

Michael Vollmer

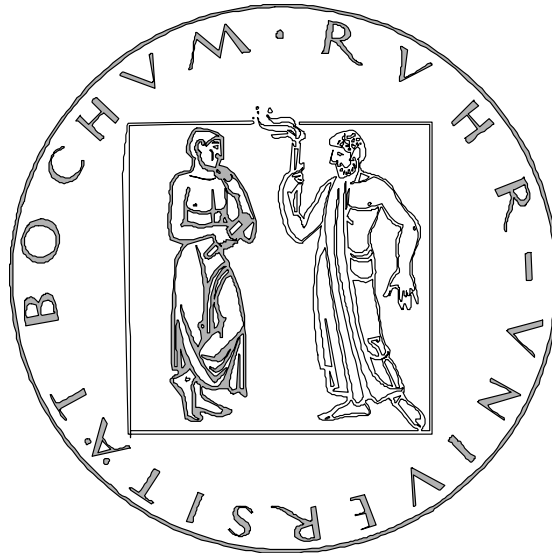
**Automatische Code-Erzeugung zur numerischen
Integration partieller Differentialgleichungen
für sicherheitskritische Anwendungen**



Cuvillier Verlag Göttingen

RUHR-UNIVERSITÄT BOCHUM

Fakultät für Elektrotechnik und Informationstechnik



**Automatische Code-Erzeugung zur numerischen Integration partieller
Differentialgleichungen für sicherheitskritische Anwendungen**

DISSERTATION

zur
Erlangung des Grades eines
Doktor-Ingenieurs

vorgelegt von
MICHAEL VOLLMER

Hamm

Bochum, den 20.04.2004

Bibliografische Information Der Deutschen Bibliothek

Die Deutsche Bibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.ddb.de> abrufbar.

1. Aufl. - Göttingen : Cuvillier, 2004
Zugl.: Bochum, Univ., Diss., 2004
ISBN 3-86537-221-X

Dissertation eingereicht:

20.04.2004

Referent:

Prof. Dr.-Ing. habil. H. D. Fischer

Korreferent:

Prof. Dr.-Ing. R. Martin

Tag der mündlichen Prüfung:

16.07.2004

© CUVILLIER VERLAG, Göttingen 2004
Nonnenstieg 8, 37075 Göttingen
Telefon: 0551-54724-0
Telefax: 0551-54724-21
www.cuvillier.de

Alle Rechte vorbehalten. Ohne ausdrückliche Genehmigung des Verlages ist es nicht gestattet, das Buch oder Teile daraus auf fotomechanischem Weg (Fotokopie, Mikrokopie) zu vervielfältigen.

1. Auflage, 2004

Gedruckt auf säurefreiem Papier

ISBN 3-86537-221-X

Vorwort

Die vorliegende Dissertation entstand während meiner Tätigkeit als wissenschaftlicher Mitarbeiter am Lehrstuhl für Nachrichtentechnik der Ruhr-Universität Bochum.

Mein besonderer Dank gilt Herrn Prof. Dr.-Ing. habil. H. D. Fischer für die Unterstützung und Betreuung meiner Arbeit. Auch danke ich ihm für die interessante Aufgabenstellung, seine wertvollen Anregungen und die sorgfältige und kritische Durchsicht dieses Manuskripts. Vor allem danke ich ihm jedoch für den mir zur Verfügung gestellten Freiraum und die angenehme Arbeitsatmosphäre.

Für die Übernahme des Korreferats und der damit verbundenen Mühe bedanke ich mich bei Herrn Prof. Dr.-Ing. R. Martin.

Weiterhin bedanke ich mich bei allen Mitarbeiterinnen und Mitarbeitern des Lehrstuhls für Nachrichtentechnik und der Arbeitsgruppe Digitale Signalverarbeitung sowie den Studien- und Diplomarbeitern und studentischen Hilfskräften, die alle zum Gelingen dieser Arbeit in freundlicher Atmosphäre beigetragen haben. Herrn Prof. Dr. sc. techn. Dr. h.c. mult. A. Fettweis danke ich für das Halten von Vorlesungen zur numerischen Integration partieller Differentialgleichungen mit dem Wellendigital-Konzept über den Zeitraum von vier Jahren. Besonders bedanken möchte ich mich bei meiner Kollegin Katrin Luhmann. Ihre wertvollen Ratschläge bei kleineren und größeren Problemen ersparten mir oft viel Zeit.

Zu guter Letzt möchte ich mich herzlich bei meiner Familie für die Unterstützung während der Erstellung meiner Arbeit bedanken.

Bochum, August 2004

MICHAEL VOLLMER

Inhaltsverzeichnis

1	Einführung	1
2	Mehrdimensionale Wellendigitalfilter	3
2.1	Abhängige und unabhängige Variablen	3
2.2	Koordinatentransformation	4
2.3	Abtastung	10
2.4	Berechnungsgebiete	13
2.5	Anordnung der Gitterpunkte und Abmessungen des Berechnungsgebietes	15
2.6	Mehrdimensionale Kirchhoff'sche Netze	15
2.7	Eigenschaften mehrdimensionaler Wellendigitalfilter	21
2.8	Wellendigitalfilter Bauelemente	21
2.8.1	Energieneutrale Bauelemente	21
2.8.2	Verallgemeinerte Verbindungsnetze	24
2.8.3	Quellen	34
2.8.4	Dissipative Bauelemente	37
2.8.5	Dynamische Bauelemente	37
2.9	Zweckmäßige Beschreibung des Wellendigitalfilters	40
2.10	Randbehandlung	44
2.11	Von einem mehrdimensionalen Wellendigitalfilter zu einem eindimensionalen	47
3	Das Engineeringsystem SPACE	53
3.1	Der Aufbau des vom Engineeringsystem SPACE erzeugten C-Codes	53
3.2	Die Einbindung der Wellendigitalfilter in das Engineeringsystem SPACE	54
3.3	Zur Wahl der Abtastperioden	57
3.4	Realisierungsaspekte	57
3.5	Verknüpfung der Gleichungen des MDWDFs mit dem Funktionsbaustein	58
4	Syntheseverfahren	65
4.1	Der behandelte Typ von PDGLn	65
4.1.1	Einschränkungen an die PDGLn	65
4.1.2	Energiebetrachtungen	66
4.1.3	Systemklassifizierung	70
4.1.4	Eingeschwungener Zustand	70
4.2	Synthese der Referenzschaltung	71
4.2.1	Vorbetrachtungen	71
4.2.2	Synthese des reaktiven Teils	73
4.2.3	Synthese des konstanten Teils	78
4.2.4	Gesamtreferenzschaltung	79
4.3	Umsetzung der elektrischen Bauelemente in WD-Elemente	80

4.4	Umsetzung der Gesamtschaltung in ein MDWDF	82
4.5	Nachweis der Berechenbarkeit	86
4.6	Alternative Methoden	88
4.7	Randbehandlung	88
4.8	Zusammenfassung	92
5	Verifikation von Software mittels formaler Methoden	93
5.1	Definitionen, Notationen und Vereinbarungen	94
5.2	Das Konzept der formalen Methoden	95
5.3	Eine axiomatische Definition von Elementaranweisungen der Programmiersprache C . . .	97
5.3.1	Einzelanweisungen	98
5.3.2	Blockanweisungen	101
5.4	Behandlung von Ausdrücken	101
5.4.1	Einzelanweisungen	102
5.4.2	Blockanweisungen	103
5.5	Einige Sätze	103
5.5.1	Skalarprodukt mit Variablen	103
5.5.2	Skalarprodukt mit Feldvariablen	103
6	Der Algorithmus	105
6.1	Gesamtanweisung und Nachbedingung	105
6.2	Berechnung der Verzögererwerte im Randfall	118
6.3	Normalberechnung der Verzögererwerte	125
6.4	Berechnung der Eingangswellen der Verzögerer	134
6.5	Berechnung der Ausgangssignale	140
6.6	Berechnung der Wellengrößen der nichtdynamischen Elemente	144
6.7	Variablendeklarationen	148
6.8	Zusammenfassung	151
7	Zusammenfassung	153
A	Die Operatoren min und max	155
B	Unitär beschränkte Matrizen	157
C	Beschränktheit der Wellengrößen	159
D	Auswirkungen endlicher Rechengenauigkeit	165
D.1	Zahlendarstellungen und Rundungsfehler	165
D.2	Idealer und realer Zweitor-Parallel-Adaptor	170
D.3	Direkte Umsetzung einer Zweitor-Adaptor-Streumatrix	177

Systematik, Formelzeichen und Abkürzungen

Systematik

In der vorliegenden Arbeit werden Vektoren mit fetten Kleinbuchstaben gekennzeichnet, Großbuchstaben in Fettdruck stehen für Matrizen, und skalare Größen werden allgemein durch Buchstaben in Normalschrift dargestellt. Variablen werden i.d.R. schräg gestellt (Schriftart *Italic*) und Konstanten werden i.d.R. gerade gesetzt (Schriftart *Roman*).

Allgemeine Festlegungen

\mathbf{M}_m^n	Matrix mit m Zeilen und n Spalten (m und n meist unterdrückt)
m_{ij}	Element in der Zeile i und der Spalte j der Matrix \mathbf{M}
$\mathbf{m}_{i,\bullet}$	Zeile i der Matrix \mathbf{M} ; Zeilenvektor
$\mathbf{m}_{\bullet,j}$	Spalte j der Matrix \mathbf{M} ; Spaltenvektor
\mathbf{M}^T	transponierte Matrix, das Element in der Zeile i und der Spalte j ist m_{ji}
\mathbf{M}^H	transjugierte Matrix, das Element in der Zeile i und der Spalte j ist m_{ji}^*
\mathbf{M}^{-T}	transponiert inverse Matrix, $\mathbf{M}^{-T} = [\mathbf{M}^T]^{-1} = [\mathbf{M}^{-1}]^T$
$\mathbf{S}_*(\mathbf{p}) = [\mathbf{S}(-\mathbf{p}^*)]^H$	Parakonjugierte der Matrix \mathbf{S}
$\mathbf{D} = \mathbf{diag}(d_1, \dots, d_m)$	Diagonalmatrix
$\overset{-j}{\underset{-i}{\mathbf{M}}}$	Matrix \mathbf{M} , in der die Zeile i und die Spalte j gestrichen wurde
$\mathbf{0}$	Nullvektor
\mathbf{e}_ν^T	Standardeinheits-Zeilenvektor $\begin{bmatrix} 1 & 2 & \dots & \nu & \dots & n \\ 0 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \end{bmatrix}$
\mathbf{e}	Vektor, deren sämtliche Koordinaten 1 sind
$\mathbf{1}_m$	Einheitsmatrix der Dimension m
$\mathbf{0}_m^n$	Nullmatrix mit m Zeilen und n Spalten (m und n meist unterdrückt)
$\lceil a \rceil$	größte ganze Zahl, die kleiner gleich a ist
$\ \mathbf{x}\ $	euklidische Vektornorm von \mathbf{x}
$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	Kronecker Produkt mit $c_{ij} = a_{ij}\mathbf{B}$
P	Bezeichner eines prädikatenlogischen Ausdrucks

S	Bezeichner eines Programmsegments
\equiv	Identität
$\sum_{\mu=N}^M \mathbf{A}_{\mu} =$	$\begin{cases} \mathbf{A}_N + \mathbf{A}_{N+1} + \cdots + \mathbf{A}_{M-1} + \mathbf{A}_M & \text{für } M \geq N \\ \mathbf{0} & \text{für } M < N \end{cases}$
$\prod_{\mu=N}^M \mathbf{A}_{\mu} =$	$\begin{cases} \mathbf{A}_N \mathbf{A}_{N+1} \cdots \mathbf{A}_{M-1} \mathbf{A}_M & \text{für } M \geq N \\ 1 & \text{für } M < N \end{cases}$
$\bigwedge_{\mu=N}^M \mathbf{P}_{\mu} \equiv$	$\begin{cases} \mathbf{P}_N \wedge \mathbf{P}_{N+1} \wedge \cdots \wedge \mathbf{P}_{M-1} \wedge \mathbf{P}_M & \text{für } M \geq N \\ \text{wahre Aussage} & \text{für } M < N \end{cases}$
$\bigwedge_{l=a}^b \mathbf{P} \equiv$	$\bigwedge_{l_1=a_1}^{b_1} \bigwedge_{l_2=a_2}^{b_2} \cdots \bigwedge_{l_n=a_n}^{b_n} \mathbf{P}$
$\mathbf{x} \geq \mathbf{0} \equiv$	$\bigwedge_{\nu=1}^n x_{\nu} \geq 0$
$\mathbf{A} \geq 0$	symbolische Darstellung für eine nicht negativ definite (n. n. d.) Matrix \mathbf{A}
$\mathbf{A} > 0$	symbolische Darstellung für eine positiv definite (p. d.) Matrix \mathbf{A}

Ausgewählte Formelzeichen

Skalare

Anzahlen des Wellendigitalfilters

n_e	Gesamtzahl der Tore der nichtdynamischen Bauelemente
$n_g = n_q + n_v + n_e$	Gesamttorzahl des Wellendigitalfilters
n_q	Anzahl der Quellen
n_v	Anzahl der dynamischen Bauelemente
n_v^{κ}	Anzahl der bzgl. t_{κ} dynamischen Bauelemente
N_e	Anzahl der nichtdynamischen Bauelemente

Anzahlen des Funktionsbausteins

n_{ea}	Anzahl der analogen Eingänge des Funktionsbausteins
----------	---

n_{aa} Anzahl der analogen Ausgänge des Funktionsbausteins

m Anzahl der Zustandsgrößen

Vektoren

abhängige Variablen des Wellendigitalfilters

\mathbf{a}_q Vektor der einfallenden Wellen der Quellen

\mathbf{a}_v Vektor der einfallenden Wellen der dynamischen Elemente

\mathbf{a}_e Vektor der einfallenden Wellen der nichtdynamischen Elemente

$\mathbf{a} = [\mathbf{a}_q^T, \mathbf{a}_v^T, \mathbf{a}_e^T]^T$ Vektor der einfallenden Wellen aller Elemente

\mathbf{b}_q Vektor der ausfallenden Wellen der Quellen

\mathbf{b}_v Vektor der ausfallenden Wellen der dynamischen Elemente

\mathbf{b}_e Vektor der ausfallenden Wellen der nichtdynamischen Elemente

$\mathbf{b} = [\mathbf{b}_q^T, \mathbf{b}_v^T, \mathbf{b}_e^T]^T$ Vektor der ausfallenden Wellen aller Elemente

$b_{v\mu}^\kappa$ Koordinate μ des Vektors \mathbf{b}_v^κ

Vektoren des Funktionsbausteins

\mathbf{a}_{FB} Vektor der Ausgangssignale des Funktionsbausteins

\mathbf{e}_{FB} Vektor der Eingangssignale des Funktionsbausteins

\mathbf{m}_{FB} Vektor der Zustandsgrößen des Funktionsbausteins

Matrizen

Matrizen des Wellendigitalfilters

\mathbf{P} Matrix, die die Verknüpfungen der einzelnen Tore des Wellendigitalfilters beschreibt

$$\mathbf{P} = \begin{matrix} & \begin{matrix} q & v & e \end{matrix} \\ \begin{matrix} q \\ v \\ e \end{matrix} & \begin{bmatrix} \mathbf{P}_{qq} & \mathbf{P}_{qv} & \mathbf{P}_{qe} \\ \mathbf{P}_{vq} & \mathbf{P}_{vv} & \mathbf{P}_{ve} \\ \mathbf{P}_{eq} & \mathbf{P}_{ev} & \mathbf{P}_{ee} \end{bmatrix} \end{matrix}$$

\mathbf{P}_b Permutationsmatrix zur Festlegung der Berechnungsreihenfolge

\mathbf{S} Streumatrix der nichtdynamischen Bauelemente

unabhängige Variablen Referenzschaltung/Wellendigitalfilter

$D_{\kappa'} = \frac{\partial}{\partial t_{\kappa'}}$	Differentialoperatoren nach den neuen Koordinaten
$\mathbf{D}_t = [D_1, \dots, D_{k'}]^T$	Gradient in den Koordinaten \mathbf{t}
$\mathbf{D}_x = [\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_{k-1}}, \frac{1}{v_k} \frac{\partial}{\partial t}]^T$	Gradient in den Koordinaten \mathbf{x}
$\mathbf{D} = [\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_{k-1}}]^T$	Gradient in den Ortskoordinaten
\mathbf{H} , $\mathbf{x} = v_0 \mathbf{H} \mathbf{t}$	Transformationsmatrix
k	Dimension des ursprünglichen Koordinatensystems
k'	Dimension des neuen Koordinatensystems
$\boldsymbol{\mu} = [\mu_1, \dots, \mu_k]^T$	Diskrete Variablen des ursprünglichen Koordinatensystems
$\Delta \boldsymbol{\mu} = [\Delta \mu_1, \dots, \Delta \mu_k]^T$	Verschiebe-Vektor in den diskreten Variablen des ursprünglichen Koordinatensystems
$\boldsymbol{\nu} = [\nu_1, \dots, \nu_{k'}]^T$	Diskrete Variablen des neuen Koordinatensystems
$\mathbf{t} = [t_1, \dots, t_{k'}]^T$	Vektor mit den unabhängigen Variablen des neuen Koordinatensystems
$\Delta \mathbf{t} = [\Delta t_1, \dots, \Delta t_{k'}]^T$	Verschiebe-Vektor des neuen Koordinatensystems
\mathbf{T}_A	Abtastmatrix des neuen Koordinatensystems
$\mathbf{x} = [x_1, \dots, x_{k-1}, v_k t]^T$	Vektor mit den unabhängigen Variablen des ursprünglichen Koordinatensystems
$_{-k} \mathbf{x} = [x_1, \dots, x_{k-1}]^T$	Vektor mit den unabhängigen Ortskoordinaten des ursprünglichen Koordinatensystems (Echtraumkoordinaten)
\mathbf{X}_A	Abtastmatrix des ursprünglichen Koordinatensystems
$\Delta \mathbf{x} = [\Delta x_1, \dots, v_k \Delta t]^T$	Verschiebe-Vektor des ursprünglichen Koordinatensystems

Abmessungen

\mathcal{G}	Berechnungsgebiet (nur Ortskoordinaten)
\mathcal{G}_x	Gebiet im ursprünglichen Koordinatensystem
\mathcal{G}_t	Gebiet im neuen Koordinatensystem

\mathcal{G}_0	Gebiet außerhalb des Berechnungsgebietes
$\partial\mathcal{G}$	Rand des Berechnungsgebietes
\mathcal{G}_n	Teilberechnungsgebiet $n > 0$
l_κ	Länge des Berechnungsgebietes in Richtung x_κ
$l_{\kappa n}$	Länge des Teilberechnungsgebietes n in Richtung x_κ
N	Anzahl der Teilberechnungsgebiete
N_κ	Anzahl der Teilberechnungsgebiete in Richtung der unabhängigen Variablen x_κ
P	Anzahl der Abtastpunkte einer Abtastschicht
P_{x_κ}	Anzahl der Abtastpunkte in Richtung x_κ
$P_{x_{\kappa n}}$	Anzahl der Abtastpunkte in Richtung x_κ des Teilberechnungsgebietes n
ΔT_κ	Abstand zweier Abtastpunkte in Richtung t_κ
$x_{\kappa \min}$	untere Grenze des Berechnungsgebietes in Richtung x_κ
$x_{\kappa \max}$	obere Grenze des Berechnungsgebietes in Richtung x_κ
ΔX_κ	Abstand zweier Abtastpunkte in Richtung x_κ

Kapitel 1

Einführung

Ein Großteil naturwissenschaftlicher und technischer Prozesse wird mithilfe von partiellen Differentialgleichungen (PDGLn) beschrieben. Da für PDGLn i.Allg. keine Verfahren zur analytischen Lösung bekannt sind, setzt man numerische Verfahren zur Berechnung von Näherungslösungen mittels eines Digitalrechners ein. Sowohl die weiter ansteigende Rechenleistung der Digitalrechner, als auch neue, effizientere Algorithmen erschließen weitere Anwendungsgebiete dieses wachsenden Teilbereichs der Mathematik. Die numerischen Verfahren sollten dabei die wesentlichen Eigenschaften des ursprünglichen Systems auf das digitale System übertragen, insbesondere die Stabilität, die Lokalität und die Passivität. Die Übertragung dieser Eigenschaften leisten die von Fettweis entwickelten Wellendigitalfilter (WDF), die ursprünglich zur digitalen Nachbildung analoger Filter dienten [Fett70], [Fett86]. Es hat sich gezeigt, dass diese Filter hervorragende Stabilitätseigenschaften auch unter realen Bedingungen (endliche Wortlänge) besitzen. Diese günstigen Eigenschaften übertragen sich auch auf die mehrdimensionalen Wellendigitalfilter. Die mehrdimensionalen Wellendigitalfilter sind digitale Nachbildungen analoger (auch nichtlinearer) mehrdimensionaler Kirchhoffscher Schaltungen und wurden erstmals in [Fisc84] zur numerischen Integration von Differentialgleichungen vorgeschlagen. Bekanntermaßen kann eine mehrdimensionale Kirchhoffsche Schaltung äquivalent durch partielle Differentialgleichungen beschrieben werden. Somit lässt sich die Wellendigitalmethode, welche in dieser Arbeit ausschließlich verwendet wird, zur numerischen Integration von partiellen Differentialgleichungen nutzen, [FN90a], [FN90b], [FN91a], [FN91b], [Nits93], [Heme95], [Feld95], [Frie95], [Krau97], [Pott98]. Mittlerweile existiert eine Vielzahl weiterer Arbeiten, die über die erfolgreiche Anwendung des Verfahrens berichten. Eine zusammenfassende Darstellung in englischer Sprache findet sich in [Bilb01] und [Bilb04]. Die praktische Anwendbarkeit der Wellendigitalmethode erfordert aber die effiziente Erzeugung eines Codes zur Simulation des Wellendigitalfilters. Der Ausgangspunkt der Erzeugung eines Codes sollte das Ausgangsproblem selber sein, also die zu lösende PDGL. In dieser Arbeit wird ein Verfahren zur automatischen Codeerzeugung vorgestellt, welches auf PDGLn anwendbar ist, die lineare, zeitinvariante, symmetrisch hyperbolische Systeme beschreiben.

Eine automatische Codeerzeugung hat nicht nur den Vorteil geringeren Herstellungsaufwands, sondern liefert auch -geeignete Implementierung vorausgesetzt- eine höhere Qualität der erzeugten Software, da der Anteil der manuellen Software-Herstellung reduziert wird. Es eröffnet sich zudem die Möglichkeit des Einsatzes der Algorithmen in sicherheitskritischen Bereichen. Wie oben bereits angedeutet, sind für eine derartige Anwendung auch strengere Qualitätsansprüche an die Implementierung zu stellen. Um diesen nachzukommen, verwenden wir die so genannten formalen Methoden der Softwaretechnik zur Verifikation der entwickelten Codes gegenüber der Spezifikation. Die Theorie der Programmverifikation wurde durch McCarthy angeregt, [McCa62], [McCa63]. Seine Intention war es, die gewünschten Eigenschaften von Programmen mittels mathematischer Methoden nachzuweisen, anstatt die Programme durch Testläufe auf Fehlerlosigkeit zu prüfen. Floyd nahm die Idee auf und schlug ein Verfahren vor mit dem ein gegebenes Software-System analysiert werden konnte, [Floy67]. Dieses Konzept wurden in

[Dijk68] erweitert und zur Synthese beweisbarer korrekter Programme genutzt. In den weiteren Jahren folgten Arbeiten, die im Wesentlichen auf den zuvor genannten Artikeln aufbauen. Hervorzuheben hieraus ist [Hoar69], in der eine Programmanweisung bzw. eine zusammengesetzte Anweisung als Transformation von Prädikaten aufgefasst wird. Weiterhin sind dort die wichtigsten Semantikregeln zu finden. Auf dieser Basis wird auch der in dieser Arbeit angegebene Beweis geführt.

Die wesentlichen Ziele dieser Arbeit sind zum einen die Entwicklung eines Verfahrens zur Synthese einer Referenzschaltung und anschließender Umsetzung in ein mehrdimensionales Wellendigitalfilter, welches die formale Spezifikation der Software festlegt. Zum anderen wird ein Code angegeben, der das mehrdimensionale Wellendigitalfilter simuliert und mittels eines formalen Korrektheitsbeweises gegenüber der Spezifikation verifiziert. Der Code soll dabei so beschaffen sein, dass er sich in das -für Steuerungs- und Regelungsfunktionen mit sicherheitsrelevanter Bedeutung in Kernkraftwerken entwickelte- digitale Leitsystem TELEPERM XS einbinden lässt.

Intention des Kapitels 2 ist es, die notwendigen Grundlagen der Theorie mehrdimensionaler Wellendigitalfilter zu rekapitulieren und zwar in einer auf unsere Aufgabenstellung angepassten Form.

Anschließend erfolgt im Kapitel 3 eine Einarbeitung in das Programmpaket SPACE und zwar einerseits aus Sicht des Benutzers (i.d.R. der Leittechniker) und andererseits aus Sicht des Entwicklers. Nach Analyse und Darstellung der erarbeiteten Erkenntnisse wird die Möglichkeit der Einbindung der Wellendigitalfilter in das Programmpaket SPACE untersucht. Dabei auftretende Probleme und deren Lösung werden aufgezeigt.

Kapitel 4 bildet den ersten Hauptteil dieser Arbeit. Dort werden wir zunächst die Klasse der in dieser Arbeit behandelten Systeme einschränken und daraus Eigenschaften der PDGLn ableiten. Im Anschluss daran werden wir das neu entwickelte Syntheseverfahren zur Gewinnung einer mehrdimensional passiven Referenzschaltung vorstellen. Ferner werden wir zu der systematisch gewonnenen Referenzschaltung ein mehrdimensionales Wellendigitalfilter angeben und deren Berechenbarkeit aufzeigen.

In Kapitel 5 werden die aus der Literatur bekannten Grundlagen zur Anwendung formaler Methoden in der Softwaretechnik behandelt. Dabei werden wir uns auf die Inhalte beschränken, die für den weiteren Verlauf der Arbeit relevant sind. Zudem werden in diesem Kapitel die für den Algorithmus notwendigen Programmanweisungen axiomatisch definiert.

Kapitel 6 beinhaltet den zweiten Hauptteil dieser Arbeit. In diesem Kapitel werden wir den Übergang zur formalen Spezifikation des Algorithmus durchführen. Zudem wird die Implementierung des Algorithmus in der Programmiersprache C unter ausschließlicher Verwendung der axiomatisch definierten Programmanweisungen durchgeführt. Weiterhin führen wir den formalen Korrektheitsbeweis in diesem Kapitel. In formalen Korrektheitsbeweisen für Algorithmen der Signalverarbeitung kommt der Verhinderung eines Überlaufs des Darstellungsbereiches eine besondere Bedeutung zu, die ebenfalls in diesem Kapitel Berücksichtigung findet.

Die Ergebnisse der Kapitel 4 und 6 werden mit eigenen Zusammenfassungen gewürdigt. Die Zusammenfassung der gesamten Arbeit befindet sich im Kapitel 7.

Kapitel 2

Mehrdimensionale Wellendigitalfilter

In diesem Kapitel werden wir die für das Verständnis der Arbeit notwendigen Grundlagen der Theorie mehrdimensionaler Wellendigitalfilter zusammenfassen. Die Zusammenfassung der Grundlagen beginnt mit der Einführung von unabhängigen und abhängigen Variablen, wobei die unabhängigen Variablen danach einer Koordinatentransformation unterworfen und diskretisiert werden. Anschließend werden die Form des Berechnungsgebietes, also der Bereich der unabhängigen Variablen für die die Lösung der PDGL erfolgen soll, festgelegt. In diesem Zusammenhang wird zudem die Lage der Abtastpunkte diskutiert.

Nachdem in den ersten Unterkapiteln die unabhängigen Variablen im Vordergrund standen, widmet sich der weitere Teil den abhängigen Variablen. Diese abhängigen Variablen einer PDGL werden als Spannungen und Ströme von Toren eines Kirchhoff'schen Netzes interpretiert. Ähnlich den unabhängigen Variablen erfolgt bei den abhängigen Variablen eine Koordinatentransformation. Diese Transformation vollzieht sich in der Form, dass anstelle von Spannung und Strom Wellengrößen genutzt werden. Von zentraler Bedeutung für Stabilitätsfragen Kirchhoff'scher Netze sind die energetischen Eigenschaften ihrer Bauelemente. Aufbauend auf der torweisen Betrachtung der Bauelemente, werden aus Energiebetrachtungen heraus Eigenschaften wie MD-Passivität und MD-Energieneutralität erläutert. Gemäß diesen Eigenschaften werden dann typische Bauelemente der Theorie elektrischer Netze, die in dieser Arbeit Verwendung finden, eingeführt, qualifiziert und in MDWDF-Bauelemente überführt. Auf den Nachweis der Übertragung der Eigenschaften der elektrischen Bauelemente auf die WD-Elemente verzichten wir und verweisen auf die existierende Literatur, [Meer79], [MF92], [Fett92].

Weiterhin wird von den bekannten Verfahren zur Randwertbehandlung eines, welches sich besonders gut in die Zielsetzung dieser Arbeit einpasst und auch später Verwendung findet, ausführlich erläutert und kompakt beschrieben.

Ergänzend wird für die torweise Verschaltung der einzelnen WD-Elemente eine übersichtliche Darstellung, in Form eines Differenzen-Gleichungssystems, angegeben und mit dessen Hilfe die Berechnungsreihenfolge der einzelnen Wellengrößen systematisch bestimmt.

Zum Abschluss wird eine Transformation der unabhängigen Variablen dargelegt, die es ermöglicht, die endliche Anzahl örtlicher Abtastpunkte auf eine natürliche Zahl eineindeutig abzubilden. Dieses Konzept wird später bei der Einbindung der MDWDF in das Programmpaket SPACE angewendet.

Um eine einfachere Recherche zu ermöglichen, findet man die Quellenangaben in den einzelnen Unterkapiteln.

2.1 Abhängige und unabhängige Variablen

Die Variablen eines Differentialgleichungssystems unterscheiden wir nach abhängigen und unabhängigen Größen. Die unabhängigen Größen sind die Ortsvariablen und die Zeit. Die abhängigen Größen werden

wir als die Feldgrößen bezeichnen. Beispiele hierfür sind Druck oder Temperatur. Die Feldgrößen sind proportional zu einer Zweigspannung oder einem Zweigstrom der Referenzschaltung. Wir fassen die unabhängigen Variablen in dem Vektor

$$\mathbf{x} = [x_1, \dots, x_{k-1}, v_k t]^T, \quad v_k > 0, \quad (2.1)$$

zusammen. Hierin ist t die physikalische Zeit und die restlichen Variablen sind beispielsweise physikalische Ortskoordinaten. Die mit der Einheit einer Geschwindigkeit behaftete Größe v_k kann eine Funktion der Zeit sein, wird aber in dieser Arbeit als konstant vorausgesetzt. Somit sind die Einheiten aller Koordinaten des Vektors \mathbf{x} gleich. Mit $\mathbf{D}\mathbf{x}$ bezeichnen wir den Vektor, der die partiellen Ableitungen bzgl. x_k enthält, d. h.

$$\mathbf{D}\mathbf{x} = \left[\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_k} \right]^T = \left[\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_{k-1}}, \frac{1}{v_k} \frac{\partial}{\partial t} \right]^T. \quad (2.2)$$

Weiterhin definieren wir einen Vektor, der als Koordinaten die komplexen Frequenzen enthält

$$\mathbf{p}_\mathbf{x} = [p_{x_1}, \dots, p_{x_k}]^T, \quad p_{x_k} = p_t / v_k. \quad (2.3)$$

2.2 Koordinatentransformation

Die in der Theorie der mehrdimensionalen Wellendigitalfilter oft verwendete Koordinatentransformation werden wir zunächst durch Überlegungen zur Passivität¹ stützen. Die wesentlichen Quellen sind [FN91b] und [Nits93]. Im eindimensionalen Fall bezieht sich die Passivität nur auf die Zeit als unabhängige Variable. In der klassischen Literatur zu mehrdimensionalen elektrischen Schaltungen und in der mehrdimensionalen Signalverarbeitung werden alle unabhängigen Variablen als untereinander gleichberechtigt aufgefasst [Koga69], [Rao69], [Bose79], [Fett79], [Huan81], [Bose82], [FB87], [BF87], [Kumm88]. Die Passivität muss bezüglich aller unabhängigen Variablen gewährleistet sein. Diese Passivität wird auch als mehrdimensionale Passivität bezeichnet. Die durch Energiebetrachtungen abgeleitete Definition der Passivität physikalischer Systeme bezieht sich hingegen nur auf die Zeit. Die Ortskoordinaten finden hierbei keine Berücksichtigung. Insofern ist zu erwarten, dass eine mehrdimensionale elektrische Schaltung, die aus einer partiellen Differentialgleichung eines passiven Systems direkt hergeleitet wurde, i.Allg. nicht mehrdimensional passiv ist.

Nichtsdestotrotz ist es wünschenswert, eine mehrdimensional passive elektrische Schaltung aus der partiellen Differentialgleichung eines passiven Systems anzugeben. Dies hat den Vorteil, dass die Erkenntnisse der mehrdimensionalen elektrischen Schaltungen und der mehrdimensionalen Signalverarbeitung genutzt werden können. Um dieses Ziel zu erreichen, führen wir eine Transformation der unabhängigen Variablen durch.

Im engen Zusammenhang mit der Passivität steht die Kausalität. Im eindimensionalen Fall ist die Kausalität so definiert, dass der aktuelle Zustand nur von Zuständen zurückliegender Zeitpunkte und vom Eingangssignal des aktuellen Zeitpunktes und zurückliegender Zeitpunkte abhängt.

In der mehrdimensionalen Signalverarbeitung wird ebenfalls eine Ablauffrichtung festgelegt. Diese Ablauffrichtung teilt das Gebiet der unabhängigen Variablen in einen Abhängigkeitsbereich (Ursachenbereich), einen Wirkungsbereich und einen Bereich, der zu keinem der beiden gehört. Bei einer bestimmten Ablauffrichtung könnte beispielsweise der Zustand im Punkt \mathbf{t}_0 nur von Zuständen und Eingangssignalen des Abhängigkeitsbereiches $\mathbf{t} \leq \mathbf{t}_0$ abhängen. Der Zustand im Punkt \mathbf{t}_0 hingegen beeinflusst nur Zustände im Wirkungsbereich $\mathbf{t} \geq \mathbf{t}_0$.

¹Die Passivitätsbegriffe werden im Kapitel 2.7 genau definiert.

Bei physikalischen Systemen bezieht sich die Kausalität hingegen nur auf die Zeit. Ein Zusammenhang zu dem gerade beschriebenen Sachverhalt in der mehrdimensionalen Signalverarbeitung ist daher nicht unmittelbar erkennbar. Insbesondere existiert kein Bereich, der nicht dem Wirkungs- oder dem Abhängigkeitsbereich zugeordnet werden kann. Allerdings zeigen sich Parallelen auf, wenn das System nur lokale Abhängigkeiten besitzt. In dem Fall bilden sich durch die endliche Ausbreitungsgeschwindigkeit des Vorgangs auch Bereiche aus, die weder dem Abhängigkeits- noch dem Wirkungsbereich zugeordnet werden können. Wir halten als Fazit dieser Überlegungen fest, dass die Darstellung eines physikalischen Systems mittels einer mehrdimensional passiven Schaltung eine endliche Ausbreitungsgeschwindigkeit voraussetzt. Im weiteren Verlauf dieses Kapitels werden wir sehen, welchen Einfluss die maximale Ausbreitungsgeschwindigkeit auf die Koordinatentransformation hat.

Da wir jeder der neuen Variablen die Bedeutung einer Zeit beimessen wollen, wählen wir als Bezeichner für die Variablen des neuen Koordinatensystems

$$\mathbf{t} = [t_1, \dots, t_{k'}]^T. \quad (2.4)$$

Die Variablen \mathbf{t} sind über

$$\mathbf{x} = v_0 \mathbf{H} \mathbf{t}, \quad v_0 > 0 \quad (2.5)$$

mit den ursprünglichen unabhängigen Variablen verknüpft ($\dim\{\mathbf{H}\} = k \times k'$). Die Variable v_0 hat die Einheit einer Geschwindigkeit und wird in dieser Arbeit ebenfalls als konstant vorausgesetzt. Der Vektor \mathbf{t} findet dann als eine mehrdimensionale Zeit eine willkommene Interpretation. Die Dimension des neuen Koordinatensystems soll nicht kleiner als die Dimension des alten Koordinatensystems sein, d.h. $k' \geq k$. Wir sprechen von einer verallgemeinerten Koordinatentransformation, wenn das neue Koordinatensystem eine höhere Dimension als das ursprüngliche Koordinatensystem hat. Weiterhin fordern wir, dass \mathbf{H} zeilenregulär ist, was zur Folge hat, dass der Rang von \mathbf{H} gleich k ist.

Nach den eingangs gemachten Bemerkungen bzgl. der Kausalität erscheint es sinnvoll, die Koordinatentransformation so durchzuführen, dass der Abhängigkeitsbereich und der Wirkungsbereich der mehrdimensionalen Zeit auf den Abhängigkeitsbereich und den Wirkungsbereich der physikalischen Zeit abgebildet werden. Zudem sollten Abhängigkeits- und Wirkungsbereich der Zeit auf den Abhängigkeits- und den Wirkungsbereich der mehrdimensionalen Zeit abgebildet werden. Aufgrund der i. Allg. rechteckförmigen Matrix \mathbf{H} kann nicht von einer bijektiven Abbildung gesprochen werden. Die bislang an die Koordinatentransformation gestellten Forderungen können wie folgt formuliert werden :

- Nimmt keine der Variablen des neuen Koordinatensystems ab und mindestens eine der neuen Variablen zu, so muss dies eine Zunahme der Zeit t bewirken, d. h.

$$[\Delta t_\kappa \geq 0 \text{ für } \kappa = 1 \dots k'] \wedge [\text{mind. ein } \Delta t_\kappa > 0] \implies \Delta t > 0. \quad (2.6)$$

- Nimmt die Zeit t zu und sind die restlichen Variablen des ursprünglichen Koordinatensystems konstant, so muss dies eine Zunahme jeder Variablen des neuen Koordinatensystems bewirken, d. h.

$$\Delta t > 0 \wedge \Delta x_\kappa = 0 \text{ für } \kappa = 1 \dots k - 1 \implies \Delta t_\kappa > 0 \text{ für } \kappa = 1 \dots k'. \quad (2.7)$$

Die Forderungen erscheinen notwendig für die gewünschte mehrdimensionale Passivität zu sein, falls das System passiv bzgl. t ist. Aus diesen Forderungen resultieren Bedingungen an die Elemente der Matrix \mathbf{H} . Wir werden später genauer darauf eingehen.

Normiert man den letzten Zeilenvektor von \mathbf{H} auf eine euklidische Norm von 1 und wählt seine Koordinaten gleich, d. h. zu $\sqrt{1/k'}$, so lässt sich zeigen (siehe z.B. [Fett99]), dass die Ungleichung

$$v_k \geq \sqrt{k-1} v_{\max} \quad (2.8)$$

notwendig und hinreichend zur Erfüllung der Forderungen Gleichung (2.6) und Gleichung (2.7) ist. Die Geschwindigkeit v_{\max} stellt die maximal mögliche Ausbreitungsgeschwindigkeit des physikalischen Vorgangs dar.

Die Differentialoperatoren im neuen Koordinatensystem sind in dem Vektor

$$\mathbf{D}_{\mathbf{t}} = \left[\frac{\partial}{\partial t_1}, \dots, \frac{\partial}{\partial t_{k'}} \right]^T \quad (2.9)$$

zusammengefasst. Wir werden im Folgenden die Beziehung zwischen den Ableitungen zu bestimmen haben und gehen dazu von der Kettenregel

$$\mathbf{J}_{\mathbf{t}}^{\mathbf{f}} = \mathbf{J}_{\mathbf{g}}^{\mathbf{f}} \cdot \mathbf{J}_{\mathbf{t}}^{\mathbf{g}}, \quad (2.10)$$

aus, die der Differentiation der beliebigen differenzierbaren Funktion $\mathbf{f}(\mathbf{g}(\mathbf{t}))$ nach \mathbf{t} entspricht. Die drei Matrizen sind die Jacobi-Matrizen, z.B.

$$\mathbf{J}_{\mathbf{t}}^{\mathbf{f}} = \left[\frac{\partial \mathbf{f}}{\partial t_1}, \dots, \frac{\partial \mathbf{f}}{\partial t_{k'}} \right] = [\mathbf{D}_{\mathbf{t}} \mathbf{f}^T]^T. \quad (2.11)$$

In unserem Fall gilt $\mathbf{g} = \mathbf{x} = v_0 \mathbf{H} \mathbf{t}$ und \mathbf{f} ist eine skalare Funktion. Die Differentiation der skalaren Funktion f nach \mathbf{t} lautet somit

$$\mathbf{J}_{\mathbf{t}}^{\mathbf{f}} = \left[\frac{\partial f}{\partial t_1}, \dots, \frac{\partial f}{\partial t_{k'}} \right] = \underbrace{\left[\frac{\partial}{\partial t_1}, \dots, \frac{\partial}{\partial t_{k'}} \right]}_{=\mathbf{D}_{\mathbf{t}}^T} f = \underbrace{\left[\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_k} \right]}_{=\mathbf{D}_{\mathbf{x}}^T} f \cdot \underbrace{\left[\frac{\partial \mathbf{x}}{\partial t_1}, \dots, \frac{\partial \mathbf{x}}{\partial t_{k'}} \right]}_{=\mathbf{J}_{\mathbf{t}}^{\mathbf{x}}}. \quad (2.12)$$

Die letzte Jacobi-Matrix berechnet man mit $\mathbf{x} = v_0 \mathbf{H} \mathbf{t}$ und

$$\frac{\partial \mathbf{x}}{\partial t_{\nu}} = \frac{\partial}{\partial t_{\nu}} v_0 \mathbf{H} \mathbf{t} = v_0 \mathbf{H} \frac{\partial \mathbf{t}}{\partial t_{\nu}} = v_0 \mathbf{H} \mathbf{e}_{\nu} \quad (2.13)$$

zu

$$\mathbf{J}_{\mathbf{t}}^{\mathbf{x}} = v_0 \mathbf{H}. \quad (2.14)$$

Die Transposition von Gleichung (2.12) liefert letztendlich das gewünschte Ergebnis

$$\mathbf{D}_{\mathbf{t}} f = v_0 \mathbf{H}^T \mathbf{D}_{\mathbf{x}} f, \quad (2.15)$$

wobei wir im Folgenden nur die Operatoren ohne die Funktion f verwenden werden, d.h.

$$\mathbf{D}_{\mathbf{t}} = v_0 \mathbf{H}^T \mathbf{D}_{\mathbf{x}}. \quad (2.16)$$

Wir werden im Folgenden Bedingungen an die Matrix \mathbf{H} aus der ersten Forderung an die Koordinatentransformation herleiten. Aus dem totalen Differential

$$dx_{\kappa} = \sum_{\mu=1}^{k'} \frac{\partial x_{\kappa}}{\partial t_{\mu}} dt_{\mu} \quad (2.17)$$

ergibt sich näherungsweise für die Zuwächse

$$\Delta x_{\kappa} = \sum_{\mu=1}^{k'} \frac{\partial x_{\kappa}}{\partial t_{\mu}} \Delta t_{\mu}. \quad (2.18)$$

Die Ableitung $\frac{\partial x_\kappa}{\partial t_\mu}$ ist das Element der Zeile κ und der Spalte μ der Jacobi-Matrix \mathbf{J}_t^x und lautet $\frac{\partial x_\kappa}{\partial t_\mu} = v_0 h_{\kappa\mu}$. Zusammenfassen von Gleichung (2.18) für alle κ in Vektorform liefert

$$\Delta \mathbf{x} = \mathbf{J}_t^x \Delta \mathbf{t} = v_0 \mathbf{H} \Delta \mathbf{t} . \quad (2.19)$$

Aus der ersten Forderung Gleichung (2.6) folgt, dass die letzte Zeile von \mathbf{H} positive Elemente besitzen muss.

Es empfiehlt sich, auch für die neuen Koordinaten einen Vektor der komplexen Frequenzen einzuführen

$$\mathbf{p}_t = [p_1, \dots, p_{k'}]^T . \quad (2.20)$$

Dieser Vektor soll der Beziehung $\mathbf{p}_x^T \mathbf{x} = \mathbf{p}_t^T \mathbf{t}$ genügen, was durch die Interpretation von Gleichung (2.16) als Operatorgleichung sichergestellt ist. Bemerkenswert ist an dieser Stelle, dass sich für $\text{Re}\{p_{x_\kappa}\} = 0$, $\kappa = 1 \dots k-1$ die linke offene p_t -Halbebene in die linke offene \mathbf{p}_t -Polyhalbebene abbildet, d. h.

$$\text{Re}\{p_t\} < 0 \wedge \text{Re}\{p_{x_\kappa}\} = 0 \text{ für } \kappa = 1 \dots k-1 \implies \text{Re}\{\mathbf{p}_t\} < \mathbf{0} . \quad (2.21)$$

Nachdem die alten Koordinaten durch die neuen Koordinaten und die neuen Ableitungen durch die alten Ableitungen ausgedrückt wurden, sollen nun die umgekehrten Beziehungen hergeleitet werden.

Wir bezeichnen eine Matrix \mathbf{H}^{-R} als die Rechtsinverse von \mathbf{H} , wenn sie eine Lösung der Gleichung

$$\mathbf{H} \mathbf{H}^{-R} = \mathbf{1}_k . \quad (2.22)$$

ist.

Von den Rechtsinversen gibt es im Falle $k' > k$ unendlich viele. Genauer gesagt, der Freiheitsgrad jedes Spaltenvektors von \mathbf{H}^{-R} ist $k' - k$. Geht man davon aus, dass Gleichung (2.22) ein konsistentes Gleichungssystem ist, so lautet die allgemeine Lösung des Gleichungssystems

$$\mathbf{H}^{-R} = \mathbf{H}_1^S + [\mathbf{1}_{k'} - \mathbf{H}_1^S \mathbf{H}] \mathbf{M}_1 \quad (2.23)$$

wobei \mathbf{M}_1 eine beliebige Matrix konsistenter Dimension ist, und die so genannte Semiinverse \mathbf{H}_1^S ist eine Matrix, die der Gleichung

$$\mathbf{H} = \mathbf{H} \mathbf{H}_1^S \mathbf{H} \quad (2.24)$$

genügt. $[\mathbf{1}_{k'} - \mathbf{H}_1^S \mathbf{H}] \mathbf{M}_1$ beschreibt die homogene Lösung von Gleichung (2.22). \mathbf{H}_1^S ist hingegen die inhomogene Lösung. Sie ist allerdings i. Allg. durch $\mathbf{H} = \mathbf{H} \mathbf{H}_1^S \mathbf{H}$ allein nicht eindeutig bestimmt. Wir fordern zusätzlich $\mathbf{H}_{1,2}^S \mathbf{H} = [\mathbf{H}_{1,2}^S \mathbf{H}]^T$. Zudem gelten $\mathbf{H} \mathbf{H}_1^S = [\mathbf{H} \mathbf{H}_1^S]^T$ und $\mathbf{H}_1^S = \mathbf{H}_1^S \mathbf{H} \mathbf{H}_1^S$, da \mathbf{H} eine zeilenreguläre Matrix ist. Die Semiinverse $\mathbf{H}_{1,2}^S$ ist somit zur Moore-Penrose-Inversen \mathbf{H}^+ festgelegt, die bekanntermaßen eindeutig ist und sich im Fall einer zeilenregulären Matrix \mathbf{H} zu $\mathbf{H}^+ = \mathbf{H}^T [\mathbf{H} \mathbf{H}^T]^{-1}$ berechnet. Der volle Lösungsraum von \mathbf{H}^{-R} bleibt aber trotz dieser Festlegung erhalten.

Wir fragen uns nun, welche Dimension der Raum hat, den die Spaltenvektoren der Matrix $[\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}]$ aufspannen. Zunächst stellen wir fest, dass $\mathbf{H}^+ \mathbf{H}$ eine idempotente Matrix vom Rang k ist. Sie hat k Eigenwerte bei 1 und $k' - k$ Eigenwerte bei 0. Da zur Einheitsmatrix jeder Vektor ein Eigenvektor ist, sind die Eigenvektoren von $\mathbf{H}^+ \mathbf{H}$ auch Eigenvektoren der Matrix $[\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}]$. Somit gilt für die Eigenwerte $\lambda\{[\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}]\} = 1 - \lambda\{[\mathbf{H}^+ \mathbf{H}]\}$. Die Matrix $[\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}]$ hat also k Eigenwerte bei 0 und $k' - k$ Eigenwerte bei 1 und der aufgespannte Spaltenraum hat damit die Dimension $k' - k$, was dem Spaltendefekt von \mathbf{H} entspricht.

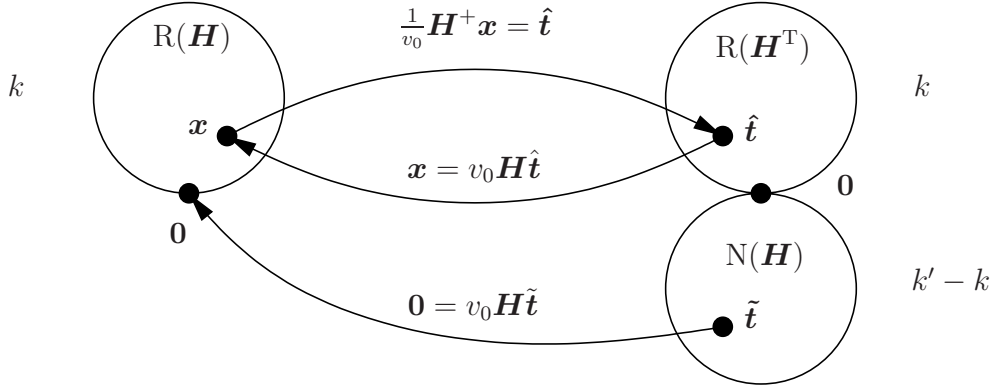


Bild 2.1: Koordinatentransformation

Nach diesen Vorbetrachtungen werden wir nun die Eigenschaften der Koordinatentransformation diskutieren. Setzt man in $\mathbf{x} = v_0 \mathbf{H} \mathbf{t}$ den Vektor

$$\mathbf{t} = \underbrace{\frac{1}{v_0} \mathbf{H}^{-R} \mathbf{x}}_{\hat{\mathbf{t}}} + \underbrace{[\mathbf{1}_{k'} - \mathbf{H}^{-R} \mathbf{H}]}_{\tilde{\mathbf{t}}} \tilde{\mathbf{t}}, \quad \tilde{\mathbf{t}} \text{ beliebig} \quad (2.25)$$

ein, so findet man

$$\mathbf{x} = v_0 \mathbf{H} \mathbf{t} = \mathbf{H} \mathbf{H}^{-R} \mathbf{x} + v_0 \underbrace{[\mathbf{H} - \mathbf{H} \mathbf{H}^{-R} \mathbf{H}]}_{=0} \tilde{\mathbf{t}} = \mathbf{x}, \quad (2.26)$$

d.h. jeder über Gleichung (2.25) aus \mathbf{x} gebildete Vektor \mathbf{t} wird über Gleichung (2.5) wieder eindeutig auf den ursprünglichen Vektor \mathbf{x} abgebildet. Dies heißt aber nicht, dass der Vektor \mathbf{t} zu einem Vektor \mathbf{x} eindeutig ist, denn zum einen ist die Rechtsinverse \mathbf{H}^{-R} nicht eindeutig, und zum anderen ist $\tilde{\mathbf{t}}$ ein beliebiger Vektor konsistenter Dimension. Als Konsequenz daraus halten wir fest, dass sich zwar die alten Koordinaten eindeutig aus den neuen gewinnen lassen, aber mehrere verschiedene neue Koordinaten auf den gleichen Vektor im alten Koordinatensystem führen. Setzen wir in Gleichung (2.25) für \mathbf{H}^{-R} die allgemeine Lösung (Gleichung (2.23)) mit $\mathbf{H}_1^S = \mathbf{H}^+$ ein, so erhalten wir

$$\mathbf{t} = \underbrace{\frac{1}{v_0} \mathbf{H}^+ \mathbf{x}}_{\hat{\mathbf{t}}} + [\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}] \underbrace{\left[\frac{1}{v_0} \mathbf{M}_1 \mathbf{x} + \tilde{\mathbf{t}} - \mathbf{M}_1 \mathbf{H} \tilde{\mathbf{t}} \right]}_{\mathbf{t}'} = \hat{\mathbf{t}} + \tilde{\mathbf{t}}, \quad \tilde{\mathbf{t}} \text{ beliebig.} \quad (2.27)$$

Der Vektor \mathbf{t}' wird durch Matrix $[\mathbf{1}_{k'} - \mathbf{H}^+ \mathbf{H}]$ in den Nullraum von \mathbf{H} projiziert, somit liegt der Vektor $\tilde{\mathbf{t}}$ im Nullraum von \mathbf{H} , d.h. $\tilde{\mathbf{t}} \in N(\mathbf{H})$. Der Vektor $\hat{\mathbf{t}}$ hingegen liegt im Spaltenraum von \mathbf{H}^T , d.h. $\hat{\mathbf{t}} \in R(\mathbf{H}^T)$. Der Vektor \mathbf{x} liegt im Spaltenraum von \mathbf{H} . Der Nullraum $N(\mathbf{H}^T)$ hat die Dimension 0, da wir die Zeilenregularität von \mathbf{H} vorausgesetzt haben. Eineindeutigkeit zwischen den Vektoren \mathbf{t} und \mathbf{x} besteht nur für den im Spaltenraum von \mathbf{H}^T liegenden Teil des Vektors \mathbf{t} . Zur Verdeutlichung des Sachverhalts dient Bild 2.1 ([Fisc99], [Stra80]).

Wir werden nun aus der zweiten Forderung an die Koordinatentransformation Bedingungen an die Matrix \mathbf{H}^{-R} herleiten. Näherungsweise gilt

$$\Delta \mathbf{t} = \mathbf{J}_{\mathbf{x}}^{\mathbf{t}} \Delta \mathbf{x} = \frac{1}{v_0} \mathbf{H}^{-R} \Delta \mathbf{x}. \quad (2.28)$$

Aus der zweiten Forderung an die Koordinatentransformation Gleichung (2.6) folgt unmittelbar, dass die letzte Spalte von \mathbf{H}^{-R} ausschließlich positive Elemente besitzen muss. Bei einer orthogonalen Matrix \mathbf{H} fallen die beiden Bedingungen an \mathbf{H} , wegen $\mathbf{H}^{-R} = \mathbf{H}^T$, zusammen.

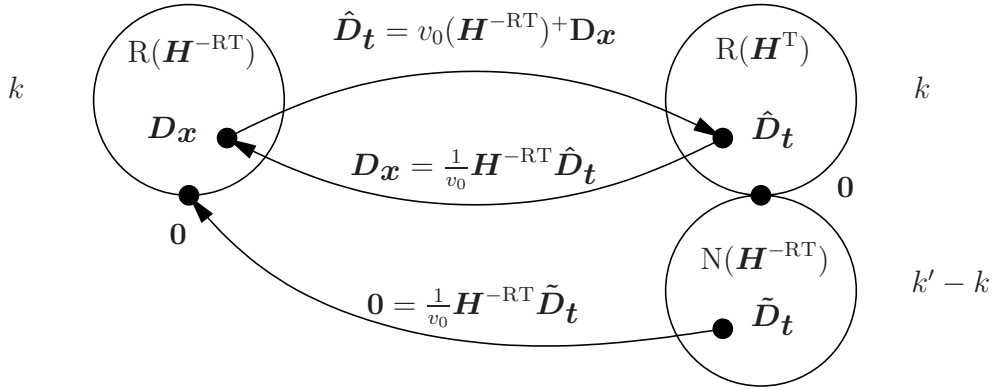


Bild 2.2: Transformation der Gradienten

Wir werden nun auf die Transformation der Gradienten im allgemeinen Fall eingehen. Dazu sei zunächst festgestellt, dass $\mathbf{D_t} = v_0 \mathbf{H^T D_x}$ ein überbestimmtes Gleichungssystem mit Spaltendefekt 0 ist, von dem wir fordern, dass es konsistent ist. Aus $\mathbf{t} = \mathbf{H^{-R} x} / v_0$ gewinnt man wiederum die Beziehung zwischen den Differentialoperatoren

$$\mathbf{D_x} = \frac{1}{v_0} [\mathbf{H^{-R}}]^T \mathbf{D_t}. \quad (2.29)$$

Ein konsistentes, überbestimmtes Gleichungssystem mit Spaltendefekt 0 hat zwar eine eindeutige Lösung, aber es existieren im Fall $k' > k$ verschiedene Möglichkeiten, diese Lösung darzustellen. Dies äußert sich wiederum dadurch, dass die transponierte Rechtsinverse $\mathbf{H^{-RT}} = [\mathbf{H^{-R}}]^T$ nicht eindeutig ist. Man verifiziert leicht, dass Gleichung (2.16) eine Lösung von Gleichung (2.29) ist und erhält durch Einsetzen von Gleichung (2.29) in Gleichung (2.16) die für eine speziell gewählte Rechtsinverse gültige Beziehung für die Differentialoperatoren im neuen Koordinatensystem

$$\mathbf{D_t} = [\mathbf{H^{-R} H}]^T \mathbf{D_t}. \quad (2.30)$$

Die allgemeine Lösung lautet

$$\mathbf{D_t} = v_0 (\mathbf{H^{-RT}})^+ \mathbf{D_x} + \underbrace{[\mathbf{1}_k - (\mathbf{H^{-RT}})^+ \mathbf{H^{-RT}}] \mathbf{D'_t}}_{=\tilde{\mathbf{D_t}}} = \hat{\mathbf{D_t}} + \tilde{\mathbf{D_t}}. \quad (2.31)$$

Die Matrix $\mathbf{H^T}$ ist selber Moore-Penrose-Inverse einer transponierten Rechtsinversen, d. h. $\mathbf{H^T} = (\mathbf{H^{-RT}})^+$. Somit sind die Spaltenräume gleich $\mathbf{R}(\mathbf{H^T}) = \mathbf{R}((\mathbf{H^{-RT}})^+)$.

Die Ableitungsoperatoren des neuen Koordinatensystems liegen in einem durch die Spaltenvektoren von $\mathbf{H^{-R}}$ aufgespannten k -dimensionalen Unterraum des $\mathbb{R}^{k'}$. Eineindeutigkeit besteht hier nur zwischen $\mathbf{D_x}$ und den Ableitungsoperatoren des neuen Koordinatensystems $\mathbf{D_t}$, die im Spaltenraum von $\mathbf{H^T}$ liegen, siehe Bild 2.2.

In der Praxis werden beim Übergang von der Differentialgleichung im ursprünglichen Koordinatensystem in das neue Koordinatensystem i. d. R. mehrere verschiedene Rechtsinverse benutzt, um so möglichst einfache Referenznetze zu erhalten.

Als Beispiel sei

$$\mathbf{H} = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad (2.32)$$

gegeben und dazu die drei Rechtsinversen

$$\mathbf{H}_1^{-R} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{H}_2^{-R} = \begin{bmatrix} 0 & 1/3 \\ -1 & 1/3 \\ 1 & 1/3 \end{bmatrix}, \quad \mathbf{H}_3^{-R} = \begin{bmatrix} -1 & 1/3 \\ -2 & 1/3 \\ 3 & 1/3 \end{bmatrix} \quad (2.33)$$

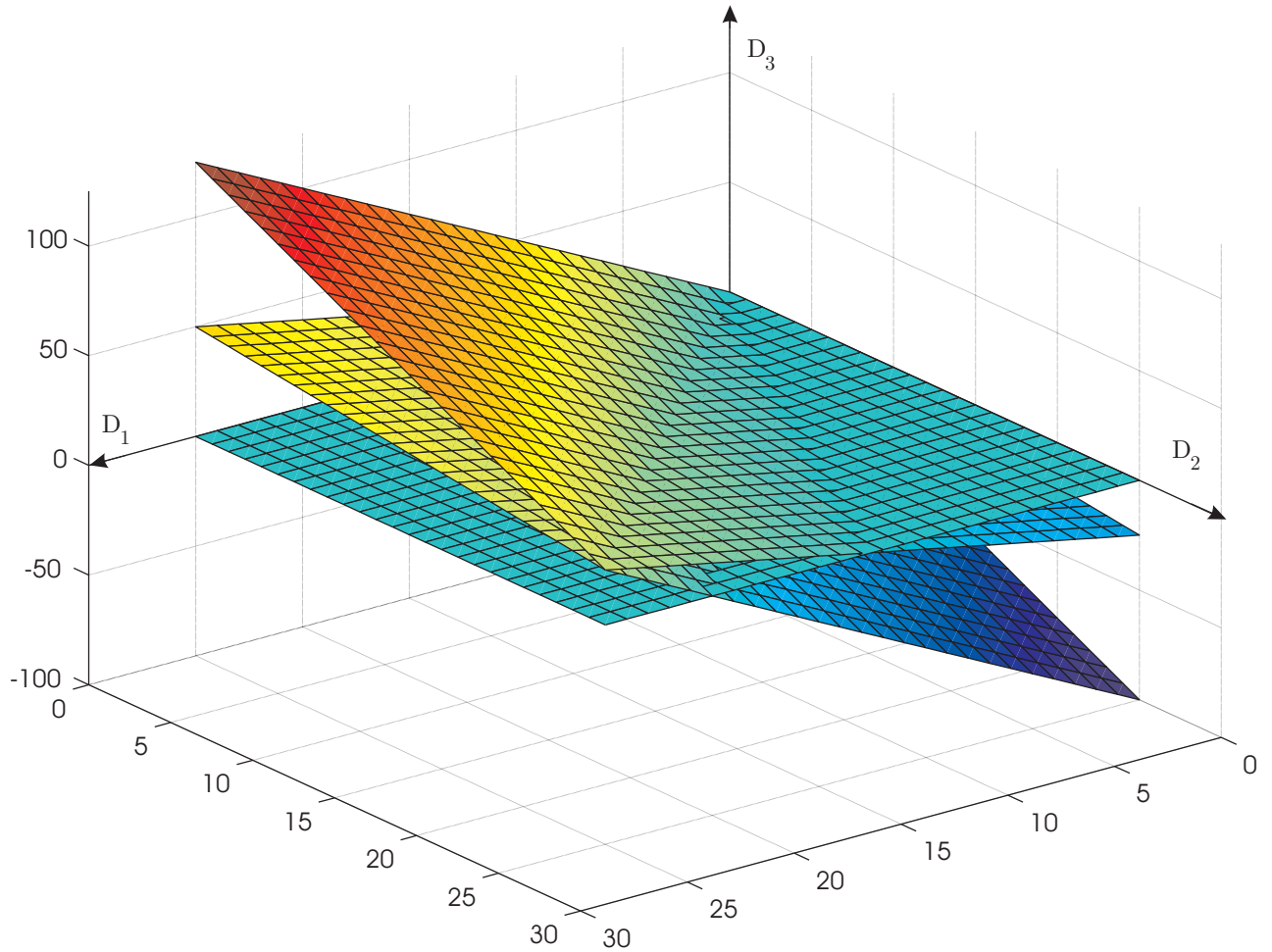


Bild 2.3: Durch die Spaltenvektoren der drei Rechtsinversen aufgespannten Ebenen (Spaltenräume von \mathbf{H}^{-R})

Bild 2.3 zeigt die Ebenen, auf denen die Lösungen \mathbf{D}_t von Gleichung (2.29) liegen.

2.3 Abtastung

Zur Erläuterung des Abtastvorgangs betrachten wir ein von der vektoriellen, kontinuierlichen Variablen \mathbf{x} abhängiges Signal $f(\mathbf{x})$. In der Regel erfolgt die Diskretisierung der Variablen \mathbf{x} durch

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{X}_A \boldsymbol{\mu} \text{ mit } \boldsymbol{\mu} = [\mu_1, \dots, \mu_{k-1}, \mu_k]^T \in \mathbb{Z}^k. \quad (2.34)$$

Weiterhin benötigen wir später

$$\boldsymbol{\mu}' = [\mu_1, \dots, \mu_{k-2}, \mu_{k-1}]^T. \quad (2.35)$$

Die quadratische Matrix \mathbf{X}_A bestimmt die Form des Abtastgitters. Wir bezeichnen sie als Abtastmatrix und setzen sie stets als regulär voraus. Sämtliche Ortpunkte zu einem konstanten Zeitpunkt μ_k wollen wir als eine Abtastschicht bezeichnen. Die durch die Abtastung entstandene Folge bezeichnen wir mit

$$\hat{f}(\boldsymbol{\mu}) = f(\mathbf{x}_0 + \mathbf{X}_A \boldsymbol{\mu}). \quad (2.36)$$

Der Abtastvorgang wird als linear bezeichnet, wenn $\mathbf{x}_0 = \mathbf{0}$ gilt, [Linn84]. Für die weitere Arbeit nehmen wir den Abtastvorgang als linear an. Die Fouriertransformierte der Funktion $f(\mathbf{x})$ lautet

$$F(j\boldsymbol{\omega}_x) = \mathcal{F}\{f(\mathbf{x})\} = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\mathbf{x}) e^{-j\boldsymbol{\omega}_x^T \mathbf{x}} d\mathbf{x}. \quad (2.37)$$

Die zugehörige Formel zur Rücktransformation ist

$$f(\mathbf{x}) = \mathcal{F}^{-1}\{F(j\boldsymbol{\omega}_x)\} = \frac{1}{(2\pi)^k} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} F(j\boldsymbol{\omega}_x) e^{j\boldsymbol{\omega}_x^T \mathbf{x}} d\boldsymbol{\omega}_x. \quad (2.38)$$

Die orts-/zeitdiskrete Fouriertransformierte der Folge $f(\mathbf{X}_A \boldsymbol{\mu})$ berechnet sich folgendermaßen

$$F_D(j\boldsymbol{\omega}_x) = \mathcal{F}_D\{f(\mathbf{X}_A \boldsymbol{\mu})\} = \sum_{\boldsymbol{\mu} \in \mathbb{Z}^k} f(\mathbf{X}_A \boldsymbol{\mu}) e^{-j\boldsymbol{\omega}_x^T \mathbf{X}_A \boldsymbol{\mu}}. \quad (2.39)$$

Mit

$$f(\mathbf{X}_A \boldsymbol{\mu}) = \frac{|\det(\mathbf{X}_A)|}{(2\pi)^k} \int_{\mathbf{A}} F_D(j\boldsymbol{\omega}_x) e^{j\boldsymbol{\omega}_x^T \mathbf{X}_A \boldsymbol{\mu}} d\boldsymbol{\omega}_x \quad (2.40)$$

erhält man die Folge aus dem Frequenzspektrum, wobei \mathbf{A} das Grundintervall ist. Das Frequenzspektrum der abgetasteten Funktion lässt sich durch $F(j\boldsymbol{\omega}_x)$ darstellen

$$\hat{f}(\boldsymbol{\mu}) \stackrel{d}{\circ} \frac{1}{|\det \mathbf{X}_A|} \sum_{\mathbf{n} \in \mathbb{Z}^k} F(j\boldsymbol{\omega}_x - j2\pi \mathbf{X}_A^{-T} \mathbf{n}). \quad (2.41)$$

Offenbar beeinflusst die Abtastmatrix \mathbf{X}_A auch das Frequenzspektrum des abgetasteten Signals. Die orts-/zeitkontinuierliche Funktion $f(\mathbf{x})$ lässt sich durch lineare Filterung zurückgewinnen, wenn sich die einzelnen Summanden nicht überlappen. Bei der Wahl der Abtastmatrix haben wir dies zu berücksichtigen. Die neue kontinuierliche unabhängige Variable \mathbf{t} wird zu

$$\mathbf{t} = \mathbf{t}_0 + \mathbf{T}_A \boldsymbol{\nu} \text{ mit } \boldsymbol{\nu} = [\nu_1, \dots, \nu_{k'-1}, \nu_{k'}]^T \in \mathbb{Z}^{k'} \quad (2.42)$$

diskretisiert. Weiterhin soll der Koordinatenursprung des Koordinatensystems $\boldsymbol{\nu}$ mit dem des Koordinatensystems $\boldsymbol{\mu}$ zusammenfallen, d. h. es gilt

$$\mathbf{x}_0 = v_0 \mathbf{H} \mathbf{t}_0. \quad (2.43)$$

Setzt man die Abtastpunkte in die Transformationsvorschrift ein, so ergibt sich

$$\mathbf{X}_A \boldsymbol{\mu} = v_0 \mathbf{H} \mathbf{T}_A \boldsymbol{\nu}. \quad (2.44)$$

Diese Beziehung offenbart die Tatsache, dass bei einer gegebenen Matrix \mathbf{H} die Abtastmatrizen \mathbf{X}_A und \mathbf{T}_A nicht unabhängig voneinander wählbar sind. Die Abtastgitter in beiden Koordinatensystemen sind somit auch miteinander verbunden.

Der einfachste Fall einer Abtastung ist die Rechteckabtastung [Linn84]. Die Abtastmatrix lautet

$$\mathbf{X}_A = \mathbf{diag}(\Delta X_1, \dots, \Delta X_k), \quad (2.45)$$

wobei ΔX_κ den Abstand zweier Abtastpunkte in Richtung x_κ bezeichnet. Die Dreiecksabtastung im Zusammenhang mit Wellendigitalfiltern ist in [LF90] beschrieben. Im Fall $k = 3$ ist die Abtastmatrix durch

$$\mathbf{X}_A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & \frac{-2}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{bmatrix} \quad (2.46)$$

gegeben. Mit der Dreiecksabtastung erreicht man die höchste spektrale Packungsdichte, wenn das Frequenzspektrum $F(j\omega_{\mathbf{x}})$ (hyper-)kugelförmig bandbegrenzt ist. Allerdings wird dieser Vorteil dadurch erkauft, dass die Implementierung eines derartigen Abtastmusters einen erhöhten Realisierungsaufwand (insbesondere bei Hardware) erfordert [Frie00], [Lisc91], [Bose01]. Führen wir die Abtastung im Koordinatensystem \mathbf{t} durch und wählen die Abtastmatrix \mathbf{X}_A proportional zu \mathbf{H} , d.h.

$$\mathbf{X}_A = v_0 \mathbf{H} \mathbf{T}_A, \quad \mathbf{T}_A = \mathbf{1}_{k'} T, \quad (2.47)$$

so können wir bei geeigneter Wahl von \mathbf{H} insbesondere Offset- und Rechteckraster erzeugen. Im Folgenden schränken wir uns auf $k = 4$ ein. Nachdem die üblicherweise verwendeten Abtastgitter kurz erläutert wurden, stellen wir nun sinnvolle Forderungen an die Abtastung.

- Jeder Verschiebe-Vektor, der von einem Abtastpunkt startet, muss wieder auf einem Abtastpunkt enden.
- Es darf nur eine Informationsübertragung von einer Abtastschicht zur unmittelbar darauf folgenden stattfinden. Diese Forderung findet ihre Begründung im Kapitel 3.
- Die Randbehandlung sollte möglichst einfach durchführbar sein.
- Eine systematische Synthese der Referenzschaltung sollte möglich sein.

Die erste Forderung lässt sich leicht erreichen, wenn die Abtastmatrix zu

$$\mathbf{X}_A = v_0 T \mathbf{H} \quad (2.48)$$

gewählt wird. Der zweiten Forderung kann dadurch genügt werden, dass die letzte Zeile von \mathbf{H} gleiche Elemente enthält. Die letzten beiden Forderungen lassen sich durch Rechteckabtastung in \mathbf{x} erreichen. Wir wählen die Transformationsmatrix zu

$$\mathbf{H} = \mathbf{diag} \left(\frac{\Delta X_1}{v_0 T}, \dots, \frac{\Delta X_k}{v_0 T} \right) \begin{bmatrix} 1 & 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & -1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}. \quad (2.49)$$

Zur prinzipiellen Erläuterung des Verfahrens wählt man üblicherweise $\Delta X_\kappa = v_0 T, \kappa = 1, 2, \dots, k$, um die Darstellung zu entlasten. Die Abtastmatrix liefert bei dieser Wahl nicht das Frequenzspektrum mit der höchsten Packungsdichte. Insofern liegt es nahe, dass wir nicht den effizientesten Algorithmus erhalten. Es soll aber nochmals darauf hingewiesen werden, dass dies nicht Ziel dieser Arbeit ist. Vielmehr steht in dieser Arbeit automatische fehlervermeidende Codegenerierung im Vordergrund.

Zum Schluss sei noch erwähnt, dass wir im weiteren Verlauf der Arbeit den Begriff Raster anstatt Abtastgitter verwenden werden, da wir im Grunde kein Signal abtasten, sondern zu bestimmten Rasterpunkten Werte durch numerische Integration berechnen.

2.4 Berechnungsgebiete

In der Regel beschreiben die ersten $k-1$ Koordinaten des Vektors \mathbf{x} Ortskoordinaten des physikalischen Raumes. Der Vektor der unabhängigen Ortsvariablen

$$_{-k}\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_{k-1} \end{bmatrix} \quad (2.50)$$

liegt im $k-1$ dimensionalen Raum über dem Körper der reellen Zahlen, d. h. $_{-k}\mathbf{x} \in \mathbb{R}^{k-1}$ und beschreibt dann einen Ortspunkt im so genannten Echtraum.² Dabei setzen wir voraus, dass die ersten $k-1$ unabhängigen Variablen des Koordinatensystems \mathbf{x} berandet sind.

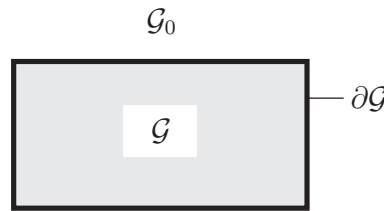


Bild 2.4: Das Berechnungsgebiet \mathcal{G} berandet durch $\partial\mathcal{G}$ mit Außengebiet \mathcal{G}_0

Die Variable t setzen wir in positiver Richtung als unbegrenzt voraus, da die Laufzeit des Betriebsmittels nicht bekannt ist. Die obere Grenze der Koordinate x_κ lautet $x_{\kappa \max}$. Die untere Grenze bezeichnen wir mit $x_{\kappa \min}$. Im Kapitel 2.10 wird deutlich, dass bei dem in dieser Arbeit eingesetzten Verfahren zur Randbehandlung, die untere Grenze nicht grundsätzlich zu $x_{\kappa \min} = 0$ gewählt werden kann. Das durch die zuvor eingeführten Grenzen beschriebene Berechnungsgebiet bezeichnen wir mit \mathcal{G} . Wir nehmen an, dass es sich bei \mathcal{G} um ein einfach zusammenhängendes Gebiet handelt. Ein Vektor \mathbf{x} liegt genau dann in dem Berechnungsgebiet \mathcal{G} falls

$$x_{\kappa \min} < x_\kappa < x_{\kappa \max} \quad \forall \kappa = 1, 2, \dots, k-1 \quad (2.51)$$

erfüllt ist. Durch die hier eingeführte untere und obere Grenze in Richtung x_κ wird in kartesischen Koordinaten für $k = 4$ eine rechteckige Box beschrieben. Im Falle allgemeiner Berechnungsgebiete wären die Grenzen selber noch Funktionen des Ortes. Der Einfachheit halber werden wir im Folgenden feste Grenzen annehmen.

Leider kann bei allgemeiner Abtastung nur mit expliziter Kenntnis der Abtastmatrix von den Grenzen in den kontinuierlichen Variablen auf die Anzahl Abtastpunkte in einer Richtung geschlossen werden. Die folgenden Ausführungen gelten daher nur für den Fall eines Rechteckgitters im Koordinatensystem \mathbf{x} . Die Anzahl der Abtastpunkte in Richtung x_κ bezeichnen wir mit P_{x_κ} . Wir definieren den Vektor

$$\mathbf{P}\mathbf{x} = [P_{x_1}, P_{x_2}, \dots, P_{x_{k-1}}]^T. \quad (2.52)$$

Im Fall eines Rechteckgitters gibt es einen kleinsten Abtastpunkt μ_κ , der zu null gewählt wird und einen größten, der zu $P_{x_\kappa} - 1$ festgelegt ist. Wohlgedenkt korrespondiert der kleinste (größte) Abtastpunkt einer Richtung nicht mit $x_{\kappa \min}$ ($x_{\kappa \max}$), sondern stellt lediglich den grenznächsten Abtastpunkt innerhalb \mathcal{G} dar. Die Gesamtanzahl der Abtastpunkte des Berechnungsgebietes \mathcal{G} lautet

$$P = \prod_{\kappa=1}^{k-1} P_{x_\kappa}. \quad (2.53)$$

²Wenn wir (unkorrektweise) sagen, x ist eine reelle Zahl, so ist dies so zu verstehen, dass der Zahlenwert der physikalischen Größe x reell ist. Soweit die konkreten Zusammenhänge nichts anderes ergeben, gilt diese Vereinbarung für die gesamte Arbeit.

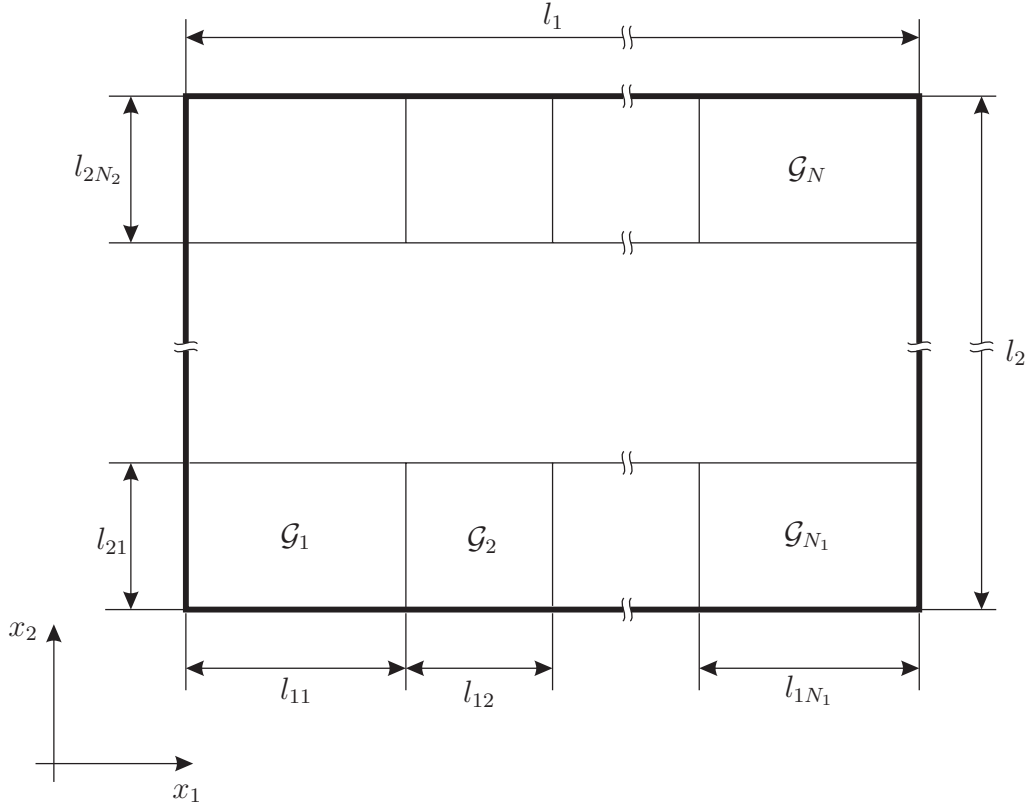


Bild 2.5: Abmessungen der Teilgebiete

Das Berechnungsgebiet \mathcal{G} ist eine echte Teilmenge des \mathbb{R}^{k-1} , d. h.

$$\mathcal{G} \subset \mathbb{R}^{k-1}. \quad (2.54)$$

Den nicht zum Berechnungsgebiet gehörenden Teil des \mathbb{R}^{k-1} bezeichnen wir mit \mathcal{G}_0 , d.h.

$$\mathcal{G}_0 = \mathbb{R}^{k-1} \setminus \mathcal{G} \iff \mathbb{R}^{k-1} = \mathcal{G}_0 \cup \mathcal{G}, \quad (2.55)$$

vgl. Bild 2.4. Teilgebiete homogener Materialverteilung bezeichnen wir mit \mathcal{G}_n . Das Berechnungsgebiet \mathcal{G} setzt sich aus N disjunkten Teilgebieten \mathcal{G}_n zusammen, d.h.

$$\mathcal{G} = \bigcup_{n=1}^N \mathcal{G}_n, \quad (2.56)$$

vgl. Bild 2.5.

Die Länge des Berechnungsgebietes in Richtung κ bezeichnen wir mit l_κ , $\kappa = 1, 2, \dots, k-1$. Die Länge der Teilgebiete \mathcal{G}_n in Richtung κ bezeichnen wir mit $l_{\kappa n}$, $n = 1, 2, \dots, N_\kappa$. Die Anzahl der Teilgebiete \mathcal{G}_n in Richtung κ lautet N_κ . Die Gesamtanzahl der Teilgebiete \mathcal{G}_n ergibt sich zu

$$\prod_{\kappa=1}^{k-1} N_\kappa = N. \quad (2.57)$$

Zur Verdeutlichung der Situation im Fall $k = 3$ dient Bild 2.5. Wir unterscheiden zwischen \mathcal{G} und \mathcal{G}_0 , da uns die Lösung der Differentialgleichung nur für das Gebiet \mathcal{G} interessiert. Die Parameter der Differentialgleichung innerhalb des Gebietes \mathcal{G} unterscheiden sich oft sehr stark von denen des Gebietes \mathcal{G}_0 . Der einfachste Fall liegt dann vor, wenn nur der Rand $\partial\mathcal{G}$ selber und nicht das Gebiet \mathcal{G}_0 Einfluss auf das Verhalten der Lösung innerhalb des Gebietes \mathcal{G} nimmt.

2.5 Anordnung der Gitterpunkte und Abmessungen des Berechnungsgebietes

Die Lage der Abtastpunkte resultiert i. W. aus der Lage der Randpunkte. Diese haben den Abstand einer halben Abtastlänge zum Rand. Die Begründung hierfür liefern wir in 2.10.

Wir betrachten nun den Fall gleich großer Teilberechnungsgebiete und rechteckförmiger Abtastung innerhalb der Berechnungsgebiete. Zwischen den Abtastpunkten $\mu_\kappa = 0$ und $\mu_\kappa = P_{x_\kappa} - 1$ beträgt der Abstand $(P_{x_\kappa} - 1)\Delta X_\kappa$. Wir berücksichtigen noch, dass aufgrund obiger Argumentation zwischen dem Rand $\partial\mathcal{G}$ und den ihm nächsten Gitterpunkten ein Abstand von einer halben Abtastlänge $\Delta X_\kappa/2$ liegt (siehe Bild 2.6). Die Längen l_κ betragen

$$l_\kappa = x_{\kappa \max} - x_{\kappa \min} = \Delta X_\kappa P_{x_\kappa} \iff \Delta X_\kappa = \frac{l_\kappa}{P_{x_\kappa}}. \quad (2.58)$$

Sind die Teilgebiete unterschiedlich groß, so ergibt sich die Länge des Teilgebietes n in Richtung μ_κ zu

$$l_{\kappa n} = \begin{cases} \Delta X_\kappa \left[P_{x_\kappa n} - \frac{1}{2} \right] & \text{für } n = 1, N_\kappa \\ \Delta X_\kappa [P_{x_\kappa n} - 1] & \text{für } n = 2, \dots, N_\kappa - 1, \end{cases} \quad (2.59)$$

wobei $P_{x_\kappa n}$ die Anzahl Abtastpunkte des Teilgebietes n in Richtung x_κ ist. Die Summe der Längen der Teilgebiete in Richtung κ muss die Gesamtlänge des Berechnungsgebietes in Richtung κ ergeben, d.h.

$$l_\kappa = \sum_{n=1}^{N_\kappa} l_{\kappa n}. \quad (2.60)$$

Setzen wir Gleichung (2.59) ein, so erhalten wir zunächst

$$l_\kappa = \Delta X_\kappa \left[-1 - (N_\kappa - 2) + \sum_{n=1}^{N_\kappa} P_{x_\kappa n} \right] = \Delta X_\kappa \left[\sum_{n=1}^{N_\kappa} P_{x_\kappa n} - (N_\kappa - 1) \right] \quad (2.61)$$

und durch Vergleich mit Gleichung (2.58)

$$P_{x_\kappa} = \sum_{n=1}^{N_\kappa} P_{x_\kappa n} - (N_\kappa - 1), \quad (2.62)$$

was durch die Tatsache bestätigt wird, dass genau die $N_\kappa - 1$ inneren Abtastpunkte in Richtung κ zu zwei Gebieten gemeinsam gehören.

2.6 Mehrdimensionale Kirchhoff'sche Netze

In diesem Abschnitt werden die für diese Arbeit benötigten Definitionen und Eigenschaften mehrdimensionaler Kirchhoff'scher Netze angegeben. Hierzu führen wir zunächst den Torbegriff ein. Wir verstehen unter einem Tor ein Klemmenpaar $1 - 1'$ gemäß Bild 2.7, welches durch die Spannung $u(\mathbf{t})$ und den Strom $i(\mathbf{t})$ (diese sind proportional zu den Feldgrößen in dem im Kapitel 2.1 definierten Sinne) charakterisiert ist. Zu beachten ist, dass hineinfließender und herausfließender Strom gleich sind. Neben der Beschreibung mit Spannung und Strom verwenden wir die Beschreibung mit Wellengrößen. Die Verwendung der Wellengrößen stellt eine bijektive Koordinatentransformation der unabhängigen Größen dar.

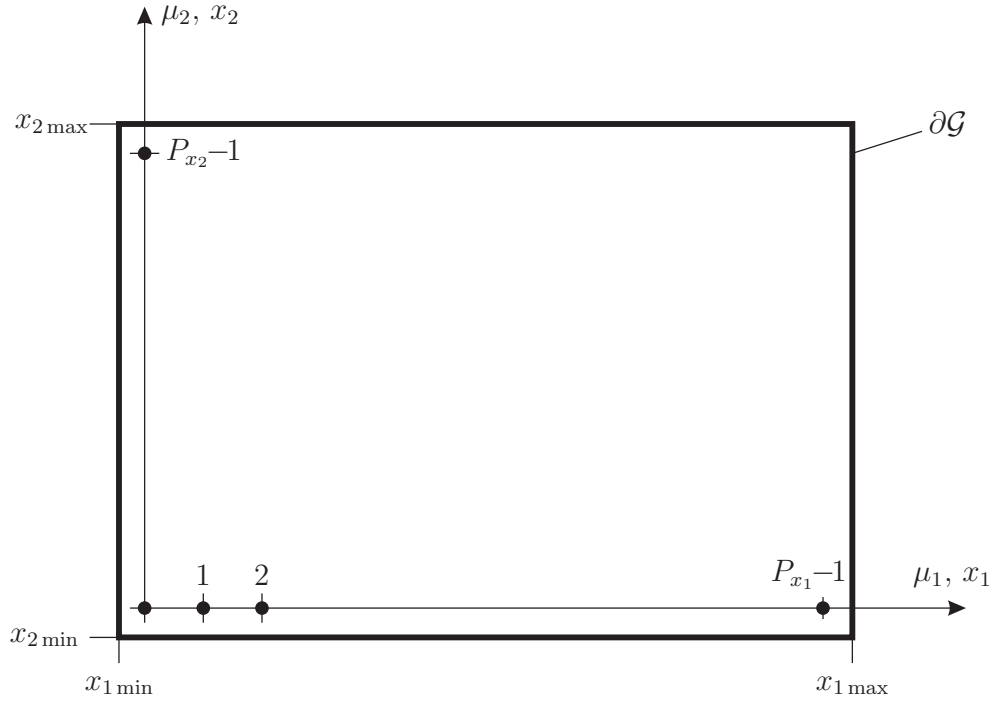


Bild 2.6: Zur Lage der Abtastpunkte

Diese Koordinatentransformation ist durch den positiv angenommenen so genannten Torwiderstand R bestimmt. Die gebräuchlichsten Wellengrößen sind die Spannungswellen

$$a' = u + Ri, \quad b' = u - Ri \quad (2.63)$$

und die Leistungswellen

$$a = \frac{u + Ri}{2\sqrt{R}}, \quad b = \frac{u - Ri}{2\sqrt{R}}. \quad (2.64)$$

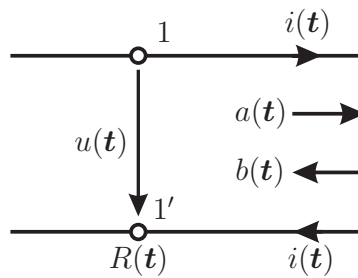


Bild 2.7: Zum Torbegriff

Nun betrachten wir ein n -Tor. Zweckmäßigerweise fassen wir zur Beschreibung des n -Tors alle Spannungen, Ströme und Wellengrößen in Vektoren zusammen. Weiterhin definieren wir eine reelle, p.d., symmetrische Torwiderstandsmatrix \mathbf{R} . Sofern nichts anderes gesagt ist, gehen wir von einer Diagonalmatrix aus, d.h. $\mathbf{R} = \mathbf{G}^{-1} = \mathbf{diag}(R_1, \dots, R_n)$. Die inverse Matrix \mathbf{G} bezeichnen wir als Torleitwert-

matrix. Wir erhalten die kompakte Darstellung

$$\begin{bmatrix} \mathbf{b}' \\ \mathbf{a}' \end{bmatrix} = \begin{bmatrix} \mathbf{1} & -\mathbf{R} \\ \mathbf{1} & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{i} \end{bmatrix} \iff \begin{bmatrix} \mathbf{u} \\ \mathbf{i} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{1} & \mathbf{1} \\ -\mathbf{G} & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{b}' \\ \mathbf{a}' \end{bmatrix} \quad (2.65)$$

$$\begin{bmatrix} \mathbf{b} \\ \mathbf{a} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{G}^{1/2} & -\mathbf{R}^{1/2} \\ \mathbf{G}^{1/2} & \mathbf{R}^{1/2} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{i} \end{bmatrix} \iff \begin{bmatrix} \mathbf{u} \\ \mathbf{i} \end{bmatrix} = \begin{bmatrix} \mathbf{R}^{1/2} & \mathbf{R}^{1/2} \\ -\mathbf{G}^{1/2} & \mathbf{G}^{1/2} \end{bmatrix} \begin{bmatrix} \mathbf{b} \\ \mathbf{a} \end{bmatrix}. \quad (2.66)$$

Falls die einfallenden Wellen des n -Tores unabhängig voneinander gewählt werden können, gelten die Beziehungen

$$\mathbf{b} = \mathbf{S}\mathbf{a} \quad \text{und} \quad \mathbf{b}' = \mathbf{S}'\mathbf{a}', \quad (2.67)$$

wobei \mathbf{S} die Leistungswellenstreumatrix und \mathbf{S}' die Spannungswellenstreumatrix darstellt. Die beiden Arten von Streumatrizen sowie die Wellengrößen lassen sich ineinander umrechnen

$$\mathbf{S}' = \mathbf{R}^{1/2} \mathbf{S} \mathbf{G}^{1/2}, \quad \mathbf{b}' = 2 \mathbf{R}^{1/2} \mathbf{b}, \quad \mathbf{a}' = 2 \mathbf{R}^{1/2} \mathbf{a}. \quad (2.68)$$

Leistungswellen garantieren die volle Robustheit des Wellendigitalalgorithmus im nicht konstanten und nicht linearen Fall. Hingegen ist bei linearen, konstanten Bauelementen die volle Robustheit auch schon durch Verwendung von Spannungswellen sichergestellt, [Fett98].

Die vom n -Tor aufgenommene Leistungsdichte ist die Summe der über die n Tore übertragene Leistungsdichte

$$p(\mathbf{t}) = \mathbf{u}^T(\mathbf{t}) \mathbf{i}(\mathbf{t}) = \|\mathbf{a}\|^2 - \|\mathbf{b}\|^2 = \mathbf{a}'^T \mathbf{G} \mathbf{a}' - \mathbf{b}'^T \mathbf{G} \mathbf{b}' = \sum_{\mu=1}^n [a_\mu^2 - b_\mu^2] = \sum_{\mu=1}^n \sum_{\nu=1}^n [a'_\mu G_{\mu\nu} a'_\nu - b'_\mu G_{\mu\nu} b'_\nu]. \quad (2.69)$$

Korrekterweise können wir im mehrdimensionalen Fall nicht mehr von Spannungen und Strömen sprechen, da es sich hier auch um Dichten handelt. Wir wollen aber trotzdem an den eingeführten Begriffen Spannung und Strom festhalten. In [Fett99] und [Heme95] wird die Leistungsdichte auch als mehrdimensionale Leistung bezeichnet.

