

# Formalsprachliche Aspekte von XML

Matthias Wendlandt

5. Februar 2008

Ich möchte mich bei meinen Eltern Elke und Günter Wendlandt bedanken. Sie gaben mir in jeder Phase der Promotion den Mut und die Kraft zielstrebig zu arbeiten.

Ganz besonders möchte ich mich bei meinem Betreuer Herrn Prof. Dr. Martin Kutrib bedanken, dessen Ratschläge, Anregungen und kritische Äußerungen außerordentlich wertvoll für diese Arbeit waren.

Natürlich möchte ich mich auch bei allen anderen Personen bedanken, die an den Diskussionen über die vorliegende Arbeit beteiligt waren, insbesondere bei Katja Meckel.

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
1.1	Problemgegenstand . . . . .	1
1.2	Problembereich . . . . .	6
1.3	Fragestellung und Aufbau . . . . .	10
<b>2</b>	<b>Grundlagen</b>	<b>12</b>
2.1	Allgemeine Definitionen . . . . .	12
2.2	k-lookahead deterministische reguläre Ausdrücke . . . . .	17
2.3	Balancierte Grammatiken . . . . .	18
2.4	Untertypen der balancierten Grammatiken . . . . .	23
<b>3</b>	<b>Durchschnitt von regulären Sprachen und Dyckprimes</b>	<b>27</b>
3.1	Konstruktion von RegDyck Automaten . . . . .	27
3.2	RegDyck Grammatiken . . . . .	43
<b>4</b>	<b>Einordnung der Modelle</b>	<b>54</b>
4.1	XML Ausdrücke . . . . .	54
4.2	Einordnung in die Chomsky Hierarchie . . . . .	57
4.2.1	XML-artige Sprachen und reguläre Sprachen . . . . .	57
4.2.2	Auflösung der Struktur von balancierten Grammatiken	58
4.2.3	Balancierte Grammatiken und die Chomsky Hierarchie	60
4.3	Hierarchie der XML-artigen Grammatiken . . . . .	73
4.4	Parser Generierung . . . . .	78
4.4.1	Sprachen und Methoden . . . . .	78
4.4.2	Konstruktion von codeterministischen Grammatiken aus balancierten Grammatiken . . . . .	78
4.4.3	Konstruktion eines deterministischen Kellerautomaten aus einer balancierten Grammatik . . . . .	81
4.4.4	Konstruktion eines Parsebaums für XML Grammatiken	82
4.4.5	Konstruktion eines regulären Ausdrucks und einer Dyck Grammatik aus einer RegDyck Grammatik . . . . .	86
4.5	Parsing Algorithmen . . . . .	87

---

4.5.1	Parsen mittels codeterministischer Grammatiken . . .	88
4.5.2	Parsen mittels Parsebäumen . . . . .	91
<b>5</b>	<b>Eigenschaften von XML-artigen Sprachen</b>	<b>94</b>
5.1	Abschlusseigenschaften . . . . .	95
5.1.1	Vereinigung . . . . .	95
5.1.2	Durchschnitt . . . . .	98
5.1.3	Komplement . . . . .	99
5.1.4	Iteration . . . . .	101
5.1.5	Konkatenation . . . . .	103
5.1.6	Spiegelung . . . . .	107
5.1.7	Substitution . . . . .	109
5.1.8	Homomorphismus . . . . .	111
5.1.9	Beispiele für Abschlusseigenschaften von XML-artigen Grammatiken . . . . .	114
5.2	Entscheidbarkeitsprobleme . . . . .	116
5.2.1	Balancierte Grammatiken und ihre Untertypen . . . .	116
5.2.2	Inklusion . . . . .	119
5.2.3	Mehrdeutigkeit . . . . .	120
5.2.4	Beispiele zu Mehrdeutigkeit . . . . .	124