In diesem Zusammenhang sprechen wir von externen Eigenschaften, wenn wir das Verhalten des Systems nur von außen betrachten. Die inneren Zustände werden dabei nicht berücksichtigt, insbesondere kann der Fall auftreten, dass das externe System eine Eigenschaft besitzt, die ein Teilsystem nicht besitzt. Ein System besteht in der Regel aus mehreren (Elementar-)elementen. Als Element bezeichnen wir ein System, welches nicht weiter zerlegt werden kann. Falls ein System keine Quellen besitzt und alle Elemente passiv (verlustfrei, energieneutral) sind, dann heißt das System intern passiv (verlustfrei, energieneutral). Zur geeigneten Definition der Passivität greifen wir auf [Nits93], [Fett99] und [Heme95] zurück. Dabei werden wir zunächst Energiebetrachtungen im Koordinatensystem \mathbf{x} durchführen. Wir wollen uns dabei auf den Fall $k = 4$ einschränken und betrachten ein 4-dimensionales Gebiet $\mathcal{G}_{\mathbf{x}}$ mit dem Rand $\partial\mathcal{G}_{\mathbf{x}}$, einem Flächenelement $dA_{\mathbf{x}}$ und dem nach außen gerichteten Normalenvektor $\mathbf{n}_{\mathbf{x}}$. Im Folgenden wollen wir die auf ein System wirkenden energetischen Einflüsse in einer Energiebilanz quantitativ beschreiben. Dazu führen wir die im System gespeicherte Energiedichte

$$\mathbf{W}_s = [W_x, W_y, W_z, W_t]^T \quad \text{mit} \quad W_t \geq 0. \quad (2.70)$$

ein. Die physikalische Einheit ist für alle Koordinaten Joule pro Kubikmeter, d. h. $[\mathbf{W}_s] = [\text{J}/\text{m}^3]$. Auf die Energiebilanz haben die folgenden Größen einen Einfluss. Zunächst lassen sich die Einflüsse unterteilen nach Ursachen innerhalb des Gebietes $\mathcal{G}_{\mathbf{x}}$ und nach Ursachen, die ihren Ursprung außerhalb des Gebiets

haben. Von Außen geht die durch die Gebietsoberfläche (diese räumliche und zeitliche Begrenzung wird auch als Segmentgrenze bezeichnet) transportierte Energiedichte $\mathbf{W}_s^T \mathbf{n}_x$ in integraler Form als die Energie

$$\int_{\partial \mathcal{G}_x} \mathbf{W}_s^T \mathbf{n}_x dA_x \quad (2.71)$$

in die Energiebilanz ein. Falls das Integral positiv ist, sprechen wir von einer nach Außen abgeführten Energie. Innerhalb des Gebietes kann ein Energieaustausch mit anderen Systemen über die so genannte Prozessgrenze (das sind die Tore des Systems) erfolgen. Die über die Tore zugeführte Leistungsdichte ist p , welche die Einheit Watt pro Kubikmeter hat, d. h. $[p] = [\text{W}/\text{m}^3]$.

Die über die Tore zugeführte Energie lautet mit $dV_x = dx_1 dx_2 dx_3 dx_4 = dx dy dz v_4 dt$

$$\frac{1}{v_4} \int_{\mathcal{G}_x} p dV_x = \int_{\mathcal{G}_x} p dx dy dz dt . \quad (2.72)$$

Die durch die Quellen über die Quellengrenze zugeführte Leistungsdichte ist p^Q . Die dem Gebiet \mathcal{G}_x durch Quellen zugeführte Energie ist

$$\frac{1}{v_4} \int_{\mathcal{G}_x} p^Q dV_x . \quad (2.73)$$

Letztendlich entweicht über die Dissipationsgrenze die nicht negative Leistungsdichte $p^V \geq 0$, sodass die im Gebiet \mathcal{G}_x dissipierte Energie

$$\frac{1}{v_4} \int_{\mathcal{G}_x} p^V dV_x \quad (2.74)$$

ist. Alle energetischen Einflüsse auf das System sind im Bild 2.8 dargestellt. Für die Energiebilanz

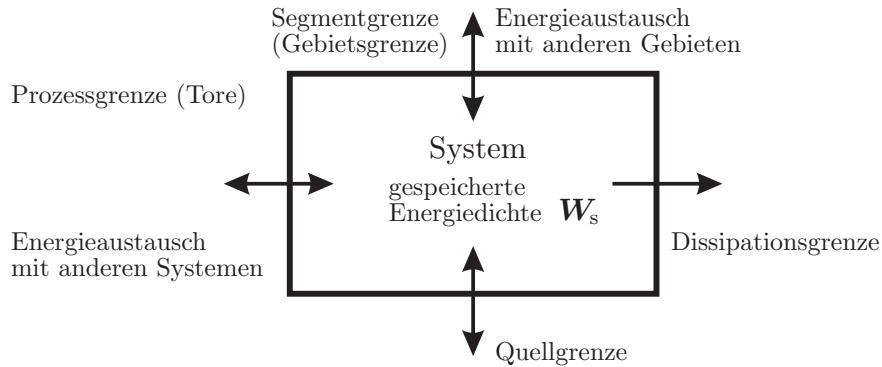


Bild 2.8: Energiebilanz eines Systems (z.B. n -Tor)

vereinbaren wir die folgende Vorzeichenkonvention

$$\int_{\partial \mathcal{G}_x} \mathbf{W}_s^T \mathbf{n} dA_x + \int_{\mathcal{G}_x} \frac{p^V}{v_4} dV_x = \int_{\mathcal{G}_x} \frac{p^Q}{v_4} dV_x + \int_{\mathcal{G}_x} \frac{p}{v_4} dV_x . \quad (2.75)$$

Mit dem Gauß'schen Satz gilt unter sehr allgemeinen Bedingungen (Stetigkeit der Integranden und endliche Grenzen sind bereits hinreichend)

$$\int_{\mathcal{G}_x} \mathbf{D}_x^T \mathbf{W}_s dV_x + \int_{\mathcal{G}_x} \frac{p^V}{v_4} dV_x = \int_{\mathcal{G}_x} \frac{p^Q}{v_4} dV_x + \int_{\mathcal{G}_x} \frac{p}{v_4} dV_x . \quad (2.76)$$

Da wir keine Einschränkungen bzgl. des Gebietes $\mathcal{G}_{\mathbf{x}}$ gemacht haben, muss die Energiebilanz für beliebige $\mathcal{G}_{\mathbf{x}}$ gültig sein. Hieraus folgt die Gleichheit der Integranden, die als Energieerhaltung in differentieller Form oder Leistungsdichtebilanz bezeichnet wird

$$v_4 \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s + p^V = p^Q + p . \quad (2.77)$$

Passive Systeme definieren wir nun durch Quellenfreiheit $p^Q = 0$. Dann gilt mit $p^V \geq 0$ die Ungleichung

$$\int_{\mathcal{G}_{\mathbf{x}}} p \, dV_{\mathbf{x}} - v_4 \int_{\mathcal{G}_{\mathbf{x}}} \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s dV_{\mathbf{x}} = \int_{\mathcal{G}_{\mathbf{x}}} p^V \, dV_{\mathbf{x}} \geq 0 , \quad (2.78)$$

bzw. in differentieller Form

$$p - v_4 \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s = p^V \geq 0 \implies p \geq v_4 \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s . \quad (2.79)$$

Verlustfreie Systeme sind passive Systeme, die keine Leistung dissipieren, d. h.

$$p^V = 0 \iff p = v_4 \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s . \quad (2.80)$$

Energieneutrale Systeme sind dadurch gekennzeichnet, dass sie verlustfrei sind und keine Energie speichern, d. h. $\mathbf{W}_s = \mathbf{0}$. Sie nehmen daher auch keine Energie auf $p = 0$.

Um eine andere Interpretation der Passivitätsbedingung zu erhalten, teilen wir die Energiedichte in einen räumlichen und einen zeitlichen Anteil auf

$$\mathbf{W}_s = [\mathbf{W}_{\text{sp}}^T, W_t]^T . \quad (2.81)$$

Zudem definieren wir die in einem örtlichen Gebiet \mathcal{G} mit Rand $\partial\mathcal{G}$, einem Flächenelement dA und dem nach außen gerichteten Normalenvektor \mathbf{n} gespeicherte Energie zu

$$E = \int_{\mathcal{G}} W_t \, dV \quad , \quad dV = dx \, dy \, dz . \quad (2.82)$$

Die Energie hat die Einheit Joule $[E] = [\text{J}]$. Nehmen wir stetige Funktionen und endliche Integralgrenzen an, so ist die Vertauschbarkeit der Integrationsreihenfolge gegeben. Dann können wir die transportierte Energie durch

$$\begin{aligned} \int_{\mathcal{G}_{\mathbf{x}}} \mathbf{D}_{\mathbf{x}}^T \mathbf{W}_s \, dV_{\mathbf{x}} &= \int_{\mathcal{G}_{\mathbf{x}}} \text{div } \mathbf{W}_{\text{sp}} \, dV_{\mathbf{x}} + \frac{1}{v_4} \int_{\mathcal{G}_{\mathbf{x}}} D_t W_t \, dV_{\mathbf{x}} \\ &= v_4 \int_{t_0}^{t_1} \int_{\mathcal{G}} \text{div } \mathbf{W}_{\text{sp}} \, dV \, dt + \int_{\mathcal{G}} \int_{t_0}^{t_1} D_t W_t \, dt \, dV \\ &= v_4 \int_{t_0}^{t_1} \int_{\mathcal{G}} \text{div } \mathbf{W}_{\text{sp}} \, dV \, dt + \int_{\mathcal{G}} [W_t(t_1) - W_t(t_0)] \, dV \\ &= v_4 \int_{t_0}^{t_1} \int_{\partial\mathcal{G}} \mathbf{W}_{\text{sp}}^T \mathbf{n} \, dA \, dt + E(t_1) - E(t_0) \end{aligned} \quad (2.83)$$

ausdrücken. Hierin ist

$$v_4 \int_{t_0}^{t_1} \int_{\partial\mathcal{G}} \mathbf{W}_{\text{sp}}^T \mathbf{n} \, dA \, dt \quad (2.84)$$

die im Zeitintervall $t_0 < t < t_1$ über die örtliche Segmentgrenze zugeführte Energie. Die im 4-dimensionalen Raum transportierte Energie entspricht somit der örtlich transportierten Energie und der Differenz der gespeicherten Energie.

Wir wollen nun eine notwendige Bedingung für passive Systeme herleiten. Dazu nehmen wir an, dass zwischen den Zeitpunkten t_0 und t_1 dem System über die örtliche Segmentgrenze keine Energie zugeführt wird. Dies ist dann der Fall, wenn die Normalkomponente der Energiedichte \mathbf{W}_{sp} auf $\partial\mathcal{G}$ verschwindet, d. h. $\mathbf{W}_{\text{sp}}^T \mathbf{n} = 0$ auf $\partial\mathcal{G}$. Nehmen wir noch $E(t_0) = 0$ an, berücksichtigen $W_t \geq 0$ und Gleichung (2.82), so folgt aus der Passivitätsbedingung Gleichung (2.78) dass die zugeführte Energie nicht negativ ist.

$$0 \leq \frac{1}{v_4} \int_{\mathcal{G}_{\mathbf{x}}} p \, dV_{\mathbf{x}} . \quad (2.85)$$

Um nun eine geeignete Definition für mehrdimensional passive Systeme in den Koordinaten \mathbf{t} anzugeben, betrachten wir ein k' -dimensionales Gebiet $\mathcal{G}_{\mathbf{t}}$ mit dem Rand $\partial\mathcal{G}_{\mathbf{t}}$, einem Flächenelement $dA_{\mathbf{t}}$ und dem nach außen gerichteten Normalenvektor $\mathbf{n}_{\mathbf{t}}$. Die vom n -Tor aufgenommene Energie ergibt sich durch Integration der zugeführten Leistungsdichte über $\mathcal{G}_{\mathbf{t}}$, d. h.

$$W = \int_{\mathcal{G}_{\mathbf{t}}} p' \, dt_1 \dots dt_{k'} . \quad (2.86)$$

Diese Energie soll nach Definition die Einheit einer physikalischen Energie besitzen. Zusätzlich zur aufgenommenen Energie definieren wir eine (vektorielle) gespeicherte Energiedichte

$$\hat{\mathbf{W}}_s(\mathbf{t}) = [W'_1, W'_2, \dots, W'_{k'}]^T . \quad (2.87)$$

Mit diesen Vorbetrachtungen können wir nun die MD-Passivität eines Systems näher beschreiben. Ein System heißt MD-extern passiv, falls eine (vektorielle) Funktion $\hat{\mathbf{W}}_s$ existiert, für die gilt

1. $\hat{\mathbf{W}}_s(\mathbf{t}) \geq \mathbf{0}$ (koordinatenweise ≥ 0) und
2. $W \geq \int_{\partial\mathcal{G}_{\mathbf{t}}} \hat{\mathbf{W}}_s^T \mathbf{n}_{\mathbf{t}} \, dA_{\mathbf{t}}$ (für beliebige Gebiete) .

(2.88)

Ein System heißt MD-extern verlustfrei, falls das System passiv ist und für beliebige Gebiete $\mathcal{G}_{\mathbf{t}}$ gilt

$$W = \int_{\partial\mathcal{G}_{\mathbf{t}}} \hat{\mathbf{W}}_s^T \mathbf{n}_{\mathbf{t}} \, dA_{\mathbf{t}} \quad (2.89)$$

Ein System heißt MD-extern energieneutral, falls das System verlustfrei ist und für beliebige Zeitpunkte \mathbf{t} gilt

$$\hat{\mathbf{W}}_s(\mathbf{t}) = \mathbf{0} . \quad (2.90)$$

Ein System heißt MD-extern dissipativ, falls das System passiv ist und ein Vorgang existiert, sodass

$$W > \int_{\partial\mathcal{G}} \hat{\mathbf{W}}_s^T \mathbf{n}_{\mathbf{t}} \, dA_{\mathbf{t}} \quad (2.91)$$

gilt.

Wir betrachten nun den Sonderfall eines linearen konstanten Systems. Im stationären Zustand folgen die Spannungen und Ströme den Beziehungen

$$\mathbf{u}(\mathbf{t}) = \text{Re}\{\mathbf{U} e^{\mathbf{p}_t^T \mathbf{t}}\} , \quad \mathbf{i}(\mathbf{t}) = \text{Re}\{\mathbf{I} e^{\mathbf{p}_t^T \mathbf{t}}\} . \quad (2.92)$$

Wir definieren die auf die gesamte \mathbf{p}_t -Polyebene verallgemeinerte Wirkleistung zu $P = \frac{1}{2} \text{Re}\{\mathbf{I}^H \mathbf{U}\}$. Wohlgedenkt ist P aber nur für $\text{Re}\mathbf{p}_t = \mathbf{0}$ die im Mittel über die n Tore übertragene Leistung. Es ergeben sich die folgenden notwendigen und hinreichenden Bedingungen für P

- Passives System : $P \geq 0$ für $\text{Re } \mathbf{p}_t > 0$
- Verlustfreies System : passiv und $P = 0$ für $\text{Re } \mathbf{p}_t = 0$
- Energieneutrales System : verlustfrei und $P = 0 \forall \mathbf{p}_t$.

2.7 Eigenschaften mehrdimensionaler Wellendigitalfilter

Ein System heißt diskret MD-extern passiv, falls eine (vektorielle) Funktion der gespeicherten Energie

$$\hat{\mathbf{W}}_s = [\hat{W}_{s1}, \dots, \hat{W}_{sk'}]^T \quad (2.93)$$

existiert, die koordinatenweise nicht negativ ist. Zudem darf die am Rasterpunkt $\mathbf{T}\nu$ zugeführte Energie $Tp(\mathbf{T}\nu)$ nicht kleiner als die Summe der Energiedifferenzen des aktuellen und des unmittelbar zurückliegenden Rasterpunktes über alle Raumkoordinaten sein, d. h. es muss gelten

$$\begin{aligned} 1. \quad & \hat{\mathbf{W}}_s(\mathbf{T}\nu) \geq \mathbf{0} \quad (\text{koordinatenweise } \geq 0) \text{ und} \\ 2. \quad & Tp(\mathbf{T}\nu) \geq \sum_{\kappa=1}^{k'} \hat{W}_{s\kappa}(\mathbf{T}\nu) - \sum_{\kappa=1}^{k'} \hat{W}_{s\kappa}(\mathbf{T}\nu - \Delta t) . \end{aligned} \quad (2.94)$$

Ein System heißt diskret MD-extern verlustfrei, falls das System passiv ist und die zum Zeitpunkt $\mathbf{T}\nu$ zugeführte Energie $Tp(\mathbf{T}\nu)$ gleich der Summe der Energiedifferenzen des aktuellen und des unmittelbar zurückliegenden Rasterpunktes über alle Raumkoordinaten ist

$$Tp(\mathbf{T}\nu) = \sum_{\kappa=1}^{k'} \hat{W}_{s\kappa}(\mathbf{T}\nu) - \sum_{\kappa=1}^{k'} \hat{W}_{s\kappa}(\mathbf{T}\nu - \Delta t) . \quad (2.95)$$

Ein System heißt diskret MD-extern energieneutral, falls das System verlustfrei ist und darüberhinaus die gespeicherte Energie an allen Abtastpunkten null ist, d. h.

$$\hat{\mathbf{W}}_s(\mathbf{T}\nu) = \mathbf{0} . \quad (2.96)$$

Die in 2.6 gemachten Aussagen bzgl. interner Eigenschaften gelten auch hier.

2.8 Wellendigitalfilter Bauelemente

2.8.1 Energieneutrale Bauelemente

Energieneutrale n -Tore sind gemäß Gleichung (2.90) dadurch gekennzeichnet, dass die aufgenommene MD-Leistung zu jedem Zeitpunkt null ist. Die zum n -Tor zugehörige Streumatrix genügt der Eigenschaft

$$\mathbf{S}^H \mathbf{S} = \mathbf{1}_n . \quad (2.97)$$

Man beachte, dass die Streumatrix eine Funktion der Zeit und der Feldgrößen sein darf.

Die direkte Implementierung der Streumatrix auf einem Digitalrechner ist nicht sinnvoll, da bei endlicher Wortlänge die Sicherstellung der Passivität einen erheblichen Aufwand erfordert. Zur Lösung des Problems greifen wir auf die bekannte Tatsache zurück, dass die Leistungswellen-Streumatrix eines energieneutralen n -Tors mittels QR-Zerlegung nach Givens durch Streumatrizen von $(n^2 - n)/2 = \binom{n}{2}$

Zweitoren und n skalaren unimodularen Multiplizierern ρ_ν dargestellt werden kann [MF96], [Frän97], d. h.

$$\mathbf{S} = \left[\prod_{\mu=1}^{\binom{n}{2}} \mathbf{H}_\mu \right] \mathbf{R} \quad , \quad \binom{n}{2} = \frac{n(n-1)}{2} \quad , \quad \text{wobei} \quad (2.98)$$

$$\begin{aligned} \mathbf{H}_\mu &= \mathbf{1}_n + \mathbf{e}_\alpha \mathbf{e}_\beta^T s_\mu + \mathbf{e}_\beta \mathbf{e}_\alpha^T s_\mu^* + [\mathbf{e}_\alpha \mathbf{e}_\alpha^T][c_\mu - 1] - [\mathbf{e}_\beta \mathbf{e}_\beta^T][c_\mu + 1] \\ &= \begin{bmatrix} \mathbf{1}_{\alpha-1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & c_\mu & \mathbf{0} & s_\mu & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1}_{\beta-\alpha-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & s_\mu^* & \mathbf{0} & -c_\mu & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1}_{n-\beta} \end{bmatrix} \end{aligned} \quad (2.99)$$

die so genannten Givens-Matrizen sind mit geeignet gewählten α, β für ein bestimmtes μ und reelles c_μ . Die Zuordnung von α, β zu einem speziellem μ kann auf verschiedene Weisen erfolgen. Im reellen Fall können die Matrix-Elemente als $c_\mu = \cos(\varphi_\mu)$ und $s_\mu = \sin(\varphi_\mu)$ interpretiert werden. Die Givens-Matrizen können in reduzierter Form als

$$\begin{aligned} \hat{\mathbf{H}}_\mu &= \mathbf{e}_1 \mathbf{e}_2^T s_\mu + \mathbf{e}_2 \mathbf{e}_1^T s_\mu^* + [\mathbf{e}_1 \mathbf{e}_1^T - \mathbf{e}_2 \mathbf{e}_2^T] c_\mu \\ &= \begin{bmatrix} c_\mu & s_\mu \\ s_\mu^* & -c_\mu \end{bmatrix} = \mathbf{1}_2 - \gamma_\mu \gamma_\mu^H \end{aligned} \quad (2.100)$$

geschrieben werden mit $\gamma_\mu^H \gamma_\mu = 2$. Wegen der Unitarität der Matrix \mathbf{S} ist die untere oder obere Dreiecksmatrix \mathbf{R} eine Diagonalmatrix,

$$\mathbf{R} = \text{diag}(\rho_1, \dots, \rho_n) \quad , \quad (2.101)$$

mit unimodularen Skalaren ρ_1, \dots, ρ_n , [Frän97]. Wählen wir die gleiche Faktorisierung wie in [Frän97] auf Seite 55, erhalten wir die Wellen-Digital-Struktur für $n = 4$ im Bild 2.9. Ein mögliche Referenzschal-

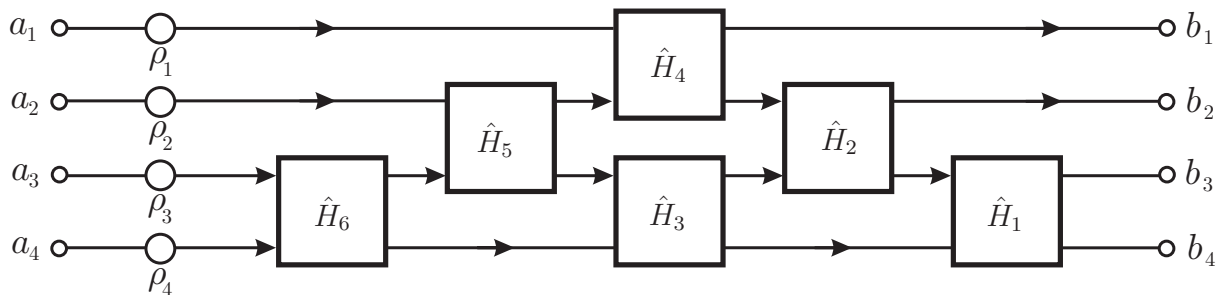


Bild 2.9: Realisierung eines energieneutralen M -Tors mit QR-Zerlegung nach Givens für $M = n = 4$

tung findet sich in [Voll04c].

Es folgen einige spezielle energieneutrale Bauelemente.

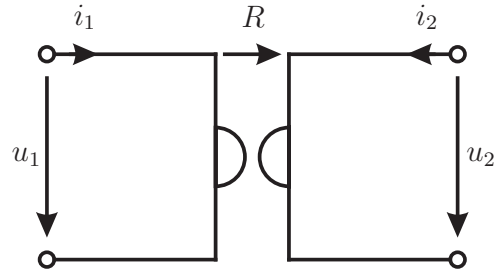


Bild 2.10: Idealer Gyrator

Idealer Gyrator

Der ideale Gyrator ist ein nichtreziprokes Bauelement, das durch die Gleichungen

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 0 & -R \\ R & 0 \end{bmatrix} \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \quad (2.102)$$

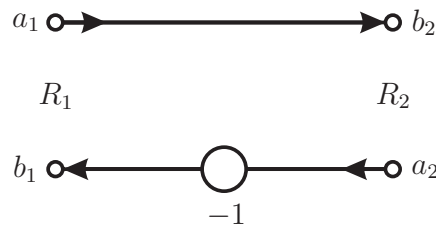
bestimmt ist. Die Streumatrizen lauten

$$\mathbf{S} = \begin{bmatrix} \frac{R^2 G_1 G_2 - 1}{1 + R^2 G_1 G_2} & \frac{-2R\sqrt{G_1 G_2}}{1 + R^2 G_1 G_2} \\ \frac{2R\sqrt{G_1 G_2}}{1 + R^2 G_1 G_2} & \frac{R^2 G_1 G_2 - 1}{1 + R^2 G_1 G_2} \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} \frac{R^2 G_1 G_2 - 1}{1 + R^2 G_1 G_2} & \frac{-2RG_2}{1 + R^2 G_1 G_2} \\ \frac{2RG_1}{1 + R^2 G_1 G_2} & \frac{R^2 G_1 G_2 - 1}{1 + R^2 G_1 G_2} \end{bmatrix}. \quad (2.103)$$

Wählt man die Torwiderstände beider Tore gleich der Gyrationkonstanten R , so berechnen sich die Streumatrizen zu

$$\mathbf{S} = \mathbf{S}' = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}. \quad (2.104)$$

Bei dieser Wahl der Torwiderstände hängt die Streumatrix von keinem Parameter ab. Das Bild 2.11 zeigt das Wellenflussdiagramm.

Bild 2.11: Wellenflussdiagramm eines idealen Gyrators für $R = R_1 = R_2$

n -Tor-Zirkulator

Ein weiteres nichtreziprokes Bauelement ist der Zirkulator, siehe Bild 2.12 a). Den Zirkulator werden wir über seine Streumatrix definieren, wobei wir annehmen, dass alle Torwiderstände gleich sind. Im Gegensatz zu [Bele68] definieren wir das Element S_{1n} zu $S_{1n} = 1$. Die gesamte Streumatrix lautet

$$\mathbf{S} = \mathbf{S}' = \begin{bmatrix} 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & 1 \\ 1 & 0 & 0 & 0 & \cdot & \cdot & \cdot & 0 & 0 \\ 0 & 1 & 0 & 0 & \cdot & \cdot & \cdot & 0 & 0 \\ 0 & 0 & 1 & 0 & \cdot & \cdot & \cdot & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & 1 & 0 \end{bmatrix}. \quad (2.105)$$

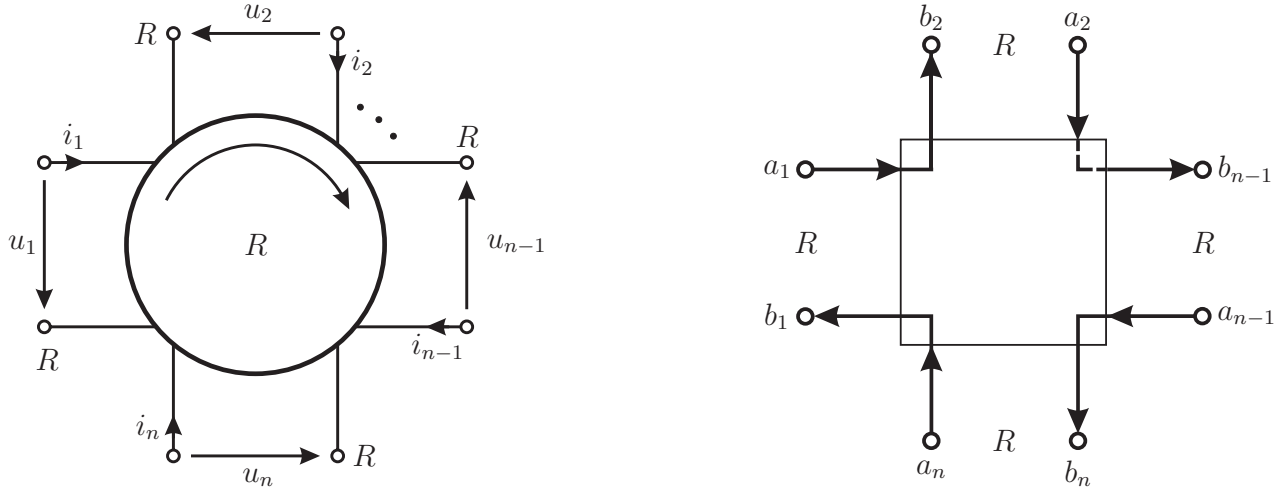


Bild 2.12: n -Tor Zirkulator: a) Kirchhoff'sche Schaltung b) Wellenflussdiagramm

Dass die Streumatrizen gleich sind, erkennt man an der Beziehung Gleichung (2.68) unter Berücksichtigung der Gleichheit der Torwiderstände. Bild 2.12 b) zeigt das Wellenflussdiagramm.

2.8.2 Verallgemeinerte Verbindungsnetze

Die Spannungen und Ströme verallgemeinerter Verbindungsnetze sind unabhängig voneinander Lösungen der homogenen Gleichungssysteme

$$\mathbf{A}\mathbf{i} = \mathbf{0} \quad \text{und} \quad \mathbf{B}\mathbf{u} = \mathbf{0}. \quad (2.106)$$

Die zeilenreguläre Matrix \mathbf{A} hat die Dimension $r \times e$. \mathbf{B} ist ebenfalls zeilenregulär und hat die Dimension $m \times e$. Weiterhin gilt $r + m = e$ und

$$\mathbf{A}\mathbf{B}^T = \mathbf{0} \quad \Leftrightarrow \quad \mathbf{B}\mathbf{A}^T = \mathbf{0}. \quad (2.107)$$

Die r (bzw. m) Spalten von \mathbf{A}^T (\mathbf{B}^T) bilden eine vollständige Basis des Nullraumes von \mathbf{B} (\mathbf{A}). Somit können durch $\mathbf{u} = \mathbf{A}^T \mathbf{u}_r$ ($\mathbf{i} = \mathbf{B}^T \mathbf{i}_m$) mit geeigneten Vektoren \mathbf{u}_r und \mathbf{i}_m sämtliche Lösungen des homogenen Gleichungssystems dargestellt werden [Fisc00a]. Aus der Orthogonalität von \mathbf{u} und \mathbf{i} folgt unmittelbar die Energieneutralität des verallgemeinerten Verbindungsnetzes. Ein verallgemeinertes Verbindungsnetz ist folglich ein spezielles energieneutrales n -Tor. Betrachten wir

$$\mathbf{b} = \mathbf{S}\mathbf{a} \quad \Leftrightarrow \quad \mathbf{G}^{1/2} \mathbf{A}^T \mathbf{u}_r - \mathbf{R}^{1/2} \mathbf{B}^T \mathbf{i}_m = \mathbf{S}[\mathbf{G}^{1/2} \mathbf{A}^T \mathbf{u}_r + \mathbf{R}^{1/2} \mathbf{B}^T \mathbf{i}_m] \quad (2.108)$$

so wird durch spezielle Wahl von \mathbf{u}_r und \mathbf{i}_m die Eigenschaft der Unitarität der Streumatrix durch

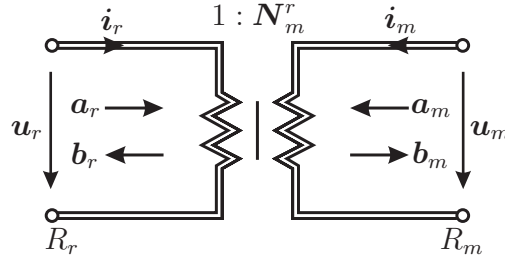
$$\mathbf{S}\mathbf{S} = \mathbf{1}_n \quad \text{und} \quad \mathbf{S}^H = \mathbf{S} \quad (2.109)$$

weiter eingeschränkt, wobei die dritte der drei Eigenschaften aus den jeweils anderen beiden hervorgeht [Frän97].

Die Streumatrizen eines verallgemeinerten Verbindungsnetzes berechnen sich zu

$$\mathbf{S} = 2\mathbf{G}^{1/2} \mathbf{A}^T [\mathbf{A}\mathbf{G}\mathbf{A}^T]^{-1} \mathbf{A}\mathbf{G}^{1/2} - \mathbf{1}_n = \mathbf{1}_n - 2\mathbf{R}^{1/2} \mathbf{B}^T [\mathbf{B}\mathbf{R}\mathbf{B}^T]^{-1} \mathbf{B}\mathbf{R}^{1/2} \quad (2.110)$$

$$\mathbf{S}' = 2\mathbf{A}^T [\mathbf{A}\mathbf{G}\mathbf{A}^T]^{-1} \mathbf{A}\mathbf{G} - \mathbf{1}_n = \mathbf{1}_n - 2\mathbf{R}\mathbf{B}^T [\mathbf{B}\mathbf{R}\mathbf{B}^T]^{-1} \mathbf{B} \quad (2.111)$$

Bild 2.13: Idealer e -Tor Übertrager

vgl. [Shek74], [Meer79], [Frie95]. Aufgrund der Zeilenregularität von \mathbf{A} ist $\mathbf{G}^{1/2} \mathbf{A}^T [\mathbf{A} \mathbf{G} \mathbf{A}^T]^{-1}$ offenbar die Moore-Penrose-Inverse von $\mathbf{A} \mathbf{G}^{1/2}$. Für den Fall gleicher Torwiderstände ergibt sich

$$\mathbf{S} = \mathbf{S}' = 2\mathbf{A}^+ \mathbf{A} - \mathbf{1} = \mathbf{1} - 2\mathbf{B}^+ \mathbf{B}. \quad (2.112)$$

Die Bestimmung der Matrix \mathbf{A} aus der Streumatrix kann wie folgt durchgeführt werden, [Ochs02]. Bekanntermaßen können die zulässigen Zweigspannungen und Zweigströme unabhängig von einander gewählt werden. Mit der Wahl der offenbar zulässigen Zweigspannungen $\mathbf{u} = \mathbf{0}$ geht $\mathbf{b}' = \mathbf{S}' \mathbf{a}'$ in $-\mathbf{R} \mathbf{i} = \mathbf{S}' \mathbf{R} \mathbf{i}$ über. Die (ggf. noch zu normierende) Matrix $[\mathbf{S}' + \mathbf{1}] \mathbf{R}$, bzw. jede andere durch reguläre Transformation daraus erzeugbare Matrix, stellt eine gültige Matrix \mathbf{A} dar. Bei reinen Verbindungsnetzen ist der maximale Rang der Matrix \mathbf{A} durch die Anzahl Eigenwerte von \mathbf{S}' mit dem Wert 1 festgelegt. Die Matrix \mathbf{A} besitzt allerdings erst nach geeigneter regulärer Transformation die Eigenschaften einer Inzidenzmatrix, wie z.B. nur Elemente und Minoren der Form $\{-1, 0, 1\}$.

Ein Bauelement, welches den Eigenschaften von verallgemeinerten Verbindungsnetzen genügt, ist der ideale reelle e -Tor Übertrager, vgl. Bild 2.13. Der ideale Übertrager soll nach Definition durch die Matrizen

$$\mathbf{A} = \begin{bmatrix} r & m \\ \mathbf{1}_r & \mathbf{N}^T \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} r & m \\ -\mathbf{N} & \mathbf{1}_m \end{bmatrix} \quad (2.113)$$

beschrieben werden. Wir nehmen die folgenden Partitionierungen vor

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_r \\ \mathbf{u}_m \end{bmatrix}, \quad \mathbf{i} = \begin{bmatrix} \mathbf{i}_r \\ \mathbf{i}_m \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}_r \\ \mathbf{b}_m \end{bmatrix}, \quad \mathbf{a} = \begin{bmatrix} \mathbf{a}_r \\ \mathbf{a}_m \end{bmatrix}, \quad \mathbf{G} = \text{diag}(\mathbf{G}_r, \mathbf{G}_m). \quad (2.114)$$

Auswerten von Gleichung (2.110) liefert die Streumatrix

$$\mathbf{S} = \left[\begin{array}{c|c} 2\mathbf{G}_r^{1/2} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{G}_r^{1/2} - \mathbf{1}_r & 2\mathbf{G}_r^{1/2} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{N}^T \mathbf{G}_m^{1/2} \\ \hline 2\mathbf{G}_m^{1/2} \mathbf{N} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{G}_r^{1/2} & 2\mathbf{G}_m^{1/2} \mathbf{N} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{N}^T \mathbf{G}_m^{1/2} - \mathbf{1}_m \end{array} \right]. \quad (2.115)$$

Im weiteren Verlauf der Arbeit werden wir den idealen Übertrager benötigen, dessen zugehörige Streumatrix die Form

$$\mathbf{S} = \begin{bmatrix} \mathbf{0} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{0} \end{bmatrix} \quad (2.116)$$

besitzt. Hierbei setzen wir \mathbf{N} als quadratisch und regulär voraus. Aus Gleichung (2.115) erhalten wir die simultan zu erfüllenden Bedingungen

$$\begin{aligned} 2\mathbf{G}_r^{1/2} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{G}_r^{1/2} &= \mathbf{1}_r & \iff & 2\mathbf{G}_r = \mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N} \\ 2\mathbf{G}_m^{1/2} \mathbf{N} [\mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{G}_r^{1/2} &= \mathbf{1}_m & \iff & 2\mathbf{N}^T \mathbf{G}_m \mathbf{N} = \mathbf{G}_r + \mathbf{N}^T \mathbf{G}_m \mathbf{N}. \end{aligned} \quad (2.117)$$

Diese Bedingungen zerfallen -offenbar aufgrund der Eigenschaften von \mathbf{S} - zu der einen Bedingung

$$\mathbf{G}_r = \mathbf{N}^T \mathbf{G}_m \mathbf{N} \iff \mathbf{R}_r \mathbf{N}^T = \mathbf{N}^{-1} \mathbf{R}_m \iff \mathbf{N} \mathbf{R}_r = \mathbf{R}_m \mathbf{N}^{-T}. \quad (2.118)$$

Die zugehörige Spannungswellenstreumatrix wollen wir auf einem alternativen Weg herleiten. Offenbar ist \mathbf{u} (\mathbf{i}^H) Rechtseigenvektor (Linkseigenvektor) zum Eigenwert 1 (-1) der Matrix

$$\begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix}, \quad (2.119)$$

d. h. es gilt

$$\mathbf{u} = \begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix} \mathbf{u} \quad , \quad \mathbf{i} = - \begin{bmatrix} \mathbf{0} & \mathbf{N}^T \\ \mathbf{N}^{-T} & \mathbf{0} \end{bmatrix} \mathbf{i}. \quad (2.120)$$

Somit ergibt sich

$$\mathbf{b}' = \mathbf{u} - \mathbf{R}\mathbf{i} = \begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix} \mathbf{u} + \begin{bmatrix} \mathbf{0} & \mathbf{R}_r \mathbf{N}^T \\ \mathbf{R}_m \mathbf{N}^{-T} & \mathbf{0} \end{bmatrix} \mathbf{i}. \quad (2.121)$$

Mit Gleichung (2.118) erhalten wir

$$\mathbf{b}' = \mathbf{u} - \mathbf{R}\mathbf{i} = \begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix} [\mathbf{u} + \mathbf{R}\mathbf{i}] = \begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix} \mathbf{a}' \quad (2.122)$$

und somit die gewünschte Spannungswellenstreumatrix

$$\mathbf{S}' = \begin{bmatrix} \mathbf{0} & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0} \end{bmatrix}. \quad (2.123)$$

n -Tor-Paralleladaptor

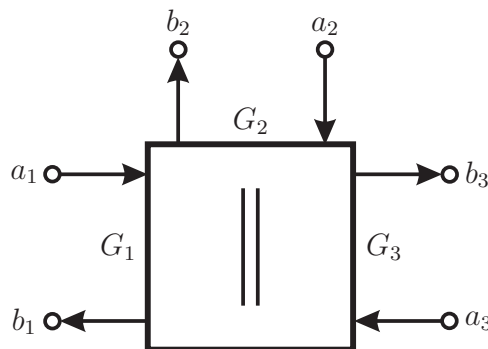


Bild 2.14: 3-Tor-Paralleladaptor

Der n -Tor-Paralleladaptor ist durch die Gleichheit seiner n Torspannungen und durch die Tatsache, dass die n Torströme eine Nullsumme bilden, gekennzeichnet, d. h.

$$u = u_\nu, \nu = 1, 2, \dots, n \quad \wedge \quad \sum_{\nu=1}^n i_\nu = 0. \quad (2.124)$$

Die Torwiderstände $R_\nu = 1/G_\nu$ des Adaptors bestimmen das Verhalten des Wellendigitalfilter-Bausteins. Drückt man nun die $2n$ Variablen in den oben angegebenen n linear unabhängigen Kirchhoff'schen Gleichungen durch die $2n$ Wellengrößen aus und führt die Hilfsvariablen

$$G_0 = \sum_{\nu=1}^n G_\nu, \quad \gamma_\nu = \sqrt{\frac{2G_\nu}{G_0}} \quad \text{und} \quad \gamma'_\nu = \gamma_\nu^2 \quad (2.125)$$

ein, so erhalten wir durch Anwendung von Gleichung (2.110) mit $\mathbf{A} = [1 \ 1 \ \dots \ 1 \ 1]$ die Streumatrizen des Paralleladaptors

$$\mathbf{S} = \begin{bmatrix} \gamma_1^2 - 1 & \gamma_1 \gamma_2 & \dots & \gamma_1 \gamma_n \\ \gamma_1 \gamma_2 & \gamma_2^2 - 1 & & \gamma_2 \gamma_n \\ \vdots & & \ddots & \vdots \\ \gamma_1 \gamma_n & \gamma_2 \gamma_n & \dots & \gamma_n^2 - 1 \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} \gamma'_1 - 1 & \gamma'_2 & \dots & \gamma'_n \\ \gamma'_1 & \gamma'_2 - 1 & & \gamma'_n \\ \vdots & & \ddots & \vdots \\ \gamma'_1 & \gamma'_2 & \dots & \gamma'_n - 1 \end{bmatrix}. \quad (2.126)$$

Werden die Koeffizienten γ_ν und γ'_ν in den Vektoren $\boldsymbol{\gamma} = [\gamma_1 \ \gamma_2 \ \dots \ \gamma_n]^T$ und $\boldsymbol{\gamma}' = [\gamma'_1 \ \gamma'_2 \ \dots \ \gamma'_n]^T$ zusammengefasst, so kann man die Streumatrizen kompakt durch

$$\mathbf{S} = \boldsymbol{\gamma} \boldsymbol{\gamma}^T - \mathbf{1}_n \quad \text{und} \quad \mathbf{S}' = \mathbf{e} \boldsymbol{\gamma}'^T - \mathbf{1}_n \quad (2.127)$$

beschreiben.

Berücksichtigt man, dass $\|\boldsymbol{\gamma}\|^2 = 2$ gilt, so erkennen wir wieder, dass \mathbf{S} eine Householdermatrix ist. Der nichtgebundene Paralleladapter besitzt zudem die Eigenschaft, dass alle Torwiderstände unabhängig voneinander vorgebar sind, wobei wir uns auf positive, endliche Torwiderstände beschränken wollen. Das zugehörige Wellendigitalfiltersymbol ist für $n = 3$ im Bild 2.14 dargestellt.

Gebundener n -Tor-Paralleladapter

Der gebundene n -Tor-Adapter ist nur ein Spezialfall des allgemeinen Leistungswellen-Paralleladaptors, der aber eine so herausragende Bedeutung besitzt, dass er als ein eigenes Wellendigitalfilter-Bauelement aufgefasst wird. Das Besondere ist, dass ein Tor des Adaptors reflexionsfrei ist und sich somit zum Aufbrechen von verzögerungsfreien gerichteten Schleifen eignet. Wir nehmen nun an, dass das reflexionsfreie Tor die Nummer ν hat. Folglich gilt für das Diagonalelement $s_{\nu\nu}$ und für den Adaptorkoeffizienten γ_ν

$$s_{\nu\nu} = 0 \Leftrightarrow \gamma_\nu^2 = 1. \quad (2.128)$$

Mit der Berechnungsvorschrift der Adaptorkoeffizienten

$$\gamma_\nu^2 = 2 \frac{G_\nu}{G_0} \quad \text{folgt} \quad G_0 = 2G_\nu. \quad (2.129)$$

Aus

$$G_0 = G_\nu + \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n G_\mu \quad \text{erhalten wir} \quad G_\nu = \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n G_\mu. \quad (2.130)$$

Die Hilfsgröße G_0 , die Summe aller Torleitwerte, kann somit durch die Torleitwerte der nicht reflexionsfreien Tore berechnet werden, d.h.

$$G_0 = 2 \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n G_\mu. \quad (2.131)$$

Die Adaptorkoeffizienten ergeben sich zu

$$\gamma_\mu = \begin{cases} \sqrt{2G_\mu/G_0} & \text{für } \mu \neq \nu \\ 1 & \text{für } \mu = \nu. \end{cases} \quad \text{und } \gamma'_\mu = \gamma_\mu^2. \quad (2.132)$$

Von $\gamma_\nu^2 = 1$ wird nur die Lösung mit dem positive Vorzeichen akzeptiert, da γ_μ über Gleichung (2.125) als positiv definiert ist. Die Streumatrizen lauten

$$\mathbf{S} = \begin{bmatrix} \gamma_1^2 - 1 & \dots & \gamma_1 & \dots & \gamma_1 \gamma_n \\ \vdots & \ddots & & & \vdots \\ \gamma_1 & \dots & 0 & \dots & \gamma_n \\ \vdots & & & \ddots & \vdots \\ \gamma_1 \gamma_n & \dots & \gamma_n & \dots & \gamma_n^2 - 1 \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} \gamma'_1 - 1 & \dots & 1 & \dots & \gamma'_n \\ \vdots & \ddots & & & \vdots \\ \gamma'_1 & \dots & 0 & \dots & \gamma'_n \\ \vdots & & & \ddots & \vdots \\ \gamma'_1 & \dots & 1 & \dots & \gamma'_n - 1 \end{bmatrix}. \quad (2.133)$$

Das zugehörige Wellendigitalfiltersymbol ist für $n = 3$ im Bild 2.15 dargestellt, wobei der Übersicht halber das reflexionsfreie Tor (hier Tor 3) besonders gekennzeichnet ist, damit auf einem Plan mit Wellendigitalfilterelementen die reflexionsfreien Tore auch ohne Kenntnis der Werte der Torwiderstände sofort erkennbar sind.

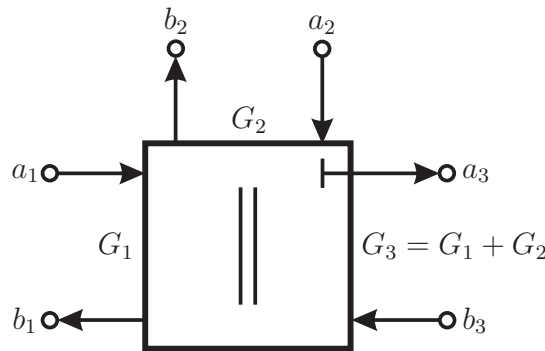


Bild 2.15: Gebundener 3-Tor-Paralleladaptor (hier: Tor 3 reflexionsfrei)

Zum Abschluss wollen wir noch den Fall einer positiv definiten, symmetrischen Torwiderstandsmatrix untersuchen. Dazu definieren wir $\bar{\mathbf{G}} = \mathbf{G}^{1/2}$ und $\boldsymbol{\beta} = [\beta_1, \dots, \beta_n]^T = \bar{\mathbf{G}} \mathbf{A}^T = \bar{\mathbf{G}}[1, \dots, 1]^T$ mit $\beta_l = \sum_{k=1}^n \bar{g}_{lk}$. Nach Gleichung (2.110) lautet die Streumatrix dann

$$\mathbf{S} = 2\boldsymbol{\beta}[\boldsymbol{\beta}^T \boldsymbol{\beta}]^{-1} \boldsymbol{\beta}^T - \mathbf{1}_n. \quad (2.134)$$

Die Forderung nach Reflexionsfreiheit am Tor ν , d. h. $s_{\nu\nu} = 0$ liefert

$$2\beta_\nu^2 = \boldsymbol{\beta}^T \boldsymbol{\beta} = \sum_{\mu=1}^n \beta_\mu^2 \iff \beta_\nu^2 = \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^n \beta_\mu^2. \quad (2.135)$$

Die Frage, die sich nun aufdrängt, ist, ob sich mehr als ein Tor reflexionsfrei wählen lässt. Wir beantworten die Frage für einen später noch relevanten Spezialfall. Und zwar nehmen wir $n = 4$ an und fordern, dass die ersten beiden Tore eine diagonale Torwiderstandsmatrix besitzen. Der Vektor $\boldsymbol{\beta}$ lautet dann

$$\boldsymbol{\beta} = [\bar{g}_{11}, \bar{g}_{22}, \bar{g}_{33} + \bar{g}_{34}, \bar{g}_{43} + \bar{g}_{44}]^T. \quad (2.136)$$

Die Forderungen $s_{11} = 0$ und $s_{22} = 0$ gehen in die simultan zu erfüllenden Beziehungen

$$\bar{g}_{11}^2 = \bar{g}_{22}^2 + \beta_3^2 + \beta_4^2 \quad \text{und} \quad \bar{g}_{22}^2 = \bar{g}_{11}^2 + \beta_3^2 + \beta_4^2 \quad (2.137)$$

über, die letztendlich in

$$\beta_3^2 + \beta_4^2 = 0 \iff (\bar{g}_{33} + \bar{g}_{43})^2 + (\bar{g}_{34} + \bar{g}_{44})^2 = 0 \quad (2.138)$$

münden. Eine Forderung, die aufgrund der Reellwertigkeit von $\mathbf{G}^{1/2}$ nur durch $\bar{g}_{33} = -\bar{g}_{43}$ und $\bar{g}_{34} = -\bar{g}_{44}$ erfüllt werden kann. Diese einzige Lösung bedeutet aber eine lineare Abhängigkeit der Zeilenvektoren von $\mathbf{G}^{1/2}$, was im Widerspruch zur vorausgesetzten Regularität von \mathbf{G} steht. In dem hier betrachteten Fall existieren somit keine 2 reflexionsfreien Tore.

n -Tor-Serienadaptor

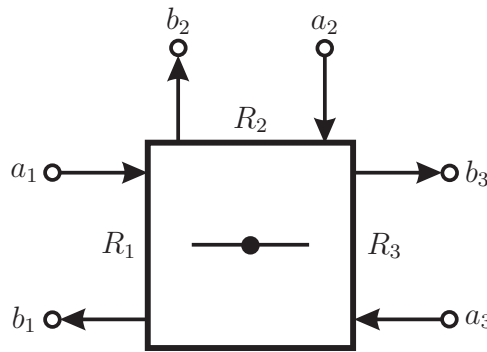


Bild 2.16: 3-Tor-Serienadaptor

Der n -Tor-Serienadaptor ist durch die Gleichheit seiner n Torströme und durch die Tatsache, dass die n Torspannungen eine Nullsumme bilden, gekennzeichnet, d. h.

$$i = i_\nu, \nu = 1, 2, \dots, n \quad \wedge \quad \sum_{\nu=1}^n u_\nu = 0. \quad (2.139)$$

Drückt man nun wiederum die $2n$ Variablen in den oben angegebenen n linear unabhängigen Kirchhoff'schen Gleichungen durch die $2n$ Wellengrößen aus und führt die Hilfsvariablen

$$R_0 = \sum_{\nu=1}^n R_\nu, \quad \gamma_\nu = \sqrt{\frac{2R_\nu}{R_0}} \quad \text{und} \quad \gamma'_\nu = \gamma_\nu^2 \quad (2.140)$$

ein, so erhalten wir durch Anwendung von Gleichung (2.110) die Streumatrizen des Serienadaptors zu

$$\mathbf{S} = \begin{bmatrix} 1-\gamma_1^2 & -\gamma_1\gamma_2 & \cdots & -\gamma_1\gamma_n \\ -\gamma_1\gamma_2 & 1-\gamma_2^2 & & -\gamma_2\gamma_n \\ \vdots & & \ddots & \vdots \\ -\gamma_1\gamma_n & -\gamma_2\gamma_n & \cdots & 1-\gamma_n^2 \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} 1-\gamma'_1 & -\gamma'_1 & \cdots & -\gamma'_1 \\ -\gamma'_2 & 1-\gamma'_2 & & -\gamma'_2 \\ \vdots & & \ddots & \vdots \\ -\gamma'_n & -\gamma'_n & \cdots & 1-\gamma'_n \end{bmatrix}, \quad (2.141)$$

bzw.

$$\mathbf{S} = \mathbf{1}_n - \boldsymbol{\gamma} \boldsymbol{\gamma}^T \quad \text{und} \quad \mathbf{S}' = \mathbf{1}_n - \boldsymbol{\gamma}' \mathbf{e}^T. \quad (2.142)$$

Das zugehörige Wellendigitalfiltersymbol findet sich für $n = 3$ im Bild 2.16.

Gebundener n -Tor-Serienadaptor

Wir nehmen für den gebundenen n -Tor-Serienadaptor auch an, dass sein reflexionsfreies Tor die Nummer ν hat. Folglich gilt für das Diagonalelement $s_{\nu\nu}$ und für den Adaptorkoeffizienten γ_ν

$$s_{\nu\nu} = 0 \Leftrightarrow \gamma_\nu^2 = 1. \quad (2.143)$$

Mit der Berechnungsvorschrift der Adaptorkoeffizienten

$$\gamma_\mu^2 = 2 \frac{R_\mu}{R_0} \quad \text{folgt} \quad R_0 = 2R_\nu. \quad (2.144)$$

Aus

$$R_0 = R_\nu + \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n R_\mu \quad \text{erhalten wir} \quad R_\nu = \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n R_\mu. \quad (2.145)$$

Die Hilfsgröße R_0 , die Summe aller Torwiderstände, kann somit durch die Torwiderstände der nicht reflexionsfreien Tore berechnet werden, d.h.

$$R_0 = 2 \sum_{\substack{\mu \neq \nu \\ \mu=1}}^n R_\mu. \quad (2.146)$$

Die Adaptorkoeffizienten ergeben sich zu

$$\gamma_\mu = \begin{cases} \sqrt{2R_\mu/R_0} & \text{für } \mu \neq \nu \\ 1 & \text{für } \mu = \nu. \end{cases}, \quad \gamma'_\mu = \gamma_\mu^2. \quad (2.147)$$

Die negative Lösung für γ_μ entfällt aus den selben Gründen wie im Falle des gebundenen Paralleladaptors. Die Streumatrizen lauten

$$\mathbf{S} = \begin{bmatrix} 1-\gamma_1^2 & \dots & -\gamma_1 & \dots & -\gamma_1\gamma_n \\ \vdots & \ddots & & & \vdots \\ -\gamma_1 & \dots & 0 & \dots & -\gamma_n \\ \vdots & & & \ddots & \vdots \\ -\gamma_1\gamma_n & \dots & -\gamma_n & \dots & 1-\gamma_n^2 \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} 1-\gamma'_1 & \dots & -\gamma'_1 & \dots & -\gamma'_1 \\ \vdots & \ddots & & & \vdots \\ -1 & \dots & 0 & \dots & -1 \\ \vdots & & & \ddots & \vdots \\ -\gamma'_n & \dots & -\gamma'_n & \dots & 1-\gamma'_n \end{bmatrix}. \quad (2.148)$$

Das zugehörige Wellendigitalfiltersymbol findet sich für $n = 3$ im Bild 2.17. Zu bemerken ist auch hier, dass der Torwiderstand R_ν nicht in die Streumatrix eingeht.

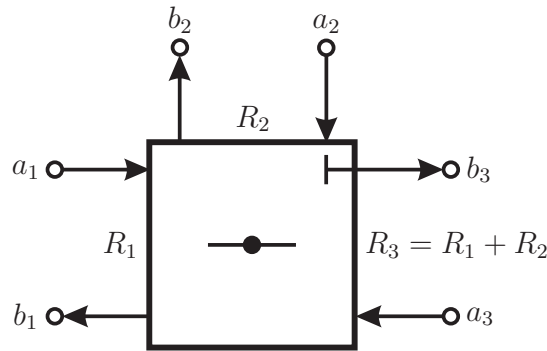


Bild 2.17: Gebundener-3-Tor-Serienadaptor (hier: Tor 3 reflexionsfrei)

2-Tor-Übertrager

Der 2-Tor-Übertrager wird durch

$$u_2 = nu_1 \quad \text{und} \quad i_1 = -ni_2 \quad (2.149)$$

definiert (vgl. Bild 2.18). Die Streumatrizen lauten

$$\mathbf{S} = \begin{bmatrix} \frac{R_2 - n^2 R_1}{R_2 + n^2 R_1} & \frac{2n\sqrt{R_1 R_2}}{R_2 + n^2 R_1} \\ \frac{2n\sqrt{R_1 R_2}}{R_2 + n^2 R_1} & \frac{n^2 R_1 - R_2}{R_2 + n^2 R_1} \end{bmatrix}, \quad \mathbf{S}' = \begin{bmatrix} \frac{R_2 - n^2 R_1}{R_2 + n^2 R_1} & \frac{2nR_1}{R_2 + n^2 R_1} \\ \frac{2nR_2}{R_2 + n^2 R_1} & \frac{n^2 R_1 - R_2}{R_2 + n^2 R_1} \end{bmatrix}. \quad (2.150)$$

Wählen wir das Übersetzungsverhältnis zu $n = \sqrt{R_2/R_1}$, so vereinfachen sich die Streumatrizen zu

$$\mathbf{S} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{und} \quad \mathbf{S}' = \begin{bmatrix} 0 & \frac{1}{n} \\ n & 0 \end{bmatrix}. \quad (2.151)$$

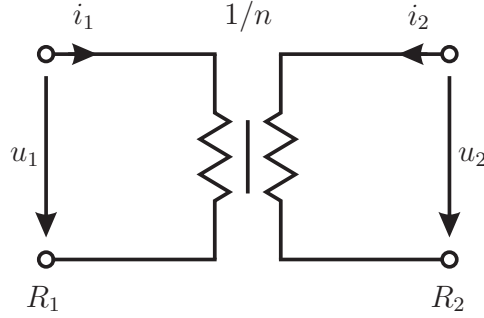


Bild 2.18: 2-Tor-Übertrager

Jaumann-Adaptor

Der im Bild 2.19 dargestellte Jaumann-Adaptor stellt ein verallgemeinertes Verbindungsnetz gemäß Gleichung (2.113) mit

$$\mathbf{N} = \frac{1}{1+n} \begin{bmatrix} 1 & -1 \\ n & 1 \end{bmatrix}, \quad \mathbf{N}^{-1} = \begin{bmatrix} 1 & 1 \\ -n & 1 \end{bmatrix} \quad (2.152)$$

dar. Der freie Parameter ist das Übersetzungsverhältnis des Übertragers. Wir nehmen an, es gilt $n \neq -1$, d.h. wir schließen den Fall aus, dass der Übertrager eine reine Parallelverbindung ist. \mathbf{N} ist somit regulär. Die Streumatrix des Jaumann-Adaptors sollte eine möglichst einfache Struktur besitzen. Aus diesem

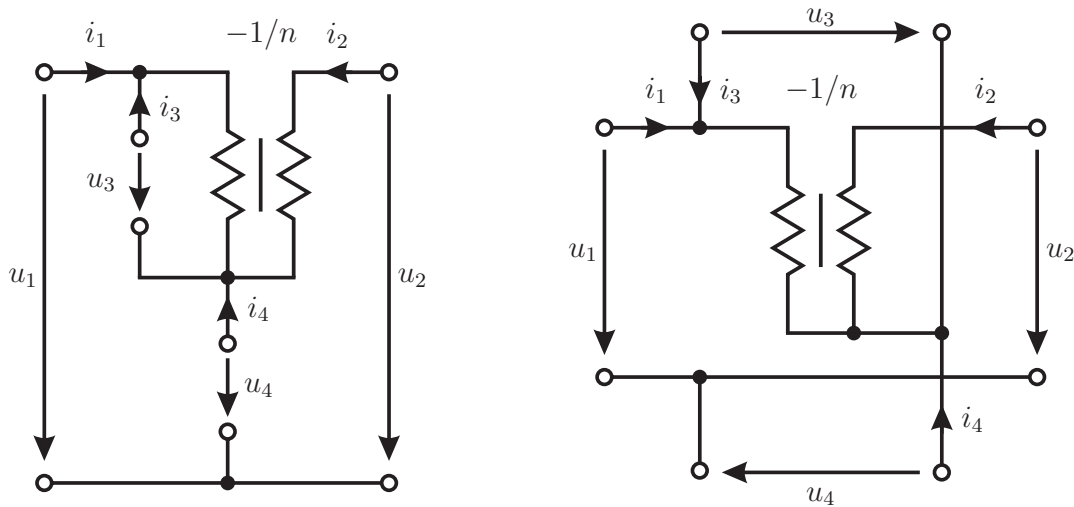


Bild 2.19: Jaumann-Adaptor

Grund wählen wir die Torwiderstandsmatrix gemäß Gleichung (2.118) zu

$$\mathbf{R} = R_1 \mathbf{diag} \left(1, n, \frac{1}{n+1}, \frac{n}{n+1} \right). \quad (2.153)$$

Durch Anwendung von Gleichung (2.123) lässt sich die Spannungswellen-Streumatrix zu

$$\mathbf{S}' = \begin{bmatrix} \mathbf{0}_2 & \mathbf{N}^{-1} \\ \mathbf{N} & \mathbf{0}_2 \end{bmatrix} \quad (2.154)$$

ablesen. Mit Gleichung (2.68) ermitteln wir die Leistungswellen-Streumatrix zu

$$\mathbf{S} = \begin{bmatrix} \mathbf{0} & \mathbf{S}_1 \\ \mathbf{S}_1^T & \mathbf{0} \end{bmatrix}, \quad (2.155)$$

wobei sich die Untermatrix \mathbf{S}_1 zu

$$\mathbf{S}_1 = \begin{bmatrix} \sqrt{\frac{1}{n+1}} & \sqrt{\frac{n}{n+1}} \\ -\sqrt{\frac{n}{n+1}} & \sqrt{\frac{1}{n+1}} \end{bmatrix} \quad (2.156)$$

ergibt. Im Bild 2.20 ist das in dieser Arbeit verwendete WDF-Symbol des Jaumann-Adaptors dargestellt. Die dickere Kante zwischen den Toren 2 und 3 der Raute trägt der Tatsache Rechnung, dass nur die Matricelemente s_{23} und s_{32} negativ sind.

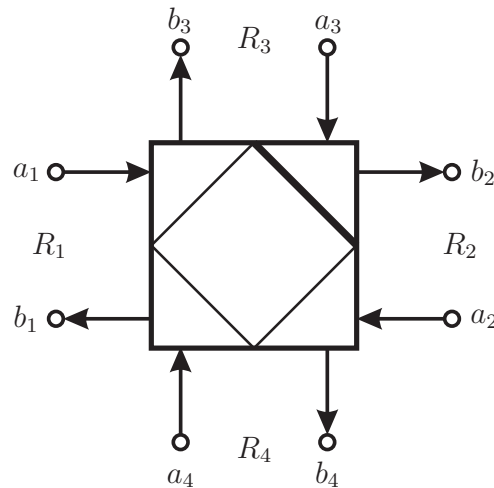


Bild 2.20: Wellendigital-Realisierung des Jaumann-Adaptors

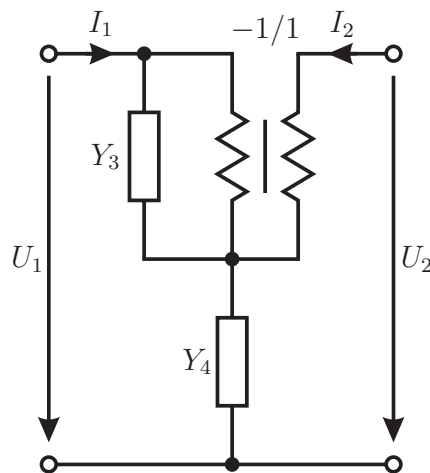


Bild 2.21: Beschlalteter Jaumann-Adaptor

In der Regel nutzt man einen Jaumann-Adaptor in Verbindung mit zwei Impedanzen Z_3 und Z_4 , mit denen die Tore 3 und 4 abgeschlossen werden (Bild 2.21), d. h. es gilt

$$\begin{bmatrix} U_3 \\ U_4 \end{bmatrix} = \begin{bmatrix} -Z_3 & 0 \\ 0 & -Z_4 \end{bmatrix} \begin{bmatrix} I_3 \\ I_4 \end{bmatrix}. \quad (2.157)$$

Durch Verwendung von Gleichung (2.152) i.V.m. Gleichung (2.120) erhält man zunächst

$$\begin{bmatrix} U_1 \\ U_2 \end{bmatrix} = \mathbf{N}^{-1} \begin{bmatrix} -Z_3 & 0 \\ 0 & -Z_4 \end{bmatrix} (-\mathbf{N}^{-T}) \begin{bmatrix} I_1 \\ I_2 \end{bmatrix}, \quad (2.158)$$

woraus man die Impedanzmatrix bezogen auf die Tore 1 und 2 zu

$$\mathbf{Z} = \begin{bmatrix} Z_4 + Z_3 & Z_4 - nZ_3 \\ Z_4 - nZ_3 & Z_4 + n^2Z_3 \end{bmatrix} \quad (2.159)$$

ermittelt. Zur Berechnung der Admittanzmatrix schreibt man zunächst

$$\begin{bmatrix} I_1 \\ I_2 \end{bmatrix} = -\mathbf{N}^T \begin{bmatrix} -\frac{1}{Z_3} & 0 \\ 0 & -\frac{1}{Z_4} \end{bmatrix} \mathbf{N} \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} \quad (2.160)$$

und erhält die Admittanzmatrix bezogen auf die Tore 1 und 2 zu

$$\mathbf{Y} = \frac{1}{(1+n)^2} \begin{bmatrix} \frac{n^2}{Z_4} + \frac{1}{Z_3} & \frac{n}{Z_4} - \frac{1}{Z_3} \\ \frac{n}{Z_4} - \frac{1}{Z_3} & \frac{1}{Z_4} + \frac{1}{Z_3} \end{bmatrix}. \quad (2.161)$$

2.8.3 Quellen

Wir betrachten die resistive Spannungsquelle im Bild 2.22 a) mit Innenwiderstand R_i , Quellenspannung e und Ausgangsspannung u_q . Wir ersetzen in der Maschengleichung $u_q = e + R_i i_q$ die auftretenden Spannungen und Ströme durch die Wellengrößen und erhalten

$$a'_q + b'_q = 2e + R_i G_q [a'_q - b'_q] \iff b'_q = \frac{2e}{1 + R_i G_q} + \rho a'_q \quad \text{mit} \quad \rho = \frac{R_i - R_q}{R_i + R_q}. \quad (2.162)$$

Das zugehörige Wellenflussdiagramm zeigt Bild 2.22 b). Das Wellenflussdiagramm vereinfacht sich, wenn wir den Torwiderstand gleich dem Innenwiderstand wählen $R_q = R_i$. Die ausfallende Welle ist nun unabhängig von der einfallenden Welle $b'_q = 2\sqrt{R_q}b_q = e$, vgl. Bild 2.22 c). Die ideale Spannungsquelle ($R_i = 0$) wird in den Wellengrößen durch $b'_q = 2e - a'_q$ beschrieben.

Analoge Ergebnisse erhalten wir bei der Untersuchung von Stromquellen. Für die im Bild 2.23 a) dargestellte resistive Stromquelle lautet die Knotengleichung $j + i_q = G_i u_q$. Wir erhalten

$$2j + G_q [a'_q - b'_q] = G_i [a'_q + b'_q] \iff b'_q = \frac{2j}{G_i + G_q} + \rho a'_q \quad \text{mit} \quad \rho = \frac{R_i - R_q}{R_i + R_q}. \quad (2.163)$$

Als Wellendigitalfiltersymbol verwenden wir das gleiche Symbol wie bei den Spannungswellen, nur, dass in diesem Fall sich die Quellenwelle anders berechnet, vgl. Bild 2.23 b) und c).

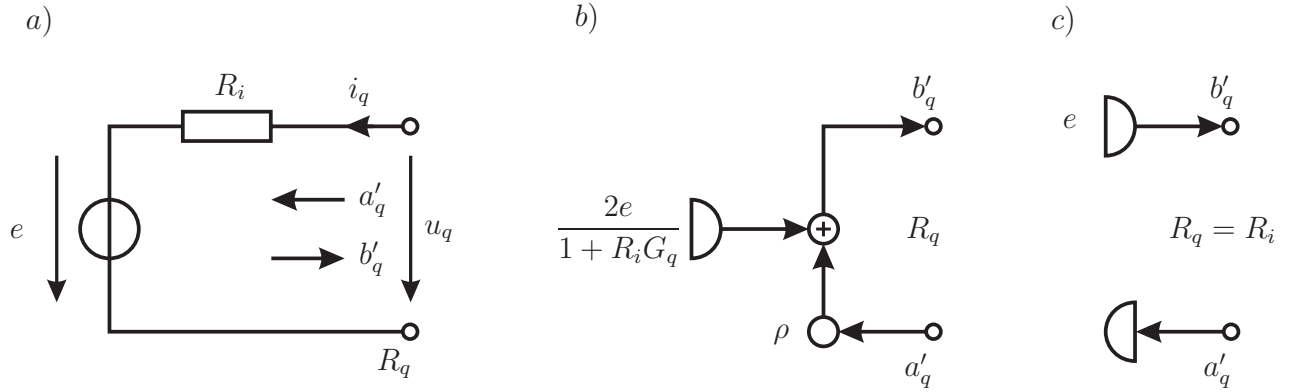


Bild 2.22: Resistive Spannungsquelle a) Kirchhoff'sche Schaltung b) und c) Wellenflussdiagramm

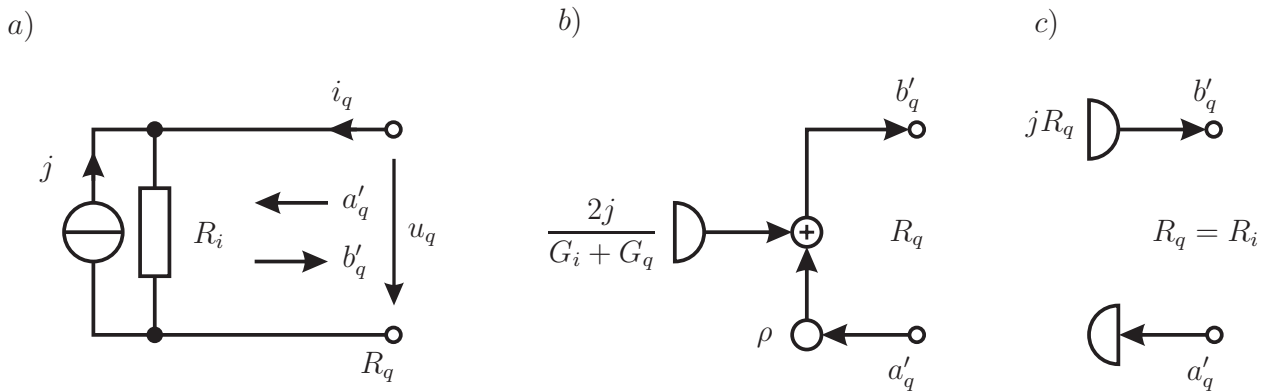
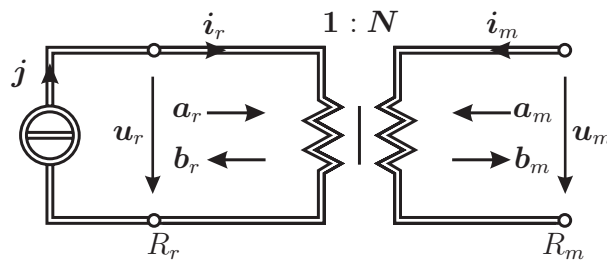


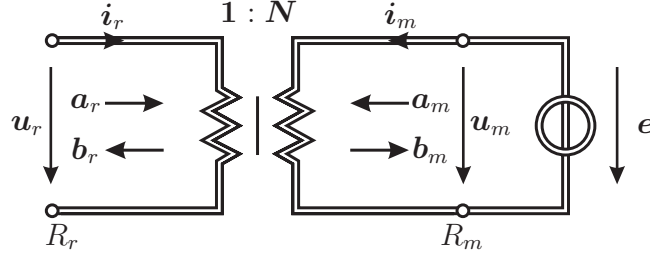
Bild 2.23: Resistive Stromquelle a) Kirchhoff'sche Schaltung b) und c) Wellenflussdiagramm

Neben den idealen Quellen kann eine Vielzahl weiterer Quellen betrachtet werden. Von besonders praktischer Relevanz sind jedoch ideale Quellen, die in einem verallgemeinerten Verbindungsnetz eingebettet sind. Um dies zu untersuchen, betrachten wir den mit einer vektoriellen idealen Stromquelle abgeschlossenen idealen e -Tor-Übertrager im Bild 2.24. Im Folgenden ist die Beschreibung durch die Wellengrößen gesucht. Der Übertrager soll Gleichung (2.113) mit den Bezeichnungen aus Gleichung (2.114) folgen. Wir nehmen $m \geq r$ an, und \mathbf{N}^T sei zeilenregulär. Obwohl die Ströme \mathbf{i}_m unabhängig sind, wollen wir die Stromquellen in die r -Zweige legen. Der Vektor $\mathbf{j} = \mathbf{i}_r$ liegt dann im Spaltenraum von \mathbf{N}^T . Gesucht ist nun eine Beziehung der Form

$$\mathbf{b}'_m = \mathbf{S}'\mathbf{a}'_m + \mathbf{Q}'\mathbf{j}. \quad (2.164)$$

Wir multiplizieren dazu die Gleichung $\mathbf{a}'_m = \mathbf{u}_m + \mathbf{R}_m\mathbf{i}_m$ von links mit $\mathbf{N}^T\mathbf{G}_m$, drücken \mathbf{u}_m , \mathbf{i}_m durch

Bild 2.24: e -Tor Übertrager mit idealen Stromquellen

Bild 2.25: e -Tor Übertrager mit idealen Spannungsquellen

\mathbf{u}_r , \mathbf{i}_r aus und lösen anschließend nach \mathbf{u}_r auf

$$\mathbf{N}^T \mathbf{G}_m \mathbf{a}'_m = \mathbf{N}^T \mathbf{G}_m \mathbf{N} \mathbf{u}_r - \mathbf{i}_r \iff \mathbf{u}_r = [\mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} [\mathbf{N}^T \mathbf{G}_m \mathbf{a}'_m + \mathbf{i}_r]. \quad (2.165)$$

Dieses setzen wir in $\mathbf{b}'_m = 2\mathbf{u}_m - \mathbf{a}'_m = 2\mathbf{N}\mathbf{u}_r - \mathbf{a}'_m$ ein und erhalten

$$\mathbf{b}'_m = \underbrace{[2\mathbf{N}[\mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{N}^T \mathbf{G}_m - \mathbf{1}_m]}_{\mathbf{S}'} \mathbf{a}'_m + \underbrace{2\mathbf{N}[\mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{j}}_{\mathbf{Q}'}. \quad (2.166)$$

Wir wollen nun die Freiheitsgrade diskutieren. Die Spannungen \mathbf{u}_m müssen im Spaltenraum von \mathbf{N} liegen, d.h. r Spannungen sind frei vorgebar. Die Ströme \mathbf{i}_m müssen $\mathbf{N}^T \mathbf{i}_m = -\mathbf{j}$ bei gegebener Quellenverteilung genügen, d.h.

$$\begin{aligned} \mathbf{i}_m &= -\mathbf{N}[\mathbf{N}^T \mathbf{N}]^{-1} \mathbf{j} + [\mathbf{1}_m - \mathbf{N}[\mathbf{N}^T \mathbf{N}]^{-1} \mathbf{N}^T] \mathbf{z} = -(\mathbf{N}^T)^+ \mathbf{j} + [\mathbf{1}_m - (\mathbf{N}^T)^+ \mathbf{N}^T] \mathbf{z} \\ &= -(\mathbf{N}^T)^+ \mathbf{j} + [\mathbf{1}_m - \mathbf{N} \mathbf{N}^+]^T \mathbf{z} = -(\mathbf{N}^T)^+ \mathbf{j} + [\mathbf{1}_m - \mathbf{N} \mathbf{N}^+] \mathbf{z}, \quad \mathbf{z} \text{ beliebig.} \end{aligned} \quad (2.167)$$

Die Ströme \mathbf{i}_m liegen daher im r -dimensionalen Unterraum des \mathbb{R}^m . Dieser Teil ist durch die Quellenverteilung festgelegt. Somit verbleiben noch die $m - r$ Freiheitsgrade des orthogonalen Raumes. Insgesamt haben wir m Freiheitsgrade, die durch die äußere Beschaltung am rechten Tor festgelegt werden. Wird z.B. \mathbf{a}'_m durch die äußere Beschaltung vorgegeben, so liegt \mathbf{b}'_m eindeutig fest. (Vgl.: in einem Netz lassen sich die r Baumzweigspannungen und die $m = e - r$ Cobaumströme unabhängig vorgeben, also e Größen insgesamt. Durch die Quellenvorgabe sind hier r Größen festgelegt. Es verbleiben $e - r = m$ Freiheitsgrade.)

Der Vollständigkeit halber wollen wir noch die Beziehung

$$\mathbf{b}_m = \mathbf{S} \mathbf{a}_m + \mathbf{Q} \mathbf{j}. \quad (2.168)$$

für die Leistungswellen herleiten. Mit Gleichung (2.68) folgt aus Gleichung (2.166)

$$\mathbf{b}_m = \underbrace{[2\mathbf{G}_m^{1/2} \mathbf{N}[\mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{N}^T \mathbf{G}_m^{1/2} - \mathbf{1}_m]}_{\mathbf{S}} \mathbf{a}_m + \underbrace{\mathbf{G}_m^{1/2} \mathbf{N}[\mathbf{N}^T \mathbf{G}_m \mathbf{N}]^{-1} \mathbf{j}}_{\mathbf{Q}}. \quad (2.169)$$

Bei Leerlauf der r -Tore reduziert sich Gleichung (2.169) auf Gleichung (2.110) mit $\mathbf{A} = \pm \mathbf{N}^T$.

Nun betrachten wir den mit einer vektoriellen idealen Spannungsquelle abgeschlossenen idealen e -Tor-Übertrager aus Bild 2.25. Es gilt $\mathbf{u}_m = \mathbf{e}$. Wir nehmen $m \leq r$ an, und \mathbf{N} sei zeilenregulär. Wir erhalten in den Wellengrößen

$$\mathbf{b}'_r = \underbrace{[\mathbf{1}_r - 2\mathbf{R}_r \mathbf{N}^T [\mathbf{N} \mathbf{R}_r \mathbf{N}^T]^{-1} \mathbf{N}]}_{\mathbf{S}'} \mathbf{a}'_r + \underbrace{2\mathbf{R}_r \mathbf{N}^T [\mathbf{N} \mathbf{R}_r \mathbf{N}^T]^{-1}}_{\mathbf{Q}'} \mathbf{e}. \quad (2.170)$$

Bei Kurzschluss der m -Tore, d.h. $\mathbf{e} = \mathbf{0}$ reduziert sich Gleichung (2.170) zu Gleichung (2.111) mit $\mathbf{B} = \pm \mathbf{N}$.

2.8.4 Dissipative Bauelemente

Widerstand

Der Widerstand R kann als reale Spannungsquelle mit der Quellenspannung $E = 0$ aufgefasst werden, womit der Widerstand schon erfasst wäre. Wir werden aber dennoch ein eigenes Bauelement einführen, da die Quellen bei uns die Eingänge eines Funktionsbausteins (der Begriff wird später noch erläutert) sind. Die Streumatrix ist im Falle des Widerstandes ein Skalar, da der Widerstand nur ein Tor besitzt

$$\mathbf{S} = \mathbf{S}' = \frac{R - R_T}{R + R_T}. \quad (2.171)$$

Üblicherweise bezeichnet man diese Streumatrix als Reflektanz ρ . Das Wellendigitalfiltersymbol ist erneut aus dem Signalflussgraph abgeleitet und ist ein Multiplizierer (siehe Bild 2.26).

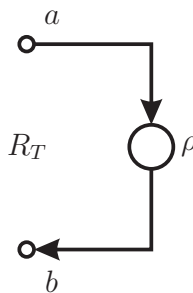


Bild 2.26: Wellenflussdiagramm eines Widerstands

2.8.5 Dynamische Bauelemente

Das einzige hier verwendete dynamische Bauelement ist der ideale Kondensator. Ideale Spulen lassen sich aus Gyrotoren und idealen Kondensatoren realisieren. Die idealen Spulen haben gegenüber den idealen Kondensatoren den Nachteil einer zusätzlichen Negation in Wellenflussdiagramm. Der ideale Kondensator wird durch

$$\hat{i} = \sqrt{C} D\{\sqrt{C}\hat{u}\} \quad (2.172)$$

beschrieben, wobei C eine Funktion der Zustandsgrößen und der unabhängigen Größen ist. Der Ableitungsoperator ist $D = \boldsymbol{\alpha}^T \mathbf{D}_t$ mit $\boldsymbol{\alpha} = \text{const.}$ und $\|\boldsymbol{\alpha}\| = 1$. Diese aus [Fett92] stammende Darstellung ist völlig gleichwertig zur ursprünglichen in [MF92] vorgeschlagenen Definition des nichtlinearen, zeitinvarianten idealen Kondensators, bietet aber den Vorteil übersichtlicherer Stromlaufpläne.³

Wir wollen zunächst eine diskrete Approximation für einen konstanten idealen Kondensator herleiten, der nur bzgl. der Richtung t_κ reaktiv ist, d. h. $\boldsymbol{\alpha} = \mathbf{e}_\kappa$, siehe Bild 2.27. Integrieren wir den Strom entlang einer Geraden vom Punkt \mathbf{t} nach $\mathbf{t} + T_0 \mathbf{e}_\kappa = \mathbf{t} + \Delta \mathbf{t}_\kappa$ und werten das Integral mittels der Trapezregel⁴ aus, so erhalten wir einen Näherungswert für die Spannung am Punkt $\mathbf{t} + T_0 \mathbf{e}_\kappa$

$$\hat{u}(\mathbf{t} + T_0 \mathbf{e}_\kappa) \approx \hat{u}(\mathbf{t}) + R \left[\hat{i}(\mathbf{t} + T_0 \mathbf{e}_\kappa) + \hat{i}(\mathbf{t}) \right] \quad \text{mit} \quad R = \frac{\|\Delta \mathbf{t}_\kappa\|}{2C} = \frac{T_0}{2C}. \quad (2.173)$$

³Genau genommen haben die angegebenen Literaturstellen anstelle des idealen Kondensators eine ideale Spule benutzt.

⁴Wir verwenden hier die Trapezregel, weil sie einfach anzuwenden ist. Zudem hat die Trapezregel die größte Konsistenzordnung und den kleinsten lokalen Fehler aller A-stabilen (passiven) linearen Mehrschritt-Verfahren, [Dahl63], [Ochs01a]. Die Verwendung von Runge-Kutta-Verfahren unter Beibehaltung der Passivität wurde im eindimensionalen Fall in [Ochs01b] gezeigt. Dem Autor der hier vorliegenden Arbeit sind jedoch keine mehrstufigen passiven Runge-Kutta-Verfahren bekannt, die auf allgemeine PDGLn anwendbar sind.

Von nun an werden nicht mehr die kontinuierlichen unabhängigen Variablen verwendet, sondern die diskreten. Außerdem werden die mit einem Dach versehenen Größen durch Näherungsgrößen ersetzt

$$u(\boldsymbol{\nu} + \mathbf{e}_\kappa) = u(\boldsymbol{\nu}) + R [i(\boldsymbol{\nu} + \mathbf{e}_\kappa) + i(\boldsymbol{\nu})]. \quad (2.174)$$

Nach Umsortieren der letzten Gleichung und Verwendung der Wellengrößen mit dem Torwiderstand R erhalten wir

$$b(\boldsymbol{\nu} + \mathbf{e}_\kappa) = a(\boldsymbol{\nu}), \quad (2.175)$$

vgl. Bild 2.28. Wir fragen uns nun, wie sich die gefundenen Zusammenhänge im ursprünglichen Koordi-

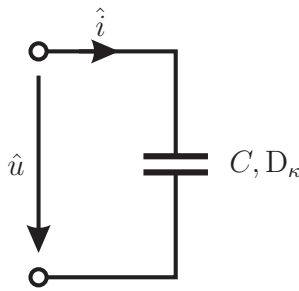


Bild 2.27: Idealer Kondensator, der bzgl. t_κ reaktiv ist

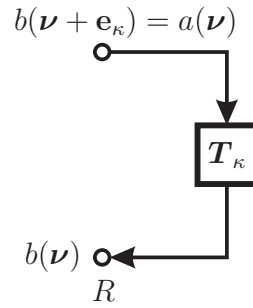


Bild 2.28: Näherungsweise Nachbildung durch ein diskretes System

natensystem \mathbf{x} darstellen. Die Beschreibung des idealen Kondensators im ursprünglichen Koordinatensystem gewinnen wir durch Anwendung von Gleichung (2.16) auf Gleichung (2.172) zu

$$\hat{i} = C D\{\hat{u}\} = \tilde{C} \frac{\mathbf{h}_\kappa^T}{\|\mathbf{h}_\kappa\|} \mathbf{D}\mathbf{x}\{\hat{u}\} \quad \text{mit} \quad \tilde{C} = C v_0 \|\mathbf{h}_\kappa\|. \quad (2.176)$$

Der zuvor durchgeführten Integration im Koordinatensystem \mathbf{t} entspricht die Integration vom Punkt \mathbf{x} nach $\mathbf{x} + \Delta\mathbf{x}_\kappa$ mit $\Delta\mathbf{x}_\kappa = v_0 T_0 \mathbf{h}_\kappa$ über das Kurvenstück

$$K = \{\boldsymbol{\xi}(\lambda) = \mathbf{x} + \lambda \Delta\mathbf{x}_\kappa \mid 0 \leq \lambda \leq 1\}. \quad (2.177)$$

Nach Umformung des Integrals in Parameterform und Berechnung mittels der Trapezregel erhalten wir einen Näherungswert für die Spannung am Punkt $\mathbf{x} + \Delta\mathbf{x}_\kappa$

$$\tilde{u}(\mathbf{x} + \Delta\mathbf{x}_\kappa) = \tilde{u}(\mathbf{x}) + \|\Delta\mathbf{x}_\kappa\| \frac{1}{2\tilde{C}} [\tilde{i}(\mathbf{x} + \Delta\mathbf{x}_\kappa) + \tilde{i}(\mathbf{x})]. \quad (2.178)$$

Der Torwiderstand lautet hier

$$R = \frac{\|\Delta\mathbf{x}_\kappa\|}{2\tilde{C}} = \frac{T_0}{2C}. \quad (2.179)$$

In den in dieser Arbeit verwendeten Zeichnungen bezieht sich der Kapazitätswert immer auf die Darstellung im neuen Koordinatensystem, vgl. Bild 2.27. Im Sinne der obigen Diskussion ist darunter also immer C und nicht \tilde{C} zu verstehen. Wie gezeigt werden kann, sind sowohl der ideale Kondensator als auch die diskrete Nachbildung mehrdimensional intern verlustfreie Bauelemente, siehe z. B. [Fett98].

Der Verschiebevektor $\Delta\boldsymbol{\nu}_\kappa = \mathbf{e}_\kappa$ geht im Koordinatensystem $\boldsymbol{\mu}$ zu

$$\Delta\boldsymbol{\mu}_\kappa = \mathbf{X}_A^{-1} v_0 \mathbf{H} T_0 \mathbf{e}_\kappa = \mathbf{X}_A^{-1} v_0 T_0 \mathbf{h}_\kappa \quad (2.180)$$

über. Da dieser Vektor im Zusammenhang mit der Randbehandlung eine wichtige Rolle spielt, führen wir für den aus den ersten 3 Koordinaten gebildeten Untervektor die besondere Bezeichnung

$$\Delta \boldsymbol{\mu}_\kappa = \begin{bmatrix} \Delta \boldsymbol{\mu}'_\kappa \\ 1 \end{bmatrix} \quad (2.181)$$

ein. Ebenso teilen wir die Spalten von \mathbf{H} wie folgt auf

$$\mathbf{h}_{\bullet, \kappa} = \begin{bmatrix} \mathbf{h}'_\kappa \\ h_{k, \kappa} \end{bmatrix}. \quad (2.182)$$

Im Folgenden soll eine Beziehung zwischen der Spannung des idealen Kondensators und den Wellengrößen hergeleitet werden, die im Zusammenhang mit der Randbehandlung erhebliche Bedeutung besitzt, vgl. [Frie95], [Fett98] und [Sera00]. Dazu entwickeln wir zunächst die Funktion $f(\mathbf{t}) = f[(\mathbf{t} \pm \mathbf{T}/2) \mp \mathbf{T}/2]$ bezüglich der Variablen $\mp \mathbf{T}/2$ um den Punkt $\mathbf{t} \pm \mathbf{T}/2$ in eine Taylor-Reihe

$$f(\mathbf{t}) = f(\mathbf{t} \pm \mathbf{T}/2) \mp \frac{\mathbf{T}^T}{2} \mathbf{D}_\tau f(\tau) |_{\tau=\mathbf{t} \pm \mathbf{T}/2} + \mathcal{O}(\|\mathbf{T}\|^2), \quad (2.183)$$

wobei

$$\mathbf{D}_\tau f(\tau) |_{\tau=\mathbf{t} \pm \mathbf{T}/2} = \mathbf{D}_\mathbf{t} f(\mathbf{t} \pm \mathbf{T}/2) \quad (2.184)$$

gilt. Aus dieser Beziehung gewinnen wir zum einen durch Wahl der oberen Vorzeichen und der Substitution $\mathbf{t} \Rightarrow \mathbf{t} - \mathbf{T}/2$

$$f(\mathbf{t} - \mathbf{T}/2) = f(\mathbf{t}) - \frac{\mathbf{T}^T}{2} \mathbf{D}_\mathbf{t} f(\mathbf{t}) + \mathcal{O}(\|\mathbf{T}\|^2) \quad (2.185)$$

und zum anderen durch Wahl der unteren Vorzeichen und der Substitution $\mathbf{t} \Rightarrow \mathbf{t} + \mathbf{T}/2$

$$f(\mathbf{t} + \mathbf{T}/2) = f(\mathbf{t}) + \frac{\mathbf{T}^T}{2} \mathbf{D}_\mathbf{t} f(\mathbf{t}) + \mathcal{O}(\|\mathbf{T}\|^2). \quad (2.186)$$

Wir betrachten nun die Definition der Leistungswellen an einem idealen Kondensator gemäß Gleichung (2.172) mit dem Torwiderstand $R = \frac{T_0}{2C}$

$$\hat{b}(\mathbf{t}) = \frac{1}{2\sqrt{R}} [\hat{u}(\mathbf{t}) - R \hat{i}(\mathbf{t})] \quad , \quad \hat{a}(\mathbf{t}) = \frac{1}{2\sqrt{R}} [\hat{u}(\mathbf{t}) + R \hat{i}(\mathbf{t})] \quad (2.187)$$

Durch Einsetzen des Stromes $\hat{i} = CD\hat{u}$ erhalten wir

$$\hat{b}(\mathbf{t}) = \frac{1}{2\sqrt{R}} \left[\hat{u}(\mathbf{t}) - \frac{T_0}{2} D \hat{u}(\mathbf{t}) \right] \quad , \quad \hat{a}(\mathbf{t}) = \frac{1}{2\sqrt{R}} \left[\hat{u}(\mathbf{t}) + \frac{T_0}{2} D \hat{u}(\mathbf{t}) \right] \quad (2.188)$$

Mit $D = \boldsymbol{\alpha}^T \mathbf{D}_\mathbf{t}$ und $\mathbf{T} = \boldsymbol{\alpha} T_0$ ergibt sich dann

$$\hat{b}(\mathbf{t}) = \frac{1}{2\sqrt{R}} \left[\hat{u}(\mathbf{t}) - \frac{\mathbf{T}^T}{2} \mathbf{D}_\mathbf{t} \hat{u}(\mathbf{t}) \right] \quad , \quad \hat{a}(\mathbf{t}) = \frac{1}{2\sqrt{R}} \left[\hat{u}(\mathbf{t}) + \frac{\mathbf{T}^T}{2} \mathbf{D}_\mathbf{t} \hat{u}(\mathbf{t}) \right] \quad (2.189)$$

Durch Nutzung der Taylorreihenentwicklungen Gleichung (2.185) und Gleichung (2.186) erhalten wir die gewünschten approximativen Beziehungen

$$\hat{b}(\mathbf{t}) = \sqrt{\frac{C}{2T_0}} u(\mathbf{t} - \mathbf{T}/2) \approx \sqrt{\frac{C}{2T_0}} \hat{u}(\mathbf{t} - \mathbf{T}/2) \quad \text{und} \quad \hat{a}(\mathbf{t}) = \sqrt{\frac{C}{2T_0}} u(\mathbf{t} + \mathbf{T}/2) \approx \sqrt{\frac{C}{2T_0}} \hat{u}(\mathbf{t} + \mathbf{T}/2),$$

(2.190)

bzw. im diskreten ursprünglichen Koordinatensystem

$$\tilde{b}(\boldsymbol{\mu}) \approx \sqrt{\frac{C}{2T_0}} \tilde{u}(\boldsymbol{\mu} - \Delta\boldsymbol{\mu}/2) \quad \text{und} \quad \tilde{a}(\boldsymbol{\mu}) \approx \sqrt{\frac{C}{2T_0}} \tilde{u}(\boldsymbol{\mu} + \Delta\boldsymbol{\mu}/2). \quad (2.191)$$

Die Ergebnisse können im linearen, konstanten Fall auch durch Betrachtung des Frequenzbereichs gewonnen werden. Dazu wenden wir die Approximationen

$$z^{1/2} = e^{pT_0/2} = \sum_{\nu=0}^{\infty} \frac{1}{\nu!} \left(\frac{pT_0}{2} \right)^{\nu} \approx 1 + pT_0/2 \iff pT_0/2 \approx z^{1/2} - 1 \quad (2.192)$$

$$z^{-1/2} = e^{-pT_0/2} = \sum_{\nu=0}^{\infty} \frac{1}{\nu!} \left(-\frac{pT_0}{2} \right)^{\nu} \approx 1 - pT_0/2 \iff pT_0/2 \approx 1 - z^{-1/2}$$

auf $CpU = I$ an. Wir erhalten näherungsweise mit dem Torwiderstand $R = \frac{T_0}{2C}$

$$[z^{1/2} - 1]U = RI \quad \text{und} \quad [1 - z^{-1/2}]U = RI. \quad (2.193)$$

Hieraus gewinnen wir durch Verwendung von Leistungswellen

$$\frac{1}{2\sqrt{R}} z^{-1/2}U = A \quad \text{und} \quad \frac{1}{2\sqrt{R}} z^{1/2}U = B. \quad (2.194)$$

Durch Übergang zu den zeitabhängigen Größen ergibt sich Gleichung (2.191).

2.9 Zweckmäßige Beschreibung des Wellendigitalfilters

Das Ziel dieses Unterkapitels ist es, ein Wellendigitalmodell, welches aus den einzelnen WDF-Bauelementen des Unterkapitels 2.8 und deren Verbindungen besteht, in Form eines Gleichungssystems darzustellen. Das Gleichungssystem soll dabei so geartet sein, dass nach einer bloßen Festlegung der Berechnungsreihenfolge eine explizite Berechnung erfolgen kann. Durch die Festlegung der Berechnungsreihenfolge erhalten wir einen Algorithmus, der uns später als Grundlage für die formale Spezifikation und die Bestimmung eines Codes zur Simulation dienen wird.