# Abbildungsverzeichnis

3.1	Kellerinhalt des Kellerautomaten $\mathcal{D}$ nach Abarbeitung des Eingabezeichens $a$ . . . . .	32
3.2	Kellerinhalt des Kellerautomaten $\mathcal{K}$ nach Abarbeitung des Eingabezeichens $a$ . . . . .	33
3.3	Kellerinhalt des Kellerautomaten $\mathcal{D}$ nach Abarbeitung des Eingabezeichens $\bar{a}$ . . . . .	34
3.4	Kellerinhalt des Kellerautomaten $\mathcal{K}$ nach Abarbeitung des Eingabezeichens $\bar{a}$ . . . . .	34
3.5	Automat $\mathcal{R}$ im Zustand $s_i$ mit der Eingabe $a$ . . . . .	36
3.6	Automat $\mathcal{K}$ im Zustand $s_i$ mit der Eingabe $a$ . Die Beschriftung einer Kante mit $a, \bar{x} / \overline{ax}$ bedeutet, dass, wenn der Automat ein Zeichen $a$ einliest und das Zeichen $\bar{x}$ auf dem Keller liegt, dann wechselt er in den Zustand zu dem die beschriftete Kante führt und legt $\overline{ax}$ auf dem Keller ab. . . . .	36
3.7	Automat $\mathcal{R}$ im Zustand $s_i$ mit der Eingabe $\bar{a}$ . . . . .	37
3.8	Automat $\mathcal{K}$ im Zustand $s_i$ mit der Eingabe $\bar{a}$ . . . . .	38
4.1	Genaue Einordnung in die Chomsky Hierarchie. . . . .	70
4.2	XML-artige Grammatiken und $LL(1)$ Grammatiken. . . . .	72
4.3	Hierarchie der XML-artigen Grammatiken. . . . .	77
4.4	Beispiel: Parsebaum zur Validierung eines XML Dokuments. . . . .	84

# Kapitel 1

## Einleitung

### 1.1 Problemgegenstand

*„I guess that XML will, by itself, solve the information interchange problem. In reality, XML just clears away some of the syntactical distractions so that we can get down to the big problem: how we arrive at common understandings about knowledge representation. That's the biggie. XML says, let's stop arguing about how we're going to represent trees and how we're going to represent attribute/value pairs. We'll just decide, let's do it as XML“.*

- Jon Bosak, Entwickler von XML [Bos05] -

Unternehmen tauschen Daten, wie Rechnungs- und Personaldaten, zwischen ihren Standorten aus. Die Raumfahrt liefert große Mengen an Informationen aus dem Weltraum und Konstruktionspläne von Maschinen werden archiviert.

Die zunehmende Globalisierung fordert eine einfache und schnelle Austauschbarkeit von Daten. Daher sind Firmenstandorte über Intra- oder Internet verbunden. Auch Universitäten arbeiten landesübergreifend zusammen und tauschen ihre Daten aus. Privatpersonen schicken Daten an entfernte Verwandte und Freunde.

Ummengen von Daten müssen in ihrer Darstellung klar strukturiert werden. Es muss nicht nur möglich sein das Design, sondern auch die inneren Beziehungen der Daten zu speichern, Hierarchien darzustellen und Eigenschaften Werte zuzuweisen. Vielmehr brauchen die Daten darüber hinaus ein einheitliches Format, das dennoch präzise genug ist, um die Daten mit verschiedenen Anwendungen bearbeiten zu können. Es muss ein einfacher Standard sein, für den sich interpretierende Software leicht entwickeln lässt.

Die herkömmlichen Datenformate dienten vor allem zur Beschreibung der

Darstellung der Daten und weniger zur Beschreibung ihrer Struktur. Ein einfacher Austausch war, bedingt durch die unterschiedlichen Formate der Anbieter, nicht möglich [Bos05].

In den 70er Jahren wurde SGML entwickelt. Dieses Format unterlag der generischen Kodierung, in der, statt der Formatierungscodes, deskriptive Tags verwendet wurden. Es fand vor allem Anwendung im Verlagswesen sowie im US Verteidigungsministerium. Für kleinere Firmen war es allerdings ungeeignet, da nur Spezialisten mit dieser Sprache umgehen konnten [Bos05].

In den 90er Jahren wurde HTML entwickelt. Es war der Anfang der generischen Kodierung für jedermann. HTML ist eine Anwendung von SGML, die eine feste Menge von Tags (deutsch: Abgrenzer) zur Beschreibung des Inhalts besitzt. Sie wurde für das Internet entwickelt und verbreitete sich der Einfachheit wegen sehr schnell. Es war kein Problem, interpretierende Software zu schreiben oder selber Dokumente mit ihr zu erzeugen.

Die Sprache entspricht jedoch nicht den allgemeinen Grundregeln der generischen Kodierung. Die feste Menge von Tags wird mit wenigen Ausnahmen nur zum Design des Dokumentes verwendet.

Es wurde eine Sprache benötigt, die der Mächtigkeit von SGML genügt und gleichzeitig die Einfachheit von HTML beibehält. C. M. Sperberg McQueen, John Bosak, Tim Bray und James Clark arbeiteten 1996 an einer SGML Lite Version [GP00]. Sie wurde 1998 als XML 1.0 (eXtensible Markup Language) [WWW] vom W3 Konsortium standardisiert [XML05].

XML hat sich seit der Veröffentlichung 1998 stark verbreitet. So ist es vom W3 Konsortium als Metasprache für Auszeichnungssprachen im Internet festgelegt worden und dient als Austauschformat jeglicher Daten, die über das Internet verschickt werden.