Die zweckmäßige Beschreibung eines mehrdimensionalen Wellendigitalfilters geschieht in Anlehnung an die Darstellung im Eindimensionalen [Rumm98]. Diese Beschreibung ist besonders geeignet für die Führung eines formalen Korrektheitsbeweises. Wir nehmen zunächst an, alle verwendeten Quellen sind reflexionsfrei, d. h. die Torwiderstände sind gleich den Innenwiderständen der Quellen. Der Vektor \mathbf{b}_q beinhaltet die skalaren Quellenwellen der Strom- und Spannungsquellen und hat die Länge n_q

$$\mathbf{b}_q = [b_{q1}, b_{q2}, \dots, b_{qn_q}]^T. \quad (2.195)$$

Die einfallenden Wellen der Quellen lauten

$$\mathbf{a}_q = [a_{q1}, a_{q2}, \dots, a_{qn_q}]^T. \quad (2.196)$$

Die Wellengrößen der dynamikfreien Elemente sind die Koordinaten der Vektoren

$$\mathbf{b}_e = [b_{e1}, b_{e2}, \dots, b_{en_e}]^T \quad \text{und} \quad \mathbf{a}_e = [a_{e1}, a_{e2}, \dots, a_{en_e}]^T. \quad (2.197)$$

Die Wellengrößen \mathbf{a}_e und \mathbf{b}_e sind über die Streumatrix \mathbf{S} der dynamikfreien Elemente miteinander verknüpft, d. h.

$$\mathbf{b}_e = \mathbf{S} \mathbf{a}_e. \quad (2.198)$$

Hierin ist die Blockdiagonalmatrix \mathbf{S} die direkte Summe der N_e Streumatrizen \mathbf{S}_ν eines jeden dynamikfreien Elementes ν

$$\mathbf{S} = \text{diag}(\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_\nu, \dots, \mathbf{S}_{N_e}). \quad (2.199)$$

Die Anzahl der Tore des dynamikfreien Elementes ν ist $n_{e\nu}$. Die Anzahl aller Tore aller dynamikfreien Elemente ist $n_e = \sum_{\nu=1}^{N_e} n_{e\nu}$. Die Verzögerer bzgl. t_κ werden durch

$$\mathbf{b}_v^\kappa([\nu_1, \dots, \nu_{\kappa-1}, \nu_\kappa + 1, \nu_{\kappa+1}, \dots, \nu_{k'}]^\text{T}) = \mathbf{a}_v^\kappa(\boldsymbol{\nu}) \quad (2.200)$$

beschrieben, wobei alle Ausgangswellen der Verzögerer bzgl. t_κ in dem Vektor

$$\mathbf{b}_v^\kappa = [b_{v1}^\kappa, b_{v2}^\kappa, \dots, b_{vn_v^\kappa}^\kappa]^\text{T} \quad (2.201)$$

und alle Eingangswellen der Verzögerer bzgl. t_κ in dem Vektor

$$\mathbf{a}_v^\kappa = [a_{v1}^\kappa, a_{v2}^\kappa, \dots, a_{vn_v^\kappa}^\kappa]^\text{T} \quad (2.202)$$

zusammengefasst sind. Die Anzahl Verzögerer bzgl. t_κ bezeichnen wir mit n_v^κ . Die Anzahl aller Verzögerer ist $n_v = \sum_{\kappa=1}^{k'} n_v^\kappa$. Mit Hilfe des Standardeinheitsvektors erhalten wir die kompakte Darstellung

$$\mathbf{b}_v^\kappa(\boldsymbol{\nu} + \mathbf{e}_\kappa) = \mathbf{a}_v^\kappa(\boldsymbol{\nu}) \quad , \quad \kappa = 1, 2, \dots, k'. \quad (2.203)$$

Weiterhin fassen wir die Vektoren der Verzögerer zu den Vektoren

$$\mathbf{a}_v = \begin{bmatrix} \mathbf{a}_v^1 \\ \mathbf{a}_v^2 \\ \vdots \\ \mathbf{a}_v^{k'} \end{bmatrix} \quad \text{und} \quad \mathbf{b}_v = \begin{bmatrix} \mathbf{b}_v^1 \\ \mathbf{b}_v^2 \\ \vdots \\ \mathbf{b}_v^{k'} \end{bmatrix} \quad (2.204)$$

zusammen.

Um das gesamte Wellendigitalfilter zu beschreiben, erweist es sich als günstig die Übervektoren

$$\mathbf{a} = \begin{bmatrix} \mathbf{a}_q \\ \mathbf{a}_v \\ \mathbf{a}_e \end{bmatrix} \quad \text{und} \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ \mathbf{b}_e \end{bmatrix} \quad (2.205)$$

zu bilden. Die Übervektoren \mathbf{a} und \mathbf{b} besitzen die Länge $n_g = n_q + n_v + n_e$.

Die Verschaltung der einzelnen Bauelemente wird durch die Permutationsmatrix \mathbf{P} beschrieben. Die einzelnen Wellengrößen sind miteinander über

$$\mathbf{a} = \mathbf{P} \mathbf{b} \quad (2.206)$$

verknüpft. Ausführlich sind die Verbindungen der Bauelemente durch

$$\begin{bmatrix} \mathbf{a}_q \\ \mathbf{a}_v^1 \\ \vdots \\ \mathbf{a}_v^{k'} \\ \mathbf{a}_e \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{qq} & \mathbf{P}_{qv^1} & \cdots & \mathbf{P}_{qv^{k'}} & \mathbf{P}_{qe} \\ \mathbf{P}_{v^1q} & \mathbf{P}_{v^1v^1} & \cdots & \mathbf{P}_{v^1v^{k'}} & \mathbf{P}_{v^1e} \\ \vdots & & \ddots & & \vdots \\ \mathbf{P}_{v^{k'}q} & \mathbf{P}_{v^{k'}v^1} & \cdots & \mathbf{P}_{v^{k'}v^{k'}} & \mathbf{P}_{v^{k'}e} \\ \mathbf{P}_{eq} & \mathbf{P}_{ev^1} & \cdots & \mathbf{P}_{ev^{k'}} & \mathbf{P}_{ee} \end{bmatrix} \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_{v^1} \\ \vdots \\ \mathbf{b}_{v^{k'}} \\ \mathbf{b}_e \end{bmatrix} \quad (2.207)$$

definiert. Durch Zusammenfassen von Matrizen zu Übermatrizen erhalten wir

$$\begin{bmatrix} \mathbf{a}_q \\ \mathbf{a}_v \\ \mathbf{a}_e \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{qq} & \mathbf{P}_{qv} & \mathbf{P}_{qe} \\ \mathbf{P}_{vq} & \mathbf{P}_{vv} & \mathbf{P}_{ve} \\ \mathbf{P}_{eq} & \mathbf{P}_{ev} & \mathbf{P}_{ee} \end{bmatrix} \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ \mathbf{b}_e \end{bmatrix} \quad (2.208)$$

mit den entsprechend zu bildenden Übermatrizen $\mathbf{P}_{qv}, \mathbf{P}_{vq}, \mathbf{P}_{ev}, \mathbf{P}_{ve}$ und \mathbf{P}_{vv} . Die Matrix \mathbf{P} ist eine quadratische Matrix der Dimension $n_g = n_q + n_v^1 + \dots + n_v^{k'} + n_e$.

Permutationsmatrizen haben die Eigenschaft orthogonal zu sein. Hier ist \mathbf{P} sogar symmetrisch, da die Verbindung der Bauelemente torweise erfolgt. Es gilt somit $\mathbf{P}^T \mathbf{P} = \mathbf{1}_{n_g}$ und $\mathbf{P}^T = \mathbf{P}$.

Das Wellendigitalkonzept geht von einer Referenzschaltung aus. Dies wollen wir beibehalten, d. h. jedes Bauelement der Referenzschaltung ist durch das Verbindungsnetz mit einem anderen verbunden. Für das digitale Verbindungsnetz bedeutet dies, dass mindestens ein WDF-Bauelement einer Verbindung ein dynamikfreies WDF-Bauelement ist, d. h. $\mathbf{P}_{vv} = \mathbf{0}, \mathbf{P}_{vq} = \mathbf{0}, \mathbf{P}_{qv} = \mathbf{0}, \mathbf{P}_{qq} = \mathbf{0}$. Dies gilt nach Definition für die weitere Arbeit. Die Permutationsmatrix lautet dann

$$\mathbf{P} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{P}_{qe} \\ \mathbf{0} & \mathbf{0} & \mathbf{P}_{ve} \\ \mathbf{P}_{eq} & \mathbf{P}_{ev} & \mathbf{P}_{ee} \end{bmatrix} = \begin{bmatrix} \mathbf{P}_q \\ \mathbf{P}_v \\ \mathbf{P}_e \end{bmatrix}. \quad (2.209)$$

Im Bild 2.29 ist die Beschreibung des MDWDFs graphisch dargestellt. Eine Detailansicht in Form eines Signalfussdiagramms zeigt Bild 2.30.

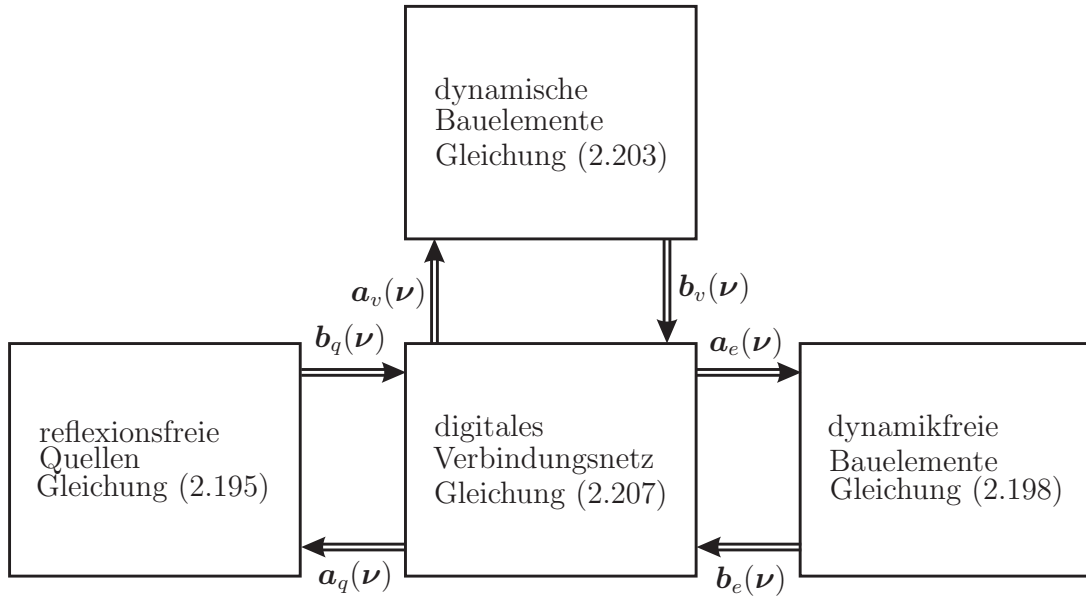


Bild 2.29: Algorithmus gültig in \mathcal{G}

Die Ausgangswellen der Quellen $\mathbf{b}_q(\nu)$ und der Verzögerer $\mathbf{b}_v^k(\nu)$ sind zu einem Zeitpunkt für alle Gitterpunkte innerhalb \mathcal{G} bekannt. Die nötige Berechnung von $\mathbf{b}_e(\nu)$ erfolgt durch Einsetzen der letzten Zeile von Gleichung (2.207) in Gleichung (2.198)

$$\mathbf{b}_e(\nu) = \mathbf{S} [\mathbf{P}_{eq} \mathbf{b}_q(\nu) + \mathbf{P}_{ev^1} \mathbf{b}_{v^1}(\nu) + \dots + \mathbf{P}_{ev^{k'}} \mathbf{b}_{v^{k'}}(\nu) + \mathbf{P}_{ee} \mathbf{b}_e(\nu)]. \quad (2.210)$$

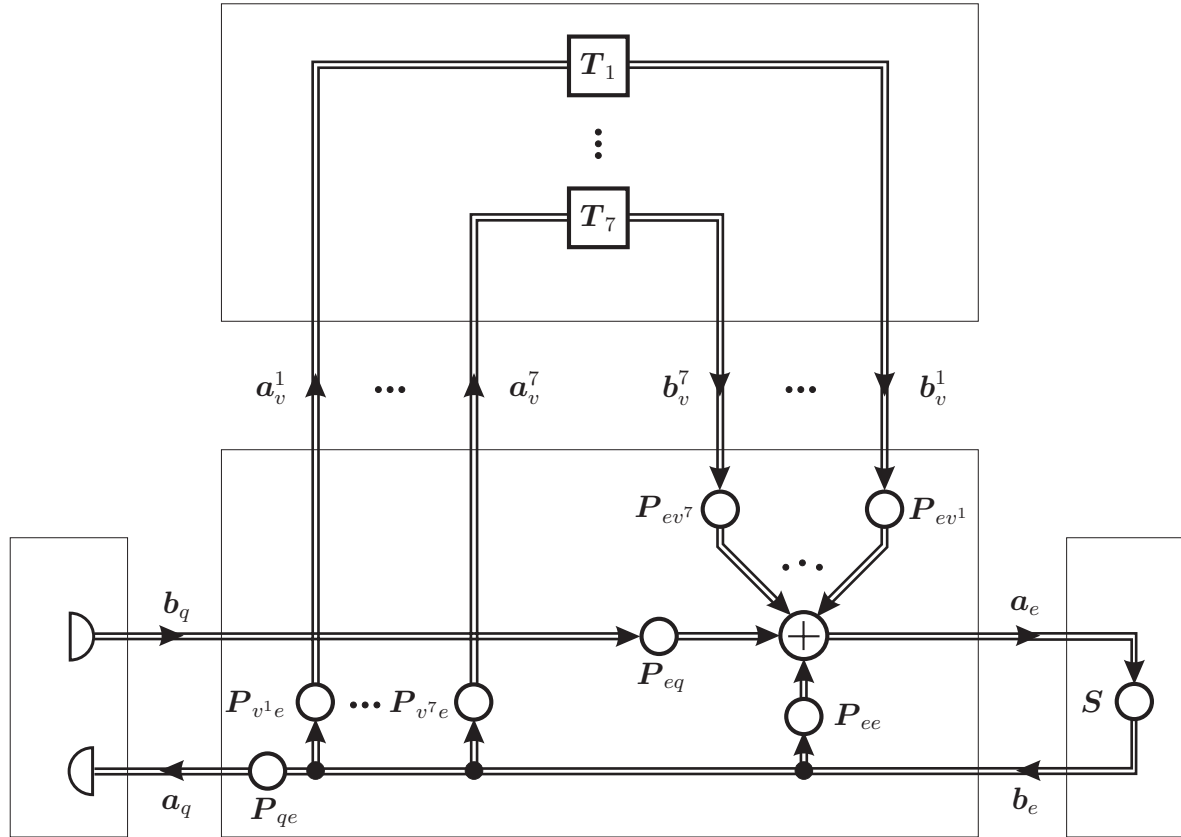


Bild 2.30: Detailansicht vom Bild 2.29

Die Matrix $\mathbf{S} \mathbf{P}_{ee}$ muss nun durch symmetrische Zeilen- und Spaltenpermutationen auf untere Dreiecksgestalt gebracht werden, bei der zudem noch die Hauptdiagonalelemente null sind. Eine untere Dreiecksmatrix, deren Hauptdiagonalelemente null sind, werden wir im weiteren Verlauf der Arbeit als strikte untere Dreiecksmatrix bezeichnen. Existiert eine Permutationsmatrix, die die Matrix $\mathbf{S} \mathbf{P}_{ee}$ auf strikte untere Dreiecksgestalt transformiert, so liegt keine verzögerungsfreie gerichtete Schleife in dem Wellendigitalfilter vor und die Berechenbarkeit ist gegeben.

Es stellt sich also die Frage, ob wenigstens eine Permutationsmatrix \mathbf{P}_b existiert, sodass

$$\mathbf{P}_b \mathbf{S} \mathbf{P}_{ee} \mathbf{P}_b^T = \begin{bmatrix} 0 & \cdots & \cdots & 0 \\ \bullet & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \bullet & \cdots & \bullet & 0 \end{bmatrix} \quad (2.211)$$

gilt. Wir nehmen im Folgenden an, dass ein korrekt spezifiziertes MDWDF vorliegt. Wir führen die permutierten Vektoren

$$\mathbf{b}'_e = \mathbf{P}_b \mathbf{b}_e \quad \text{und} \quad \mathbf{b}' = [\mathbf{b}'_q \quad \mathbf{b}'_v \quad \mathbf{b}'_e]^T \quad (2.212)$$

ein und erhalten die in ihrer Reihenfolge getauschten Wellengrößen durch

$$\mathbf{b}'_e = \mathbf{P}_b \mathbf{b}_e = \mathbf{P}_b \mathbf{S} \mathbf{P}_b^T [\mathbf{P}_b \mathbf{P}_{eq} \quad \mathbf{P}_b \mathbf{P}_{ev^1} \quad \cdots \quad \mathbf{P}_b \mathbf{P}_{ev^{k'}} \quad \mathbf{P}_b \mathbf{P}_{ee} \mathbf{P}_b^T] \mathbf{b}' = \mathbf{L}' \mathbf{b}'. \quad (2.213)$$

Das weitere Einführen von permutierten Matrizen

$$\mathbf{S}' = \mathbf{P}_b \mathbf{S} \mathbf{P}_b^T, \quad \mathbf{P}'_{eq} = \mathbf{P}_b \mathbf{P}_{eq}, \quad \mathbf{P}'_{ev} = \mathbf{P}_b \mathbf{P}_{ev}, \quad \mathbf{P}'_{ee} = \mathbf{P}_b \mathbf{P}_{ee} \mathbf{P}_b^T \quad (2.214)$$

führt uns zunächst zu

$$\mathbf{S}'\mathbf{P}'_{eq} = \mathbf{P}_b\mathbf{S}\mathbf{P}_{eq} \quad , \quad \mathbf{S}'\mathbf{P}'_{ev} = \mathbf{P}_b\mathbf{S}\mathbf{P}_{ev} \quad , \quad \mathbf{S}'\mathbf{P}'_{ee} = \mathbf{P}_b\mathbf{S}\mathbf{P}_{ee}\mathbf{P}_b^T, \quad (2.215)$$

so dass wir den permutierten Vektor \mathbf{b}'_e durch

$$\begin{aligned} \mathbf{b}'_e &= \mathbf{S}'\mathbf{P}'_{eq}\mathbf{b}_q + \mathbf{S}'\mathbf{P}'_{ev}\mathbf{b}_v + \mathbf{S}'\mathbf{P}'_{ee}\mathbf{b}'_e \\ &= \underbrace{[\mathbf{S}'\mathbf{P}'_{eq} \quad \mathbf{S}'\mathbf{P}'_{ev} \quad \mathbf{S}'\mathbf{P}'_{ee}]}_{\mathbf{L}'=} \mathbf{b}' \end{aligned} \quad (2.216)$$

darstellen können. Dieses Gleichungssystem kann aufgrund der Struktur der Matrix \mathbf{L}' durch Vorwärts-substitution berechnet werden. Der neue Zustandsvektor ergibt sich zu

$$\mathbf{b}_v(\boldsymbol{\nu} + \mathbf{e}_\nu) = \mathbf{P}_{ve}\mathbf{b}_e(\boldsymbol{\nu}) = \mathbf{P}_{ve}\mathbf{P}_b^T = \mathbf{P}'_{ve}\mathbf{b}'_e(\boldsymbol{\nu}). \quad (2.217)$$

Definieren wir abschliessend

$$\mathbf{L}'_q = \mathbf{S}'\mathbf{P}'_{eq}, \quad \mathbf{L}'_v = \mathbf{S}'\mathbf{P}'_{ev}, \quad \mathbf{L}'_e = \mathbf{S}'\mathbf{P}'_{ee}, \quad (2.218)$$

so können letztendlich die zu implementierenden Gleichungen zu

$$\begin{aligned} \mathbf{b}'_e(\boldsymbol{\nu}) &= \mathbf{L}'_q\mathbf{b}_q(\boldsymbol{\nu}) + \mathbf{L}'_v\mathbf{b}_v(\boldsymbol{\nu}) + \mathbf{L}'_e\mathbf{b}'_e(\boldsymbol{\nu}) \\ \mathbf{b}_v^\kappa(\boldsymbol{\nu} + \mathbf{e}_\kappa) &= \mathbf{P}_{v^\kappa e}\mathbf{b}_e(\boldsymbol{\nu}) \quad \forall \kappa = 1, 2, \dots, k' \end{aligned} \quad (2.219)$$

notiert werden. Dabei hat die Berechnung der Vektoren $\mathbf{b}'_e(\boldsymbol{\nu})$ und $\mathbf{b}_v(\boldsymbol{\nu})$ zu allen Abtastpunkten einer Abtastschicht zu erfolgen, wobei allerdings die zweite Gleichung nur für eingeschränkte $\boldsymbol{\nu}$ gültig ist. Auf diese Problematik wird im Zusammenhang mit den Randwertfragen im Kapitel 2.10 eingegangen.

2.10 Randbehandlung

Die Verwendung von Digitalfiltern zur numerische Integration erfordert eine Berücksichtigung der Ränder. Ein Signalflussdiagramm (nach Bild 2.29) ist eine äquivalente Darstellung eines Differenzengleichungssystems (Gleichung (2.198), Gleichung (2.203) und Gleichung (2.207)). Die Differenzengleichung wird dabei aus der zu lösenden Differentialgleichung gewonnen. Der Algorithmus eines mehrdimensionalen Wellendigitalfilters zur numerischen Integration ist nun nicht allein durch die Differenzengleichung bestimmt, sondern erfordert zusätzlich eine Beschreibung der Abmessungen des Berechnungsgebietes und die Behandlung des Randes.

Bevor wir die Diskussion fortsetzen, wollen wir den Begriff Randbehandlung präzisieren. Dazu betrachten wir eine beliebige Fläche im Raum und stellen uns auf eine Seite der Fläche. Das Verhalten der Feldgrößen auf dieser Seite der Fläche wird vom gesamten Raum hinter der Fläche beeinflusst. Oft machen es bestimmte Materialübergänge möglich, das Verhalten der Feldgrößen hinter der Fläche durch ein Modell zu approximieren, welches keine Abhängigkeit in Richtung des Normalenvektors der Grenzfläche besitzt. Wir definieren eine Fläche, auf der eine Approximation des Verhaltens der Feldgrößen innerhalb Raumes hinter dieser Fläche durchgeführt wird, als eine Randfläche. Wir unterscheiden bei denen durch die Approximation entstandenen Modelle i.W. zwei Fälle. Zum einen betrachten wir Modelle, die ausschließlich algebraische Beziehungen zwischen den Feldgrößen auf dem Rand liefern und zum anderen modellieren wir das Randverhalten durch differentielle Beziehungen zwischen den Feldgrößen.

Im Folgenden soll der Grund für die besondere Behandlung des Randes eines Berechnungsgebietes bei der numerischen Integration mit mehrdimensionalen Wellendigitalfiltern erläutert werden. Dazu betrachten wir noch einmal das Vorgehen bei der Simulation eines mehrdimensionalen Wellendigitalfilters. Zunächst werden gemäß Gleichung (2.219) alle fehlenden Wellengrößen der dynamikfreien Elemente des Berechnungsgebietes ermittelt. Danach werden die ausfallenden Wellen der dynamischen Elemente der nächsten Abtastschicht, d. h.

$$\mathbf{b}_v^\kappa(\boldsymbol{\nu} + \mathbf{e}_\kappa) = \mathbf{a}_v^\kappa(\boldsymbol{\nu}) \quad \text{für } \kappa = 1, \dots, k' \quad (2.220)$$

ermittelt. Wir gehen dazu von Gleichung (2.203) aus, wobei wir zuvor $\mathbf{a}_v^\kappa(\boldsymbol{\nu})$ aus allen bereits vorliegenden Vektoren $\mathbf{b}_q(\boldsymbol{\nu})$, $\mathbf{b}_v(\boldsymbol{\nu})$ und $\mathbf{b}_e(\boldsymbol{\nu})$ gemäß

$$\mathbf{a}_v^\kappa(\boldsymbol{\nu}) = \mathbf{P}_{v^\kappa e} \mathbf{b}_e(\boldsymbol{\nu}) \quad (2.221)$$

berechnet haben. Aus dieser Gleichung lässt sich allerdings nicht ohne Weiteres eine Aussage über die Randproblematik ableiten, da das Berechnungsgebiet im neuen Koordinatensystem nicht mehr die Form eines Quaders hat. Wir ziehen daher für die weitere Diskussion eine Beschreibung der dynamischen Elemente im ursprünglichen Koordinatensystem vor, d. h.

$$\tilde{\mathbf{b}}_v^\kappa(\boldsymbol{\mu}) = \mathbf{b}_v^\kappa(\boldsymbol{\nu}) = \mathbf{a}_v^\kappa(\boldsymbol{\nu} - \mathbf{e}_\kappa) = \tilde{\mathbf{a}}_v^\kappa(\boldsymbol{\mu} - \mathbf{h}_\kappa) . \quad (2.222)$$

Über diese Beziehung ist nun der Vektor $\tilde{\mathbf{b}}_v^\kappa(\boldsymbol{\mu})$ für alle örtlichen Gitterpunkte $\boldsymbol{\mu}' \in \mathcal{G}$ zu bestimmen, siehe Bild 2.31. Dafür sind die Wellengrößen $\tilde{\mathbf{a}}_v^\kappa$ zu den Gitterpunkten $(\boldsymbol{\mu}' - \mathbf{h}_\kappa')$ notwendig. Die Wellengrößen liegen aber nur für die Gitterpunkte innerhalb des Berechnungsgebietes vor. Wir haben daher bei der Bestimmung von $\tilde{\mathbf{b}}_v^\kappa$ zwischen den Punkten $(\boldsymbol{\mu}' - \mathbf{h}_\kappa) \in \mathcal{G}$ und den Punkten $(\boldsymbol{\mu}' - \mathbf{h}_\kappa) \in \mathcal{G}_0$ zu unterscheiden. (Hier wurde selbstverständlich wieder angenommen, dass $\boldsymbol{\mu}'$ selber im Berechnungsgebiet \mathcal{G} liegt.) Ein Algorithmus zur Bestimmung der Wellengrößen $\tilde{\mathbf{b}}_v^\kappa(\boldsymbol{\mu})$ ist somit in zwei Teilalgorithmen zu zerlegen, die Normalbehandlung und die Randbehandlung.

Es stellt sich daher die Frage, wie wir die notwendigen Wellengrößen am Rand bestimmen. Wir verwenden für die Randbehandlung einen Vorschlag von Fettweis und Seraji, [FS99], [Sera00]. In [Sera00] ist die Behandlung von Rändern im $(1+1)$ D-Fall am Beispiel der Telegraphengleichung dargestellt. Es wird darauf hingewiesen, dass die vorgestellte Methode auf Problemstellungen mit beliebiger Dimension erweiterbar ist. Wir werden im Folgenden den für uns wesentlichen Teil aus dieser Arbeit rekapitulieren und das beschriebene Verfahren verallgemeinern. Hierbei beschränken wir uns auf die Transformationsmatrix und Abtastmatrix gemäß Gleichung (2.49). Untersuchungen haben ergeben, dass dies die einfachste Möglichkeit der Randbehandlung darstellt [Voll02]. Aus dieser Festlegung resultiert zudem der Vorteil einer sich ergebenden festen Codestruktur, wie sich im Kapitel 6 zeigen wird.

Die generelle Schwierigkeit der Randbehandlung resultiert daraus, dass der Rand in den Feldgrößen beschrieben wird, aber das MDWDF als Zustandsgrößen die Wellengrößen verwendet. Zwar lassen sich die Feldgrößen aus den Wellengrößen bestimmen, die umgekehrte Richtung ist allerdings nur in Sonderfällen möglich. Das Grundprinzip zur Vermeidung dieser Schwierigkeit beruht auf der zuerst in [Frie95] dargestellten Tatsache, dass sich die Wellengrößen eines idealen Kondensators bzw. einer idealen Spule durch die Spannung bzw. den Strom zu einem um den halben Verschiebe-Vektor versetzten Punkt mit einer lokalen Konsistenzordnung von eins approximieren lassen, siehe Gleichung (2.190). Da die globale Konsistenzordnung der Trapezregel auch eins beträgt [HNW93] fügt sich diese Randbehandlung in den Rest des Verfahrens ein. Diese Methode der Randbehandlung erfordert einen Abstand von einem halben Verschiebe-Vektor zwischen dem randnächstem Punkt und der Grenzebene, so wie es im Kapitel 2.5 beschrieben wurde.

In [Sera00] wird ein weiteres, ähnliches Verfahren zur Randbehandlung beschrieben. Dabei liegen die Randpunkte auf der Grenzebene. Um nun die in Gleichung (2.190) beschriebene Approximation nutzen

zu können, wird die Abtastrate am Rand geändert. Dies hat allerdings zur Folge, dass die Torwiderstände der ein- und ausfallenden Wellen unterschiedlich sind. Für ein Werkzeug, welches eine automatische Codegenerierung ermöglichen soll, erscheint dieses Verfahren zu aufwendig und wird deshalb in dieser Arbeit nicht verwendet.

Um nun das Verfahren zu beschreiben, gehen wir der Einfachheit halber, ohne Einschränkung der Allgemeinheit, davon aus, dass die behandelte Referenzschaltung keine idealen Spulen, sondern nur ideale Kondensatoren als reaktive Bauelemente besitzt. Weiterhin ist festzustellen, dass sich bei der Randbehandlung nur diejenigen idealen Kondensatoren gegenseitig beeinflussen, deren zugehörige Verschiebe-Vektoren einen gemeinsamen Punkt auf der Grenzebene besitzen. Die Verschiebe-Vektoren mit einem gemeinsamen Punkt auf der Grenzebene müssen wir zusätzlich unterscheiden und zwar zwischen Verschiebe-Vektoren, die aus dem örtlichen⁵ Berechnungsgebiet hinaus gerichtet sind, denen, die in das Berechnungsgebiet herein zeigen und denen, die orthogonal zum Normalenvektor sind. Dazu betrachten wir nun den Normalenvektor \mathbf{n} , der senkrecht auf der Grenzfläche steht und aus dem Berechnungsgebiet hinaus zeigt. Im (3+1)D-Fall haben wir, falls \mathcal{G} ein Quader ist, die 6 verschiedenen Normalenvektoren $\mathbf{n}_1 = -\mathbf{n}_4 = \mathbf{e}_1$, $\mathbf{n}_2 = -\mathbf{n}_5 = \mathbf{e}_2$, $\mathbf{n}_3 = -\mathbf{n}_6 = \mathbf{e}_3$. Ein Orts- Verschiebe-Vektor \mathbf{h}'_κ zeigt dann aus dem Berechnungsgebiet hinaus (herein), wenn das Skalarprodukt $\mathbf{n}^T \mathbf{h}'_\kappa$ positiv (negativ) ist. Ist das Skalarprodukt $\mathbf{n}^T \mathbf{h}'_\kappa$ gleich null, so ist zum einen für Zustandsgrößen mit dem Argument $\boldsymbol{\mu} + \mathbf{h}_\kappa$ keine Randbehandlung erforderlich und zum anderen liefern diese Zustandsgrößen keinen Beitrag zur Berechnung der anderen Randwerte. Die Spannungen derjenigen idealen Kondensatoren, deren zugehöriger Verschiebe-Vektor aus dem Berechnungsgebiet hinaus zeigt, sind die bekannten Feldgrößen eines Modells zur Randbehandlung. Die Spannungen derjenigen idealen Kondensatoren, deren zugehöriger Verschiebe-Vektor in das Berechnungsgebiet herein zeigt, müssen über das den Rand beschreibende Modell bestimmt werden.

Als Beispiel betrachten wir erneut ein Berechnungsgebiet in Form eines Quaders und nehmen eine Randfläche $x_1 = x_{1\max}$ an, siehe Bild 2.31. Wir betrachten zwei vektorielle ideale Kondensatoren deren diskrete Approximationen die Verschiebe-Vektoren \mathbf{h}_1 und \mathbf{h}_4 besitzen. Sie sollen den gemeinsamen Gitterpunkt $\boldsymbol{\mu} + \mathbf{h}_1/2$ haben. Der Verschiebe-Vektor \mathbf{h}_1 zeigt wegen $\mathbf{n}_1^T \mathbf{h}'_1 = 1 > 0$ aus dem Berechnungsgebiet hinaus. Der Verschiebe-Vektor \mathbf{h}_4 zeigt in das Berechnungsgebiet herein.

Als Basis für das Aufstellen der benötigten Gleichungen dient Gleichung (2.191) in vektorieller Form mit $\mathbf{C}^\kappa = \text{diag}(C_1^\kappa, \dots, C_n^\kappa)$

$$\tilde{\mathbf{b}}^\kappa(\boldsymbol{\mu}) = [\mathbf{C}^\kappa]^{1/2} \sqrt{1/2T_0} \tilde{\mathbf{u}}^\kappa(\boldsymbol{\mu} - \mathbf{h}_\kappa/2) \quad \text{und} \quad \tilde{\mathbf{a}}^\kappa(\boldsymbol{\mu}) = [\mathbf{C}^\kappa]^{1/2} \sqrt{1/2T_0} \tilde{\mathbf{u}}^\kappa(\boldsymbol{\mu} + \mathbf{h}_\kappa/2), \quad (2.223)$$

deren Argumente zu unserem Beispiel so verschoben werden müssen, dass die Spannungen als Argument den betrachteten Punkt $\boldsymbol{\mu} + \mathbf{h}_1/2$ haben. Der Vektor $\tilde{\mathbf{a}}^1(\boldsymbol{\mu})$ ist bekannt. Mit der Gleichung (2.191)

$$\tilde{\mathbf{a}}^1(\boldsymbol{\mu}) = [\mathbf{C}^1]^{1/2} \sqrt{1/2T_0} \tilde{\mathbf{u}}^1(\boldsymbol{\mu} + \mathbf{h}_1/2) \quad (2.224)$$

können wir den Vektor der Feldgrößen $\tilde{\mathbf{u}}^1(\boldsymbol{\mu} + \mathbf{h}_1/2)$ bestimmen. Die gesuchten Wellengrößen sind

$$\tilde{\mathbf{b}}^4(\boldsymbol{\mu} + \mathbf{h}_1/2 + \mathbf{h}_4/2) = [\mathbf{C}^4]^{1/2} \sqrt{1/2T_0} \tilde{\mathbf{u}}^4(\boldsymbol{\mu} + \mathbf{h}_1/2). \quad (2.225)$$

Die anderen vorhandenen idealen Kondensatoren benötigen keine Randbehandlung und tragen auch nichts zur Randbehandlung bei, da ihre Verschiebe-Vektoren \mathbf{n}_2 , \mathbf{n}_3 , \mathbf{n}_5 und \mathbf{n}_6 orthogonal zum Normalenvektor \mathbf{n}_1 sind. Wir beachten $[\mathbf{h}_1 + \mathbf{h}_4]/2 = \mathbf{e}_4$.

Für den Randpunkt $\boldsymbol{\mu} + \mathbf{h}_1/2$ ist nun ein Abbild der physikalischen Realität anzugeben. Dieses Modell muss so geartet sein, dass nur die approximierten Größen $\tilde{\mathbf{u}}^1$ und $\tilde{\mathbf{u}}^4$ miteinander verknüpft sind, z. B. in der Form

$$\tilde{\mathbf{u}}^4 = \mathbf{M} \tilde{\mathbf{u}}^1. \quad (2.226)$$

⁵Im weiteren Verlauf dieses Unterkapitels ist immer das örtliche Berechnungsgebiet gemeint.

In diesem Fall haben wir

$$\tilde{\mathbf{b}}^4(\boldsymbol{\mu} + \mathbf{e}_4) = [\mathbf{C}^4]^{1/2} \mathbf{M} [\mathbf{C}^1]^{-1/2} \tilde{\mathbf{a}}^1(\boldsymbol{\mu}). \quad (2.227)$$

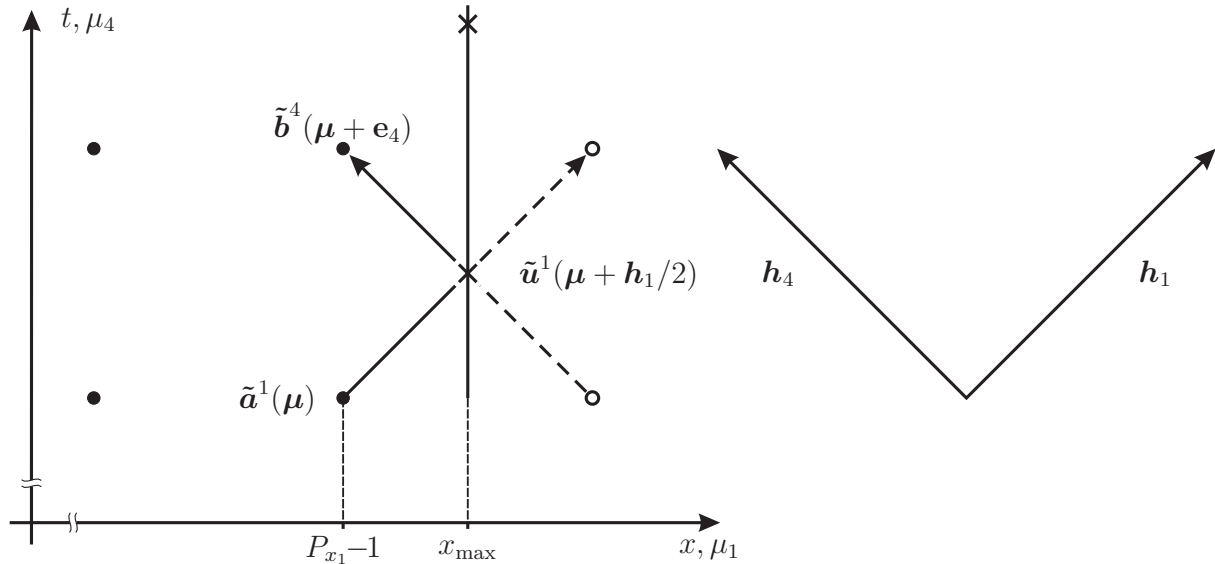


Bild 2.31: Randbehandlung

Im Falle eines nicht konstanten Modells am Rand ist das implizite Randlösungsmodell durch ein (ggf. nichtlineares) Wellendigitalfilter zu approximieren. Die Frage, die sich in jedem konkreten Einzelfall erneut stellt, ist, ob ein Randmodell gefunden werden kann, mit dem am Punkt $\boldsymbol{\mu} + \mathbf{h}_1/2$ aus den Spannungen $\tilde{\mathbf{u}}^1$ die gesuchten Spannungen $\tilde{\mathbf{u}}^4$ ermittelbar sind. Wir werden nicht weiter darauf eingehen, wie man ein derartiges Wellendigitalfilter des Randes findet, sondern es als gegeben annehmen.

Im weiteren Verlauf der Arbeit werden wir uns auf algebraische Randmodelle beschränken, da zu diesem Zeitpunkt noch keine ausreichenden Erkenntnisse über die Behandlung differentieller Ränder im (3+1)D Fall vorliegen. Im Falle algebraischer Randbedingungen können die Wellengrößen mittels Matrizen miteinander verknüpft werden

$$\tilde{\mathbf{b}}_v^\kappa(\boldsymbol{\mu} + \mathbf{e}_k) = \mathbf{R}^\kappa \tilde{\mathbf{a}}_v^{(\kappa+3)}(\boldsymbol{\mu}) \text{ für } \kappa = 1, 2, 3 \quad \text{und} \quad \tilde{\mathbf{b}}_v^\kappa(\boldsymbol{\mu} + \mathbf{e}_k) = \mathbf{R}^\kappa \tilde{\mathbf{a}}_v^{(\kappa-3)}(\boldsymbol{\mu}) \text{ für } \kappa = 4, 5, 6. \quad (2.228)$$

Die Matrizen \mathbf{R}^1 bis \mathbf{R}^6 stellen einen Teil der MDWDF-Spezifikation dar. Anhand der Spezifikation des Wellendigitalfilters für das Berechnungsgebiet kann eine Prüfung auf Vollständigkeit der Matrizen \mathbf{R}^1 bis \mathbf{R}^6 für die Randbehandlung erfolgen.

2.11 Von einem mehrdimensionalen Wellendigitalfilter zu einem eindimensionalen

Wir wollen nun auf den eigentlichen Ablauf der Berechnung eines Wellendigitalfilters eingehen. Wir betrachten dazu die notwendigen Rechenschritte zur Berechnung aller Wellengrößen zu einem festen Zeitpunkt μ_k . Zunächst haben wir die Quellenwerte $\mathbf{b}_q(\boldsymbol{\mu})$ für $\boldsymbol{\mu}' \in \mathcal{G}$ als bekannt vorausgesetzt. Die Verzögererwerte $\mathbf{b}_v(\boldsymbol{\nu})$ aller Abtastpunkte des Berechnungsgebietes der Zeitschicht μ_k können durch Auslesen der Speicherwerte aus den vorhandenen Werten der Zeitschicht $\mu_k - 1$ ermittelt werden. Liegt die Matrix \mathbf{L}' aus Gleichung (2.219) in der benötigten Form vor, so kann der Vektor $\mathbf{b}'_e(\boldsymbol{\nu})$ mittels Vorwärtssubstitution berechnet werden. Die Berechnung des Vektors $\mathbf{b}'_e(\boldsymbol{\nu})$ kann dabei für jeden Abtastpunkt der

Zeitschicht μ_k parallel erfolgen, da wie oben bereits erwähnt die benötigten Vektoren $\mathbf{b}_q(\boldsymbol{\nu})$ und $\mathbf{b}_v(\boldsymbol{\nu})$ der Abtastschicht μ_k vorliegen. Wir haben also eine Kombination aus paralleler und serieller Berechnung. Die Ausnutzung der möglichen Parallelität erfordert allerdings eine spezielle Hardware, die für jeden Abtastpunkt des Berechnungsgebietes eine Recheneinheit vorsieht, die Gleichung (2.219) implementiert. Bei der Realisierung eines mehrdimensionalen Wellendigitalfilters auf einem Universalrechner werden aber alle Rechenoperationen zumindest im Quellcode sequentiell programmiert. (Inwieweit letztendlich eine Berechnung dann parallel ausgeführt wird, hängt natürlich von dem Universalrechner selber und dem verwendeten Compiler ab.) Durch das Verzichten auf die massiv parallele Abarbeitung gehen aber keine weiteren Vorteile des Verfahrens der numerischen Integration nach dem Wellendigitalfilterprinzip verloren [Fett99].

Im Folgenden soll eine analytische Beschreibung für die sequentielle Berechnung eines mehrdimensionalen Wellendigitalfilters hergeleitet werden, deren Grundlagen auf [Kumm88] zurückgehen. Wir bilden dazu das mehrdimensionale System auf ein eindimensionales System ab. Dabei kommt uns die Tatsache zu Hilfe, dass wir ein endliches Berechnungsgebiet vorliegen haben. Wir bilden jeden Ortsvektor der P Abtastpunkte einer Abtastschicht auf eine ganze Zahl δ ab. Sollen dann alle Abtastpunkte eines Berechnungsgebietes berechnet werden, so kann dies dadurch geschehen, dass die neue Variable δ die Werte 0 bis $P - 1$ durchläuft und die Wellengrößen an diesen Punkten berechnet werden.

Die Verknüpfung von δ mit dem Vektor ${}_k\boldsymbol{\mu} = \boldsymbol{\mu}'$ erfolgt durch die eineindeutige nichtlineare Abbildung

$$\begin{aligned} \boldsymbol{\mu}' = \mathbf{g}(\delta) \quad : \quad \delta \in [0, P-1] \subset \mathbb{Z} &\longmapsto \boldsymbol{\mu}' \in \mathcal{G} \subset \mathbb{Z}^{k-1} \\ \delta = f(\boldsymbol{\mu}') \quad : \quad \boldsymbol{\mu}' \in \mathcal{G} \subset \mathbb{Z}^{k-1} &\longmapsto \delta \in [0, P-1] \subset \mathbb{Z} \end{aligned} \quad (2.229)$$

die durch

$$\begin{aligned} \mu_\kappa &= \left\lceil \frac{\delta - \sum_{n=\kappa+1}^{k-1} \mu_n c_n}{c_\kappa} \right\rceil, \quad \kappa = 1, 2, \dots, k-1 \\ \delta &= \sum_{n=1}^{k-1} \mu_n c_n = \boldsymbol{\mu}'^T \mathbf{c} \\ c_n &= \prod_{d=1}^{n-1} P_{x_d} \end{aligned} \quad (2.230)$$

definiert ist. Bei der Ermittlung des Vektors μ_κ ist man auf die μ_η , $\eta > \kappa$ angewiesen. Bei einer konkreten Berechnung beginnt man dann bei der letzten Koordinate und arbeitet sich zur ersten Koordinate durch, was einer Art Rückwärtssubstitution bei linearen Gleichungssystemen entspricht. Zur Verdeutlichung dient Bild 2.32.

Zum Beweis der Eineindeutigkeit der Abbildung benutzen wir zunächst

Satz 1

$$\sum_{n=1}^N (P_{x_n} - 1) \prod_{d=1}^{n-1} P_{x_d} = -1 + \prod_{d=1}^N P_{x_d}. \quad (2.231)$$

Beweis:

$$\sum_{n=1}^N (P_{x_n} - 1) \prod_{d=1}^{n-1} P_{x_d} = \sum_{n=1}^N P_{x_n} \prod_{d=1}^{n-1} P_{x_d} - \sum_{n=1}^N \prod_{d=1}^{n-1} P_{x_d} = \sum_{n=1}^N \prod_{d=1}^n P_{x_d} - \sum_{n=0}^{N-1} \prod_{d=1}^n P_{x_d} = \prod_{d=1}^N P_{x_d} - \prod_{d=1}^0 P_{x_d}$$

(2.232)

Wir benötigen außerdem den

Satz 2

$$\sum_{n=1}^N \mu_n c_n < c_{N+1}. \quad (2.233)$$

Beweis:

Offenbar nimmt die Summe ihren maximalen Wert an, wenn in allen Richtungen der letzte Abtastpunkt innerhalb \mathcal{G} eingesetzt wird, d. h. $\mu_n = P_n - 1, n = 1, 2, \dots, N$ gilt. Wir können somit die Beweisaufgabe auf das Zeigen der Richtigkeit von

$$\sum_{n=1}^N (P_{x_n} - 1) c_n < c_{N+1} \quad (2.234)$$

zurückführen. Wir setzen die Definition der c_n aus Gleichung (2.230) ein und erhalten

$$\sum_{n=1}^N (P_{x_n} - 1) \prod_{d=1}^{n-1} P_{x_d} < \prod_{d=1}^N P_{x_d}. \quad (2.235)$$

Durch Verwendung von Satz 1 erkennt man unmittelbar die Korrektheit von Satz 2.

Wir kommen nun zur eigentlichen Aufgabe zurück, die Eineindeutigkeit der Abbildung zu zeigen. Um $\mu' = g(f(\mu'))$ zu belegen, setzen wir

$$\delta = \sum_{n=1}^{k-1} \mu_n c_n \quad (2.236)$$

in

$$\mu_\kappa = \left\lceil \frac{\delta - \sum_{n=\kappa+1}^{k-1} \mu_n c_n}{c_\kappa} \right\rceil \quad (2.237)$$

ein und bekommen

$$\mu_\kappa = \left\lceil \frac{\sum_{n=1}^{k-1} \mu_n c_n - \sum_{n=\kappa+1}^{k-1} \mu_n c_n}{c_\kappa} \right\rceil = \left\lceil \frac{\sum_{n=1}^{\kappa} \mu_n c_n}{c_\kappa} \right\rceil. \quad (2.238)$$

Durch Aufteilen der Summe und Verwendung von Satz 2 finden wir

$$\mu_\kappa = \left\lceil \frac{\sum_{n=1}^{\kappa-1} \mu_n c_n + \mu_\kappa c_\kappa}{c_\kappa} \right\rceil = \left\lceil \underbrace{\frac{\sum_{n=1}^{\kappa-1} \mu_n c_n}{c_\kappa}}_{<1} + \mu_\kappa \right\rceil = \mu_\kappa. \quad (2.239)$$

Nachdem die Eineindeutigkeit bewiesen ist, kann die Berechnung der Punkte einer Abtastschicht auf das Berechnen der Wellengrößen zu den einzelnen Punkten der skalaren, unabhängigen Variablen δ zurückgeführt werden.

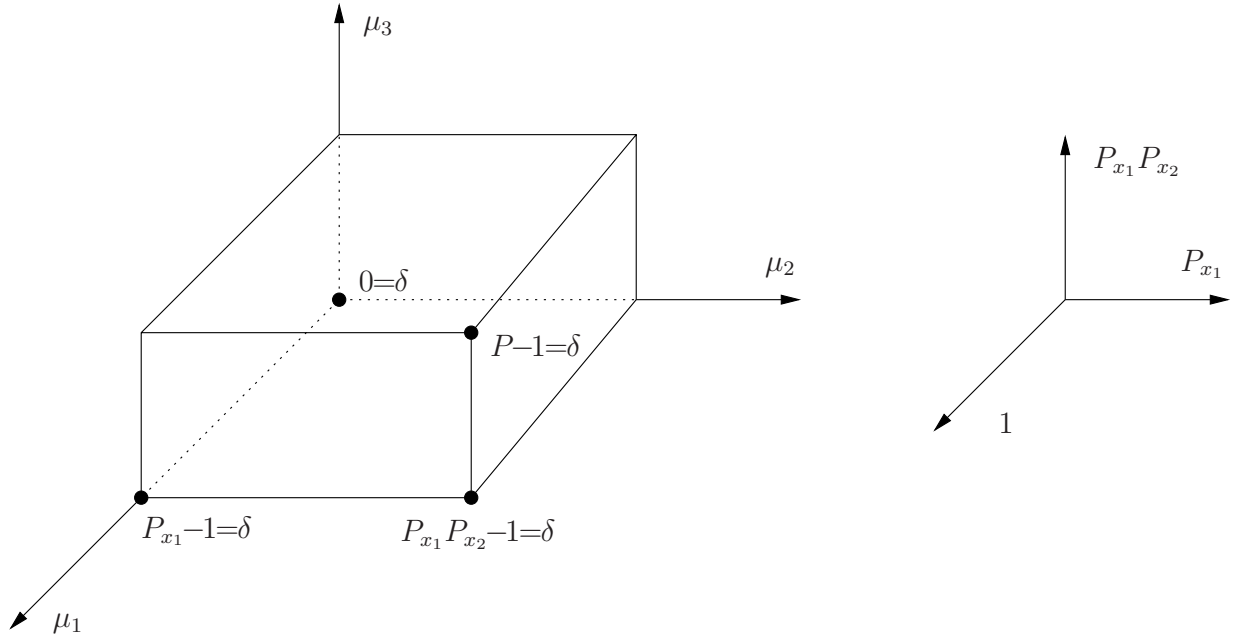


Bild 2.32: Koordinatensystem (links) und symbolische Darstellung der Gewichtungsfaktoren (rechts)

Abschließend soll noch einmal der nichtlineare Charakter der Abbildung veranschaulicht werden. Dazu nehmen wir drei konkrete Punkte δ aus dem Berechnungsgebiet an, die jeweils mit einem Punkt μ' korrespondieren sollen, d.h.

$$\delta^{(1)} = f(\mu'^{(1)}) , \quad \delta^{(2)} = f(\mu'^{(2)}) , \quad \delta^{(3)} = f(\mu'^{(3)}) . \quad (2.240)$$

Aus der Summe bzw. der Differenz zweier δ folgt i. Allg. nicht, dass diese der Summe bzw. Differenz der zugehörigen μ' entspricht, d.h.

$$\delta^{(1)} \pm \delta^{(2)} = \delta^{(3)} \implies \mu'^{(1)} \pm \mu'^{(2)} \stackrel{\text{i.Allg.}}{\neq} \mu'^{(3)} . \quad (2.241)$$

Für die umgekehrte Richtung trifft dies ebenso zu

$$\mu'^{(1)} \pm \mu'^{(2)} = \mu'^{(3)} \implies \delta^{(1)} \pm \delta^{(2)} \stackrel{\text{i.Allg.}}{\neq} \delta^{(3)} . \quad (2.242)$$

Schränkt man sich allerdings auf Vektoren $\mu'^{(1)}$ und $\mu'^{(2)}$ ein, deren Summe

$$0 \leq \mu_{\kappa}^{(1)} + \mu_{\kappa}^{(2)} \leq P_{x_{\kappa}} - 1 \quad \forall \kappa = 1, 2, \dots, k-1 \quad (2.243)$$

genügt, so gilt

$$f(\mu'^{(1)} + \mu'^{(2)}) = \delta^{(1)} + \delta^{(2)} \quad (2.244)$$

$$\iff g(\delta^{(1)} + \delta^{(2)}) = \mu'^{(1)} + \mu'^{(2)} . \quad (2.245)$$

Wir führen den Beweis der Übersicht halber nur für den relevanten Fall $k = 4$. Mit

$$\delta^{(1)} = \mu_1^{(1)} + \mu_2^{(1)} P_{x_1} + \mu_3^{(1)} P_{x_1} P_{x_2} \quad \text{und} \quad \delta^{(2)} = \mu_1^{(2)} + \mu_2^{(2)} P_{x_1} + \mu_3^{(2)} P_{x_1} P_{x_2} \quad (2.246)$$

erhält man die rechte Gleichungsseite zu

$$\delta^{(1)} + \delta^{(2)} = \mu_1^{(1)} + \mu_2^{(1)} P_{x_1} + \mu_3^{(1)} P_{x_1} P_{x_2} + \mu_1^{(2)} + \mu_2^{(2)} P_{x_1} + \mu_3^{(2)} P_{x_1} P_{x_2} . \quad (2.247)$$

Da $\boldsymbol{\mu}'^{(1)} + \boldsymbol{\mu}'^{(2)}$ ebenfalls im Berechnungsgebiet liegt, lässt sich die linke Gleichungsseite zu

$$f(\boldsymbol{\mu}'^{(1)} + \boldsymbol{\mu}'^{(2)}) = (\mu_1^{(1)} + \mu_1^{(2)}) + (\mu_2^{(1)} + \mu_2^{(2)})P_{x_1} + (\mu_3^{(1)} + \mu_3^{(2)})P_{x_1}P_{x_2}. \quad (2.248)$$

berechnen.

Die Richtigkeit der umgekehrten Beziehung werden wir im Folgenden zeigen. Nach Gleichung (2.230) berechnen wir zunächst die dritte Koordinate zu

$$\mu_3^{(3)} = \left\lceil \frac{\delta^{(1)} + \delta^{(2)}}{P_{x_1}P_{x_2}} \right\rceil = \left\lceil \frac{\mu_1^{(1)} + \mu_1^{(2)}}{P_{x_1}P_{x_2}} + \frac{\mu_2^{(1)} + \mu_2^{(2)}}{P_{x_2}} + \mu_3^{(1)} + \mu_3^{(2)} \right\rceil = \mu_3^{(1)} + \mu_3^{(2)}, \quad (2.249)$$

wobei die letzte Identität durch mehrfache Anwendung von Gleichung (C.35) gewonnen wurde, d. h.

$$\frac{\mu_1^{(1)} + \mu_1^{(2)} + P_{x_1}[\mu_2^{(1)} + \mu_2^{(2)}]}{P_{x_1}P_{x_2}} \leq \frac{\mu_1^{(1)} + \mu_1^{(2)} + P_{x_1}[P_{x_2} - 1]}{P_{x_1}P_{x_2}} = 1 + \frac{\mu_1^{(1)} + \mu_1^{(2)} - P_{x_1}}{P_{x_1}P_{x_2}} \leq 1 - \frac{1}{P_{x_1}P_{x_2}} < 1. \quad (2.250)$$

Die zweite Koordinate berechnet man zu

$$\mu_2^{(3)} = \left\lceil \frac{\delta^{(1)} + \delta^{(2)} - (\mu_3^{(1)} + \mu_3^{(2)})P_{x_1}P_{x_2}}{P_{x_1}} \right\rceil = \left\lceil \underbrace{\frac{\mu_1^{(1)} + \mu_1^{(2)}}{P_{x_1}}}_{<1} + \mu_2^{(1)} + \mu_2^{(2)} \right\rceil = \mu_2^{(1)} + \mu_2^{(2)}. \quad (2.251)$$

Die erste Koordinate kann man dann zu

$$\mu_1^{(3)} = \left\lceil \delta^{(1)} + \delta^{(2)} - (\mu_2^{(1)} + \mu_2^{(2)})P_{x_1} - (\mu_3^{(1)} + \mu_3^{(2)})P_{x_1}P_{x_2} \right\rceil = \left\lceil \mu_1^{(1)} + \mu_1^{(2)} \right\rceil = \mu_1^{(1)} + \mu_1^{(2)}. \quad (2.252)$$

ermitteln, womit Gleichung (2.245) auch bewiesen ist.

Kapitel 3

Das Engineeringsystem SPACE

Für Steuerungs- und Regelungsfunktionen mit sicherheitstechnischer Bedeutung wird das digitale Leitsystem Teleperm XS verwendet. Die typischen Anwendungen von Teleperm XS liegen im Bereich des Reaktorschutzes und der ESFAS-Funktionen (Engineered Safety Features Actuation System). Des Weiteren können sicherheitsrelevante Funktionen wie etwa die Reaktorregelung oder die Aufgaben des Steuerstabsfahrrechners übernommen werden.

Die Projektierung und Wartung von Teleperm XS Systemen erfolgt mit dem Engineeringsystem SPACE (SPEcification And Coding Environment). Das Engineeringsystem SPACE stellt eine graphische Oberfläche zur Verfügung, unter der sowohl die Leittechnikfunktionen, als auch die Hardwarearchitektur spezifiziert werden können. Aus der Spezifikation kann automatisch C-Code generiert werden, der anschließend übersetzt, gebunden und in das Zielsystem (mehrere Verarbeitungseinheiten der Teleperm XS Hardware) geladen werden kann. [KWU 98]

Zur Zeit wird am Lehrstuhl für Nachrichtentechnik der Ruhr-Universität-Bochum an einem Zusatzpaket zum Engineeringsystem SPACE gearbeitet, welches die i. W. von Fettweis entwickelten Wellendigitalfilter, unter Beibehaltung aller Qualitätsmerkmale in Bezug auf die Korrektheit der erstellten Software, in das Programmpaket SPACE einbindet. Die Wellendigitalfilter sollen zur numerischen Integration von partiellen Differentialgleichungen, insbesondere der Neutronendifusionsgleichung, genutzt werden, um das Potential digitaler Leittechnik zur weiteren Optimierung von Steuerungs- und Regelungsstrategien nutzbar zu machen.

3.1 Der Aufbau des vom Engineeringsystem SPACE erzeugten C-Codes

Das Engineeringsystem SPACE ermöglicht die graphische Projektierung einerseits der Hardwarekonfiguration und andererseits der Funktionen des zu realisierenden Leittechniksystems. Ausschlaggebend für die Funktion des Leittechniksystems ist die graphische Spezifikation der Software, was i. W. durch Verschalten von vorgegebenen Funktionsbausteinen (FB) erfolgt, wodurch Funktionspläne (FP) entstehen. Die Funktionsbausteine sind manuell codiert, qualifiziert und anlagenunabhängig. Jeder Funktionsplan eines Projektes ist einer in der Hardwarekonfiguration angegebenen Verarbeitungseinheit zuzuweisen. Auf jeder Verarbeitungseinheit läuft im Betrieb der Anlage eine Ablaufumgebung, in die eine Gruppe von Funktionsplänen eingebettet ist. Eine Funktionsplangruppe (FPG) besteht aus einem oder mehreren Funktionsplänen, die wiederum die Funktionsbausteine beinhalten. Der C-Code-Generator erzeugt zu jedem Funktionsplan und zu jeder Funktionsplangruppe automatisch eine C-Code-Datei und ein zugehöriges Header-File. Der C-Code der Funktionspläne ruft in seinem Code die fertig vorliegenden Funktionsbausteine auf, wobei die Funktionspläne selbst von der Funktionsplangruppe aufgerufen werden. Somit liegt ein kompletter C-Code für eine Verarbeitungseinheit vor, der nun compiliert und auf das

Zielsystem TELEPERM XS geladen werden kann oder von einer Simulationsumgebung zur Verifizierung compiliert wird. Die Aufrufhierarchie verdeutlicht Bild 3.1.

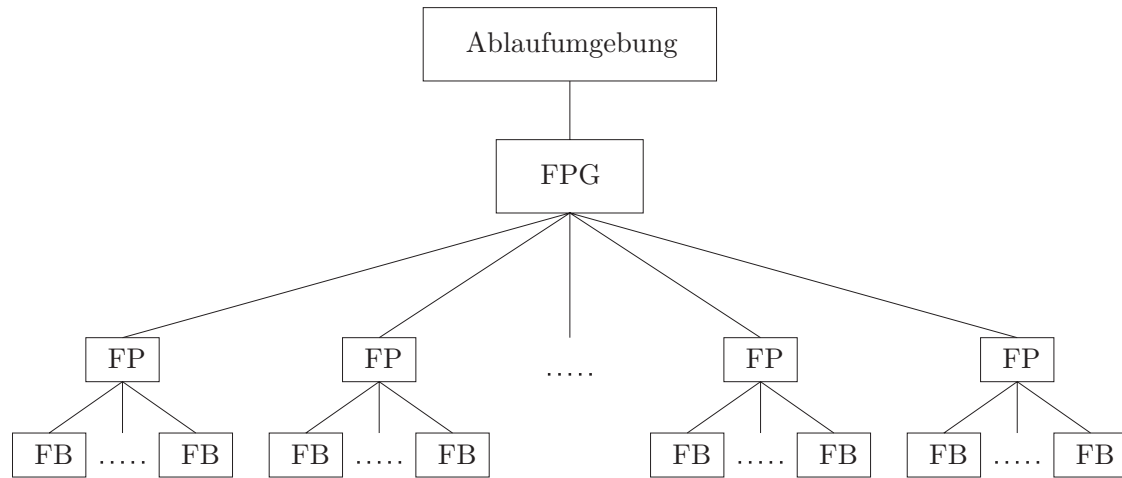


Bild 3.1: Aufrufhierarchie Funktionsplangruppe-Funktionsplan-Funktionsbaustein

3.2 Die Einbindung der Wellendigitalfilter in das Engineeringsystem SPACE

Vielen Bauelementen der klassischen Theorie elektrischer Netze, können entsprechende Wellendigitalfilter-Elemente zugeordnet werden. Um nun das numerische Integrationsverfahren nach dem Wellendigitalfilterprinzip in das Engineeringsystem SPACE zu integrieren, könnte man jedes Wellendigitalfilter-Element als einen Funktionsbaustein realisieren. Dieses Vorgehen führt allerdings zu den folgenden drei schwerwiegenden Problemen:

- Das Programmpaket SPACE detektiert verzögerungsfreie gerichtete Schleifen, obwohl keine vorhanden sind:

Die Prüfung der Funktionsplangruppen und Funktionspläne auf verzögerungsfreie gerichtete Schleifen erfolgt ohne Berücksichtigung der inneren Struktur der miteinander verschalteten Funktionsbausteine [Waed94]. Dieses Vorgehen kommt in nahezu allen Problemstellungen einer Detektion von echten verzögerungsfreien gerichteten Schleifen gleich, da davon ausgegangen wird, dass jeder Eingang eines Funktionsbausteins alle Ausgänge beeinflusst. In der Theorie der Wellendigitalfilter existieren mehrere Bauelemente, bei denen das Ausgangssignal eines Tores unabhängig von dem Eingangssignal des gleichen Tores ist. An dieser Stelle wird eine mögliche verzögerungsfreie gerichtete Schleife durch eine spezielle Wahl der inneren Struktur der Bauelemente aufgebrochen. In nahezu jeder praktisch relevanten WDF-Nachbildung einer analogen Referenzschaltung tritt eines dieser Bauelemente auf, da diese Bauelemente gerade für das Aufbrechen der verzögerungsfreien gerichteten Schleifen Verwendung finden. Bei dem Wellendigitalfilter nach Bild 3.2 würde das Programmpaket SPACE eine verzögerungsfreie gerichtete Schleifen detektieren, obwohl das linke Bauelement diese Schleife aufbricht.

- Im Programmpaket SPACE ist nur eine unabhängige Variable vorgesehen:

Das Programmpaket SPACE stellt nur Totzeitglieder in der Zeit zur Verfügung. Möchte man aber eine partielle Differentialgleichung mit der Zeit und den drei Ortskoordinaten als unabhängige Variablen lösen, so hat das zugehörige zeitdiskrete System mindestens vier unabhängige Variablen.

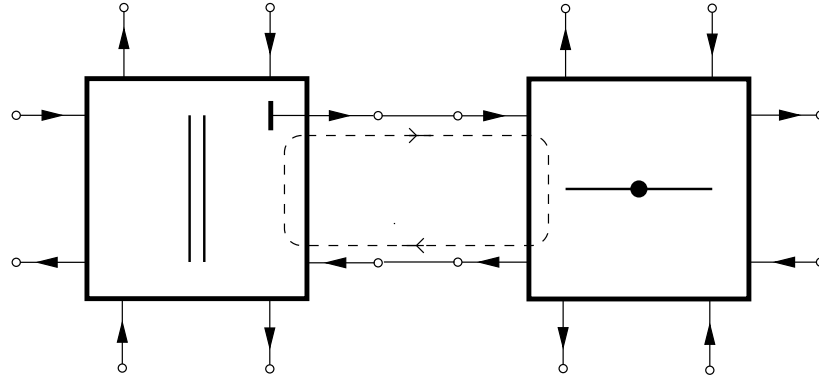


Bild 3.2: Aufbrechen von verzögerungsfreien gerichteten Schleifen

- Einem WDF-Baustein kann nicht ohne weiteres einem C-Code-Modul zugeordnet werden:

Der Grund liegt darin, dass die Berechnung eines Wellendigitalfilters nicht bausteinweise erfolgt. Stattdessen wird die Berechnungsreihenfolge nach Gleichung (2.216) genutzt.

Um diesen beiden Problemen entgegenzutreten zu können, werden wir zunächst den Aufbau eines Funktionsbausteins näher untersuchen. Im Bild 3.3 ist die Aufrufschnittstelle eines Funktionsbausteins

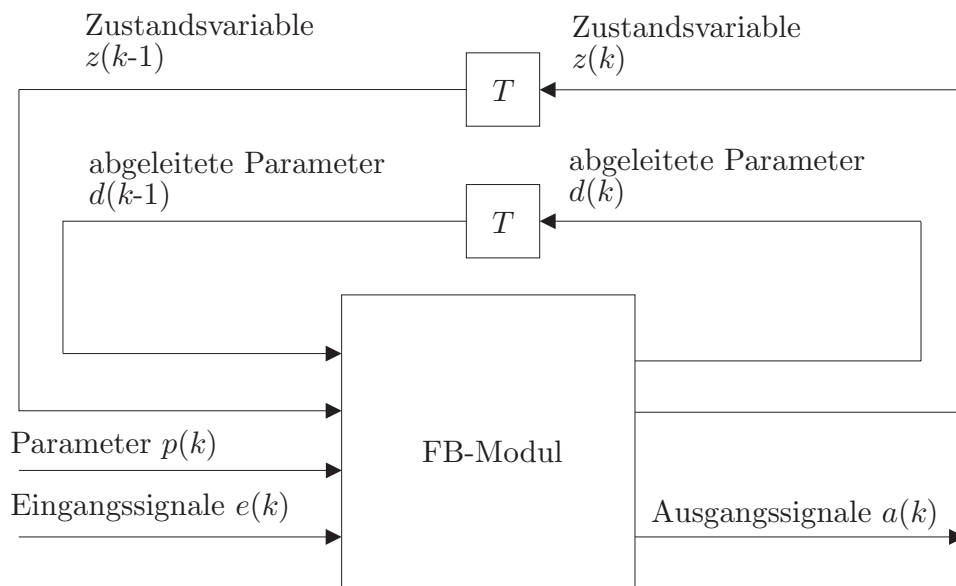


Bild 3.3: Funktionsbaustein eingebettet in einen Funktionsplan

dargestellt. Fasst man die dort auftretenden Vektoren wie folgt zusammen

$$\mathbf{w}(k+1) = \begin{bmatrix} \mathbf{z}(k) \\ \mathbf{d}(k) \end{bmatrix}, \quad \mathbf{x}(k) = \begin{bmatrix} \mathbf{p}(k) \\ \mathbf{e}(k) \end{bmatrix}, \quad \mathbf{y}(k) = \mathbf{a}(k) \quad (3.1)$$

und ändert den Laufindex ab, d.h. $k = \nu$, so erhalten wir die in der digitalen Signalverarbeitung übliche Zustandsraumdarstellung eines eindimensionalen Systems

$$\begin{aligned} \mathbf{w}(\nu+1) &= \mathbf{A} \mathbf{w}(\nu) + \mathbf{B} \mathbf{x}(\nu) \\ \mathbf{y}(\nu) &= \mathbf{C} \mathbf{w}(\nu) + \mathbf{D} \mathbf{x}(\nu). \end{aligned} \quad (3.2)$$

Folglich kann jedes zeitdiskrete eindimensionale, lineare mathematische Modell mittels eines Funktionsbaustein-Moduls in Verbindung mit den zugehörigen Zustandsspeichern gebildet werden. Ein 1D-Wellendigitalfilter ist eine i. d. R. nichtkanonische Zustandsraumdarstellung eines Digitalfilters.

Wir halten als Ergebnis dieses Abschnitts fest, dass eine Realisierung eines Wellendigitalfilters mit einzelnen Wellendigitalfilter-Elementen, ausgeführt als Funktionsbaustein, i. Allg. nicht möglich ist, ein 1D-Wellendigitalfilter aber selber immer durch ein Funktionsbaustein-Modul aufgebaut werden kann.

Wir werden von dieser Erkenntnis Gebrauch machen und ein komplettes Wellendigitalfilter als einen Funktionsbaustein realisieren. Der Leittechniker in der Praxis kann dann auf fertige Funktionsbausteine zurückgreifen, welche z. B. Löser einer Differentialgleichung repräsentieren. Dieses Vorgehen erscheint auch aus dem Grund sinnvoll, da das Auffinden eines Wellendigitalfilters zu einer gegebenen Differentialgleichung (z. B. Neutronendifusionsgleichung) einen gewissen Sachverstand voraussetzt [Luhm04].

Im Folgenden soll beispielhaft der Original C-Code des Funktionsbausteins `fb_154.c` erläutert werden. Es handelt sich hierbei i. W. um einen Integrierer, der durch verschiedene Steuersignale noch Sonderfunktionen, wie z. B. Schnellangleich, externen Halt usw. ausführen kann.

Der Original C-Code ist für einen Integrierer recht umfangreich. Dies liegt zum einen an den realisierten Sonderfunktionen und zum anderen an eingebauten Konsistenzprüfungen. Wir interessieren uns in diesem Abschnitt nur für den eigentlichen Verarbeitungsteil des Integrierers. Aus diesem Grunde haben wir den C-Code zu diesem Programmpfad kompakt dargestellt:

```
dlOut    = ddlC * (DL_t) eas1.v + mdlState ;
aas1.v   = (FL_t) dlOut ;
mdlState = ddlC * (DL_t) eas1.v + (DL_t) aas1.v;
```

Hinter `DL_t`, `FL_t` verbirgt sich `Double` und `Float`. Mit `(DL_t)` wird nicht multipliziert, sondern `DL_t` bezeichnet den so genannten Cast-(Umwandlungs-)Operator d. h. die Variablen `eas1.v`, `dlOut` und `aas1.v` werden explizit in den entsprechenden Datentyp umgewandelt.

Bevor wir mit der signaltheoretischen Untersuchung beginnen, sei bemerkt, dass jedes Aufrufen des C-Codes einem Fortschreiten um einen Abtastpunkt entspricht, d. h. die Betriebsperiode des Funktionsbausteins ist durch die Betriebsperiode des Funktionsplanes vorgegeben.

Nun zum eigentlichen Programmcode. In der ersten Programmzeile wird das Eingangssignal `eas1.v` mit einer Konstanten `ddlC` multipliziert und der aktuelle Wert der Zustandsgröße `mdlState` hinzuaddiert. Das Ergebnis wird der temporären Variablen `dlOut` zugewiesen. Die temporäre Variable `dlOut` stellt gleichzeitig das Ausgangssignal dar. Deshalb wird der Wert von `dlOut` der Ausgangsvariablen `aas1.v` in der zweiten Zeile zugewiesen. In der dritten Zeile wird dann der neue Wert der Zustandsgröße aus dem Ein- und dem Ausgangssignal berechnet.

Die temporäre Konstante `ddlC` besteht aus der Betriebsperiode T_a und einer Integrierkonstanten T_i : $ddlC = T_a / 2T_i$.

Fassen wir den Programmcode als ein lineares konstantes zeitdiskretes System auf und eliminieren wir `dlOut`, so können die Systemgleichungen in die übliche Standardnotation Gleichung (3.2) gebracht werden

$$\begin{aligned} aas1.v &= mdlState + \frac{T_a}{2T_i} eas1.v \\ mdlState &= mdlState + \frac{T_a}{T_i} eas1.v. \end{aligned} \tag{3.3}$$

An dieser Stelle ist es von größter Wichtigkeit auf Folgendes hinzuweisen. Die Variable `mdlState` ist eine skalare Variable. Möchte man die Werte dieser Variablen zu mehreren Zeitpunkten speichern, so muss die Variable ein Feld sein. In dem Fall liegt kein Problem bei der Zuordnung vor. Eine mögliche Zuordnung ist `mdlState[k] = x(k)` und `mdlState[k+1] = x(k+1)`. Erfordert die Anwendung nur den aktuellen Wert der Zustandsvariablen, so lässt sich die Zustandsgröße in einem Programm durch eine skalare Variable

realisieren. Fraglich ist nun, ob man der Programmvariablen den Wert $x(k)$ oder $x(k+1)$ zuordnet. Die Frage kann folgendermaßen beantwortet werden. Solange $x(k)$ benötigt wird, muss $\mathbf{mdlState} = x(k)$ gelten. Wird $x(k)$ nicht mehr benötigt, so kann die Zuweisung $\mathbf{mdlState} = x(k+1)$ erfolgen. Dies erfordert eine Reihenfolge der Berechnung und zwar werden erst alle Leseoperationen, dann der Wechsel von k nach $k+1$ durchgeführt und zum Schluss erfolgt die Schreiboperation. Diese Reihenfolge ist der Grund dafür, dass die Zeilen in Gleichung (3.3) gegenüber denen in Gleichung (3.2) vertauscht sind.

Die Übertragungsfunktion berechnet sich mit $\mathbf{H}(z) = \mathbf{C}[z\mathbf{1} - \mathbf{A}]^{-1}\mathbf{B} + \mathbf{D}$ zu

$$H(z) = \frac{T_a}{2T_i} \left[\frac{2}{z-1} + 1 \right] = \frac{T_a}{2T_i} \left[\frac{z+1}{z-1} \right] = \frac{T_a}{2T_i\psi}. \quad (3.4)$$

Der Übertragungsfunktion ist zu entnehmen, dass es sich um einen Integrierer nach der Trapezregel handelt.

3.3 Zur Wahl der Abtastperioden

Die Wahl der Abtastmatrix \mathbf{X}_A ist zunächst durch das Abtasttheorem festgelegt. Zusätzliche Randbedingungen sind der vorhandene Speicherplatz, die zur Verfügung stehende Rechenzeit und Forderungen des Lastenhefts. Eine Forderung des Pflichtenhefts ist die Festlegung eines Zeitintervalls, in dem die Berechnung der Lösung zu erfolgen hat. Die bisher eingeführten Zeiten sind die Zeit t als unabhängige Variable der Differentialgleichung und die mehrdimensionale Zeit \mathbf{t} . Keine der beiden Zeiten hat allerdings unmittelbar mit dem Fortschreiten der Rechenzeit des implementierten Algorithmus zu tun. Das Fortschreiten der Rechenzeit hängt unter anderem von der verwendeten Hardware ab. Die Simulationszeit und die unabhängige Variable t stimmen i. Allg. nicht überein. Bei unserer vorliegenden Echtzeitanwendung (siehe Gesamtkonzept) heißt dies nun, dass alle notwendigen Berechnungen zur Abarbeitung einer Abtastschicht innerhalb der Abtastperiode der Zeit durchzuführen sind. Natürlich kann man, wenn nicht genügend Rechenzeit zur Verfügung steht, die Abtastperiode T erhöhen oder die Anzahl der Abtastpunkte P reduzieren, so dass alle nötigen Berechnungen innerhalb der Abtastperiode T ausgeführt werden können. Dieses Vorgehen, der Erhöhung der Abtastperioden, steht allerdings diametral dem Abtasttheorem gegenüber. Eine andere Möglichkeit ist eine Abtastratenumsetzung, d. h. der Löser der Differentialgleichung hat nicht die gleiche Abtastrate wie der Versorgungsblock.

3.4 Realisierungsaspekte

Die vorliegende Beschreibung der Verzögerer bedeutet im ursprünglichen Koordinatensystem, dass der Vektor $\tilde{\mathbf{b}}_v(\boldsymbol{\mu})$ zu jedem festen Zeitpunkt μ_k und beliebigem Abtastpunkt nur von Werten eines Zeitpunktes zuvor, d. h. $\mu_k - 1$, abhängt und bietet somit den Vorteil, dass nur Werte einer Abtastschicht gespeichert werden müssen. Die dabei zum Zeitpunkt μ_k benötigten Randwerte werden aus den Werten des Berechnungsgebietes zum Zeitpunkt $\mu_k - 1$ durch Gleichung (2.228) ermittelt. Die Tatsache, dass alle zur Zeitschicht μ_k zu berechnenden Werte durch Kenntnis der Werte des physikalischen Zeitpunktes zuvor berechnet werden können, ermöglicht die in 2.9 vorgestellte sehr einfache, systematische Beschreibungsweise eines mehrdimensionalen Wellendigitalfilters. Dies hat noch den weiteren Vorteil, dass nach jedem Zeitschritt globale Größen des Berechnungsgebietes \mathcal{G} wie z. B. die Gesamtleistung, der Eigenwert, berechnet werden können, bevor der nächste Zeitschritt beginnt.

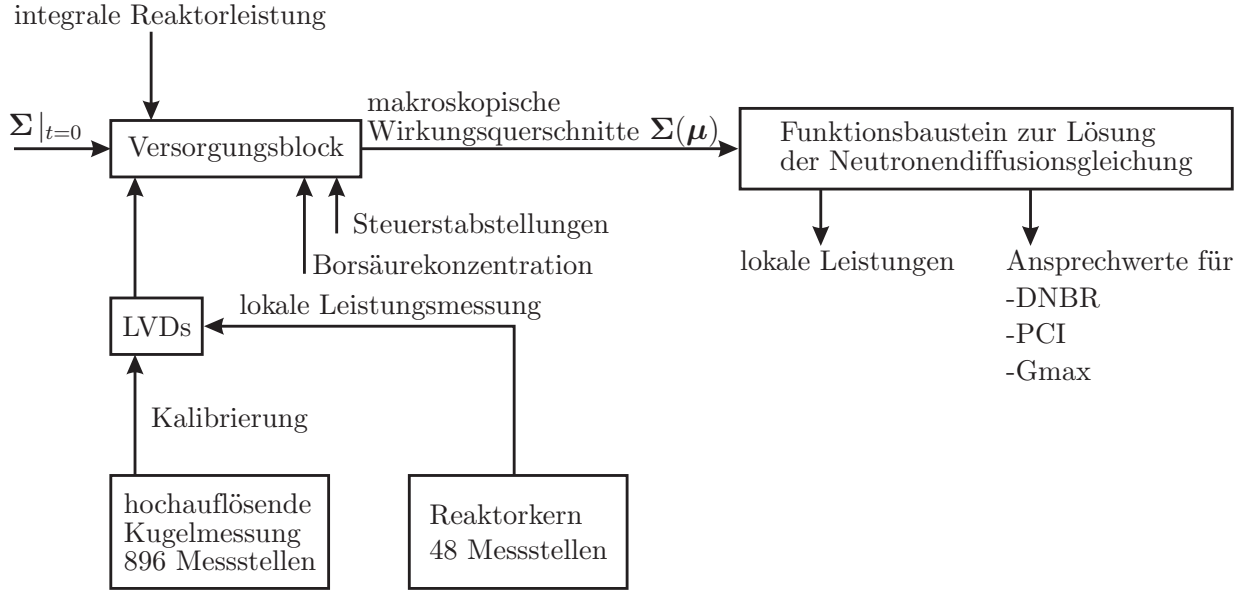


Bild 3.4: Blockschaltbild

3.5 Verknüpfung der Gleichungen des MDWDFs mit dem Funktionsbaustein

Bevor wir auf die Verknüpfung eingehen wollen, muss noch diskutiert werden, ob und an welcher Stelle der Anwender (i. d. R. der Leittechniker) Einfluss auf die Parameter des Wellendigitalfilters nehmen kann. Zunächst soll anhand eines typischen Anwendungsbeispiels die Frage geklärt werden, ob der Anwender überhaupt Parameter an das Wellendigitalfilter übergeben muss. Dazu betrachten wir Bild 3.4. Das dort angegebene Blockschaltbild erläutert u. a. die Berechnung der Leistungsverteilung im Reaktorkern, eine typische Aufgabe eines Zusatzpaketes für TELEPERM XS. Wie dem Bild zu entnehmen ist, setzt die Berechnung der lokalen Leistung die Lösung der Neutronendiffusionsgleichung voraus, was aber wiederum ein typisches Einsatzgebiet der Wellendigitalfilter ist. Legt man zwei Energiegruppen zugrunde, so berechnet sich die lokale thermische Reaktorleistung P_n innerhalb des Volumens V zu

$$P_n = \kappa \int_V (\nu \Sigma_{f1} \varphi_1 + \nu \Sigma_{f2} \varphi_2) dV. \quad (3.5)$$

Die Neutronenflussdichten φ_1 und φ_2 werden durch Lösung der Neutronendiffusionsgleichung

$$\begin{bmatrix} \text{div } \mathbf{j}_1 \\ \text{div } \mathbf{j}_2 \end{bmatrix} + \left[\text{diag} \left\{ \frac{1}{v_1}, \frac{1}{v_2} \right\} \frac{\partial}{\partial t} + \mathbf{A} \right] \begin{bmatrix} \varphi_1 \\ \varphi_2 \end{bmatrix} = \mathbf{0}, \quad \mathbf{A} = \begin{bmatrix} \Sigma_{a1} + \Sigma_{12} - \frac{\nu_1}{\lambda} \Sigma_{f1} & -\frac{\nu_2}{\lambda} \Sigma_{f2} \\ -\Sigma_{12} & \Sigma_{a2} \end{bmatrix} \quad (3.6)$$

$$(\tau_1 \frac{\partial}{\partial t} + 1) \mathbf{j}_1 + \mathcal{D}_1 \text{grad } \varphi_1 = 0$$

(Ficksches Gesetz)

$$(\tau_2 \frac{\partial}{\partial t} + 1) \mathbf{j}_2 + \mathcal{D}_2 \text{grad } \varphi_2 = 0$$

ermittelt [Pott98], [Smid75]. Die Vektoren \mathbf{j}_1 und \mathbf{j}_2 sind die Neutronenstromdichten der Energiegruppen 1 und 2. Die Parameter Σ beschreiben die makroskopischen Wirkungsquerschnitte: Σ_a ist der Absorptionsquerschnitt, Σ_f ist der Wirkungsquerschnitt für die Kernspaltung und Σ_{12} beschreibt die Streuung von der Energiegruppe 1 in die Energiegruppe 2. \mathcal{D}_1 und \mathcal{D}_2 stellen die Diffusionskonstanten dar. Die in der Matrix \mathbf{A} auftretenden Größen ν_1 und ν_2 sind in diesem Zusammenhang die mittleren Anzahlen der Spaltneutronen einer Energiegruppe und nicht die diskreten Variablen des neuen Koordinatensystems.

Die Größen τ_1 , τ_2 , v_1 und v_2 dienen der Hyperbolisierung der üblicherweise verwendeten Form der Differentialgleichung. Alle auftretenden Parameter sind i. Allg. Funktionen des Ortes, werden aber lokal als konstant angenommen. Räume mit konstanten Parametern werden den Teilberechnungsgebieten \mathcal{G}_n zugeordnet.

Zur weiteren Erläuterung des Blocksaltbildes betrachten wir eine Betriebsperiode des Reaktors in der keine Änderung des Beladungszustands vorgenommen wird (i.d.R. 1 Jahr). Vor dem Reaktorstart ($t=0$) wird dieser mit Brennstäben beladen. Aus den bekannten Daten dieser Brennstäbe lassen sich die makroskopischen Wirkungsquerschnitte berechnen $\Sigma|_{t=0}$. Mit diesen Parametern lässt sich nun die Leistungsverteilung berechnen. Während des Betriebs werden sich die Materialeigenschaften des Brennstoffes ändern, was sich in der Differentialgleichung durch Veränderung der Parameter äußert. Die makroskopischen Wirkungsquerschnitte sind somit nicht nur Funktionen des Ortes, sondern auch der Zeit. Die Berechnung der makroskopischen Wirkungsquerschnitte übernimmt der Versorgungsblock. Der Versorgungsblock benötigt dazu die integrale Reaktorleistung, die makroskopischen Wirkungsquerschnitte zur Zeit der Neubeladung und die lokalen Leistungen an (z.B. 48) Messstellen. Die lokalen Leistungen werden durch LVDs (Leistungsverteilungsdetektoren) gemessen. Die Leistungsverteilungsdetektoren wiederum werden ca. alle 14 Tage durch hochauflösende Kugelmessungen an den $28 \times 32 = 896$ Messstellen kalibriert. Durch Vergleich der gemessenen und der für $t = 0$ berechneten Leistungsverteilung kann der Versorgungsblock eine Adaption der makroskopischen Wirkungsquerschnitte $\Sigma(\mathbf{x})$ vornehmen [Fisc00b].

Als Ergebnis der Diskussion halten wir fest, dass zur Lösung der Differentialgleichung die makroskopischen Wirkungsquerschnitte, die Anzahlen der Spaltneutronen und die Diffusionskonstanten erforderlich sind. Ein entsprechender Funktionsbaustein, der den Löser der Neutronendifusionsgleichung repräsentiert, muss demnach von außen die Parameter des zugehörigen Wellendigitalfilters übergeben bekommen.

Nun stellt sich die Frage, auf welche Art man generell Parameter an die Funktionsbausteine übergibt. Betrachten wir erneut Bild 3.3, so stellen wir fest, dass die Übergabe prinzipiell über den dort angegebenen Eingangsvektor $\mathbf{e}(k)$ und den Parametervektor $\mathbf{p}(k)$ erfolgen kann. Die Werte des Parametervektors $\mathbf{p}(k)$ werden über eine so genannten Parametriermaske bei der Erstellung eines Funktionsplans eingegeben und können somit während des Betriebs nicht mehr verändert werden. Dies liegt darin begründet, dass die Parameter des Vektors $\mathbf{p}(k)$ zwar Übergabeparameter eines FB-Softwaremoduls sind, aber keine Eingangssignale eines Funktionsbausteins, wie er im Schaltplaneditor auftritt. Da unsere Aufgabenstellung, wie oben bereits diskutiert, eine Änderung der Parameter des Wellendigitalfilters im Betrieb erfordert, müssen diese Parameter als Eingänge des Funktionsbausteins aufgefasst werden.

Eine Änderung der Parameter des Wellendigitalfilters hat aber eine Änderung der Streumatrix \mathbf{S} zur Folge, was wiederum eine Änderung der Matrizen \mathbf{L}'_q , \mathbf{L}'_v und \mathbf{L}'_e in Gleichung (2.219) bewirkt. Aus Gleichung (2.219) wird aber unmittelbar der C-Code der Funktionsbausteine generiert, d. h. der gemäß Gleichung (2.219) generierte C-Code hängt von den Übergabeparametern des Funktionsbausteins ab. Die sich daraus ergebenden Änderungen zur Generierung des C-Codes ist Gegenstand der Arbeit [Juss00].

Die während des Betriebs sich nicht ändernden Parameter des Wellendigitalfilters können über eine Parametriermaske in SPACE eingegeben werden. Zudem sind bei der eigentlichen Erstellung eines Funktionsbausteins natürlich beliebig viele Parameter und die Struktur selber änderbar.

Syntaktische Grundlage für die Einbindung eines Wellendigitalfilters in einen Funktionsbaustein ist die Übergabeschnittstelle. Beim Aufruf eines Funktionsbausteins durch einen Funktionsplan wird grundsätzlich die Adresse einer Variablen vom Typ `fb_<log_id>_t` übergeben. Der Datentyp dieser Variablen wird mithilfe des C-Schlüsselwortes `struct` zur Definition von Datenstrukturen festgelegt. Die Datenstruktur lautet [KWU 99] :

```
typedef struct
{ const TX_t      *txExtIdent_p;          /* externe Identifikation
```



```

const <E/A-Typ> *e<e/a-Typ>1_p;      /* Zeiger auf 1. Eingangssignal      */
const <E/A-Typ> *e<e/a-Typ>2_p;      /* Zeiger auf 2. Eingangssignal      */
/* ... */
const <E/A-Typ> *e<e/a-Typ><i>_p;      /* Zeiger auf i. Eingangssignal      */
<E/A-Typ> *a<e/a-Typ><i+1>_p;          /* Zeiger auf 1. Ausgangssignal      */
<E/A-Typ> *a<e/a-Typ><i+2>_p;          /* Zeiger auf 2. Ausgangssignal      */
/* ... */
<E/A-Typ> *a<e/a-Typ><i+o>_p;          /* Zeiger auf o. Ausgangssignal      */
const <Par-Typ> *p<par-Typ>1_p;        /* Zeiger auf 1. Parameter           */
const <Par-Typ> *p<par-Typ>2_p;        /* Zeiger auf 2. Parameter           */
/* ... */
const <Par-Typ> *p<par-Typ><p>_p;        /* Zeiger auf p. Parameter           */
<Mem-Typ> *d<mem-Typ>1_p;              /* Zeiger auf 1. abgeleiteten Parameter */
<Mem-Typ> *d<mem-Typ>2_p;              /* Zeiger auf 2. abgeleiteten Parameter */
/* ... */
<Mem-Typ> *d<mem-Typ><d>_p;              /* Zeiger auf d. abgeleiteten Parameter */
<Mem-Typ> *m<mem-Typ>1_p;              /* Zeiger auf 1. Zustandsspeicher     */
<Mem-Typ> *m<mem-Typ>2_p;              /* Zeiger auf 2. Zustandsspeicher     */
/* ... */
<Mem-Typ> *m<mem-Typ><m>_p;              /* Zeiger auf m. Zustandsspeicher     */
}fb_log_id>_t ;

```

Da nach Verlassen des Funktionsbausteins alle internen Variablen wieder gelöscht werden, müssen die Zustände in den Zustandsspeichern `m<mem-Typ><nr>_p` abgelegt werden. Die Anzahl der Zustandsspeicher ergibt sich als Produkt aus den Zustandsgrößen des mehrdimensionalen Wellendigitalfilters und der Anzahl der Abtastpunkte einer Abtastschicht: $n_v P$. Legen wir im Fall $k = 4$ eine Zahl von 100 Berechnungspunkten in jede Richtung zugrunde, so benötigen wir $n_v 10^6$ skalare Übergabeparameter, was nicht sinnvoll zu realisieren ist. Um eine derartige Menge an Variablen kompakt beschreiben zu können, gibt es in der Programmiersprache C das Konzept der vektoriellen Datentypen. Die Deklaration von Feldern ist aber im Programmpaket SPACE nicht vorgesehen. *Die Einbindung von vektoriellen Datentypen sollte zukünftig in das Programmpaket SPACE integriert sein.* Dazu sind entsprechende Erweiterungen an dem Codegenerator des Programmpaketes SPACE vorzunehmen [Waed00].

Wir werden im weiteren Verlauf dieser Arbeit davon ausgehen, dass innerhalb des Codegenerators auch vektorielle Signale definiert werden können. Für die Quellen, Verzögerer und Ausgänge sind die an einen Funktionsbaustein übergebenen Variablen Vektoren, deren Länge gleich der Anzahl Abtastpunkte einer Abtastschicht ist. Der Wert eines Feldelementes entspricht `*(args->mdl1_p)[elementnr]`. Um eine bessere Übersicht in den Quellcodes zu ermöglichen, werden wir in den Quellcodes anstatt `*(args->mdl1_p)[elementnr]` Bezeichner nutzen, die den verwendeten Variablen des Wellendigitalfilters entsprechen. Diese Bezeichner werden dann mittels den in den Header-Dateien definierten Präprozessordirektiven `#define` aufgelöst, in diesem Fall durch `#define b_v_1_1 (*(args->mdl1_p))`. Näheres dazu findet sich in [Voll02].

Auf die einzelnen Koordinaten kann dann mittels der Adressierung `[delta]` zugegriffen werden. Die Länge der Vektoren der variablen Parameter ist gleich der Anzahl der Teilberechnungsgebiete, wobei wir der Einfachheit halber annehmen, dass die Länge der Vektoren der variablen Parameter auch P ist. Eine Koordinate des übergebenen Vektors von Zeigern ist dabei ein Zeiger, der auf einen Vektor mit den ortsabhängigen Parametern des Wellendigitalfilters zeigt. Die neue Datenstruktur mit vektoriellen Signalen lautet

```

typedef struct{
    const TX_t *txExtIdent_p;          /* externe Identifikation            */
    const <E/A-Typ> *e<e/a-Typ>1_p[P]; /* Zeiger auf 1. Quellensignal       */

```

```

const <E/A-Typ> *e<e/a-Typ>2_p[P];      /* Zeiger auf 2. Quellensignal      */
/* ...                                  */
const <E/A-Typ> *e<e/a-Typ><n_q>_p[P];    /* Zeiger auf n_q. Quellensignal    */
const <E/A-Typ> *e<e/a-Typ><n_q+1>_p[P]; /* Zeiger auf 1. variablen Parameter */
/* ...                                  */
const <E/A-Typ> *e<e/a-Typ><i>_p[P];      /* Zeiger auf n_vP. variablen Parameter */
<E/A-Typ> *a<e/a-Typ><i+1>_p[P];          /* Zeiger auf 1. Ausgangssignal      */
<E/A-Typ> *a<e/a-Typ><i+2>_p[P];          /* Zeiger auf 2. Ausgangssignal      */
/* ...                                  */
<E/A-Typ> *a<e/a-Typ><i+o>_p[P];          /* Zeiger auf o. Ausgangssignal      */
const <Par-Typ> *p<par-Typ>1_p;           /* Zeiger auf 1. Parameter           */
const <Par-Typ> *p<par-Typ>2_p;           /* Zeiger auf 2. Parameter           */
/* ...                                  */
const <Par-Typ> *p<par-Typ><p>_p;          /* Zeiger auf p. Parameter           */
<Mem-Typ> *d<mem-Typ>1_p;                 /* Zeiger auf 1. abgeleiteten Parameter */
<Mem-Typ> *d<mem-Typ>2_p;                 /* Zeiger auf 2. abgeleiteten Parameter */
/* ...                                  */
<Mem-Typ> *d<mem-Typ><d>_p;                /* Zeiger auf d. abgeleiteten Parameter */
<Mem-Typ> *m<mem-Typ>1_p[P];              /* Zeiger auf 1. Zustandsspeicher     */
<Mem-Typ> *m<mem-Typ>2_p[P];              /* Zeiger auf 2. Zustandsspeicher     */
/* ...                                  */
<Mem-Typ> *m<mem-Typ><m>_p[P];             /* Zeiger auf m. Zustandsspeicher     */
}fb_log_id_t ;

```

Die Eingangs-, Ausgangs- und Zustandssignale bekommen in Übereinstimmung mit der Header-Datei der Funktionsbausteine fb_t.h des Programmpaketes SPACE fortlaufende Nummern, und zwar in der Reihenfolge, dass die ersten i Nummern von den Eingangssignalen belegt sind und die Ausgangssignale von den darauf folgenden Nummern. Die Zustandsgrößen erhalten neue Nummern. Wir werden die Variablennamen der Übersicht halber für die theoretische Beschreibung in Vektoren zusammenfassen:

$$\begin{aligned}
\mathbf{e}_{\text{FB}} &= \left[\begin{array}{c} e_{\langle e/a\text{-Typ} \rangle 1_p[P]} \\ \vdots \\ e_{\langle e/a\text{-Typ} \rangle \langle i \rangle_p[P]} \end{array} \right] \Bigg\}_i, & \mathbf{a}_{\text{FB}} &= \left[\begin{array}{c} a_{\langle e/a\text{-Typ} \rangle \langle i+1 \rangle_p[P]} \\ \vdots \\ a_{\langle e/a\text{-Typ} \rangle \langle i+o \rangle_p[P]} \end{array} \right] \Bigg\}_o \\
\mathbf{m}_{\text{FB}} &= \left[\begin{array}{c} m_{\langle \text{mem-Typ} \rangle 1_p[P]} \\ \vdots \\ m_{\langle \text{mem-Typ} \rangle \langle m \rangle_p[P]} \end{array} \right] \Bigg\}_m,
\end{aligned} \tag{3.7}$$

hierbei steht der erste Buchstabe für den Eingang (e), Ausgang (a) oder Zustand (m).

Wir werden im Folgenden die Dimensionen der oben definierten Vektoren diskutieren. Die Eingangssignale des Funktionsbausteins seien Spannungen und Ströme von idealen Quellen und die Parameter für das Wellendigitalfilter. Der Vektor \mathbf{e}_{FB} mit den Eingangssignalen bildet sich folglich aus einem Vektor mit den Quellenwellen und einem Vektor der Länge n_{vP} mit variablen Parametern. Die Anzahl der Eingänge des Funktionsbausteins entspricht der Summe der Quellen und der änderbaren Parameter des Wellendigitalfilters, d. h.

$$i = n_q + n_{vP} . \tag{3.8}$$

Die Anzahl der Zustandsvariablen des Funktionsbausteins entspricht natürlich der Anzahl der dynamischen Bauelemente des Wellendigitalfilters, d. h.

$$m = n_v . \tag{3.9}$$

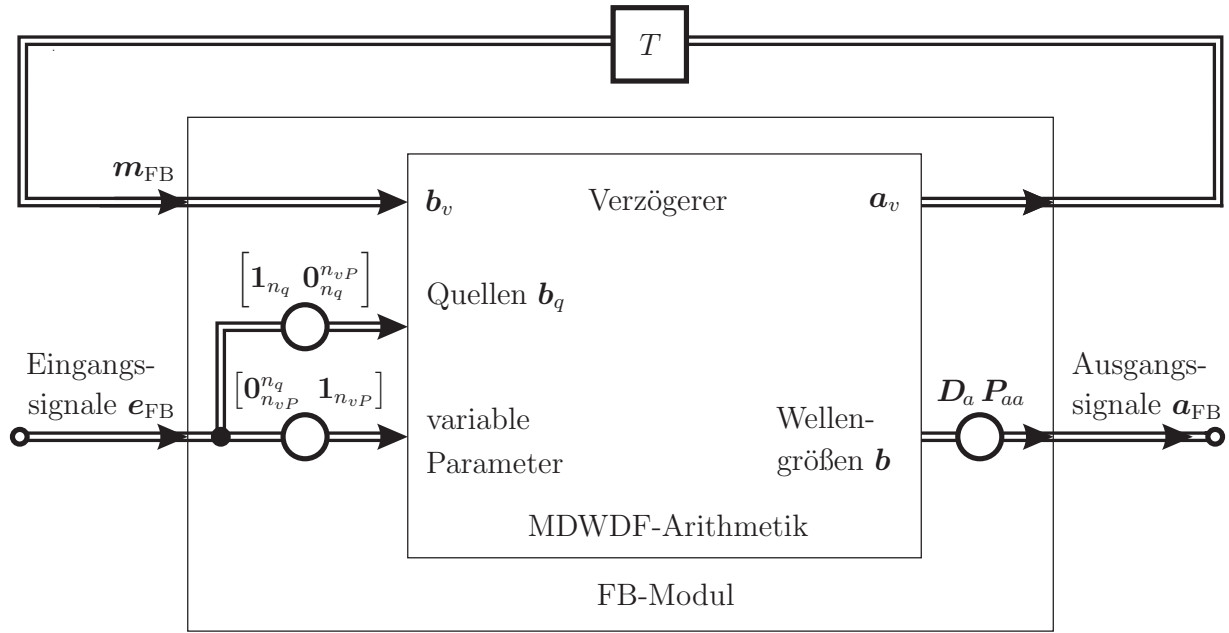


Bild 3.5: Einbettung der Wellendigitalfilter in das Programmpaket SPACE

Da die Anzahl Tore $n_g = n_q + n_v + n_e$ beträgt (wobei sich jeweils zwei Tore nur von der Richtung des Stromes und der Wellengrößen unterscheiden), können maximal n_g verschiedene Signale die Ausgangssignale des Funktionsbausteins darstellen, d.h.

$$n_{aa} \leq n_g. \quad (3.10)$$

Da das Wellendigitalfilter integraler Bestandteil des Funktionsbausteins ist, besteht nun die Notwendigkeit, diese Einbettung mathematisch zu beschreiben. Zweckmäßigerweise erfolgt die Verknüpfung der Wellengrößen a und b mit den Ein- und Ausgangssignalen e_{FB} und a_{FB} des Funktionsbausteins durch Matrizen. Wir bezeichnen die Verknüpfungsmatrix mit P_{aa} . Das Wellendigitalfilter wird vollständig durch

$$b_q = \begin{bmatrix} 1_{n_q} & 0_{n_{vP}} \end{bmatrix} e_{FB} \quad , \quad a_{FB} = D_a P_{aa} \begin{bmatrix} a \\ b \end{bmatrix} \quad , \quad m_{FB} = b_v \quad (3.11)$$

in einen Funktionsbaustein eingebunden. Die Einbindung der Wellendigitalfilter-Arithmetik in einen Funktionsbaustein ist im Bild 3.5 dargestellt. In dieser Arbeit wollen wir nur lineare zeit-/ortsinvariante Systeme betrachten. Ergebnisse zur Behandlung von linearen zeit-/ortsvarianten Systemen finden sich in [Juss00].

Wir werden nun die Diagonalmatrix D_a näher erläutern. Ist der Ausgang eine Spannung, so ergibt sich das Ausgangssignal als gewichtete Summe der zugehörigen Wellengrößen

$$u = (a + b)\sqrt{R} = (a' + b')/2. \quad (3.12)$$

Ist hingegen das Ausgangssignal ein Strom, so ergibt sich das Ausgangssignal als gewichtete Differenz der zugehörigen Wellengrößen

$$i = (a - b)/\sqrt{R} = (a' - b')/(2R). \quad (3.13)$$

Der Vektor der Ausgangssignale a_{FB} berechnet sich somit als Produkt einer Diagonalmatrix mit einem Vektor, der die Summen oder Differenzen der Wellengrößen beinhaltet. Bei Verwendung von Leistungswellen ist ein Element der Diagonalmatrix D_a durch die Wurzel des Torwiderstandes gegeben, wenn das

entsprechende Ausgangssignal eine Spannung ist oder im Falle des Stromes durch das Reziproke der Wurzel des Torwiderstandes. Bei Verwendung von Spannungswellen sind die Elemente von \mathbf{D}_a entsprechend durch $1/2$ oder $1/(2R)$ gegeben.

Die Matrix \mathbf{P}_{aa} kann durch $\mathbf{P}_{aa} = [\mathbf{P}_{aaa} \mathbf{P}_{aab}]$ dargestellt werden, wobei die Zeilenanzahl beider Untermatrizen gleich n_{aa} ist und die Spaltenanzahl gleich n_g . Die Elemente der Matrix \mathbf{P}_{aaa} besitzen nur die Werte 0 und 1

$$p_{aaa\mu\nu} = \begin{cases} 1 & \text{falls der Ausgang } \mu \text{ des Funktionsbausteins eine Spannung} \\ & \text{oder einen Strom des WDF-Tores } \nu \text{ darstellt} \\ 0 & \text{sonst.} \end{cases} \quad (3.14)$$

Dagegen nehmen die Elemente von \mathbf{P}_{aab} nicht nur 0 oder 1 an, sondern auch -1, und zwar in Abhängigkeit davon, ob das Ausgangssignal eine Spannung oder ein Strom ist

$$p_{aab\mu\nu} = \begin{cases} 1 & \text{falls der Ausgang } \mu \text{ des Funktionsbausteins eine Spannung} \\ & \text{des WDF-Tores } \nu \text{ ist} \\ -1 & \text{falls der Ausgang } \mu \text{ des Funktionsbausteins ein Strom des} \\ & \text{WDF-Tores } \nu \text{ ist} \\ 0 & \text{sonst.} \end{cases} \quad (3.15)$$

Jede Zeile der Matrizen \mathbf{P}_{aaa} und \mathbf{P}_{aab} hat nur einen Eintrag, aber eine Spalte kann mehrere Einträge besitzen, nämlich dann, wenn sowohl die Spannung als auch der Strom eines Tores ein Ausgangssignal ist.

Sämtliche einfallenden Wellen $\mathbf{a}_q, \mathbf{a}_v$ und \mathbf{a}_e lassen sich durch die Permutationsmatrix \mathbf{P} aus den Wellen $\mathbf{b}_q, \mathbf{b}_v$ und \mathbf{b}_e berechnen. Folglich können die Ausgangsgrößen allein durch Kenntnis der ausfallenden Wellen $\mathbf{b}_q, \mathbf{b}_v$ und \mathbf{b}_e ermittelt werden. Mit Gleichung (2.206) erhalten wir somit

$$\mathbf{a}_{\text{FB}} = \mathbf{D}_a[\mathbf{P}_{aaa}\mathbf{a} + \mathbf{P}_{aab}\mathbf{b}] = \mathbf{D}_a[\mathbf{P}_{aaa}\mathbf{P} + \mathbf{P}_{aab}]\mathbf{b} = \underbrace{[\mathbf{D}_a\mathbf{P}_{aaa}\mathbf{P} + \mathbf{D}_a\mathbf{P}_{aab}]}_{\mathbf{A}}\mathbf{b}. \quad (3.16)$$

Kapitel 4

Syntheseverfahren

Dieses Kapitel widmet sich der systematischen Generierung mehrdimensional intern passiver Wellendigitalfilter zur numerischen Integration einer Klasse von PDGLn. In [Kumm88] wurde der Beweis erbracht, dass i. Allg. keine intern passive Schaltung zu einem drei- und höherdimensionalen Reaktanzmehrtor existiert. Da die Reaktanzmehrtore eine Untermenge passiver Mehrtore sind, existiert somit auch für die passiven Mehrtore i. Allg. keine intern passive Schaltung für drei- und mehr Dimensionen. Es ist allerdings gelungen zu vielen PDGLn eine Referenzschaltung zu finden. Da ein systematisches Vorgehen zur Durchführung der Synthese aber nicht bekannt ist, ist viel Erfahrung zur Durchführung erforderlich. Die Synthese erfordert das Aufsuchen von mehreren geeigneten Rechtsinversen zur Koordinatentransformationsmatrix unter gleichzeitiger Berücksichtigung der mehrdimensionalen konkreten Passivität einer möglichen Referenzschaltung. Anschließend wird aus der Referenzschaltung ein mehrdimensionales Wellendigitalfilter gewonnen. Die soeben erläuterte Vorgehensweise stellt die erste Phase des Softwareerstellungsprozesses dar. Diese Phase ist manuell durchzuführen. In der zweiten Phase wird das mehrdimensionale Wellendigitalfilter als graphische Spezifikation der zur Simulation verwendeten Software genutzt. Die Umsetzung der graphischen Spezifikation in eine Anweisungssequenz kann mittels Rechnerunterstützung automatisiert durchgeführt werden. Diese zweistufige Vorgehensweise bringt i.W. zwei Nachteile mit sich. Zum einen wird der Kreis der Anwender auf Sachkundige stark eingeschränkt. Zum zweiten ist die manuelle Phase mit hohen Kosten und erhöhter Gefahr eines Fehlers verbunden. Wünschenswert wäre daher ein komplett automatisierter Softwareerstellungsprozess ausgehend von der PDGL. Dieses ist aber im allgemeinen Fall nicht möglich, wie bereits erläutert wurde. Möglicherweise kann man aber die Klasse der PDGLn so einschränken, dass einerseits eine sinnvolle Klasse physikalischer Probleme behandelbar ist und andererseits ein komplett automatisierter Softwareerstellungsprozess möglich ist. Es soll daher hier ein Verfahren zur automatisierten Synthese vorgestellt werden, welches auf eine eingeschränkte Klasse PDGLn anwendbar ist, [Voll04a].

4.1 Der behandelte Typ von PDGLn

4.1.1 Einschränkungen an die PDGLn

Ausgangspunkt der Diskussion ist das folgende System nichtlinearer, nicht konstanter partieller Differentialgleichungen

$$\sum_{\xi=x,y,z,t} \sum_{\nu=1}^N \operatorname{sgn}(c_{\mu\nu}^{\xi}) |c_{\mu\nu}^{\xi}|^{1/2} \frac{\partial}{\partial \xi} [|c_{\mu\nu}^{\xi}|^{1/2} u_{\nu}] + \sum_{\nu=1}^N y_{\mu\nu}^k u_{\nu} = j_{\mu}, \quad \mu = 1, 2, \dots, N, \quad (4.1)$$

wobei die $c_{\mu\nu}^{\xi}$ und $y_{\mu\nu}^k$ reellwertige Funktionen von u_1, u_2, \dots, u_N und \mathbf{x} sind. Ein derartiges System wird auch als quasilinear bezeichnet, da die Ableitungen der abhängigen Variablen nur linear auftreten.

Man beachte die Beziehung

$$|c_{\mu\nu}^\xi|^{1/2} \frac{\partial}{\partial \xi} [|c_{\mu\nu}^\xi|^{1/2} u_\nu] = \frac{1}{2} \left[|c_{\mu\nu}^\xi| \frac{\partial u_\nu}{\partial \xi} + \frac{\partial}{\partial \xi} (|c_{\mu\nu}^\xi| u_\nu) \right]. \quad (4.2)$$

Im Falle eines Kondensators stellt diese Definition den Mittelwert aus der globalen und der differentiellen Kapazität dar, [Fett98].

Die unabhängigen Variablen des partiellen Differentialgleichungs-Systems sind \mathbf{x} . Weiter definieren wir \mathbf{u} als den Vektor der abhängigen Feldgrößen u_ν . Die eingepägten Größen j_μ sind in dem Vektor \mathbf{j} zusammengefasst. Im Falle konstanter $c_{\mu\nu}^\xi$ lautet das PDGL-System

$$[\mathbf{C}^x \mathbf{D}_x + \mathbf{C}^y \mathbf{D}_y + \mathbf{C}^z \mathbf{D}_z + \mathbf{C}^t \mathbf{D}_t + \mathbf{Y}^k] \mathbf{u} = \mathbf{j}, \quad (4.3)$$

wobei $c_{\mu\nu}^\xi$, $y_{\mu\nu}^k$ die reellen Elemente der Matrizen \mathbf{C}^ξ und \mathbf{Y}^k sind. Im weiteren Verlauf der Arbeit werden wir nur diesen Fall behandeln. Die Teilströme der einzelnen Elemente definieren wir durch $\mathbf{i}_c^x = \mathbf{C}^x \mathbf{D}_x \mathbf{u}$, \dots , $\mathbf{i}_c^t = \mathbf{C}^t \mathbf{D}_t \mathbf{u}$, $\mathbf{i}_k = \mathbf{Y}^k \mathbf{u}$.

Die Eigenschaften der Matrizen \mathbf{C}^ξ und \mathbf{Y}^k werden wir im Folgenden durch Energiebetrachtungen weiter einschränken.

4.1.2 Energiebetrachtungen

Der quellenfreie Teil der PDGL, das ist ihr homogener Teil,

$$\mathbf{i} = [\mathbf{C}^x \mathbf{D}_x + \mathbf{C}^y \mathbf{D}_y + \mathbf{C}^z \mathbf{D}_z + \mathbf{C}^t \mathbf{D}_t + \mathbf{Y}^k] \mathbf{u}. \quad (4.4)$$

soll ein passives System beschreiben. Voraussetzung für die folgenden Energiebetrachtungen ist, dass die Koordinaten mit gleichem Index der Vektoren \mathbf{u} und \mathbf{i} ein Paar komplementärer Feldgrößen ist. Mit anderen Worten, jeder Summand des Skalarprodukts $p = \mathbf{u}^T \mathbf{i}$ muss eine Energiedichte sein. Die Tatsache, dass jeder Summand die Dimension einer Energiedichte besitzt, reicht nicht aus. Zur Erläuterung betrachten wir 2 konzentrierte Elementarkörper mit den Dichten ρ_1 und ρ_2 , die sich mit den Geschwindigkeiten v_1 bzw. v_2 bewegen. Die kinetische Energiedichte ist $\rho_1 v_1^2/2 + \rho_2 v_2^2/2$. Hingegen ist $\rho_2 v_1^2/2 + \rho_1 v_2^2/2$ keine Energiedichte, obwohl die Dimensionen gleich sind.

Fassen wir die Größen \mathbf{u} und \mathbf{i} als Spannung und Strom in einer elektrischen Schaltung auf, so wird implizit angenommen, dass es sich hier um ein Paar komplementärer Feldgrößen handelt. Die Energiebetrachtungen beziehen sich dann aber auf die elektrische Schaltung und nicht auf das zugrunde liegende System.

Satz 1 *Der quellenfreie Teil des durch Gleichung (4.3) beschriebenen Systems ist genau dann passiv, wenn \mathbf{C}^x , \mathbf{C}^y , \mathbf{C}^z , \mathbf{C}^t symmetrisch sind und sowohl \mathbf{C}^t als auch der symmetrische Teil der Matrix \mathbf{Y}^k n. n. d. ¹ sind.*

Beweis 1

Der Beweis ist dem aus dem eindimensionalen bekannten ähnlich, vgl. [Tell54], [FH92].

Notwendigkeit

Für jeden beliebigen Vorgang ist nach Gleichung (2.85) zu zeigen, dass die über die Systemgrenze aufgenommene Energie nicht negativ ist, wenn die über die örtliche Segmentgrenze des Systems transportierte Energiedichte null ist, d. h.

$$\mathbf{W}_{\text{sp}}^T \mathbf{n} = 0 \text{ auf } \partial \mathcal{G}. \quad (4.5)$$

¹nicht negativ definit

Weiterhin soll zum Zeitpunkt t_0 im System keine Energie gespeichert sein, d. h. $E(t_0) = 0$. Mit $p = \mathbf{u}^T \mathbf{i}$ lautet die Bedingung

$$0 \leq \int_{t_0}^{t_1} \int_{\mathcal{G}} p \, dV dt \quad \text{für} \quad t_1 \geq t_0. \quad (4.6)$$

Da die Passivitätsforderungen für beliebige Vorgänge gelten, müssen sie auch für einen von uns speziell gewählten Vorgang gelten. Dazu nehmen für \mathcal{G} einen Quader an. Setzen wir zudem $\mathbf{i}^{ct} = \mathbf{i}^{cy} = \mathbf{i}^{cz} = \mathbf{i}_k = \mathbf{0}$, so muss $\mathbf{i}^{cx} = \mathbf{C}^x \mathbf{D}^x \mathbf{i}$ schon allein passiv sein.

Für die Spannungen und die Ströme \mathbf{i}^{cx} nehmen wir 2 spezielle Vorgänge an. Für den 1. Vorgang soll $u_\sigma = 0$ und $i_\sigma = 0$ für $\mu \neq \sigma \neq \nu$ gelten. Die anderen beiden Spannungen und Ströme folgen dem Kreisprozess (mit als stetig angenommenen Leistungsdichten)

Intervall	Spannung u_μ	Spannung u_ν	$D_x u_\mu$	$D_x u_\nu$
$x_0 \leq x \leq x_1$	$u_\mu = \begin{cases} 0 & \text{für } x = x_0 \\ \text{beliebig} & \text{für } x_0 < x < x_1 \\ u_{\mu 0} > 0 & \text{für } x = x_1 \end{cases}$	$u_\nu = 0$		0
$x_1 \leq x \leq x_2$	$u_\mu = u_{\mu 0}$	$u_\nu = \begin{cases} 0 & \text{für } x = x_1 \\ \text{beliebig} & \text{für } x_1 < x < x_2 \\ u_{\nu 0} > 0 & \text{für } x = x_2 \end{cases}$	0	
$x_2 \leq x \leq x_3$	$u_\mu = \begin{cases} u_{\mu 0} & \text{für } x = x_2 \\ \text{beliebig} & \text{für } x_2 < x < x_3 \\ 0 & \text{für } x = x_3 \end{cases}$	$u_\nu = u_{\nu 0}$		0
$x_3 \leq x \leq x_4$	$u_\mu = 0$	$u_\nu = \begin{cases} u_{\nu 0} & \text{für } x = x_3 \\ \text{beliebig} & \text{für } x_3 < x < x_4 \\ 0 & \text{für } x = x_4 \end{cases}$	0	

der im Bild 4.1 a) graphisch dargestellt ist.

Die eingangs gestellte Bedingung Gleichung (4.5) ist wegen $\mathbf{u}(x_0) = \mathbf{u}(x_4) = \mathbf{0}$ erfüllt. Aufgrund des quaderförmigen örtlichen Berechnungsgebietes hat das Volumenintegral mit $t_1 = t_0 + T$ feste Grenzen

$$\int_{t_0}^{t_1} \int_V p \, dV dt = \int_{t_0}^{t_1} \int_0^Z \int_0^Y \int_{x_0}^{x_4} p \, dx \, dy \, dz \, dt = [t_1 - t_0] ZY \int_{x_0}^{x_4} p = TZY \int_{x_0}^{x_4} p \, dx. \quad (4.7)$$

Im Folgenden wird das innere Integral

$$\int_{x_0}^{x_4} \mathbf{u}^T \mathbf{i} \, dx = \int_{x_0}^{x_4} c_{\mu\mu}^x u_\mu D_x u_\mu \, dx + \int_{x_0}^{x_4} c_{\mu\nu}^x u_\mu D_x u_\nu \, dx + \int_{x_0}^{x_4} c_{\nu\mu}^x u_\nu D_x u_\mu \, dx + \int_{x_0}^{x_4} c_{\nu\nu}^x u_\nu D_x u_\nu \, dx \quad (4.8)$$

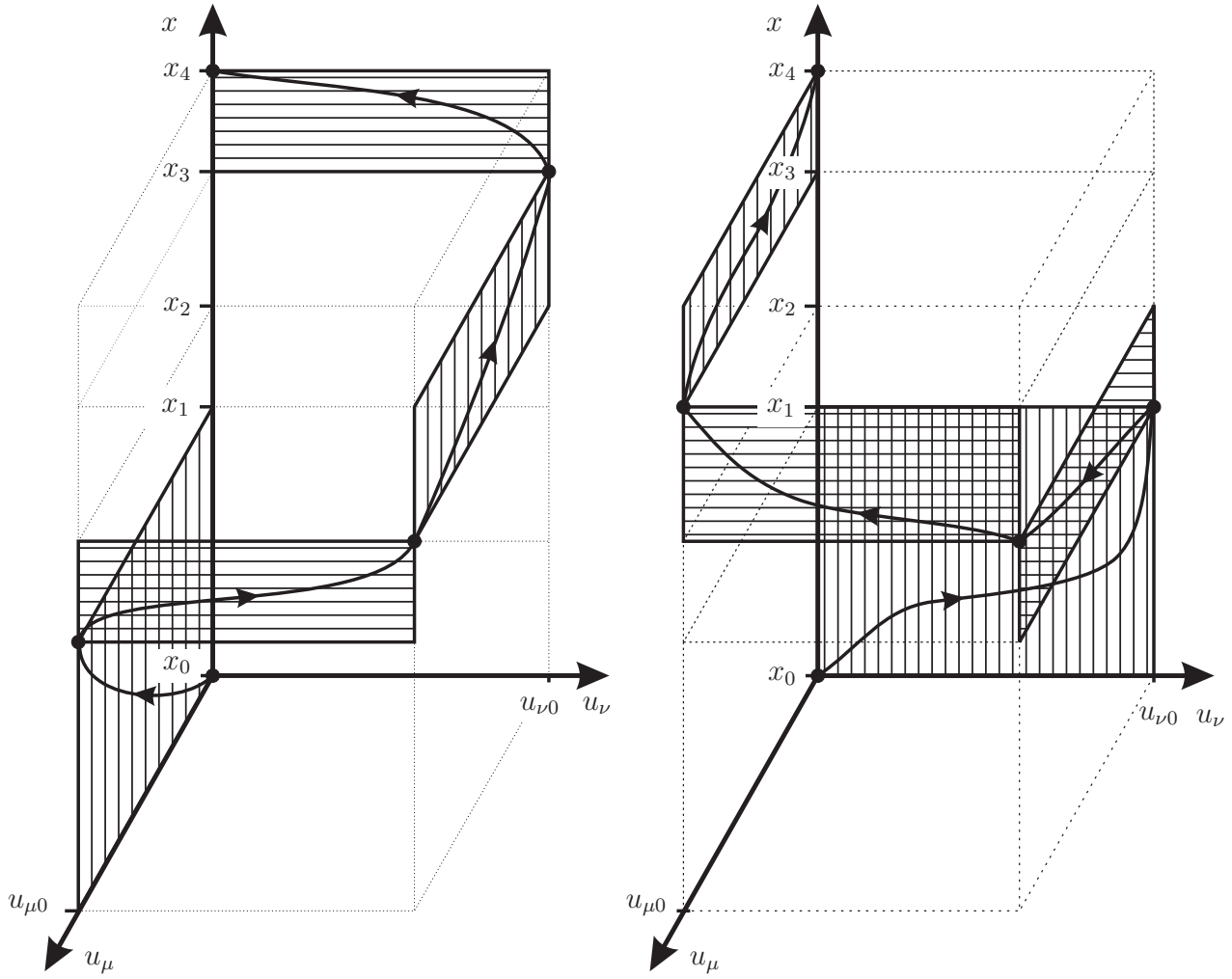


Bild 4.1: a) Erster Kreisprozess

b) Zweiter Kreisprozess

untersucht. Die Auswertung der Integrale über die einzelnen Teilintervalle ergibt

$$\begin{aligned}
 \int_{x_0}^{x_1} \mathbf{u}^T \mathbf{i} \, dx &= c_{\mu\mu}^x \frac{u_{\mu 0}^2}{2} \\
 \int_{x_1}^{x_2} \mathbf{u}^T \mathbf{i} \, dx &= c_{\mu\nu}^x u_{\mu 0} u_{\nu 0} + c_{\nu\nu}^x \frac{u_{\nu 0}^2}{2} \\
 \int_{x_2}^{x_3} \mathbf{u}^T \mathbf{i} \, dx &= -c_{\mu\mu}^x \frac{u_{\mu 0}^2}{2} - c_{\nu\mu}^x u_{\mu 0} u_{\nu 0} \\
 \int_{x_3}^{x_4} \mathbf{u}^T \mathbf{i} \, dx &= -c_{\nu\nu}^x \frac{u_{\nu 0}^2}{2} .
 \end{aligned} \tag{4.9}$$

Für den 1. speziellen Vorgang erhalten wir schließlich

$$\int_{x_0}^{x_4} p \, dx = \int_{x_0}^{x_4} \mathbf{u}^T \mathbf{i} \, dx = [c_{\mu\nu}^x - c_{\nu\mu}^x] u_{\mu 0} u_{\nu 0} \geq 0 . \tag{4.10}$$

Für den 2. Vorgang gilt $u_\sigma = 0$ und $i_\sigma = 0$ für $\mu \neq \sigma \neq \nu$. Die anderen beiden Spannungen und Ströme folgen dem ersten Kreisprozess in umgekehrter Richtung, siehe Bild 4.1 b).

Für den 2. speziellen Vorgang erhalten wir

$$\int_{x_0}^{x_4} p \, dx = \int_{x_0}^{x_4} \mathbf{u}^T \mathbf{i} \, dx = [c_{\nu\mu}^x - c_{\mu\nu}^x] u_{\mu 0} u_{\nu 0} \geq 0. \quad (4.11)$$

Dieses Ergebnis hätten wir auch durch Vertauschung von μ und ν ermitteln können, da Gleichung (4.6) für beliebige Vorgänge gelten muss und die ursprüngliche Bezeichnung ja willkürlich gewählt wurde.

Mit $t_1 > t_0$, $Z > 0$, $Y > 0$, $u_{\mu 0} > 0$ und $u_{\nu 0} > 0$ erhalten wir aus den beiden Vorgängen die simultan zu erfüllenden Ungleichungen

$$[c_{\nu\mu}^x - c_{\mu\nu}^x] \geq 0 \quad \text{und} \quad [c_{\mu\nu}^x - c_{\nu\mu}^x] \geq 0. \quad (4.12)$$

Die einzige Lösung ist $c_{\mu\nu}^x = c_{\nu\mu}^x$. Da μ und ν beliebig gewählt wurden, muss $c_{\mu\nu}^x = c_{\nu\mu}^x$ für alle μ und ν gelten. Somit muss \mathbf{C}^x symmetrisch sein. Durch ähnliches Vorgehen ergibt sich als notwendige Bedingung an \mathbf{C}^y , \mathbf{C}^z und \mathbf{C}^t ebenfalls die Symmetrie. Aufgrund der Symmetrie lassen sich die aufgenommenen Teilleistungsdichten durch die partielle Ableitung einer quadratischen Form ausdrücken

$$\mathbf{u}^T \mathbf{C}^\xi \mathbf{D}_\xi \mathbf{u} = \frac{1}{2} \mathbf{D}_\xi \{ \mathbf{u}^T \mathbf{C}^\xi \mathbf{u} \} \quad \text{für} \quad \xi = x, y, z \text{ und } t. \quad (4.13)$$

Die gesamte aufgenommene Leistungsdichte kann daher durch

$$p = \mathbf{u}^T \mathbf{i} = \frac{1}{2} [\mathbf{D}_x \mathbf{u}^T \mathbf{C}^x \mathbf{u} + \mathbf{D}_y \mathbf{u}^T \mathbf{C}^y \mathbf{u} + \mathbf{D}_z \mathbf{u}^T \mathbf{C}^z \mathbf{u} + \mathbf{D}_t \mathbf{u}^T \mathbf{C}^t \mathbf{u}] + \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \quad (4.14)$$

dargestellt werden. Die Energiedichte \mathbf{W}_s ist folglich durch

$$\mathbf{W}_s = \frac{1}{2} \left[\mathbf{u}^T \frac{\mathbf{C}^x}{v_4} \mathbf{u}, \mathbf{u}^T \frac{\mathbf{C}^y}{v_4} \mathbf{u}, \mathbf{u}^T \frac{\mathbf{C}^z}{v_4} \mathbf{u}, \mathbf{u}^T \mathbf{C}^t \mathbf{u} \right]^T \quad (4.15)$$

gegeben. Das Integral

$$\begin{aligned} \int_{\mathcal{G}_x} p \, dV dt &= v_4 \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{D}_x^T \mathbf{W}_s \, dV dt + \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \, dV dt \\ &= v_4 \int_{t_0}^{t_1} \int_{\mathcal{G}} \operatorname{div} \mathbf{W}_{\text{sp}} \, dV dt + \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{D}_t W_t \, dt \, dV + \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \, dV dt \\ &= v_4 \int_{t_0}^{t_1} \int_{\partial \mathcal{G}} \mathbf{W}_{\text{sp}}^T \mathbf{n} \, dA \, dt + \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{D}_t W_t \, dt \, dV + \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \, dV dt \\ &= \int_{\mathcal{G}} \left[\frac{1}{2} \mathbf{u}^T(t_1) \mathbf{C}^t \mathbf{u}(t_1) + \int_{t_0}^{t_1} \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \, dt \right] dV \end{aligned} \quad (4.16)$$

darf nach Gleichung (2.85) für beliebige Vorgänge \mathbf{u} und $t_1 \geq t_0$ nicht negativ werden. Dies ist nur dann der Fall, wenn die Matrizen \mathbf{C}^t und \mathbf{Y}^k n. n. d. sind.

Hinlänglichkeit

Nach Gleichung (2.78) müssen wir zeigen, dass die Differenz aus der über die Tore aufgenommenen Energie und der über die Gebietsgrenze aufgenommenen Energie nicht kleiner als null ist. Bilden dieser Differenz liefert

$$\int_{\mathcal{G}_x} p \, dV \, dt - v_4 \int_{\mathcal{G}_x} \mathbf{D}_x^T \mathbf{W}_s \, dV \, dt = \int_{\mathcal{G}} \int_{t_0}^{t_1} \mathbf{u}^T \mathbf{Y}^k \mathbf{u} \, dt \, dV. \quad (4.17)$$

Der Integrand ist aufgrund der nicht negativen Definitheit von \mathbf{Y}^k nicht negativ. Wegen $t_1 \geq t_0$ und der Tatsache, dass über ein Volumen integriert wird, ist die Differenz für alle Vorgänge nicht negativ. Damit ist der Satz bewiesen.

Abschließend wollen wir noch einmal die Ergebnisse zusammentragen. Notwendig und hinreichend für die Passivität sind symmetrische Matrizen \mathbf{C}^t , \mathbf{C}^x , \mathbf{C}^y und \mathbf{C}^z sowie nicht negativ definite Matrizen \mathbf{C}^t und \mathbf{Y}^k .

4.1.3 Systemklassifizierung

Ein partielles Differentialgleichungssystem Gleichung (4.3) wird nach [Frie54] symmetrisch hyperbolisches System bezeichnet, wenn die Matrizen \mathbf{C}^x , \mathbf{C}^y , \mathbf{C}^z und \mathbf{C}^t symmetrisch sind und eine von ihnen positiv definit ist.

Wie wir später sehen werden, benötigt das in diesem Kapitel vorgestellte Syntheseverfahren eine positiv definite Matrix \mathbf{C}^t . Es liegt daher nahe, die aus der Passivität resultierende Forderung an \mathbf{C}^t von der nicht negativen Definitheit auf die positive Definitheit zu verstärken. Das vorliegende partielle Differentialgleichungssystem gehört dann bei beliebigen symmetrischen Matrizen \mathbf{C}^x , \mathbf{C}^y und \mathbf{C}^z zur Klasse der symmetrisch hyperbolischen.

Wir machen noch eine weitere Einschränkung an die Matrizen \mathbf{C}^x , \mathbf{C}^y und \mathbf{C}^z . Wir fordern, dass deren Hauptdiagonalelemente verschwinden. Eine Synthese und anschließende Umsetzung in eine WD-Struktur ließe sich zwar auch bei nichtverschwindenden Hauptdiagonalen erreichen, aber bislang konnte keine Methode zur Randbehandlung dieser Anteile gefunden werden. Eine erwähnenswerte Eigenschaft einer reellen, symmetrischen Matrix \mathbf{M} mit verschwindender Hauptdiagonale ist es, indefinit zu sein (abgesehen von der Nullmatrix). Die Summe der Eigenwerte ist nämlich gleich der Spur der Matrix und diese verschwindet, d. h. $\text{spur } \mathbf{M} = \sum \lambda_\nu = 0$. Es existieren daher mindestens zwei (reelle) Eigenwerte mit unterschiedlichem Vorzeichen, d.h. die Matrix ist indefinit.

4.1.4 Eingeschwungener Zustand

Aufgrund der Linearität existiert bei Anregung der Form

$$\mathbf{j} = \text{Re}\{\mathbf{J} e^{\mathbf{p}_x^T \mathbf{x}}\}, \quad (4.18)$$

bis auf endlich viele Frequenzen \mathbf{p}_x , eine Lösung der Form

$$\mathbf{u} = \text{Re}\{\mathbf{U} e^{\mathbf{p}_x^T \mathbf{x}}\}, \quad (4.19)$$

wobei \mathbf{J} und \mathbf{U} Vektoren mit komplexen konstanten Koordinaten sind. Bei Betrachtung des stationären Zustands und Verwendung der komplexen Größen anstatt der zeit- und ortsabhängigen Größen tritt die Frequenzvariable \mathbf{p}_x an die Stelle des Differentialoperators \mathbf{D}_x . Im stationärem Zustand geht Gleichung (4.3) in

$$[\mathbf{C}^x p_x + \mathbf{C}^y p_y + \mathbf{C}^z p_z + \mathbf{C}^t p_t + \mathbf{Y}^k] \mathbf{U} = \mathbf{J} \quad (4.20)$$

über. Mit der Interpretation der Koordinaten von \mathbf{u} als Spannungen und der Koordinaten von \mathbf{j} als eingeprägte Ströme stellen

$$\mathbf{Y}^{cx}(p_x) = p_x \mathbf{C}^x, \mathbf{Y}^{cy}(p_y) = p_y \mathbf{C}^y, \mathbf{Y}^{cz}(p_z) = p_z \mathbf{C}^z \text{ und } \mathbf{Y}^{ct}(p_t) = p_t \mathbf{C}^t \quad (4.21)$$

Admittanzmatrizen dar. Der reaktive Teil besitzt die Admittanzmatrix

$$\mathbf{Y}^c(\mathbf{p}_x) = \mathbf{Y}^{cx}(p_x) + \mathbf{Y}^{cy}(p_y) + \mathbf{Y}^{cz}(p_z) + \mathbf{Y}^{ct}(p_t). \quad (4.22)$$

Die durch $\mathbf{Y}(\mathbf{p}_x)\mathbf{U} = \mathbf{J}$ festgelegte gesamte Admittanzmatrix $\mathbf{Y}(\mathbf{p}_x)$ ist die Summe aus dem reaktiven und dem konstanten Teil $\mathbf{Y}(\mathbf{p}_x) = \mathbf{Y}^c(\mathbf{p}_x) + \mathbf{Y}^k$.

Wir wollen nun untersuchen, wie die Einschränkungen an die partielle Differentialgleichung sich im Frequenzbereich auswirken und welche Eigenschaften daraus für die zuvor definierten Admittanzmatrizen resultieren. Da die partielle Differentialgleichung ein reelles System beschreiben soll, gilt $\mathbf{C}^x, \mathbf{C}^y, \mathbf{C}^z, \mathbf{C}^t, \mathbf{Y}^k \in \mathbb{R}^{N \times N}$. Folglich ist jedes Element der Admittanzmatrix gleich dem bikonjugierten Element, d. h. $\mathbf{Y}(\mathbf{p}_x) = \mathbf{Y}^*(\mathbf{p}_x^*)$. Die Admittanzmatrix \mathbf{Y}^c soll ein verlustfreies System beschreiben. Somit ist die Summe aus der Matrix und ihrer parakonjugierten Matrix gleich der Nullmatrix, d. h. $\mathbf{Y}^c(\mathbf{p}_x) + \mathbf{Y}_*^c(\mathbf{p}_x) = \mathbf{0}$, [Bele68], [Bose82].

Offenbar ist die Admittanzmatrix des reaktiven Teils eine ungerade Funktion in \mathbf{p}_x . Aus dieser Beziehung $\mathbf{Y}^c(\mathbf{p}_x) = -\mathbf{Y}^c(-\mathbf{p}_x)$, den Eigenschaften eines reellen Systems und den Eigenschaften eines verlustfreien Systems, folgt das schon aus dem Orts-/Zeitbereich bekannte Ergebnis, dass der reaktive Teil durch eine symmetrische Matrix beschrieben werden kann

$$\mathbf{Y}^c(\mathbf{p}_x) = \mathbf{Y}^{cT}(\mathbf{p}_x). \quad (4.23)$$

Die partielle Differentialgleichung soll ein extern passives, symmetrisch hyperbolisches System beschreiben. Im Koordinatensystem \mathbf{x} ist \mathbf{Y}^k eine positive Matrix in der Frequenzvariablen p_t , d. h. \mathbf{Y}^k ist analytisch und der symmetrische Teil ist n.n.d. für $\text{Re}\{p_t\} > 0$

$$\begin{aligned} \mathbf{Y}^k & \text{ analytisch in } \text{Re}\{p_t\} > 0 \\ \mathbf{Y}^k + \mathbf{Y}^{kT} & \geq 0 \quad \text{in } \text{Re}\{p_t\} > 0. \end{aligned} \quad (4.24)$$

Für \mathbf{Y}^{ct} gilt hingegen die schärfere Bedingung

$$\begin{aligned} \mathbf{Y}^{ct} & \text{ analytisch in } \text{Re}\{p_t\} > 0 \\ \mathbf{Y}^{ct} + \mathbf{Y}^{ctT} & > 0 \quad \text{in } \text{Re}\{p_t\} > 0. \end{aligned} \quad (4.25)$$

Die Matrizen $\mathbf{Y}^{cx}, \mathbf{Y}^{cy}$ und \mathbf{Y}^{cz} unterliegen bis auf die schon festgelegte Symmetrie keinen weiteren Restriktionen.

Im Koordinatensystem \mathbf{t} ergeben sich für \mathbf{Y} die folgenden notwendigen und hinreichenden Eigenschaften für die MD-Passivität. \mathbf{Y} ist eine positive Matrix in der Variablen \mathbf{p}_t , d. h. \mathbf{Y} ist analytisch und $\mathbf{Y} + \mathbf{Y}^H$ ist n.n.d. in der rechten offenen Polyhalbebene $\text{Re}\{\mathbf{p}_t\} > \mathbf{0}$, d. h.

$$\begin{aligned} \mathbf{Y} & \text{ analytisch in } \text{Re}\{\mathbf{p}_t\} > \mathbf{0} \\ \mathbf{Y} + \mathbf{Y}^H & \geq 0 \quad \text{in } \text{Re}\{\mathbf{p}_t\} > \mathbf{0}. \end{aligned} \quad (4.26)$$

Die letzte Eigenschaft setzt allerdings die Wahl einer geeigneten Konstanten v_4 voraus. Wir werden später auf die Wahl der Konstanten und den Beweis der Äquivalenz der Aussagen in den beiden Koordinatensystemen zu sprechen kommen.

4.2 Synthese der Referenzschaltung

4.2.1 Vorbetrachtungen

Jede quadratische Matrix \mathbf{M} mit den Elementen $m_{\mu\nu}$ und N Zeilen kann durch

$$\mathbf{M} = \sum_{\mu=1}^N \sum_{\nu=1}^N \mathbf{e}_\mu \mathbf{e}_\nu^T m_{\mu\nu} \quad (4.27)$$

dargestellt werden. Aus der Summe spalten wir die Summe mit den Hauptdiagonalelementen ab

$$\mathbf{M} = \sum_{\mu=1}^N \mathbf{e}_{\mu} \mathbf{e}_{\mu}^T m_{\mu\mu} + \sum_{\mu=1}^N \sum_{\substack{\nu=1 \\ \nu \neq \mu}}^N \mathbf{e}_{\mu} \mathbf{e}_{\nu}^T m_{\mu\nu}. \quad (4.28)$$

Den zweiten Teil zerlegen wir in die Summe einer strikten linken unteren Dreiecksmatrix und einer strikten rechten oberen Dreiecksmatrix

$$\sum_{\mu=1}^N \sum_{\substack{\nu=1 \\ \nu \neq \mu}}^N \mathbf{e}_{\mu} \mathbf{e}_{\nu}^T m_{\mu\nu} = \underbrace{\sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{e}_{\mu} \mathbf{e}_{\nu}^T m_{\mu\nu}}_{\text{strikte linke untere Dreiecksmatrix}} + \underbrace{\sum_{\nu=2}^N \sum_{\mu=1}^{\nu-1} \mathbf{e}_{\mu} \mathbf{e}_{\nu}^T m_{\mu\nu}}_{\text{strikte rechte obere Dreiecksmatrix}}. \quad (4.29)$$

Die innere Summe der linken unteren Dreiecksmatrix hat $\mu - 1$ Summanden. Somit lautet die Zahl der Summanden der Doppelsumme

$$\sum_{\mu=2}^N \mu - 1 = \sum_{\mu=1}^{N-1} \mu = \frac{N-1}{2} [(N-1) + 1] = \frac{N^2 - N}{2} = \frac{N(N-1)}{2} = \frac{N!}{2!(N-2)!} = \binom{N}{2}. \quad (4.30)$$

Dies entspricht der Hälfte aus der Differenz der Zahl der Matrixelemente und der Zahl der Hauptdiagonalelemente.

Wir nehmen im Folgenden eine symmetrische Matrix \mathbf{M} an. Die strikte rechte obere Dreiecksmatrix formen wir unter Ausnutzung von $m_{\mu\nu} = m_{\nu\mu}$ und Tauschen der (bzgl. der Summen lokalen) Indizes um, d.h.

$$\sum_{\nu=2}^N \sum_{\mu=1}^{\nu-1} \mathbf{e}_{\mu} \mathbf{e}_{\nu}^T m_{\mu\nu} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{e}_{\nu} \mathbf{e}_{\mu}^T m_{\mu\nu}. \quad (4.31)$$

Die Matrix \mathbf{M} lautet dann

$$\mathbf{M} = \sum_{\mu=1}^N \mathbf{e}_{\mu} \mathbf{e}_{\mu}^T m_{\mu\mu} + \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_{\mu} \mathbf{e}_{\nu}^T + \mathbf{e}_{\nu} \mathbf{e}_{\mu}^T] m_{\mu\nu}. \quad (4.32)$$

Wir werden im Folgenden die Summe der Dyaden $\mathbf{e}_{\mu} \mathbf{e}_{\nu}^T$ und $\mathbf{e}_{\nu} \mathbf{e}_{\mu}^T$ aufbereiten

$$\mathbf{e}_{\mu} \mathbf{e}_{\nu}^T + \mathbf{e}_{\nu} \mathbf{e}_{\mu}^T = \begin{bmatrix} \mathbf{e}_{\mu} & \mathbf{e}_{\nu} \end{bmatrix} \begin{bmatrix} \mathbf{e}_{\nu}^T \\ \mathbf{e}_{\mu}^T \end{bmatrix} = \begin{bmatrix} \mathbf{e}_{\mu} & \mathbf{e}_{\nu} \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{e}_{\mu}^T \\ \mathbf{e}_{\nu}^T \end{bmatrix}. \quad (4.33)$$

Wir definieren nun eine Matrix $\mathbf{T}_{\mu\nu} = [\mathbf{e}_{\mu} \mathbf{e}_{\nu}]$, sodass

$$\mathbf{e}_{\mu} \mathbf{e}_{\nu}^T + \mathbf{e}_{\nu} \mathbf{e}_{\mu}^T = \mathbf{T}_{\mu\nu} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \mathbf{T}_{\mu\nu}^T \quad (4.34)$$

gilt. In diesem Zusammenhang nutzen wir die folgende Zerlegung, die einer Eigenwertzerlegung ähnlich ist

$$2 \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}}_{\mathbf{S}_0} \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad (4.35)$$

sodass

$$\mathbf{e}_\mu \mathbf{e}_\nu^T + \mathbf{e}_\nu \mathbf{e}_\mu^T = \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} -1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T \quad (4.36)$$

resultiert. Die mit \mathbf{S}_0 definierte Hadamard-Matrix wird im weiteren Verlauf noch eine große Rolle spielen. Weiterhin werden wir später die leicht verifizierbare Beziehung der speziellen Diagonalmatrix

$$\mathbf{e}_\mu \mathbf{e}_\mu^T + \mathbf{e}_\nu \mathbf{e}_\nu^T = \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T \quad (4.37)$$

benötigen.

4.2.2 Synthese des reaktiven Teils

Ziel der weiteren Überlegungen ist es, die Matrix \mathbf{Y}^c als Knotenleitwertmatrix einer mehrdimensional konkret passiven Schaltung darzustellen, d.h.

$$\mathbf{Y}^c = \mathbf{A} \mathbf{Y}_D^c \mathbf{A}^T. \quad (4.38)$$

Hierin ist \mathbf{Y}_D^c eine Diagonalmatrix mit den Elementen $p_\kappa C_\sigma$, $\kappa = 1, 2, \dots, 7$, $C_\sigma \geq 0$. \mathbf{A} ist die Knoten-Zweig-Inzidenzmatrix eines verallgemeinerten Verbindungsnetzes. Es wird sich später herausstellen, dass das verallgemeinerte Verbindungsnetz mittels eines Verbindungsnetzes und Zweitorübertragern mit Übersetzungsverhältnis $-1/1$ realisiert werden kann. Da die Elemente von \mathbf{Y}_D^c positive Funktionen sind, stellt \mathbf{Y}_D^c eine positive Matrix dar. Weiterhin ist sichergestellt, dass \mathbf{Y}^c eine positive Matrix ist, weil die nicht negative Definitheit invariant bzgl. der Kongruenztransformation ist.

Unter den gemachten Voraussetzungen lautet die Matrix \mathbf{Y}^{cx}

$$\mathbf{Y}^{cx} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_\mu \mathbf{e}_\nu^T + \mathbf{e}_\nu \mathbf{e}_\mu^T] p_x c_{\mu\nu}^x. \quad (4.39)$$

Mit den Vorbemerkungen kann nun \mathbf{Y}^{cx} durch Nutzen von Gleichung (4.36) und $c_{\mu\nu}^x = |c_{\mu\nu}^x| \operatorname{sgn}(c_{\mu\nu}^x)$ wie folgt dargestellt werden

$$\mathbf{Y}^{cx} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} -p_x \frac{|c_{\mu\nu}^x| \operatorname{sgn}(c_{\mu\nu}^x)}{2} & 0 \\ 0 & p_x \frac{|c_{\mu\nu}^x| \operatorname{sgn}(c_{\mu\nu}^x)}{2} \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T. \quad (4.40)$$

Nun definieren wir die Diagonalmatrizen $\mathbf{Y}^{ct1,4}$ und $\mathbf{C}^{t1,4}$

$$\mathbf{Y}^{ct1,4} = p_t \mathbf{C}^{t1,4} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_\mu \mathbf{e}_\mu^T + \mathbf{e}_\nu \mathbf{e}_\nu^T] \frac{p_t}{v_4} |c_{\mu\nu}^x| = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T \frac{p_t}{v_4} |c_{\mu\nu}^x|, \quad (4.41)$$

wobei die letzte Identität über Gleichung (4.37) zu gewinnen ist². Die Matrix $\mathbf{Y}^{ct1,4}$ wird nun zu \mathbf{Y}^{cx} aus Gleichung (4.40) addiert

$$\mathbf{Y}^{c1,4} = \mathbf{Y}^{cx} + \mathbf{Y}^{ct1,4} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} \frac{p_t}{v_4} - p_x \operatorname{sgn}(c_{\mu\nu}^x) & 0 \\ 0 & \frac{p_t}{v_4} + p_x \operatorname{sgn}(c_{\mu\nu}^x) \end{bmatrix} \frac{|c_{\mu\nu}^x|}{2} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T. \quad (4.42)$$

²Die Indizes 1 und 4 stellen einen Bezug zu den Koordinaten 1 und 4 der mehrdimensionalen Zeit her.

Mögliche Darstellungen von $p_t/v_4 - p_x$ und $p_t/v_4 + p_x$ sind $p_t/v_4 - p_x = p_4/v_0$ bzw. $p_t/v_4 + p_x = p_1/v_0$. Diese Darstellungen entsprechen der Wahl zweier Spaltenvektoren einer Rechtsinversen von \mathbf{H} . Damit haben wir in den neuen Koordinaten

$$\mathbf{Y}^{c1,4} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{|c_{\mu\nu}^x|}{2v_0} & 0 \\ 0 & \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{|c_{\mu\nu}^x|}{2v_0} \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T. \quad (4.43)$$

Hier und im Folgenden sind die geschweiften Klammern so zu interpretieren, dass die obere (untere) Alternative zu nehmen ist, wenn die Signum-Funktion $\text{sgn}(c_{\mu\nu}^x)$ positiv (negativ) ist. Eine formale Darstellung kann mithilfe der Einheitssprungfunktion erreicht werden

$$\left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} = p_4 u(c_{\mu\nu}^x) + p_1 u(-c_{\mu\nu}^x). \quad (4.44)$$

Wir definieren

$$C_{\mu\nu}^\xi = \frac{2|c_{\mu\nu}^\xi|}{v_0} \quad \text{für} \quad \xi = x, y, z \quad (4.45)$$

die später Kapazitäten von Kondensatoren sind. Dann haben wir

$$\mathbf{Y}^{c1,4} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} & 0 \\ 0 & \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T, \quad (4.46)$$

vgl. dazu Gleichung (2.160). Ebenso verfährt man mit den anderen Richtungen, die die Matrizen $\mathbf{Y}^{c2,5}$ und $\mathbf{Y}^{c3,6}$ definieren. Da wir die ursprüngliche Matrix \mathbf{Y}^c erhalten wollen, erfolgt eine Kompensation der Additionen der Matrix \mathbf{Y}^{ct} , und zwar wie folgt

$$\mathbf{Y}^c = \mathbf{Y}^{c1,4} + \mathbf{Y}^{c2,5} + \mathbf{Y}^{c3,6} + \mathbf{Y}^{ct} - \underbrace{\sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_\mu \mathbf{e}_\mu^T + \mathbf{e}_\nu \mathbf{e}_\nu^T] \frac{p_t}{v_4} [|c_{\mu\nu}^x| + |c_{\mu\nu}^y| + |c_{\mu\nu}^z|]}_{\mathbf{Y}^{c7} = p_7 \mathbf{C}^7 = p_t \frac{v_0}{v_4} \mathbf{C}^7}. \quad (4.47)$$

Durch Wahl einer hinreichend großen Konstanten v_4 kann die nicht negative Definitheit von

$$\mathbf{C}^7 = \frac{v_4}{v_0} \mathbf{C}^t - \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_\mu \mathbf{e}_\mu^T + \mathbf{e}_\nu \mathbf{e}_\nu^T] \frac{1}{v_0} [|c_{\mu\nu}^x| + |c_{\mu\nu}^y| + |c_{\mu\nu}^z|], \quad (4.48)$$

die für die MD-Passivität notwendig und hinreichend ist, erreicht werden. Zum Beweis berücksichtigen wir, dass die Eigenwerte einer symmetrischen reellen p.d. Matrix \mathbf{C}^t reell und positiv sind. Die Eigenwerte einer symmetrischen reellen n.n.d. Matrix

$$\mathbf{C} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} [\mathbf{e}_\mu \mathbf{e}_\mu^T + \mathbf{e}_\nu \mathbf{e}_\nu^T] \frac{1}{v_0} [|c_{\mu\nu}^x| + |c_{\mu\nu}^y| + |c_{\mu\nu}^z|] \quad (4.49)$$

sind reell und nicht negativ. Weiterhin besitzt der Rayleigh-Quotient einer reellen symmetrischen Matrix als untere (obere) Schranke ihren kleinsten (größten) Eigenwert. Der kleinste (größte) Eigenwert der Matrix \mathbf{C}^t (\mathbf{C}) soll $\lambda_{\min}^{C^t}$ (λ_{\max}^C) sein. Dann gilt für $\mathbf{x} \neq \mathbf{0}$

$$\frac{1}{\lambda_{\min}^{C^t}} \frac{\mathbf{x}^H \mathbf{C}^t \mathbf{x}}{\|\mathbf{x}\|^2} \geq 1 \geq \frac{1}{\lambda_{\max}^C} \frac{\mathbf{x}^H \mathbf{C} \mathbf{x}}{\|\mathbf{x}\|^2} \iff \mathbf{x}^H \left[\frac{\lambda_{\max}^C}{\lambda_{\min}^{C^t}} \mathbf{C}^t - \mathbf{C} \right] \mathbf{x} \geq 0. \quad (4.50)$$

Wir müssen somit

$$v_4 \geq v_0 \frac{\lambda_{\max}^C}{\lambda_{\min}^{C^t}} \quad (4.51)$$

wählen um die mehrdimensionale Passivität zu garantieren.

Die gesamte Admittanzmatrix des reaktiven Teils kann nun in Abhängigkeit der Frequenzvariablen des neuen Koordinatensystems durch

$$\mathbf{Y}^c = \sum_{\kappa=1}^7 \mathbf{C}^\kappa p_\kappa \quad (4.52)$$

ausgedrückt werden.

Durch Darstellen der Doppelsumme in Gleichung (4.46) mittels eines Skalarprodukts

$$\mathbf{Y}^{c1,4} = \underbrace{[\mathbf{T}_{21}\mathbf{S}_0, \dots, \mathbf{T}_{N,N-1}\mathbf{S}_0]}_{\mathbf{A}} \underbrace{\text{diag} \left(\left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{C_{21}^x}{4}, \dots, \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{C_{N,N-1}^x}{4} \right)}_{\mathbf{Y}_D^{c1,4}} [\mathbf{T}_{21}\mathbf{S}_0, \dots, \mathbf{T}_{N,N-1}\mathbf{S}_0]^T \quad (4.53)$$

haben wir für die Matrix $\mathbf{Y}^{c1,4}$ die Struktur als Knotenleitwertmatrix gewonnen, die für die weiteren Überlegungen günstig ist. In der Regel sind \mathbf{C}^x , \mathbf{C}^y , \mathbf{C}^z , \mathbf{C}^t und \mathbf{Y}^k lichte Matrizen. Aus diesem Grunde kann eine Reduktion der Matrizen durch Streichen der verschwindenden Einträge in $\mathbf{Y}_D^{c1,4}$ durchgeführt werden, d. h.

$$\mathbf{Y}^{c1,4} = \mathbf{A}^x \mathbf{Y}_{D,\text{red}}^{c1,4} \mathbf{A}^{xT}, \quad (4.54)$$

mit regulärer Diagonalmatrix $\mathbf{Y}_{D,\text{red}}^{c1,4}$. Wir werden nun die Matrix des verallgemeinerten Verbindungsnetzes wie folgt aufteilen

$$\mathbf{A} = [\mathbf{T}_{21}, \dots, \mathbf{T}_{N,N-1}] \text{diag}(\mathbf{S}_0, \dots, \mathbf{S}_0) = \mathbf{A}_{\text{rvn}} \text{diag}(\mathbf{S}_0, \dots, \mathbf{S}_0). \quad (4.55)$$

Die Matrix \mathbf{A}_{rvn} hat N Zeilen, $N^2 - N$ Spalten und ist Knoten-Zweig-Inzidenzmatrix eines topologischen Verbindungsnetzes. (Da $\mathbf{Y}^{c1,4}$ als Summe von Matrizen ausdrückbar ist, handelt es sich um eine Parallelschaltung.) Die zweite Matrix $\text{diag}\{\mathbf{S}_0, \dots, \mathbf{S}_0\}$ wird als das Übersetzungsverhältnis von Jaumann-Adaptoren interpretiert, die wiederum selber aufgrund der speziellen Struktur des Übersetzungsverhältnisses durch $\binom{N}{2}$ 2-Tor-Übertrager mit Übersetzungsverhältnis $-1/1$ realisiert werden können. Um dies zu erläutern stellen wir $\mathbf{Y}^{c1,4}$ aus Gleichung (4.46) durch

$$\mathbf{Y}^{c1,4} = \mathbf{A}_{\text{rvn}} \text{diag}(\mathbf{Y}_{21}^{c1,4}, \dots, \mathbf{Y}_{N,N-1}^{c1,4}) \mathbf{A}_{\text{rvn}}^T \quad \text{mit} \quad \mathbf{Y}_{\mu\nu}^{c1,4} = \mathbf{S}_0 \begin{bmatrix} \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} & 0 \\ 0 & \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} \end{bmatrix} \mathbf{S}_0^T \quad (4.56)$$

dar. Die Admittanzmatrix $\mathbf{Y}_{\mu\nu}^{c1,4}$ entspricht der Admittanzmatrix aus Gleichung (2.160), bei Wahl von $\frac{1}{Z_3} = \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} C_{\mu\nu}^x$, $\frac{1}{Z_4} = \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} C_{\mu\nu}^x$, $n = 1$ und $\mathbf{S}_0 = 2\mathbf{N}^T$. Die Realisierung der Jaumann-Schaltung ist im Bild 2.21 angegeben.

Es verbleibt, eine Realisierung der Matrix \mathbf{Y}^{c7} zu finden. Wir suchen ebenfalls eine Zerlegung der Form

$$\mathbf{Y}^{c7} = \mathbf{A}^7 \mathbf{Y}_D^{c7} \mathbf{A}^{7T}. \quad (4.57)$$

Um eine Aufwandsreduzierung zu erreichen ist es sinnvoll, einen möglicherweise vorhandenen Diagonaltteil der Matrix \mathbf{Y}^{c7} abzuspalten. Zu diesem Zweck sortieren wir die Zeilen und Spalten \mathbf{C}^7 mithilfe der Permutationsmatrix \mathbf{P}^7 zu

$$\mathbf{C}^7 = \mathbf{P}^7 \begin{bmatrix} \bar{\mathbf{C}}^7 & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{C}}^7 \end{bmatrix} \mathbf{P}^{7T} \quad (4.58)$$

um, wobei $\tilde{\mathbf{C}}^7$ eine Diagonalmatrix mit Kapazitäten ist. Die verbleibende Matrix $\bar{\mathbf{C}}^7$ hat die Dimension $\bar{n} \times \bar{n}$. Da $\bar{\mathbf{C}}^7$ reell und symmetrisch ist, existiert eine zugehörige reelle, orthogonale Matrix $\mathbf{N} = \mathbf{N}^{-T}$ der Eigenwertzerlegung

$$\bar{\mathbf{C}}^7 = \mathbf{N}^T \text{diag}(\bar{C}_1, \dots, \bar{C}_{\bar{n}}) \mathbf{N} = \mathbf{N}^T \bar{\mathbf{C}}_D^7 \mathbf{N}. \quad (4.59)$$

Hierbei sind aufgrund der nicht negativen Definitheit die Eigenwerte nicht negativ, $\bar{C}_\nu \geq 0$. Die Eigenwerte sind Kapazitäten von Kondensatoren. Wir fassen nun \mathbf{N} als das Übersetzungsverhältnis eines idealen $2\bar{n}$ -Tor Übertragers gemäß Gleichung (2.113) auf, der mit \bar{n} idealen Kondensatoren abgeschlossen ist, d. h. es gilt $\mathbf{I}_m = -p_7 \bar{\mathbf{C}}^7 \mathbf{U}_m$. Mit $\mathbf{I}_r = -\mathbf{N}^T \mathbf{I}_m$ und $\mathbf{U}_m = \mathbf{N} \mathbf{U}_r$ haben wir

$$\mathbf{I}_r = \mathbf{N}^T p_7 \bar{\mathbf{C}}^7 \mathbf{N} \mathbf{U}_r \quad (4.60)$$

und somit eine geeignete Synthese von $\bar{\mathbf{Y}}^{c7} = p_7 \bar{\mathbf{C}}^7$ gefunden. Kompakt lässt sich \mathbf{Y}^{c7} durch

$$\mathbf{C}^7 = \mathbf{P}^7 \begin{bmatrix} \mathbf{N}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \text{diag}(\bar{\mathbf{C}}_D^7, \tilde{\mathbf{C}}^7) \begin{bmatrix} \mathbf{N} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \mathbf{P}^{7T} \quad (4.61)$$

beschreiben. Hierin ist

$$\tilde{\mathbf{N}} = \begin{bmatrix} \mathbf{N} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \mathbf{P}^{7T} \quad (4.62)$$

das Übersetzungsverhältnis eines idealen Übertragers.

Ziel unserer Überlegungen ist es, eine konkret mehrdimensional passive Schaltung für \mathbf{Y}^c zu finden. Aus den erarbeiteten Darstellungen können wir zwei wesentliche Resultate festhalten. Zum einen haben wir eine mehrdimensional konkret passive Realisierung für $\mathbf{Y}^{c1,4}$, $\mathbf{Y}^{c2,5}$ und $\mathbf{Y}^{c3,6}$ gefunden. Zum anderen haben wir für eine n. n. d. Matrix \mathbf{C}^7 eine mehrdimensional konkret passive Realisierung von \mathbf{Y}^{c7} angegeben. Wir haben gezeigt, dass \mathbf{C}^7 durch Wahl eines genügend großen v_4 n. n. d. ist.

Es folgt der nachzutragende Beweis für die im Kapitel 4.1.4 gemachte Aussage in Gestalt von Gleichung (4.26). Zur Beweisvorbereitung gehen wir von

$$\mathbf{C}^x = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} -1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^{c^x} \quad (4.63)$$

und

$$\mathbf{C}^{t1,4} = \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu} \frac{|c_{\mu\nu}^x|}{v_4} \quad (4.64)$$

aus. Damit können wir durch

$$\mathbf{Y}^{c1,4} = p_x \mathbf{C}^x + p_t \mathbf{C}^{t1,4} = \frac{1}{2v_0} [p_1 - p_4] \mathbf{C}^x + \frac{v_4}{2v_0} [p_1 + p_4] \mathbf{C}^{t1,4} = p_1 \underbrace{\left[\frac{v_4 \mathbf{C}^{t1,4}}{2v_0} + \frac{\mathbf{C}^x}{2v_0} \right]}_{\mathbf{C}^1} + p_4 \underbrace{\left[\frac{v_4 \mathbf{C}^{t1,4}}{2v_0} - \frac{\mathbf{C}^x}{2v_0} \right]}_{\mathbf{C}^4}$$

(4.65)

die Trennung nach \mathbf{C}^1 und \mathbf{C}^4 durchführen. \mathbf{C}^1 berechnet sich zu

$$\mathbf{C}^1 = \frac{1}{4v_0} \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} |c_{\mu\nu}^x| - c_{\mu\nu}^x & 0 \\ 0 & |c_{\mu\nu}^x| + c_{\mu\nu}^x \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T \quad (4.66)$$

und ist n. n. d., da die Diagonalmatrix keine negativen Elemente enthält und die Kongruenztransformation die nicht negative Definitheit nicht beeinflusst. Offenbar ist auch

$$\mathbf{C}^4 = \frac{1}{4v_0} \sum_{\mu=2}^N \sum_{\nu=1}^{\mu-1} \mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} |c_{\mu\nu}^x| + c_{\mu\nu}^x & 0 \\ 0 & |c_{\mu\nu}^x| - c_{\mu\nu}^x \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T \quad (4.67)$$

nicht negativ definit. Ebenso können wir zeigen, dass \mathbf{C}^2 , \mathbf{C}^3 , \mathbf{C}^5 und \mathbf{C}^6 n. n. d. sind.

Nach diesen Vorbetrachtungen kommen wir nun zum eigentlichen Beweis. Wir zeigen zunächst

$$[\mathbf{Y} + \mathbf{Y}^H \geq 0 \text{ für } \operatorname{Re}\{\mathbf{p}_t\} > \mathbf{0}] \implies [\mathbf{Y}^{ct} + \mathbf{Y}^{ctH} > 0 \text{ für } \operatorname{Re}\{p_t\} > 0] . \quad (4.68)$$

Wir gehen dazu von der verallgemeinerten Wirkleistung aus

$$2P = \operatorname{Re}\{\mathbf{I}^H \mathbf{U}\} = \mathbf{U}^H \frac{1}{2} [\mathbf{Y} + \mathbf{Y}^H] \mathbf{U} = \sum_{\kappa=1}^7 \mathbf{U}^H \mathbf{C}^\kappa \operatorname{Re}\{p_\kappa\} \mathbf{U} + \mathbf{U}^H \frac{1}{2} [\mathbf{Y}^k + \mathbf{Y}^{kT}] \mathbf{U} \geq 0 . \quad (4.69)$$

Da die Ungleichung für alle $\operatorname{Re}\{\mathbf{p}_t\} > \mathbf{0}$ gilt, folgt aus der Prämisse insbesondere die nicht negative Definitheit von \mathbf{C}^7 . Die Matrix $\mathbf{C}^7 = \frac{v_4}{v_0} \mathbf{C}^t - \sum_{\kappa=1}^6 \mathbf{C}^\kappa$ ist also für beliebige symmetrische n. n. d. \mathbf{C}^κ nicht negativ definit. Dies kann nur dann der Fall sein, wenn \mathbf{C}^t p.d. ist, die nicht negative Definitheit von \mathbf{C}^t reicht nicht aus. Denn bei nur n. n. d. Matrix \mathbf{C}^t ist nicht auszuschließen, dass ein Vektor $\mathbf{v} \neq \mathbf{0}$ existiert für den

$$\mathbf{v}^T \mathbf{C}^7 \mathbf{v} = \underbrace{\mathbf{v}^T \mathbf{C}^t \mathbf{v}}_{=0} - \frac{v_0}{v_4} \sum_{\kappa=1}^6 \mathbf{v}^T \mathbf{C}^\kappa \mathbf{v} < 0 \quad (4.70)$$

gilt und somit \mathbf{C}^7 der unmittelbar aus der Prämisse abgeleiteten Aussage widerspräche.

Eine weitere Darstellung der verallgemeinerten Wirkleistung erhalten wir mit

$$\hat{\mathbf{W}}_s = [\mathbf{U}^H \mathbf{C}^1 \mathbf{U}, \mathbf{U}^H \mathbf{C}^2 \mathbf{U}, \dots, \mathbf{U}^H \mathbf{C}^7 \mathbf{U}]^T \quad (4.71)$$

zu

$$\begin{aligned} 2P &= \operatorname{Re}\{\mathbf{p}_t^T\} \hat{\mathbf{W}}_s + \mathbf{U}^H \frac{1}{2} [\mathbf{Y}^k + \mathbf{Y}^{kT}] \mathbf{U} \\ &= \operatorname{Re}\{\mathbf{p}_x^T\} \mathbf{H} v_0 \hat{\mathbf{W}}_s + \mathbf{U}^H \frac{1}{2} [\mathbf{Y}^k + \mathbf{Y}^{kT}] \mathbf{U} = \operatorname{Re}\{\mathbf{p}_x^T\} v_0 \begin{bmatrix} \mathbf{U}^H [\mathbf{C}^1 - \mathbf{C}^4] \mathbf{U} \\ \mathbf{U}^H [\mathbf{C}^2 - \mathbf{C}^5] \mathbf{U} \\ \mathbf{U}^H [\mathbf{C}^3 - \mathbf{C}^6] \mathbf{U} \\ \mathbf{U}^H \sum_{\nu=1}^7 \mathbf{C}^\nu \mathbf{U} \end{bmatrix} + \mathbf{U}^H \frac{1}{2} [\mathbf{Y}^k + \mathbf{Y}^{kT}] \mathbf{U} . \end{aligned} \quad (4.72)$$

Drücken wir nun wieder \mathbf{C}^1 bis \mathbf{C}^7 durch \mathbf{C}^x , \mathbf{C}^y , \mathbf{C}^z und \mathbf{C}^t aus, so erhalten wir

$$2P = \mathbf{U}^H \mathbf{C}^x \operatorname{Re}\{p_x\} \mathbf{U} + \mathbf{U}^H \mathbf{C}^y \operatorname{Re}\{p_y\} \mathbf{U} + \mathbf{U}^H \mathbf{C}^z \operatorname{Re}\{p_z\} \mathbf{U} + \mathbf{U}^H \mathbf{C}^t \operatorname{Re}\{p_t\} \mathbf{U} + \mathbf{U}^H \frac{1}{2} [\mathbf{Y}^k + \mathbf{Y}^{kT}] \mathbf{U} .$$

(4.73)

Nun zeigen wir

$$\exists v_4 \quad [\mathbf{Y}^{ct} + \mathbf{Y}^{ctH} > 0 \text{ für } \operatorname{Re}\{p_t\} > 0] \implies [\mathbf{Y} + \mathbf{Y}^H \geq 0 \text{ für } \operatorname{Re}\{\mathbf{p}_t\} > \mathbf{0}] . \quad (4.74)$$

Wie bereits nachgewiesen wurde, folgt aus $\mathbf{Y}^{ct} + \mathbf{Y}^{ctH} > 0$ für $\operatorname{Re}\{p_t\} > 0$ eine positiv definite Matrix \mathbf{C}^t . Die nicht negative Definitheit von \mathbf{C}^1 bis \mathbf{C}^6 ist aus Gleichung (4.66) und Gleichung (4.67) ersichtlich. Aus Gleichung (4.52) i. V. m. Gleichung (4.64) folgt zudem, dass bei hinreichend großer Wahl von v_4 die Matrix \mathbf{C}^7 n. n. d. ist. Mit

$$\mathbf{Y} = \sum_{\nu=1}^7 \mathbf{C}^\nu p_\nu + \mathbf{Y}^k \quad (4.75)$$

folgt

$$[\mathbf{Y} + \mathbf{Y}^H] = 2 \sum_{\nu=1}^7 \mathbf{C}^\nu \operatorname{Re}\{p_\nu\} + [\mathbf{Y}^k + \mathbf{Y}^{kT}] . \quad (4.76)$$

Da nach Voraussetzung $\operatorname{Re}\mathbf{p}_t > \mathbf{0}$ gilt und die Matrizen in der Summe n. n. d. sind, ist die gesamte Summe n. n. d., was zu zeigen war.

4.2.3 Synthese des konstanten Teils

Zunächst werden wir wie auch im Fall der Matrix \mathbf{Y}^{c7} einen schon vorhandenen Diagonaleil der Matrix \mathbf{Y}^k abspalten. Zu diesem Zweck sortieren wir die Zeilen und Spalten von \mathbf{Y}^k mithilfe der Permutationsmatrix \mathbf{P}^k um

$$\mathbf{Y}^k = \mathbf{P}^k \begin{bmatrix} \bar{\mathbf{Y}}^k & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{Y}}^k \end{bmatrix} \mathbf{P}^{kT} , \quad (4.77)$$

wobei $\bar{\mathbf{Y}}^k$ eine reelle Matrix der Dimension $\tilde{n} \times \tilde{n}$ ist und einen n. n. d. symmetrischen Teil besitzt.

Die verbleibende reelle Diagonalmatrix $\tilde{\mathbf{Y}}^k$ hat dann die Dimension $N - \tilde{n} \times N - \tilde{n}$ und ist ebenfalls n. n. d.. Die Matrix $\bar{\mathbf{Y}}^k$ wird in einen symmetrischen und einen schiefsymmetrischen Teil zerlegt

$$\bar{\mathbf{Y}}^k = \underbrace{[\bar{\mathbf{Y}}^k + \bar{\mathbf{Y}}^{kT}]}_{\mathbf{Y}_R} \frac{1}{2} + \underbrace{[\bar{\mathbf{Y}}^k - \bar{\mathbf{Y}}^{kT}]}_{\mathbf{Y}_G} \frac{1}{2} . \quad (4.78)$$

Der Index R steht für einen Widerstand und der Index G für einen Gyrator. Der symmetrische Teil \mathbf{Y}_R kann mittels eines mit positiven Widerständen abgeschlossenen Mehrtorübertragers synthetisiert werden. Dieses Vorgehen entspricht dem, welches zur Synthese von $\bar{\mathbf{C}}^7$ angewendet wurde. Der schiefsymmetrische Teil \mathbf{Y}_G kann durch Gyratoren synthetisiert werden. Aufgrund der Addition von \mathbf{Y}_R und \mathbf{Y}_G müssen beide Teile mittels Parallelschaltung miteinander verbunden werden. Ohne dies im Detail erläutern zu wollen, sei gesagt, dass dieses Vorgehen zu verzögerungsfreien gerichteten Schleifen in den Wellengrößen führt. Wir suchen daher nach einer Synthese des konstanten Teils, welche in eine Wellen-Digital-Struktur mündet, die keine verzögerungsfreien Schleifen enthält. Zudem fordern wir, dass die Wellen-Digital-Struktur nur aus Elementarbausteinen mit 2 Toren im Wellenbereich besteht, da es so einfacher ist die Passivität sicherzustellen [Meer79].

Zur Lösung des Problems greifen wir auf die Erkenntnisse aus Abschnitt 2.8.1 zurück. Wir wählen nun $M = 2\tilde{n}$. Schließen wir das M -Tor an den letzten \tilde{n} Toren mit positiven Widerständen ab, so stellt die Gesamtschaltung bzgl. der ersten \tilde{n} Tore eine passive dynamikfreie, nichtreziproke Schaltung dar und lässt sich unter Einhaltung der obigen Forderungen realisieren. Zur Einhaltung der Forderung nach Nichtauftreten von verzögerungsfreien Schleifen, wählen wir die Torwiderstände gleich den Abschlusswiderständen, sodass die Schaltung im Bild 4.2 resultiert. Die Frage ist nun, wie man sich die Streu- bzw.

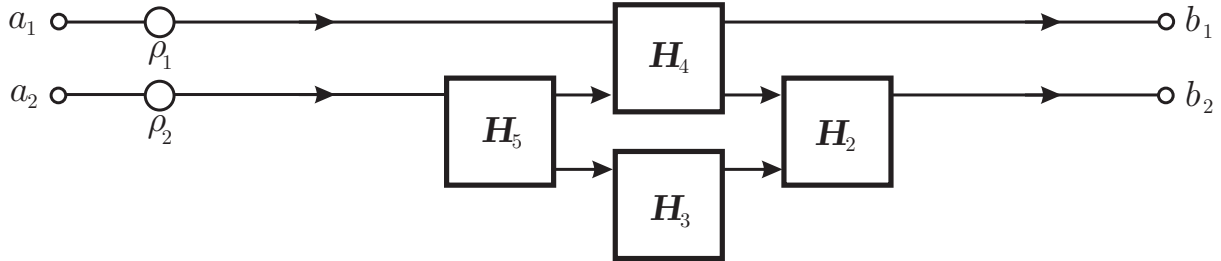


Bild 4.2: Realisierung eines passiven nichtreziproken \tilde{n} -Tors für $\tilde{n} = 2$

die Impedanzmatrix der energieneutralen Schaltung beschafft, sodass die gewünschte resultierende dissipative Schaltung bzgl. der ersten \tilde{n} Tore entsteht. Geht man von der partitionierten Streumatrix aus, so entstehen kompliziert zu lösende Beziehungen. Wir gehen deshalb von der Admittanzmatrix aus

$$\begin{bmatrix} \mathbf{I}_1 \\ \mathbf{I}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{Y}_{11} & \mathbf{Y}_{12} \\ \mathbf{Y}_{21} & \mathbf{Y}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{U}_1 \\ \mathbf{U}_2 \end{bmatrix}. \quad (4.79)$$

Der Abschluss der letzten \tilde{n} Tore wird durch $\mathbf{U}_2 = -\mathbf{Z}_t \mathbf{I}_2$ beschrieben, sodass

$$\mathbf{I}_1 = [\mathbf{Y}_{11} - \mathbf{Y}_{12}[\mathbf{Y}_t + \mathbf{Y}_{22}]^{-1}\mathbf{Y}_{21}]\mathbf{U}_1 \quad (4.80)$$

resultiert. Nun bringen wir die Bedingungen der Energieneutralität für die Admittanzmatrix mit ein $\mathbf{Y}_{11} = -\mathbf{Y}_{11}^H$, $\mathbf{Y}_{12} = -\mathbf{Y}_{21}^H$, $\mathbf{Y}_{22} = -\mathbf{Y}_{22}^H$. Setzen wir dies ein, so haben wir

$$\bar{\mathbf{Y}}_k = \mathbf{Y}_{11} + \mathbf{Y}_{21}^H[\mathbf{Y}_t + \mathbf{Y}_{22}]^{-1}\mathbf{Y}_{21}. \quad (4.81)$$

Wir setzen $\mathbf{Y}_{22} = \mathbf{0}$ und $\mathbf{Y}_t = G\mathbf{1}_{\tilde{n}}$. Dadurch ist der symmetrische Teil von $\bar{\mathbf{Y}}_k$ durch $\mathbf{Y}_{21}^H \mathbf{Y}_{21} / G$ bestimmt und der schiefsymmetrische Teil durch \mathbf{Y}_{11} , d.h.

$$\mathbf{Y}_{11} = \frac{1}{2}[\bar{\mathbf{Y}}_k - \bar{\mathbf{Y}}_k^T] \quad \text{und} \quad \mathbf{Y}_{21}^H \mathbf{Y}_{21} = \frac{G}{2}[\bar{\mathbf{Y}}_k + \bar{\mathbf{Y}}_k^T]. \quad (4.82)$$

Wobei wir nun \mathbf{Y}_{21} z.B. durch die Cholesky-Zerlegung ermitteln können. Da $\bar{\mathbf{Y}}_k$ reell ist, folgt aus der nicht negativen Definitheit des symmetrischen Teils auch die Reellwertigkeit von \mathbf{Y}_{21} . Es ist offensichtlich, dass die Cholesky-Zerlegung nicht die optimale Lösung bzgl. des Rechenaufwands liefert. Insbesondere könnte eine Optimierung dahingehend durchgeführt werden, dass die Anzahl der Multiplizierer der Wellen-Digital-Struktur minimiert wird.

4.2.4 Gesamtreferenzschaltung

Nachdem die Teilschaltungen synthetisiert wurden, suchen wir nun nach einer geeigneten Realisierung der Gesamtschaltung. Insbesondere ist darauf zu achten, dass das MDWDF keine verzögerungsfreien gerichteten Schleifen enthält. Von den untersuchten Gesamtschaltungen wird im Folgenden eine vorgestellt. Zur Herleitung der Schaltung gehen wir von Gleichung (4.20) in der Form

$$\mathbf{I}^{c1,4} + \mathbf{I}^{c2,5} + \mathbf{I}^{c3,6} + \mathbf{I}^{c7} + \mathbf{I}_k = \mathbf{J} \quad (4.83)$$

beschrieben werden. Die Synthese erfolgte durch eine mit $\frac{1}{Z_3} = \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} C_{\mu\nu}^x$ und $\frac{1}{Z_4} = \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} C_{\mu\nu}^x$ abgeschlossene Jaumann-Schaltung mit $n = 1$. Die Nachbildung der idealen Kondensatoren durch Verzögererelemente im Wellenbereich erfolgt gemäß Kapitel 2.8.5 durch Wahl der Torwiderstände nach Gleichung (2.173) zu

$$R_3 = R_4 = \frac{T}{2C_{\mu\nu}^x} . \quad (4.89)$$

Fassen wir die Torleitwerte aller $N^2 - N$ idealen Kondensatoren der Matrix $\mathbf{Y}^{c1,4}$ in der Diagonalmatrix

$$\mathbf{G}^{c1,4} = \frac{2}{T} \mathbf{diag} (C_{21}^x, \dots, C_{N,N-1}^x) \quad (4.90)$$

zusammen, so erhalten wir durch Anwendung der Trapezregel auf die reaktiven Bauelemente mit

$$z_\nu = e^{p_\nu T} , \quad z_\nu = \frac{1 + \psi_\nu}{1 - \psi_\nu} , \quad p_\nu = \frac{2}{T} \operatorname{artanh}(\psi_\nu) \approx \frac{2\psi_\nu}{T} \quad (4.91)$$

dann näherungsweise

$$\mathbf{R}^{c1,4} \mathbf{Y}_D^{c1,4} \approx \mathbf{diag} \left(\left\{ \begin{smallmatrix} \psi_4 \\ \psi_1 \end{smallmatrix} \right\}, \dots, \left\{ \begin{smallmatrix} \psi_1 \\ \psi_4 \end{smallmatrix} \right\} \right) . \quad (4.92)$$

Wir werden im Folgenden das Näherungszeichen durch das Gleichheitszeichen ersetzen. Die zu $\mathbf{Y}_D^{c1,4}$ gehörige Streumatrix lautet dann

$$\mathbf{S}_D^{c1,4} = [\mathbf{1} - \mathbf{R}^{c1,4} \mathbf{Y}_D^{c1,4}] [\mathbf{1} + \mathbf{R}^{c1,4} \mathbf{Y}_D^{c1,4}]^{-1} = \mathbf{diag} \left(\left\{ \begin{smallmatrix} z_4^{-1} \\ z_1^{-1} \end{smallmatrix} \right\}, \dots, \left\{ \begin{smallmatrix} z_1^{-1} \\ z_4^{-1} \end{smallmatrix} \right\} \right) . \quad (4.93)$$

Die Matrizen $\mathbf{Y}_D^{c2,5}$ und $\mathbf{Y}_D^{c3,6}$ werden in gleicher Weise in eine WD-Struktur umgesetzt.

Nun suchen wir eine WDF-Realisierung für \mathbf{Y}^{c7} . Diese Matrix wurde in Gleichung (4.61) aufbereitet. Um für die darin enthaltenen idealen Kondensatoren $\mathbf{diag}(\bar{\mathbf{C}}_D, \tilde{\mathbf{C}}^7)$ einfache Signalflussdiagramme zu erhalten, legen wir die Torleitwertmatrizen zu

$$\bar{\mathbf{G}}^{c7} = \frac{2}{T} \bar{\mathbf{C}}_D . \quad (4.94)$$

und

$$\tilde{\mathbf{G}}^{c7} = \frac{2}{T} \tilde{\mathbf{C}}^7 . \quad (4.95)$$

fest. Die idealen Kondensatoren haben somit i. S. der Trapezregel die Streumatrix $\mathbf{1}_N z_7^{-1}$.

Umsetzung der Jaumann-Adaptoren in WD-Elemente

Die einzelnen Jaumann-Strukturen, die durch

$$\mathbf{S}_0 \begin{bmatrix} \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} & 0 \\ 0 & \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} \end{bmatrix} \mathbf{S}_0^T \quad (4.96)$$

beschrieben werden, sind im Kapitel 2.8.2 erläutert worden. Durch die Festlegungen der Torwiderstände der idealen Kondensatoren an den jeweiligen Toren 3 und 4 gemäß Gleichung (4.90) ist der erste Teil

der Torwiderstände der Jaumann-Adaptoren bestimmt. Durch die Forderung nach Reflexionsfreiheit resultieren die Torleitwerte der jeweils anderen Tore 1 und 2 aus Gleichung (2.153) mit $n = 1$ zu

$$\frac{1}{2} \mathbf{G}^{c1,4}. \quad (4.97)$$

Nachdem alle Torwiderstände festliegen, ergibt sich die Streumatrix der Jaumann-Adaptoren als direkte Summe der Streumatrizen eines Jaumann-Adaptors gemäß Gleichung (2.155) zu

$$\mathbf{S}^{jau} = \text{diag} \left(\begin{bmatrix} \mathbf{0} & \mathbf{S}_1 \\ \mathbf{S}_1^T & \mathbf{0} \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{0} & \mathbf{S}_1 \\ \mathbf{S}_1^T & \mathbf{0} \end{bmatrix} \right). \quad (4.98)$$

Umsetzung des idealen Übertragers in WD-Elemente

Für die WD-Struktur des $\mathbf{1} : \mathbf{N}$ -Übertragers fordern wir eine reflexionsfreie Realisierung. Die Torleitwerte der rechten Seite liegen gemäß Gleichung (4.94) fest. Die Torleitwerte $\hat{\mathbf{G}}^{c7}$ der linken Seite des Übertragers ergeben sich dann gemäß Gleichung (2.118) zu $\hat{\mathbf{G}}^{c7} = \mathbf{N}^T \text{diag} \bar{\mathbf{G}}^{c7} \mathbf{N}$. Wir bemerken, dass hier die Torwiderstandsmatrix $\hat{\mathbf{G}}^{c7}$ keine Diagonalmatrix mehr ist. Welche Auswirkungen dies hat, werden wir später sehen.

Umsetzung des konstanten Teils in WD-Elemente

Die Torleitwertmatrix von $\mathbf{P}^k \mathbf{Y}^k \mathbf{P}^{kT}$ ist

$$\mathbf{G}^k = \text{diag}(\bar{\mathbf{G}}^k, \tilde{\mathbf{G}}^k). \quad (4.99)$$

Die Torleitwerte der Widerstände werden zu $\tilde{\mathbf{G}}^k = \tilde{\mathbf{Y}}^k$ gewählt. Dann sind die Streumatrizen der Widerstände $\tilde{\mathbf{S}}^k = \mathbf{0}$. Die Torwiderstände $\bar{\mathbf{G}}^k$ sind frei wählbar.

4.4 Umsetzung der Gesamtschaltung in ein MDWDF

Bei der Gesamtschaltung handelt es sich i. W. um N Parallelschaltungen von $2N$ -Toren, die mit idealen Stromquellen abgeschlossen sind. Die zur Umsetzung in ein MDWDF notwendigen N Paralleladaptoren mit angeschlossenen idealen Stromquellen wurden im Abschnitt 2.8.3 hergeleitet.

Im Bild 4.4 ist das MDWDF der direkten Umsetzung der Referenzschaltung aus Bild 4.3 dargestellt. Hierbei wurde die Tatsache ausgenutzt, dass sich ein n -Tor-Paralleladaptor durch einen $(n-1)$ -Tor-Paralleladaptor und einen 3-Tor-Paralleladaptor darstellen lässt. Dies ist aus dem Grund nötig, da die gemeinsamen Tore von Paralleladaptoren und idealen Übertragern keine diagonale Torwiderstandsmatrix mehr besitzen. Die rechten N skalaren Dreitor-Paralleladaptoren müssen als ein (vektorieller) Paralleladaptor aufgefasst werden. Im Kapitel 2.8.2 wurde gezeigt, dass nur ein Tor eines Paralleladaptors reflexionsfrei sein kann. Demnach sind zwischen dem rechten vektoriellen Paralleladaptor und dem energieneutralen $2\tilde{n}$ -Tor verzögerungsfreie gerichtete Schleifen vorhanden.