Die XML Dateien bestehen aus Tags und Fließtext zwischen den Tags. Es gibt zwei verschiedene Sorten von Tags: öffnende und schließende Tags. Ein öffnender Tag besteht aus einer öffnenden, spitzen Klammer <, einem Namen und einer schließenden, spitzen Klammer >. Ein schließender Tag besitzt zusätzlich nach der öffnenden, spitzen Klammer einen Slash /. Zu jedem öffnenden Tag muss ein schließender Tag existieren. Innerhalb eines Tag-paars können neue Tags oder Fließtext stehen. Es gibt genau einen Start- und genau einen Endtag.

### Beispiel 1.1.1 Ein XML Dokument

```
<Dokument>
  <Buch>
    <Titel> Don Quijote </Titel>
    <Kapitel>
```

```

    <Ktitel> Junker Don Quijote </Ktitel>
    <Paragraph>...</Paragraph>
  </Kapitel>
</Buch>
</Dokument>

```

□

XML ermöglicht es, Hierarchien darzustellen und einzelnen Elementen Attribute zuzuordnen. Zudem besitzt sie keine feste Menge von Tags, sondern die Menge kann selbst definiert werden. Um ein Standardformat von Dokumenten zu erfassen, können Dokumentklassen durch ein übergeordnetes Beschreibungssystem angelegt werden.

Diese Beschreibungssysteme sind ähnlich aufgebaut wie Grammatiken der Chomsky Hierarchie. Als Besonderheit verlangen sie einen Starttag zum Einleiten der rechten Seite einer Produktion und einen Endtag zum Beenden der Produktion.

Es existiert eine große Anzahl verschiedener Beschreibungssysteme, wie zum Beispiel die DTD [DTD07], das XML Schema [XML07], TREX [Claa] oder RELAX NG [Clab].

Die Beschreibungssysteme erläutern, wie die Tagpaare eines XML Dokumentes angeordnet sein müssen.

Eine DTD besteht aus einer besonderen Art von Tags. Diese Tags werden auch als Elemente bezeichnet. Jedes Element beschreibt ein Tagpaar aus dem später entstehendem XML Dokument. Sie beginnen mit einer spitzen Klammer <, darauf folgen ein Ausrufezeichen!, der Elementname, eingeklammert und durch Komma getrennt die Namen, der in diesem Element enthaltenen Tags, bzw. mit #PCDATA gekennzeichnet, die Möglichkeit Fließtext zu erlauben. Zum Schluss steht die schließende, spitze Klammer >.

### Beispiel 1.1.2 Die DTD *Buch*

```

<!ELEMENT Buch (Titel, ISBN, Kommentare*, Autor+)>
<!ELEMENT Titel (#PCDATA)>
<!ELEMENT ISBN (#PCDATA)>
<!ELEMENT Kommentare (#PCDATA)>
<!ELEMENT Autor (#PCDATA)>

```

□

Alles was in den Elementen nach dem Namen in den Klammern aufgeführt wird, heißt Inhaltsmodell und kann aus einem regulären Ausdruck über den Namen der Elemente bestehen. Ein Stern bedeutet dabei, dass das Element in dem XML Dokument mehrere Male oder auch keinmal hintereinander aufgeführt werden kann. Das Plus bedeutet, dass es mindestens einmal enthalten sein muss, genauso gut aber auch mehrere Male vorkommen kann.

Ein großer Anwendungsbereich von XML findet sich im Datenaustausch. Durch die Beschreibungssysteme kann eine Struktur der Daten festgelegt werden. Die Beschreibungssysteme geben an, welche Tags existieren und in welcher Form sie angeordnet werden dürfen. Die Verarbeitung der Daten durch Programme wird hiermit vereinfacht.

Der Datenaustausch kann durch die allgemeine Struktur der Daten zwischen zwei gleichen Systemen, zum Beispiel zwei Datenbanken oder zwischen zwei unterschiedlichen Systemen, zum Beispiel einem Textverarbeitungssystem und einem Internet Content Management System, geschehen. Die Programme brauchen lediglich einen Algorithmus, der angibt, wie die Struktur, zum Beispiel eine hierarchische Beziehung, verarbeitet wird.

Ein weiterer Vorteil von XML im Datenaustausch ist, dass komplizierte Datenstrukturen, wie Objekte, Listen und Bäume, dargestellt und verschickt werden können [See03]. Die standardisierte Struktur macht es verarbeitenden Programmen einfach, die Struktur zu erkennen, obwohl sie womöglich den Inhalt der Daten nicht kennen.

XML kann beim Festlegen von Standarddokumenten hilfreich sein. Briefe und Firmendokumente erhalten damit eine einheitliche Form, die von verschiedenen Softwareprodukten bearbeitet werden kann. So ist der Firmenauftritt standardisiert.

In modernen Programmen werden XML Dateien verwendet, um Daten abzuspeichern. Spielstände von Computerspielen werden darin abgelegt oder Maße von Bauteilen gespeichert.

Die DOM (Document Object Model) (Spezifikation, [DOM07]) (nähere Informationen über den Gebrauch, [GP00]) und SAX (Simple API for XML) (Spezifikation, [SAX07]) (nähere Informationen über den Gebrauch, [Meg]) Schnittstellen erleichtern es den Programmierern, Daten aus den Dateien auszulesen. Die Programmierung eines komplizierten Dateiinhaltsverzeichnisses gehört der Vergangenheit an. Die DOM Schnittstelle [Har99] liefert die Daten als Baum und bietet über Methoden einfache Zugriffsmöglichkeiten auf die Elemente. Die SAX Schnittstelle liefert die Daten seriell.

Die programmiertechnischen Vorteile können auch in der einfacheren Planung der Software genutzt werden. Es werden keine komplizierten Schnittstellen definiert, Daten werden durch XML Beschreibungsmodelle standardisiert und können somit einfacher verarbeitet werden.

Die Speicherung von Daten in Datenbanken kann durch XML modernisiert werden. So gibt es Ansätze, XML als eine Art Datenbank [See03] zu verwenden. Weil XML mächtigere semantische Beziehungen als relationale Datenbanken darstellen kann, lassen sich beispielsweise Hierarchien abbilden.

Im Bereich der Datenbankenforschung nimmt XML einen wichtigen Platz ein. In diesen Bereich fallen vor allem Verbindungen zu den relationalen Datenbanken, die die häufigste Speicherform von Daten in Datenbanken sind. Transformationsalgorithmen, welche dafür sorgen, dass eine XML Datei in ein relationales Datenbankkonzept umgewandelt wird, wurden in [DFS99] [FK99] [STZ<sup>+</sup>99] erforscht. Die Arbeit [STZ<sup>+</sup>99] gibt drei Algorithmen an, mit denen aus einer DTD ein relationales Modell erzeugt wird. In ihr wird vor allem die Anzahl der entstehenden Tabellen untersucht. [FK99] stellt acht verschiedene Algorithmen vor, in denen er auch DTDs in ein relationales Datenbankmodell umwandelt und berücksichtigt dabei vor allem die Leistungsmerkmale in Bezug auf Attribut-Werte Paare.

[DFS99] gibt einen Algorithmus an, der ein XML Dokument direkt in ein relationales Modell umwandelt, ohne eine DTD zu benötigen.

Alle Algorithmen setzen die Struktur der XML Daten in ein relationales Schema um, doch die Semantik der Daten bleibt unbehandelt.

Das folgende Beispiel zeigt eine Überführung einer DTD:

**Beispiel 1.1.3** Überführung DTD *Tagung* in ein relationales Schema

```
<!ELEMENT Tagung (Titel,Gesellschaft,Jahr,Papier+)>
<!ELEMENT Papier (Pid,Titel,Zusammenfassung?)>
```

Es sei angenommen, der *Titel* und das *Jahr* in *Tagung* identifizieren eindeutig jedes Tupel in *Tagung*.

Wenn der Hybrid Inlining Algorithmus aus [STZ<sup>+</sup>99] auf die DTD angewendet wird, entsteht folgendes relationales Modell:

```
Tagung (Titel,Gesellschaft,Jahr)
Papier (Pid,Titel,TagungTitel,TagungJahr,Zusammenfassung)
```

Die Struktur wird richtig übertragen, aber die Algorithmen fordern keine korrekte Semantik. Es könnte ein Tupel

$$t_0 : (100, \text{'DTD...'}, \text{'ER'}, 3000, \text{'...'})$$

in *Papier* eingefügt werden. Die DTD verlangt, dass ein *Papier* nur in Verbindung mit einer *Tagung* bestehen darf. Hierbei ist es möglich, dass eine *Tagung* nicht existiert. Das führt zu einem semantischen Fehler.

In relationalen Datenbankbedingungen könnte dies durch eine Teilmengenabhängigkeit gelöst werden.

□

Die Arbeiten [LMM00] [MLM01a] [LMCC02] beschäftigen sich mit diesen Problemen. Es werden Algorithmen vorgestellt, mit denen die XML Modelle semantisch korrekt dargestellt werden. Sie beschränken sich jedoch ausschließlich auf DTD Beschreibungsmodelle. Dazu werden Praxisstudien beschrieben.