Um die verzögerungsfreien Schleifen zu beseitigen, fassen wir den rechten vektoriellen Paralleladaptor, die MDWDF-Realisierung des Übertragers und das energieneutrale $2\tilde{n}$ -Tor zu einem energieneutralen $(\tilde{n} + \tilde{n} + \tilde{n})$ -Tor zusammen. Die Toranzahl \tilde{n} wird später erläutert. Das Bild 4.5 zeigt das MDWDF. Zu diesem energieneutralen $(\tilde{n} + \tilde{n} + \tilde{n})$ -Tor berechnen wir nun die Streumatrix. Dazu betrachten wir die Matrix

$$\begin{aligned} \mathbf{Y}_r &= \mathbf{P}^{c7} \text{diag}(p_7 \bar{\mathbf{C}}^7, \mathbf{0}) \mathbf{P}^{c7T} + \mathbf{P}^k \text{diag}(\bar{\mathbf{Y}}^k, \mathbf{0}) \mathbf{P}^{kT} \\ &= \mathbf{P}^{c7} \begin{bmatrix} \bar{\mathbf{C}}^7 p_7 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{P}^{c7T} + \mathbf{P}^k \begin{bmatrix} \bar{\mathbf{Y}}^k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{P}^{kT}, \quad \bar{\mathbf{C}}^7 = \mathbf{N}^T \bar{\mathbf{C}}_D^7 \mathbf{N}. \end{aligned} \quad (4.100)$$

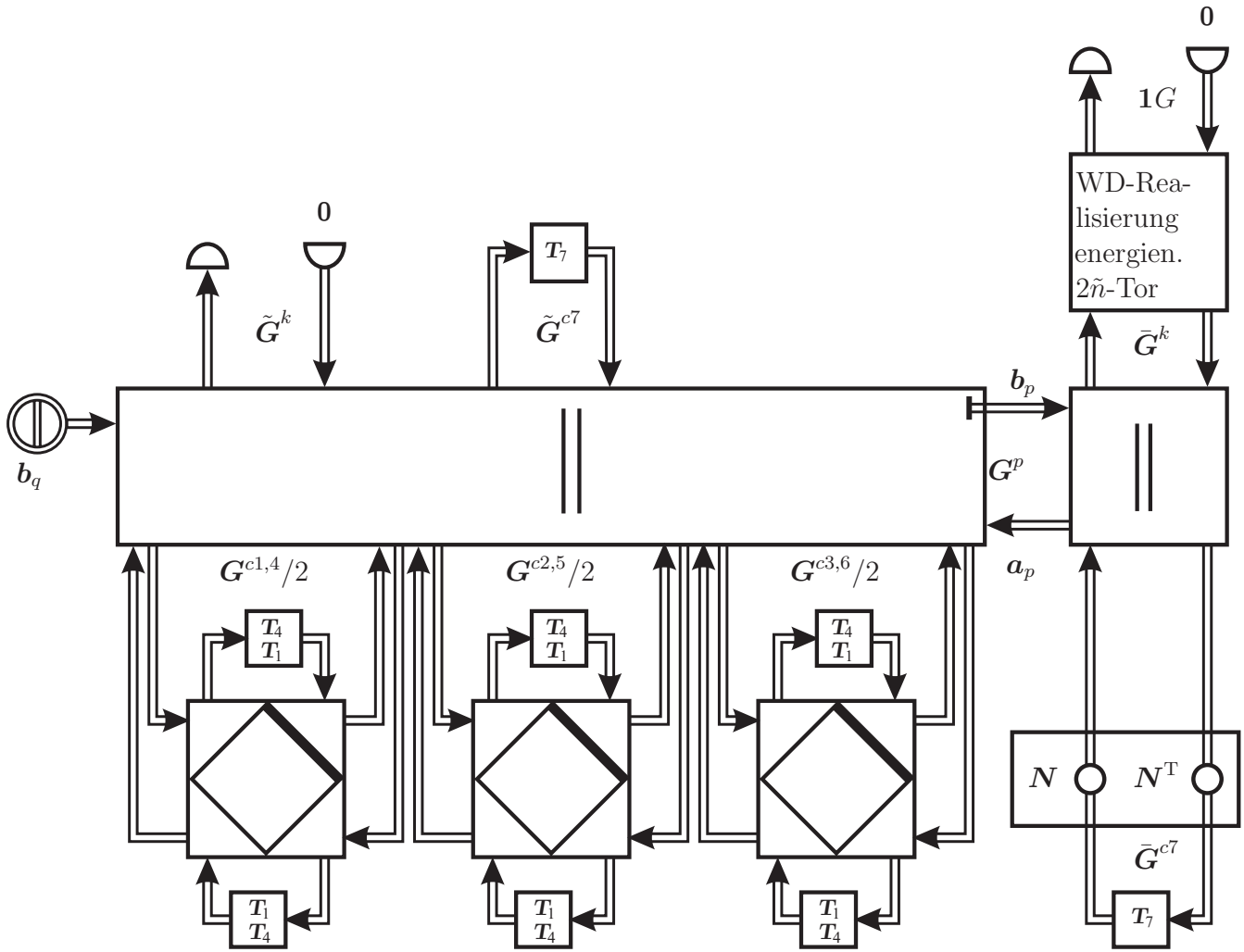


Bild 4.4: Wellendigitalrealisierung

Die Matrix besitzt möglicherweise Nullzeilen und Nullspalten mit gleichem Index. Ihre Anzahl soll $N - \tilde{n}$ sein. Sie werden im Folgenden eliminiert. Dazu multiplizieren wir die Matrix von links und von rechts mit einer Permutationsmatrix, sodass die Matrix

$$\begin{aligned}
 \hat{P}^T Y_r \hat{P} &= \hat{P}^T P^{c7} \text{diag}(p_7 \bar{C}^7, 0) P^{c7T} \hat{P} + \hat{P}^T P^k \text{diag}(\bar{Y}^k, 0) P^{kT} \hat{P} \\
 &= \hat{P}^T P^{c7} \begin{bmatrix} \bar{C}^7 p_7 & 0 \\ 0 & 0 \end{bmatrix} P^{c7T} \hat{P} + \hat{P}^T P^k \begin{bmatrix} \bar{Y}^k & 0 \\ 0 & 0 \end{bmatrix} P^{kT} \hat{P}
 \end{aligned} \tag{4.101}$$

resultiert. Die letzten $N - \tilde{n}$ Zeilen und Spalten der zu addierenden Matrizen haben ebenfalls nur Elemente, die null sind. Wir betrachten nun nur noch die ersten \tilde{n} Zeilen und Spalten von $\hat{P}^T Y_r \hat{P}$, d. h. die Matrix

$$Y_1 = \begin{bmatrix} \mathbf{1}_{\tilde{n}} & 0 \end{bmatrix} \hat{P}^T Y_r \hat{P} \begin{bmatrix} \mathbf{1}_{\tilde{n}} \\ 0 \end{bmatrix}, \quad I_1 = Y_1 U_1. \tag{4.102}$$

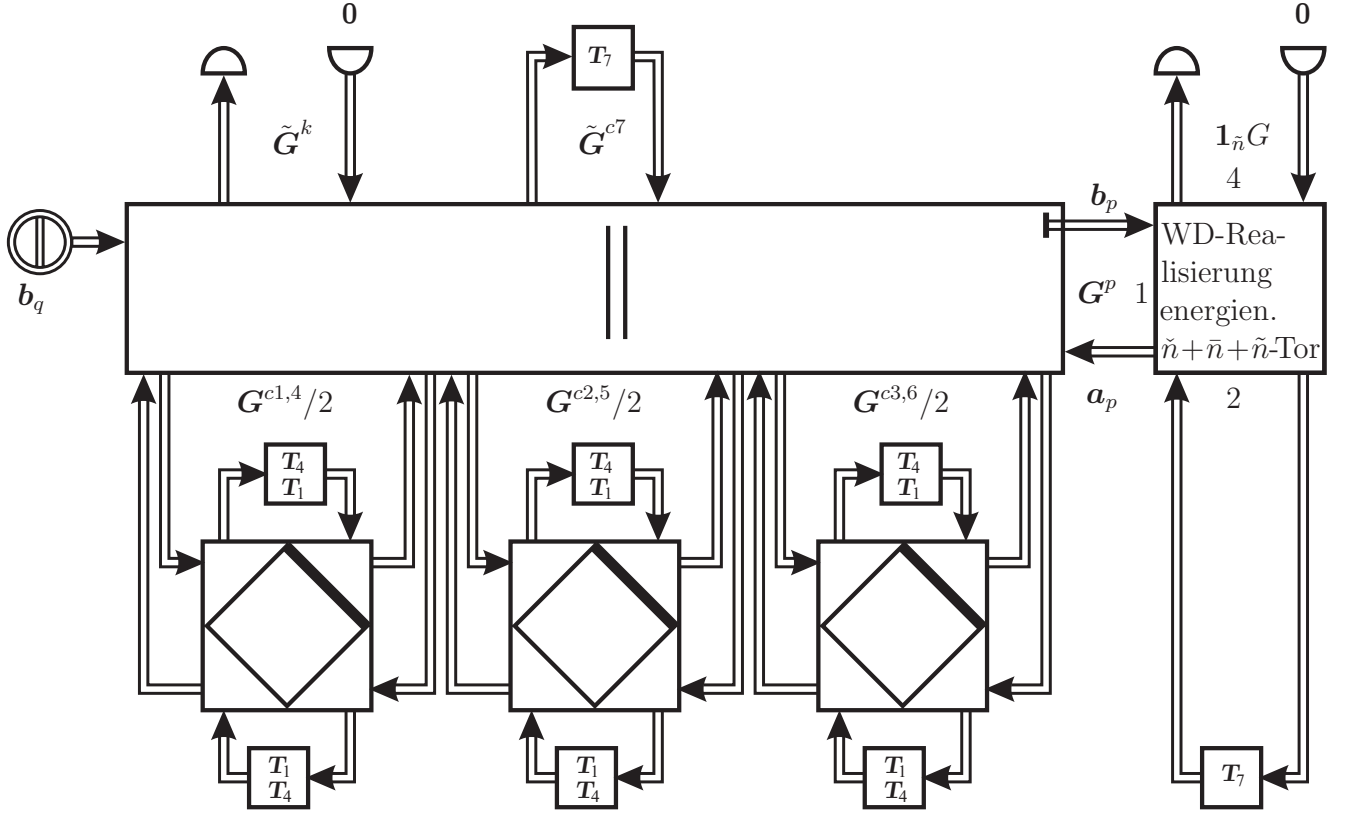


Bild 4.5: Wellendigitalrealisierung ohne verzögerungsfreie gerichtete Schleifen

Nun schreiben wir die Beziehungen zwischen Strom und Spannung neu auf

$$\begin{aligned} I_1 = Y_1 U_1 = \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} \hat{P}^T P^{c7} \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} N^T p_7 \bar{C}_D^7 N \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} P^{c7T} \hat{P} \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} U_1 + \\ \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} \hat{P}^T P^k \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} \bar{Y}^k \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} P^{kT} \hat{P} \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} U_1. \end{aligned} \quad (4.103)$$

Führen wir die Matrizen

$$C_1 = \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} \hat{P}^T P^{c7} \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} N^T \quad \text{und} \quad C_2 = \begin{bmatrix} 1_{\tilde{n}} & 0 \end{bmatrix} \hat{P}^T P^k \begin{bmatrix} 1_{\tilde{n}} \\ 0 \end{bmatrix} \quad (4.104)$$

als wesentlichen Teil von Inzidenzmatrizen ein, so lautet dies

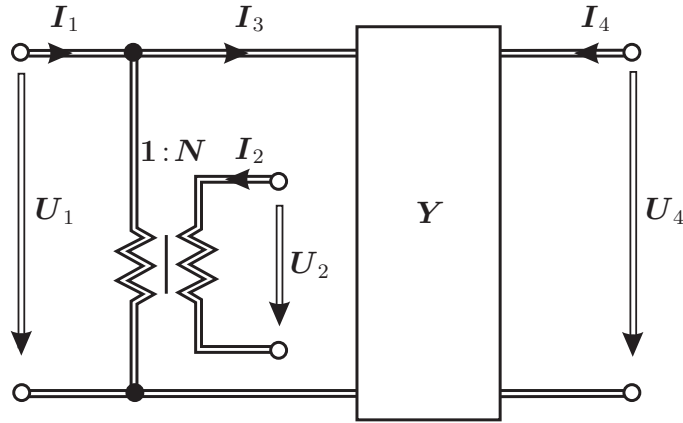
$$I_1 = Y_1 U_1 = C_1 p_7 \bar{C}_D^7 C_1^T U_1 + C_2 \bar{Y}^k C_2^T U_1 = C_1 p_7 \bar{C}_D^7 U_2 + C_2 \bar{Y}^k U_3 = C_1 (-I_2) + C_2 I_3. \quad (4.105)$$

Das Verbindungsnetz kann somit durch

$$\begin{bmatrix} 1_{\tilde{n}} & C_1 & C_2 \end{bmatrix} \begin{bmatrix} I_1 \\ I_2 \\ -I_3 \end{bmatrix} = 0 \quad \text{und} \quad \begin{bmatrix} -C_1^T & 1_{\tilde{n}} & 0 \\ -C_2^T & 0 & 1_{\tilde{n}} \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \\ U_3 \end{bmatrix} = 0 \quad (4.106)$$

beschrieben werden, wobei wir uns von der Orthogonalitätsbeziehung zwischen den Inzidenzmatrizen überzeugen

$$\begin{bmatrix} 1_{\tilde{n}} & C_1 & C_2 \end{bmatrix} \begin{bmatrix} -C_1 & -C_2 \\ 1_{\tilde{n}} & 0 \\ 0 & 1_{\tilde{n}} \end{bmatrix} = 0. \quad (4.107)$$

Bild 4.6: Symbolische Darstellung zur Berechnung der Streumatrix des energieneutralen $\tilde{n} + \tilde{n} + \tilde{n}$ -Tors

Die Bezugsrichtung des Stromes I_3 wurde entgegen der Bezugsrichtung der Spannung U_3 gewählt. Nun wird Tor 3 mit dem energieneutralen $2\tilde{n}$ -Tor

$$\begin{bmatrix} I_3 \\ I_4 \end{bmatrix} = \begin{bmatrix} Y_{11} & Y_{12} \\ Y_{21} & 0 \end{bmatrix} \begin{bmatrix} U_3 \\ U_4 \end{bmatrix} \quad (4.108)$$

entsprechend Bild 4.5 verbunden. Es ist die Streumatrix bezogen auf die Tore 1,2 und 4 gesucht, vgl. Bild 4.6. Diese Zeichnung ist wieder symbolisch zu verstehen, d.h. die auftretenden Knoten sind keine im üblichen Sinne, sondern werden durch Gleichung (4.106) beschrieben. Es werden zunächst I_3 und U_3 durch Nutzen von Gleichung (4.106) eliminiert

$$\begin{bmatrix} C_2 I_3 \\ I_4 \end{bmatrix} = \begin{bmatrix} C_2 Y_{11} & C_2 Y_{12} \\ Y_{21} & 0 \end{bmatrix} \begin{bmatrix} U_3 \\ U_4 \end{bmatrix} = \begin{bmatrix} I_1 + C_1 I_2 \\ I_4 \end{bmatrix} = \begin{bmatrix} C_2 Y_{11} & C_2 Y_{12} \\ Y_{21} & 0 \end{bmatrix} \begin{bmatrix} C_2^T U_1 \\ U_4 \end{bmatrix}. \quad (4.109)$$

Zu diesen 2 Gleichungen nehmen wir noch das Spannungsübertragungsverhältnis $C_1^T U_1 = U_2$ in der Form $0 = -G_2 C_1^T U_1 + G_2 U_2$ mit G_2 als Torleitwertmatrix der Tore 2 hinzu und erhalten

$$\underbrace{\begin{bmatrix} 1_{\tilde{n}} & C_1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1_{\tilde{n}} \end{bmatrix}}_Z \begin{bmatrix} I_1 \\ I_2 \\ I_4 \end{bmatrix} = \underbrace{\begin{bmatrix} C_2 Y_{11} C_2^T & 0 & C_2 Y_{12} \\ -G_2 C_1^T & G_2 & 0 \\ Y_{21} C_2^T & 0 & 0 \end{bmatrix}}_{\check{Y}} \begin{bmatrix} U_1 \\ U_2 \\ U_4 \end{bmatrix}. \quad (4.110)$$

Übrigens äußert sich die Energieneutralität in $z \check{Y}^H + \check{Y} z^H = 0$. Die eigentlich gesuchte Streumatrix berechnet sich durch

$$S = G^{1/2} [zG + \check{Y}]^{-1} [zG - \check{Y}] R^{1/2}. \quad (4.111)$$

Im Folgenden werden wir zeigen, dass die Inverse von

$$[zG + \check{Y}] = \begin{bmatrix} G_1 + C_2 Y_{11} C_2^T & C_1 G_2 & C_2 Y_{12} \\ -G_2 C_1^T & G_2 & 0 \\ Y_{21} C_2^T & 0 & G_4 \end{bmatrix} \quad (4.112)$$

existiert. Dazu rechnen wir ihren symmetrischen Teil aus

$$[zG + \check{Y}] + [zG + \check{Y}]^T = \begin{bmatrix} G_1 + G_1 + C_2 (Y_{11} + Y_{11}^T) C_2^T & C_1 (G_2 - G_2) & C_2 (Y_{12} + Y_{21}^T) \\ -(G_2 - G_2) C_1^T & G_2 + G_2 & 0 \\ (Y_{12}^T + Y_{21}) C_2^T & 0 & G_4 + G_4 \end{bmatrix} = 2G$$

(4.113)

also die p.d. Matrix der doppelten Torleitwerte. Wenn der symmetrische Teil einer reellen Matrix \mathbf{M} p.d. ist, ist die Matrix selbst regulär. Zu dieser Erkenntnis gelangt man durch Einsetzen von $\mathbf{x}^T \mathbf{M} \mathbf{x} = \mathbf{x}^T \mathbf{M}^T \mathbf{x}$ in die Definition der positiven Definitheit des symmetrischen Teils $\mathbf{x}^T [\mathbf{M}/2 + \mathbf{M}^T/2] \mathbf{x} > 0 \forall \mathbf{x} \neq \mathbf{0}$ d.h. $\mathbf{x}^T \mathbf{M} \mathbf{x} > 0 \forall \mathbf{x} \neq \mathbf{0}$. Die Matrix \mathbf{M} muss regulär sein, andernfalls würde ein Vektor $\mathbf{x} \neq \mathbf{0}$ existieren, der $\mathbf{M} \mathbf{x} = \mathbf{0}$ erfüllt und die Ungleichung verletzt.

Wir berechnen nun die Streumatrix der N Paralleladaptoren aus Bild 4.5. Jeder Paralleladaptor hat maximal $M = \{3[N^2 - N] + 3\}$ Tore. Die Torleitwerte der Paralleladaptoren liegen bis auf \mathbf{G}^p fest. Dies Torleitwerte werden nun so gewählt, dass die zugehörigen Tore reflexionsfrei sind. Den Torleitwert des Tores ν bezeichnen wir mit $g_{\nu\nu}^p$. Nach Gleichung (2.130) berechnet sich dieser Torleitwert durch die Summe der Torleitwerte der restlichen Tore, die an diesen Knoten inzidieren. Diese Summe kann sehr einfach mittels Skalarprodukt aus den Inzidenzmatrizen dargestellt werden. Die Torleitwertmatrix lautet

$$\mathbf{G}^p = \frac{1}{2} \mathbf{A}_{\text{rvn}} [\mathbf{G}^{c1,4} + \mathbf{G}^{c2,5} + \mathbf{G}^{c3,6}] \mathbf{A}_{\text{rvn}}^T + \mathbf{P}^k \text{diag}(\mathbf{0}_{\tilde{n}}, \tilde{\mathbf{G}}^k) \mathbf{P}^{kT} + \mathbf{P}^7 \text{diag}(\mathbf{0}_{\tilde{n}}, \tilde{\mathbf{G}}^{c7}) \mathbf{P}^{7T}. \quad (4.114)$$

Diese Torleitwerte bestimmen dann auch die des energieneutralen $\tilde{n} + \tilde{n} + \bar{n}$ -Tors.

4.5 Nachweis der Berechenbarkeit

Zum Nachvollziehen der folgenden Überlegungen empfiehlt sich die Betrachtung von Bild 4.5. Zunächst stellen wir fest, dass die ausfallenden Wellengrößen der Quellen, der Verzögererelemente und der Widerstände bekannt sind. Aufgrund der Reflexionsfreiheit der Jaumann-Adaptoren lassen sich die in Richtung der Paralleladaptoren (mit idealer Stromquelle) aus den Jaumann-Adaptoren ausfallenden Wellen (wir wollen diese hier $\mathbf{b}_{\nu\nu}$ nennen) unmittelbar bestimmen. Die Paralleladaptoren nach Gleichung (2.166) können durch

$$\begin{bmatrix} \mathbf{a}_{\nu\nu} \\ \mathbf{b}_{p\nu} \end{bmatrix} = \mathbf{S}_\nu \begin{bmatrix} \mathbf{b}_{\nu\nu} \\ \mathbf{a}_{p\nu} \end{bmatrix} + \mathbf{b}_{q\nu} \quad (4.115)$$

beschrieben werden. Die Vektoren $\mathbf{b}_{\nu\nu}$ und $\mathbf{b}_{q\nu}$ sind bekannt. Aufgrund der Reflexionsfreiheit erfährt $\mathbf{a}_{p\nu}$ keine Berücksichtigung bei der Berechnung von $\mathbf{b}_{p\nu}$. Somit kann $\mathbf{b}_{p\nu}$ für $\nu = 1, 2, \dots, N$ parallel berechnet werden. Der Vektor der ausfallenden Wellen der Paralleladaptoren \mathbf{b}_p ist nun bekannt. Da die von den Widerständen ausfallenden Wellen null sind und die aus den Verzögerern bzgl. t_7 ausfallenden Wellen bekannt sind, können sämtliche aus dem energieneutralen $(\tilde{n} + \tilde{n} + \bar{n})$ -Tor ausfallenden Wellen \mathbf{a}_p bestimmt werden. Diese Berechnung stellt die einzige Kopplung der jeweiligen WD-Realisierungen, der an die N Netzknoten inzidierenden Kirchhoff'schen Elemente dar. Diese Berechnung kann nicht knotenweise parallel ausgeführt werden. Da \mathbf{a}_p nun bekannt ist, können die restlichen einfallenden Wellen der Jaumann-Adaptoren bestimmt werden und damit $\mathbf{a}_{\nu\nu}$ für $\nu = 1, 2, \dots, N$. Zusammengefasst ergibt sich somit die Reihenfolge

1. berechne alle in Richtung der Paralleladaptoren aus den Jaumann-Adaptoren ausfallenden Wellen,
2. berechne $\mathbf{b}_{p\nu}$ für jeden Knoten $\nu = 1, 2, \dots, N$ parallel,
3. verknüpfe die Knoten durch Berechnung von \mathbf{a}_p und gleichzeitig Berechnung der restlichen aus dem $(\tilde{n} + \tilde{n} + \bar{n})$ -Tor ausfallenden Wellen,
4. berechne die ausfallenden Wellen der Paralleladaptoren,
5. berechne $\mathbf{a}_{\nu\nu}$ für jeden Knoten $\nu = 1, 2, \dots, N$ parallel.

Somit ist gezeigt, dass die Berechenbarkeit sichergestellt ist. Die Tatsache, dass die einzelnen Ortspunkte parallel berechnet werden können (bis auf den Rand) bleibt hiervon unberührt. Durch die vorliegende Darstellung wird deutlich, dass innerhalb der Parallelberechnung der Ortspunkte auch noch die Berechnung der Knoten (des elektrischen Netzes) zum Großteil parallel erfolgen kann.

festgelegt. Die Streumatrix der Jaumann-Adaptoren lautet

$$\begin{bmatrix} \mathbf{b}_{e3} \\ \mathbf{b}_{e12} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{S}_{12}^j \\ \mathbf{S}_{21}^j & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{e3} \\ \mathbf{a}_{e12} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{S}_{12}^j \\ \mathbf{S}_{21}^j & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_{e11} \\ \mathbf{b}_{v1} \end{bmatrix}. \quad (4.118)$$

Die Streumatrizen der resistiven Abschlüsse sind gleich der Nullmatrix, d. h.

$$\mathbf{b}_{e1} = \mathbf{0} \mathbf{a}_{e1} \quad , \quad \mathbf{b}_{e2} = \mathbf{0} \mathbf{a}_{e2} . \quad (4.119)$$

Auf eine Beschreibung der torweisen Verbindungen mittels der Permutationsmatrix \mathbf{P} verzichten wir, da sich die Matrix \mathbf{L} unmittelbar aus Bild 4.7 und den zuvor definierten Matrizen zu

$$\mathbf{L} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & | & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & | & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{12}^j & \mathbf{0} & \mathbf{0} & | & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}_{12}^p & \mathbf{0} & \mathbf{S}_{14}^p & \mathbf{0} & | & \mathbf{S}_{13}^p & \mathbf{0} & \mathbf{S}_{15}^p & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{S}_{13}^{en} & | & \mathbf{0} & \mathbf{S}_{12}^{en} & \mathbf{0} & \mathbf{S}_{11}^{en} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{S}_{23}^{en} & | & \mathbf{0} & \mathbf{S}_{22}^{en} & \mathbf{0} & \mathbf{S}_{21}^{en} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{S}_{33}^{en} & | & \mathbf{0} & \mathbf{S}_{32}^{en} & \mathbf{0} & \mathbf{S}_{31}^{en} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}_{22}^p & \mathbf{0} & \mathbf{S}_{24}^p & \mathbf{0} & | & \mathbf{S}_{23}^p & \mathbf{0} & \mathbf{S}_{25}^p & \mathbf{0} & \mathbf{S}_{21}^p & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}_{32}^p & \mathbf{0} & \mathbf{S}_{34}^p & \mathbf{0} & | & \mathbf{S}_{33}^p & \mathbf{0} & \mathbf{S}_{35}^p & \mathbf{0} & \mathbf{S}_{31}^p & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}_{42}^p & \mathbf{0} & \mathbf{S}_{44}^p & \mathbf{0} & | & \mathbf{S}_{43}^p & \mathbf{0} & \mathbf{S}_{45}^p & \mathbf{0} & \mathbf{S}_{41}^p & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{S}_{52}^p & \mathbf{0} & \mathbf{S}_{54}^p & \mathbf{0} & | & \mathbf{S}_{53}^p & \mathbf{0} & \mathbf{S}_{55}^p & \mathbf{0} & \mathbf{S}_{51}^p & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & | & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{S}_{21}^j & \mathbf{0} \end{bmatrix}, \quad \mathbf{b}_e = \mathbf{L} \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ \mathbf{b}_e \end{bmatrix} \quad (4.120)$$

ablesen lässt. Die strikte untere Dreiecksgestalt der Matrix \mathbf{SP}_{ee} liegt offenbar vor, sodass die Berechenbarkeit gegeben und unmittelbar erkennbar ist.

4.6 Alternative Methoden

Als Alternative können wir Referenzschaltungen herleiten, die auf einer Zustandstransformation basieren. Dazu transformieren wir \mathbf{Y}^{c7} auf Diagonalgestalt. Die meisten Eigenschaften von \mathbf{Y} bleiben durch die reguläre Transformation erhalten. Der wesentliche Nachteil dieser Synthese ist der rapide ansteigende Aufwand, da die i.d.R. lichten Matrizen i. Allg. zu voll besetzten entarten. Aus diesem Grunde sollen diese Referenzschaltungen hier nicht weiter betrachtet werden, stattdessen wird auf die Untersuchungsergebnisse in [Voll02] verwiesen.

4.7 Randbehandlung

Wie im Kapitel 2.10 dargestellt, benötigen die Verzögerer in Richtung \mathbf{T}_1 bis \mathbf{T}_6 am Rand eine spezielle Behandlung. Die Randbehandlung erfolgt jedoch nicht für jeden Verzögerer einzeln, sondern immer für eine Jaumann-Struktur. Wir erläutern die Randbehandlung an Hand

$$\mathbf{T}_{\mu\nu} \mathbf{S}_0 \begin{bmatrix} \left\{ \begin{smallmatrix} p_4 \\ p_1 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} & 0 \\ 0 & \left\{ \begin{smallmatrix} p_1 \\ p_4 \end{smallmatrix} \right\} \frac{C_{\mu\nu}^x}{4} \end{bmatrix} \mathbf{S}_0^T \mathbf{T}_{\mu\nu}^T. \quad (4.121)$$

Neben den schon festliegenden Knotenspannungen definieren wir noch die zugehörigen Torströme $I_\mu^{C_{\mu\nu}^x}$ und $I_\nu^{C_{\mu\nu}^x}$ wie folgt

$$\mathbf{S}_0 \begin{bmatrix} \{p_4\} C_{\mu\nu}^x/4 & 0 \\ 0 & \{p_1\} C_{\mu\nu}^x/4 \end{bmatrix} \mathbf{S}_0^T \begin{bmatrix} U_\mu^{C_{\mu\nu}^x} \\ U_\nu^{C_{\mu\nu}^x} \end{bmatrix} = \begin{bmatrix} I_\mu^{C_{\mu\nu}^x} \\ I_\nu^{C_{\mu\nu}^x} \end{bmatrix}. \quad (4.122)$$

Im Folgenden betrachten wir nur die obere Alternative die untere Alternative erfährt eine ähnliche Behandlung. Zudem soll auf das Abdrucken des hochgestellten $C_{\mu\nu}^x$ verzichtet werden. Ferner ersetzen wir für die restlichen Überlegungen dieses Kapitels die Indizes der unabhängigen Größen durch 1 und 2. Die Tornummerierungen und die Bezugspfeilrichtungen in diesem Kapitel entsprechen denen des Jaumann-Adaptors im Kapitel 2.8.2.

Zunächst wollen wir das Problem konkretisieren, indem wir feststellen, welcher Verzögerer an welchem Rand eine Randbehandlung benötigt. Dazu formen wir die letzte Gleichung zu

$$\begin{bmatrix} p_4 C_{\mu\nu}^x & 0 \\ 0 & p_1 C_{\mu\nu}^x \end{bmatrix} \frac{1}{2} \mathbf{S}_0^T \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} = \mathbf{S}_0^T \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} \quad (4.123)$$

um. Mit den aus Kapitel 2.8.2 bekannten Beziehungen

$$\begin{bmatrix} U_3 \\ U_4 \end{bmatrix} = \frac{1}{2} \mathbf{S}_0^T \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}, \quad \begin{bmatrix} I_3 \\ I_4 \end{bmatrix} = -\mathbf{S}_0^T \begin{bmatrix} I_1 \\ I_2 \end{bmatrix} \quad (4.124)$$

haben wir

$$\begin{bmatrix} p_4 C_{\mu\nu}^x & 0 \\ 0 & p_1 C_{\mu\nu}^x \end{bmatrix} \begin{bmatrix} U_3 \\ U_4 \end{bmatrix} = - \begin{bmatrix} I_3 \\ I_4 \end{bmatrix}. \quad (4.125)$$

Die Anwendung der Trapezregel entspricht der Näherung $p_\nu \approx 2\psi_\nu/T$. Mit dem Torwiderstand $R = T/(2C_{\mu\nu}^x)$ an den idealen Kondensatoren haben wir

$$\begin{bmatrix} \psi_4/R & 0 \\ 0 & \psi_1/R \end{bmatrix} \begin{bmatrix} U_3 \\ U_4 \end{bmatrix} \approx - \begin{bmatrix} I_3 \\ I_4 \end{bmatrix}. \quad (4.126)$$

Im Folgenden werden wir anstatt des Näherungszeichens ein Gleichheitszeichen verwenden. Mit $\psi_\nu = \frac{z_\nu - 1}{z_\nu + 1}$ erhalten wir in den Wellengrößen

$$\begin{bmatrix} A_3 \\ A_4 \end{bmatrix} = \mathbf{diag}(z_4^{-1}, z_1^{-1}) \begin{bmatrix} B_3 \\ B_4 \end{bmatrix}. \quad (4.127)$$

Mit den aus Kapitel 2.8.2 bekannten Beziehungen

$$\begin{bmatrix} A_3 \\ A_4 \end{bmatrix} = \mathbf{S}_1^T \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, \quad \begin{bmatrix} B_3 \\ B_4 \end{bmatrix} = \mathbf{S}_1^T \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad (4.128)$$

schreiben wir diese Gleichungen als

$$\begin{bmatrix} B_1 \\ B_2 \end{bmatrix} = \mathbf{S}_1 \mathbf{diag}(z_4^{-1}, z_1^{-1}) \mathbf{S}_1^T \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}. \quad (4.129)$$

Nun lässt sich das Signalflussdiagramm leicht angeben, vgl. Bild 4.8. Durch Betrachtung der Gleichungen $A_3 = z_4^{-1} B_3$ und $A_4 = z_1^{-1} B_4$ im Zeitbereich, d. h. $a_3(\boldsymbol{\nu} + \mathbf{e}_4) = b_3(\boldsymbol{\nu})$ und $a_4(\boldsymbol{\nu} + \mathbf{e}_1) = b_4(\boldsymbol{\nu})$ erkennen

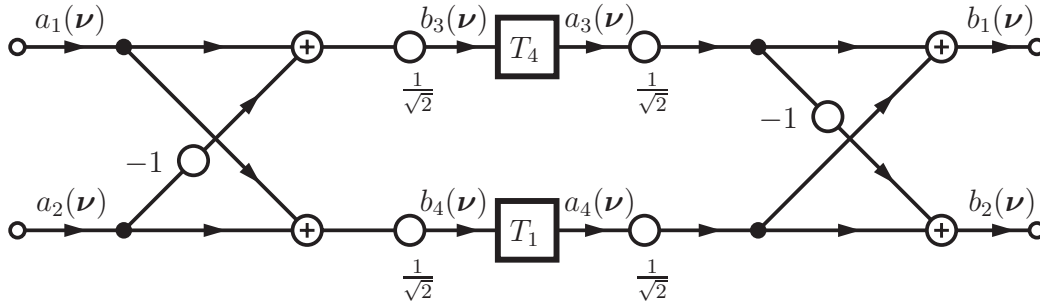


Bild 4.8: Signalflussdiagramm bei Normalberechnung

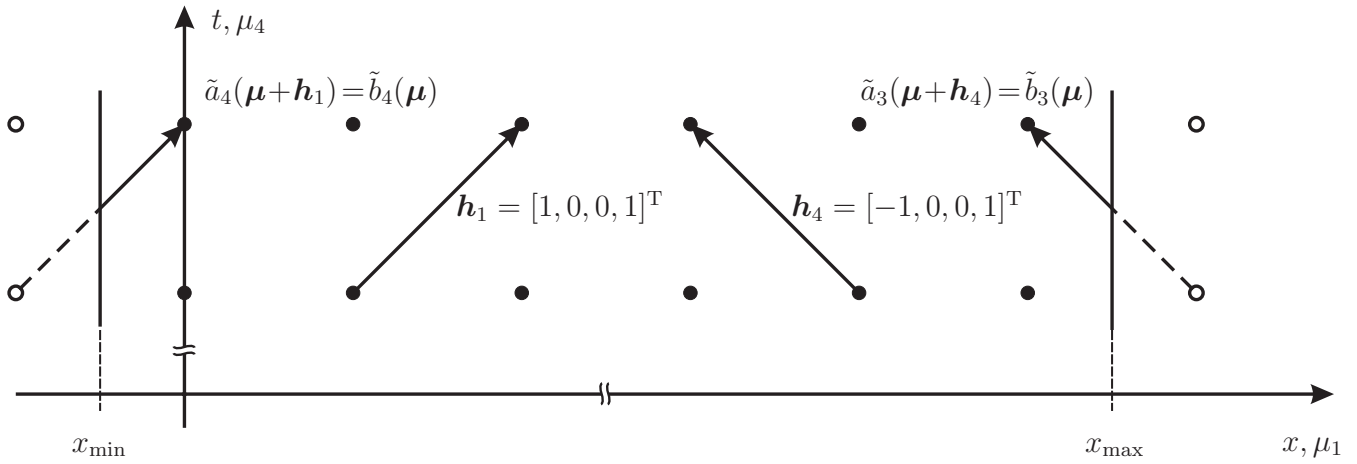


Bild 4.9: Verschiebevektoren bei Normalberechnung

wir, dass der erste Kondensator eine Randbehandlung an $x = x_{\max}$ und der zweite Kondensator eine Randbehandlung an $x = x_{\min}$ erfordert, vgl. Bild 4.9.

Zur eigentlichen Randbehandlung wenden wir Gleichung (2.191) auf Gleichung (4.125) im Zeitbereich an. Hierbei haben wir die mit negativen Vorzeichen versehenen Ströme durch Vertauschen der ein- und ausfallenden Wellen zu berücksichtigen. Mit Wahl der Torwiderstände zu $R = T/(2C_{\mu\nu}^x)$ erhalten wir für den ersten Kondensator

$$a_3(\nu) = \frac{1}{2\sqrt{R}}u_3(\nu - \mathbf{e}_4/2) \quad , \quad b_3(\nu) = \frac{1}{2\sqrt{R}}u_3(\nu + \mathbf{e}_4/2) \quad (4.130)$$

und für den zweiten Kondensator

$$a_4(\nu) = \frac{1}{2\sqrt{R}}u_4(\nu - \mathbf{e}_1/2) \quad , \quad b_4(\nu) = \frac{1}{2\sqrt{R}}u_4(\nu + \mathbf{e}_1/2) . \quad (4.131)$$

Von diesen 4 Gleichungen werden jeweils 2 zur Randbehandlung an $x = x_{\min}$ und $x = x_{\max}$ herangezogen. Um die Randfläche $x = x_{\max}$ zu behandeln, nutzen wir die 2 Gleichungen

$$a_3(\nu) = \frac{1}{2\sqrt{R}}u_3(\nu - \mathbf{e}_4/2) \quad \text{und} \quad b_4(\nu) = \frac{1}{2\sqrt{R}}u_4(\nu + \mathbf{e}_1/2) . \quad (4.132)$$

Um die Randfläche $x = x_{\min}$ zu behandeln, nutzen wir die verbleibenden 2 Gleichungen

$$a_4(\nu) = \frac{1}{2\sqrt{R}}u_4(\nu - \mathbf{e}_1/2) \quad \text{und} \quad b_3(\nu) = \frac{1}{2\sqrt{R}}u_3(\nu + \mathbf{e}_4/2) . \quad (4.133)$$

Die am Rand $x = x_{\max}$ durch die normale Berechnung nicht ermittelbare Welle $a_3(\nu)$ können wir uns nun aus der bekannten Welle $b_4(\nu)$ beschaffen, wenn wir einen Zusammenhang zwischen den Feldgrößen u_1 und u_2 auf dem Rand haben. Dabei greifen wir zwei Spezialfälle heraus. Zum einen soll eine konstante Beziehung zwischen u_1 und u_2 auf dem Rand und zum anderen eine lineare differentielle Beziehung bzgl. t zwischen u_1 und u_2 auf dem Rand bestehen. Beide Fälle sind mit dem Ansatz $U_1 = F(x, p_t)U_2$ behandelbar. In Gleichung (4.132) können nun die Feldgrößen eliminiert werden und wir erhalten die gewünschte Beziehung zwischen den Zustandsgrößen auf der Randfläche $x = x_{\max}$

$$A_3 = \underbrace{\frac{F(x_{\max}, p_t) - 1}{F(x_{\max}, p_t) + 1}}_{H_1(p_t)} z_1^{-1/2} z_4^{-1/2} B_4 \quad , \quad (4.134)$$

bzw. im Zeitbereich

$$a_3(\nu + \mathbf{e}_7) = h_1(\nu_7) * b_4(\nu) \quad (4.135)$$

mit dem Faltungsoperator bzgl. ν_7 . Im Bild 4.10 ist das WDF für die Randfläche $x = x_{\max}$ dargestellt, wobei das System S_1 die Impulsantwort $h_1(\nu_7)$ besitzt.

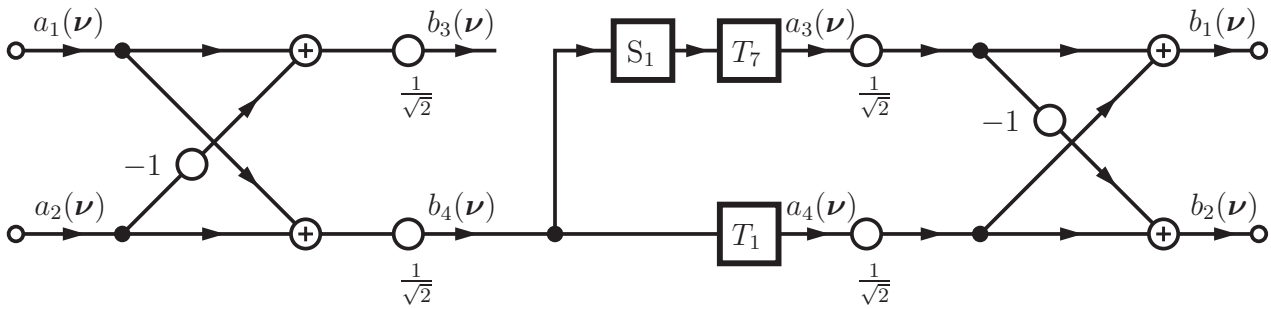


Bild 4.10: Zur Erläuterung der Randbehandlung an $x = x_{\max}$

Für $x = x_{\min}$ wird aus Gleichung (4.133)

$$A_4 = \underbrace{\frac{F(x_{\min}, p_t) + 1}{F(x_{\min}, p_t) - 1}}_{H_2(p_t)} z_1^{-1/2} z_4^{-1/2} B_3 \quad , \quad (4.136)$$

bzw. im Zeitbereich

$$a_4(\nu + \mathbf{e}_7) = h_2(\nu_7) * b_3(\nu) \quad . \quad (4.137)$$

Im Bild 4.11 ist das WDF für die Randfläche $x = x_{\min}$ dargestellt, wobei das System S_2 die Impulsantwort $h_2(\nu_7)$ besitzt.

Aufgrund des Verzögerers $z_1^{-1/2} z_4^{-1/2} = z_t^{-1} = z_7^{-1}$, ist das MDWDF auch am Rand explizit berechenbar. Bild 4.12 zeigt die Verschiebevektoren bei Randbehandlung.

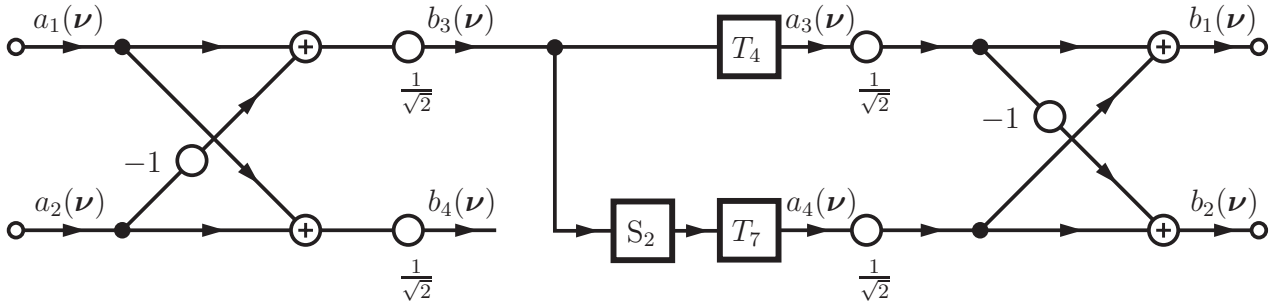
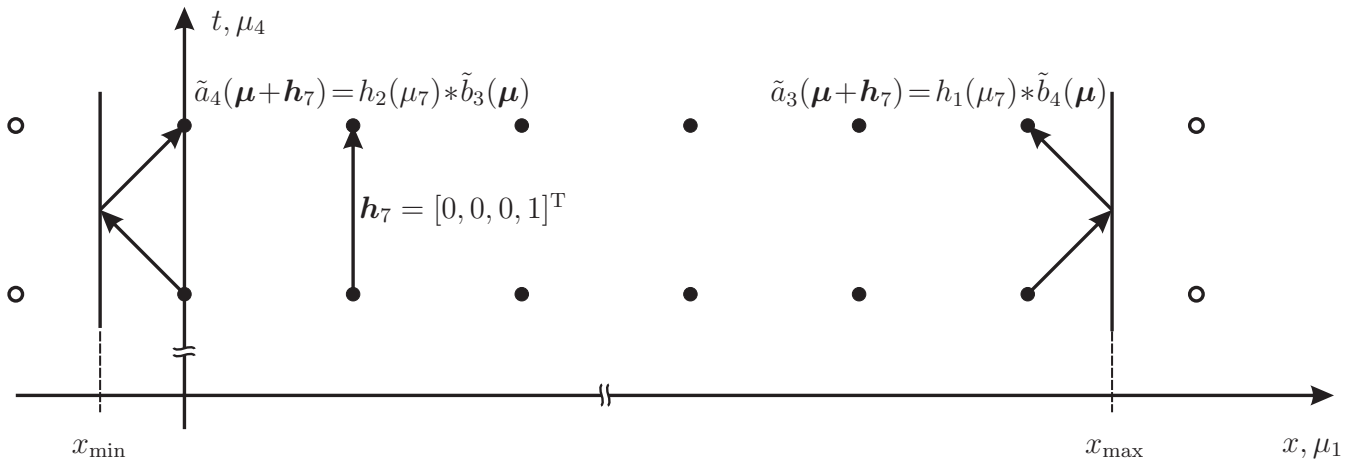
Bild 4.11: Zur Erläuterung der Randbehandlung an $x = x_{\min}$ 

Bild 4.12: Verschiebevektoren bei Randbehandlung

4.8 Zusammenfassung

Zunächst wurde in diesem Kapitel ein System aus linearen partiellen Differentialgleichungen betrachtet und dieses auf die Eigenschaft der Passivität untersucht. Aus der Forderung nach Passivität des quellenfreien Teils haben wir gezeigt, dass die Symmetrie der Matrizen, die den reaktiven Teil beschreiben, notwendig ist. Zudem forderte die Passivität die nicht negative Definitheit des konstanten Teils und des reaktiven Teils bzgl. der Zeit, wobei wir letztere Eigenschaft zur positiven Definitheit stärkten. Systeme mit diesen Eigenschaften werden symmetrisch hyperbolische Systeme genannt. Hyperbolische Systeme folgen insbesondere dem Nahwirkungsprinzip und ihre Vorgänge pflanzen sich mit endlicher Ausbreitungsgeschwindigkeit fort. Die vorhandene Beschreibung des bzgl. der Zeit passiven PDGLn-Systems wurde geeignet umgeformt. Anschließend wurde dieses System durch eine geeignete Koordinatentransformation in ein MD-passives System überführt. Dabei überträgt sich die Eigenschaft der Lokalität auf das System in den neuen Koordinaten. Diese vorliegende Beschreibung wurde zur Synthese einer MD-passiven Kirchhoff'schen Schaltung (der so genannten Referenzschaltung) genutzt. Das kontinuierliche System in Form einer MD-passiven Kirchhoff'schen Schaltung setzten wir durch Anwendung der verallgemeinerten Trapezregel und Verwendung von Wellengrößen in ein MDWDF um, welches die Eigenschaft der MD-Passivität im Diskreten besitzt. Es folgte das Aufbrechen von verzögerungsfreien gerichteten Schleifen und der anschließende Nachweis der Rekursierbarkeit. Zum Abschluss wurde, ausgehend von bekannten Ergebnissen im 1+1D-Fall, eine kompakte und übersichtliche Behandlung von Rändern, für die durch die Synthese entstandenen Strukturen, aufgezeigt. Hierbei haben wir uns auf konstante Randmodelle eingeschränkt.

Kapitel 5

Verifikation von Software mittels formaler Methoden

Ziel dieses Kapitels ist es, die Anwendung formaler Methoden zur Erstellung hochzuverlässiger Software darzulegen. Weiterhin werden die notwendigen Grundlagen und verwendeten Notationen für den im Kapitel 6 geführten Korrektheitsbeweis erläutert. Die herkömmliche Softwareerstellung erfolgt i.d.R. mittels des Wasserfallmodells, Spiralmodells oder des Phasenmodells [Royc87], [Boeh76], [Boeh81], [Boeh86], [Bun96]. Eine Verifikation wird zwischen den einzelnen Stufen dieser Modelle durchgeführt. Obwohl die Schwachstelle der Software-Zuverlässigkeit im Übergang von der formalen Anforderungsspezifikation (das ist Teil des Lastenheftes) zur Software-Spezifikation (das ist Teil des Pflichtenheftes) liegt, [Fisc03], steht in dieser Arbeit die Verifikation der Implementierung gegenüber der formalen Software-Spezifikation im Vordergrund.

Als Maß für die Zuverlässigkeit technischer Systeme wird oft die Ausfallrate, d.h. die Anzahl der Ausfälle innerhalb einer Zeitspanne gewählt. Je nach Anwendung ergeben sich unterschiedliche Anforderungen an die Ausfallrate. Als Methoden zur Bestimmung der Ausfallrate kommen insbesondere Tests und die Schätzung von Parametern eines Zuverlässigkeitsmodells in Betracht. Die Schätzung der Parameter erfolgt durch mehrfaches Durchlaufen eines Zyklusses, der auch zur Qualitätsverbesserung eingesetzt wird. Innerhalb eines Zyklusses wird zunächst ein Test durchgeführt. Die auftretenden Fehler werden behoben und die korrigierte Software erneut dem Test unterzogen. Dieses Vorgehen ist aber mit prinzipiellen Nachteilen behaftet. Aus dem Durchlauf eines Zyklusses muss nicht notwendigerweise eine niedrigere Ausfallrate folgen. Möglicherweise entstehen sogar durch das Beheben des Fehlers neue Fehler. Zudem werden durch den Test nicht alle Fehler aufgedeckt. Ein fehlerfreier Test hat immer nur für die gewählten Testzustände und Testeingänge Gültigkeit. Schon bei Systemen geringer Komplexität kann die Software nicht für ihr gesamtes Definitionsgebiet getestet werden. Aus den eingeschränkten Testeingaben resultiert eine nicht unterschreitbare Grenze der Ausfallrate. Weiterer Nachteil ist, dass eine nicht formale Spezifikation unvollständig, widersprüchlich oder gar falsch sein kann. Abhilfe schafft hier nur die Verwendung anderer Konzepte, z.B. die Verwendung von formalen Methoden zur Verifikation. Schon die Spezifikation der Software erfolgt mittels formaler Methoden, [Floy67], [Dijk68], [Hoar69]. Der Grundgedanke dieses Verfahrens besteht darin, eine mathematische Beschreibung für jede Teilanweisung zu finden, wie es auch in anderen ingenieurwissenschaftlichen Disziplinen üblich ist. Vergleichbare Konzepte sind z.B. Signalflussdiagramme oder Schaltpläne. In diesen Fällen ist jeweils die symbolische Darstellung in Form von Blöcken oder elektrischen Schaltelementen völlig äquivalent zu einer mathematischen Beschreibung, z.B. mittels Differentialgleichungen. Bei den formalen Methoden ist nun jeder symbolischen Darstellung in Form des Quellcodes eine mathematische Beschreibung (die Semantik des Programmcodes) zugeordnet. Ebenso wie für Signalflussdiagramme und Schaltpläne, Regeln für die Beschreibung zusammengesetzter Elemente existieren, sind auch Regeln für die Verknüpfung von Elementaranweisungen vorhanden. Dem gut dokumentierten, systematischen Erstellen von möglichst

allgemeinen wiederverwendbaren Modulen kommt bei Anwendung formaler Methoden im Hinblick auf Qualität und schnelle Entwicklungszeiten noch größere Bedeutung zu. Der formale Korrektheitsbeweis muss für jedes Modul geführt werden. Aus diesen Modulen kann ein System erzeugt werden. Für den Korrektheitsbeweis des Systems kann auf die Beweise der Module zurückgegriffen werden.

Obwohl die Theorie der Programm-Korrektheits-Beweisführung hinsichtlich sequentieller Programme weitgehend entwickelt worden ist, findet sie in der Praxis nur wenig Anwendung. Die wesentlichen Kritikpunkte, die in [Babe95] aufgelistet sind, wollen wir kurz zusammenfassen. Formale Verifikation ist zu kompliziert, erfordert zu viel hochqualifiziertes Personal, ist zu teuer in der Anwendung, wenig verständlich und nur auf kleine, meist künstliche Problemen anwendbar. Zudem bezieht sich die Beweisführung nur auf den Quellcode. Übersetzer, Betriebssystem, andere systemnahe Software, Beweiswerkzeug und die Hardware können weiterhin fehlerhaft sein. Die angesprochenen Kritikpunkte sind die Gründe dafür, dass formale Verifikationsmethoden bisher fast nur in sicherheitsrelevanten Anwendungen eingesetzt werden. Als Beitrag zur Umsetzung der bekannten theoretischen Konzepte in die Praxis ist [Babe95] zu verstehen. Eine weitere wichtige Quelle, die auch für die vorliegende Arbeit verwendet wurde, ist [Rumm98]. Die letztgenannte Arbeit beinhaltet die Programm-Korrektheits-Beweisführung für eindimensionale Wellendigitalfilter.

Das vorliegende Kapitel stellt die Grundlage für den im Kapitel 6 geführten Korrektheitsbeweis dar. Wir erläutern die Grundbegriffe und Grundkonzepte der Programm-Korrektheits-Beweisführung. Es folgt die Festlegung einer axiomatischen Basis für die in dieser Arbeit notwendigen Anweisungen der Programmiersprache C. Auf eine Berücksichtigung der Laufzeit wie sie in [Rumm98] durchgeführt wurde, wird in dieser Arbeit verzichtet, da die Übersichtlichkeit des Beweises stark darunter leiden würde. Die Berücksichtigung der Zeit kann leicht ergänzt werden. Dies läuft darauf hinaus, dass sich die Laufzeit als eine Summe der Ausführungszeiten für Additionen, Multiplikationen, For-Schleifen und Vergleichsoperationen ergibt.

5.1 Definitionen, Notationen und Vereinbarungen

In diesem Abschnitt werden wir die für den Korrektheitsbeweis benötigten Begriffe erläutern. Hierbei geben wir allerdings keine allgemeine Einführung, sondern beschränken uns auf die Dinge, die im weiteren Verlauf der Arbeit verwendet werden. Die wichtigsten Quellen sind [Rumm98], [Babe95], [HA28], [Novi73] und [Meye88].

Die Menge der Variablennamen bezeichnen wir mit \mathcal{N} . Die Menge der durch die Programmiersprache C und von uns benötigten Datentypen ist $\mathcal{D} = \{ \text{int}, \text{float} \}$. Die Menge \mathcal{B} beinhaltet die vom Programm schon genutzten Variablenbezeichner. Als Abkürzungen für die kleinsten und größten positiven Zahlen nutzen wir

$$\begin{aligned} \text{Max}_i &: \text{Darstellbare ganze Zahl mit dem größten Betrag} \\ \text{Max}_f &: \text{Darstellbare float-Zahl mit dem größten Betrag.} \end{aligned} \tag{5.1}$$

Eine skalare Programmvariable `var` ist semantisch ein Quadrupel, bestehend aus den vier PL-Variablen¹ Variablenbezeichner `var.name`, Variablentyp `var.typ`, und Variablenwert (Value) `var.v` und `var.l = -1`. Letzte kennzeichnet die Eigenschaft einer skalaren Variablen. Ist die Programmvariable ein Feld, so liegt auf semantischer Ebene ein Quadrupel vor, wobei für die Feldlänge `var.l > 0` gilt. Zudem ist `var.v` nun ein Vektor. Sofern sich durch konkrete Umstände nichts anderes ergibt, soll zur Vereinfachung im

¹Das Akronym PL steht für Prädikatenlogik. Die Prädikatenlogik untersucht den Zusammenhang zwischen den Axiomen und Sätzen mathematischer Theorien. Die Prädikatenlogik versteht sich also als eine formale Sprache, die zur Darstellung mathematischer Aussagen benutzt wird. Die logischen Sachverhalte, die zwischen Urteilen, Begriffen usw. bestehen, finden ihre Darstellung durch Formeln, deren Interpretation frei ist von Unklarheiten, die beim sprachlichen Ausdruck leicht auftreten können [HA28].

Folgenden die Programmvariable *var* abkürzend nur noch durch die PL-Variablen *name.typ*, *name.v* und *name.l* vertreten werden. Zwischen den Bezeichnern der PL-Variablen und den Bezeichnern der Programmvariablen ist streng zu unterscheiden. Nichtsdestotrotz besteht per Konvention zwischen den beiden Bezeichnern ein eindeutiger Zusammenhang. Für die Bezeichner soll gelten, dass die Programmvariablen in Schreibmaschinenschrift gesetzt werden. Die PL-Variablen werden normal kursiv gesetzt. Die PL-Variable *name.v* beschreibt immer den Ist-Wert. I.d.R. sind in Bedingungen Variablen enthalten, z.B. soll die Bedingung garantieren, dass die Variable *f* den so genannten Sollwert *x* annimmt, d.h. $f.v = x$. Der Übersicht halber wollen wir vereinbaren, dass sich die Bezeichner der PL-Variablen aus den Bezeichnern der Sollwerte und Anfügen von *.v* ergeben. In diesem Fall hätten wir $f.v = f$. Somit sind aus Kenntnis eines der Bezeichner von Sollwert, PL-Variable oder Programm-Variablen-Bezeichner die anderen zwei jeweils eindeutig bestimmbar. Wir greifen also bei der Wahl der Programm-Variablen-Bezeichner de facto auf die in der Spezifikation enthaltenen Bezeichner der Sollwerte zurück.

Bei der Betrachtung von Variablen mit Indizes stellt sich heraus, dass das angegebene Vorgehen nicht möglich ist, da die Syntax der Programm-Variablen-Bezeichner der Programmiersprache C keine Indizes erlaubt. Ist die PL-Variable ein Matrixelement, so soll der zugehörige Programm-Variablen-Bezeichner

`<Matrixkleinbuchstabe><Matrixindexoben><Matrixindexunten>_<Zeilenindex>_<Spaltenindex>`

lauten. Als Beispiel betrachten wir die PL-Variable $\mathbf{R}_a.v$. Für das Element der Zeile 5 und Spalte 2 lautet die PL-Variable $r_{a52}.v$ und die Programmvariable `ra_5_2`.

Im Falle von Feldern muss sowohl aus den PL-Variablen als auch aus den Programmvariablen hervorgehen, welches Feldelement angesprochen wird.

Für die Programmiersprache C liegt die Programmvariable zu `a[k]` fest. Die zugehörige PL-Variable trägt die Bezeichnung $a(k)$. Der Feldindex $k.v$ liegt im Intervall $0 \leq k.v \leq a.l - 1$.

Wir besprechen nun noch die Behandlung griechischer Buchstaben, die als PL-Variablen Verwendung finden. Der Programm-Variablen-Bezeichner ist der in kleinen arabischen Lettern ausgeschriebene griechische Buchstabe. Der Ist-Wert der PL-Variablen ist der griechische Buchstabe mit der zusätzlichen Endung *.v*.

Als abschließendes Beispiel betrachten wir das Element der Zeile *m* und der Spalte *n* der Matrix \mathbf{A}_R^e , deren Elemente Felder sind und vom Feldindex `delta` abhängen. Der Sollwert lautet $a_{Rmn}^e(\delta)$, die PL-Variable $a_{Rmn}^e(\delta.v).v$ und die Programmvariable `aeR_m_n[delta]`.

Im Zusammenhang mit prädikatenlogischen Ausdrücken werden wir die Identität

$$P_1 \equiv P_2 \left[\begin{smallmatrix} a.v \\ b.v \end{smallmatrix} \right] \quad (5.2)$$

verwenden, die so zu verstehen ist, dass der Ausdruck P_1 sich aus dem Ausdruck P_2 ergibt, indem dort jede auftretende prädikatenlogische Variable *a.v* durch *b.v* ersetzt wird.

5.2 Das Konzept der formalen Methoden

Um das Konzept der formalen Methoden zu erläutern, betrachten wir eine Folge von Programmvariablen, die wir als Datenumgebung bezeichnen wollen. Als Vorbedingung wird ein Zustand einer Datenumgebung bezeichnet, in dem sich die Datenumgebung vor der Ausführung der Anweisung befinden muss. Die Nachbedingung gibt den Zustand an, in dem sich die Datenumgebung nach der Ausführung der Anweisung befinden muss. Die Nachbedingung ist i.d.R. abhängig vom Datenzustand vor Ausführung der Anweisung. Das Prinzip der formalen Methoden basiert nun darauf, das Verhalten einer Anweisung auf den Datenzustand vor und nach der Anweisung durch die PL-Formel

$$\{V\}S\{P\} \quad (5.3)$$

zu beschreiben. Hierin sind die Vorbedingung V und die Nachbedingung P selber PL-Formeln. An Stelle des Begriffs Bedingung wird auch der Begriff Zusicherung verwandt.

Der fundamentale Satz der formalen Methoden ist der folgende: eine Bedingung V ist eine (gewöhnliche) Vorbedingung der Nachbedingung P bezüglich der Programmanweisung S , falls die Wahrheit von V vor der Ausführung von S die Wahrheit von P nach Ausführung sicherstellt, d.h. die PL-Formel

$$V \implies \{V\}S\{P\} \quad (5.4)$$

wahr ist. In dem Fall erfüllt die Anweisung die Spezifikation und wird als korrekt bezeichnet. In der Literatur wird noch zwischen partieller Korrektheit und totaler Korrektheit differenziert. Eine Korrektheitsformel $\{V\}S\{P\}$ heißt partiell korrekt, wenn die Anweisung aus jeder zulässigen Vorbedingung in eine beliebige Nachbedingung $\{P\}$ terminiert. Eine Korrektheitsformel $\{V\}S\{P\}$ heißt total korrekt, wenn die Anweisung aus jeder zulässigen Vorbedingung in eine vorgegebene Nachbedingung $\{P\}$ terminiert. In dieser Arbeit interessieren wir uns nur für die totale Korrektheit. Aus diesem Grund werden wir auf das Adjektiv total in Zukunft verzichten.

Für die Beweisführung sind prinzipiell 2 Varianten anwendbar. Ausgehend von einer gegebenen Vorbedingung kann man den Nachweis führen, dass die Anweisung der gegebenen Nachbedingung genügt. Andererseits kann man von der Nachbedingung im Rückschluss den Nachweis führen, dass die Anweisung der Vorbedingung genügt. Oft ist es auch so, dass nur die Nachbedingung vorgeschrieben ist. Durch Anwendung der formalen Methoden wird dann eine Vorbedingung gefunden, die einen Teil der Spezifikation bildet.

In der Praxis sind die Bedingungen durch logische Verknüpfungen einzelner Formeln definiert. Die Formeln können abhängig von reellen Zahlen sein, wie das folgende Beispiel zeigt. Wir betrachten eine Anweisung, die die Wurzel einer Zahl vom Datentyp `float` ziehen soll.

Anweisung S : $y = \text{sqrt}(x)$
 Vorbedingung V : $\text{Max}_f \geq x.v \geq 0 \wedge \alpha = \sqrt{x.v}$
 Nachbedingung P : $y.v = \alpha$

In der Kompaktnotation

$$\{\text{Max}_f \geq x.v \geq 0 \wedge \alpha = \sqrt{x.v}\}y = \text{sqrt}(x)\{y.v = \alpha\}. \quad (5.5)$$

In diesem einführenden Beispiel wurden einige Bedingungen unterdrückt.

Wie jede mathematische Theorie, benötigt auch die Verifikation von Software als Basis Axiome. Diese Basis bildet sich aus Elementaranweisungen, Axiomen zur Verknüpfung von Elementaranweisungen und Axiomen zur Vereinfachung der Beweisführung.

Elementaranweisungen

Elementaranweisungen sind Anweisungen, die axiomatisch definiert sind. Man unterscheidet nach Einzelanweisungen und Blockanweisungen. Blockanweisungen sind dadurch gekennzeichnet, dass sich innerhalb der Blockanweisung noch Anweisungen befinden.

Die von uns benötigten Elementaranweisungen sind Variablendeklarationsanweisungen, Zuweisungen, algebraische Operatoren (Addition, Subtraktion, Multiplikation), Boolesche Operatoren (Und, Oder, Negation), Klammeroperatoren und for-Schleifenanweisungen. Eine axiomatische Definition dieser Elementaranweisungen in der Programmiersprache C erfolgt im Kapitel Gleichung (5.3).

Anweisungssequenzen

Mittels einer Elementaranweisung lassen sich keine komplexen Aufgaben lösen. Ein Programm wird daher aus einer Sequenz von Elementaranweisungen bestehen. Ein Axiom welches die PL-Ausdrücke zweier Anweisungen zu einem PL-Ausdruck der Anweisungssequenz verknüpft, wird nun eingeführt.

Dazu betrachten wir die beiden Anweisungen S_1 und S_2 , die nacheinander abgearbeitet werden sollen. Die Frage ist, wie die Korrektheit der zusammengesetzten Anweisung $S_1; S_2$ bewiesen werden kann. Die Vorbedingung dieser Anweisungssequenz sei V und die Nachbedingung P , d.h.

$$\{V\} S_1; S_2 \{P\} . \quad (5.6)$$

Die Nachbedingung von S_1 bezeichnen wir mit R . Sie stellt gleichzeitig die Vorbedingung von S_2 dar. Für den Nachweis der Korrektheit kann man auf das folgende Axiom zurückgreifen

$$\{V\} S_1 \{R\} \wedge \{R\} S_2 \{P\} \implies \{V\} S_1; S_2 \{P\} , \quad (5.7)$$

d.h. der Nachweis kann auf das Zeigen der Richtigkeit von $\{V\} S_1 \{R\}$ und $\{R\} S_2 \{P\}$ zurückgeführt werden.

Durch Sequenzen von Elementaranweisungen kann man neue Anweisungen definieren. Für den Nachweis der Korrektheit der Anweisung darf man nur auf Axiome oder als korrekt bewiesene Anweisungen zurückgreifen. Für die effiziente Durchführung von Korrektheitsbeweisen bietet es sich daher an, eine Bibliothek neuer Anweisungen zu erstellen oder zu erwerben. Diese Module sollten einen möglichst allgemein gültigen Charakter besitzen.

Stärkung, Schwächung

Die Axiome „Stärkung der Vorbedingung“ und „Schwächung der Nachbedingung“, welche auch gelegentlich als Konsequenzaxiome bezeichnet werden, können zur Vereinfachung von Beweisen genutzt werden.

Zur Erläuterung der Axiome gehen wir von einer Vorbedingung V_1 und einer Nachbedingung P_1 bzgl. der Anweisung S aus. Das Axiom über die „Stärkung der Vorbedingung“ besagt, dass V auch eine gültige Vorbedingung bzgl. S und P_1 ist, wenn V schärfer als V_1 ist. Anders ausgedrückt, dass die Vorbedingung V die Vorbedingung V_1 impliziert, d.h.

$$V \implies V_1 . \quad (5.8)$$

Die Menge der Zustände, die V erfüllen, ist dann nicht größer als die Menge der Zustände, die V_1 erfüllen. Ist z z.B. eine komplexe Zahl und die Vorbedingung lautet $V_1 \equiv |z| < 3$, dann ist mit der schärferen Vorbedingung $V \equiv |z| < 2$ auch V_1 erfüllt.

Das Axiom über die „Schwächung der Nachbedingung“ besagt, dass P auch eine gültige Nachbedingung bzgl. S und $\{V_1\}$ ist, wenn P schwächer als P_1 ist. Anders ausgedrückt, dass die Bedingung P_1 die Bedingung P impliziert, d.h.

$$P_1 \implies P . \quad (5.9)$$

Ist z z.B. eine komplexe Zahl und die Nachbedingung lautet $P_1 \equiv |x| < 1$, dann ist die schwächere Nachbedingung $P \equiv |x| < 2$ auch eine gültige Nachbedingung.

Beide Axiome lassen sich wie folgt kompakt darstellen

$$[(\{V_1\} S \{P_1\}) \wedge (V \implies V_1) \wedge (P_1 \implies P)] \implies \{V\} S \{P\} . \quad (5.10)$$

5.3 Eine axiomatische Definition von Elementaranweisungen der Programmiersprache C

Selbst bei so populären Sprachen wie C oder C++ ist die Syntax im Detail beschrieben. Eine formale Spezifikation der Semantik existiert aber nicht. Die Bedeutung der Anweisungen wird meist intuitiv den

Anweisungsbezeichnern entnommen und nicht einer formalen Definition. Wir werden daher die Semantik axiomatisch, formal und praxisgerecht definieren. Als Programmiersprache verwenden wir dazu C im ANSI-C-Standard, da die entwickelte Software in das Programmpaket SPACE eingebunden werden soll, das wir im Kapitel 3 vorgestellt haben. Die Definition erfolgt auf Grundlage von [Inte99]. Dabei werden wir uns auf die Definition derjenigen Elementaranweisungen beschränken, die wir benötigen, um unsere Aufgabenstellung zu lösen. Da in unserem Fall der Sicherheitsaspekt im Vordergrund steht, werden wir nicht jede Anweisung, die zur Verkürzung der Programmlaufzeit führt, verwenden.

5.3.1 Einzelanweisungen

Deklaration

Die Deklarationsanweisung einer skalaren Variablen lautet

$$\text{vartyp } \mathbf{a}; . \quad (5.11)$$

Die Vorbedingung ergibt sich aus der abgeänderten Nachbedingung. Konjunktiv kommt die Existenz des angeforderten Datentyps und die Tatsache, dass der Bezeichner gültig und noch nicht verwendet worden ist, hinzu. Die PL-Variable $a.\text{typ}$ ist durch vartyp zu ersetzen. Zusätzlich wird $a.l$ durch -1 ersetzt, um zu kennzeichnen, dass es sich um eine skalare Variable handelt. Formal erhalten wir das Deklarationsaxiom zu

$$\left\{ P \left[\begin{array}{c} a.\text{typ} \\ \text{vartyp} \end{array} \right] \left[\begin{array}{c} a.l \\ -1 \end{array} \right] \wedge \text{vartyp} \in \mathcal{D} \wedge \mathbf{a} \notin \mathcal{N} \wedge \mathbf{a} \in \mathcal{B} \right\} \text{vartyp } \mathbf{a}; \{P\} . \quad (5.12)$$

Deklaration von Feldern

Die Deklarationsanweisung einer Feldvariablen lautet

$$\text{vartyp } \mathbf{x}[\text{laenge}]; . \quad (5.13)$$

Ein eindimensionales Feld stellt einen Vektor dar. Wir verwenden für die PL-Variablen die Bezeichnungen

$$\mathbf{x}.v = [x_{0.v}, x_{1.v}, \dots, x_{x.l-2.v}, x_{x.l-1.v}]^T \quad \text{und} \quad x_\nu.v = \mathbf{x}(\nu).v . \quad (5.14)$$

Jede vektorielle PL-Variable trägt neben den Attributen der Variablen zusätzlich die Länge $x.l$. Die vektorielle PL-Variable stellt dann ein Quadrupel dar, wobei $\mathbf{x}.v$ selber ein Vektor mit den Koordinaten $x.v_\nu$, $\nu = 1, 2, \dots, x.l - 1$ ist. Die Vorbedingungen sind denen der Deklaration von skalaren Variablen ähnlich. Hinzukommt noch die Bedingung, dass die Länge des Vektors positiv sein muss. Das Axiom lautet daher

$$P \left[\begin{array}{c} x.\text{typ} \\ \text{vartyp} \end{array} \right] \left[\begin{array}{c} x.l \\ \text{laenge} \end{array} \right] \wedge \text{vartyp} \in \mathcal{D} \wedge x \notin \mathcal{N} \wedge x \in \mathcal{B} \wedge \text{laenge} > 0 \text{vartyp } \mathbf{x}[\text{laenge}]; \{P\} . \quad (5.15)$$

Zuweisung

Die Zuweisungsanweisung der Form

$$\mathbf{a} = \mathbf{b}; . \quad (5.16)$$

Hierbei muss unterschieden werden zwischen einer Programmvariablen \mathbf{b} und einer Konstanten b . Gehen wir zunächst davon aus, dass \mathbf{b} eine Programmvariable ist, lautet das Axiom

$$\{P_{[b.v]}^{a.v}\} \wedge \text{Max}_{a.\text{typ}} \geq |b.v| \wedge a.l = b.l = -1 \wedge a.\text{typ} \in \mathcal{D} \wedge b.\text{typ} = a.\text{typ}\} \mathbf{a} = \mathbf{b}; \{P\} . \quad (5.17)$$

Im Falle einer Konstanten b lautet das Axiom für den Datentyp `int`

$$\{P_{[b]}^{a.v} \wedge \text{Max}_i \geq |b| \wedge a.l = -1 \wedge a.typ = \text{int}\} a = b; \{P\} . \quad (5.18)$$

Die Zuweisung gilt nicht für float-Variablen. Im Falle von float-Variablen ist der Konstanten das Postfix f hinzuzufügen

$$\{P_{[b.f]}^{a.v} \wedge \text{Max}_f \geq |b| \wedge a.l = -1 \wedge a.typ = \text{float}\} a = b.f; \{P\} . \quad (5.19)$$

Zuweisungen bei Feldern

Die Zuweisung der Form $x[k]=y[m]$; folgt dem Axiom

$$\begin{aligned} & \{P_{[y_{m.v.v}]}^{x_{k.v.v}} \wedge \text{Max}_{a.typ} \geq |y_{m.v.v}| \wedge 0 \leq k.v \leq x.l - 1 \wedge 0 \leq m.v \leq y.l - 1 \wedge x.typ = y.typ \wedge x.typ \in \mathcal{D} \wedge \\ & 1 \leq k.v \leq \text{Max}_i \wedge 1 \leq m.v \leq \text{Max}_i \wedge k.l = -1 \wedge m.l = -1 \wedge k.typ = \text{int} \wedge m.typ = \text{int}\} \\ & x[k]=y[m]; \{P\} . \end{aligned} \quad (5.20)$$

Algebraische Grundoperationen

Die konsequente Erweiterung der Zuweisung entsteht, wenn die rechte Seite nun nicht mehr durch eine Konstante oder Variable gegeben ist, sondern durch mehrere mit algebraischen Grundoperationen verknüpfte Variablen bestimmt ist. Weiterhin ist es so, dass die Variablen der rechten Seite auch durch Konstanten ersetzt werden dürfen.

Unter den Variablen verstehen wir hier sowohl skalare Variablen, als auch Komponenten von Feldern, wobei dann in der Vorbedingung noch konjunktive Verknüpfungen bezüglich der PL-Variablen der Längen und der Typen hinzukommen. Dies soll nicht wiederholt werden, da es bei der Zuweisung schon erläutert wurde.

In der Programmiersprache C sind Operationen und Zuweisungen zwischen 'gemischten' Datentypen erlaubt. Insofern könnte man daraus weitere Axiome für die algebraischen Grundoperationen ableiten. Da wir diese im Verlauf der Arbeit nicht benötigen, verzichten wir darauf. Ebenso verzichten wir auf die Division, da der WDF-Algorithmus keine enthält.

Addition / Subtraktion von Ganzzahlen

Rundungsfehler treten nur dann auf, wenn eine Bereichsüberschreitung vorliegt

$$\{P_{[(b.v \pm c.v)]}^{a.v} \wedge \text{Max}_i \geq |b.v \pm c.v| \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b \pm c \{P\} . \quad (5.21)$$

Addition / Subtraktion von Gleitkommazahlen des Typs float

Rundungsfehler treten i.Allg. auf. Zudem ist eine Bereichsüberschreitung möglich

$$\begin{aligned} & \{P_{[(b.v \pm c.v + \varepsilon(b.v \pm c.v))]}^{a.v} \wedge \text{Max}_f \geq |b.v \pm c.v| \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{float}\} \\ & a = b \pm c \{P\} . \end{aligned} \quad (5.22)$$

Multiplikation von Ganzzahlen

Rundungsfehler treten nur dann auf, wenn eine Bereichsüberschreitung vorliegt

$$\{P_{[(b.v \cdot c.v)]}^{a.v}\} \wedge \text{Max}_i \geq |b.v \cdot c.v| \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} \quad (5.23)$$

$$a = b * c \{P\}.$$

Multiplikation von Gleitkommazahlen des Typs float

Sowohl Rundungsfehler als auch eine Bereichsüberschreitung treten i.Allg. auf.

$$\{P_{[(b.v \cdot c.v) + \varepsilon(b.v \cdot c.v)]}^{a.v}\} \wedge \text{Max}_f \geq |b.v \cdot c.v| \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{float}\}$$

$$a = b * c \{P\}.$$

(5.24)

Die im Axiom (5.22) und im Axiom (5.24) auftretende Funktion ε stellt den Rundungsfehler dar. Dieser ist zwar kaum explizit auszudrücken, jedoch kann eine betragsmäßige Obergrenze angegeben werden [Rumm98], [Babe95], [Inte99] und Anhang D.

Logische Operatoren

Im Zusammenhang mit den hier definierten logischen Operatoren wird der Variablen $a.v$ der Wert 1 zugewiesen, falls die Aussage wahr ist, andernfalls 0

$$\{P_{[(b.v \wedge c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b \ \&\& \ c; \{P\} \quad (5.25)$$

$$\{P_{[(b.v \vee c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b \ || \ c; \{P\} \quad (5.26)$$

$$\{P_{[(b.v < c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b < c; \{P\} \quad (5.27)$$

$$\{P_{[(b.v \leq c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b \leq c; \{P\} \quad (5.28)$$

$$\{P_{[(b.v > c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b > c; \{P\} \quad (5.29)$$

$$\{P_{[(b.v \geq c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b \geq c; \{P\} \quad (5.30)$$

$$\{P_{[(b.v = c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b == c; \{P\} \quad (5.31)$$

$$\{P_{[(b.v \neq c.v)]}^{a.v}\} \wedge a.l = b.l = c.l = -1 \wedge a.typ = b.typ = c.typ = \text{int}\} a = b != c; \{P\} \quad (5.32)$$

$$\{P_{[-b.v]}^{a.v}\} \wedge a.l = b.l = -1 \wedge a.typ = b.typ = \text{int}\} a = !b; \{P\}. \quad (5.33)$$

5.3.2 Blockanweisungen

For-Schleifenanweisung

Die For-Schleifenanweisung lautet

$$\text{for } (S_1 ; b ; S_3) \{ S_4 \} ; . \quad (5.34)$$

Hierin stellt S_1 die Anweisung zur Initialisierung einer Schleifenvariablen, $b.v \neq 0$ die Schleifenbedingung und S_3 die Anweisung zur Reinitialisierung der Schleifenvariablen dar. S_4 ist die Anweisung, die innerhalb der Schleife ausgeführt wird.

Der Programmablaufplan der For-Schleifenanweisung ist mit Bedingungen im Bild 5.1 dargestellt. Mithilfe der so genannten Schleifeninvarianten I gelingt die Beweisführung bei Schleifenanweisungen. Wie der Name schon sagt, ändert sich die Schleifeninvariante während eines Schleifendurchlaufs nicht. Diese Eigenschaft macht man sich beim Beweis der Korrektheit zu nutze. Die konjunktive Verknüpfung der Invarianten und der negierten Schleifenbedingung stellt die Nachbedingung dar. Zur Erfüllung der gestellten Aufgabe (Erreichung der Nachbedingung P aus der Vorbedingung V) sind mehrere Schleifeninvarianten I geeignet. Die Schleifeninvariante muss vor der Programmausführung bekannt sein. Zur Förderung des Verständnisses dient die Darstellung der Anweisung mit den Bedingungen

$ \begin{array}{ll} \{V\} & \\ & \text{for } (S_1 ; b ; S_3) \{ \\ \{I\} & \\ \{I \wedge (b.v \neq 0)\} & \\ & S_4 \\ \{Z\} & \\ & \}; \\ \{P\} & \end{array} $.
--	---

Das Axiom lautet in der Kompaktnotation

$$\begin{aligned}
 & \{V\} S_1 \{I\} \wedge \{I \wedge (b.v \neq 0)\} S_4 ; S_3 \{I\} \wedge b.typ = \text{int} \wedge b.l = -1 \\
 \Rightarrow & \{V\} \text{for}(S_1 ; b ; S_3) \{S_4\} \{I \wedge (b.v = 0)\}
 \end{aligned} \quad (5.35)$$

und es gilt $P \equiv I \wedge b.v = 0$. Zu beachten ist, dass die Klammern um S_4 syntaktisch zum C-Code gehören und nicht zu einer Bedingung.

5.4 Behandlung von Ausdrücken

In die bislang eingeführten Axiome lassen sich nur Elementaranweisungen einsetzen, welche keine arithmetischen Ausdrücke, sondern nur einfache Programmvariablen besitzen. Um nun auch arithmetische und logische Ausdrücke, wie sie im Abschnitt 5.3.1 beschrieben sind verwenden zu können, werden noch weitere Axiome nötig. Mithilfe dieser Axiome kann eine Programmanweisung, welche einen arithmetische Ausdruck besitzt, in mehrere einfachere Anweisungen zerlegt werden. Die in diesen Axiomen auftretenden Variablen b und c stellen PL-Variablen dar, welche nicht in den PL-Formeln V , P , oder I vorkommen dürfen.

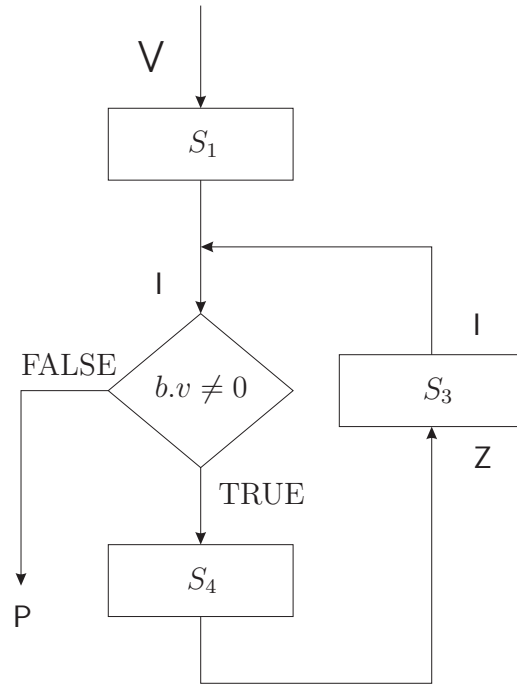


Bild 5.1: for-schleife

5.4.1 Einzelanweisungen

Wendet man den unären Negations-Operator auf einen Ausdruck $expr$ an, so kann der Ausdruck vor der Anwendung des Operators ausgewertet werden

$$\{V \wedge b.l = -1 \wedge b.typ = \text{int}\} \mathbf{b} = expr ; \mathbf{a} = !\mathbf{b}; \{P\} \Rightarrow \{V\} \mathbf{a} = !expr \{P\} . \quad (5.36)$$

Ebenso können die Ausdrücke $expr_1$ und $expr_2$ ausgewertet werden, bevor man einen beliebigen der definierten binären logischen Operatoren op auf die Ausdrücke anwendet. Wir erhalten dabei drei Fälle. Zunächst zwei Ausdrücke als Operanden

$$\begin{aligned} & \{V \wedge b.typ = \text{vartype1} \wedge c.typ = \text{vartype2} \wedge b.l = -1 \wedge c.l = -1\} \\ & \mathbf{b} = expr_1; \mathbf{c} = expr_2; \mathbf{a} = \mathbf{b} \text{ op } \mathbf{c} \{P\} \\ & \Rightarrow \{V\} \mathbf{a} = (expr_1) \text{ op } (expr_2) \{P\} . \end{aligned} \quad (5.37)$$

Für vartype1 bzw. vartype2 darf hier ein beliebiger Variablentyp aus \mathcal{D} eingesetzt werden. Die Typüberprüfung erfolgt dann später bei der Auflösung von $expr_1$ und $expr_2$ sowie bei der Zuweisung $\mathbf{a} = \mathbf{b} \text{ op } \mathbf{c}$. Die anderen beiden Fälle sind durch jeweils einen Operanden als Ausdruck und einen als Variable gekennzeichnet, d. h.

$$\{V \wedge b.typ = \text{vartype} \wedge b.l = -1\} \mathbf{b} = expr; \mathbf{a} = \mathbf{b} \text{ op } \mathbf{c} \{P\} \Rightarrow \{V\} \mathbf{a} = (expr) \text{ op } \mathbf{c} \{P\} \quad (5.38)$$

und

$$\{V \wedge c.typ = \text{vartype} \wedge c.l = -1\} \mathbf{c} = expr; \mathbf{a} = \mathbf{b} \text{ op } \mathbf{c} \{P\} \Rightarrow \{V\} \mathbf{a} = \mathbf{b} \text{ op } (expr) \{P\} . \quad (5.39)$$

5.4.2 Blockanweisungen

5.5 Einige Sätze

5.5.1 Skalarprodukt mit Variablen

Der folgende Satz stammt aus [Rumm98]. Eine Anweisungssequenz zur Berechnung des Skalarprodukts der Vektoren $\mathbf{b}.v$ und $\mathbf{c}.v$ mit der Länge N und den reellen Koordinaten $b_1.v, b_2.v, \dots, b_N.v, c_1.v, c_2.v, \dots, c_N.v$ lautet

$$\begin{array}{l} \{\mathbf{V}\} \\ \quad \mathbf{a} = ((\dots (\mathbf{b_1} * \mathbf{c_1}) \dots) + (\mathbf{b_N-1} * \mathbf{c_N-1})) + (\mathbf{b_N} * \mathbf{c_N}) \\ \{\mathbf{P}\}. \end{array} \quad (5.40)$$

Satz: Eine gültige Vorbedingung ist

$$\begin{aligned} \mathbf{V} \equiv & \mathbf{P} \left[\begin{array}{c} a.v \\ (\mathbf{b}^T.v \mathbf{c}.v) \end{array} \right] \wedge \bigwedge_{n=1}^N |b_n.v \ c_n.v| \leq \text{Max}_f \bigwedge_{n=2}^N \left| \sum_{\nu=1}^n b_\nu.v \ c_\nu.v \right| \leq \text{Max}_f \wedge a.l = -1 \\ & \bigwedge_{n=1}^N b_n.l = -1 \bigwedge_{n=1}^N c_n.l = -1 \wedge a.typ = \text{float} \wedge \bigwedge_{n=1}^N a.typ = b_n.typ \bigwedge_{n=1}^N a.typ = c_n.typ \end{aligned} \quad (5.41)$$

wobei

$$\mathbf{b}^T.v \mathbf{c}.v = \sum_{n=1}^N b_n.v \ c_n.v$$

gilt. Der angegebene Satz lässt sich mit Hilfe von Axiom (5.22) und Axiom (5.24) beweisen [Rumm98], S.121-124. Wichtig ist in diesem Zusammenhang, dass für den Programmausdruck weder das Kommutativgesetz noch das Assoziativgesetz gelten. Die Klammern dürfen nicht weggelassen werden, da die Reihenfolge der Auswertung von Unterausdrücken gleicher Priorität in C nicht definiert ist [Brey96], [Hero99].

5.5.2 Skalarprodukt mit Feldvariablen

Eine Anweisungssequenz zur Berechnung des Skalarprodukts der Vektoren $\mathbf{b}(x).v$ und $\mathbf{c}(y).v$ mit der Länge N und den reellen Koordinaten $b_1(x).v, b_2(x).v, \dots, b_N(x).v, c_1(y).v, c_2(y).v, \dots, c_N(y).v$ lautet

$$\begin{array}{l} \{\mathbf{V}\} \\ \quad \mathbf{a}[\mathbf{z}] = ((\dots (\mathbf{b_1}[\mathbf{x}] * \mathbf{c_1}[\mathbf{y}]) \dots) + (\mathbf{b_N}[\mathbf{x}] * \mathbf{c_N}[\mathbf{y}])) \\ \{\mathbf{P}\}. \end{array} \quad (5.42)$$

Satz: Eine gültige Vorbedingung ist

$$\begin{aligned}
 V \equiv & \text{P} \left[\mathbf{b}_{(x.v).v}^{\mathbf{T}} \mathbf{c}_{(y.v).v} \right] \wedge \bigwedge_{n=1}^N |b_n(x.v).v \ c_n(y.v).v| \leq \text{Max}_f \bigwedge_{n=2}^N \left| \sum_{\nu=1}^n b_\nu(x.v).v \ c_\nu(y.v).v \right| \leq \text{Max}_f \\
 & \wedge \ 0 \leq z.v \leq a.l - 1 \wedge 0 \leq x.v \leq b.l - 1 \wedge 0 \leq y.v \leq c.l - 1 \\
 & \wedge \ a.typ = \mathbf{float} \wedge \bigwedge_{n=1}^N a.typ = b_n.typ \bigwedge_{n=1}^N a.typ = c_n.typ \\
 & \wedge \ 0 \leq z.v \leq \text{Max}_i \wedge 0 \leq x.v \leq \text{Max}_i \wedge 0 \leq y.v \leq \text{Max}_i \\
 & \wedge \ z.l = y.l = x.l = -1 \wedge z.typ = x.typ = y.typ = \mathbf{int}
 \end{aligned} \tag{5.43}$$

wobei

$$\mathbf{b}^{\mathbf{T}}(x.v).v \mathbf{c}(y.v).v = \sum_{n=1}^N b_n(x.v).v \ c_n(y.v).v$$

gilt. Einer der beiden Vektoren kann auch durch einen Vektor mit Konstanten ersetzt werden. In dem Fall erübrigen sich die Forderungen an den ersetzten Vektor.

Kapitel 6

Der Algorithmus

Ziel dieses Kapitels ist es, einen durch formale Methoden verifizierten Algorithmus durch Verknüpfen von Elementaranweisungen der Programmiersprache C anzugeben, der ein lineares, konstantes MDWDF für ein quaderförmiges Berechnungsgebiet simuliert. Die Anweisungssequenz soll eine Abtastschicht berechnen und einen, gemäß Kapitel 3.5 spezifizierten, Funktionsbaustein des Programmpaketes SPACE darstellen, [Voll04b].

Inhalt des Kapitels 6.1 ist es, Vor- und Nachbedingungen des Algorithmus anzugeben. Dies entspricht im Softwareentwicklungsprozess dem Übergang von der Anforderungsspezifikation zur formalen Spezifikation. Für den formalen Korrektheitsbeweis verzichten wir allerdings aus Gründen der Übersichtlichkeit auf eine Berücksichtigung der Laufzeit. In den Kapiteln 6.2 bis 6.7 erfolgt zu den Teilanweisungen die manuelle Softwaresynthese. Zudem wird die Erfüllung der in Kapitel 6.1 festgelegten Vor- und Nachbedingungen jeweils nachgewiesen. Im Kapitel 6.8 erfolgt schließlich eine Zusammenfassung und Bewertung der Ergebnisse.

6.1 Gesamtanweisung und Nachbedingung

In diesem Unterkapitel werden wir zunächst die relevanten Ergebnisse der vorherigen Kapitel rekapitulieren, insbesondere die Voraussetzungen für die Gültigkeit der folgenden Ergebnisse wiederholen. Um einen leichteren Übergang zur formalen Spezifikation zu ermöglichen, werden anschließend die Gleichungen aufbereitet. Zum einen werden wir die Wellengrößen zu dimensionslosen Größen normieren und zum anderen wird die Darstellung der Gültigkeitsbereiche mittels Allquantoren auf konjunktiv verknüpfte Prädikate umgeformt. Es folgen Definitionen von Prädikaten, die eine einfacherer Schreibweise ermöglichen. Danach geben wir einen Algorithmus an und illustrieren diesen anhand von Programmablaufplänen. Diesen Grobentwurf zerlegen wir wiederum in sequentiell abzuarbeitende, syntaktisch zusammengehörige Teilanweisungen. Abschließend legen wir die im Kapitel 6.2 bis 6.7 benötigten Nach-, Zwischen- und Vorbedingungen fest.

Das Vorgehen und die Voraussetzungen zur Einbindung mehrdimensionaler linearer konstanter Wellendigitalfilter bei quaderförmigen Berechnungsgebiet in das Programmpaket SPACE nach Kapitel 3 wird wie folgt zusammengefasst :

- Ein Funktionsbaustein stellt den algebraischen Teil eines mehrdimensionalen Wellendigitalfilters dar und jeder Aufruf eines Funktionsbausteins entspricht der Berechnung einer Abtastschicht des Wellendigitalfilters und somit dem Fortschreiten der physikalischen Zeit.
- Die Eingänge des Funktionsbausteins sind die Quellen-Wellen des Wellendigitalfilters.
- Die Ausgänge des Funktionsbausteins sind die Spannungen und Ströme an den Toren des Referenznetzes.

- Wir gehen davon aus, dass das Verhalten der Feldgrößen am Rand algebraisch modelliert wird.
- Die abgeänderte Schnittstelle für Funktionsbausteine des Programmpaketes SPACE nach Kapitel 3.5 wird genutzt.

Die Gleichungen (2.219), (2.221), (2.228) und (3.16) aus den Kapiteln 2.9, 2.10 und 3.5, die ein mehrdimensionales Wellendigitalfilter beschreiben, lauten in den Variablen μ

$$\mathbf{b}_e(\mu) = \mathbf{L}_q \mathbf{b}_q(\mu) + \mathbf{L}_v \mathbf{b}_v(\mu) + \mathbf{L}_e \mathbf{b}_e(\mu) \quad \forall \quad \mu \in \mathcal{G} \quad (6.1)$$

$$\mathbf{a}_{\text{FB}}(\mu) = \mathbf{A} \mathbf{b}(\mu) \quad \forall \quad \mu \in \mathcal{G}, \quad (6.2)$$

$$\mathbf{a}_v^\kappa(\mu) = \mathbf{P}_{v^\kappa e} \mathbf{b}_e(\mu) \quad \forall \quad (\kappa = 1, 2, \dots, k' - 1) \wedge (\mu \in \mathcal{G}) \wedge (\mu + \mathbf{h}_\kappa \in \mathcal{G}_0) \quad (6.3)$$

$$\mathbf{b}_v^\kappa(\mu + \mathbf{h}_\kappa) = \mathbf{P}_{v^\kappa e} \mathbf{b}_e(\mu) \quad \forall \quad (\kappa = 1, 2, \dots, k') \wedge (\mu \in \mathcal{G}) \wedge (\mu + \mathbf{h}_\kappa \in \mathcal{G}) \quad (6.4)$$

$$\mathbf{b}_v^\kappa(\mu + \mathbf{e}_7) = \mathbf{R}^\kappa \mathbf{a}_v^{(\kappa+3)}(\mu) \quad \forall \quad (\kappa = 1, 2, 3) \wedge (\mu \in \mathcal{G}) \wedge (\mu + \mathbf{h}_\kappa \in \mathcal{G}_0) \quad \text{und} \quad (6.5)$$

$$\mathbf{b}_v^\kappa(\mu + \mathbf{e}_7) = \mathbf{R}^\kappa \mathbf{a}_v^{(\kappa-3)}(\mu) \quad \forall \quad (\kappa = 4, 5, 6) \wedge (\mu \in \mathcal{G}) \wedge (\mu + \mathbf{h}_\kappa \in \mathcal{G}_0),$$

wobei wir der Übersicht halber die Tilde weglassen und die gestrichenen Größen in Gleichung (6.1) durch ungestrichene Größen ersetzen. Zu beachten ist, dass die Struktur für ein beliebiges mehrdimensionales Wellendigitalfilter mit der in Gleichung (2.49) gewählten Transformationsmatrix \mathbf{H} gleich ist. Diese Regelmäßigkeit eröffnet gerade die Möglichkeit einer strukturierten Programmsynthese und darüber hinaus eine systematische und allgemeingültige Verifizierbarkeit der Code-Struktur gegenüber der Spezifikation.