In einer weiteren Arbeit [SSB<sup>+</sup>01] wird beschrieben, wie aus einem relationalen Modell Daten abgefragt werden und in einem XML Dokument präsentiert werden.

Dies dient vor allem der Darstellung bestehender Daten im Internet oder hilft bei der Übertragung von Daten via Internet. Dabei werden SQL Kommandos mit Zusätzen ergänzt. Das Ergebnis einer Anfrage ist dann eine XML Datei.

## 1.2 Problembereich

Produktionen der XML-artigen Grammatiken bestehen auf der rechten Seite aus regulären Ausdrücken. Die meisten praktischen Beschreibungsmodelle, wie zum Beispiel DTDs, dürfen keine mehrdeutigen regulären Ausdrücke enthalten. Jeder mehrdeutige reguläre Ausdruck kann zwar durch den in [BEGO71] vorgestellten Algorithmus in einen eindeutigen regulären Ausdruck umgewandelt werden, jedoch ergibt sich daraus eine maximale Laufzeit von  $O(2^n)$ . Algorithmen mit solchen Laufzeiten wären in der Praxis nicht einsetzbar.

Darüber hinaus ist die Eindeutigkeit, wie sie in [BEGO71] beschrieben wird, so definiert, dass es für einen eindeutigen regulären Ausdruck  $r$  für ein beliebiges Wort  $w$  der Sprache nur genau eine Ableitung geben darf. So ist der reguläre Ausdruck  $r_1 = (a|b)^*aa^*$  mehrdeutig. Eine Indizierung eines regulären Ausdrucks bedeutet, dass jedem Zeichen des Ausdrucks ein Index beigelegt wird. Die Indizierung beginnt bei 1 und wird bei jedem Auftreten des Zeichens um eins erhöht. Wird  $r_1$  indiziert, so wird er zu  $(a_1|b_1)^*a_2a_3^*$ . Daraus lassen sich für das Wort  $aaa$  folgende drei Ableitungen erzeugen  $a_1a_1a_2$ ,  $a_1a_2a_3$  und  $a_2a_3a_3$ .

Für einen deterministischen Kellerautomaten, der als Basis für Parser genutzt wird, ist es wichtig, dass er anhand der nächsten  $k$  Zeichen entscheiden kann, welche Überführung angewendet werden muss. In der Praxis wird normalerweise der Fall  $k = 1$  gewählt.

Betrachten wir den regulären Ausdruck  $r = (a|b)^*a$ . Der Ausdruck  $r$  ist eindeutig. Wird jedoch der Ausdruck indiziert  $(a_1|b_1)^*a_2$ , so ist es für ein Wort  $w = baaa, w \in L(r)$  nach Einlesen des Buchstabens  $b$  nicht klar, ob nach  $a_1$  abgeleitet werden muss oder nach  $a_2$ .

Für die Inhaltsmodelle einiger XML-artiger Grammatiken sind nur solche regulären Ausdrücke erlaubt, bei denen anhand des nächsten Eingabezeichens (1-lookahead) entscheidbar ist, welcher Teil des regulären Ausdrucks benutzt werden muss, um das Eingabewort zu erzeugen. Diese regulären Ausdrücke wurden von Brüggemann-Klein definiert und untersucht [BK93b] [BKW92] [BK93a] [BKW98]. Sie nennen sich 1-eindeutige reguläre Ausdrücke. Sie sind eine echte Teilmenge der Menge der regulären Ausdrücke.

In einer Arbeit von 2006 [HW07] wird die 1-Eindeutigkeit erweitert, mit einem neuen Namen versehen und zu  $k$ -lookahead Determinismus und  $k$ -block Determinismus. Bei der  $k$ -lookahead Variante werden die nächsten  $k$  Zeichen gelesen, jedoch nur das nächste Zeichen verarbeitet, wohingegen bei der anderen Variante  $k$  Zeichen gelesen und alle verarbeitet werden.

Es wird gezeigt, dass eine echte Hierarchie innerhalb der  $k$ -block deterministischen Sprachen existiert und dass die  $k$ -block deterministischen Sprachen eine echte Teilmenge der  $k$ -lookahead deterministischen Sprachen sind. Es fehlt der Beweis, dass die  $k$ -lookahead Determinismus Hierarchie echt ist.

Im formalsprachlichen Bereich gibt es verschiedene Forschungsansätze, die sich mit den Beschreibungsmodellen von XML befassen. Die Beschreibungsmodelle werden als Grammatik aufgefasst. Im Wesentlichen gibt es dabei zwei verschiedene Ansatzpunkte:

- Beschreibungsmodelle werden als Zeichengrammatiken aufgefasst.
- Beschreibungsmodelle werden als Baumgrammatiken aufgefasst.