Für den Übergang von dimensionsbehafteten zu dimensionslosen Größen definieren wir

$$\begin{aligned} \mathbf{u}_N &= \frac{\mathbf{u}}{u_B}, & u_B &\dots \text{Bezugsspannung} \\ \mathbf{i}_N &= \frac{\mathbf{i}}{i_B}, & i_B &\dots \text{Bezugsstrom} \\ \mathbf{R}_N &= \frac{\mathbf{R}}{R_B}, & R_B &\dots \text{Bezugswiderstand} \\ \mathbf{G}_N &= \mathbf{G} R_B \\ \mathbf{b}_N &= \frac{\mathbf{b}}{b_B}, & b_B &\dots \text{Bezugswelle} \\ \mathbf{a}_N &= \frac{\mathbf{a}}{b_B}. \end{aligned} \quad (6.6)$$

Die Wellengrößen sind mit Kirchhoff'schen Größen über Gleichung (2.66) miteinander verknüpft. Wir fordern, dass diese Beziehung für normierte Größen erhalten bleiben soll, d. h.

$$\begin{bmatrix} \mathbf{a}_N \\ \mathbf{b}_N \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{G}_N^{1/2} & \mathbf{R}_N^{1/2} \\ \mathbf{G}_N^{1/2} & -\mathbf{R}_N^{1/2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_N \\ \mathbf{i}_N \end{bmatrix}, \quad (6.7)$$

$$\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} = b_B \begin{bmatrix} \mathbf{a}_N \\ \mathbf{b}_N \end{bmatrix} = b_B \frac{1}{2} \begin{bmatrix} \mathbf{G}_N^{1/2} & \mathbf{R}_N^{1/2} \\ \mathbf{G}_N^{1/2} & -\mathbf{R}_N^{1/2} \end{bmatrix} \begin{bmatrix} \mathbf{u}_N \\ \mathbf{i}_N \end{bmatrix} = b_B \frac{1}{2} \begin{bmatrix} \mathbf{G}_N^{1/2}/u_B & \mathbf{R}_N^{1/2}/i_B \\ \mathbf{G}_N^{1/2}/u_B & -\mathbf{R}_N^{1/2}/i_B \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{i} \end{bmatrix}, \quad (6.8)$$

$$\mathbf{G}^{1/2} = \mathbf{G}_N^{1/2} \frac{b_B}{u_B} \quad \text{und} \quad \mathbf{R}^{1/2} = \mathbf{R}_N^{1/2} \frac{b_B}{i_B}. \quad (6.9)$$

Ihre Multiplikation liefert

$$\mathbf{1} = \mathbf{G}^{1/2} \mathbf{R}^{1/2} = \mathbf{G}_N^{1/2} \frac{b_B}{u_B} \mathbf{R}_N^{1/2} \frac{b_B}{i_B} \Rightarrow b_B^2 = u_B i_B \quad (6.10)$$

und mit invertierter 2. Gleichung (6.9) liefert die Multiplikation

$$\mathbf{G} = \mathbf{G}_N \frac{i_B}{u_B} \Rightarrow R_B = \frac{u_B}{i_B}. \quad (6.11)$$

Somit sind nur 2 der 4 Bezugsgrößen frei wählbar. Im Folgenden werden wir den Algorithmus in den normierten Größen entwickeln. Um die Übersicht zu wahren, wollen wir allerdings den Zusatz N für die normierten Größen in dem Algorithmus weglassen.

Nun werden die Gleichungen (6.3) - (6.5) so umgeformt, dass die Allquantoren verschwinden. Insbesondere ist es das Ziel der folgenden Überlegungen, festzustellen, in welchen Intervallen die Variablen μ_1, μ_2 und μ_3 liegen, wenn die Vektoren $\boldsymbol{\mu}$ den beiden folgenden Restriktionen unterliegen:

- a) $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}_0)$
- b) $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G})$.

Diese Bereiche geben uns an, für welche Werte der Variablen μ_1, μ_2 und μ_3 wir welche Gleichung zu erfüllen haben.

Gemäß Kapitel 2.4 sind \mathcal{G} und \mathcal{G}_0 durch

$\boldsymbol{\mu} \in \mathcal{G} : 0 \leq \mu_1 \leq P_{x_1} - 1 \wedge 0 \leq \mu_2 \leq P_{x_2} - 1 \wedge 0 \leq \mu_3 \leq P_{x_3} - 1$ und

$\boldsymbol{\mu} \in \mathcal{G}_0 : (0 > \mu_1 \vee \mu_1 > P_{x_1} - 1) \vee (0 > \mu_2 \vee \mu_2 > P_{x_2} - 1) \vee (0 > \mu_3 \vee \mu_3 > P_{x_3} - 1)$

definiert. Zur weiteren Beantwortung der Frage betrachten wir als Vorüberlegung die Koordinate σ des Vektors $\boldsymbol{\mu}$, wobei später die Koordinaten für $\sigma = 1, 2, 3, 4$ miteinander verknüpft werden müssen um das Gebiet zu beschreiben. Hierbei unterscheiden wir danach, ob das Element der Matrix \mathbf{H} den Wert 1, -1 oder 0 besitzt.

1. $h_{\sigma\kappa} = 1$

$$\begin{aligned} \text{a) } & (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \wedge (0 > \mu_\sigma + 1 \vee \mu_\sigma + 1 > P_{x_\sigma} - 1) \equiv \mu_\sigma = P_{x_\sigma} - 1 \\ \text{b) } & (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \wedge (0 \leq \mu_\sigma + 1 \leq P_{x_\sigma} - 1) \equiv 0 \leq \mu_\sigma \leq P_{x_\sigma} - 2 \end{aligned} \quad (6.12)$$

2. $h_{\sigma\kappa} = -1$

$$\begin{aligned} \text{a) } & (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \wedge (0 > \mu_\sigma - 1 \vee \mu_\sigma - 1 > P_{x_\sigma} - 1) \equiv \mu_\sigma = 0 \\ \text{b) } & (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \wedge (0 \leq \mu_\sigma - 1 \leq P_{x_\sigma} - 1) \equiv 1 \leq \mu_\sigma \leq P_{x_\sigma} - 1 \end{aligned} \quad (6.13)$$

3. $h_{\sigma\kappa} = 0$

$$\text{b) } (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \wedge (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1) \equiv (0 \leq \mu_\sigma \leq P_{x_\sigma} - 1)$$

a) Dieser Fall ist wesentlich komplizierter. Wir müssen hier die Fälle unterscheiden, in denen die anderen beiden der ersten 3 Spaltenelemente gleich null sind und die, in denen sie es nicht sind.

κ	$h_{\sigma\kappa}$	a) $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}_0)$	b) $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G})$
1	$h_{11} = 1$ $h_{21} = 0$ $h_{31} = 0$	$\mu_1 = P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$	$0 \leq \mu_1 \leq P_{x_1} - 2$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$
2	$h_{12} = 0$ $h_{22} = 1$ $h_{32} = 0$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $\mu_2 = P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 2$ $0 \leq \mu_3 \leq P_{x_3} - 1$
3	$h_{13} = 0$ $h_{23} = 0$ $h_{33} = 1$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $\mu_3 = P_{x_3} - 1$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 2$
4	$h_{14} = -1$ $h_{24} = 0$ $h_{34} = 0$	$\mu_1 = 0$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$	$1 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$
5	$h_{15} = 0$ $h_{25} = -1$ $h_{35} = 0$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $\mu_2 = 0$ $0 \leq \mu_3 \leq P_{x_3} - 1$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $1 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$
6	$h_{16} = 0$ $h_{26} = 0$ $h_{36} = -1$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $\mu_3 = 0$	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $1 \leq \mu_3 \leq P_{x_3} - 1$
7	$h_{17} = 0$ $h_{27} = 0$ $h_{37} = 0$	<i>FALSE</i>	$0 \leq \mu_1 \leq P_{x_1} - 1$ $0 \leq \mu_2 \leq P_{x_2} - 1$ $0 \leq \mu_3 \leq P_{x_3} - 1$ $0 \leq \delta \leq P_{x_1}P_{x_2}P_{x_3} - 1$

Tabelle 6.1: Grenzen

1) Wir wollen dies am Beispiel $\kappa = 1$ diskutieren. Hier gilt $h_{11} = 1$ und für die anderen beiden Spalten-elemente $h_{12} = h_{13} = 0$. Für das Prädikat $\boldsymbol{\mu} + \mathbf{h}_1 \in \mathcal{G}_0$ ergibt sich als konjunktive Bedingung

$$(0 \leq \mu_2 \leq P_{x_2} - 1) \wedge (0 \leq \mu_3 \leq P_{x_3} - 1) . \quad (6.14)$$

Das Prädikat $\boldsymbol{\mu} + \mathbf{h}_1 \in \mathcal{G}_0$ ist mit $\boldsymbol{\mu} \in \mathcal{G}$ konjunktiv zu verknüpfen und liefert für die Spalten 2 und 3

$$(0 \leq \mu_2 \leq P_{x_2} - 1) \wedge (0 \leq \mu_3 \leq P_{x_3} - 1) . \quad (6.15)$$

Hinzu kommt die konjunktive Verknüpfung mit den in Gleichung (6.12) ermittelten Ergebnissen für die erste Spalte, d.h. $\mu_1 = P_{x_1} - 1$.

2) Nun betrachten wir das Beispiel $\kappa = 7$. Hier gilt neben $h_{17} = 0$ auch für die anderen beiden Spalten-elemente $h_{72} = h_{73} = 0$.

Da \mathcal{G} und \mathcal{G}_0 für ein konstantes μ_4 disjunkt sind, gilt

$$(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_7 \in \mathcal{G}_0) \equiv \text{FALSE} \quad (6.16)$$

Die Bedeutung dieses Resultates ist, dass für den in 1) genannten Bereich die Aussage wahr ist und für den in 2) genannten Bereich die Aussage falsch ist.

Für jedes κ und der Kenntnis von \mathbf{H} muss jeweils für ein Spaltenelement $0, -1, 1$ unterschieden werden. Das ergibt die in der Tabelle 6.1 dargestellten Grenzen.

Bemerkung : Sowohl an $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}_0) \vee (\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}) \equiv \boldsymbol{\mu} \in \mathcal{G}$ (die disjunktive Verknüpfung wirkt wie eine Vereinigung der Mengen), als auch an den Werten in der Tabelle ist zu

erkennen, dass sämtliche Punkte, für die $\boldsymbol{\mu} \in \mathcal{G}$ gilt, abgearbeitet werden. Mit der Tabelle sind wir später in der Lage, den Schleifenanfang und die Schleifenabbruchbedingung in den Variablen μ_1, μ_2, μ_3 festzulegen.

Wir wollen nun die Auswirkungen einer Verschiebung des Vektors $\boldsymbol{\mu}$ nach $\boldsymbol{\mu} + \mathbf{h}_\kappa$ auf δ untersuchen. Konkret möchten wir die wie folgt definierten δ_κ bestimmen

$$\delta(\boldsymbol{\mu}) + \delta_\kappa = \delta(\boldsymbol{\mu} + \mathbf{h}_\kappa) . \quad (6.17)$$

Nach Gleichung (2.230) gilt

$$\delta(\boldsymbol{\mu}) = \mu_1 + \mu_2 P_{x_1} + \mu_3 P_{x_1} P_{x_2} = \boldsymbol{\mu}^T \mathbf{c} \quad , \quad \mathbf{c}^T = [1, P_{x_1}, P_{x_1} P_{x_2}] \quad (6.18)$$

und es folgt

$$\delta_\kappa = \delta(\boldsymbol{\mu} + \mathbf{h}_\kappa) - \delta(\boldsymbol{\mu}) = \boldsymbol{\mu}^T \mathbf{c} + \mathbf{h}_\kappa^T \mathbf{c} - \boldsymbol{\mu}^T \mathbf{c} = \mathbf{h}_\kappa^T \mathbf{c} . \quad (6.19)$$

Im Einzelnen ergeben sich die folgenden Werte

$$\delta_1 = -\delta_4 = 1 \quad , \quad \delta_2 = -\delta_5 = P_{x_1} \quad , \quad \delta_3 = -\delta_6 = P_{x_1} P_{x_2} \quad , \quad \delta_7 = 0 . \quad (6.20)$$

Wir führen nun eine Reihe von Prädikaten ein, deren Gegenstände der Ortspunkt des Sollwertes, ausgedrückt durch μ_1, μ_2, μ_3 und der Ortspunkt des Ist-Wertes δ ist. Die Wahrheit eines Prädikates entspricht der Tatsache, dass der Wert der Programmvariablen (der Ist-Wert) identisch mit dem geforderten Wert (dem Sollwert) ist. Wir benötigen jeweils ein Prädikat für die ausfallenden Wellen der nichtdynamischen Bauelemente

$$P_{ek}(\mu_1, \mu_2, \mu_3, \delta) \equiv b_{ek}(\boldsymbol{\mu}) = b_{ek}(\delta).v \quad , \quad (6.21)$$

die Ausgangswellen der Quellen

$$P_{qk}(\mu_1, \mu_2, \mu_3, \delta) \equiv b_{qk}(\boldsymbol{\mu}) = b_{qk}(\delta).v \quad , \quad (6.22)$$

die Ausgangssignale des Funktionsbausteins

$$P_{aFBk}(\boldsymbol{\mu}, \delta) \equiv a_{FBk}(\boldsymbol{\mu}) = a_{FBk}(\delta).v \quad , \quad (6.23)$$

die einfallenden Wellen der Verzögerer am Rand

$$P_{ak}^\kappa(\mu_1, \mu_2, \mu_3, \delta) \equiv a_{vk}^\kappa(\boldsymbol{\mu}) = a_{vk}^\kappa(\delta).v \quad (6.24)$$

und die ausfallenden Wellen der Verzögerer

$$P_{bk}^\kappa(\mu_1, \mu_2, \mu_3, \delta) \equiv b_{vk}^\kappa(\boldsymbol{\mu}) = b_{vk}^\kappa(\delta).v . \quad (6.25)$$

Diese Prädikate nutzen wir, um weitere Prädikate zu definieren. Wir verknüpfen konjunktiv die soeben definierten Prädikate, ausgewertet an Ortspunkten eines Gebietes, um die Gleichheit von Soll- und Istwert in einem Gebiet auszudrücken. Zunächst definieren wir Prädikate für die ausfallenden Wellen der

Verzögerer innerhalb des Gebietes $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G})$

$$\begin{aligned}
P_{\text{norm}}^1 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-2} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1+1, \mu_2, \mu_3, \delta+\delta_1) \\
P_{\text{norm}}^2 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-2} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^2} P_{bk}^2(\mu_1, \mu_2+1, \mu_3, \delta+\delta_2) \\
P_{\text{norm}}^3 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-2} \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1, \mu_2, \mu_3+1, \delta+\delta_3) \\
P_{\text{norm}}^4 &\equiv \bigwedge_{\mu_1=1}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^4} P_{bk}^4(\mu_1-1, \mu_2, \mu_3, \delta+\delta_4) \\
P_{\text{norm}}^5 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=1}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^5} P_{bk}^5(\mu_1, \mu_2-1, \mu_3, \delta+\delta_5) \\
P_{\text{norm}}^6 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=1}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^6} P_{bk}^6(\mu_1, \mu_2, \mu_3-1, \delta+\delta_6) \\
P_{\text{norm}}^7 &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^7} P_{bk}^7(\mu_1, \mu_2, \mu_3, \delta+\delta_7) ,
\end{aligned} \tag{6.26}$$

oder kompakt

$$\begin{aligned}
P_{\text{norm}} &\equiv \bigwedge_{\kappa=1}^3 \bigwedge_{-4\boldsymbol{\mu}=\mathbf{0}}^{\mathbf{Px}-\mathbf{e}-4\mathbf{h}_\kappa} \bigwedge_{k=1}^{n_v^\kappa} P_{bk}^\kappa(\mu_1+h_{1\kappa}, \mu_2+h_{2\kappa}, \mu_3+h_{3\kappa}, \delta+\delta_\kappa) \wedge \\
&\bigwedge_{\kappa=4}^7 \bigwedge_{-4\boldsymbol{\mu}=-4\mathbf{h}_\kappa}^{\mathbf{Px}-\mathbf{e}} \bigwedge_{k=1}^{n_v^\kappa} P_{bk}^\kappa(\mu_1+h_{1\kappa}, \mu_2+h_{2\kappa}, \mu_3+h_{3\kappa}, \delta+\delta_\kappa) .
\end{aligned} \tag{6.27}$$

Weiterhin definieren wir Prädikate für die ausfallenden Wellen der Verzögerer des Gebietes $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}_0)$

$$\begin{aligned}
P_x(\mu_2, \mu_3) &\equiv \bigwedge_{k=1}^{n_v^1} P_{bk}^1(0, \mu_2, \mu_3, \delta([0, \mu_2, \mu_3]^T)) \wedge \bigwedge_{k=1}^{n_v^4} P_{bk}^4(P_{x_1}-1, \mu_2, \mu_3, \delta([P_{x_1}-1, \mu_2, \mu_3]^T)) \\
P_y(\mu_1, \mu_3) &\equiv \bigwedge_{k=1}^{n_v^2} P_{bk}^2(\mu_1, 0, \mu_3, \delta([\mu_1, 0, \mu_3]^T)) \wedge \bigwedge_{k=1}^{n_v^5} P_{bk}^5(\mu_1, P_{x_2}-1, \mu_3, \delta([\mu_1, P_{x_2}-1, \mu_3]^T)) \\
P_z(\mu_1, \mu_2) &\equiv \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1, \mu_2, 0, \delta([\mu_1, \mu_2, 0]^T)) \wedge \bigwedge_{k=1}^{n_v^6} P_{bk}^6(\mu_1, \mu_2, P_{x_3}-1, \delta([\mu_1, \mu_2, P_{x_3}-1]^T)) .
\end{aligned} \tag{6.28}$$

Der nächste große Block von Prädikatsdefinitionen gilt den einfallenden Wellen der Verzögerer des Ge-

bietes $(\boldsymbol{\mu} \in \mathcal{G}) \wedge (\boldsymbol{\mu} + \mathbf{h}_\kappa \in \mathcal{G}_0)$

$$\begin{aligned}
P_{ax}(\mu_2, \mu_3) &\equiv \bigwedge_{k=1}^{n_v^1} P_{ak}^1(P_{x_1} - 1, \mu_2, \mu_3, \delta([P_{x_1} - 1, \mu_2, \mu_3]^T)) \wedge \bigwedge_{k=1}^{n_v^4} P_{ak}^4(0, \mu_2, \mu_3, \delta([0, \mu_2, \mu_3]^T)) \\
P_{ay}(\mu_1, \mu_3) &\equiv \bigwedge_{k=1}^{n_v^2} P_{ak}^2(\mu_1, P_{x_2} - 1, \mu_3, \delta([\mu_1, P_{x_2} - 1, \mu_3]^T)) \wedge \bigwedge_{k=1}^{n_v^5} P_{ak}^5(\mu_1, 0, \mu_3, \delta([\mu_1, 0, \mu_3]^T)) \\
P_{az}(\mu_1, \mu_2) &\equiv \bigwedge_{k=1}^{n_v^3} P_{ak}^3(\mu_1, \mu_2, P_{x_3} - 1, \delta([\mu_1, \mu_2, P_{x_3} - 1]^T)) \wedge \bigwedge_{k=1}^{n_v^6} P_{ak}^6(\mu_1, \mu_2, 0, \delta([\mu_1, \mu_2, 0]^T)) .
\end{aligned} \tag{6.29}$$

Über die ganzen Randflächen ergibt sich

$$\begin{aligned}
P'_{avx} &\equiv \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} P_{ax}(\mu_2, \mu_3) \\
P'_{avy} &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} P_{ay}(\mu_1, \mu_3) \\
P'_{avz} &\equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1, \mu_2) .
\end{aligned} \tag{6.30}$$

Außerdem benötigen wir Prädikate, die genau dann wahr sind, wenn die Werte der Variablen der Quellenwellen

$$P_q \equiv \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_q} P_{qk}(\mu_1(\delta), \mu_2(\delta), \mu_3(\delta), \delta) \tag{6.31}$$

und die Werte der Verzögerer

$$P_v \equiv \bigwedge_{\delta=0}^{P-1} \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} P_{bk}^\kappa(\mu_1(\delta), \mu_2(\delta), \mu_3(\delta), \delta) \tag{6.32}$$

den tatsächlichen Werten im gesamten Berechnungsgebiet entsprechen. Im Folgenden werden wir eine Reihe von häufig verwendeten logischen Prädikaten definieren, die uns später bei dem eigentlichen Korrektheitsbeweis eine einfachere Schreibweise ermöglichen. Hierbei dient

$$P_{dec} \equiv P_{decbq} \wedge P_{decbv} \wedge P_{decFB} \wedge P_{decuv} \wedge P_{decbe} \wedge P_{decav} \tag{6.33}$$

mit

$$\begin{aligned}
P_{\text{dec bq}} &\equiv \bigwedge_{n=1}^{n_q} [b_{qn}.typ = \text{float} \wedge P \leq b_{qn}.l \leq \text{Max}_i] \\
P_{\text{dec bv}} &\equiv \bigwedge_{\kappa=1}^7 \bigwedge_{n=1}^{n_v^\kappa} [b_{vn}^\kappa.typ = \text{float} \wedge P \leq b_{vn}^\kappa.l \leq \text{Max}_i] \\
P_{\text{dec be}} &\equiv \bigwedge_{n=1}^{n_{aa}} [b_{en}.typ = \text{float} \wedge P \leq b_{en}.l \leq \text{Max}_i] \\
P_{\text{deca FB}} &\equiv \bigwedge_{n=1}^6 [a_{FBn}.typ = \text{float} \wedge P \leq a_{FBn}.l \leq \text{Max}_i] \\
P_{\text{decav}} &\equiv \bigwedge_{\kappa=1}^6 \bigwedge_{n=1}^{n_v^\kappa} [a_{vn}^\kappa.typ = \text{float} \wedge P \leq a_{vn}^\kappa.l \leq \text{Max}_i] \\
P_{\text{decuv}} &\equiv \delta.typ = \text{int} \wedge \delta.l = -1 \wedge \bigwedge_{\kappa=1}^3 [\mu_\kappa.typ = \text{int} \wedge \mu_\kappa.l = -1]
\end{aligned} \tag{6.34}$$

für die Deklariertheit von vektoriellen Variablen der richtigen Länge. Dabei unterscheiden wir zwischen denen, die vor der Ausführung von S_{alg} schon deklariert sind und denen, die noch deklariert werden müssen. Das wichtigste Prädikat ist das korrekte WDF, welches aus Gleichung (6.1) bis Gleichung (6.5) entsteht

$$\begin{aligned}
P_{\text{WDF}} &\equiv \bigwedge_{\mu=0}^{P_{x-e}} b_e(\mu) = L_q b_q(\mu) + L_v b_v(\mu) + L_e b_e(\mu) \wedge a_{\text{FB}}(\mu) = Ab(\mu) \\
&\bigwedge_{\kappa=1}^3 \bigwedge_{-4\mu=0}^{P_{x-e}-4h_\kappa} \bigwedge_{k=1}^{n_v^\kappa} b_v^\kappa(\mu + h_\kappa) = P_{v^\kappa e} b_e(\mu) \wedge \bigwedge_{\kappa=4}^7 \bigwedge_{-4\mu=-4h_\kappa}^{P_{x-e}} \bigwedge_{k=1}^{n_v^\kappa} b_v^\kappa(\mu + h_\kappa) = P_{v^\kappa e} b_e(\mu) \\
&\bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} a_v^1(\mu) = P_{v^1 e} b_e(\mu) \wedge b_v^1(\mu + e_7) = R^1 a_v^4(\mu) \wedge \mu_1 = P_{x_1} - 1 \\
&\bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} a_v^2(\mu) = P_{v^2 e} b_e(\mu) \wedge b_v^2(\mu + e_7) = R^2 a_v^5(\mu) \wedge \mu_2 = P_{x_2} - 1 \\
&\bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} a_v^3(\mu) = P_{v^3 e} b_e(\mu) \wedge b_v^3(\mu + e_7) = R^3 a_v^6(\mu) \wedge \mu_3 = P_{x_3} - 1 \\
&\bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} a_v^6(\mu) = P_{v^6 e} b_e(\mu) \wedge b_v^6(\mu + e_7) = R^6 a_v^3(\mu) \wedge \mu_3 = 0
\end{aligned}$$

$$\begin{aligned}
& \bigwedge_{\mu=1}^{n_e} \bigwedge_{\nu=\mu}^{n_e} l_{e\mu\nu} = 0 \wedge \|x\| \geq \|Lx\| \forall x \wedge \bigwedge_{\kappa=1}^6 \|x\| \geq \|R^\kappa x\| \forall x \wedge \\
& \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} a_v^4(\mu) = P_{v^4e} b_e(\mu) \wedge b_v^4(\mu + e_7) = R^4 a_v^1(\mu) \wedge \mu_1 = 0 \\
& \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} a_v^5(\mu) = P_{v^5e} b_e(\mu) \wedge b_v^5(\mu + e_7) = R^5 a_v^2(\mu) \wedge \mu_2 = 0 \\
& \bigwedge_{\kappa=1}^3 [\text{Max}_i \geq P_{x_\kappa} > 0 \wedge P_{x_\kappa} \in \mathbb{N}] \wedge P = P_{x_1} P_{x_2} P_{x_3} \wedge \text{Max}_i \geq P .
\end{aligned} \tag{6.35}$$

Weiterhin definieren wir ein Prädikat, welches im Zusammenhang mit der Überschreitung des Wertebereichs der abhängigen Variablen verwendet wird

$$P_{\max} \equiv \bigwedge_{\delta=0}^{P-1} \sum_{\nu=1}^{n_q} |b_{q\nu}(\delta).v|^2 + \sum_{\nu=1}^{n_v} |b_{v\nu}(\delta).v|^2 < \text{Max}_f^2 / (\beta^2 4) , \tag{6.36}$$

wobei die Konstante $\beta > 0$ sich aus

$$\beta^2 = \max_{\mu} \{ \max \{ R_{N\mu}, G_{N\mu} \} \} \tag{6.37}$$

ergibt und $R_{N\mu}, G_{N\mu}$ die (normierten) Torwiderstände bzw. Torleitwerte sind. Unter Berücksichtigung von Gleichung (A.4) gilt $\beta \geq 1$. Somit kann das Radizieren in den Operator hineingezogen werden, so dass

$$\beta = \max_{\mu} \left\{ \max \left\{ \sqrt{R_{N\mu}}, \sqrt{G_{N\mu}} \right\} \right\} \tag{6.38}$$

entsteht. Zudem benötigen wir

$$P_{\text{avx max}} \equiv \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \left[\bigwedge_{\nu=1}^{n_v^1} |a_{v\nu}^1(\delta([0, \mu_2, \mu_3]^T))| < \text{Max}_f \wedge \bigwedge_{\nu=1}^{n_v^4} |a_{v\nu}^4(\delta([P_{x_1}-1, \mu_2, \mu_3]^T))| < \text{Max}_f \right] ,$$

$$P_{\text{avy max}} \equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \left[\bigwedge_{\nu=1}^{n_v^2} |a_{v\nu}^2(\delta([\mu_1, 0, \mu_3]^T))| < \text{Max}_f \wedge \bigwedge_{\nu=1}^{n_v^5} |a_{v\nu}^5(\delta([\mu_1, P_{x_2}-1, \mu_3]^T))| < \text{Max}_f \right] ,$$

$$P_{\text{avz max}} \equiv \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \left[\bigwedge_{\nu=1}^{n_v^3} |a_{v\nu}^3(\delta([\mu_1, \mu_2, 0]^T))| < \text{Max}_f \wedge \bigwedge_{\nu=1}^{n_v^6} |a_{v\nu}^6(\delta([\mu_1, \mu_2, P_{x_3}-1]^T))| < \text{Max}_f \right] ,$$

$$P_{\text{av max}} \equiv P_{\text{avx max}} \wedge P_{\text{avy max}} \wedge P_{\text{avz max}} \text{ und}$$

$$P_{\text{be max}} \equiv \bigwedge_{\delta=0}^{P-1} \bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta).v| 2\beta < \text{Max}_f .$$

(6.39)

Zum Abschluss benötigen wir für die korrekte Abbildung des Vektors μ ins Eindimensionale das Prädikat

$$P_{\text{MDED}} \equiv \mu_1 + \mu_2 P_{x_1} + \mu_3 P_{x_1} P_{x_2} = \delta(\mu). \quad (6.40)$$

Im engen Zusammenhang damit steht der folgende

Satz 2 Sei $\{V\} \text{ delta} = \mu_1 + \mu_2 * P_{x_1} + \mu_3 * P_{x_1} * P_{x_2} ; \{P\}$.

Wenn dazu die Nachbedingung $P \equiv P' \wedge 0 \leq \delta.v \leq P - 1 \wedge P_{\text{dec}}$ ist, dann ist

$$0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1 \wedge P' \left[\begin{smallmatrix} \delta.v \\ \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \end{smallmatrix} \right] \wedge P_{\text{dec}} \quad (6.41)$$

eine gültige Vorbedingung $\{V\}$.

Beweis 2

Durch Verwendung von Gleichung (2.233) i.V.m. Gleichung (2.230) erhalten wir

$$0 \leq \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \leq P - 1 \wedge P' \left[\begin{smallmatrix} \delta.v \\ \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \end{smallmatrix} \right] \wedge$$

$$P_{\text{dec}} \wedge \delta.l = -1 \wedge \delta.\text{typ} = \text{int} \wedge |\mu_1.v| \leq \text{Max}_i \wedge |\mu_1.v + P_{x_1} \mu_2.v| \leq \text{Max}_i \wedge \quad (6.42)$$

$$|\mu_1.v + P_{x_1} \mu_2.v + P_{x_1} P_{x_2} \mu_3.v| \leq \text{Max}_i \wedge \bigwedge_{l=1}^3 [\mu_l.l = -1 \wedge \mu_l.\text{typ} = \text{int}] .$$

Zunächst stellen wir fest, dass die auftretenden Typ- und Längenbedingungen abgetrennt werden können, da sie von P_{dec} impliziert werden.

Hinreichend für $0 \leq \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \leq P - 1$ mit $P = P_{x_1} P_{x_2} P_{x_3}$ ist, wie im Kapitel 2.11 gezeigt

$$0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1 \quad (6.43)$$

Zudem erübrigt sich die weitere Betrachtung der Teilsommen, da die Summanden dieses Skalarproduktes nicht negativ sind. P_{dec} impliziert die Aussage $P - 1 \leq \text{Max}_i$, so dass mit $0 \leq \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \leq P - 1$ die Aussage $\mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} \leq \text{Max}_i$ auch aus Gleichung (6.42) abgetrennt werden kann.

Wir kommen nun zur eigentlichen Problemstellung zurück, der Formulierung eines Algorithmus zur Berechnung eines mehrdimensionalen Wellendigitalfilters, welcher sich in das Programmpaket SPACE einfügt. Hierbei berücksichtigen wir, dass der Aufruf der Funktionsbausteine durch die Ablaufumgebung erfolgt. Bild 6.1 zeigt den Teil der Ablaufumgebung, der die Funktionsbausteine aufruft. Die äußere nicht-terminierte Schleife stellt die Ablaufumgebung dar, die über die Funktionsplangruppen und den Funktionsplänen die einzelnen Funktionsbausteine aufruft. Die diskrete Variable μ_k des Wellendigitalfilters entspricht der diskreten Variablen der Ablaufumgebung, d.h. das Fortschreiten von einer Abtastschicht zur nächsten entspricht dem Aufruf des MDWDF-Funktionsbausteins, wie im Kapitel 3 gefordert. Der Funktionsplan umfasst den MDWDF-Funktionsbaustein und weitere Funktionsbausteine. Der Teil des Programmablaufplans innerhalb der gestrichelten Linie ist eine detailliertere Darstellung des mehrdimensionalen Wellendigitalfilters. Hinter dem Anweisungsblock zur Ermittlung der Verzögererwerte der neuen Abtastschicht verbergen sich die im Bild 6.2 dargestellten Anweisungsblöcke, die wiederum später im Detail exemplarisch im Bild 6.5 dargestellt werden. Beim Aufruf des MDWDF-Funktionsbausteins

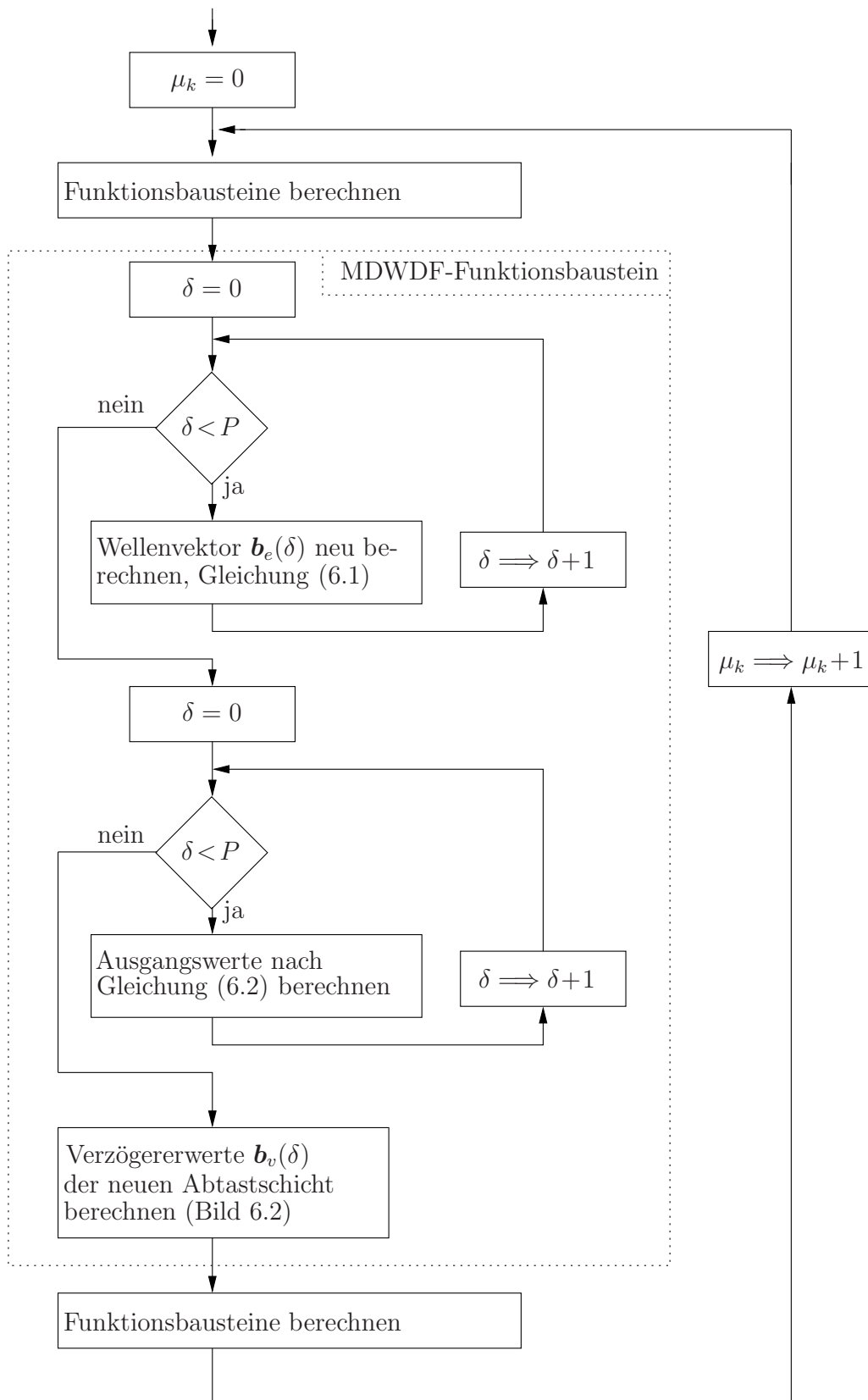


Bild 6.1: Einbindung eines MDWDF-Funktionsbausteins in eine Ablaufumgebung

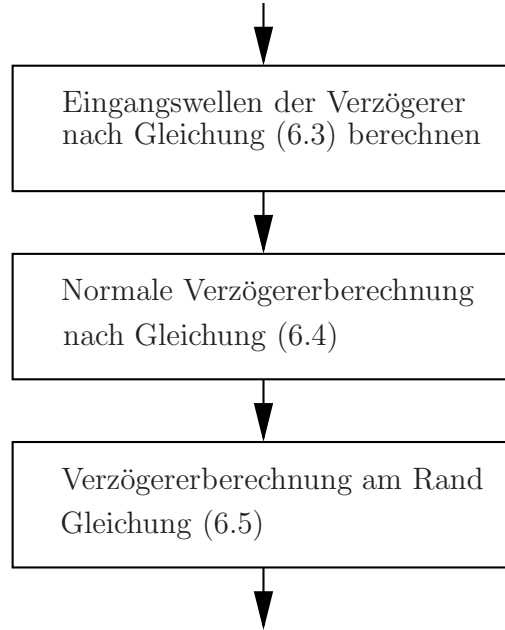


Bild 6.2: Verzögererwerte $\mathbf{b}_v(\delta)$ der neuen Abtastschicht ermitteln

werden zunächst die Wellengrößen der nichtdynamischen Elemente nach Gleichung (6.1) berechnet. Im Anschluss können die Ausgangswerte des Funktionsbausteins nach Gleichung (6.2) ermittelt werden. Darauf folgend berechnen wir die Eingangswellen der Verzögerer am Rand gemäß Gleichung (6.3).

Wir wollen an dieser Stelle noch einen wichtigen Aspekt erläutern. Bis einschließlich Gleichung (6.3) wurde nur lesend auf die Ausgangswellen der Verzögerer $\mathbf{m}_{\text{FB}} = \mathbf{b}_v([\bullet, \bullet, \bullet, \mu_k.v]^T)$ zugegriffen. Nach Gleichung (6.3) wird nicht mehr lesend darauf zugegriffen, d. h., ab hier werden die Verzögererwerte der aktuellen Abtastschicht nicht mehr benötigt. Bei der Berechnung der Verzögererwerte der nächsten Abtastschicht starten wir mit den Verzögererwellen, die keine Randbehandlung benötigen, d. h. also mit Gleichung (6.4). Hierbei betrachten wir nun \mathbf{m}_{FB} als die Verzögererwerte der nächsten Abtastschicht, d. h. es gilt $\mathbf{m}_{\text{FB}} = \mathbf{b}_v([\bullet, \bullet, \bullet, \mu_k.v + 1]^T)$. Der Zugriff auf \mathbf{m}_{FB} erfolgt nun schreibend. Dass wir die Verzögererwerte der nächsten Abtastschicht berechnen, erkennen wir daran, dass $\mu_k.v$ immer noch den gleichen Wert hat und im Argument von \mathbf{b}_v der Wert $\mu_k.v + 1$ steht. Insbesondere ist zu beachten, dass erst nach Inkrementierung von $\mu_k.v$ die Beziehung $\mathbf{m}_{\text{FB}} = \mathbf{b}_v([\bullet, \bullet, \bullet, \mu_k.v]^T)$ gilt und erst dann \mathbf{m}_{FB} wieder benutzt werden darf.

Zum Abschluss berechnet man sequentiell die Verzögererwellen des Randes nach Gleichung (6.5). Außerhalb des Funktionsbausteins wird $\mu_k.v$ inkrementiert. Für den lesenden Zugriff auf \mathbf{b}_v gilt nun $\mathbf{m}_{\text{FB}} = \mathbf{b}_v([\bullet, \bullet, \bullet, \mu_k.v]^T)$.

Das erläuterte Vorgehen entspricht einem Grobentwurf, also einer Zerlegung der Gesamtanweisung S_{alg} in aufeinanderfolgende Teilanweisungen, d. h.

$$S_{\text{alg}} \equiv S_{\text{dekl}}; S_{\text{be}}; S_{\text{aus}}; S_{\text{av}}; S_{\text{norm}}; S_{\text{rand}}. \quad (6.44)$$

Aus bereits im Kapitel 5 diskutierten Gründen werden wir den Korrektheitsbeweis führen, indem wir eine Vorbedingung aus einer vorgegebenen Nachbedingung ermitteln.¹ Insofern ist es sinnvoll, an dieser Stelle die Nachbedingung des Algorithmus formal zu spezifizieren, um auf dieser Basis dann die

¹Die Bedingungen stellen Prädikate dar, die von ungebundenen Variablen abhängen. Streng genommen müssten die Bezeichner der Prädikate mit diesen Gegenständen erfolgen. Der Übersicht halber wollen wir dennoch darauf verzichten. Beispiel : Sei $P(x) \equiv x > 0$. Dann schreiben wir anstatt $P(x)$ nur P . Dieses unpräzise Vorgehen wirft aber insbesondere dann Probleme auf, wenn die Ersetzungsregel genutzt wird. Falls im weiteren Verlauf der Arbeit auf ein Prädikat die Ersetzungsregel angewendet wird, werden wir abweichend von der obigen Vereinbarung die präzise Darstellung wählen.

Teilanweisungen zu entwerfen. Die Nachbedingung an den Algorithmus stellt den Datenzustand des Programms dar, der nach Ausführung des Algorithmus erreicht werden soll. Wir führen hierfür weitere Zwischenbedingungen ein

Bedingung	Anweisung	Kommentar
$\{V\}$	S_{dekl}	Variablendeklarationen
$\{P_{\text{dekl}}\}$	S_{be}	Berechnung des Vektors \mathbf{b}_e
$\{P_{\text{be}}\}$	S_{aus}	Berechnung der Ausgangssignale \mathbf{a}_{FB}
$\{P_{\text{aus}}\}$	S_{av}	Berechnung der Verzögererwerte \mathbf{a}_v am Rand
$\{P_{\text{av}}\}$	S_{norm}	Berechnung der Verzögerer \mathbf{b}_v im Normalfall
$\{P_{\text{norm}}\}$	S_{rand}	Berechnung der Verzögererwerte \mathbf{b}_v am Rand
$\{P\}$		

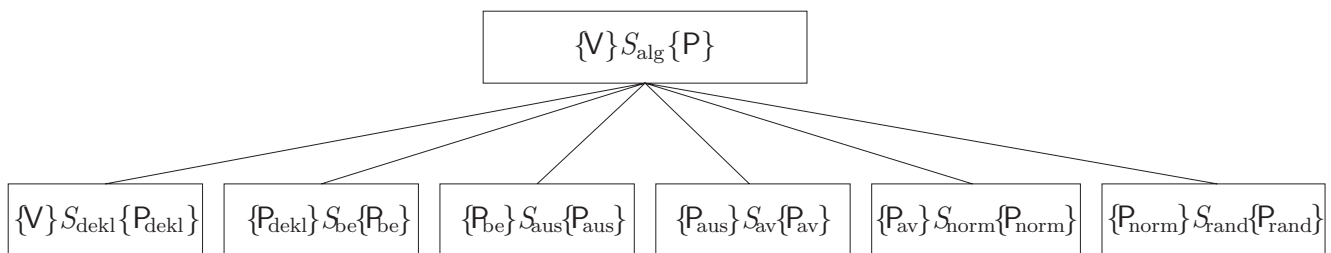


Bild 6.3: Übersicht über den Korrektheitsbeweis

Im Bild 6.3 ist der Algorithmus und die Zerlegung in sequentielle Teilanweisungen dargestellt. Dieses Diagramm und die folgenden Diagramme sind so zu verstehen, dass aus der Korrektheit aller untergeordneten Prädikate die Korrektheit des übergeordneten Prädikats folgt. Die einzelnen Teilanweisungen werden in den Kapiteln 6.2-6.7 festgelegt. Jede Nachbedingung dieser Teilanweisungen soll aus der konjunktiven Verknüpfung der Vorbedingung und einem bestimmten zusätzlichem Teil bestehen. Um diese zusätzlichen Nachbedingungen von den Nachbedingungen zu unterscheiden, kennzeichnen wir die Bezeichner der Prädikate dieser zusätzlichen Nachbedingungen mit einem '.

Mithilfe der definierten Prädikate können wir nun die zusätzlichen Nachbedingungen der 6 Teilanweisungen aus Gleichung (6.44) in der folgenden Tabelle darstellen

Anweisung	zusätzliche Nachbedingung
S_{dekl}	$P'_{\text{be}} \equiv P_{\text{decbe}} \wedge P_{\text{decav}} \wedge P_{\text{decuv}}$
S_{be}	$P'_{\text{be}} \equiv \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_e} P_{ek}(\mu_1(\delta), \mu_2(\delta), \mu_3(\delta), \delta) \wedge P_{\text{be max}}$
S_{aus}	$P'_{\text{aus}} \equiv \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_{aa}} P_{a\text{FB}k}(\mu(\delta), \delta)$
S_{av}	$P'_{\text{av}} \equiv P'_{\text{avx}} \wedge P'_{\text{avx}} \wedge P'_{\text{avx}} \wedge P_{\text{av max}}$
S_{norm}	$P'_{\text{norm}} \equiv \bigwedge_{\kappa=1}^7 P_{\text{norm}}^{\kappa}$
S_{rand}	$P'_{\text{rand}} \equiv \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} P_x(\mu_2, \mu_3) \wedge \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} P_y(\mu_1, \mu_3) \wedge \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1, \mu_2)$

Die Nachbedingung des Gesamtalgorithmus ist nun die konjunktive Verknüpfung mehrerer dieser Bedingungen, d.h.

$$P \equiv P'_{\text{aus}} \wedge P'_{\text{norm}} \wedge P'_{\text{randx}} \wedge P'_{\text{randy}} \wedge P'_{\text{randz}} . \quad (6.45)$$

Die in der Tabelle auftretenden und nicht in P berücksichtigten Bedingungen können beispielsweise die Speicherung von Zwischenwerten fordern.

Bevor wir zur eigentlichen Beweisführung übergehen, geben wir die Grobübersicht noch einmal an, diesmal aber mit den Zwischenbedingungen

Bedingung	Anweisung
$\{V\} \equiv P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{decbq}} \wedge P_{\text{decbv}} \wedge P_{\text{decaFB}} \wedge P_{\text{name}}$	S_{dekl}
$\{P_{\text{dekl}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_q \wedge P_v$	S_{be}
$\{P_{\text{be}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_q \wedge P_v \wedge P'_{\text{be}}$	S_{aus}
$\{P_{\text{aus}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_q \wedge P'_{\text{be}} \wedge P'_{\text{aus}}$	S_{av}
$\{P_{\text{av}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_q \wedge P'_{\text{be}} \wedge P'_{\text{aus}} \wedge P'_{\text{av}}$	S_{norm}
$\{P_{\text{norm}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P'_{\text{aus}} \wedge P'_{\text{av}} \wedge P'_{\text{norm}}$	S_{rand}
$\{P\} \equiv P_{\text{dec}} \wedge P'_{\text{aus}} \wedge P'_{\text{norm}} \wedge P'_{\text{randx}} \wedge P'_{\text{randy}} \wedge P'_{\text{randz}}$	

6.2 Berechnung der Verzögererwerte im Randfall

Die Anweisung zur Berechnung der Verzögererwerte im Randfall S_{rand} zerlegen wir in

$$S_{\text{rand}} \equiv S_{\text{randx}}; S_{\text{randy}}; S_{\text{randz}} . \quad (6.46)$$

Wir werden nur die Anweisung S_{randz} ausführlich behandeln. Die anderen zwei Anweisungen sind zu S_{randz} analog, vgl. Bild 6.4. Die Nachbedingung der Anweisung S_{randz} entspricht der des gesamten Algorithmus S_{alg} , also P gemäß Gleichung (6.45). Der Beitrag, den die Anweisung S_{randz} zu P leistet, ist die Berechnung der Verzögererwerte am Rand $z = \text{const.}$, formal repräsentiert durch die zusätzliche Nachbedingung P'_{randz} . Um das geforderte Ziel zu erreichen, nutzen wir den im Bild 6.5 dargestellten

Algorithmus. Die Implementierung in C-Code lautet :

Anweisungssequenz $\{P_{\text{randy}}\}S_{\text{randz}}\{P_{\text{randz}}\}$ <pre> for(mu1 = 0 ; mu1 < P_{x₁} ; mu1=mu1+1){ for(mu2 = 0 ; mu2 < P_{x₂} ; mu2=mu2+1){ mu3 = 0; delta = mu1 + mu2 * P_{x₁} + mu3 * P_{x₁} * P_{x₂} ; S_{b_{v1}³} : S_{b_{vn₃³}} mu3 = P_{x₁} - 1; delta = mu1 + mu2 * P_{x₁} + mu3 * P_{x₁} * P_{x₂} ; S_{b_{v1}⁶} : S_{b_{vn₆⁶}} } }</pre>

wobei die Zuweisungsoperationen durch

$$\begin{aligned}
S_{b_{v1}^3} &\equiv b_{v_3_1}[\text{delta}] = r_{11}^3 * a_{v_6_1}[\text{delta}] + \dots + r_{1n_v^6}^3 * a_{v_6_n_v}[\text{delta}] ; \\
&\vdots \\
S_{b_{vn_3^3}^3} &\equiv b_{v_3_n_v^3}[\text{delta}] = r_{n_v^3 1}^3 * a_{v_6_1}[\text{delta}] + \dots + r_{n_v^3 n_v^6}^3 * a_{v_6_n_v}[\text{delta}] ; \\
S_{b_{v1}^6} &\equiv b_{v_6_1}[\text{delta}] = r_{11}^6 * a_{v_3_1}[\text{delta}] + \dots + r_{1n_v^3}^6 * a_{v_1_n_v^3}[\text{delta}] ; \\
&\vdots \\
S_{b_{vn_6^6}^6} &\equiv b_{v_6_n_v^6}[\text{delta}] = r_{n_v^6 1}^6 * a_{v_3_1}[\text{delta}] + \dots + r_{n_v^6 n_v^3}^6 * a_{v_3_n_v^3}[\text{delta}] ;
\end{aligned}$$

gegeben sind.

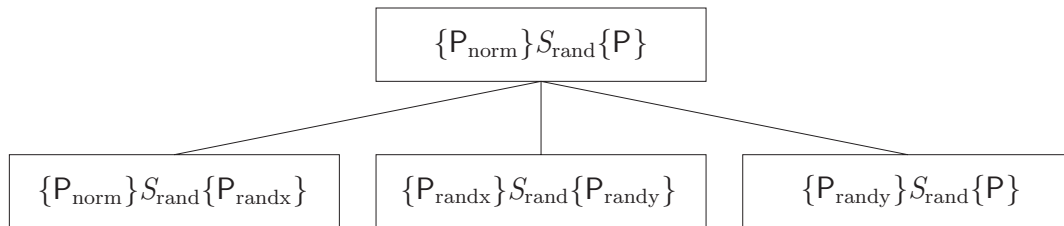


Bild 6.4: Die Berechnung der Verzögererwerte im Randfall

Wir ermitteln nun die Vorbedingung P_{randy} . Im Bild 6.6 sind die notwendigen Schritte zum Nachweis der Korrektheit der Anweisung S_{randz} dargestellt. Bei S_{randz} handelt es sich um zwei verschachtelte for-Schleifen mit den Schleifenbedingungen $b_a \equiv \mu_1.v < P_{x_1}$ und $b_i \equiv \mu_2.v < P_{x_2}$. Dem Bild 6.6 ist weiter zu entnehmen, dass jeweils 2 Mal Axiom (5.35) und Axiom (5.7) zur Anwendung kommen. Die verbleibende Aufgabe besteht nun darin, die Korrektheit der im Bild 6.6 dargestellten Anweisungen nachzuweisen. Vorab wollen wir uns Gedanken über die Reihenfolge der Beweisführung machen. Wir werden zuerst die Invariante der äußeren Schleife I_a in einem Satz angeben. Im Anschluss daran weisen wir durch die Initialisierungsanweisung die Vorbedingung

$$P_{\text{randy}} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P'_{\text{aus}} \wedge P'_{\text{av}} \wedge P'_{\text{norm}} \wedge P'_{\text{randx}} \wedge P'_{\text{randy}} \quad (6.47)$$

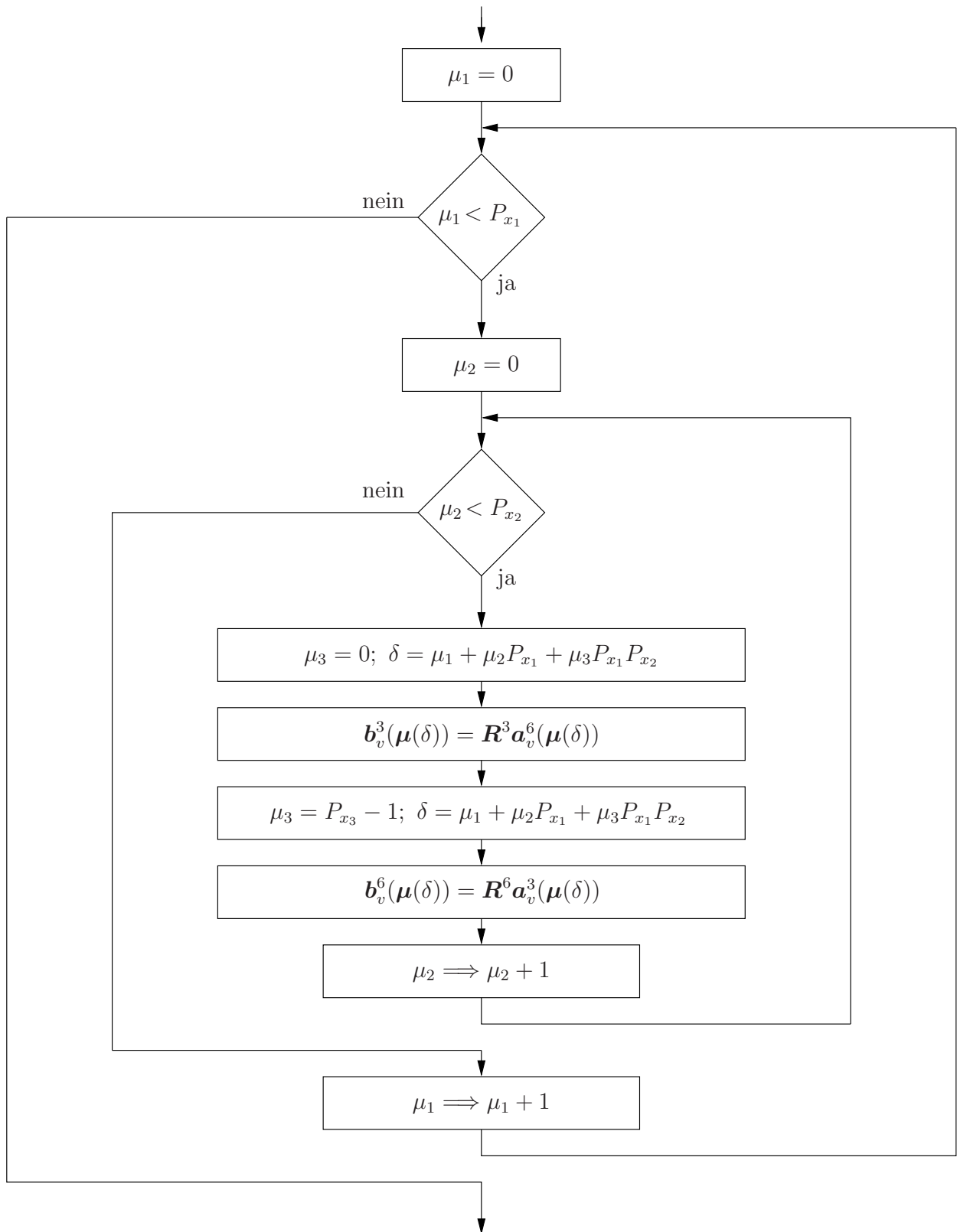


Bild 6.5: Berechnung der neuen Verzögererwerte \mathbf{b}_v^3 und \mathbf{b}_v^6 an den Grenzflächen $\mu_3 = 0$ und $\mu_3 = P_{x_3} - 1$

und über die Abbruchbedingung die Nachbedingung P nach. Das Zeigen der Invarianz von I_a bzgl. $S_{4a}; S_{3a}$ führt uns schließlich zur inneren Schleife.

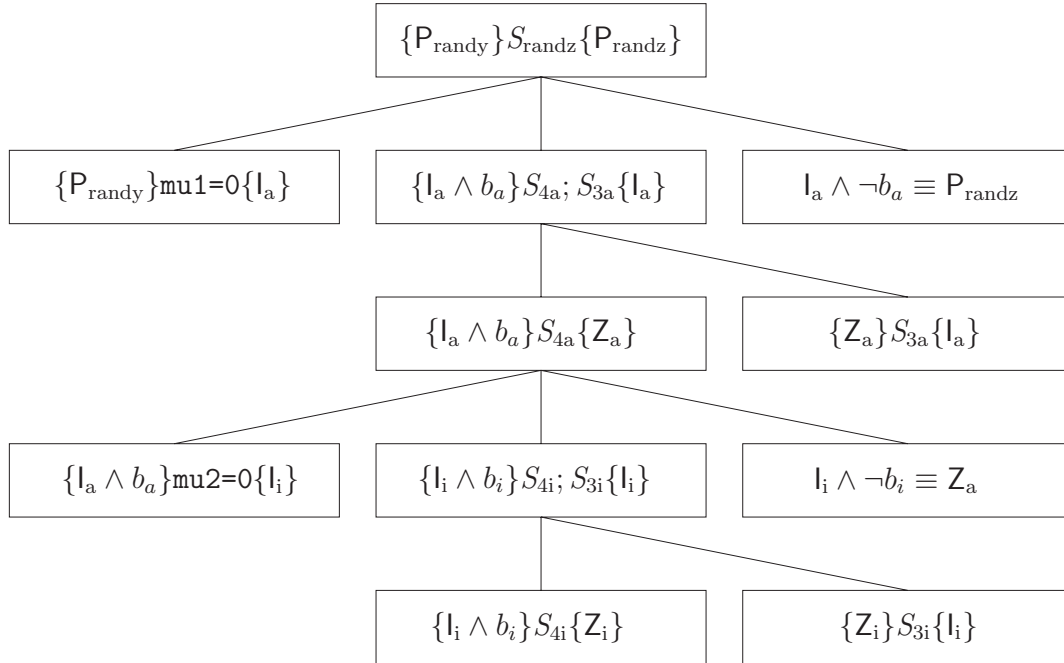


Bild 6.6: Übersicht zum Korrektheitsbeweis der Anweisung S_{randz} .

Satz 3 *Der Ausdruck*

$$I_a \equiv 0 \leq \mu_1.v \leq P_{x_1} \wedge P_{\text{randy}} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1, \mu_2) \quad (6.48)$$

ist eine gültige Schleifeninvariante der äußeren Schleife.

Beweis 3

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\text{mu1}=0$ liefert die Vorbedingung

$$0 \leq P_{x_1} \wedge P_{\text{randy}} \wedge \bigwedge_{\mu_1=0}^{-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1, \mu_2) \wedge 0 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.\text{typ} = \text{int} \equiv P_{\text{randy}}, \quad (6.49)$$

da P_{randy} den Ausdruck $\mu_1.\text{typ} = \text{int} \wedge \mu_1.l = -1$ impliziert. Die Nachbedingung bestimmt sich zu

$$I_a \wedge \mu_1.v \geq P_{x_1} \equiv P_{\text{randy}} \wedge \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1, \mu_2) \wedge \mu_1.v = P_{x_1} \equiv P_{\text{randz}} \wedge \mu_1.v = P_{x_1} \quad (6.50)$$

und kann gemäß Gleichung (5.10) zu P_{randz} geschwächt werden, weil der Wert der Variablen mu1 nicht weiter benötigt wird. Die Zwischenbedingung der äußeren Schleife bestimmt sich zu

$$Z_a \equiv 0 \leq \mu_1.v+1 \leq P_{x_1} \wedge P_{\text{randy}} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1, \mu_2) \wedge \mu_1.v+1 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.\text{typ} = \text{int}.$$

(6.51)

Wir stärken die Zwischenbedingung um $0 \leq \mu_1.v$ zu

$$Z'_a \equiv Z_a \wedge 0 \leq \mu_1.v \equiv l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1.v, \mu_2). \quad (6.52)$$

Der Beweis von $\{l_a \wedge b_a\} S_{4a} \{Z'_a\}$ wird erneut gemäß Axiom (5.35) durchgeführt, da die Anweisung S_{4a} die innere for-Schleife ist.

Satz 4 *Der Ausdruck*

$$l_i \equiv 0 \leq \mu_2.v \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v-1} P_z(\mu_1.v, \mu_2) \quad (6.53)$$

ist eine gültige Schleifeninvariante der inneren Schleife.

Beweis 4

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\text{mu2}=0$ liefert die behauptete Vorbedingung von S_{4a}

$$0 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{-1} P_z(\mu_1.v, \mu_2) \wedge 0 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.\text{typ} = \text{int} \equiv l_a \wedge b_a. \quad (6.54)$$

Die Nachbedingung der inneren Schleife lautet

$$Z_a \equiv l_i \wedge \mu_2.v \geq P_{x_2} \equiv l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_z(\mu_1.v, \mu_2). \quad (6.55)$$

Die Zwischenbedingung der inneren Schleife lautet

$$Z_i \equiv 0 \leq \mu_2.v + 1 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v} P_z(\mu_1.v, \mu_2) \wedge \mu_2.v + 1 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.\text{typ} = \text{int}. \quad (6.56)$$

Wir stärken Sie durch $0 \leq \mu_2.v$ zu

$$Z'_i \equiv Z_i \wedge 0 \leq \mu_2.v \equiv 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v} P_z(\mu_1.v, \mu_2) \equiv l_i \wedge b_i \wedge P_z(\mu_1.v, \mu_2.v). \quad (6.57)$$

Es verbleibt, die Richtigkeit von $\{l_i \wedge b_i\} S_{4i} \{Z'_i\}$ zu zeigen. Wir ermitteln dazu zunächst die Vorbedingung der Anweisung $S_{b_6}^{v n_v^6}$. Mit Satz (5.43) erhalten wir die Vorbedingung zu

$$\begin{aligned} & l_i \wedge b_i \wedge P_z(\mu_1.v, \mu_2.v) \left[\begin{matrix} b_v^6 n_v^6(\delta.v).v \\ \mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v \end{matrix} \right] \wedge |\mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v| \leq \text{Max}_f \wedge 0 \leq \delta.v \leq b_v^6 n_v^6.l - 1 \wedge \\ & \bigwedge_{\nu=1}^{n_v^3} 0 \leq \delta.v \leq a_{v\nu}^3.l - 1 \wedge \bigwedge_{\nu=1}^{n_v^3} b_v^6 n_v^6.\text{typ} = a_{v\nu}^3.\text{typ} \wedge \delta.l = -1 \wedge \delta.\text{typ} = \text{int}. \end{aligned} \quad (6.58)$$

Die Bedingungen an die Beschränktheit der Teilsummen des Skalarprodukts treten nicht auf, da nur eine Koordinate von $\mathbf{r}_{n_v^6, \bullet}^6$ von null verschieden ist. Das Bestreben ist es im Folgenden, die Teile dieser Aussage abzuspalten, die P_{randy} impliziert. Wir werden dabei mit den hinteren Teilaussagen anfangen, d. h. mit den Bedingungen an den Datentyp. Die Verwendung von Gleichung (6.48) i.V.m. Gleichung (6.33) und Gleichung (6.34) liefert

$$P_{\text{randy}} \Rightarrow b_{v n_v^6}^6 \cdot \text{typ} = \text{float} \wedge \bigwedge_{l=1}^{n_v^3} a_{v l}^3 \cdot \text{typ} = \text{float} \Rightarrow \bigwedge_{l=1}^{n_v^3} a_{v l}^3 \cdot \text{typ} = b_{v n_v^6}^6 \cdot \text{typ}. \quad (6.59)$$

Die Aussage

$$\bigwedge_{l=1}^{n_v^3} a_{v l}^3 \cdot \text{typ} = b_{v n_v^6}^6 \cdot \text{typ}$$

in Gleichung (6.58) ist somit wahr und kann gemäß $[(A \rightarrow B) \wedge A] \rightarrow B$ abgetrennt werden.

Um in der Vorbedingung Gleichung (6.58), die Bedingungen an die Länge der Felder abtrennen zu können, müssen wir uns eine Aussage über $\delta.v$ verschaffen. Dazu nehmen wir zunächst eine Stärkung von Gleichung (6.58) um $0 \leq \delta.v \leq P - 1$ vor. Wir fordern also zusätzlich, dass $\delta.v$ innerhalb eines Intervalls liegt. Zudem nutzen wir Gleichung (6.48) i.V.m. Gleichung (6.33) und Gleichung (6.34) zur Gewinnung von

$$P_{\text{randy}} \Rightarrow 0 \leq P - 1 \leq b_{v n_v^6}^6 \cdot l - 1 \wedge \bigwedge_{l=1}^{n_v^3} 0 \leq P - 1 \leq a_{v l}^3 \cdot l - 1, \quad (6.60)$$

d. h. die Feldlängen betragen mindestens P . Konjunktive Verknüpfung von P_{randy} und $0 \leq \delta.v \leq P - 1$ liefert

$$P_{\text{randy}} \wedge 0 \leq \delta.v \leq P - 1 \Rightarrow 0 \leq \delta.v \leq b_{v n_v^6}^6 \cdot l - 1 \wedge \bigwedge_{l=1}^{n_v^3} 0 \leq \delta.v \leq a_{v l}^3 \cdot l - 1. \quad (6.61)$$

Die Aussage

$$0 \leq \delta.v \leq b_{v n_v^6}^6 \cdot l - 1 \wedge \bigwedge_{l=1}^{n_v^3} 0 \leq \delta.v \leq a_{v l}^3 \cdot l - 1 \quad (6.62)$$

in Gleichung (6.58) ist somit schon in $P_{\text{randy}} \wedge 0 \leq \delta.v \leq P - 1$ enthalten und kann entfernt werden. Zudem gilt

$$P_{\text{randy}} \Rightarrow \delta.l = -1 \wedge \delta.\text{typ} = \text{int}. \quad (6.63)$$

Es verbleibt die Richtigkeit von $0 \leq \delta.v \leq P - 1$ zu zeigen. Die Vorbedingung von S_{4i} impliziert die Ungleichungen $0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1$ und P_{MED} . Daher ist nach Satz 6.4 die Richtigkeit von $0 \leq \delta.v \leq P - 1$ schon gezeigt, wenn $0 \leq \mu_3.v \leq P_{x_3} - 1$ gilt. An dieser Stelle kann allerdings keine weitere Vereinfachung des Ausdrucks vorgenommen werden. $0 \leq \mu_3.v \leq P_{x_3} - 1$ wird in die Vorbedingung übernommen.

Die Zuweisungsanweisungen in Gleichung (6.58) fordern, dass $|\mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v| < \text{Max}_f$ gilt. Da die Matrizen \mathbf{R}^κ unitär beschränkt sind, ist dies dann der Fall, wenn $\bigwedge_{\nu=1}^{n_v^3} |a_{v \nu}^3(\delta.v).v| < \text{Max}_f$ gilt.

Um den in Gleichung (6.58) auftretenden Ausdruck

$$P_z(\mu_1.v, \mu_2.v) \left[\begin{array}{c} b_{v n_v^6}^6(\delta.v).v \\ \mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v \end{array} \right] \quad (6.64)$$

zu vereinfachen, schreiben wir $P_z(\mu_1.v, \mu_2.v)$ ausführlich auf

$$P_z(\mu_1.v, \mu_2.v) \equiv \bigwedge_{k=1}^{n_v^3} P_{b_k}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) \wedge \bigwedge_{k=1}^{n_v^6-1} P_{b_k}^6(\mu_1.v, \mu_2.v, P_{x_3} - 1, \delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)) \wedge P_{b_{n_v^6}}^6(\mu_1.v, \mu_2.v, P_{x_3} - 1, \delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)), \quad (6.65)$$

wobei

$$P_{b_{n_v^6}}^6(\mu_1.v, \mu_2.v, P_{x_3}-1, \delta([\mu_1.v, \mu_2.v, P_{x_3}-1]^T)) \equiv b_{v n_v^6}^6(\delta([\mu_1.v, \mu_2.v, P_{x_3}-1]^T)).v = b_{v n_v^6}^6([\mu_1.v, \mu_2.v, P_{x_3}-1]^T). \quad (6.66)$$

Offenbar erfährt nur der letzte Konjunktionsterm von $P_z(\mu_1.v, \mu_2.v)$ eine Änderung in Gleichung (6.64). Wir werden im Folgenden diesen untersuchen, d.h.

$$b_{v n_v^6}^6(\delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)).v = b_{v n_v^6}^6([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) \begin{bmatrix} b_{v n_v^6}^6(\delta.v).v \\ \mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v \end{bmatrix} \quad (6.67)$$

bestimmen. Dazu stärken wir allerdings noch die Vorbedingung von $S_{b_{v n_v^6}^6}$ durch $\delta.v = \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1)P_{x_1}P_{x_2}$. Aus dem Teilprädikat Gleichung (6.67) wird durch Nutzen der Stärkung und Ersetzen des δ durch $\delta.v$ auf der linken Seite

$$\mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta.v).v = b_{v k}^6([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \quad (6.68)$$

Nutzen wir nun die Stärkung in anderer Richtung, um in $\mathbf{a}_v^3(\delta.v).v$ das $\delta.v$ durch $\delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)$ zu ersetzen, so haben wir

$$\mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3(\delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)).v = b_{v k}^6([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \quad (6.69)$$

Wir werden nun die Istwerte der abhängigen Variablen durch die Sollwerte ersetzen. Dazu nutzen wir Gleichung (6.48) i.V.m. Gleichung (6.33) und Gleichung (6.35) und ermitteln

$$P_{\text{randy}} \Rightarrow \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{\text{avz}}(\mu_1, \mu_2) \Rightarrow \mathbf{a}_v^3(\delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)).v = \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \quad (6.70)$$

Aus Gleichung (6.69) wird dann

$$\mathbf{r}_{n_v^6, \bullet}^6 \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) = b_{v k}^6([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \quad (6.71)$$

Diese Beziehung ist offenbar in P_{WDF} enthalten und kann somit auch aus der Vorbedingung Gleichung (6.58) abgetrennt werden. Die gestärkte Vorbedingung aus Gleichung (6.58) ergibt sich also zu

$$l_i \wedge b_i \wedge \delta.v = \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1)P_{x_1}P_{x_2} \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1 \wedge \bigwedge_{\nu=1}^{n_v^3} |a_{v\nu}^3(\delta.v).v| \leq \text{Max}_f$$

$$\bigwedge_{k=1}^{n_v^3} P_{b_k}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) \wedge \bigwedge_{k=1}^{n_v^6-1} P_{b_k}^6(\mu_1.v, \mu_2.v, P_{x_3} - 1, \delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)). \quad (6.72)$$

Die aus den Anweisungen $S_{b_{v_1}^6}$ bis einschließlich $S_{b_{v_1}^6}$ resultierenden Vorbedingungen können ebenso umgeformt und abgetrennt werden. Als Vorbedingung von $S_{b_{v_1}^6}$ erhalten wir

$$\begin{aligned} & l_i \wedge b_i \wedge \delta.v = \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2} \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1 \wedge \\ & \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) \wedge \bigwedge_{\nu=1}^{n_v^3} |a_{v\nu}^3(\delta.v).v| < \text{Max}_f . \end{aligned} \quad (6.73)$$

Die Vorbedingung von **delta** = **mu1** + **mu2** * P_{x_1} + **mu3** * P_{x_1} * P_{x_2} ; lautet mit Gleichung (6.41)

$$\begin{aligned} & l_i \wedge b_i \wedge \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} = \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2} \wedge \\ & \bigwedge_{\nu=1}^{n_v^3} |a_{v\nu}^3(\mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2}).v| \leq \text{Max}_f \wedge \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) \wedge \\ & 0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1 . \end{aligned} \quad (6.74)$$

Die Vorbedingung von **mu3** = $P_{x_3} - 1$; ergibt sich leicht aus Gleichung (6.74) zu

$$\begin{aligned} & l_i \wedge b_i \wedge \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2} = \mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2} \wedge 0 \leq P_{x_3} - 1 \wedge \\ & \bigwedge_{\nu=1}^{n_v^3} |a_{v\nu}^3(\mu_1.v + \mu_2.v P_{x_1} + (P_{x_3} - 1) P_{x_1} P_{x_2}).v| \leq \text{Max}_f \wedge \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) \wedge \\ & 0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 . \end{aligned} \quad (6.75)$$

Sowohl die wahren Aussagen, als auch die von $l_i \wedge b_i$ implizierten Aussagen können abgetrennt werden und es verbleibt

$$l_i \wedge b_i \wedge \bigwedge_{k=1}^{n_v^3} P_{bk}^3(\mu_1.v, \mu_2.v, 0, \delta([\mu_1.v, \mu_2.v, 0]^T)) . \quad (6.76)$$

Es folgen die Anweisungen $S_{b_{v_1}^3}$ bis einschließlich **mu3**=0; deren Beweis ähnlich dem vorherigen ist. Für die Vorbedingung der Anweisung S_{4i} verbleibt, wie behauptet, $l_i \wedge b_i$.

Die Programmablaufpläne und Anweisungssequenzen für S_{randx} und S_{randy} entsprechen bis auf kleine Änderungen denen der Anweisungssequenz S_{randz} und werden hier nicht angegeben. Die Beweise werden ebenfalls analog zu denen der Anweisung S_{randz} geführt.

6.3 Normalberechnung der Verzögererwerte

Die Anweisung S_{norm} dient der Berechnung der neuen Verzögererwerte innerhalb des Berechnungsgebietes, d. h., es erfolgt die Berechnung von Gleichung (6.4). Die zu erfüllende Nachbedingung lautet

$$P_{\text{norm}} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P'_{\text{aus}} \wedge P'_{\text{av}} \wedge P'_{\text{norm}} . \quad (6.77)$$

Die Anweisung S_{norm} erfolgt durch Zerlegung in 7 sequentielle Anweisungen von denen jeweils die Anweisung $S_{\kappa\text{norm}}$ den Vektor \mathbf{b}_v^κ berechnet. Dazu die folgende Übersicht

Bedingung	Anweisung
$\{P_{av}\}$	S_{1norm}
$\{P_{norm}^1\}$	S_{2norm}
	\vdots
	S_{6norm}
$\{P_{norm}^6\} \equiv P_{WDF} \wedge P_{dec} \wedge P'_{aus} \wedge P'_{av} \wedge \bigwedge_{\kappa=1}^6 P_{norm}^{\kappa}$	S_{7norm}
$\{P_{norm}\} \equiv P_{WDF} \wedge P_{dec} \wedge P'_{aus} \wedge P'_{av} \wedge P'_{norm}$	

Im Bild 6.7 ist der Programmablaufplan für die Berechnung des Vektors \mathbf{b}_v^1 dargestellt. Die Grenzen der Schleifen der Anweisungen S_{2norm} bis S_{7norm} ergeben sich aus Tabelle 6.1. Im Bild 6.8 sind die für die Beweisführung nötigen Schritte für die Anweisung S_{1norm} dargestellt.² Die Schleifenbedingungen sind $b_a \equiv \mu_1.v < P_{x_1} - 1 \equiv \mu_1.v + 1 \leq P_{x_1} - 1$, $b_m \equiv \mu_2.v < P_{x_2} \equiv \mu_2.v + 1 \leq P_{x_2}$ und $b_i \equiv \mu_3.v < P_{x_3} \equiv \mu_3.v + 1 \leq P_{x_3}$. Aus dem Programmablaufplan ermitteln wir den C-Code von S_{1norm} zu

Bedingung	Anweisung	Kommentar
$\{P_{av}\}$	<pre> for(mu1 = 0 ; mu1 < P_{x1} - 1 ; mu1=mu1+1){ for(mu2 = 0 ; mu2 < P_{x2} ; mu2=mu2+1){ for(mu3 = 0 ; mu3 < P_{x3} ; mu3=mu3+1){ delta = mu1 + mu2 * P_{x1} + mu3 * P_{x1} * P_{x2} ; S_{bv11} : S_{bv1n_v¹} } } } </pre>	<p>Wellen b_{v1}^1 berechnen</p> <p>Wellen $b_{vn_v^1}^1$ berechnen</p>
$\{P_{norm}^1\}$		

mit den Anweisungen

$$\begin{aligned}
S_{bv1} &\equiv \mathbf{b}_{v_1_1}[\text{delta}+1] = p_{v^1e11}^1 * \mathbf{b}_{e_1}[\text{delta}] + \dots + p_{v^1e1n_e}^1 * \mathbf{b}_{e_n_e}[\text{delta}] ; \\
&\vdots \\
S_{bv n_v^1} &\equiv \mathbf{b}_{v_1_n_v^1}[\text{delta}+1] = p_{v^1e n_v^1 1}^1 * \mathbf{b}_{e_1}[\text{delta}] + \dots + p_{v^1e n_v^1 n_e}^1 * \mathbf{b}_{e_n_e}[\text{delta}] ;
\end{aligned}$$

Der Beweis wird nur für die erste der 7 Teilanweisungen geführt, da die Teilanweisungen einander ähnlich sind.

Satz 5 *Der Ausdruck*

$$I_a \equiv 0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge P_{av} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1+1, \mu_2, \mu_3, \delta(\boldsymbol{\mu}) + \delta_1) \quad (6.78)$$

²Eigentlich müssten für die hier auftretenden Invarianten I_a , I_i , und Bedingungen b_a , b_i neue Bezeichnungen eingeführt werden, da sie schon im Kapitel 6.2 verwendet wurden. Um die Übersicht zu wahren, verzichten wir jedoch darauf und wollen die Gültigkeit dieser Bezeichner auf die Unterkapitel beschränken.

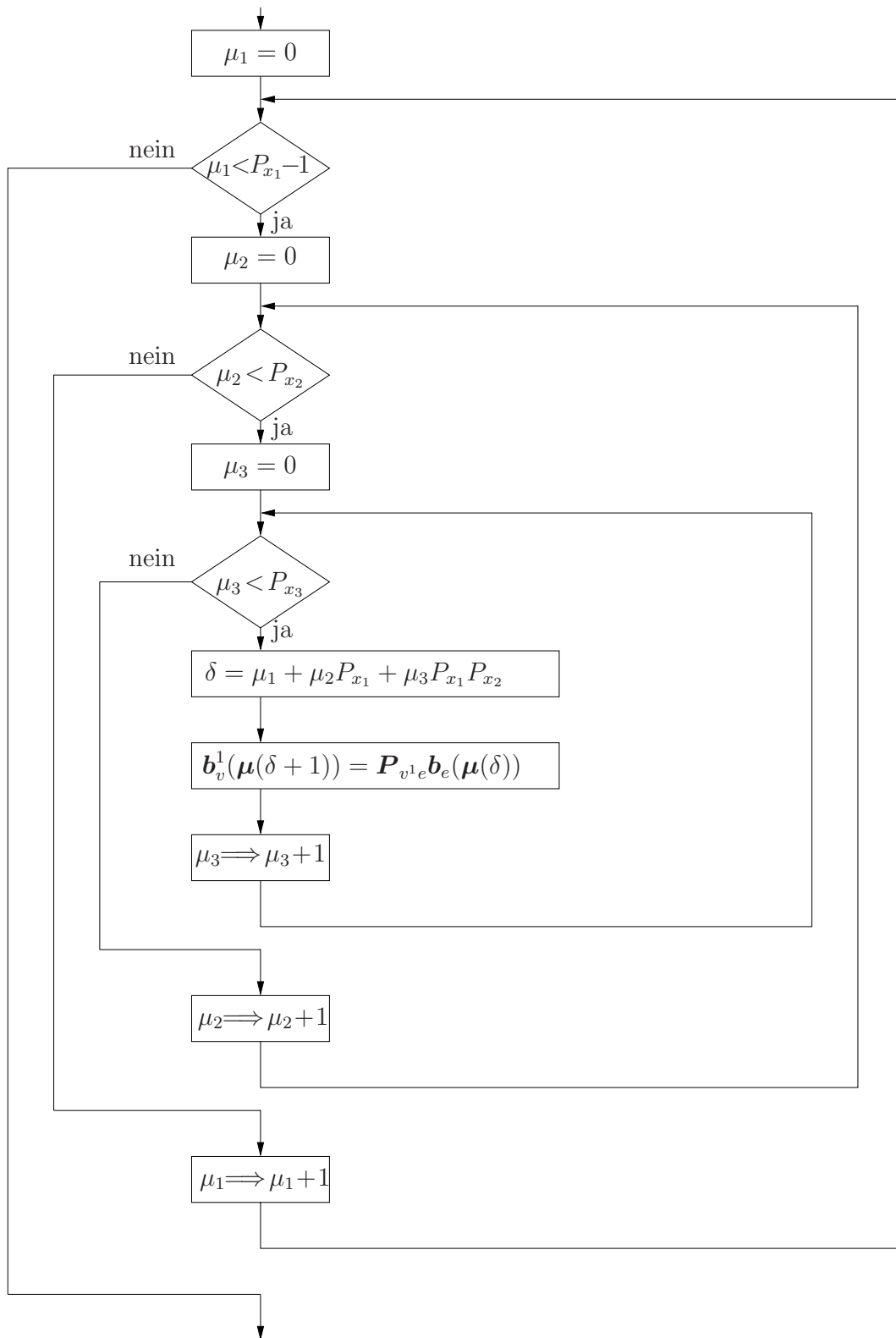


Bild 6.7: Berechnung der neuen Verzögererwerte innerhalb des Berechnungsgebietes (Gleichung (6.4))

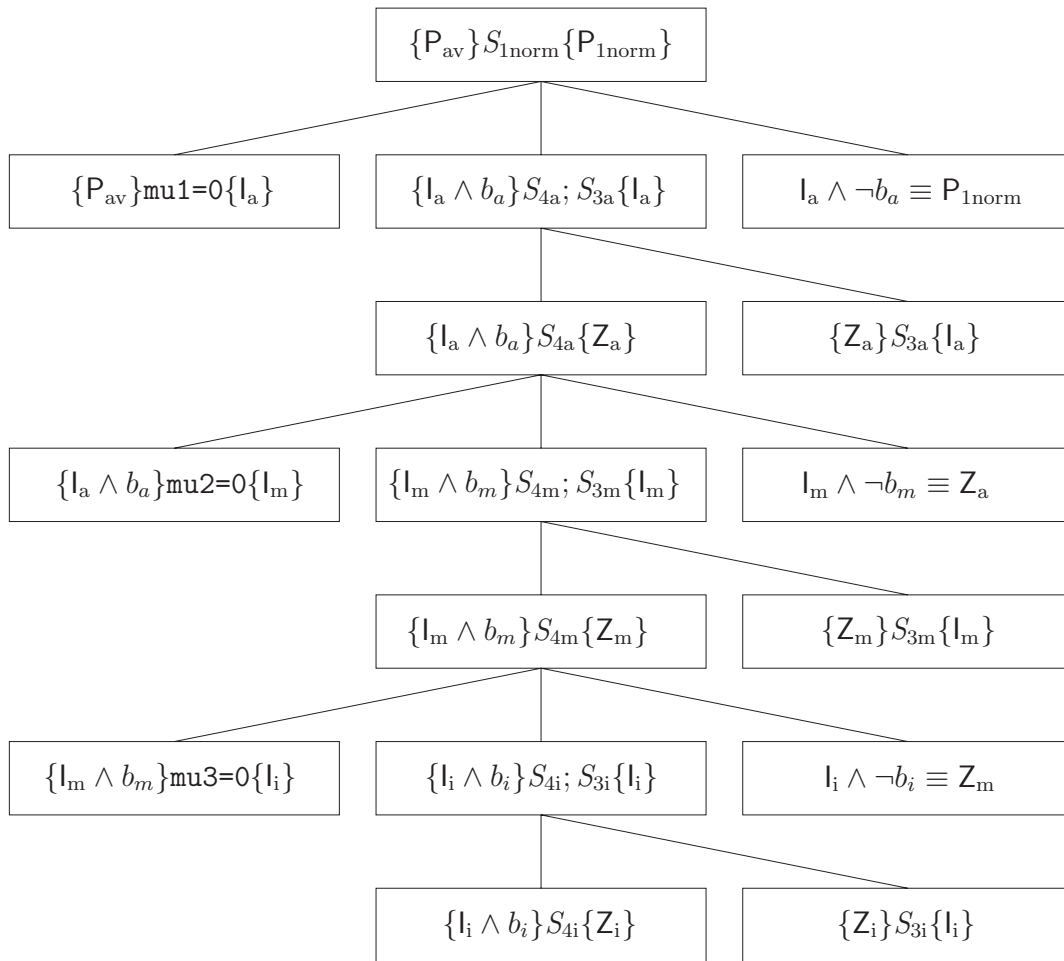


Bild 6.8: Berechnung der Verzögererwerte \mathbf{b}_v^1 im Normalfall

mit

$$P_{av} \equiv P'_{aus} \wedge P_{WDF} \wedge P'_{av} \wedge P'_{be} \wedge P_q \wedge P_{dec} . \quad (6.79)$$

ist eine gültige Schleifeninvariante der äußeren Schleife.

Beweis 5

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung liefert die Vorbedingung

$$P_{av} \equiv 0 \leq P_{x_1} - 1 \wedge P_{av} \wedge \bigwedge_{\mu_1=0}^{-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1+1, \mu_2, \mu_3, \delta(\boldsymbol{\mu})+\delta_1) \wedge \quad (6.80)$$

$$0 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.type = \text{int} \equiv P_{av} ,$$

da P_{av} die Bedingung $0 \leq P_{x_1} - 1$ impliziert. Die Nachbedingung bestimmt sich (unter Beachtung von $P_{x_1} - 1 \leq \mu_1.v \wedge 0 \leq \mu_1.v \leq P_{x_1} - 1 \equiv 0 \leq \mu_1.v = P_{x_1} - 1$) zu

$$P_{\text{norm}}^1 \equiv I_a \wedge P_{x_1} - 1 \leq \mu_1.v \equiv 0 \leq \mu_1.v = P_{x_1} - 1 \wedge P_{av} \wedge \bigwedge_{\mu_1=0}^{P_{x_1}-2} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1+1, \mu_2, \mu_3, \delta+\delta_1) \quad (6.81)$$

und kann zu $P_{av} \wedge P_{\text{norm}}^1$ geschwächt werden. Die Zwischenbedingung bestimmt sich zu

$$Z_a \equiv 0 \leq \mu_1.v + 1 \leq P_{x_1} - 1 \wedge P_{av} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v} \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1+1, \mu_2, \mu_3, \delta(\boldsymbol{\mu})+\delta_1) \wedge \quad (6.82)$$

$$\mu_1.v + 1 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.type = \text{int} .$$

Stärkung um $0 \leq \mu_1.v$ liefert

$$Z'_a \equiv Z_a \wedge 0 \leq \mu_1.v \equiv I_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2, \mu_3]^T)+\delta_1) . \quad (6.83)$$

Der Beweis von $\{I_a \wedge b_a\} S_{4a} \{Z'_a\}$ wird gemäß Axiom (5.35) durchgeführt, da S_{4a} die mittlere Schleife darstellt.

Satz 6 Der Ausdruck

$$I_m \equiv 0 \leq \mu_2.v \leq P_{x_2} \wedge I_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2, \mu_3]^T)+\delta_1) \quad (6.84)$$

ist eine gültige Schleifeninvariante der mittleren Schleife.

Beweis 6

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\mu_2=0$ liefert

$$\begin{aligned} 0 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2, \mu_3]^T) + \delta_1) \wedge \\ 0 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.typ = \text{int} \equiv l_a \wedge b_a, \end{aligned} \quad (6.85)$$

da l_a die Aussage $0 \leq P_{x_2}$ impliziert, sodass wir, wie zu zeigen ist, die Vorbedingung von S_{4a} erhalten. Die Nachbedingung bestimmt sich zu

$$l_m \wedge P_{x_2} \leq \mu_2.v \equiv 0 \leq \mu_2.v = P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2, \mu_3]^T) + \delta_1) \quad (6.86)$$

und stellt nach Abtrennung von $0 \leq \mu_2.v = P_{x_2}$, wie behauptet, Z'_a aus Gleichung (6.83) dar. Die Zwischenbedingung bestimmt sich zu

$$\begin{aligned} Z_m \equiv 0 \leq \mu_2.v + 1 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2, \mu_3]^T) + \delta_1) \\ \wedge \mu_2.v + 1 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.typ = \text{int}. \end{aligned} \quad (6.87)$$

Stärkung um $0 \leq \mu_2.v$ und Vereinfachung liefert

$$Z'_m \equiv Z_m \wedge 0 \leq \mu_2.v \equiv l_m \wedge b_m \wedge \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1). \quad (6.88)$$

Der Beweis von $\{l_m \wedge b_m\} S_{4m} \{Z'_m\}$ wird gemäß Axiom (5.35) durchgeführt, da S_{4m} die innere Schleife darstellt.

Satz 7 *Der Ausdruck*

$$l_i \equiv 0 \leq \mu_3.v \leq P_{x_3} \wedge l_m \wedge b_m \wedge \bigwedge_{\mu_3=0}^{\mu_3.v-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1). \quad (6.89)$$

ist eine gültige Schleifeninvariante der inneren Schleife.

Beweis 7

Zunächst bemerken wir $l_m \wedge b_m \Rightarrow 0 \leq \mu_2.v \leq P_{x_2} - 1$. Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\mu_3=0$ liefert

$$0 \leq P_{x_3} \wedge l_m \wedge b_m \wedge \bigwedge_{\mu_3=0}^{-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1) \wedge \quad (6.90)$$

$$0 \leq \text{Max}_i \wedge \mu_3.l = -1 \wedge \mu_3.typ = \text{int} \equiv l_m \wedge b_m.$$

Die Nachbedingung der inneren Schleife lautet

$$l_i \wedge P_{x_3} \leq \mu_3.v \equiv 0 \leq \mu_3.v = P_{x_3} \wedge l_m \wedge b_m \wedge \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1). \quad (6.91)$$

Einsetzen von l_m aus Gleichung (6.84) und b_m liefert

$$0 \leq \mu_3.v = P_{x_3} \wedge 0 \leq \mu_2.v+1 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v} \bigwedge_{\mu_3=0}^{P_{x_3}-1} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1) \quad (6.92)$$

und ist wegen $P_{x_3} > 0$ die Bedingung Z_m aus Gleichung (6.87). Die Zwischenbedingung der inneren Schleife lautet

$$Z_i \equiv 0 \leq \mu_3.v+1 \leq P_{x_3} \wedge l_m \wedge b_m \wedge \bigwedge_{\mu_3=0}^{\mu_3.v} \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3, \delta([\mu_1.v, \mu_2.v, \mu_3]^T) + \delta_1) \wedge \mu_3.v+1 \leq \text{Max}_i \wedge \mu_3.l = -1 \wedge \mu_3.typ = \text{int}. \quad (6.93)$$

Stärkung um $0 \leq \mu_3.v$ und Vereinfachung liefert

$$Z'_i \equiv Z_i \wedge 0 \leq \mu_3.v \equiv l_i \wedge b_i \wedge \bigwedge_{k=1}^{n_v^1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta([\mu_1.v, \mu_2.v, \mu_3.v]^T) + \delta_1). \quad (6.94)$$

Es verbleibt die Richtigkeit von $\{l_i \wedge b_i\} S_{4i} \{Z'_i\}$ zu zeigen. Zu diesem Zweck bereiten wir Z'_i weiter auf

$$Z'_i \equiv l_i \wedge b_i \wedge \bigwedge_{k=1}^{n_v^1-1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta(\mu.v) + \delta_1) \wedge P_{bn_v^1}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta(\mu.v) + \delta_1). \quad (6.95)$$

Die Vorbedingung von $S_{b_{vn_v^1}}$ erhalten wir durch Anwenden von Gleichung (5.43) und Stärkung zu

$$\begin{aligned} & l_i \wedge b_i \wedge P_{bn_v^1}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta(\mu.v) + \delta_1) \left[\mathbf{p}_{e_{n_v^1}, \bullet}^{b_{vn_v^1}^1(\delta.v+\delta_1).v} \mathbf{b}_{e(\delta.v).v} \right] \wedge \\ & \bigwedge_{k=1}^{n_v^1-1} P_{bk}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta(\mu.v) + \delta_1) \wedge \bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| \leq \text{Max}_f \wedge \\ & 0 \leq \delta.v + \delta_1 \leq b_{vn_v^1}^1.l - 1 \wedge \bigwedge_{l=1}^{n_e} 0 \leq \delta.v \leq b_{el}.l - 1 \wedge \bigwedge_{k=1}^{n_e} b_{vn_v^1}^1.typ = b_{ek}.typ = \text{float} \wedge \\ & 0 \leq \delta.v + \delta_1 \leq \text{Max}_i \wedge 0 \leq \delta.v \leq \text{Max}_i \wedge \delta.l = -1 \wedge \delta.typ = \text{int}. \end{aligned} \quad (6.96)$$

Hierbei wurde berücksichtigt, dass $n_e - 1$ Elemente des Zeilenvektors $\mathbf{p}_{ve\mu, \bullet}^\kappa$ den Wert null haben und das verbleibende Element den Wert eins hat. Weitere Vereinfachungen können wegen $l_i \Rightarrow \delta.typ =$

$\text{int} \wedge \delta.l = -1$ vorgenommen werden, indem die in l_i enthaltenen Prädikate abgetrennt werden. Gleiches Vorgehen für die Anweisungen $S_{b_{v_{n_v^1-1}}}, \dots, S_{b_{v_2}}, S_{b_{v_1}}$ führt zur Vorbedingung von $S_{b_{v_1}}$

$$\begin{aligned} l_i \wedge b_i \wedge \bigwedge_{k=1}^{n_v^1} P_{b_k}^1(\mu_1.v+1, \mu_2.v, \mu_3.v, \delta(\mu.v)+\delta_1) \left[\begin{matrix} b_{v_k}^1(\delta.v+\delta_1).v \\ \mathbf{p}_{e_k, \bullet}^1 \mathbf{b}_e(\delta.v).v \end{matrix} \right] \wedge \bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| \leq \text{Max}_f \wedge \\ \bigwedge_{l=1}^{n_v^1} \bigwedge_{k=1}^{n_e} b_{vl}^1.\text{typ} = b_{ek}.\text{typ} = \text{float} \wedge \bigwedge_{k=1}^{n_v^1} 0 \leq \delta.v + \delta_1 \leq b_{vk}^1.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \wedge \\ 0 \leq \delta.v + \delta_1 \leq \text{Max}_i \wedge 0 \leq \delta.v \leq \text{Max}_i. \end{aligned} \quad (6.97)$$

Diese Vorbedingung stärken wir durch $\delta.v = \delta(\mu.v)$ und erhalten durch Auswertung und Nutzung von Gleichung (6.25)

$$\begin{aligned} l_i \wedge b_i \wedge \delta.v = \delta(\mu.v) \wedge \bigwedge_{k=1}^{n_v^1} \mathbf{p}_{e_k, \bullet}^1 \mathbf{b}_e(\delta.v).v = b_{vk}^1(\mu.v + \mathbf{h}_1) \wedge \bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| \leq \text{Max}_f \wedge \\ \bigwedge_{l=1}^{n_v^1} \bigwedge_{k=1}^{n_e} b_{vl}^1.\text{typ} = b_{ek}.\text{typ} = \text{float} \wedge \bigwedge_{k=1}^{n_v^1} 0 \leq \delta.v + \delta_1 \leq b_{vk}^1.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \wedge \\ 0 \leq \delta.v + \delta_1 \leq \text{Max}_i \wedge 0 \leq \delta.v \leq \text{Max}_i. \end{aligned} \quad (6.98)$$

Wir werden zunächst zeigen, dass l_i die in Gleichung (6.98) auftretenden Aussagen bzgl. der Datentypen impliziert. Die Verwendung von Gleichung (6.89) i.V.m. Gleichung (6.33), Gleichung (6.34), Gleichung (6.78), Gleichung (6.79) und Gleichung (6.84), liefert

$$l_i \Rightarrow \bigwedge_{k=1}^{n_v^1} b_{vk}^1.\text{typ} = \text{float} \wedge \bigwedge_{k=1}^{n_e} b_{ek}.\text{typ} = \text{float}. \quad (6.99)$$

Alle Datentypen sind somit gleich. Insbesondere folgt daraus -wie behauptet- die Richtigkeit von

$$l_i \Rightarrow \bigwedge_{l=1}^{n_v^1} \bigwedge_{k=1}^{n_e} b_{vl}^1.\text{typ} = b_{ek}.\text{typ}, \quad (6.100)$$

sodass die Forderungen an die Datentypen aus Gleichung (6.97) abgetrennt werden können. Wie im Folgenden gezeigt wird, impliziert l_i auch die in Gleichung (6.98) auftretenden Ausdrücke bzgl. der Feldlängen. Um dies zu zeigen verwenden wir erneut Gleichung (6.89) i.V.m. Gleichung (6.33), Gleichung (6.34), Gleichung (6.78), Gleichung (6.79) und Gleichung (6.84). Wir erhalten

$$l_i \wedge b_i \Rightarrow 0 \leq \mu_1.v \leq P_{x_1} - 2 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1. \quad (6.101)$$

Wegen Gleichung (6.40) folgt für $\delta.v$ die Beziehung

$$l_i \wedge b_i \Rightarrow 0 \leq \delta.v \leq P - 2. \quad (6.102)$$

Durch mehrfache Anwendung des Kettenschlusses und Gleichung (6.34) folgen die Implikationen

$$l_i \Rightarrow \bigwedge_{k=1}^{n_v^1} P \leq b_{vk}^1.l \leq \text{Max}_i, \quad l_i \Rightarrow \bigwedge_{k=1}^{n_e} P \leq b_{ek}.l \leq \text{Max}_i. \quad (6.103)$$

Konjunktive Verknüpfung mit 6.102 liefert

$$l_i \wedge b_i \Rightarrow \bigwedge_{k=1}^{n_v^1} 0 \leq \delta.v \leq b_{vk}^1.l - 2, \quad l_i \wedge b_i \Rightarrow \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 2. \quad (6.104)$$

Anschaulich bedeutet diese Aussage, dass der Wert der Variablen δ innerhalb eines Intervalls liegt. Nun betrachten wir die folgenden -offenbar richtigen- Aussagen

$$\left[\bigwedge_{k=1}^{n_v^1} 1 \leq \delta.v + 1 \leq b_{vk}^1.l - 1 \right] \Rightarrow \left[\bigwedge_{k=1}^{n_v^1} 0 \leq \delta.v + 1 \leq b_{vk}^1.l - 1 \right], \quad (6.105)$$

$$\left[\bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 2 \right] \Rightarrow \left[\bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \right].$$

Offenbar sind auch die Aussagen bzgl. der Feldlänge in Gleichung (6.98) richtig und werden abgetrennt. Aufgrund von $l_i \Rightarrow b_{ek}.l \leq \text{Max}_i$ folgt aus l_i auch $l_i \Rightarrow 0 \leq \delta.v + \delta_1 \leq \text{Max}_i \wedge 0 \leq \delta.v \leq \text{Max}_i$, was ebenfalls in Gleichung (6.98) abgetrennt werden kann.