Die erste Strukturierung aus den bestehenden Beschreibungsmodellen in Grammatiken wurde in [MLM01b] und [LMM00] getroffen. Hierbei wurden die bestehenden Beschreibungsmodelle in vier verschiedene Baumgrammatiktypen eingeteilt: Local Baumgrammatiken, Single-Type Baumgrammatiken, Restrained Competition Baumgrammatiken und reguläre Baumgrammatiken.

Die Grammatiken bestehen dabei aus einem 4-Tupel  $G = (N, T, S, P)$ . Die Produktionen haben die Form

$$X \longrightarrow a(m), m \in \text{Reg}(N), a \in T, X \in N.$$

Die Produktionen werden hierbei als Baumproduktionen aufgefasst, wobei  $a$  ein Knoten ist,  $(m)$  die Nachfolgebäume sind und  $\text{Reg}(N)$  ein regulärer Ausdruck über Nichtterminalzeichen ist.

Es wird eine Hierarchie zwischen den Grammatiken gezeigt, wobei die regulären Baumgrammatiken die ausdrucksmächtigsten Grammatiken sind und die Local Baumgrammatiken die ausdruckschwächsten.

Desweiteren werden Algorithmen zum Parsen der entstehenden Modelle erzeugt. Dabei werden für die drei schwächeren Grammatiken deterministische Baumautomaten verwendet. Für die regulären Baumgrammatiken wird ein nichtdeterministischer Baumautomat angegeben.

Nähere Analysen von Baumgrammatiken werden in [CDGJ97] durchgeführt. Es werden nur die regulären Baumgrammatiken behandelt. Hier werden Abchlusseigenschaften, Entscheidbarkeitsfragen, Automatenmodelle und weiterführende Ergebnisse der regulären Baumgrammatiken präsentiert.

Es gibt in [MLM01b] und [LMM00] jedoch Unterschiede zwischen ihrem Modell und dem Modell von [CDGJ97]. In [MLM01b] und [LMM00] dürfen die regulären Ausdrücke der Inhaltsmodelle einen Kleenestern enthalten. Die regulären Baumgrammatiken sind zwar unter Iteration abgeschlossen, jedoch nur unter geschachtelter Iteration und nicht unter einer konkatenierten Iteration.

In einer weiteren Arbeit [LC00] werden sechs Beschreibungsmodelle auf ihre Zugehörigkeit zu den Grammatiktypen überprüft.

Murata [MLM01b] beschreibt einen Grammatiktyp, der den Baumgrammatiken ähnelt. Die Grammatiken heißen Hedge Grammatiken. Dabei muss eine Sprache nicht genau einen Wurzelknoten besitzen, sondern kann auch aus mehreren Bäumen bestehen.

Es werden sowohl Hedge Automaten vorgestellt, als auch die booleschen Abschlusseigenschaften ihrer Sprache. So sind sie sowohl unter Durchschnitt als auch unter Vereinigung und Komplement abgeschlossen.

In der Praxis finden die Hedge Grammatiken jedoch keine Anwendung. In der XML Spezifikation ist eindeutig definiert, dass der Anfang eines XML Dokumentes aus einem Starttag und das Ende aus dem dazugehörigen Endtag besteht [WWW]. Genauso sind HTML Dokumente aufgebaut, die in neuen Versionen eine Dokumentklasse von XML sind. Mit dieser Struktur wird ein Baum beschrieben. In der Spezifikation von XML wird dies als Wohlgeformtheit bezeichnet.

Ein wichtiger Forschungsbereich, mit dem sich diese Arbeit vor allem beschäftigt, sind XML Beschreibungsmodelle aus Sicht von Wortgrammatiken. Den Anstoß geben zwei Arbeiten von Boasson und Berstel [BB02b] [BB00], in denen ein Grammatiktyp namens XML Grammatik definiert wird. Diese Art von Grammatik ist vergleichbar mit DTDs, die durch die Formalisierung abstrahiert werden.

XML Grammatiken besitzen Produktionen mit einem Nichtterminalzeichen auf der linken Seite und einem Tagpaar auf der rechten Seite, das einen regulären Ausdruck über Nichtterminalen einschließt:

$$X \longrightarrow xReg(N)\bar{x}.$$

Die Produktionen haben die spezielle Form, dass zu jedem Terminalzeichen  $x \in T, T = A \cup \bar{A}$ , genau ein Nichtterminalzeichen existiert.

Für diese Form der Grammatiken wurde gezeigt, dass es zu jeder XML Sprache eine eindeutige minimale Grammatik gibt. Weiter wurde gezeigt, dass sie unter Durchschnitt, nicht aber unter Vereinigung und somit auch nicht unter Komplement abgeschlossen sind.

Eine weitere Arbeit befasst sich mit einer Verallgemeinerung der XML Grammatiken. Die Grammatiken heißen balancierte Grammatiken [BB02b]. Sie unterliegen nicht mehr der Einschränkung, dass jedem Terminalzeichen  $x \in T$  genau ein Nichtterminalzeichen zugeordnet sein muss.

Wichtige Erkenntnisse ergeben sich durch den Beweis, dass es zu jeder balancierten Grammatik eine eindeutige Minimalgrammatik gibt. Die balancierten Sprachen sind unter allen booleschen Operationen abgeschlossen.

Darüber hinaus ist es entscheidbar, ob eine balancierte Grammatik eine reguläre bzw. eine kontextfreie Sprache erzeugt. Eine wichtige Eigenschaft, auch in Bezug auf die folgende Arbeit, ist, dass es für jede balancierte Grammatik eine codeterministische äquivalente Grammatik gibt. In älterer Literatur wird diese Eigenschaft als *backwards deterministic* bezeichnet. Die Eigenschaft legt zugrunde, dass es keine zwei Produktionen,  $X \rightarrow m$  und  $X' \rightarrow m$ , in  $G$  gibt, ohne  $X = X'$  zur Folge zu haben.

Es gibt einige ältere Ansätze, die sich mit ähnlichen Grammatiktypen befassen:

- Parenthesis Grammatiken
- Bracketed Grammatiken
- Chomsky-Schützenberger Grammatiken

Die Parenthesis Grammatiken wurden von Knuth [Knu67] und McNaughton [McN67] untersucht. Parenthesis Grammatiken sind Grammatiken  $G = (N, T, S, P)$ , wobei jede Produktion die Form

$$X \rightarrow am\bar{a}$$

hat. Das Terminalalphabet  $T$  besteht aus einem Tagpaar  $\{a, \bar{a}\}$  und einer zu  $T$  disjunkten Menge von Zeichen  $B$ , also  $T = \{a, \bar{a}\} \cup B$ . Die Inhaltsmodelle  $m$  in den Produktionen bestehen aus Nichtterminalen und Zeichen aus  $B$ ,  $m \in (N \cup B)^*$ .

Es existiert nur ein Klammertyp und die Inhaltsmodelle sind endlich.

Ist  $B = \emptyset$ , wird die Grammatik als „strenge“ bezeichnet.

Die Bracketed Grammatiken wurden von Ginsburg und Harrison [GH67] erforscht. Die Terminalzeichen bestehen aus der Menge der Zeichen  $T = A \cup \bar{B} \cup C$ . Die einzelnen Mengen sind disjunkt. Jede Produktion hat die Form

$$X \rightarrow am\bar{b}, a \in A, \bar{b} \in \bar{B}.$$

Für das Inhaltsmodell  $m$  gilt  $m \in (N \cup C)^*$ . Es existiert eine Bijektion zwischen den Terminalen in  $A$  und den Produktionen. So kann in jedem Ableitungsschritt eindeutig erkannt werden, welche Produktion benutzt werden

muss.

Der Unterschied zu den XML Grammatiken besteht darin, dass die Inhaltsmodelle endlich sind.

Die Chomsky-Schützenberger Grammatiken werden in dem Chomsky-Schützenberger Theorem benutzt. Hierbei besteht das Terminalalphabet aus den Zeichen  $T = A \cup \bar{A} \cup B, m \in N^*$ . Die Produktionen haben die Form

$$X \longrightarrow am\bar{a}.$$

Die Beziehung zwischen den Zeichen  $a \in A$  und den Produktionen ist bijektiv.

### 1.3 Fragestellung und Aufbau

Diese Arbeit beschäftigt sich in zwei verschiedenen Richtungen mit XML-artigen Grammatiken.

Der erste Ansatz ist rein theoretisch und beschäftigt sich mit den formalsprachlichen Eigenschaften der Beschreibungsmodelle von XML. Schwerpunkt bilden die genaue Einordnung in die Chomsky Hierarchie, die Abschlusseigenschaften und die Entscheidbarkeitsfragen.

Die von Berstel und Boasson definierten Grammatiken werden, durch die in [MLM01b] gezeigten Grammatiken, ergänzt und bilden mit einem neuen Typ von Grammatiken, den RegDyck Grammatiken, die Basis der Arbeit:

- XML Grammatiken
- einfache balancierte Grammatiken
- konkurrenzverhindernde balancierte Grammatiken
- RegDyck Grammatiken
- balancierte Grammatiken.

Der zweite Ansatz dient dazu, einen praktischen Nutzen aus diesen Erkenntnissen zu ziehen. Es werden Parsergenerierungsalgorithmen und Parser konstruiert und deren Laufzeit bestimmt.