Da l_i die Beziehung $\bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| \leq \text{Max}_f$ impliziert, verbleibt von Gleichung (6.98) noch

$$l_i \wedge b_i \wedge \delta.v = \delta(\mu.v) \wedge \bigwedge_{k=1}^{n_v^1} p_{ek,\bullet}^1 b_e(\delta.v).v = b_{vk}^1(\mu.v + h_1). \quad (6.106)$$

Hierin ersetzen wir $\delta.v$ durch $\delta(\mu.v)$. Zudem werden wir im Folgenden die Istwerte durch die Sollwerte ersetzen. Dazu ermitteln wir zunächst aus Gleichung (6.89), Gleichung (6.84), Gleichung (6.78), P_{av} , P_e und P'_{be}

$$l_i \Rightarrow \bigwedge_{\delta=0}^{P-1} b_e(\delta).v = b_e(\mu(\delta)). \quad (6.107)$$

Da diese Gleichung für alle δ des Berechnungsgebietes gilt, gilt sie auch für $\delta.v = \delta(\mu.v)$, d. h.

$$l_i \wedge \delta.v = \delta(\mu.v) \Rightarrow b_e(\delta.v).v = b_e(\mu(\delta.v)). \quad (6.108)$$

Somit lautet die um $\delta.v = \delta(\mu.v)$ gestärkte Vorbedingung von S_{bv1}

$$l_i \wedge b_i \wedge \delta.v = \delta(\mu.v) \wedge \bigwedge_{k=1}^{n_v^1} p_{ek,\bullet}^1 b_e(\mu(\delta.v)) = b_{vk}^1(\mu.v + h_1). \quad (6.109)$$

Mit $\delta.v = \delta(\mu.v) \wedge P_{MDED} \equiv \mu(\delta.v) = \mu.v \wedge P_{MDED}$ haben wir

$$l_i \wedge b_i \wedge \delta.v = \delta(\mu.v) \wedge \bigwedge_{k=1}^{n_v^1} p_{ek,\bullet}^1 b_e(\mu.v) = b_{vk}^1(\mu.v + h_1). \quad (6.110)$$

Mit $l_i \Rightarrow P_{WDF}$ folgt die Wahrheit der letzten Aussage. Durch Abtrennung der wahren Aussage erhalten wir die Vorbedingung von S_{bv1} letztendlich zu $l_i \wedge b_i \wedge \delta.v = \delta(\mu.v)$. Die Vorbedingung von δ lautet mit Gleichung (6.41)

$$l_i \wedge b_i \wedge \mu_1.v + \mu_2.v P_{x_1} + \mu_3.v P_{x_1} P_{x_2} = \delta(\mu.v) \wedge \quad (6.111)$$

$$0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1.$$

Wegen $l_i \Rightarrow P_{\text{MDED}} \wedge P_{\text{decuv}}$ und $l_i \wedge b_i \implies 0 \leq \mu_1.v \leq P_{x_1} - 1 \wedge 0 \leq \mu_2.v \leq P_{x_2} - 1 \wedge 0 \leq \mu_3.v \leq P_{x_3} - 1$ folgt schließlich, wie behauptet,

$$l_i \wedge b_i. \quad (6.112)$$

Die anderen 6 Anweisungen wollen wir nicht beweisen, da sie ähnlich sind. Die Unterschiede in den Nachbedingungen liegen darin, dass jeweils P_{norm}^κ hinzukommt.

6.4 Berechnung der Eingangswellen der Verzögerer

Die Nachbedingung von S_{av} lautet

$$P_{\text{av}} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_{\text{q}} \wedge P'_{\text{be}} \wedge P'_{\text{aus}} \wedge P'_{\text{av}} \quad (6.113)$$

und beinhaltet die korrekt berechneten Werte der Wellengrößen \mathbf{a}_v am Rand. Diese Berechnung leistet die Anweisung S_{av} . Wie zuvor erfolgt zunächst eine Zerlegung S_{av} in sequentielle Teilanweisungen nach $S_{\text{av}} \equiv S_{\text{avx}}; S_{\text{avy}}; S_{\text{avz}}; .$ Dazu ist die Übersicht

Bedingung	Anweisung
$\{P_{\text{aus}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_{\text{q}} \wedge P'_{\text{be}} \wedge P'_{\text{aus}}$	S_{avx}
$\{P_{\text{avx}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_{\text{q}} \wedge P'_{\text{be}} \wedge P'_{\text{aus}} \wedge P'_{\text{avx}}$	S_{avy}
$\{P_{\text{avy}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_{\text{q}} \wedge P'_{\text{be}} \wedge P'_{\text{aus}} \wedge P'_{\text{avx}} \wedge P'_{\text{avy}}$	S_{avz}
$\{P_{\text{av}}\} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_{\text{q}} \wedge P'_{\text{be}} \wedge P'_{\text{aus}} \wedge P'_{\text{avx}} \wedge P'_{\text{avy}} \wedge P'_{\text{avz}}$	

hilfreich. Die Anweisungen S_{avx} , S_{avy} und S_{avz} sind jeweils zweifach verschachtelte for-Schleifen in den Variablen μ_1, μ_2 und μ_3 , von denen jeweils eine konstant ist. Die verbliebenen zwei unabhängigen Variablen bilden dann eine Grenzfläche. Innerhalb der Schleifen wird Gleichung (6.3) ausgewertet. Der zugehörige Programmablaufplan für S_{avz} ist im Bild 6.9 dargestellt. Der C-Code für S_{avz} lautet

Bedingung	Anweisung	Kommentar
$\{P_{\text{avy}}\}$	<pre> for(mu1 = 0 ; mu1 < P_{x1} ; mu1=mu1+1){ for(mu2 = 0 ; mu2 < P_{x2} ; mu2=mu2+1){ mu3 = 0; delta = mu1 + mu2 * P_{x1} + mu3 * P_{x1} * P_{x2} ; S_{av6} mu3 = P_{x3} - 1; delta = mu1 + mu2 * P_{x1} + mu3 * P_{x1} * P_{x2} ; S_{av3} } } </pre>	Wellen \mathbf{a}_{v6} berechnen
$\{P_{\text{avz}}\}$		Wellen \mathbf{a}_{v3} berechnen

wobei die Berechnungsanweisungen durch

$$\begin{aligned}
S_{\text{av1}} &\equiv \mathbf{a}_{v1}[\text{delta}] = p_{v1e11} * \mathbf{b}_{e1}[\text{delta}] + \dots + p_{v1en_e} * \mathbf{b}_{en_e}[\text{delta}]; \\
&\vdots \\
\mathbf{a}_{v1n_v^1}[\text{delta}] &= p_{v1en_v^1} * \mathbf{b}_{e1}[\text{delta}] + \dots + p_{v1en_v^1 n_e} * \mathbf{b}_{en_e}[\text{delta}];
\end{aligned}$$

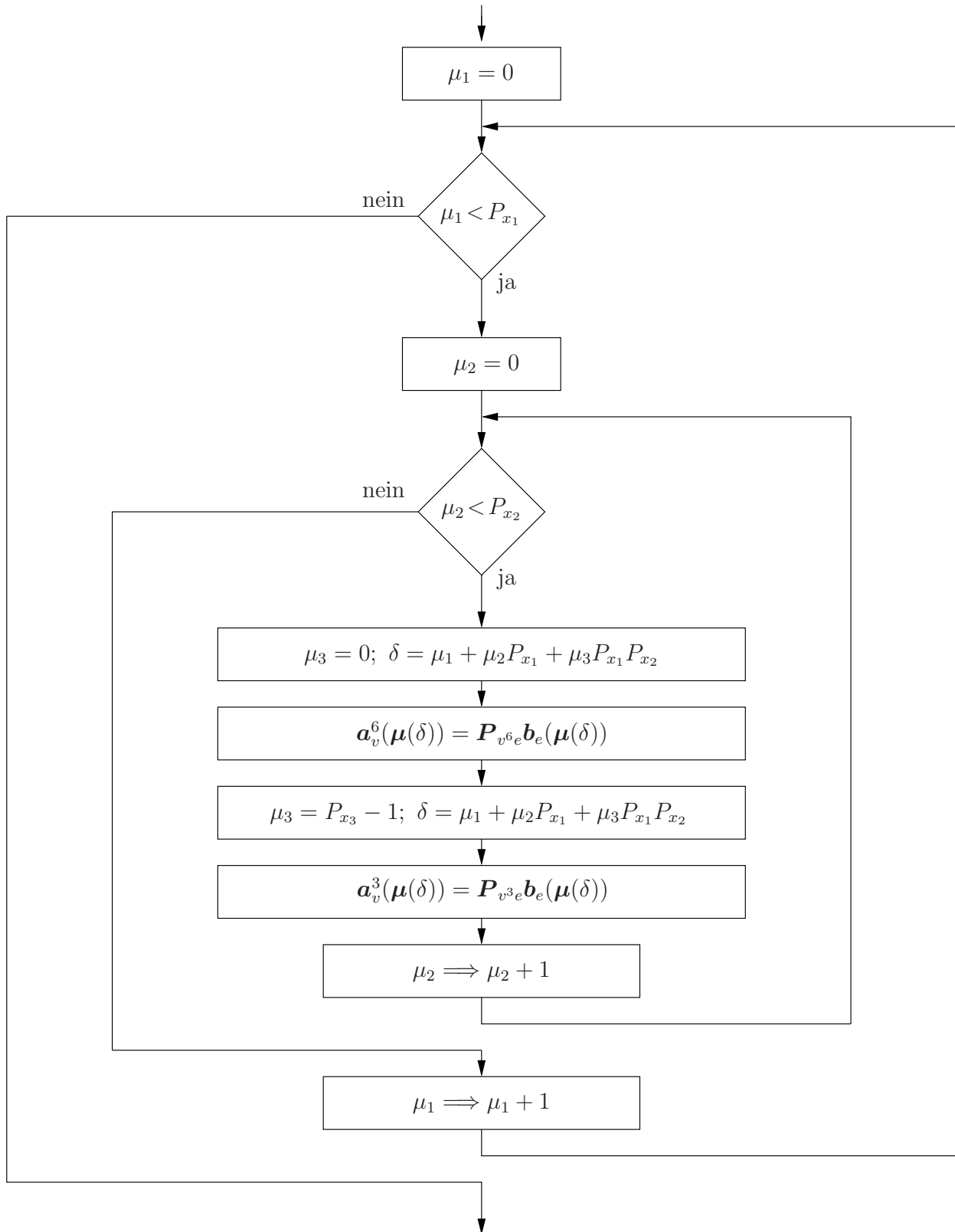


Bild 6.9: Berechnung der Wellen, die in der nächsten Abtastschicht durch die Grenzflächen $\mu_3 = 0$ und $\mu_3 = P_{x_3} - 1$ das Gebiet \mathcal{G} verlassen (Gleichung (6.3) für $\kappa = 3, 6$)

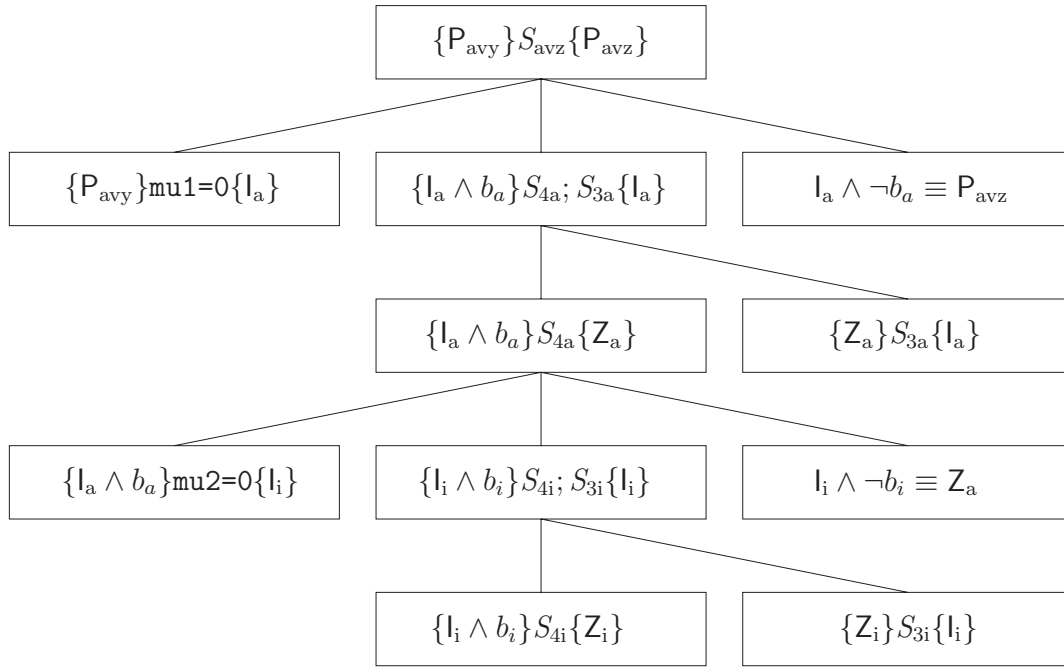


Bild 6.10: Berechnung der Eingangswellen der Verzögerer

bis

$$\begin{aligned}
 S_{av6} &\equiv a_{v_6_1}[\delta] = p_{v^6e_{11}} * b_{e_1}[\delta] + \dots + p_{v^6e_{1n_e}} * b_{e_n_e}[\delta]; \\
 &\vdots \\
 a_{v_6_n_v^6}[\delta] &= p_{v^6e_{n_v^6 1}} * b_{e_1}[\delta] + \dots + p_{v^6e_{n_v^6 n_e}} * b_{e_n_e}[\delta];
 \end{aligned}$$

gegeben sind. Zur Ermittlung der Vorbedingung von S_{avz} ist der Überblick im Bild 6.10 nützlich. Hierin ist $b_a \equiv \mu_1.v < P_{x_1} \equiv \mu_1.v + 1 \leq P_{x_1}$ und $b_i \equiv \mu_2.v < P_{x_2} \equiv \mu_2.v + 1 \leq P_{x_2}$. Wir werden nun die Beweisführung für S_{avz} vornehmen und beginnen mit der äußeren Schleife.

Satz 8 *Der Ausdruck*

$$I_a \equiv 0 \leq \mu_1.v \leq P_{x_1} \wedge P_{avy} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1, \mu_2) \quad (6.114)$$

ist eine gültige Schleifeninvariante der äußeren Schleife.

Beweis 8

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\mu_1=0$ liefert

$$0 \leq P_{x_1} \wedge P_{avy} \wedge \bigwedge_{\mu_1=0}^{-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1, \mu_2) \wedge 0 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.typ = \text{int} \equiv P_{avy}. \quad (6.115)$$

Die Nachbedingung bestimmt sich zu

$$I_a \wedge \mu_1.v \geq P_{x_1} \equiv \mu_1.v = P_{x_1} \wedge P_{avy} \wedge \bigwedge_{\mu_1=0}^{P_{x_1}-1} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1, \mu_2) \quad (6.116)$$

und kann um $\mu_1.v = P_{x_1}$ zu P_{avz} geschwächt werden. Die Zwischenbedingung bestimmt sich zu

$$Z_a \equiv 0 \leq \mu_1.v + 1 \leq P_{x_1} \wedge P_{avy} \wedge \bigwedge_{\mu_1=0}^{\mu_1.v} \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1, \mu_2) \wedge \quad (6.117)$$

$$\mu_1.v + 1 \leq \text{Max}_i \wedge \mu_1.l = -1 \wedge \mu_1.typ = \text{int}.$$

Stärkung um $0 \leq \mu_1.v$ liefert

$$Z'_a \equiv Z_a \wedge 0 \leq \mu_1.v \equiv l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1.v, \mu_2). \quad (6.118)$$

Der Beweis von $\{l_a \wedge b_a\}S_{4a}\{Z_a\}$ wird wiederum gemäß Axiom (5.35) durchgeführt, da S_{4a} die innere Schleife darstellt.

Satz 9 *Der Ausdruck*

$$l_i \equiv 0 \leq \mu_2.v \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v-1} P_{az}(\mu_1.v, \mu_2) \quad (6.119)$$

ist eine gültige Schleifeninvariante der inneren Schleife.

Beweis 9

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung $\text{mu2}=0$ liefert

$$0 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{-1} P_{az}(\mu_1.v, \mu_2) \wedge \mu_2.v + 1 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.typ = \text{int} \equiv l_a \wedge b_a. \quad (6.120)$$

Die Nachbedingung der inneren Schleife lautet

$$l_i \wedge \mu_2.v \geq P_{x_2} \equiv \mu_2.v = P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{P_{x_2}-1} P_{az}(\mu_1.v, \mu_2) \quad (6.121)$$

und ist nach Schwächung um $\mu_2.v = P_{x_2}$, wie behauptet, Z'_a aus Gleichung (6.118). Die Zwischenbedingung der inneren Schleife lautet

$$Z_i \equiv 0 \leq \mu_2.v + 1 \leq P_{x_2} \wedge l_a \wedge b_a \wedge \bigwedge_{\mu_2=0}^{\mu_2.v} P_{az}(\mu_1.v, \mu_2) \wedge \mu_2.v + 1 \leq \text{Max}_i \wedge \mu_2.l = -1 \wedge \mu_2.typ = \text{int}. \quad (6.122)$$

Stärkung um $0 \leq \mu_2.v$ liefert

$$Z'_i \equiv Z_i \wedge 0 \leq \mu_2.v \equiv l_i \wedge b_i \wedge P_{az}(\mu_1.v, \mu_2.v). \quad (6.123)$$

Es verbleibt die Richtigkeit von $\{l_i \wedge b_i\}S_{4i}\{Z'_i\}$ zu zeigen. Dazu zerlegen wir Z'_i wie folgt

$$\begin{aligned} Z'_i &\equiv l_i \wedge b_i \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) \wedge \\ &\quad \mathbf{a}_v^3(\delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T)).v = \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \end{aligned} \quad (6.124)$$

Zudem wird die Nachbedingung von S_{av3} durch

$$\delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \equiv \delta.v = \delta([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) \wedge \mu_3.v = P_{x_3} - 1 \quad (6.125)$$

gestärkt. Für die gestärkte Bedingung Z'_i notieren wir

$$\begin{aligned} Z'_i \wedge \delta.v &= \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \equiv l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \\ \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v &= \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) \wedge \mathbf{a}_v^3(\delta.v).v = \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T). \end{aligned} \quad (6.126)$$

Nach diesen Vorbereitungen gewinnen wir die Vorbedingung von S_{av3} zu

$$\begin{aligned} Z'_i \left[\begin{array}{c} \mathbf{a}_v^3(\delta.v).v \\ \mathbf{P}_{v_3e} \mathbf{b}_e(\delta.v).v \end{array} \right] \wedge \delta.v &= \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| < \text{Max}_f \wedge \\ \bigwedge_{k=1}^{n_v^3} 0 \leq \delta.v &\leq a_{vk}^3.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \wedge \bigwedge_{l=1}^{n_v^3} \bigwedge_{k=1}^{n_e} a_{vl}^3.typ = b_{ek}.typ \wedge \\ \delta.v &\leq \text{Max}_i \wedge \delta.l = -1 \wedge \delta.typ = \text{int} . \end{aligned} \quad (6.127)$$

Zur Verdichtung des obigen Prädikats werden wir zeigen, dass $l_i \wedge b_i$ die Prädikate bzgl. der Datentypen und der Feldlänge impliziert.

Durch Gleichung (6.119) , Gleichung (6.114) , Gleichung (6.113), Gleichung (6.34) und Gleichung (6.33) erhalten wir für die Datentypen

$$l_i \wedge b_i \Rightarrow \bigwedge_{k=1}^{n_v^3} a_{vk}^3.typ = \text{float} \wedge \bigwedge_{k=1}^{n_e} b_{ek}.typ = \text{float} . \quad (6.128)$$

Mit dem Prädikatenkalkül folgt wie behauptet

$$l_i \Rightarrow \bigwedge_{l=1}^{n_v^3} \bigwedge_{k=1}^{n_e} a_{vl}^3.typ = b_{ek}.typ . \quad (6.129)$$

Zudem gilt $l_i \Rightarrow \delta.typ = \text{int} \wedge \delta.l = -1$, so dass die Forderungen bzgl. der Datentypen aus Gleichung (6.127) abgetrennt werden können. Nun werden wir zeigen, dass die Bedingungen an die Feldlängen abgetrennt werden können. Einerseits folgt aus Gleichung (6.119) , Gleichung (6.114) , Gleichung (6.113), Gleichung (6.34) und Gleichung (6.33)

$$l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \Rightarrow \bigwedge_{k=1}^{n_v^3} P - 1 \leq a_{vk}^3.l - 1 \wedge \bigwedge_{k=1}^{n_e} P - 1 \leq b_{ek}.l - 1 \quad (6.130)$$

und andererseits gilt wegen $Z_i \Rightarrow P_{\text{MEDD}}$ die Beziehung

$$Z'_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \Rightarrow 0 \leq \delta.v \leq P - 1 . \quad (6.131)$$

Somit erhalten wir

$$l_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \Rightarrow \bigwedge_{k=1}^{n_v^3} 0 \leq \delta.v \leq a_{vk}^3.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 . \quad (6.132)$$

Dies ist der in Gleichung (6.127) vorhandene Teil der Vorbedingung, der somit abgetrennt werden kann. Die Bedingung l_i impliziert

$$\bigwedge_{\nu=1}^{n_e} |b_{e\nu}(\delta.v).v| \leq \text{Max}_f . \quad (6.133)$$

Von Gleichung (6.127) verbleibt somit als Vorbedingung

$$\begin{aligned} l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) \wedge \\ \mathbf{P}_{v^3e} \mathbf{b}_e(\delta.v).v = \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) . \end{aligned} \quad (6.134)$$

Den letzten Konjunktionsterm ersetzen wir mit

$$\mathbf{b}_e(\delta.v).v = \mathbf{b}_e(\boldsymbol{\mu}(\delta.v)) = \mathbf{b}_e([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) \quad (6.135)$$

zu

$$\begin{aligned} l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) \wedge \\ \mathbf{P}_{v^3e} \mathbf{b}_e([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) = \mathbf{a}_v^3([\mu_1.v, \mu_2.v, P_{x_3} - 1]^T) . \end{aligned} \quad (6.136)$$

Berücksichtigen wir $l_i \Rightarrow \mathbf{P}_{\text{WDF}}$, so haben wir schließlich

$$l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) . \quad (6.137)$$

Die Vorbedingung der Anweisung $\text{delta} = \text{mu1} + \text{mu2} * P_{x_1} + \text{mu3} * P_{x_1} * P_{x_2}$; lautet mit Gleichung (6.41)

$$\begin{aligned} l_i \wedge b_i \wedge \mu_1.v + P_{x_1}\mu_2.v + P_{x_1}P_{x_2}\mu_3.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = P_{x_3} - 1 \wedge \\ \mu_1.v \leq P_{x_1} - 1 \wedge \mu_2.v \leq P_{x_2} - 1 \wedge \mu_3.v \leq P_{x_3} - 1 \wedge \\ \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) . \end{aligned} \quad (6.138)$$

Trennen wir zudem die von $l_i \wedge b_i$ implizierten Konjunktionsterme ab, so ist die Vorbedingung von $\text{mu3} = P_{x_3} - 1$

$$l_i \wedge b_i \wedge P_{x_3} - 1 = P_{x_3} - 1 \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) \wedge P_{x_3} - 1 \leq \text{Max}_i . \quad (6.139)$$

Berücksichtigung von $l_i \Rightarrow P_{x_3} \leq \text{Max}_i$ (nach Gleichung (6.34)) liefert schließlich

$$l_i \wedge b_i \wedge \mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T) . \quad (6.140)$$

Durch analoge Überlegungen erreichen wir durch S_{av6} eine Abtrennung von $\mathbf{a}_v^6(\delta([\mu_1.v, \mu_2.v, 0]^T)).v = \mathbf{a}_v^6([\mu_1.v, \mu_2.v, 0]^T)$. Wir bekommen als Vorbedingung von S_{av6} im Wesentlichen

$$l_i \wedge b_i \wedge \delta.v = \delta(\boldsymbol{\mu}.v) \wedge \mu_3.v = 0 . \quad (6.141)$$

Durch die beiden zuvor ausgeführten Anweisungen können $\delta.v = \delta(\boldsymbol{\mu}.v)$ und $\mu_3.v = 0$ abgetrennt werden. Die Vorbedingung von $\text{mu3}=0$; und somit auch die von S_{4i} ist, wie behauptet, durch $l_i \wedge b_i$ gegeben.

Die Vorbedingungen der Anweisungen S_{avx} und S_{avy} werden in gleicher Weise ermittelt. Auf die Beweisführung wird an dieser Stelle verzichtet.

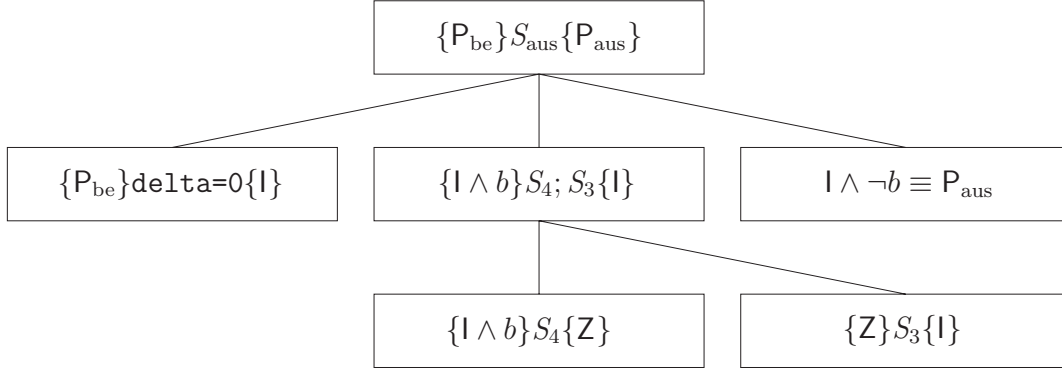


Bild 6.11: Berechnung der Eingangswellen der Verzögerer

6.5 Berechnung der Ausgangssignale

Die Berechnung der Ausgangssignale \mathbf{a}_{FB} des Funktionsbausteins für jeden räumlichen Gitterpunkt gemäß Gleichung (6.2) erfolgt mittels der Anweisung S_{aus} . Die Nachbedingung lautet

$$P_{aus} \equiv P_{WDF} \wedge P_{dec} \wedge P_q \wedge P'_{be} \wedge P'_{aus}. \quad (6.142)$$

Die Berechnung über alle Gitterpunkte erfolgt ebenfalls durch Verwendung einer for-Schleife der Länge P und der Laufvariablen δ . Wir verwenden dazu den folgenden C-Code

Bedingung	Anweisung	Kommentar
$\{P_{be}\}$	<code>for(delta=0 ; delta < P ; delta=delta+1){</code>	
$\{I\}$	<code>S_{aFB1}</code>	Invariante
	<code>⋮</code>	Berechnung von a_{FB1}
	<code>S_{aFBn_{aa}}</code>	
	<code>}</code>	Berechnung von a_{FBnaa}
$\{P_{aus}\}$		

wobei die Berechnungsanweisungen durch

$$S_{aFB1} \equiv \mathbf{a}_{FB_1}[\delta] = \begin{matrix} a_{11} & *b_{q_1}[\delta] & + \dots + a_{1n_q} & *b_{q_n_q}[\delta] + \\ a_{1n_q+1} & *b_{v_1_1}[\delta] + \dots + a_{1n_q+n_v} & *b_{v_7_n_v}[\delta]; \\ a_{1n_q+n_v+1} & *b_{e_1}[\delta] & + \dots + a_{1n_q+n_v+n_e} & *b_{e_n_e}[\delta]; \end{matrix}$$

bis

$$S_{aFBn_{aa}} \equiv \mathbf{a}_{FB_n_{aa}}[\delta] = \begin{matrix} a_{n_{aa}1} & *b_{q_1}[\delta] & + \dots + a_{n_{aa}n_q} & *b_{q_n_q}[\delta] + \\ a_{n_{aa}n_q+1} & *b_{v_1_1}[\delta] + \dots + a_{n_{aa}n_q+n_v} & *b_{v_7_n_v}[\delta]; \\ a_{n_{aa}n_q+n_v+1} & *b_{e_1}[\delta] & + \dots + a_{n_{aa}n_q+n_v+n_e} & *b_{e_n_e}[\delta]; \end{matrix}$$

gegeben sind. Im Bild 6.11 ist eine Übersicht des Beweises (der Ermittlung der Vorbedingung) dargestellt. Hierin ist die Schleifenbedingung $b \equiv \delta.v < P$.

Satz 10 Der Ausdruck

$$I \equiv 0 \leq \delta.v \leq P \wedge P_{be} \wedge \bigwedge_{\delta=0}^{\delta.v-1} \bigwedge_{k=1}^{n_{aa}} P_{aFBk}(\mu(\delta), \delta) \quad (6.143)$$

mit

$$P_{be} \equiv P_{WDF} \wedge P_{dec} \wedge P_q \wedge P_v \wedge P'_{be} \quad (6.144)$$

ist eine gültige Schleifeninvariante der Schleife.

Beweis 10

Anwenden von Axiom (5.18) auf die Initialisierungsanweisung `delta=0` liefert

$$0 \leq P \wedge P_{be} \wedge \delta.l = -1 \wedge \delta.typ = \text{int} \wedge \bigwedge_{\delta=0}^{-1} \bigwedge_{k=1}^{n_{aa}} P_{aFBk}(\boldsymbol{\mu}(\delta), \delta) \equiv P_{be}. \quad (6.145)$$

Die Nachbedingung der Schleife lautet

$$I \wedge \delta.v \geq P \equiv \delta.v = P \wedge P_{be} \wedge \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_{aa}} P_{aFBk}(\boldsymbol{\mu}(\delta), \delta) \quad (6.146)$$

und ist nach Schwächung um $\delta.v = P$ und P_v die Nachbedingung P_{aus} nach Gleichung (6.142). Die Zwischenbedingung der Schleife lautet

$$Z \equiv 0 \leq \delta.v+1 \leq P \wedge P_{be} \wedge \bigwedge_{\delta=0}^{\delta.v} \bigwedge_{k=1}^{n_{aa}} P_{aFBk}(\boldsymbol{\mu}(\delta), \delta) \wedge \delta.v+1 \leq \text{Max}_i \wedge \delta.l = -1 \wedge \delta.typ = \text{int}. \quad (6.147)$$

Wegen $P_{be} \implies P \leq \text{Max}_i$ kann $\delta.v + 1 \leq \text{Max}_i$ abgetrennt werden. Mit der Stärkung um $0 \leq \delta.v$ und einigen Vereinfachungen erhalten wir

$$Z' \equiv Z \wedge 0 \leq \delta.v \equiv I \wedge b \wedge \bigwedge_{k=1}^{n_{aa}} P_{aFBk}(\boldsymbol{\mu}(\delta.v), \delta.v). \quad (6.148)$$

Es muss nun die Richtigkeit von $\{I \wedge b\} S_4 \{Z'\}$ mit $S_4 \equiv S_{aFB1}; S_{aFB2}; \dots S_{aFBn_{aa}}$ gezeigt werden. Wir bestimmen die Vorbedingung von S_{aFB1} mit Axiom (5.43) zu

$$\begin{aligned} & I \wedge b \wedge \bigwedge_{k=1}^{n_{aa}} a_{FBk}(\delta.v).v = a_{FBk}(\boldsymbol{\mu}(\delta.v)) \left[a_{FBk}(\delta.v).v \right. \\ & \quad \left. \left[\mathbf{a}_{k,\bullet} [\mathbf{b}_q^T(\delta.v).v, \mathbf{b}_v^{1T}(\delta.v).v, \dots, \mathbf{b}_v^{7T}(\delta.v).v, \mathbf{b}_e^T(\delta.v).v]^T \right] \right] \\ & \quad \bigwedge_{k=1}^{n_{aa}} \left[\bigwedge_{m=1}^{n_q} \left| \sum_{\mu=1}^m a_{k\mu} b_{q\mu}(\delta.v).v \right| \leq \text{Max}_f \wedge \bigwedge_{l=1}^7 \bigwedge_{m=1}^{n_v^l} \left| \sum_{\mu=1}^{n_q} a_{k\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^l \sum_{\mu=1}^m a_{k*} b_{v\mu}^\kappa(\delta.v).v \right| \leq \text{Max}_f \wedge \right. \\ & \quad \left. \bigwedge_{m=1}^{n_e} \left| \sum_{\mu=1}^{n_q} a_{k\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^7 \sum_{\mu=1}^{n_v^\kappa} a_{k*} b_{v\mu}^\kappa(\delta.v).v + \sum_{\mu=1}^m a_{k*} b_{e\mu}(\delta.v).v \right| \leq \text{Max}_f \right] \wedge \\ & \quad \bigwedge_{k=1}^{n_{aa}} \left[\bigwedge_{\mu=1}^{n_q} |a_{k\mu} b_{q\mu}(\delta.v).v| \leq \text{Max}_f \wedge \bigwedge_{l=1}^7 \bigwedge_{\mu=1}^{n_v^l} |a_{k*} b_{v\mu}^l(\delta.v).v| \leq \text{Max}_f \wedge \bigwedge_{\mu=1}^{n_e} |a_{k*} b_{e\mu}(\delta.v).v| \leq \text{Max}_f \right] \wedge \\ & \quad \bigwedge_{k=1}^{n_{aa}} 0 \leq \delta.v \leq a_{FBk}.l-1 \wedge \bigwedge_{k=1}^{n_q} 0 \leq \delta.v \leq b_{qk}.l-1 \wedge \bigwedge_{m=1}^7 \bigwedge_{\mu=1}^{n_v^m} 0 \leq \delta.v \leq b_{v\mu}^m.l-1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l-1 \wedge \\ & \quad \bigwedge_{l=1}^{n_{aa}} \left[\bigwedge_{k=1}^{n_q} a_{FBk}.typ = b_{qk}.typ \wedge \bigwedge_{k=1}^{n_e} a_{FBk}.typ = b_{ek}.typ \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} a_{FBk}.typ = b_{v\kappa}^\kappa.typ \right] \wedge \\ & \quad \delta.v \leq \text{Max}_i \wedge \delta.l = -1 \wedge \delta.typ = \text{int}, \end{aligned} \quad (6.149)$$

wobei der Wert von $*$ der Übersicht halber nicht explizit angegeben wird. Zur Verdichtung des obigen Prädikats werden wir zeigen, dass l die Prädikate bzgl. der Datentypen und der Feldlänge impliziert. Aus Gleichung (6.143), Gleichung (6.144), Gleichung (6.33) und Gleichung (6.34) folgt

$$l \Rightarrow \bigwedge_{k=1}^{n_{aa}} a_{\text{FB}k}.typ = \text{float} \wedge \bigwedge_{k=1}^{n_q} b_{qk}.typ = \text{float} \wedge \bigwedge_{k=1}^{n_e} b_{ek}.typ = \text{float} \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} b_{v\kappa}^\kappa.typ = \text{float} . \quad (6.150)$$

Durch Anwendung des Prädikatenkalküls folgt

$$\bigwedge_{l=1}^{n_{aa}} \left[\bigwedge_{k=1}^{n_q} a_{\text{FB}l}.typ = b_{qk}.typ \wedge \bigwedge_{k=1}^{n_e} a_{\text{FB}l}.typ = b_{ek}.typ \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} a_{\text{FB}l}.typ = b_{v\kappa}^\kappa.typ \right] \quad (6.151)$$

und letzteres Prädikat kann aus Gleichung (6.149) abgetrennt werden. Ebenso kann aus Gleichung (6.143), Gleichung (6.144), Gleichung (6.33) und Gleichung (6.34)

$$l \wedge b \Rightarrow \bigwedge_{k=1}^{n_{aa}} P - 1 \leq a_{\text{FB}k}.l - 1 \wedge \bigwedge_{k=1}^{n_q} P - 1 \leq b_{qk}.l - 1 \wedge \bigwedge_{k=1}^{n_e} P - 1 \leq b_{ek}.l - 1 \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} P - 1 \leq b_{v\kappa}^\kappa.l - 1 \wedge P \leq \text{Max}_i \quad (6.152)$$

gewonnen werden. Mit $\delta.v + 1 \leq P \equiv \delta.v \leq P - 1$ folgt

$$l \wedge b \Rightarrow \bigwedge_{k=1}^{n_{aa}} 0 \leq \delta.v \leq a_{\text{FB}k}.l - 1 \wedge \bigwedge_{k=1}^{n_q} 0 \leq \delta.v \leq b_{qk}.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} 0 \leq \delta.v \leq b_{v\kappa}^\kappa.l - 1 \wedge \delta.v + 1 \leq \text{Max}_i . \quad (6.153)$$

Die Prädikate bzgl. der Länge und $\delta.v \leq \text{Max}_i$ können somit auch aus Gleichung (6.149) abgetrennt werden. Ferner gilt

$$l \Rightarrow \delta.typ = \text{int} \wedge \delta.l = -1 . \quad (6.154)$$

Wir zeigen nun, dass aus der Vorbedingung P_{aus} die Beschränktheit der Summanden und der Summen folgt. Da die Matrix \mathbf{A} eine spezielle Struktur hat, ergibt sich der (Soll-)Wert eines Ausgangssignals $a_{\text{FB}k}(\delta)$ entweder durch

$$u_\mu = \sqrt{R_\mu} [a_\mu + b_\mu] \quad \text{oder durch} \quad i_\mu = \sqrt{G_\mu} [a_\mu - b_\mu] . \quad (6.155)$$

Zunächst weisen wir nach, dass jeder Summand beschränkt ist. Dazu berücksichtigen wir

$$l \Rightarrow P_{\text{be}} \Rightarrow \bigwedge_{\delta=0}^{P-1} \left[\bigwedge_{\mu=1}^{n_q} |b_{q\mu}(\delta.v).v| 2\beta \leq \text{Max}_f \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{\mu=1}^{n_v^\kappa} |b_{v\mu}^\kappa(\delta.v).v| 2\beta \leq \text{Max}_f \wedge \bigwedge_{\mu=1}^{n_e} |b_{e\mu}(\delta.v).v| 2\beta \leq \text{Max}_f \right] \quad (6.156)$$

und durch Schwächung

$$l \Rightarrow P_{\text{be}} \Rightarrow \bigwedge_{\delta=0}^{P-1} \left[\bigwedge_{\mu=1}^{n_q} |b_{q\mu}(\delta.v).v| \beta \leq \text{Max}_f \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{\mu=1}^{n_v^\kappa} |b_{v\mu}^\kappa(\delta.v).v| \beta \leq \text{Max}_f \wedge \bigwedge_{\mu=1}^{n_e} |b_{e\mu}(\delta.v).v| \beta \leq \text{Max}_f \right] .$$

(6.157)

Beachten wir die Definition von β nach Gleichung (6.38), so haben wir

$$\begin{aligned} \mathbf{I} \Rightarrow \mathbf{P}_{\text{be}} \Rightarrow \bigwedge_{\delta=0}^{P-1} \left[\bigwedge_{\mu=1}^{n_g} \left[\sqrt{R_\mu} |b_{q\mu}(\delta.v).v| \leq \text{Max}_f \wedge \sqrt{G_\mu} |b_{q\mu}(\delta.v).v| \leq \text{Max}_f \right] \wedge \right. \\ \bigwedge_{\kappa=1}^7 \bigwedge_{\mu=1}^{n_v^\kappa} \left[\sqrt{R_\mu} |b_{v\mu}^\kappa(\delta.v).v| \leq \text{Max}_f \wedge \sqrt{G_\mu} |b_{v\mu}^\kappa(\delta.v).v| \leq \text{Max}_f \right] \wedge \\ \left. \bigwedge_{\mu=1}^{n_e} \left[\sqrt{R_\mu} |b_{e\mu}(\delta.v).v| \leq \text{Max}_f \wedge \sqrt{G_\mu} |b_{e\mu}(\delta.v).v| \leq \text{Max}_f \right] \right], \end{aligned} \quad (6.158)$$

d.h. die Tatsache, dass jede Welle a_μ gleich einer Welle b_ν ist und erhalten mittels Schwächung den gesuchten Nachweis der Beschränktheit jedes Summanden.

Wir wollen im Folgenden die Beschränktheit der Summen nachweisen, d.h. $|u_\mu| \leq \text{Max}_f$ und $|i_\mu| \leq \text{Max}_f$. Mit der Dreiecksungleichung erhalten wir aus Gleichung (6.155)

$$|u_\mu| \leq \sqrt{R_\mu} [|a_\mu| + |b_\mu|] \quad , \quad |i_\mu| \leq \sqrt{G_\mu} [|a_\mu| + |b_\mu|] . \quad (6.159)$$

Wenden wir darauf die Gleichung (6.158) an, so offenbart dies, dass die Beträge der Torspannungen und Torströme, die den Ausgangssignalen $a_{\text{FB}k}(\delta)$ entsprechen, durch

$$\bigwedge_{\delta=0}^{P-1} \bigwedge_{\mu=1}^{n_g} \left[|a_{\text{FB}k}(\delta)| \leq \sqrt{R_\mu} \left[\frac{\text{Max}_f}{2\beta} + \frac{\text{Max}_f}{2\beta} \right] = \frac{\sqrt{R_\mu}}{\beta} \text{Max}_f \wedge |a_{\text{FB}k}(\delta)| \leq \sqrt{G_\mu} \left[\frac{\text{Max}_f}{2\beta} + \frac{\text{Max}_f}{2\beta} \right] = \frac{\sqrt{G_\mu}}{\beta} \text{Max}_f \right] \quad (6.160)$$

begrenzt sind. Aus der Definition von β folgt nun

$$\sqrt{R_\mu} \leq \beta \iff \frac{\sqrt{R_\mu}}{\beta} \text{Max}_f \leq \text{Max}_f \quad \wedge \quad \sqrt{G_\mu} \leq \beta \iff \frac{\sqrt{G_\mu}}{\beta} \text{Max}_f \leq \text{Max}_f . \quad (6.161)$$

Demnach gilt auch

$$\bigwedge_{\delta=0}^{P-1} \bigwedge_{\mu=1}^{n_g} \left[\frac{\sqrt{R_\mu}}{\beta} \text{Max}_f \leq \text{Max}_f \wedge \frac{\sqrt{G_\mu}}{\beta} \text{Max}_f \leq \text{Max}_f \right] . \quad (6.162)$$

Konjunktive Verknüpfung von Gleichung (6.160), Gleichung (6.162) und Beachtung von $0 \leq \delta.v \leq P-1$ ergibt

$$\bigwedge_{\mu=1}^{n_g} [|a_{\text{FB}k}(\delta.v)| \leq \text{Max}_f \wedge |a_{\text{FB}k}(\delta.v)| \leq \text{Max}_f] . \quad (6.163)$$

Es verbleibt somit von Gleichung (6.149)

$$\mathbf{I} \wedge \mathbf{b} \wedge \bigwedge_{k=1}^{n_{aa}} \mathbf{a}_{k,\bullet} [\mathbf{b}_q^T(\delta.v).v, \mathbf{b}_v^{1T}(\delta.v).v, \dots, \mathbf{b}_v^{7T}(\delta.v).v, \mathbf{b}_e^T(\delta.v).v]^T = a_{\text{FB}k}(\boldsymbol{\mu}(\delta.v)) . \quad (6.164)$$

Mit

$$\mathbf{I} \Rightarrow \mathbf{b}_q(\delta.v).v = \mathbf{b}_q(\boldsymbol{\mu}(\delta.v)) \wedge \mathbf{b}_e(\delta.v).v = \mathbf{b}_e(\boldsymbol{\mu}(\delta.v)) \wedge \bigwedge_{\kappa=1}^7 \mathbf{b}_v^\kappa(\delta.v).v = \mathbf{b}_v^\kappa(\boldsymbol{\mu}(\delta.v)) \quad (6.165)$$

können wir in Gleichung (6.164) die prädikatenlogischen Variablen $\mathbf{b}(\delta.v).v$ durch die Sollwerte $\mathbf{b}(\boldsymbol{\mu}(\delta.v))$ ersetzen. Anschließend nutzen wir $\mathbf{I} \Rightarrow \mathbf{P}_{\text{WDF}}$ und trennen die Skalarprodukte ab. Wir erhalten, wie behauptet, $\mathbf{I} \wedge \mathbf{b}$.

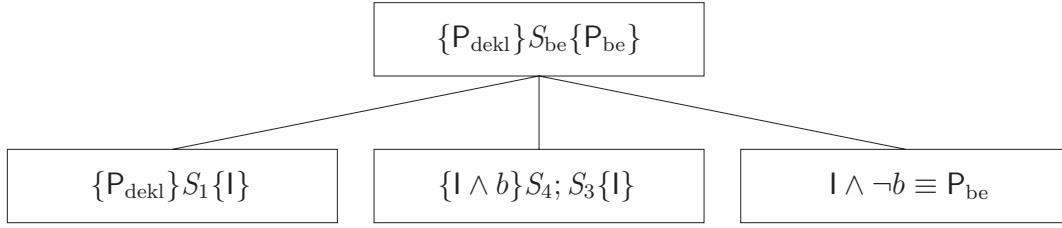


Bild 6.12: Berechnung der Wellengrößen der nichtdynamischen Elemente

6.6 Berechnung der Wellengrößen der nichtdynamischen Elemente

Wir behandeln nun die Anweisung S_{be} , die die Berechnung der Wellengrößen der nichtdynamischen Elemente \mathbf{b}_e durchführt. Nach Gleichung (6.1) sind die Wellengrößen \mathbf{b}_e für alle Gitterpunkte des Berechnungsgebietes zu ermitteln. Wir gehen davon aus, dass vor dem Aufruf des Funktionsbausteins zu jedem Gitterpunkt der Abtastschicht die Wellengrößen \mathbf{b}_v und \mathbf{b}_q bekannt sind. Die Nachbedingung von S_{be} lautet

$$P_{be} \equiv P_{WDF} \wedge P_{dec} \wedge P_q \wedge P_v \wedge P'_{be}. \quad (6.166)$$

Die Auswertung von Gleichung (6.1) erfolgt innerhalb einer for-Schleife mit der Laufvariablen $\delta.v$ von 0 bis $P - 1$, sodass sich die Berechnung über alle Gitterpunkte der Abtastschicht erstreckt.

Die zusätzliche Nachbedingung P'_{be} besagt, dass die Werte der Variablen für die nichtdynamischen Elemente denen aus Gleichung (6.1) für den aktuellen Zeitpunkt μ_4 und für $\delta.v = 0, 1, \dots, P - 1$ entsprechen müssen. Anders formuliert lautet die Nachbedingung P_{be}

$$P_{be} \equiv P_{dekl} \wedge \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v| 2\beta \leq \text{Maxf}] . \quad (6.167)$$

Zur Erfüllung der Bedingungen nutzen wir die folgende Anweisung in C-Code

Bedingung	Anweisung	Kommentar
$\{P_{dekl}\}$		
$\{I\}$	for(delta=0 ; delta < P ; delta=delta+1){	Invariante
	S_{be1}	Berechnung von \mathbf{b}_{e1}
	\vdots	
	S_{ben_e}	Berechnung von \mathbf{b}_{en_e}
	}	
$\{P_{be}\}$		

mit den Zuweisungsanweisungen

$$S_{be1} \equiv \mathbf{b}_{e1}[\text{delta}] = l_{11} * \mathbf{b}_{q1}[\text{delta}] + \dots + l_{1n_q} * \mathbf{b}_{qn_q}[\text{delta}] + l_{1(n_q+1)} * \mathbf{b}_{v1_1}[\text{delta}] + \dots + l_{1(n_q+n_v)} * \mathbf{b}_{v6_n_v}[\text{delta}]; \quad (6.168)$$

und für $k = 2 \dots n_e$

$$S_{bek} \equiv \mathbf{b}_{ek}[\text{delta}] = l_{k1} * \mathbf{b}_{q1}[\text{delta}] + \dots + l_{kn_q} * \mathbf{b}_{qn_q}[\text{delta}] + l_{k(n_q+1)} * \mathbf{b}_{v1_1}[\text{delta}] + \dots + l_{k(n_q+n_v)} * \mathbf{b}_{v6_n_v}[\text{delta}] + l_{k(n_q+n_v+1)} * \mathbf{b}_{e1}[\text{delta}] + \dots + l_{k(n_q+n_v+k-1)} * \mathbf{b}_{ek-1}[\text{delta}]; . \quad (6.169)$$

Wir führen nun den Nachweis der Korrektheit von $\{P_{\text{dekl}}\} S_{\text{be}} \{P_{\text{be}}\}$, vgl. dazu die Übersicht im Bild 6.12 mit der Schleifenbedingung $b \equiv \delta.v < P \equiv \delta.v \leq P - 1$.

Satz 11 *Der Ausdruck*

$$I \equiv 0 \leq \delta.v \leq P \wedge P_{\text{dekl}} \wedge \bigwedge_{\delta=0}^{\delta.v-1} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v|2\beta \leq \text{Max}_f] \quad (6.170)$$

mit

$$P_{\text{dekl}} \equiv P_{\text{WDF}} \wedge P_{\text{dec}} \wedge P_q \wedge P_v \quad (6.171)$$

ist eine gültige Invariante der Schleife.

Beweis 11

Gemäß Axiom (5.35) können wir den Nachweis auf das Zeigen der Initialisierung $\{P_{\text{dekl}}\} S_1 \{I\}$, des Schleifenendes $P_{\text{be}} \equiv I \wedge \neg b$ und der Invarianz $\{I \wedge b\} S_4; S_3 \{I\}$ zurückführen.

Die Vorbedingung von $\text{delta}=0$; ermittelt sich zu

$$0 \leq P \wedge P_{\text{dekl}} \wedge \bigwedge_{\delta=0}^{-1} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v|2\beta \leq \text{Max}_f] \wedge \delta.l = -1 \wedge \delta.\text{typ} = \text{int} . \quad (6.172)$$

Dies entspricht der Vorbedingung P_{dekl} , da P_{dekl} die Bedingung $0 \leq P$ impliziert. Die Abbruchbedingung lautet $\neg b \equiv \delta.v \geq P$, sodass wir als Nachbedingung

$$I \wedge \delta.v \geq P \equiv \delta.v = P \wedge P_{\text{dekl}} \wedge \bigwedge_{\delta=0}^{P-1} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v|2\beta \leq \text{Max}_f] \quad (6.173)$$

erhalten, was nach Schwächung um $\delta.v = P$, den Forderungen gemäß Gleichung (6.167) entspricht.

Für den Nachweis von $\{I \wedge b\} S_4; S_3 \{I\}$ bestimmen wir zunächst die durch $\{Z\} S_3 \{I\}$ definierte Zwischenbedingung Z und bereiten diese durch Nutzen von Gleichung (6.170) zu

$$\begin{aligned} Z &\equiv 0 \leq \delta.v + 1 \leq P \wedge P_{\text{dekl}} \wedge \bigwedge_{\delta=0}^{\delta.v} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v|2\beta \leq \text{Max}_f] \wedge \\ &\quad \delta.v \leq \text{Max}_i \wedge \delta.\text{typ} = \text{int} \wedge \delta.l = -1 \\ &\equiv 0 \leq \delta.v + 1 \leq P \wedge P_{\text{dekl}} \wedge \bigwedge_{\delta=0}^{\delta.v-1} \bigwedge_{k=1}^{n_e} [b_{ek}(\delta).v = b_{ek}(\boldsymbol{\mu}(\delta)) \wedge |b_{ek}(\delta).v|2\beta \leq \text{Max}_f] \wedge \\ &\quad \bigwedge_{k=1}^{n_e} [b_{ek}(\delta.v).v = b_{ek}(\boldsymbol{\mu}(\delta.v)) \wedge |b_{ek}(\delta.v).v|2\beta \leq \text{Max}_f] \wedge \delta.v \leq \text{Max}_i \wedge \delta.\text{typ} = \text{int} \wedge \delta.l = -1 \end{aligned} \quad (6.174)$$

auf. Wir stärken die Vorbedingung zudem durch $0 \leq \delta.v$ und erhalten nach einigen Vereinfachungen

$$Z' \equiv Z \wedge 0 \leq \delta.v \equiv I \wedge b \wedge \bigwedge_{k=1}^{n_e} [b_{ek}(\delta.v).v = b_{ek}(\boldsymbol{\mu}(\delta.v)) \wedge |b_{ek}(\delta.v).v|2\beta \leq \text{Max}_f] . \quad (6.175)$$

Es verbleibt $\{l \wedge b\} S_4 \{Z'\}$ zu zeigen. Zunächst behandeln wir die Anweisungen S_{bek} , $k = 1, \dots, n_e$. Den Beweis haben wir mit Axiom (5.7) erbracht, wenn wir mit

$$\begin{aligned} V_{S_{bek}} &\equiv l \wedge b \wedge \bigwedge_{\kappa=1}^{k-1} [b_{e\kappa}(\delta.v).v = b_{e\kappa}(\mu(\delta.v)) \wedge |b_{e\kappa}(\delta.v).v|2\beta \leq \text{Max}_f] \\ P_{S_{bek}} &\equiv l \wedge b \wedge \bigwedge_{\kappa=1}^k [b_{e\kappa}(\delta.v).v = b_{e\kappa}(\mu(\delta.v)) \wedge |b_{e\kappa}(\delta.v).v|2\beta \leq \text{Max}_f] \\ &\equiv V_{S_{bek}} \wedge b_{ek}(\delta.v).v = b_{ek}(\mu(\delta.v)) \wedge |b_{ek}(\delta.v).v|2\beta \leq \text{Max}_f. \end{aligned} \quad (6.176)$$

die Korrektheit von $l \wedge b \equiv V_{S_{be1}}$, $Z' \equiv P_{S_{ben_e}}$ und $\{V_{S_{bek}}\} S_{bek} \{P_{S_{bek}}\}$ nachweisen. Das entspricht dem Zeigen der Richtigkeit von

$$\{V_{S_{be1}}\} \underbrace{S_{be1}; \dots S_{ben_e}}_{S_4} \{P_{S_{ben_e}}\}. \quad (6.177)$$

Die Gültigkeit von $l \wedge b \equiv V_{S_{be1}}$ ist offensichtlich. Weiterhin gilt

$$P_{S_{ben_e}} \equiv l \wedge b \wedge \bigwedge_{k=1}^{n_e} [b_{ek}(\delta.v).v = b_{ek}(\mu(\delta.v)) \wedge |b_{ek}(\delta.v).v|2\beta \leq \text{Max}_f] \quad (6.178)$$

und ist Z' wie zu zeigen war. Es verbleibt $\{V_{S_{bek}}\} S_{bek} \{P_{S_{bek}}\}$ zu zeigen. Die Vorbedingung von S_{bek} ergibt sich zu

$$\begin{aligned} &V_{bek} \wedge b_{ek}(\delta.v).v = b_{ek}(\mu(\delta.v)) \left[b_{ek}(\delta.v).v \right. \\ &\quad \left. |b_{ek}(\delta.v).v|2\beta \leq \text{Max}_f \right] \wedge \\ &\quad \bigwedge_{m=0}^{n_q} \left| \sum_{\mu=1}^m l_{qk\mu} b_{q\mu}(\delta.v).v \right| \leq \text{Max}_f \wedge \\ &\quad \bigwedge_{l=1}^7 \bigwedge_{m=1}^{n_v^l} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^l \sum_{\mu=1}^m l_{vk\kappa} b_{v\mu}^\kappa(\delta.v).v \right| \leq \text{Max}_f \wedge \\ &\quad \bigwedge_{m=1}^{n_e} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^7 \sum_{\mu=1}^{n_v^\kappa} l_{vk\kappa} b_{v\mu}^\kappa(\delta.v).v + \sum_{\mu=1}^m l_{ek\kappa} b_{e\mu}(\delta.v).v \right| \leq \text{Max}_f \wedge \\ &\quad \bigwedge_{\mu=1}^{n_q} |l_{qk\mu} b_{q\mu}(\delta.v).v| \leq \text{Max}_f \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{\mu=1}^{n_v^\kappa} |l_{vk\kappa} b_{v\mu}^\kappa(\delta.v).v| \leq \text{Max}_f \wedge \bigwedge_{\mu=1}^{n_e} |l_{ek\kappa} b_{e\mu}(\delta.v).v| \leq \text{Max}_f \wedge \\ &\quad \bigwedge_{k'=1}^{n_q} 0 \leq \delta.v \leq b_{qk'}.l - 1 \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k'=1}^{n_v^\kappa} 0 \leq \delta.v \leq b_{vk'}^\kappa.l - 1 \wedge \bigwedge_{k'=1}^{n_e} 0 \leq \delta.v \leq b_{ek'}.l - 1 \wedge \\ &\quad \bigwedge_{k'=1}^{n_q} b_{ek}.typ = b_{qk'}.typ \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k'=1}^{n_v^\kappa} b_{ek}.typ = b_{vk'}^\kappa.typ \wedge \bigwedge_{k'=1}^{k-1} b_{ek}.typ = b_{ek'}.typ \wedge \\ &\quad \delta.v + 1 \leq \text{Max}_i \wedge \delta.l = -1 \wedge \delta.typ = \text{int}. \end{aligned} \quad (6.179)$$

Zur Verdichtung des obigen Prädikats werden wir zeigen, dass $l \wedge b$ die Prädikate bzgl. der Datentypen und der Feldlänge impliziert. Mit Gleichung (6.170), Gleichung (6.171), Gleichung (6.33) und Gleichung

(6.34) erhalten wir zum einen

$$l \wedge b \Rightarrow \bigwedge_{k=1}^{n_q} b_{qk}.typ = \text{float} \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} b_{v\kappa}^\kappa.typ = \text{float} \wedge \bigwedge_{k=1}^{n_e} b_{ek}.typ = \text{float}, \quad (6.180)$$

was somit aus Gleichung (6.179) abgetrennt werden kann und zum anderen

$$l \wedge b \Rightarrow \bigwedge_{k=1}^{n_q} P - 1 \leq b_{qk}.l - 1 \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} P - 1 \leq b_{v\kappa}^\kappa.l - 1 \wedge \bigwedge_{k=1}^{n_e} P - 1 \leq b_{ek}.l - 1 \wedge P \leq \text{Max}_i. \quad (6.181)$$

Daraus gewinnen wir mit $l \wedge b \Rightarrow 0 \leq \delta.v \leq P - 1$

$$l \wedge b \Rightarrow \bigwedge_{k=1}^{n_q} 0 \leq \delta.v \leq b_{qk}.l - 1 \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} 0 \leq \delta.v \leq b_{v\kappa}^\kappa.l - 1 \wedge \bigwedge_{k=1}^{n_e} 0 \leq \delta.v \leq b_{ek}.l - 1 \wedge \delta.v \leq \text{Max}_i \wedge \delta.typ = \text{int} \wedge \delta.l = -1, \quad (6.182)$$

was wiederum aus Gleichung (6.179) abgetrennt wird. Wir zeigen nun, dass die Beschränktheit der Teilskalarprodukte in Gleichung (6.179) durch $P_{\text{dekl}} \Rightarrow P_{\text{WDF}} \wedge P_{\text{max}}$ impliziert wird. Zum Beweis betrachten wir die Ergebnisse von Gleichung (C.11), Gleichung (C.15) und Gleichung (C.35) (für die PL-Variablen der Programmvariablen), die wiederum durch l impliziert werden. Wir haben somit

$$\bigwedge_{\delta.v=0}^{P-1} \bigwedge_{k=1}^{n_e} \left[\bigwedge_{m=0}^{n_q} \left| \sum_{\mu=1}^m l_{qk\mu} b_{q\mu}(\delta.v).v \right|^2 \beta^2 4 \leq \text{Max}_f^2 \wedge \bigwedge_{l=1}^7 \bigwedge_{m=1}^{n_v^l} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^l \sum_{\mu=1}^m l_{v\kappa*} b_{v\mu}^\kappa(\delta.v).v \right|^2 \beta^2 4 \leq \text{Max}_f^2 \wedge \bigwedge_{m=1}^{n_e} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^7 \sum_{\mu=1}^{n_v^\kappa} l_{v\kappa*} b_{v\mu}^\kappa(\delta.v).v + \sum_{\mu=1}^m l_{ek*} b_{e\mu}(\delta.v).v \right|^2 \beta^2 4 \leq \text{Max}_f^2 \right]. \quad (6.183)$$

Wir radizieren jeweils auf beiden Seiten, beachten, dass β nicht kleiner als 1 ist, berücksichtigen Gleichung (A.5) und erhalten dann

$$\bigwedge_{\delta.v=0}^{P-1} \bigwedge_{k=1}^{n_e} \left[\bigwedge_{m=0}^{n_q} \left| \sum_{\mu=1}^m l_{qk\mu} b_{q\mu}(\delta.v).v \right| \leq \text{Max}_f \wedge \bigwedge_{l=1}^7 \bigwedge_{m=1}^{n_v^l} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^l \sum_{\mu=1}^m l_{v\kappa*} b_{v\mu}^\kappa(\delta.v).v \right| \leq \text{Max}_f \wedge \bigwedge_{m=1}^{n_e} \left| \sum_{\mu=1}^{n_q} l_{qk\mu} b_{q\mu}(\delta.v).v + \sum_{\kappa=1}^7 \sum_{\mu=1}^{n_v^\kappa} l_{v\kappa*} b_{v\mu}^\kappa(\delta.v).v + \sum_{\mu=1}^m l_{ek*} b_{e\mu}(\delta.v).v \right| \leq \text{Max}_f \right], \quad (6.184)$$

sodass dieser Teil und

$$|l_{qk,\bullet} b_{q\bullet}(\delta.v).v + l_{v\kappa,\bullet} b_{v\bullet}(\delta.v).v + l_{ek,\bullet} b_{e\bullet}(\delta.v).v| 2\beta \leq \text{Max}_f \quad (6.185)$$

ebenfalls von Gleichung (6.179) abgetrennt werden kann. Aufgrund der von \mathbf{l} implizierten unitären Beschränktheit der Matrix \mathbf{L} ist jedes Element der Matrix dem Betrage nach nicht größer als eins. Folglich können auch Forderungen an die Beschränktheit der Summanden des Skalarprodukts aus Gleichung (6.179) abgetrennt werden. Die Vorbedingung ist somit äquivalent zu

$$\mathbf{V}_{\text{bek}} \wedge \mathbf{l}_{qk,\bullet} \mathbf{b}_q(\delta.v).v + \mathbf{l}_{vk,\bullet} \mathbf{b}_v(\delta.v).v + \mathbf{l}_{ek,\bullet} \mathbf{b}_e(\delta.v).v = b_{ek}(\boldsymbol{\mu}(\delta.v)) . \quad (6.186)$$

Aufgrund der strikten unteren Dreiecksgestalt von \mathbf{L} hat das dritte Skalarprodukt i. Allg. nur $k - 1$ Summanden

$$\mathbf{l}_{ek,\bullet} \mathbf{b}_e(\delta.v).v = \sum_{\nu=1}^{n_e} \mathbf{l}_{ek\nu} \mathbf{b}_{e\nu}(\delta.v).v = \sum_{\nu=1}^{k-1} \mathbf{l}_{ek\nu} \mathbf{b}_{e\nu}(\delta.v).v . \quad (6.187)$$

Mit

$$\mathbf{V}_{\text{Sbek}} \Rightarrow \bigwedge_{k=1}^{n_q} b_{qk}(\delta.v).v = b_{qk}(\boldsymbol{\mu}(\delta.v)) \wedge \bigwedge_{\kappa=1}^7 \bigwedge_{k=1}^{n_v^\kappa} b_{v\kappa}^\kappa(\delta.v).v = b_{v\kappa}^\kappa(\boldsymbol{\mu}(\delta.v)) \wedge \bigwedge_{\nu=1}^{k-1} b_{e\nu}(\delta.v).v = b_{e\nu}(\boldsymbol{\mu}(\delta.v)) \quad (6.188)$$

kann die Vorbedingung durch

$$\mathbf{V}_{\text{Sbek}} \wedge \mathbf{l}_{qk,\bullet} \mathbf{b}_q(\boldsymbol{\mu}(\delta.v)) + \mathbf{l}_{vk,\bullet} \mathbf{b}_v(\boldsymbol{\mu}(\delta.v)) + \mathbf{l}_{ek,\bullet} \mathbf{b}_e(\boldsymbol{\mu}(\delta.v)) = b_{ek}(\boldsymbol{\mu}(\delta.v)) \quad (6.189)$$

oder vereinfacht durch

$$\mathbf{V}_{\text{Sbek}} \wedge \mathbf{l}_{qk,\bullet} \mathbf{b}_q(\boldsymbol{\mu}(\delta.v)) + \mathbf{l}_{vk,\bullet} \mathbf{b}_v(\boldsymbol{\mu}(\delta.v)) + \sum_{\nu=1}^{k-1} \mathbf{l}_{ek\nu} \mathbf{b}_{e\nu}(\boldsymbol{\mu}(\delta.v)) = b_{ek}(\boldsymbol{\mu}(\delta.v)) \quad (6.190)$$

dargestellt werden. Mit $\mathbf{l} \Rightarrow \mathbf{P}_{\text{WDF}}$ und der Abtrennung der wahren Aussage verbleibt letztendlich von der Vorbedingung Gleichung (6.179), wie behauptet, nur \mathbf{V}_{Sbek} .

6.7 Variablendeklarationen

In den vorangegangenen Kapiteln hatten wir festgestellt, dass S_{dekl} die Nachbedingung

$$\mathbf{P}_{\text{dekl}} \equiv \mathbf{P}_{\text{WDF}} \wedge \mathbf{P}_{\text{dec}} \wedge \mathbf{P}_q \wedge \mathbf{P}_v \quad (6.191)$$

haben muss. Durch die Anweisung S_{dekl} erfolgt die Deklaration der Variablen. Für unseren Algorithmus unterscheiden wir zwei Typen von Variablen. Zum einen haben wir die Variablen, die schon vor dem Aufruf des Funktionsbausteins deklariert sind und beim Programmaufruf übergeben werden. Das sind die Ein- und Ausgangssignale und die Verzögererwerte, d.h. die Variablen, die den Vektoren \mathbf{b}_v , \mathbf{b}_q und \mathbf{a}_{FB} entsprechen. Deren Deklariertheit bildet einen Teil der Vorbedingung \mathbf{V} . Zum anderen haben wir die lokalen Variablen, d.h. Variablen, die innerhalb des Funktionsbausteins deklariert werden. Hierzu gehören diejenigen Variablen, die nicht für die nächste Abtastschicht verwendet werden, d.h. die Variablen, die den Vektoren \mathbf{b}_e , $\mathbf{a}_v^1, \dots, \mathbf{a}_v^6$ entsprechen und die Laufvariablen der for-Schleifen. Wir fassen nun die Forderungen an die gesamte Vorbedingung $\{\mathbf{V}\}$ des Algorithmus zusammen. Die Variablen der Verzögerer \mathbf{b}_v , der Quellen \mathbf{b}_q und der Ausgangssignale \mathbf{a}_{FB} müssen deklariert sein, die Feldlänge P haben und vom Typ `float` sein. Zudem müssen die Werte der Variablen für die Verzögerer und die Quellen zum aktuellen

Zeitpunkt μ_4 und zu jedem Ortspunkt mit den tatsächlichen Werten übereinstimmen. Ferner muss ein gültiges MDWDF vorliegen. Die Vorbedingung V lautet daher

$$V \equiv P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec}bq} \wedge P_{\text{dec}bv} \wedge P_{\text{deca}FB} \wedge P_{be} \wedge P_{\text{name}} \quad (6.192)$$

mit

$$P_{\text{name}} \equiv P_{\text{avname}} \wedge P_{\text{bename}} \wedge P_{\text{uvname}} \quad (6.193)$$

und

$$\begin{aligned} P_{\text{avname}} &\equiv \bigwedge_{\kappa=1}^6 \bigwedge_{k=1}^{n_v^\kappa} [a_{v_k}^\kappa.\text{name} \notin \mathcal{N} \wedge a_{v_k}^\kappa.\text{name} \in \mathcal{B}] , \\ P_{\text{bename}} &\equiv \bigwedge_{k=1}^{n_e} [b_{e_k}.\text{name} \notin \mathcal{N} \wedge b_{e_k}.\text{name} \in \mathcal{B}] , \\ P_{\text{uvname}} &\equiv \delta.\text{name} \notin \mathcal{N} \wedge \delta.\text{name} \in \mathcal{B} \wedge \bigwedge_{\kappa=1}^3 [\mu_\kappa.\text{name} \notin \mathcal{N} \wedge \mu_\kappa.\text{name} \in \mathcal{B}] . \end{aligned} \quad (6.194)$$

Die durch S_{dekl} zu erfüllenden Aufgaben können in der zusätzlichen Nachbedingung

$$P'_{\text{dekl}} \equiv P_{\text{dec}be} \wedge P_{\text{dec}av} \wedge P_{\text{dec}uv} \quad (6.195)$$

zusammengefasst werden. Zur Erfüllung der Bedingungen werden wir den (hier kommentierten) C-Code

Bedingung	Anweisung	Kommentar
$\{V\} \equiv$	$P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec}bq} \wedge P_{\text{dec}bv} \wedge$ $P_{\text{deca}FB} \wedge P_{\text{uvname}} \wedge P_{\text{bename}} \wedge P_{\text{avname}}$ <code>int delta;</code> <code>int mu1;</code> <code>int mu2;</code> <code>int mu3;</code>	Deklarationen der unabhängigen Variablen
$\{V_{\text{deklbe}}\} \equiv$	$P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec}bq} \wedge P_{\text{dec}bv} \wedge$ $P_{\text{deca}FB} \wedge P_{\text{dec}uv} \wedge P_{\text{bename}} \wedge P_{\text{avname}}$ <code>float b_e_1[P];</code> <code>:</code> <code>float b_e_n_e[P];</code>	Variablendeklarationen für den Vektor \mathbf{b}_e
$\{P_{\text{deklbe}}\} \equiv$	$P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec}bq} \wedge P_{\text{dec}bv} \wedge$ $P_{\text{deca}FB} \wedge P_{\text{dec}uv} \wedge P_{\text{dec}be} \wedge P_{\text{avname}}$ <code>float a_v_1_1[P];</code> <code>:</code> <code>float a_v_1_n_v^1[P];</code> <code>float a_v_2_1[P];</code> <code>:</code> <code>float a_v_2_n_v^2[P];</code> <code>float a_v_3_1[P];</code> <code>:</code>	Variablendeklarationen für den Vektor \mathbf{a}_v^1 Variablendeklarationen für den Vektor \mathbf{a}_v^2 Variablendeklarationen für

<code>float a_v_3_n_v³[P];</code>	den Vektor \mathbf{a}_v^3
<code>float a_v_4_1[P];</code>	Variablendeklarationen
<code>⋮</code>	für
<code>float a_v_4_n_v⁴[P];</code>	den Vektor \mathbf{a}_v^4
<code>float a_v_5_1[P];</code>	Variablendeklarationen
<code>⋮</code>	für
<code>float a_v_5_n_v⁵[P];</code>	den Vektor \mathbf{a}_v^5
<code>float a_v_6_1[P];</code>	Variablendeklarationen
<code>⋮</code>	für
<code>float a_v_6_n_v⁶[P];</code>	den Vektor \mathbf{a}_v^6
$\{P_{\text{dekl}}\} \equiv P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec bq}} \wedge P_{\text{dec bv}} \wedge P_{\text{deca FB}} \wedge P_{\text{dec uv}} \wedge P_{\text{dec be}} \wedge P_{\text{dec av}}$	

verwenden. Im Folgenden werden wir die Korrektheit des Codes nachweisen. Anwendung von Axiom (5.15) auf die Deklarationsanweisungen für \mathbf{a}_v liefert die Vorbedingung von `float a_v_1_1[P]`;

$$P_{\text{dekl}} \left[\begin{smallmatrix} a_{v1}^1 \text{.typ} \\ \text{float} \end{smallmatrix} \right] \cdots \left[\begin{smallmatrix} a_{v n_v^6}^6 \text{.typ} \\ \text{float} \end{smallmatrix} \right] \left[\begin{smallmatrix} a_{v1}^1 \text{.l} \\ P \end{smallmatrix} \right] \cdots \left[\begin{smallmatrix} a_{v n_v^6}^6 \text{.l} \\ P \end{smallmatrix} \right] \wedge \text{float} \in \mathcal{D} \wedge P_{\text{avname}} \wedge P > 0. \quad (6.196)$$

Die Ersetzungen, die an P_{dekl} vorzunehmen sind, betreffen nur P_{decav} . Der Datentyp `float` existiert, so dass $\text{float} \in \mathcal{D}$ eine wahre Aussage ist. Weiterhin berücksichtigen wir $P_{\text{dekl}} \Rightarrow P > 0$. Daher wird aus Gleichung (6.196)

$$P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec bq}} \wedge P_{\text{dec bv}} \wedge P_{\text{deca FB}} \wedge P_{\text{dec uv}} \wedge P_{\text{dec be}} \wedge P_{\text{avname}} \wedge \bigwedge_{\kappa=1}^6 \bigwedge_{l=1}^{n_v^\kappa} [\text{float} = \text{float} \wedge P \leq P]. \quad (6.197)$$

Die wahren Aussagen entfallen aufgrund des Prädikatenkalküls, so dass Gleichung (6.197), wie behauptet, P_{deklbe} darstellt.

Die Vorbedingung zur Anweisung `float b_e_1[P]`; ergibt sich durch Anwendung von Axiom (5.15) zu

$$P_{\text{deklbe}} \left[\begin{smallmatrix} b_{e1} \text{.typ} \\ \text{float} \end{smallmatrix} \right] \cdots \left[\begin{smallmatrix} b_{e n_e} \text{.typ} \\ \text{float} \end{smallmatrix} \right] \left[\begin{smallmatrix} b_{e1} \text{.l} \\ P \end{smallmatrix} \right] \cdots \left[\begin{smallmatrix} b_{e n_e} \text{.l} \\ P \end{smallmatrix} \right] \wedge \text{float} \in \mathcal{D} \wedge P_{\text{bename}} \wedge P > 0. \quad (6.198)$$

Mit gleicher Argumentation wie zuvor erhalten wir aus Gleichung (6.198)

$$P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{dec bq}} \wedge P_{\text{dec bv}} \wedge P_{\text{deca FB}} \wedge P_{\text{dec uv}} \wedge \bigwedge_{l=1}^{n_e} [\text{float} = \text{float} \wedge P \leq P \leq \text{Max}_i] \wedge P_{\text{avname}} \wedge P_{\text{bename}}. \quad (6.199)$$

Die wahren Aussagen entfallen wiederum auf Grund des Prädikatenkalküls, und es verbleibt von Gleichung (6.199), wie behauptet, V_{declbe} .

Die eigentliche Vorbedingung V des Algorithmus bestimmt sich durch Anwendung von Axiom (5.12) zu

$$V_{\text{declbe}} \left[\begin{smallmatrix} \delta \text{.typ} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_1 \text{.typ} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_2 \text{.typ} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_3 \text{.typ} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \delta \text{.l} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_1 \text{.l} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_2 \text{.l} \\ \text{int} \end{smallmatrix} \right] \left[\begin{smallmatrix} \mu_3 \text{.l} \\ \text{int} \end{smallmatrix} \right] \wedge \text{int} \in \mathcal{D} \wedge P_{\text{uvname}} \wedge P > 0. \quad (6.200)$$

Der Datentyp `int` existiert, so dass $\text{int} \in \mathcal{D}$ eine wahre Aussage ist. Die Ersetzungen, die an V_{deklbe} vorzunehmen sind, betreffen nur P_{decuv} , so dass Gleichung (6.200) ausführlich

$$V \equiv P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{decqb}} \wedge P_{\text{decbv}} \wedge P_{\text{decaFB}} \wedge P_{\text{uvname}} \wedge P_{\text{bename}} \wedge P_{\text{avname}} \wedge \text{int} = \text{int} \wedge -1 = -1 \quad (6.201)$$

ist. Die wahren Aussagen entfallen wiederum auf Grund des Prädikatenkalküls und liefern die Vorbedingung für den Algorithmus, wie behauptet, zu

$$V \equiv P_{\text{WDF}} \wedge P_q \wedge P_v \wedge P_{\text{decqb}} \wedge P_{\text{decbv}} \wedge P_{\text{decaFB}} \wedge P_{\text{uvname}} \wedge P_{\text{bename}} \wedge P_{\text{avname}} . \quad (6.202)$$

6.8 Zusammenfassung

Wir wollen nun die wesentlichen Ergebnisse dieses Kapitels zusammenfassen und bewerten. Es wurde ein Algorithmus durch Verknüpfen von Elementaranweisungen der Programmiersprache C angegeben, der ein lineares, konstantes MDWDF für ein quaderförmiges Berechnungsgebiet simuliert. Die Anweisungssequenz, die eine Abtastschicht berechnet, ist dabei so geartet, dass sie einen, gemäß Kapitel 3.5 spezifizierten, Funktionsbaustein des Programmpaketes SPACE darstellt. Weiterhin erfolgte die Niederlegung der Anforderungen mittels einer formalen Spezifikation. Durch konsequente Anwendung der aus der Literatur bekannten und im Kapitel 5 dargelegten Ergebnisse wurde aus der Nachbedingung und dem Algorithmus eine geeignete Vorbedingung ermittelt. Die Tatsache, dass auch das Nicht-Überschreiten des Rechenbereiches durch eine geeignete Vorbedingung sichergestellt werden kann, hebt dieses Verfahren gegenüber anderen hervor. Mit diesem formalen Korrektheitsbeweis ist nunmehr das Gebiet der numerischen Integration partieller Differentialgleichungen auch für sicherheitskritische Anwendungen erschlossen worden.

Kapitel 7

Zusammenfassung

Die vorliegende Arbeit befasst sich vornehmlich mit zwei bisher nicht behandelten Aufgabenstellungen aus dem Bereich der mehrdimensionalen Wellendigitalfilter. Zum einen mit der systematischen Erzeugung von Simulations-Algorithmen, ausgehend von einem System linearer, konstanter partieller Differentialgleichungen. Den zweiten bildet der Nachweis der Korrektheit des verwendeten Programmcodes, der die Simulations-Algorithmen implementiert.

Im Kapitel 2 wurden zunächst die wichtigsten für die vorliegende Arbeit notwendigen Grundlagen der mehrdimensionalen Wellendigitalfilter wiederholt. Dabei wurde auf eine Darstellungsform Wert gelegt, die den Aufgaben der folgenden Kapitel Rechnung trägt.

Im Kapitel 3 wurde das Programmpaket SPACE beschrieben und die Einbindung von mehrdimensionalen Wellendigitalfiltern erörtert. Dabei zeigte sich, dass die vorhandenen Schnittstellen der einzelnen Softwaremodule, des zur Behandlung eindimensionaler Probleme vorgesehenen Programmpaketes SPACE, ungeeignet für unser Ziel sind. Dem Problem sind wir durch Erweiterung der Signale des Programmpaketes SPACE auf vektorielle Signale (Felder) entgegengetreten. Diese erforderliche Abänderung der vorhandenen Schnittstelle und die daraus resultierende Möglichkeit der Einbindung der mehrdimensionalen Wellendigitalfilter sind ebenfalls im Kapitel 3 niedergelegt.

Im Kapitel 4 wurde zunächst durch Passivitätsbetrachtungen die Klasse der PDGLn eingeschränkt, die behandelt werden. Für diese Klasse PDGLn wurde ein Verfahren zur automatischen Synthese von Referenzschaltung und Wellendigitalfilter entwickelt, welches sich somit auch zur automatischen Codeerzeugung für eine Simulation des zu Grunde liegenden Systems mittels mehrdimensionaler Wellendigitalfilter eignet. Aspekte wie die Genauigkeit des eingesetzten numerischen Integrationsverfahrens und die Effizienz der Implementierung rückten hierbei in den Hintergrund. Vielmehr war es das Ziel ein numerisches Integrationsverfahren, welches die Möglichkeit der formalen Verifikation bietet, zu erhalten.

Die aus der Literatur bekannten, anerkannten und für die vorliegende Arbeit notwendigen Grundlagen der Software-Verifikation mittels formaler Methoden wurden im Kapitel 5 rekapituliert. Die Definitionen der notwendigen Programmanweisungen finden sich ebenfalls in diesem Kapitel.

Im Kapitel 6 wurde aus den zu Grunde liegenden Anforderungen eine formale Spezifikation entwickelt. Im Anschluss daran wurde ein Programmcode zur Simulation des mehrdimensionalen Wellendigitalfilters festgelegt. Dabei wurde der Entwurf des Codes so durchgeführt, dass dieser sich in das Programmpaket SPACE einbinden lässt. Im Anschluss daran erfolgte für diesen vorliegenden Programmcode der Nachweis der Korrektheit mittels formaler Methoden. Es zeigt sich, dass auch bei anspruchsvolleren Problemstellungen der digitalen Signalverarbeitung eine Beweisführung der Korrektheit relativ einfach durchzuführen ist. Im Zusammenhang mit der Beschränktheit der verwendeten Signale hat die Passivität des eingangs betrachteten Systems einen wesentlichen Anteil daran. Viele der für dieses Kapitel notwendigen Berechnungen und Abschätzungen in Bezug auf die Beschränktheit der verwendeten Signale, insbesondere im Falle von endlicher Signalwortlänge, wurden in den Anhang ausgelagert.

Durch die vorliegende Arbeit ergeben sich i. W. zwei neue Ergebnisse zur numerischen Integration von

PDGLn mithilfe mehrdimensionaler Wellendigitalfilter. Zum einen kommt der effizienten Codeerzeugung in der Praxis eine erhebliche Bedeutung zu, der durch das entwickelte automatische Syntheseverfahren ermöglicht wird. Auf Basis dieses Verfahrens wird zur Zeit am Lehrstuhl für Nachrichtentechnik ein Werkzeug zur Erzeugung von C-Code entwickelt [BV03]. Nach Vorliegen dieses Simulations-Werkzeugs wird sich der Anwenderkreis des numerischen Integrationsverfahrens nach der Wellendigitalmethode erheblich erweitern. Zum anderen eröffnet der geführte formale Korrektheitsbeweis die Möglichkeit der Anwendung des Integrationsverfahrens in sicherheitskritischen Systemen.

In weiterführenden Untersuchungen könnte das vorgestellte Syntheseverfahren auf quasilineare, symmetrisch hyperbolische PDGL erweitert werden. Hierbei wird man zunächst nicht den allgemeinen Fall behandeln, sondern Einschränkungen an die PDGL machen, insbesondere in Bezug auf die Abhängigkeiten der Koeffizienten von der Zeit. Ferner könnte geklärt werden, ob die Führung eines Korrektheitsbeweises für ein Werkzeug zur automatischen Codeerzeugung möglich ist. Weiterhin wäre eine Skalierung des Algorithmus, in Teilalgorithmen entsprechend der Teilberechnungsgebiete inklusive Kommunikation zwischen den Rechnern sinnvoll. In diesem Zusammenhang empfehlen sich auch Untersuchungen zur Behandlung von beliebig gearteten Berechnungsgebieten mit dem vorliegenden Verfahren. Die automatische Codeerzeugung für den nichtlinearen Fall stellt eine weitere interessante Aufgabenstellung dar. Als besondere Herausforderung dürfte sich in diesem Fall der Nachweis der Korrektheit der zu implementierenden Algorithmen erweisen, die der Lösung der algebraischen, nichtlinearen Gleichungen dienen.

Anhang A

Die Operatoren \min und \max

Die folgenden Sätze erscheinen trivial, ermöglichen aber ermüdungsfreies Arbeiten. α, β sollen reelle Zahlen sein und x, x_μ, Max_f sollen positiv sein.

$$\min \left\{ x, \frac{1}{x} \right\} = \frac{1}{\max \left\{ x, \frac{1}{x} \right\}} \quad (\text{A.1})$$

$$\min_{\mu} \left\{ \frac{1}{x_{\mu}} \right\} = \frac{1}{\max_{\mu} \{x_{\mu}\}} \quad (\text{A.2})$$

$$\min_{\mu} \left\{ \min \left\{ x_{\mu}, \frac{1}{x_{\mu}} \right\} \right\} = \frac{1}{\max_{\mu} \left\{ \max \left\{ x_{\mu}, \frac{1}{x_{\mu}} \right\} \right\}} \quad (\text{A.3})$$

$$\max \left\{ x, \frac{1}{x} \right\} \geq 1, \quad \min \left\{ x, \frac{1}{x} \right\} \leq 1 \quad (\text{A.4})$$

$$\begin{aligned} [\quad \alpha \cdot x < \text{Max}_f \wedge 1 \leq \alpha \quad] &\implies [\quad x < \text{Max}_f \quad] \\ [\quad x < \beta \cdot \text{Max}_f \wedge \beta \leq 1 \quad] &\implies [\quad x < \text{Max}_f \quad] \end{aligned} \quad (\text{A.5})$$

Anhang B

Unitär beschränkte Matrizen

Wir wollen eine $m \times n$ Matrix \mathbf{S} als unitär beschränkt (oder subunitär) bezeichnen, wenn sie

$$\|\mathbf{S}\mathbf{a}\|^2 \leq \|\mathbf{a}\|^2 \quad \forall \mathbf{a} \in \mathbb{C}^n \quad (\text{B.1})$$

genügt. Es ergeben sich für die hermitesch konjugierte Matrix \mathbf{S}^H und die um die Spalte ν reduzierte Matrix ${}^{-\nu}\mathbf{S}$ die folgenden Eigenschaften

$$\begin{array}{ccccccc}
 \boxed{\|\mathbf{S}\mathbf{a}\|^2 \leq \|\mathbf{a}\|^2 \quad \forall \mathbf{a}} & \Leftrightarrow & \boxed{\mathbf{D} = \mathbf{1} - \mathbf{S}^H \mathbf{S} \geq 0} & \Leftrightarrow & \boxed{\tilde{\mathbf{D}} = \mathbf{1} - \mathbf{S} \mathbf{S}^H \geq 0} & \Leftrightarrow & \boxed{\|\mathbf{S}^H \mathbf{b}\|^2 \leq \|\mathbf{b}\|^2 \quad \forall \mathbf{b}} \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 \boxed{\|{}^{-\nu}\mathbf{S} {}^{-\nu}\mathbf{a}\|^2 \leq \|{}^{-\nu}\mathbf{a}\|^2 \quad \forall {}^{-\nu}\mathbf{a}} & \Leftrightarrow & \boxed{{}^{-\nu}\mathbf{D} \geq 0} & & \boxed{{}^{-\nu}\tilde{\mathbf{D}} \geq 0} & \Leftrightarrow & \boxed{\|{}^{-\nu}\mathbf{S}^H {}^{-\nu}\mathbf{b}\|^2 \leq \|{}^{-\nu}\mathbf{b}\|^2 \quad \forall {}^{-\nu}\mathbf{b}} .
 \end{array} \quad (\text{B.2})$$

Anhang C

Beschränktheit der Wellengrößen

Im Kapitel 6 benötigten wir einige Sätze im Zusammenhang mit der Beschränktheit der Programmvariablen. Diese Sätze werden in diesem Anhang nachgeliefert. Um eine übersichtliche Darstellung zu gewähren, ersetzen wir die Programmvariablen durch ihre PL-Variablen. Das Postfix $.v$ entfällt somit regelmäßig.

Ziel dieses Kapitels ist es, unter der Voraussetzung der Beschränktheit von $\|\mathbf{b}_q\|^2 + \|\mathbf{b}_v\|^2$, die Beschränktheit der Wellen $b_{e\mu}$ und der zur Berechnung von $b_{e\mu}$ benutzten Zwischengrößen nachzuweisen. Zunächst wird eine Relation zwischen den oben angegebenen Größen für die Leistungswellen hergeleitet, um daraus eine Beziehung für die Spannungswellen zu bestimmen. Der Beweis ist dem in [Rumm98] für unitäre Matrizen \mathbf{L}' geführten ähnlich, wobei wir im Folgenden \mathbf{L} anstatt \mathbf{L}' schreiben wollen. In unserem Fall ist \mathbf{L} eine unitär beschränkte Matrix, was nun gezeigt wird. Dazu bereiten wir zunächst

$$\mathbf{L} \mathbf{L}^H = \mathbf{P}_b \mathbf{S} \begin{bmatrix} \mathbf{P}_{eq} & \mathbf{P}_{ev} & (\mathbf{P}_{ee} \mathbf{P}_b^T) \end{bmatrix} \begin{bmatrix} \mathbf{P}_{eq}^T \\ \mathbf{P}_{ev}^T \\ \mathbf{P}_b \mathbf{P}_{ee}^T \end{bmatrix} \mathbf{S}^H \mathbf{P}_b^T \quad (\text{C.1})$$

auf. Mit $\mathbf{P}_b^T \mathbf{P}_b = \mathbf{1}$ und $\mathbf{P}_e \mathbf{P}_e^T = \mathbf{1}$ erhalten wir

$$\mathbf{L} \mathbf{L}^H = \mathbf{P}_b \mathbf{S} \mathbf{S}^H \mathbf{P}_b^T. \quad (\text{C.2})$$

Aus der Passivität folgt die unitäre Beschränktheit von \mathbf{S} . Wegen Gleichung (B.2) ist auch \mathbf{S}^H unitär beschränkt, d. h. es gilt

$$\mathbf{x}^H [\mathbf{1} - \mathbf{S} \mathbf{S}^H] \mathbf{x} \geq 0 \quad \forall \mathbf{x}. \quad (\text{C.3})$$

Aufgrund der Unitarität von \mathbf{P}_b ist dies äquivalent zu

$$\mathbf{x}^H \mathbf{P}_b [\mathbf{1} - \mathbf{S} \mathbf{S}^H] \mathbf{P}_b^T \mathbf{x} = \mathbf{x}^H [\mathbf{1} - \mathbf{P}_b \mathbf{S} \mathbf{S}^H \mathbf{P}_b^T] \mathbf{x} \geq 0 \quad \forall \mathbf{x}. \quad (\text{C.4})$$

Nutzen von Gleichung (C.2) liefert zunächst

$$\mathbf{x}^H [\mathbf{1} - \mathbf{L} \mathbf{L}^H] \mathbf{x} \geq 0 \quad \forall \mathbf{x} \quad (\text{C.5})$$

und bestätigt mit Gleichung (B.2) die behauptete Aussage.

Die Koordinate μ aus dem Spaltenvektor $\mathbf{b}_e = \mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v + \mathbf{L}_e \mathbf{b}_e$ berechnet sich zu

$$b_{e\mu} = l_{q\mu, \bullet} \mathbf{b}_q + l_{v\mu, \bullet} \mathbf{b}_v + \sum_{\nu=1}^{n_e} l_{e\mu\nu} b_{e\nu}. \quad (\text{C.6})$$

Wir definieren nun den Vektor

$${}^m\mathbf{b}_e = \begin{bmatrix} {}^mb_{e1} \\ \vdots \\ {}^mb_{en_e} \end{bmatrix} \quad (\text{C.7})$$

dessen Koordinate μ aus den ersten $n_q + n_v + m$ Summanden von $b_{e\mu}$ bestehen soll. Dabei unterscheiden wir für m drei Fälle. Zunächst soll sich die Summation nur auf die gewichteten Quellenwellen beschränken, d. h.

$${}^mb_{e\mu} = \sum_{\nu=1}^{m+n_q+n_v} l_{q\mu\nu} b_{q\nu}, \quad (\text{C.8})$$

wobei hier $m = 1 - n_q - n_v, 2 - n_q - n_v, \dots, -1 - n_v, -n_v$ gilt. Da die Matrix \mathbf{L} unitär beschränkt ist, gilt dies auch für deren Untermatrizen, d. h.

$$|{}^mb_{e\mu}|^2 = \left| \sum_{\nu=1}^{m+n_q+n_v} l_{q\mu\nu} b_{q\nu} \right|^2 \leq \sum_{\nu=1}^{m+n_q+n_v} |b_{q\nu}|^2. \quad (\text{C.9})$$

Berücksichtigen wir nun

$$\sum_{\nu=1}^{m+n_q+n_v} |b_{q\nu}|^2 \leq \|\mathbf{b}_q\|^2 + \|\mathbf{b}_v\|^2 < \frac{\text{Max}_f^2}{4\beta^2}, \quad (\text{C.10})$$

so gilt tatsächlich

$$|{}^mb_{e\mu}|^2 < \frac{\text{Max}_f^2}{4\beta^2}. \quad (\text{C.11})$$

Im zweiten Fall erfolgt die Summation von allen gewichteten Quellenwellen und über gewichtete Verzögererwellen, d. h.

$${}^mb_{e\mu} = l_{q\mu,\bullet}\mathbf{b}_q + \sum_{\nu=1}^{m+n_v} l_{v\mu\nu} b_{v\nu}, \quad (\text{C.12})$$

wobei hier $m = 1 - n_v, 2 - n_v, \dots, -1, 0$ gilt. Da die Matrix \mathbf{L} unitär beschränkt ist, gilt dies auch für deren Untermatrizen, d. h.

$$|{}^mb_{e\mu}|^2 = |l_{q\mu,\bullet}\mathbf{b}_q + \sum_{\nu=1}^{m+n_v} l_{v\mu\nu} b_{v\nu}|^2 \leq \|\mathbf{b}_q\|^2 + \sum_{\nu=1}^{m+n_v} |b_{v\nu}|^2. \quad (\text{C.13})$$

Berücksichtigen wir nun

$$\|\mathbf{b}_q\|^2 + \sum_{\nu=1}^{m+n_v} |b_{v\nu}|^2 \leq \|\mathbf{b}_q\|^2 + \|\mathbf{b}_v\|^2 < \frac{\text{Max}_f^2}{4\beta^2}, \quad (\text{C.14})$$

so gilt tatsächlich auch im zweiten Fall

$$|{}^mb_{e\mu}|^2 < \frac{\text{Max}_f^2}{4\beta^2}. \quad (\text{C.15})$$

Der dritte Fall ist der komplizierteste. Für ihn gilt

$${}^m b_{e\mu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \sum_{\nu=1}^m l_{e\mu\nu} b_{e\nu}, \quad (\text{C.16})$$

wobei m hier die Werte $1, 2, \dots, n_e$ annehmen darf.

Wegen der strikten unteren Dreiecksgestalt von \mathbf{L}_e gilt $l_{e\mu\nu} = 0$ für $\nu \geq \mu$. Nutzen wir dies für die Berechnung von $b_{e\mu}$ aus, so lautet Gleichung (C.6)

$$b_{e\mu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \sum_{\nu=1}^{\mu-1} l_{e\mu\nu} b_{e\nu}. \quad (\text{C.17})$$

Sinnvoll ist es hier, zwei Fälle zu unterscheiden

$$b_{e\mu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \sum_{\nu=1}^{\mu-1} l_{e\mu\nu} b_{e\nu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \begin{cases} \sum_{\nu=1}^{\mu-1} l_{e\mu\nu} b_{e\nu} & \text{für } m \geq \mu - 1 \\ \sum_{\nu=1}^m l_{e\mu\nu} b_{e\nu} + \sum_{\nu=m+1}^{\mu-1} l_{e\mu\nu} b_{e\nu} & \text{für } m < \mu - 1. \end{cases} \quad (\text{C.18})$$

Für die Teilsummen gilt entsprechend

$${}^m b_{e\mu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \begin{cases} \sum_{\nu=1}^{\mu-1} l_{e\mu\nu} b_{e\nu} & \text{für } m \geq \mu - 1 \\ \sum_{\nu=1}^m l_{e\mu\nu} b_{e\nu} & \text{für } m < \mu - 1 \text{ (und für } m \geq \mu - 1 \text{).} \end{cases} \quad (\text{C.19})$$

Die Differenz aus Gleichung (C.18) und Gleichung (C.19) lautet

$$b_{e\mu} - {}^m b_{e\mu} = \begin{cases} 0 & \text{für } m \geq \mu - 1 \\ \sum_{\nu=m+1}^{\mu-1} l_{e\mu\nu} b_{e\nu} & \text{für } m < \mu - 1. \end{cases} \quad (\text{C.20})$$

Offenbar stimmen die ersten $m+1$ Koordinaten der Vektoren \mathbf{b}_e und ${}^m \mathbf{b}_e$ überein, d. h.

$$\begin{bmatrix} b_{e1} \\ \vdots \\ b_{em+1} \end{bmatrix} = \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{em+1} \end{bmatrix} \quad \text{bzw.} \quad [\mathbf{1}_{m+1} \quad \mathbf{0}] \mathbf{b}_e = [\mathbf{1}_{m+1} \quad \mathbf{0}] {}^m \mathbf{b}_e. \quad (\text{C.21})$$

Umbenennung der Indizes von μ nach ν liefert ${}^m b_{e\nu} = b_{e\nu}$ für $\nu \leq m+1$. Somit können wir in Gleichung (C.16) jede Welle $b_{e\nu}$ durch ${}^m b_{e\nu}$ ersetzen d. h.

$${}^m b_{e\mu} = \mathbf{l}_{q\mu, \bullet} \mathbf{b}_q + \mathbf{l}_{v\mu, \bullet} \mathbf{b}_v + \sum_{\nu=1}^m l_{e\mu\nu} {}^m b_{e\nu}. \quad (\text{C.22})$$

Zu beachten ist, dass diese Gleichung für alle μ und m gilt, da die Voraussetzung $\nu \leq m+1$ sich nur auf den lokalen Summenindex ν bezog. Auf den ersten Blick mag diese Formel verwundern, da auf der

rechten Gleichungsseite nun auch die Teilsummen auftreten. Dies liegt daran, dass zur Berechnung von ${}^m b_{e\mu}$ die Wellen b_{e1} bis b_{em} benötigt werden. Diese Wellen sind aber mit den Teilsummen ${}^m b_{e1}$ und ${}^m b_{em}$ identisch. Somit hängt ${}^m b_{e\mu}$ auch nur von ${}^m b_{e1}$ bis ${}^m b_{em}$ ab.

Gleichung (C.22) kann auch in Matrixschreibweise ermittelt werden. Ausgehend von Gleichung (C.16) für $\mu = 1, \dots, n_e$

$${}^m \mathbf{b}_e = \mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v + \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{n_e-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{b}_e \quad (\text{C.23})$$

wenden wir die aus Gleichung (C.21) unmittelbar folgende Gleichung $[\mathbf{1}_m \quad \mathbf{0}_m^{n_e-m}] \mathbf{b}_e = [\mathbf{1}_m \quad \mathbf{0}_m^{n_e-m}] {}^m \mathbf{b}_e$ an und erhalten

$${}^m \mathbf{b}_e = \mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v + \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{n_e-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} {}^m \mathbf{b}_e. \quad (\text{C.24})$$

Von dieser Gleichung betrachten wir nun nur die ersten μ Zeilen, d. h.

$$[\mathbf{1}_\mu \quad \mathbf{0}_\mu^{n_e-\mu}] {}^m \mathbf{b}_e = [\mathbf{1}_\mu \quad \mathbf{0}_\mu^{n_e-\mu}] [\mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v] + [\mathbf{1}_\mu \quad \mathbf{0}_\mu^{n_e-\mu}] \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{n_e-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} {}^m \mathbf{b}_e. \quad (\text{C.25})$$

Wir wollen nun zeigen, dass

$$|{}^m b_{e\mu}|^2 \leq \|\mathbf{b}_q\|^2 + \|\mathbf{b}_v\|^2 \quad (\text{C.26})$$

gilt. Dazu betrachten wir zunächst den Fall $m \leq \mu - 1 \Leftrightarrow m + 1 \leq \mu$. Für $m \leq \mu - 1$ gilt offenbar

$$\begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{n_e-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} {}^m \mathbf{b}_e = \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{n_e-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{em} \\ \vdots \\ {}^m b_{e n_e} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{\mu-1-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{em} \\ \vdots \\ {}^m b_{e \mu-1} \end{bmatrix}, \quad (\text{C.27})$$

so dass aus Gleichung (C.25)

$$\begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu} \end{bmatrix} = [\mathbf{1}_\mu \quad \mathbf{0}] [\mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v] + [\mathbf{1}_\mu \quad \mathbf{0}] \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{\mu-1-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{em} \\ \vdots \\ {}^m b_{e \mu-1} \end{bmatrix} \quad (\text{C.28})$$

wird. Um diese Gleichung übersichtlicher schreiben zu können, führen wir die Vektoren

$${}^{m,\mu} \mathbf{b}_e = \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu} \end{bmatrix}, \quad {}^{m,\mu-1} \mathbf{b} = \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ {}^{m,\mu-1} \mathbf{b}_e \end{bmatrix} \quad (\text{C.29})$$

und die Matrix ${}^{\mu,m} \mathbf{L}$ als die linke obere Untermatrix von \mathbf{L} mit μ Zeilen und $n_q + n_v + \mu - 1$ Spalten, d. h.

$${}^{\mu,m} \mathbf{L} = [\mathbf{1}_\mu \quad \mathbf{0}] \left[\mathbf{L}_q \quad \mathbf{L}_v \quad \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{\mu-1-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right] \quad (\text{C.30})$$

ein. Damit erhalten wir Gleichung (C.28) zunächst in der Form

$${}^{m,\mu}\mathbf{b}_e = [\mathbf{1}_\mu \quad \mathbf{0}][\mathbf{L}_q \mathbf{b}_q + \mathbf{L}_v \mathbf{b}_v] + [\mathbf{1}_\mu \quad \mathbf{0}]\mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{\mu-1-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} {}^{m,\mu-1}\mathbf{b}_e \quad (\text{C.31})$$

und weiter

$${}^{m,\mu}\mathbf{b}_e = [\mathbf{1}_\mu \quad \mathbf{0}] \left[\mathbf{L}_q \quad \mathbf{L}_v \quad \mathbf{L}_e \begin{bmatrix} \mathbf{1}_m & \mathbf{0}_m^{\mu-1-m} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right] {}^{m,\mu-1}\mathbf{b} = {}^{\mu,m}\mathbf{L} \quad {}^{m,\mu-1}\mathbf{b} . \quad (\text{C.32})$$

Wegen der unitären Beschränktheit von \mathbf{L} ist nach Gleichung (B.2) auch die Teilmatrix ${}^{\mu,m}\mathbf{L}$ unitär beschränkt, d. h.

$$\| {}^{\mu,m}\mathbf{L} \quad {}^{m,\mu-1}\mathbf{b} \|^2 \leq \| {}^{m,\mu-1}\mathbf{b} \|^2 \quad \forall \quad {}^{m,\mu-1}\mathbf{b} . \quad (\text{C.33})$$

Setzen wir ${}^{m,\mu}\mathbf{b}_e = {}^{\mu,m}\mathbf{L} \quad {}^{m,\mu-1}\mathbf{b}$ aus Gleichung (C.32) ein, so haben wir

$$\sum_{\nu=1}^{\mu} |{}^m b_{e\nu}|^2 \leq \| \mathbf{b}_q \|^2 + \| \mathbf{b}_v \|^2 + \sum_{\nu=1}^{\mu-1} |{}^m b_{e\nu}|^2 , \quad (\text{C.34})$$

also das gewünschte Ergebnis

$$|{}^m b_{e\mu}|^2 \leq \| \mathbf{b}_q \|^2 + \| \mathbf{b}_v \|^2 , \quad |{}^m b_{e\mu}|^2 < \frac{\text{Max}_f^2}{4\beta^2} , \quad (\text{C.35})$$

die gesuchte Obergrenze von $|{}^m b_{e\mu}|$, welche für μ, m mit $m+1 \leq \mu$ gültig ist .

Wir zeigen nun, dass sie auch für $m+1 > \mu$ gültig ist. Gleichung (C.22) lautet in Matrixform

$$\begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu} \end{bmatrix} = [\mathbf{1}_\mu \quad \mathbf{0}]\mathbf{L} \begin{bmatrix} \mathbf{1}_{n_q+n_v+\mu-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{b} = [\mathbf{1}_\mu \quad \mathbf{0}]\mathbf{L} \begin{bmatrix} \mathbf{1}_{n_q+n_v+\mu-1} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu-1} \end{bmatrix} . \quad (\text{C.36})$$

Die unitäre Beschränktheit der ersten μ Zeilen und der ersten $n_q + n_v + \mu - 1$ Spalten von \mathbf{L} liefert

$$\left\| \begin{bmatrix} {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu} \end{bmatrix} \right\|^2 = \left\| [\mathbf{1}_\mu \quad \mathbf{0}]\mathbf{L} \begin{bmatrix} \mathbf{1}_{n_q+n_v+\mu-1} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu-1} \end{bmatrix} \right\|^2 \leq \left\| \begin{bmatrix} \mathbf{b}_q \\ \mathbf{b}_v \\ {}^m b_{e1} \\ \vdots \\ {}^m b_{e\mu-1} \end{bmatrix} \right\|^2 , \quad (\text{C.37})$$

was wieder Gleichung (C.35) aber diesmal gültig für $m+1 \leq \mu$, darstellt.

Um dieses Resultat auf die Spannungswellen zu übertragen, nutzen wir die Beziehung $\mathbf{b} = \frac{1}{2} \mathbf{diag} \{ \mathbf{G}_q^{1/2}, \mathbf{G}_v^{1/2}, \mathbf{G}_e^{1/2} \} \mathbf{b}'$ und ersetzen die Leistungswellen in Gleichung (C.35)

$$|{}^m b'_{e\mu}|^2 \leq \| (R_{e\mu} \mathbf{G}_q)^{1/2} \mathbf{b}'_q \|^2 + \| (R_{e\mu} \mathbf{G}_v)^{1/2} \mathbf{b}'_v \|^2 . \quad (\text{C.38})$$

Anhang D

Auswirkungen endlicher Rechengenauigkeit

Anhang C ging von unendlicher Rechengenauigkeit aus (ideale Rechenoperationen). Bekannt ist, dass passive Bauelemente auch bei endlicher Rechengenauigkeit durch Auswahl eines geeigneten Rundungsverfahrens noch passiv sind, [FM75], [Fett78], [Fett90]. In diesem Abschnitt werden praxisgerechte Implementierungen der Arithmetik für verschiedene Zielsysteme vorgestellt. Die zentrale Frage ist, wie die Subunitarität der Streumatrizen nichtreaktiver Bauelemente unter realen Bedingungen gewährleistet werden kann. Untersuchungen zu diesem Thema, bei Benutzung von Festkomma-Arithmetik, finden sich in [Meer79]. In dieser Arbeit steht die Gleitkomma-Arithmetik im Mittelpunkt und wird in diesem Anhang ausschließlich betrachtet. Dazu werden nur WDF-Realisierungen eines (im Idealfall energie-neutralen) Zweitorts untersucht, da jede unitäre Streumatrix auf eine Realisierung zurückgeführt werden kann, die aus eben diesen Zweitort-Streumatrizen besteht, vgl. Kapitel 2.8.1. Die Untersuchungen werden für die gängigen Rundungsarten des Betragsschneidens und des Gitterpunktschneidens durchgeführt. Es werden für beide Rundungsarten Realisierungen angegeben, die unter realen Bedingungen passiv sind.

D.1 Zahlendarstellungen und Rundungsfehler

Darstellung von Gleitkommazahlen

Wir legen für diese Arbeit die Darstellung einer Gleitkommazahl (Floating-Point-Number) nach dem Standard IEEE 754-1985 zu Grunde. Die Menge aller von null verschiedenen Gleitkommazahlen wird durch

$$\tilde{z}_i = (-1)^{s_i} [1 + f_i 2^{-p+1}] 2^{e_i-v}, \quad f_i = \sum_{\nu=1}^{p-1} b_{i\nu} 2^{\nu-1}, \quad b_{i\nu} \in \{0, 1\} \quad (\text{D.1})$$

gebildet. Hierin ist das Zahlensystem durch

- die Anzahl der Bits für den Betrag der Mantisse $p - 1$,
- die Anzahl der Bits für den Exponent q ,
- die Gesamtzahl der Bits $p + q$ und
- der Verschiebung $v = 2^{q-1} - 1$

bestimmt. Für Zahlen einfacher Länge (Datentyp `float` in der Programmiersprache C) gilt $p = 24$ und $q = 8$. Eine konkrete Zahl $\neq 0$ ist dann zum einen bestimmt durch die Mantisse $(-1)^{s_i}[1 + f_i 2^{-p+1}]$ bestehend aus dem Vorzeichenbit $s_i \in \{0, 1\}$ und dem informationstragenden Teil des Betrages (der so genannten Charakteristik) f_i und zum anderen durch den Exponenten 2^{e_i-v} mit $0 < e_i < 2v+1 = 2^q - 1$. Die Null wird durch $e_i = 0 \wedge f_i = 0$ dargestellt. Der Maximalwert von f_i ist

$$f_{\max} = \sum_{\nu=1}^{p-1} 2^{\nu-1} = \sum_{\nu=0}^{p-2} 2^{\nu} = \frac{1 - 2^{p-1}}{1 - 2} = 2^{p-1} - 1. \quad (\text{D.2})$$

f_i ist somit eine ganze Zahl aus dem Intervall $0 \leq f_i \leq f_{\max}$.

Für die folgenden Überlegungen bzgl. der Rundung nehmen wir an, dass keine Überlaufeffekte auftreten. Ist z das Ergebnis einer idealen Rechenoperation und nicht durch eine Gleitkommazahl des verwendeten Zahlensystems darstellbar, so entsteht bei der realen Rechenoperation der Rundungsfehler $z - \tilde{z}$. Gebräuchliche Rundungsarten sind

- Rundung in Richtung 0 (Betragsschneiden), d. h. $\tilde{z} = \text{sgn}(z)[|z|]$, wobei sich hier die Gauß-Klammern auf das Zahlenraster und nicht auf die ganzen Zahlen beziehen,
- Rundung mit minimalem Betrag des Fehlers, d. h. $|z - \tilde{z}| \Rightarrow \min$,
- Rundung in Richtung ∞ , d. h. $\tilde{z} = \lfloor z \rfloor$,
- Rundung in Richtung $-\infty$, d. h. $\tilde{z} = \lceil z \rceil$.

Wir werden nur die ersten beiden betrachten.

Maximaler Fehler bei der Addition und Subtraktion zweier Zahlen

Wir betrachten das Ergebnis der idealen Addition zweier Gleitkommazahlen $z_1 = \tilde{z}_1$ und $z_2 = \tilde{z}_2$

$$z_3 = z_1 + z_2 \quad (\text{D.3})$$

und der realen Addition

$$\tilde{z}_3 = z_1 \oplus z_2. \quad (\text{D.4})$$

Zunächst untersuchen wir das Betragsschneiden. O.E.d.A. nehmen wir $z_1 > z_2$ und $s_1 = s_2 = 0$ an. Der maximale betragliche relative Fehler tritt genau dann auf, wenn z_2 gleich der größten Zahl ist, die keinen Beitrag zu \tilde{z}_3 liefert, d. h. $\tilde{z}_3 = z_1$ und z_3 die kleinstmögliche Mantisse hat. Dann gilt offenbar $z_3 - \tilde{z}_3 = z_1 + z_2 - z_1 = z_2 = \tilde{z}_2$. Die oben gemachten Überlegungen besagen, dass die Mantissen von z_1 durch $f_1 = 0$ und von z_2 durch $f_2 = f_{\max}$ bestimmt sind. Damit z_2 keinen Beitrag zu \tilde{z}_3 liefert, müssen die Exponenten von z_1 und z_2 gerade um p differieren. Bei Festkommazahlen tritt bei der Addition überhaupt kein Fehler auf, da die Exponenten gleich sind. Der maximale Betrag des relativen Fehlers bezogen auf z_3 ist

$$\max \left\{ \left| \frac{z_3 - \tilde{z}_3}{z_3} \right| \right\} = \frac{z_2}{z_3} = \frac{1 + f_{\max} \cdot 2^{-p+1}}{1 + 0 \cdot 2^{-p+1}} \frac{1}{2^p} = \frac{2 - 2^{-p+1}}{2^p} = \frac{1 - 2^{-p}}{2^{p-1}} < \frac{1}{2^{p-1}} = 2^{1-p} \quad (\text{D.5})$$

und bezogen auf \tilde{z}_3

$$\max \left\{ \left| \frac{z_3 - \tilde{z}_3}{\tilde{z}_3} \right| \right\} = \frac{z_2}{\tilde{z}_3} = \frac{z_2}{z_3 - z_2} = \frac{z_2/z_3}{1 - z_2/z_3} = \frac{2 - 2^{-p+1}}{2^p - 2 + 2^{-p+1}} = \frac{1 - 2^{-p}}{2^{p-1} - 1 + 2^{-p}}. \quad (\text{D.6})$$

Beispiel : Sei $p = 4$ und es gelten die folgenden Zahlenwerte

$$\begin{aligned} z_1 &= (1 + 0 \cdot 2^{-1} + 0 \cdot 2^{-2} + 0 \cdot 2^{-3}) \cdot 2^4 = (1 \cdot 2^4 + 0 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1) = 16 \\ z_2 &= (1 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3}) \cdot 2^0 = 1,875 \\ z_3 &= (1 \cdot 2^4 + 0 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3}) \\ \tilde{z}_3 &= (1 \cdot 2^4 + 0 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1) = z_1 . \end{aligned} \quad (\text{D.7})$$

Der relative Fehler ist

$$\frac{z_3 - \tilde{z}_3}{z_3} = \frac{z_2}{z_3} = \frac{1,875}{16} = \frac{1 - 2^{-4}}{2^{4-1}} \approx 0,1172 . \quad (\text{D.8})$$

Zum Vergleich der abgeschätzte relative Fehler

$$\frac{z_3 - \tilde{z}_3}{z_3} < 2^{-3} = 0,125 . \quad (\text{D.9})$$

Im Bild D.1 ist die fehlerhafte Addition dargestellt.

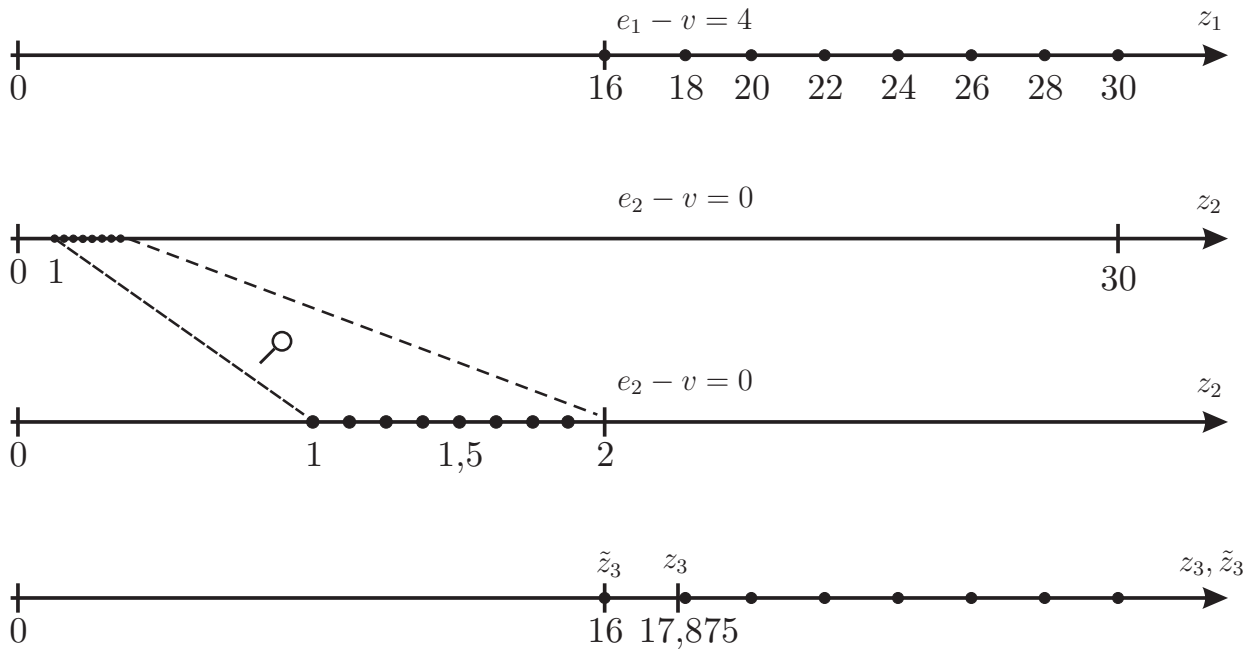


Bild D.1: Zahlenraster zu dem Beispiel der fehlerhaften Addition.

Im Fall der Rundung zum nächsten Gitterpunkt halbiert sich der relative Fehler bezogen auf z_3 zu

$$\max \left\{ \left| \frac{z_3 - \tilde{z}_3}{z_3} \right| \right\} = \frac{1 - 2^{-p}}{2^p} < 2^{-p} . \quad (\text{D.10})$$

Die unterschiedlichen Vorzeichen des Fehlers $z_3 - \tilde{z}_3$ und von z_3 legen eine Ober- und eine Untergrenze fest, welche im Folgenden unter Zuhilfenahme der Abkürzung

$$\varepsilon = \begin{cases} 2^{1-p} & \text{für Betragsschneiden,} \\ 2^{-p} & \text{für Gitterpunktschneiden.} \end{cases} \quad (\text{D.11})$$

ermittelt werden soll. Beide Schneideoperationen bewirken, dass die Vorzeichen von \tilde{z}_3 und z_3 gleich

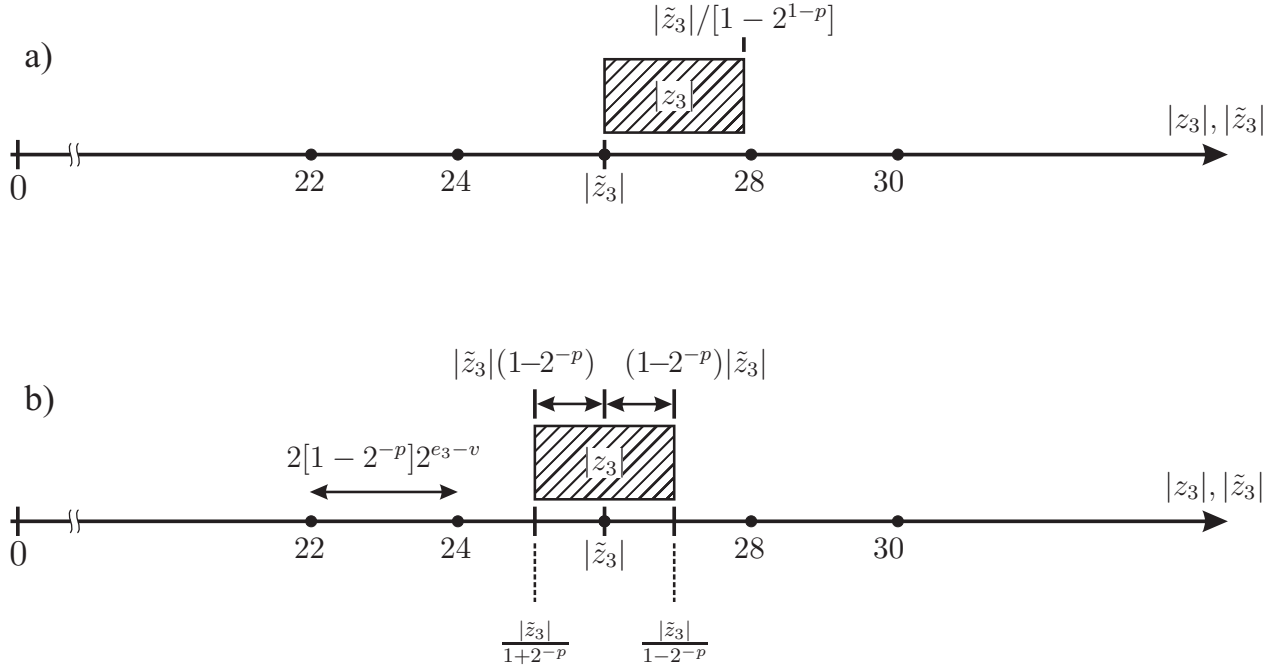


Bild D.2: Fehler bei der Addition a) Betragsschneiden b) Gitterpunktschneiden

sind, d. h. $\tilde{z}_3 z_3 \geq 0$. Dann ergibt sich

$$\frac{|z_3 - \tilde{z}_3|}{|z_3|} < \varepsilon \Leftrightarrow \pm[|z_3| - |\tilde{z}_3|] < \varepsilon |z_3| \Leftrightarrow |z_3| - |\tilde{z}_3| \stackrel{<}{>} \pm \varepsilon |z_3| \Leftrightarrow |z_3|[1 \mp \varepsilon] \stackrel{<}{>} |\tilde{z}_3|. \quad (\text{D.12})$$

Dieses Ergebnis können wir zu

$$[1 - \varepsilon] < \frac{|\tilde{z}_3|}{|z_3|} < [1 + \varepsilon] \iff \frac{1}{1 + \varepsilon} < \frac{|z_3|}{|\tilde{z}_3|} < \frac{1}{1 - \varepsilon} \quad (\text{D.13})$$

zusammenfassen. Beim Betragsschneiden gilt zusätzlich die schärfere Bedingung $\frac{|\tilde{z}_3|}{|z_3|} \leq 1$ und folglich

$$1 - 2^{1-p} < \frac{|\tilde{z}_3|}{|z_3|} \leq 1 \iff 1 \leq \frac{|z_3|}{|\tilde{z}_3|} < \frac{1}{1 - 2^{1-p}}. \quad (\text{D.14})$$

Für das Gitterpunktschneiden hingegen gilt

$$[1 - 2^{-p}] < \frac{|\tilde{z}_3|}{|z_3|} < [1 + 2^{-p}] \iff \frac{1}{1 + 2^{-p}} < \frac{|z_3|}{|\tilde{z}_3|} < \frac{1}{1 - 2^{-p}}. \quad (\text{D.15})$$

Im Bild D.2 sind die Fehler in beiden Fällen dargestellt.

Maximaler Fehler bei der Multiplikation zweier Zahlen

Das Ergebnis der idealen Multiplikation zweier Gleitkommazahlen $z_1 = \tilde{z}_1$ und $z_2 = \tilde{z}_2$ ist

$$z_3 = z_1 \cdot z_2 \quad (\text{D.16})$$

und bei realer Multiplikation

$$\tilde{z}_3 = z_1 \odot z_2. \quad (\text{D.17})$$

Wir betrachten den Fall $f_1 = f_2 = f_{\max}$ (lauter Einsen in der Binärdarstellung der Mantisse), da dies die maximale Mantisse bei z_3 liefert. O.E.d.A. nehmen wir $z_1, z_2 > 0$ an,

$$\begin{aligned}
z_3 &= z_1 z_2 = [1 + f_{\max} 2^{1-p}]^2 2^{e_1+e_2-2v} = [2 - 2^{1-p}]^2 2^{e_1+e_2-2v} = [1 - 2^{-p}]^2 2^{2+e_1+e_2-2v} \\
&= [1 - 2 \cdot 2^{-p} + 2^{-2p}] 2^{2+e_1+e_2-2v} = [1 - 2^{-p+1} + 2^{-2p}] 2^{2+e_1+e_2-2v} \\
&= [2^{1-p} \sum_{\nu=0}^{p-2} 2^\nu + 2^{-2p}] 2^{2+e_1+e_2-2v} = [\sum_{\nu=0}^{p-2} 2^{1+\nu-p} + 2^{-2p}] 2^{2+e_1+e_2-2v} \\
&= [1 + \sum_{\nu=0}^{p-3} 2^{2+\nu-p} + 2^{1-2p}] 2^{1+e_1+e_2-2v} = [1 + 2^{1-2p} \sum_{\nu=1}^{2p-1} b_{3\nu} 2^{\nu-1}] 2^{1+e_1+e_2-2v} \\
&= [1 + 2^{1-p} \sum_{\nu=1}^{2p-1} b_{3\nu} 2^{\nu-1-p}] 2^{1+e_1+e_2-2v} = [1 + 2^{1-p} \sum_{\nu=1}^p b_{3\nu} 2^{\nu-1-p} + 2^{1-p} \sum_{\nu=p+1}^{2p-1} b_{3\nu} 2^{\nu-1-p}] 2^{1+e_1+e_2-2v}.
\end{aligned} \tag{D.18}$$

Es gilt $e_3 - v = 1 + e_1 + e_2 - 2v$. Die Länge der Charakteristik ist offenbar von $p - 1$ auf $2p - 1$ angestiegen. Presst man eine Zahl mit Mantissenlänge $2p$ auf einen Gitterpunkt eines Zahlensystems mit Mantissenlänge p , so ist der Betrag des dadurch entstehenden maximalen Fehlers dann gegeben, wenn $b_{3\nu} = 1$ für $\nu = 1, 2, \dots, p$ gilt. Der maximale Fehler in der Mantisse ist also bei Rundung mittels Betragsschneiden

$$z_3 - \tilde{z}_3 = 2^{1-2p} \sum_{\nu=1}^p 2^{\nu-1} 2^{1+e_1+e_2-2v} = 2^{2-2p+e_1+e_2-2v} \sum_{\nu=0}^{p-1} 2^\nu = 2^{2-p+e_1+e_2-2v} (1 - 2^{-p}). \tag{D.19}$$

Der größte relative Fehler bezogen auf ein z_3 entsteht, wenn z_3 die kleinste Mantisse hat, zu

$$\frac{z_3 - \tilde{z}_3}{z_3} < \frac{2^{2-p+e_1+e_2-2v} (1 - 2^{-p})}{2^{1+e_1+e_2-2v}} = 2^{1-p} (1 - 2^{-p}) < 2^{1-p}. \tag{D.20}$$

Der relative Fehler ist also nicht größer als das letzte Bit der Mantisse von z_3 , d. h. $2^{1-p} 2^{e_3-v}$.

Bei Rundung zum nächsten Gitterpunkt ist hingegen der Fehler nicht größer als die Hälfte des letzten Bits von z_3 , d. h. $2^{-p} 2^{e_3-v}$. Die relativen Fehler bezogen auf z_3 sind dann bei beliebigen Zahlen durch

- bei der Rundung mittels Betragsschneidens durch

$$\left| \frac{z_3 - \tilde{z}_3}{z_3} \right| < 2^{-p+1} \tag{D.21}$$

- und bei Rundung zum nächsten Gitterpunkt durch

$$\left| \frac{z_3 - \tilde{z}_3}{z_3} \right| < 2^{-p}. \tag{D.22}$$

beschränkt. Die maximalen Fehler, die bei der Addition und der Multiplikation entstehen sind zwar nicht exakt gleich, aber die angegebenen oberen Abschätzungen Gleichung (D.5) bzw. Gleichung (D.10) und Gleichung (D.21) bzw. Gleichung (D.22) sind bei beiden Rechenoperationen gleich. Da wir ein z_3 mit kleinster Mantisse angenommen haben, gilt die Abschätzung auch für den relativen Fehler bezogen auf \tilde{z}_3 , d. h.

- bei der Rundung mittels Betragsschneidens

$$\left| \frac{z_3 - \tilde{z}_3}{\tilde{z}_3} \right| < 2^{-p+1} \quad (\text{D.23})$$

- und bei Rundung zum nächsten Gitterpunkt

$$\left| \frac{z_3 - \tilde{z}_3}{\tilde{z}_3} \right| < 2^{-p} . \quad (\text{D.24})$$

Beispiel Erneut gilt $p = 4$ und die Zahlenwerte lauten

$$\begin{aligned} z_1 &= (1 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3}) \cdot 2^3 = 2^3 + 2^1 + 2^0 = 11 \\ z_2 &= (1 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}) \cdot 2^3 = 2^3 + 2^2 + 2^0 = 13 . \end{aligned} \quad (\text{D.25})$$

Die Berechnung liefert

$$\begin{array}{rcccccccc} z_1 \cdot 2^3 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ z_1 \cdot 2^2 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ z_1 \cdot 2^0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ \hline z_3 & 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{array} \quad (\text{D.26})$$

oder in Dezimaldarstellung $z_3 = 143$. Bei Abrundung haben wir $\tilde{z}_3 = 2^7 = 128$ und bei Aufrundung $\tilde{z}_3 = 2^7 + 2^4 = 144$. Die relativen Fehler $15/143 = 0,104895$ bzw. $1/143 = 0,006993$ sind, wie behauptet, durch die Obergrenzen $2^{-p+1} = 0,125$ bzw. $2^{-p} = 0,0625$ begrenzt.

Die unterschiedlichen Vorzeichen des Fehlers von z_3 legen die Ober- und die Untergrenze fest. Bereits ermittelt wurde

$$\left| \frac{z_3 - \tilde{z}_3}{z_3} \right| < \varepsilon \implies \begin{cases} 1 - 2^{1-p} < \frac{|\tilde{z}_3|}{|z_3|} \leq 1 & \text{für Betragsschneiden,} \\ [1 - 2^{-p}] < \frac{|\tilde{z}_3|}{|z_3|} < [1 + 2^{-p}] & \text{für Gitterpunktschneiden.} \end{cases} \quad (\text{D.27})$$

D.2 Idealer und realer Zweitor-Parallel-Adaptor

Wir werden nun Untersuchungen zu einem Zweitor-Parallel-Adaptor durchführen, der energieneutral ist und im realen Fall passiv sein soll. Die Berechnungsgleichungen im idealen Fall sind

$$\mathbf{b} = [\mathbf{\Gamma}\mathbf{\Gamma}^T - \mathbf{1}_2]\mathbf{a} \quad , \quad \mathbf{\Gamma} = [\gamma_1, \gamma_2]^T \quad , \quad \gamma_1, \gamma_2 > 0 \quad , \quad \mathbf{\Gamma}^T\mathbf{\Gamma} = 2 . \quad (\text{D.28})$$

Die Bestimmung von b_1 und b_2 soll durch

$$a_0 = \gamma_1 a_1 + \gamma_2 a_2 \quad , \quad b_1 = \gamma_1 a_0 - a_1 \quad , \quad b_2 = \gamma_2 a_0 - a_2 \quad (\text{D.29})$$

erfolgen. Das Bild D.3 zeigt das Signalfuss-Diagramm des idealen und Bild D.4 das des realen Zweitor-Adaptors. Im realen Fall kann die reale Operation (Multiplikation/Addition) durch die ideale Operation und einem Reformatierungsbaustein dargestellt werden. Die Untersuchungen in Bezug auf den Fehler betreffen 4 Punkte :

1. die ersten zwei Multiplikationen,
2. die erste Addition,

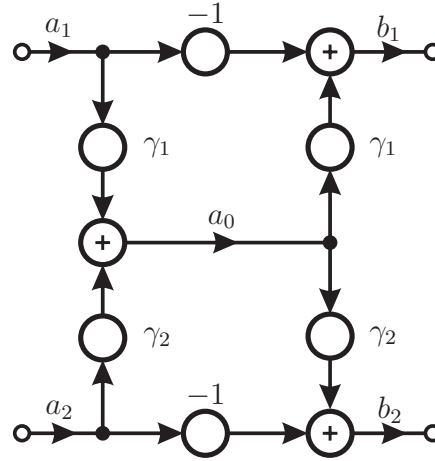


Bild D.3: Signalfluss-Diagramm eines idealen Zweitor-Parallel-Adaptors.

3. die zweiten zwei Multiplikationen und

4. die zweiten zwei Additionen.

Dabei wird jeweils die Untersuchung für die Rundungsart a) Betragsschneiden als auch für b) Gitterpunktschneiden durchgeführt.

Im weiteren Verlauf werden wir die Passivität an manchen Stellen dadurch garantieren, dass anstelle der unitären Matrix \mathbf{S} die unitär beschränkte Matrix $\mathbf{S}' = \mathbf{\Gamma}'_b \mathbf{\Gamma}'_a^T - \mathbf{1}_2$ verwendet wird, wobei sich die Vektoren $\mathbf{\Gamma}'_a$ und $\mathbf{\Gamma}'_b$ durch Kürzung des Vektors $\mathbf{\Gamma}$ ergeben, d. h. es gilt $\mathbf{\Gamma}'_a \|\mathbf{\Gamma}\| = \mathbf{\Gamma} \|\mathbf{\Gamma}'_a\|$, $\mathbf{\Gamma}'_b \|\mathbf{\Gamma}\| = \mathbf{\Gamma} \|\mathbf{\Gamma}'_b\|$, $\|\mathbf{\Gamma}'_a\| \leq \|\mathbf{\Gamma}\|$ und $\|\mathbf{\Gamma}'_b\| \leq \|\mathbf{\Gamma}\|$. Wir müssen nun noch zeigen, dass die Eigenwerte die für die unitäre Beschränktheit notwendigen und hinreichenden Bedingungen $|\lambda'_1| \leq 1 \wedge |\lambda'_2| \leq 1$ erfüllen. Dazu berechnen wir explizit die Eigenwerte

$$\begin{aligned} \det[\mathbf{1}_2 \lambda - \mathbf{S}'] &= \det[\mathbf{1}_2 \lambda - \mathbf{\Gamma}'_b \mathbf{\Gamma}'_a^T + \mathbf{1}_2] = \det[\mathbf{1}_2 (\lambda + 1) - \mathbf{\Gamma}'_b \mathbf{\Gamma}'_a^T] = (\lambda + 1)^2 \det[\mathbf{1}_2 - \mathbf{\Gamma}'_b \mathbf{\Gamma}'_a^T / (\lambda + 1)] \\ &= (\lambda + 1)^2 [1 - \mathbf{\Gamma}'_a^T \mathbf{\Gamma}'_b / (\lambda + 1)] = (\lambda + 1) [\lambda + 1 - \mathbf{\Gamma}'_a^T \mathbf{\Gamma}'_b] = (\lambda + 1) [\lambda + 1 - \|\mathbf{\Gamma}'_a\| \|\mathbf{\Gamma}'_b\|] . \end{aligned} \quad (\text{D.30})$$

Beide Eigenwerte

$$\lambda'_1 = -1 \quad \text{und} \quad \lambda'_2 = \|\mathbf{\Gamma}'_a\| \|\mathbf{\Gamma}'_b\| - 1 \quad (\text{D.31})$$

sind unimodular beschränkt. Für λ'_2 sieht man das wie folgt

$$\|\mathbf{\Gamma}'_a\| \|\mathbf{\Gamma}'_b\| \leq \|\mathbf{\Gamma}\|^2 = 2 \Leftrightarrow \lambda'_2 = \|\mathbf{\Gamma}'_a\| \|\mathbf{\Gamma}'_b\| - 1 \leq 1 . \quad (\text{D.32})$$

Da der Betrag einer Zahl nicht negativ sein kann, ist λ'_2 natürlich auch nicht kleiner als -1 .

1. Die erste Multiplikation

Die reale Multiplikation $\gamma_1 \odot a_1$ entspricht (für feste Faktoren γ_1 und a_1) der idealen Multiplikation von a_1 mit einem modifiziertem $\tilde{\gamma}_1$, d. h.

$$\tilde{z}_3 = \gamma_1 \odot a_1 = \tilde{\gamma}_1 a_1 . \quad (\text{D.33})$$

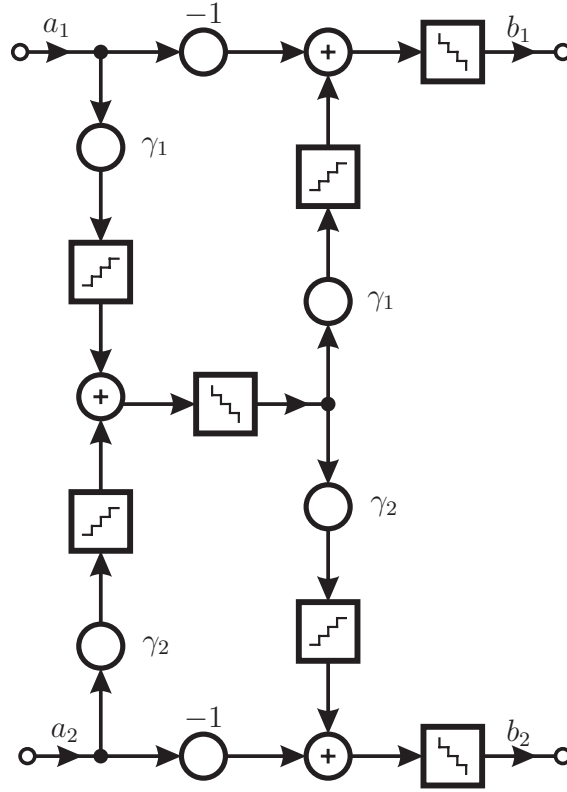


Bild D.4: Signalfluss-Diagramm eines realen Zweiter-Parallel-Adaptors.

Entsprechendes gilt für $\gamma_2 \odot a_2$. In Matrixform lautet dies $\tilde{\mathbf{b}} = [\mathbf{\Gamma} \tilde{\mathbf{\Gamma}}^T - \mathbf{1}_2] \mathbf{a} = \tilde{\mathbf{S}} \mathbf{a}$. Aus Passivitätsgründen muss $\tilde{\mathbf{S}}$ unitär beschränkt sein. Die Passivität ist gemäß Gleichung (D.32) genau dann gegeben, wenn $\|\tilde{\mathbf{\Gamma}}\| \leq \|\mathbf{\Gamma}\|$ gilt. Nehmen wir an, $\gamma_2 \odot a_2$ sei fehlerfrei, so folgt $\tilde{\gamma}_1 \leq \gamma_1$ und somit

$$\tilde{\gamma}_1 |a_1| \leq \gamma_1 |a_1|. \quad (\text{D.34})$$

Da die Multipliziererkoeffizienten positiv sein sollen, können sie in die Betragsoperation gezogen werden.

$$|\tilde{\gamma}_1 a_1| \leq |\gamma_1 a_1|. \quad (\text{D.35})$$

Mit Gleichung (D.33) haben wir

$$|\gamma_1 \odot a_1| \leq |\gamma_1 a_1|. \quad (\text{D.36})$$

Die Passivität fordert also, dass das Ergebnis der realen Multiplikation nicht größer als das der idealen Multiplikation ist.

a) Rundung mittels Betragsschneiden

Bei Verwendung des Betragsschneidens ergibt sich das Ergebnis der realen Multiplikation durch Abrundung des Betrages aus der idealen Rechnung. Ohne weitere Eingriffe ist Gleichung (D.36) erfüllt und die Passivität sichergestellt.

b) Gitterpunktschneiden

Bei Verwendung des Gitterpunktschneidens existieren hingegen Fälle, die Gleichung (D.36) verletzen. Abhilfe schafft für diese Rundungsart ein Stauchen der Multipliziererkoeffizienten. Nach Gleichung (D.27) ist mit $z_3 = \gamma_1 a_1$ und $\tilde{z}_3 = \gamma_1 \odot a_1$ der Fehlerbereich nur durch

$$[1 - 2^{-p}] < \left| \frac{\gamma_1 \odot a_1}{\gamma_1 a_1} \right| < [1 + 2^{-p}] \quad (\text{D.37})$$

eingeschränkt, was Gleichung (D.36) zur oberen Schranke hin, wie eingangs behauptet, verletzen kann. Um dem zu begegnen, wird in der realen Multiplikation der Koeffizient γ_1 durch $\gamma'_1 = \gamma_1/[1+2^{-p}]$ ersetzt. Dann haben wir

$$[1 - 2^{-p}] < \left| \frac{\gamma'_1 \odot a_1}{\gamma'_1 a_1} \right| < [1 + 2^{-p}] . \quad (\text{D.38})$$

Setzen wir $\gamma'_1 = \gamma_1/[1 + 2^{-p}]$ ein, so folgt

$$\left| \frac{[1 + 2^{-p}] \gamma'_1 \odot a_1}{\gamma_1 a_1} \right| < [1 + 2^{-p}] \iff |\gamma'_1 \odot a_1| < |\gamma_1 a_1| . \quad (\text{D.39})$$

Offenbar wird nun Gleichung (D.36) immer erfüllt.

2. Die erste Addition

a) Rundung mittels Betragsschneiden

Wir nehmen nun an, dass die ersten beiden Multiplikationen exakt ausgeführt wurden. Anstatt der idealen Addition $a_0 = \gamma_1 a_1 + \gamma_2 a_2 = \mathbf{\Gamma}^T \mathbf{a}$ haben wir im realen Fall $\tilde{a}_0 = \gamma_1 a_1 \oplus \gamma_2 a_2$. Wir werden zeigen, dass ein Betragsschneiden von \tilde{a}_0 die Passivität des realen Zweitor-Paralleladaptors sicherstellt. Dies ist ein bemerkenswertes Ergebnis (welches übrigens auch für n -Tor-Serien- und Paralleladaptoren gültig ist), da a_0 im Gegensatz zu a_1, a_2, b_1, b_2 nicht mehr als eine Leistungsgröße interpretiert werden kann.

Setzen wir zudem voraus, dass die weiteren Operationen auch ideal durchgeführt werden, so berechnen sich die Ausgangswellen im realen Fall durch

$$\tilde{\mathbf{b}} = \mathbf{\Gamma} \tilde{a}_0 - \mathbf{a} . \quad (\text{D.40})$$

Die „reflektierte Leistung“ des Adaptors ist das Normquadrat von $\tilde{\mathbf{b}}$ und lautet

$$\|\tilde{\mathbf{b}}\|^2 = \|\tilde{a}_0 \mathbf{\Gamma}\|^2 - 2 \tilde{a}_0 \mathbf{\Gamma}^T \mathbf{a} + \|\mathbf{a}\|^2 = 2 \tilde{a}_0 [\tilde{a}_0 - a_0] + \|\mathbf{a}\|^2 \Leftrightarrow \|\tilde{\mathbf{b}}\|^2 - \|\mathbf{a}\|^2 = 2 \tilde{a}_0 [\tilde{a}_0 - a_0] . \quad (\text{D.41})$$

Die Passivitätsbedingung lautet in dem vorliegenden Fall

$$\|\tilde{\mathbf{b}}\|^2 \leq \|\mathbf{a}\|^2 \Leftrightarrow \|\tilde{\mathbf{b}}\|^2 - \|\mathbf{a}\|^2 \leq 0 . \quad (\text{D.42})$$

Offenbar ist $\tilde{a}_0[a_0 - \tilde{a}_0] \geq 0$ notwendig und hinreichend für die Passivität. Diese Beziehung wird durch ein Betragsschneiden von a_0 nach \tilde{a}_0 erfüllt, wie im Folgenden gezeigt wird. Das Betragsschneiden impliziert das Nichtändern des Vorzeichens, d. h. $|\tilde{a}_0| \operatorname{sgn}(a_0) = |\tilde{a}_0| \operatorname{sgn}(\tilde{a}_0)$. Mit $\tilde{a}_0 = |\tilde{a}_0| \operatorname{sgn}(\tilde{a}_0)$ erhalten wir dann

$$\tilde{a}_0[a_0 - \tilde{a}_0] = |\tilde{a}_0|[\operatorname{sgn}(\tilde{a}_0) a_0 - \operatorname{sgn}(\tilde{a}_0) \tilde{a}_0] = |\tilde{a}_0|[\operatorname{sgn}(a_0) a_0 - \operatorname{sgn}(\tilde{a}_0) \tilde{a}_0] = |\tilde{a}_0| [|a_0| - |\tilde{a}_0|] \geq 0 \Leftrightarrow |a_0| \geq |\tilde{a}_0| \quad (\text{D.43})$$

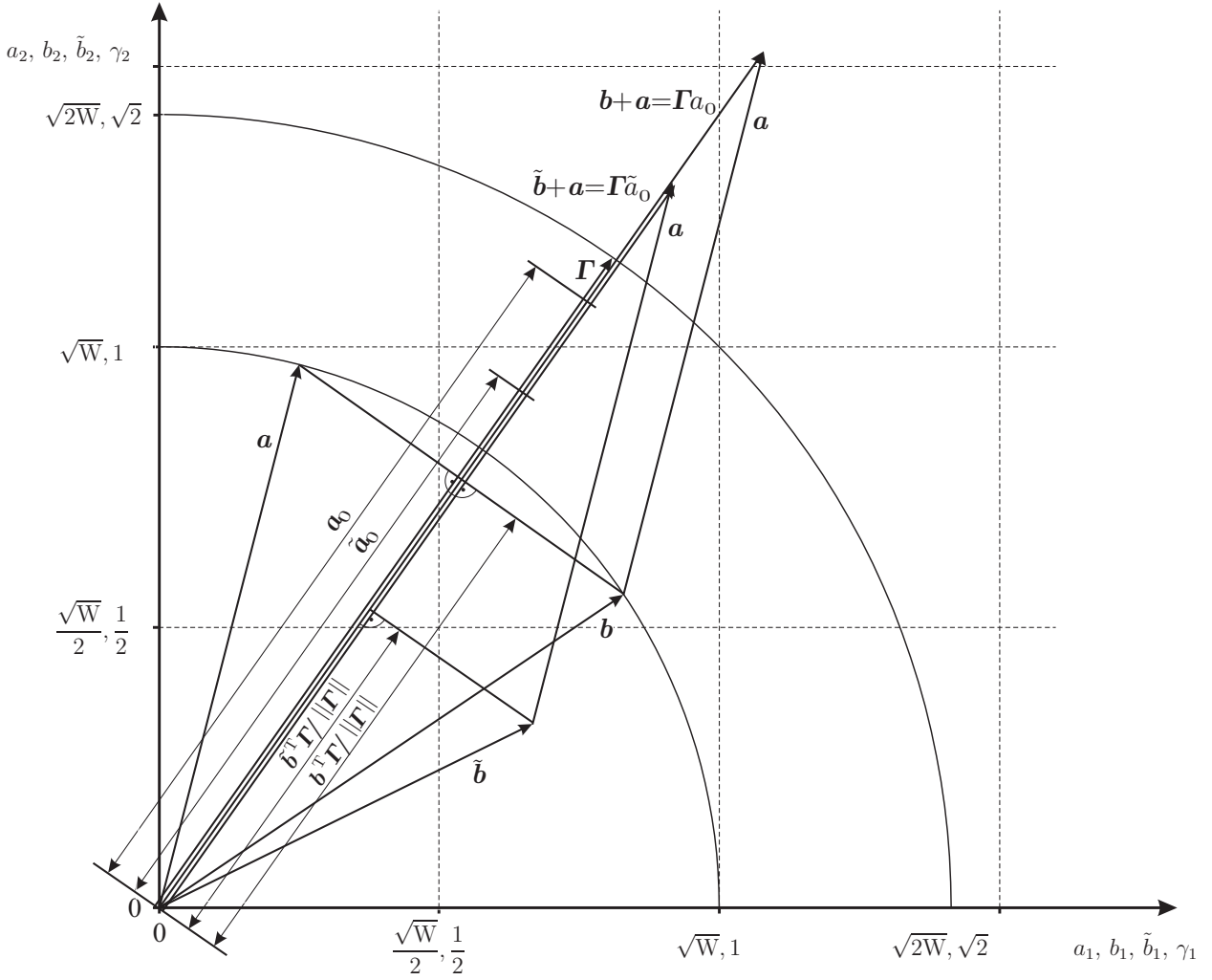


Bild D.5: Wellengrößen eines realen und eines idealen Zweitor-Paralleladaptors

Der Vorgang des Betragsschneidens und der Berechnung von \mathbf{b} und $\tilde{\mathbf{b}}$ ist im Bild D.5 für $\|\mathbf{a}\|^2 = 1W$ gezeichnet.

b) Gitterpunktschneiden

Wie wir im vorherigen Abschnitt festgestellt haben, ist $\tilde{a}_0[a_0 - \tilde{a}_0] \geq 0$ notwendig und hinreichend für die Passivität. Durch das Gitterpunktschneiden kann natürlich die Situation $\tilde{a}_0 > a_0 > 0$ auftreten, die die obige Bedingung verletzt.

Wir werden dieses Problem erneut dadurch lösen, dass wir statt der ersten Multipliziererkoeffizienten Γ die Multipliziererkoeffizienten $\Gamma' = \alpha\Gamma$ (mit den Koordinaten $\gamma'_1 = \gamma_1\alpha$ und $\gamma'_2 = \gamma_2\alpha$) mit $1 > \alpha > 0$ verwenden. Der Wert $\tilde{a}'_0 = \tilde{z}_3$ der realen Addition berechnet sich nun durch

$$\tilde{a}'_0 = \tilde{z}_3 = \gamma'_1 a_1 \oplus \gamma'_2 a_2, \quad (\text{D.44})$$

mit den angenommenen Summanden $z_1 = \gamma'_1 a_1 > \gamma'_2 a_2 = z_2 > 0$. Somit ergibt sich der Ausgangsvektor zu

$$\tilde{\mathbf{b}}' = \Gamma \tilde{a}'_0 - \mathbf{a}. \quad (\text{D.45})$$

Nun stellt sich die Frage, wie groß die Stauchung sein muss. Die Antwort ergibt sich aus der Passivitätsbedingung. Demnach darf für beliebige Eingangsvektoren \mathbf{a} die Ausgangsleistung $\|\tilde{\mathbf{b}}'\|^2$ nicht größer als

die Eingangsleistung $\|\mathbf{a}\|^2$ sein. Hinreichend und notwendig dafür ist $|\tilde{a}'_0| \leq |a_0|$. Wir müssen nun den fehlerhaften Wert \tilde{a}'_0 durch a_0 , p und α ausdrücken und die Bedingung nach α auflösen. Der Fehler wird durch die Ungleichung D.15 mit

$$z_3 = a'_0 = \gamma'_1 a_1 + \gamma'_2 a_2 = \alpha \gamma_1 a_1 + \alpha \gamma_2 a_2 = \alpha \mathbf{\Gamma}^T \mathbf{a} = \alpha a_0 \quad (\text{D.46})$$

zu

$$[1 - 2^{-p}] < \frac{|\tilde{a}'_0|}{|a'_0|} < [1 + 2^{-p}] \quad (\text{D.47})$$

eingegrenzt. In Bezug auf die Passivität ist die Aufrundung des Betrages, d. h. die Obergrenze

$$|\tilde{a}'_0| < |a'_0| [1 + 2^{-p}] \quad (\text{D.48})$$

entscheidend. Setzen wir hier $a'_0 = \alpha a_0$ ein, so haben wir

$$|\tilde{a}'_0| < \alpha |a_0| [1 + 2^{-p}] . \quad (\text{D.49})$$

Im schlimmsten Fall darf $|\tilde{a}'_0|$ gerade so groß wie $|a_0|$ sein, d. h. $|a_0| \geq |\tilde{a}'_0|$. Aus der in Gleichung (D.49) bestimmten Obergrenze von $|\tilde{a}'_0|$ erhalten wir den gesuchten Wert für α

$$|a_0| \geq \alpha |a_0| [1 + 2^{-p}] \quad \Leftrightarrow \quad \alpha \leq \frac{1}{[1 + 2^{-p}]} \approx [1 - 2^{-p}] . \quad (\text{D.50})$$

Somit wurde der Nachweis der Passivität bei Stauchung der Multiplizierer und Anwendung des Gitterpunktschneidens für die erste Addition erbracht.

3. Die zweite Multiplikation

Wir gehen hier von einem korrekt berechnetem a_0 aus und betrachten die ideale Multiplikation in $b_1 = \gamma_1 a_0 - a_1$. Das Ergebnis der realen Multiplikation $\gamma_1 \odot a_0$ ist gleich der idealen Multiplikation mit einem modifiziertem γ_1 , d. h.

$$\tilde{z}_3 = \gamma_1 \odot a_0 = \tilde{\gamma}_1 a_0 . \quad (\text{D.51})$$

Gemäß Gleichung (D.32) ist für die Passivität

$$\tilde{\gamma}_1 \leq \gamma_1 \quad (\text{D.52})$$

notwendig. Somit auch

$$\tilde{\gamma}_1 |a_0| \leq \gamma_1 |a_0| \quad (\text{D.53})$$

und mit $\gamma_1, \tilde{\gamma}_1 > 0$

$$|\tilde{\gamma}_1 a_0| \leq |\gamma_1 a_0| . \quad (\text{D.54})$$

Verwenden von Gleichung (D.51) liefert dann

$$|\gamma_1 \odot a_0| \leq |\gamma_1 a_0| . \quad (\text{D.55})$$

Dies bedeutet, dass das Ergebnis der realen Multiplikation dem Betrage nach kleiner als das Ergebnis der idealen Multiplikation sein muss, um die Passivität zu gewährleisten.

a) Rundung mittels Betragsschneiden

Bei Verwendung des Betragsschneidens ergibt sich das Ergebnis der realen Multiplikation durch Abrundung des Betrages einer idealen Rechnung. Ohne weitere Eingriffe ist Gleichung (D.55) erfüllt und die Passivität sichergestellt.

b) Gitterpunktschneiden

Bei Verwendung des Gitterpunktschneidens existieren hingegen Fälle, die Gleichung (D.55) verletzen. Abhilfe schafft für diese Rundungsart ein Stauchen der Koeffizienten. Nach Gleichung (D.27) ist der relative Fehler bei der Multiplikation mit $z_3 = \gamma_1 a_0$ und $\tilde{z}_3 = \tilde{\gamma}_1 a_0 = \gamma_1 \odot a_0$ durch

$$1 - 2^{-p} < \frac{|\gamma_1 \odot a_0|}{|\gamma_1 a_0|} < 1 + 2^{-p} \quad (\text{D.56})$$

eingeschränkt. Das bedeutet, dass der Fall $|\gamma_1 \odot a_0| > |\gamma_1 a_0|$ auftreten kann und Gleichung (D.55), wie eingangs behauptet, verletzt. Um dem zu begegnen, wird in der realen Multiplikation der Koeffizienten γ_1 durch $\gamma'_1 = \gamma_1 / [1 + 2^{-p}]$ ersetzt. Dann haben wir aus Gleichung (D.27) die Fehlerschranke

$$1 - 2^{-p} < \frac{|\gamma'_1 \odot a_0|}{|\gamma'_1 a_0|} < 1 + 2^{-p}, \quad (\text{D.57})$$

d. h.

$$\frac{|[1 + 2^{-p}] \gamma'_1 \odot a_0|}{|\gamma_1 a_0|} < 1 + 2^{-p} \iff |\gamma'_1 \odot a_0| < |\gamma_1 a_0|. \quad (\text{D.58})$$

Offenbar wird nun Gleichung (D.55) immer erfüllt.

4. Die zweite Addition

Hier wird der Rundungsfehler berücksichtigt, der bei der Berechnung von b_1 und b_2 auftritt. Die idealen Additionen lauten $b_1 = \gamma_1 a_0 + (-a_1)$ und $b_2 = \gamma_2 a_0 + (-a_2)$, die realen Additionen $\tilde{b}_1 = \gamma_1 a_0 \oplus (-a_1)$ und $\tilde{b}_2 = \gamma_2 a_0 \oplus (-a_2)$. Für die Passivität muss $||\tilde{\mathbf{b}}|| \leq ||\mathbf{a}|| = ||\mathbf{b}||$ gelten, was durch $|\tilde{b}_1| \leq |b_1|$ und $|\tilde{b}_2| \leq |b_2|$ gewährleistet ist.

a) Rundung mittels Betragsschneiden

Im Falle der Anwendung des Betragsschneidens gilt $|\tilde{b}_1| \leq |b_1|$ und $|\tilde{b}_2| \leq |b_2|$. Die Passivität an dieser Additionsstelle ist nicht gefährdet, da b_1 und b_2 Energiegrößen sind.

b) Gitterpunktschneiden

Im Falle von Gitterpunktschneiden ist dies aber nicht sichergestellt. Im Gegenteil, wir werden zeigen, dass das System selbst dann aktiv ist, wenn wir den Vektor $\mathbf{\Gamma}$ stauchen.

Der Fehler bei der Berechnung von $b_1 = -a_1 + \gamma_1 a_0$ durch $\tilde{b}_1 = -a_1 \oplus \gamma_1 a_0$ ist gemäß Gleichung (D.15) durch

$$1 - 2^{-p} < \frac{|\tilde{b}_1|}{|b_1|} < 1 + 2^{-p} \quad (\text{D.59})$$

eingeschränkt. Somit ist auch

$$|b_1| < |\tilde{b}_1| \iff -|\tilde{b}_1|^2 < -|b_1|^2 \quad (\text{D.60})$$

möglich. Die aufgenommene Leistung berechnet sich dann durch

$$\tilde{p} = \|\mathbf{a}\|^2 - \|\tilde{\mathbf{b}}\|^2 = \|\mathbf{a}\|^2 - |\tilde{b}_1|^2 - |b_2|^2 < \|\mathbf{a}\|^2 - |b_1|^2 - |b_2|^2 = \|\mathbf{a}\|^2 - \|\mathbf{b}\|^2 = 0 \quad (\text{D.61})$$

und ist negativ. Das System ist folglich aktiv.

Die Streumatrix \mathbf{S} wird nun durch $\mathbf{S}' = \mathbf{\Gamma}'\mathbf{\Gamma}'^T - \mathbf{1}_2$ ersetzt, wobei $\alpha\mathbf{\Gamma} = \mathbf{\Gamma}'$ und $\mathbf{\Gamma}^T\mathbf{\Gamma} > \mathbf{\Gamma}'^T\mathbf{\Gamma}'$ gilt. Wir legen die Bezeichnung $\tilde{\mathbf{b}}' = \mathbf{S}' \odot \mathbf{a}$ für die reale Matrixmultiplikation (hier nur mit fehlerhafter Addition bei den Skalarprodukten) fest. Wie wir eben gesehen haben, kann sich \mathbf{b}' maximal durch die zweite Addition auf $\mathbf{b}[1 + 2^{1-p}]$ vergrößern. Dieser Fehler bei der realen Multiplikation mit der Matrix \mathbf{S}' entspricht dem der idealen Multiplikation mit der Matrix $\mathbf{S}'[1 + 2^{1-p}]$, d. h.

$$\tilde{\mathbf{b}}' = \mathbf{S}'[1 + 2^{1-p}]\mathbf{a} = \mathbf{S}' \odot \mathbf{a} . \quad (\text{D.62})$$

Das Ergebnis der realen Multiplikation soll aber gerade der Passivitätsbedingung

$$\|\tilde{\mathbf{b}}'\|^2 = \|\mathbf{S}' \odot \mathbf{a}\|^2 \leq \|\mathbf{a}\|^2 \quad (\text{D.63})$$

genügen. Das heißt, der neue Vektor $\mathbf{\Gamma}'$ muss nun so geartet sein, dass $\mathbf{S}'[1 + 2^{1-p}]$ unitär ist. Die Eigenwerte dürfen dem Betrag nach nicht größer als eins sein, d. h.

$$\begin{aligned} \det[\mathbf{1}_2\lambda - [1 + 2^{1-p}]\mathbf{S}'] &= \det[\mathbf{1}_2\lambda - [1 + 2^{1-p}][\mathbf{\Gamma}'\mathbf{\Gamma}'^T - \mathbf{1}_2]] = \det[\mathbf{1}_2(\lambda + 1 + 2^{1-p}) - [1 + 2^{1-p}]\mathbf{\Gamma}'\mathbf{\Gamma}'^T] \\ &= (\lambda + 1 + 2^{1-p})^2 \det[\mathbf{1}_2 - [1 + 2^{1-p}]\mathbf{\Gamma}'\mathbf{\Gamma}'^T / (\lambda + 1 + 2^{1-p})] \\ &= (\lambda + 1 + 2^{1-p})[(\lambda + 1 + 2^{1-p}) - [1 + 2^{1-p}]\mathbf{\Gamma}'^T\mathbf{\Gamma}'] . \end{aligned} \quad (\text{D.64})$$

Offenbar ist ein Eigenwert $\lambda = -1 - 2^{1-p}$ und somit vom Betrag her größer als eins, ganz gleich wie wir $\mathbf{\Gamma}'$ wählen.

Als Ergebnis halten wir zum einen fest, dass wir für das im Bild D.4 dargestellte Signalflussdiagramm bei Verwendung des Betragsschneidens eine passive Realisierung gefunden haben. Hingegen waren wir bei Verwendung der Gitterpunkt-Rundung nicht in der Lage, ein Signalflussdiagramm passiv zu realisieren. Abschließend sei erwähnt, dass sich die gefundenen Ergebnisse auf Serienadaptoren verallgemeinern lassen.

D.3 Direkte Umsetzung einer Zweitor-Adaptor-Streumatrix

Leider sind die meisten C-Compiler (aufgrund des Zielsystems) auf Gitterpunktschneiden eingerichtet. Wie in D.2 nachgewiesen wurde, führt dies dazu, dass die Implementierung nicht mehr passiv ist. Ein Betragsschneiden müsste daher softwareseitig implementiert werden. Dies ist zwar möglich, stellt sich aber als zu rechenaufwendig heraus und ist für die Praxis völlig ungeeignet. Als Alternative bietet sich die direkte Umsetzung der Streumatrix mit gestauchten Koeffizienten an. Die Berechnung erfolgt im idealen Fall durch

$$b_1 = S_{11}a_1 + S_{12}a_2 \quad ; \quad b_2 = S_{21}a_1 + S_{22}a_2 \quad (\text{D.65})$$

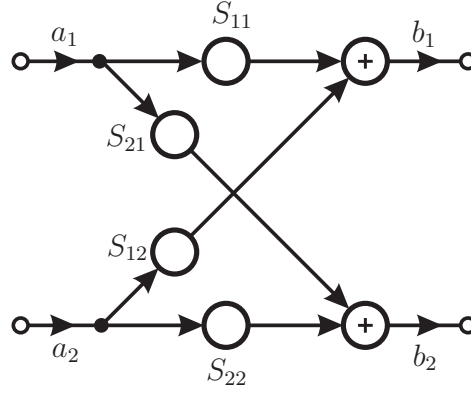


Bild D.6: Signalflussdiagramm eines realen Zweitor-Paralleladaptors bei direkter Umsetzung.

und im realen Fall durch

$$\tilde{b}_1 = S_{11} \odot a_1 \oplus S_{12} \odot a_2 \quad ; \quad \tilde{b}_2 = S_{21} \odot a_1 \oplus S_{22} \odot a_2 \quad , \quad (\text{D.66})$$

wobei auch hier Punkt- vor Strichrechnung gilt. Im Bild D.6 ist das Signalflussdiagramm eines realen Zweitor-Paralleladaptors bei direkter Umsetzung dargestellt.

Die Berechnung erfolgt im realen Fall mit gestauchten Koeffizienten $|S'_{\mu\nu}| < |S_{\mu\nu}|$ für $\mu, \nu = 1, 2$ durch

$$\tilde{b}'_1 = S'_{11} \odot a_1 \oplus S'_{12} \odot a_2 \quad ; \quad \tilde{b}'_2 = S'_{21} \odot a_1 \oplus S'_{22} \odot a_2 \quad . \quad (\text{D.67})$$

Wir müssen zunächst eine Stauchung aufgrund der fehlerhaften Multiplikation vornehmen. Der Fehler bei der Multiplikation kann gemäß Gleichung (D.27) durch

$$[1 - 2^{-p}] < \frac{|S'_{11} \odot a_1|}{|S'_{11} a_1|} < [1 + 2^{-p}] \quad (\text{D.68})$$

eingegrenzt werden. Nehmen wir den schlimmsten Fall an, dass alle 4 Multiplikationen eine Erhöhung der Leistung bewirken, so gilt für die Norm der Ausgangswellen $||\tilde{\mathbf{b}}'|| < (1 + 2^{-p})||\mathbf{b}'||$. Setzen wir $\mathbf{b}' = \mathbf{S}'\mathbf{a}$ ein, so ergibt sich

$$||\tilde{\mathbf{b}}'|| < (1 + 2^{-p})||\mathbf{S}'\mathbf{a}|| \quad . \quad (\text{D.69})$$

Für die Passivität muss $||\tilde{\mathbf{b}}'|| \leq ||\mathbf{a}||$ gelten. Diese Bedingung wird von $\mathbf{S}' = \mathbf{S}/(1 + 2^{-p})$ erfüllt, wie durch Einsetzen in Gleichung (D.69) nachprüfbar ist

$$||\tilde{\mathbf{b}}'|| < (1 + 2^{-p})||\mathbf{S}'\mathbf{a}|| = (1 + 2^{-p})||\mathbf{S}\mathbf{a}/(1 + 2^{-p})|| = ||\mathbf{S}\mathbf{a}|| = ||\mathbf{a}|| \quad . \quad (\text{D.70})$$

Desweiteren müssen wir eine Stauchung der Matrixelemente aufgrund der fehlerhaften Addition durchführen, da durch die reale Addition eine Erhöhung der Leistung auftreten kann. Mit

$$\tilde{b}'_1 = S'_{11}a_1 \oplus S'_{12}a_2 \quad (\text{D.71})$$

gilt die Eingrenzung nach Gleichung (D.15)

$$[1 - 2^{-p}] < \frac{|S'_{11}a_1 \oplus S'_{12}a_2|}{|S'_{11}a_1 + S'_{12}a_2|} < [1 + 2^{-p}] \quad . \quad (\text{D.72})$$

Der zur Erhöhung der Ausgangsleistung führende Fall ist durch $|S'_{11}a_1 \oplus S'_{12}a_2| < |[S'_{11}a_1 + S'_{12}a_2][1 + 2^{-p}]|$ beschränkt. Bewirken beide Additionen eine Erhöhung der Leistung, so gilt

$$||\tilde{\mathbf{b}}'|| < (1 + 2^{-p})||\mathbf{S}'\mathbf{a}|| \quad . \quad (\text{D.73})$$

Für die Passivität muss $\|\tilde{\mathbf{b}}'\| \leq \|\mathbf{a}\|$ gelten. Dies ist dann erfüllt, wenn sich die modifizierte Streumatrix zu $\mathbf{S}' = \mathbf{S}/[1 + 2^{-p}]$ berechnet, denn auch hier ergibt sich durch Einsetzen

$$\|\tilde{\mathbf{b}}'\| < (1 + 2^{-p})\|\mathbf{S}'\mathbf{a}\| = (1 + 2^{-p})\|\mathbf{S}\mathbf{a}/(1 + 2^{-p})\| = \|\mathbf{S}\mathbf{a}\| = \|\mathbf{a}\|. \quad (\text{D.74})$$

Fassen wir die Änderungen wegen fehlerbehafteter Multiplikationen und Additionen zusammen, so berechnet sich insgesamt die modifizierte Streumatrix zu

$$\mathbf{S}' = \frac{\mathbf{S}}{[1 + 2^{-p}]^2}. \quad (\text{D.75})$$

Zum Leistungsvergleich berechnen wir mit $\mathbf{b}' = \mathbf{S}'\mathbf{a} = \mathbf{S}\mathbf{a}/[1 + 2^{-p}]^2$ die Dämpfung in dB

$$\begin{aligned} A_{\text{dB}} &= 10 \lg \left(\frac{\|\mathbf{a}\|^2}{\|\mathbf{b}'\|^2} \right) = 20 \lg \left(\frac{[1 + 2^{-p}]^2 \|\mathbf{a}\|}{\|\mathbf{S}\mathbf{a}\|} \right) = 20 \log_{10} ([1 + 2^{-p}]^2) \approx 20 \log_{10} \left(\frac{1 + 2^{-p}}{1 - 2^{-p}} \right) \\ &= \frac{20}{\ln 10} \ln \left(\frac{1 + 2^{-p}}{1 - 2^{-p}} \right) \approx \frac{20}{\ln 10} (2 \cdot 2^{-p}) = \frac{10}{\ln 10} 2^{2-p} \approx 2^{4-p}. \end{aligned} \quad (\text{D.76})$$

Für $p = 24$ ergibt sich $A_{\text{dB}} \approx 2^{-20} \text{dB} \approx 1000^{-2} \text{dB} = 1 \mu\text{dB}$.

Wir wollen festhalten, dass sich die Ergebnisse auch hier auf Serienadaptoren verallgemeinern lassen. Zum Abschluss fassen wir das Vorgehen zur Bestimmung der Multipliziererkoeffizienten bei direkter Implementierung nach Bild D.6 zusammen

1. Berechnung von \mathbf{S}
2. Stauchung der Streumatrix zu $\mathbf{S}' = \mathbf{S}/[(1 + 2^{-p})]^2$
3. Quantisierung der Elemente von \mathbf{S}' durch Betragsschneiden.

Literaturverzeichnis

- [Babe95] R. L. Baber: *Praktische Anwendbarkeit mathematisch rigoroser Methoden zum Sicherstellen der Programmkorrektheit*, Walter de Gruyter Verlag, Berlin, 1995.
- [Bele68] V. Belevitch: *Classical network theory*, Holden-Day, Inc., 1968.
- [BF87] S. Basu, A. Fettweis: "New Results on Stable Multidimensional Polynomials-Part II: Discrete Case", *IEEE Transactions on Circuits and Systems*, Bd. 34, Nr. 11, S. 1264-1274, 1987.
- [Bilb01] S. D. Bilbao: "Wave and Scattering Methods for the Numerical Integration of Partial Differential Equations", Dissertation an der Stanford University, 2001.
- [Bilb04] S. D. Bilbao: *Wave and Scattering Methods for the Numerical Simulation*, John Wiley and Sons. Chichester, UK, 2004, 2004.
- [Boeh76] B. Boehm: "Software engineering", *IEEE Transactions on Electronic Computers* **25**(12), 1976.
- [Boeh81] B. Boehm: *Software Engineering Economics*, Prentice-Hall, 1981.
- [Boeh86] B. Boehm: "A spiral model of program development and enhancement", *Software Engineering Notes* **11**(4), 1986, S. 14-24.
- [Bose79] N. K. Bose (Hrsg.): *Multidimensional Systems: Theory and Application*, IEEE Press Selected Reprint Series, New York: IEEE Press, 1979.
- [Bose82] N. K. Bose: *Applied Multidimensional System Theory*, Van Nordstrand Reinhold Company, 1982.
- [Bose01] N. K. Bose: "Multidimensional Signal Processing", Vorlesung an der Ruhr-Universität Bochum, 2001.
- [Brey96] U. Breymann: *C++ Eine Einführung*, Carl Hanser Verlag, 1996.
- [Bun96] "IT-Phasenmodell (Vorgehensmodell V-Modell 97)", *Bundesanzeiger*, Nr.125, Seite 7722ff., 1996.
- [BV03] S. Buchholz, M. Vollmer: "Ein Werkzeug zur automatische Codegenerierung zur numerischen Integration linearer, konstanter, partieller Differentialgleichungen mit Hilfe von Wellendigitalfiltern", Interner Bericht des Lehrstuhls für Nachrichtentechnik der Ruhr-Universität Bochum, 2003.
- [Dahl63] G. Dahlquist: "A special stability problem for linear multistep methods", *BIT* **3**, S. 27-43, 1963.

- [Dijk68] E. W. Dijkstra: "A Constructive Approach to the Problem of Program Correctness", *BIT Nord. Tidskr. Inform. Bd. 8, S. 174-186*, 1968.
- [FB87] A. Fettweis, S. Basu: "New Results on Stable Multidimensional Polynomials-Part I: Continuous Case", *IEEE Transactions on Circuits and Systems, Bd. 34, Nr. 10, S. 1221-1232*, 1987.
- [Feld95] T. Felderhoff: "Digitale Simulation nichtlinearer Systeme mit Methoden der Netzwerktheorie", Dissertation an der Universität-Gesamthochschule Paderborn, 1995.
- [Fett70] A. Fettweis: "Entwurf von Digitalfiltern in Anlehnung an Verfahren der klassischen Netzwerktheorie", *NTZ-Fachtagung 15./16.10.1970 Stuttgart*, 1970.
- [Fett78] A. Fettweis: "Suppression of parasitic oscillations in multidimensional wave digital filters", *IEEE Transactions on Circuits and Systems, Bd. 25, Nr. 12, S. 1060-1066*, 1978.
- [Fett79] A. Fettweis: "Principles of multidimensional wave digital filtering". in J. K. Aggarwal (Hrsg.), *Digital Signal Processing*, S. 261-282. Western Periodicals, Point Lobos Press, North Hollywood, CA, USA, 1979.
- [Fett86] A. Fettweis: "Wave digital filters: Theory and Practice", (*invited paper*) *Proceedings of the IEEE (The Institute of Electrical and Electronics Engineers), Bd. 74, Nr. 2, S. 270-327*, Februar 1986.
- [Fett90] A. Fettweis: "On assessing robustness of recursive digital filters", *European Transactions on Telecommunications, Bd. 1, Nr. 2, S. 103-109*, 1990.
- [Fett92] A. Fettweis: "Discrete passive modelling of physical systems described by PDEs", *European Signal Processing Conference (EUSIPCO), Bd. 1, S. 55-62 Brüssel, Belgien, 24.08. - 27.08.*, 1992.
- [Fett98] A. Fettweis: "Numerische Integration nach dem Wellendigitalfilterprinzip", Vorlesung an der Ruhr-Universität Bochum, 1998.
- [Fett99] A. Fettweis: "Numerische Integration nach dem Wellendigitalfilterprinzip", Vorlesung an der Ruhr-Universität Bochum, 1999.
- [FH92] A. Fettweis, G. Hemetsberger: *Grundlagen der Theorie elektrischer Schaltungen*, Universitätsverlag Dr. N. Brockmeyer, 1992.
- [Fisc84] H. D. Fischer: "Wave digital filters for numerical integration", *ntz-Archiv, Bd. 6, Nr. 2, S. 37-40*, 1984.
- [Fisc99] H. D. Fischer: "Numerische Optimierung", Vorlesung an der Ruhr-Universität Bochum, 1999.
- [Fisc00a] H. D. Fischer: "Ergänzungen zur Berechnung elektrischer Netze im stationären Zustand", Vorlesung an der Ruhr-Universität Bochum, 2000.
- [Fisc00b] H. D. Fischer: "Persönliche Mitteilung", 2000.
- [Fisc03] H. D. Fischer: "Technische Zuverlässigkeit", Vorlesung an der Ruhr-Universität Bochum, 2003.

- [Floy67] R. W. Floyd: "Assigning Meanings to Programs", *Proceedings of the Symposium of Applied Mathematics, Bd. 19, S. 19-32, American Mathematical Society, Providence, RI, 1967.*
- [FM75] A. Fettweis, K. Meerkötter: "Suppression of parasitic oscillations in wave digital filters", *IEEE Transactions on Circuits and Systems, Bd. 22, Nr. 3, S. 239-246, 1975.*
- [FN90a] A. Fettweis, G. Nitsche: "Massively parallel algorithms for numerical integration of partial differential equations", *International Workshop on Algorithms and Parallel VLSI Architectures, Summaries of Contributions, S. 475-484, Pont-à-Mousson, France, 1990.*
- [FN90b] A. Fettweis, G. Nitsche: "Numerical integration of partial differential equations by means of multidimensional wave digital filters", *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS '90), Bd. 2, S. 954-957, New Orleans, LA, USA, 30 May - 2 June, 1990.*
- [FN91a] A. Fettweis, G. Nitsche: "Numerical integration of partial differential equations using principles of multidimensional wave digital filters", *Journal of VLSI Signal Processing, Bd. 3, S. 7-24, Kluwer Academic Publishers, Boston, MA, USA, 1991.*
- [FN91b] A. Fettweis, G. Nitsche: "Transformation approach to numerically integrating PDEs by means of WDF principles", *Multidimensional Systems and Signal Processing, Bd. 2, Nr. 2, S. 127-159, 1991.*
- [Frän97] D. Fränken: "Passive Systeme zur Verarbeitung komplexer zeitdiskreter Signale", Dissertation an der Universität-Gesamthochschule Paderborn, 1997.
- [Frie54] K. O. Friedrichs: "Symmetric hyperbolic differential equations", *Comm. Pure Appl. Math., Bd. 7, S. 354-392, 1954.*
- [Frie95] M. Fries: "Numerical integration of Euler flow by means of multidimensional wave digital principles", Dissertation an der Ruhr-Universität Bochum, 1995.
- [Frie00] M. Fries: "Persönliche Mitteilung", 2000.
- [FS99] A. Fettweis, G. A. Seraji: "New results in numerically integrating PDEs by the wave digital approach", *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS '99), Bd. 5, S. 17-20, Orlando, FL, USA, 30 May - 2 June, 1999.*
- [HA28] D. Hilbert, W. Ackermann: *Grundzüge der theoretischen Logik*, Springer-Verlag, 1928.
- [Heme95] G. Hemetsberger: "Numerische Integration hyperbolischer partieller Differentialgleichungen unter Verwendung mehrdimensionaler Wellendigitalfilter", Dissertation an der Ruhr-Universität Bochum, 1995.
- [Hero99] H. Herold: *C-Kompaktreferenz*, ADDISON-WESLEY LONGMAN, 1999.
- [HNW93] E. Hairer, S. P. Nørsett, G. Wanner: *Solving Ordinary Differential Equations I*, Springer-Verlag, 1993.
- [Hoar69] C. Hoare: "An Axiomatic Basis for Computer Programming", *CACM, Bd. 12, Nr. 10, S. 576-580+583, 1969.*
- [Huan81] Huang: *Two-Dimensional Digital Signal Processing I*, Springer-Verlag, 1981.

- [Inte99] International Electrotechnical Commission: "Programming languages – C", *ISO/IEC 9899, International Electrotechnical Commission, Genf*, 1999.
- [Juss00] J. Jussen: "Untersuchungen zur automatischen Codegenerierung einer steuerbaren, eindimensionalen Wellendigitalstruktur", Studienarbeit an der Ruhr-Universität Bochum, 2000.
- [Koga69] T. Koga: "Synthesis of passive n -ports with prescribed positive real matrices of several variables", *IEEE Transactions on Circuit Theory, Bd. 15, S. 2-23*, 1969.
- [Krau97] H. Krauß: "Simulation linearer dissipativer Wellenausbreitungsvorgänge mit mehrdimensionalen digitalen Filtern", Dissertation an der Friedrich Alexander Universität Erlangen-Nürnberg, 1997.
- [Kumm88] A. Kummert: "Beiträge zur Synthese mehrdimensionaler Reaktanzmehrtore", Dissertation an der Ruhr-Universität Bochum, 1988.
- [KWU 98] KWU NLLZ ST: "Systemübersicht zum digitalen Leitsystem Teleperm XS", 1998.
- [KWU 99] KWU NLLZ ST: "Entwicklungsunterlagen : Design FB ModuleDesign", 1999.
- [LF90] X. Liu, A. Fettweis: "Multidimensional Digital Filtering by Using Parallel Algorithms Based on Diagonal Processing", *Multidimensional Systems and Signal Processing, S. 51-66*, 1990.
- [Linn84] G. Linnenberg: "Über die diskrete Verarbeitung mehrdimensionaler Signale unter Verwendung mehrdimensionaler Wellendigitalfilter", Dissertation an der Ruhr-Universität Bochum, 1984.
- [Lisc91] J. Lischetzki: "Untersuchungen zur Genauigkeit von Wellendigitalfilterverfahren zur numerischen Integration partieller Differentialgleichungen", Diplomarbeit an der Ruhr-Universität Bochum, 1991.
- [Luhm04] K. Luhmann: "Die numerische Lösung der Neutronendifusionsgleichungen in zwei Energiegruppen mit dem Wellendigital-Konzept", Dissertation an der Ruhr-Universität Bochum, 2004.
- [McCa62] J. McCarthy: "Towards a Mathematical Science of Computation", *Proceedings of the IFIP Congress, S.21-28, North-Holland, Amsterdam*, 1962.
- [McCa63] J. McCarthy: "A Basis of Mathematical Theory of Computations", *Computer programming and formal systems, S.33-70, North-Holland, Amsterdam*, 1963.
- [Meer79] K. Meerkötter: "Beiträge zur Theorie der Wellendigitalfilter", Dissertation an der Ruhr-Universität Bochum, 1979.
- [Meye88] B. Meyer: *Objektorientierte Softwareentwicklung*, Hanser, 1988.
- [MF92] K. Meerkötter, T. Felderhoff: "Simulation of nonlinear transmission lines by wave digital filter principles". in *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS '92), Bd. 2, S. 875-878, San Diego, CA, USA*, 1992.
- [MF96] K. Meerkötter, D. Fränken: "Digital Realization of Connection Networks by Voltage-Wave Two-Port Adaptors". in *Archiv für Elektronik und Übertragungstechnik 50*, 1996.

- [Nits93] G. Nitsche: "Numerische Lösung partieller Differentialgleichungen mit Hilfe von Wellendigitalfiltern", Dissertation an der Ruhr-Universität Bochum, 1993.
- [Novi73] P. S. Novikov: *Mathematische Logik*, Friedrich Vieweg und Sohn Verlag, 1973.
- [Ochs01a] K. Ochs: "Passive Integration Methods : Fundamental Theory", *Archiv für Elektronik und Übertragungstechnik*, Bd. 55, Nr. 3, S. 153-163, 2001.
- [Ochs01b] K. Ochs: "Passive Integrationsmethoden", Dissertation an der Universität-Gesamthochschule Paderborn, 2001.
- [Ochs02] K. Ochs: "Persönliche Mitteilung", 2002.
- [Pott98] R. Pott: "Numerische Integration von Neutronendifusionsgleichungen unter Verwendung mehrdimensionaler Wellendigitalfilter auf parallelen Rechnerarchitekturen", Dissertation an der Ruhr-Universität Bochum, 1998.
- [Rao69] T. Rao: "Minimal synthesis of two-variable reactance matrices", *Bell Systems Technical Journal*, S. 163-199, 1969.
- [Royc87] W. Royce: "Managing the development of large software systems: Concepts and techniques", 1987, S. 328-338. Nachdruck aus dem Jahre 1970.
- [Rumm98] E. Rummert: "Methodik eines formalen Korrektheitsbeweises bei graphisch spezifizierter Software", Dissertation an der Ruhr-Universität Bochum, 1998.
- [Sera00] G. A. Seraji: "Randwertfragen bei der numerischen Integration mit Hilfe der Wellendigitalmethode", Dissertation an der Ruhr-Universität Bochum, 2000.
- [Shek74] J. Shekel: "The junction matrix in the analysis of scattering networks", *IEEE Transactions on Circuits and Systems*, Bd. 21, Nr. 1, S. 21-22, 1974.
- [Smid75] D. Smidt: *Reaktortechnik*, G. Braun, Karlsruhe, 1975.
- [Stra80] G. Strang: *Linear Algebra and Its Applications*, Harcourt Brace Jonanovich, 1980.
- [Tell54] B. D. H. Tellegen: *Theorie der Elektrische Netwerken*, Nordhoff, 1954.
- [Voll02] M. Vollmer: "Automatische Codegenerierung zur numerischen Integration linearer, konstanter, partieller Differentialgleichungen mit Hilfe von Wellendigitalfiltern", Interner Bericht des Lehrstuhls für Nachrichtentechnik der Ruhr-Universität Bochum, 2002.
- [Voll04a] M. Vollmer: "An approach to automatic generation of wave digital structures from PDEs", *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS '04)*, Vancouver, Canada, 23 May - 26 May, 2004.
- [Voll04b] M. Vollmer: "Numerical Integration of PDEs for Safety Critical Applications implemented by I & C Systems", *The 23rd International Conference on Computer Safety, Reliability and Security*, 21-24 September 2004, Potsdam, Germany, 2004.
- [Voll04c] M. Vollmer: "Reference circuits for nonenergetic multi-ports without directed loops in the wave domain", *zur Veröffentlichung eingereicht*, 2004.
- [Waed94] K. Waedt: "Fehlervermeidende Codegenerierung für verteilte, responsive Systeme", Dissertation an der Friedrich Alexander Universität Erlangen-Nürnberg, 1994.

[Waed00] K. Waedt: “Persönliche Mitteilung”, 2000.

Lebenslauf

Persönliche Daten

Name:	Michael Vollmer
Geburtsdatum:	14. Dezember 1971
Geburtsort:	Hamm
Familienstand:	ledig
Staatsangehörigkeit:	deutsch

Schulbildung

08/1978 – 07/1982	Grundschule in Hamm
08/1982 – 06/1988	Realschule in Hamm
08/1991 – 07/1992	Fachoberschule in Hamm

Berufsausbildung

09/1988 – 06/1991	Ausbildung zum Energieelektroniker bei der Deutschen Bundesbahn in Hamm
-------------------	---

Wehrdienst

10/1992 – 09/1993	Grundwehrdienst in Hamm
-------------------	-------------------------

Studium

10/1993 – 05/1998	Elektrotechnik mit der Vertiefungsrichtung Informationstechnik in Paderborn
-------------------	---

Berufstätigkeit

09/1998 – 09/2004	Wissenschaftlicher Mitarbeiter am Lehrstuhl für Nachrichtentechnik der Ruhr-Universität Bochum
-------------------	--