Aus diesen beiden Ansätzen ergibt sich die zentrale Fragestellung: Wie können die formalsprachlichen Eigenschaften der Beschreibungsmodelle von XML praktisch genutzt werden?

Im Einzelnen ergeben sich dadurch folgende Fragen:

- Gibt es neue, verbesserte XML-artige Grammatiken?
- Zu welchem speziellen Typ der Chomsky Hierarchie gehören die Familien der XML-artigen Sprachen?

- Gibt es eine Hierarchie zwischen den XML-artigen Sprachfamilien?
- Wie können Parser generiert werden?
- Mit welchen Methoden kann geparst werden?
- Welche Laufzeiten haben Parser und Parsergenerator?
- Welche Abschlusseigenschaften besitzen die Sprachen der Grammatiken?

Zuerst werden im zweiten Kapitel die grundlegenden Definitionen der Grammatiken beschrieben. Hierbei werden die bestehenden Definitionen der XML Grammatiken, der balancierten Grammatiken und die Inhaltsmodelleigenschaft des 1-lookahead Determinismus erläutert. Anschließend werden die bereits bestehenden Definitionen der Single-Type Baumgrammatiken und der Restrained Competition Baumgrammatiken in eine Definition der Wortgrammatiken übersetzt.

In Kapitel 3 wird ein neuer Grammatiktyp „RegDyck Grammatiken“ hergeleitet. Sie beschreiben den Durchschnitt zwischen regulären Sprachen und den Dyck Sprachen. Es wird dazu auch ein erkennender Automat angegeben.

Kapitel 4 beschäftigt sich mit der Einordnung der XML-artigen Grammatiken. Zuerst wird überprüft, in welche Klasse der Chomsky Hierarchie die Grammatiken gehören. Danach wird eine Hierarchie innerhalb der XML-artigen Grammatiken gezeigt. Der letzte Teil des Kapitels beschäftigt sich mit Parsergeneratoren und Parsingmethoden. Es werden Laufzeitanalysen durchgeführt und neue Methoden zur Parserkonstruktion und zum Parsen erzeugt. Die Grundlagen hierfür legen die bis dahin erörterten theoretischen Grundlagen.

Im letzten Kapitel werden die Abschlusseigenschaften der Grammatiken erforscht: Vereinigung, Durchschnitt, erweitertes Komplement, erweiterte Konkatenation, erweiterte Iteration, erweiterte Substitution und erweiterter Homomorphismus. Danach folgt die Klärung von Entscheidbarkeitsfragen.

# Kapitel 2

## Grundlagen

### 2.1 Allgemeine Definitionen

Die Menge der natürlichen Zahlen wird mit  $\mathbb{N} = \{1, 2, 3, \dots\}$  bezeichnet. Die Vereinigung  $\mathbb{N} \cup \{0\}$  wird mit  $\mathbb{N}_0$  gekennzeichnet. Für die Potenzmenge einer endlichen Menge  $M$  wird  $2^M$  geschrieben. Eine endliche, nichtleere Menge  $A$  von Zeichen wird als Alphabet bezeichnet. Ein Wort ist eine endliche Folge von Zeichen über einem Alphabet  $A$ .  $A^+$  ist die Menge der nichtleeren Worte über  $A$ . Die Vereinigung  $A^+ \cup \lambda$  wird durch  $A^*$  dargestellt. Die leere Menge wird mit  $\emptyset$  bezeichnet. Es gilt  $A^0 = \{\lambda\}$ .

Die Länge eines Wortes  $w \in A^*$  wird abkürzend als  $|w|$  geschrieben, wobei gilt  $|\lambda| = 0$ . Die Spiegelung  $u_n \dots u_2 u_1$  eines Wortes  $w = u_1 \dots u_n$  wird durch  $w^R$  gekennzeichnet.

Es sei  $A$  ein Alphabet, dann ist  $\bar{A}$  die zu  $A$  disjunkte Menge von Zeichen, für die gilt, wenn  $a \in A$ , dann ist  $\bar{a} \in \bar{A}$  und wenn  $a \in A, b \in A, a \neq b$ , dann ist  $\bar{a} \neq \bar{b}$ .

Ein Dyckprime ist ein Wort  $x\bar{x}$  mit  $x \in A, \bar{x} \in \bar{A}$  oder ein Wort  $x d_1 \dots d_n \bar{x}$ ,  $n \in \mathbb{N}$ , wobei  $d_1, \dots, d_n$  Dyckprimes sind. Es sei  $T$  ein Alphabet mit  $T = A \cup \bar{A}$ , dann ist  $D_T$  die Menge aller Dyckprimes über diesem Alphabet.

Ein Dyckwort über dem Alphabet  $T$  hat die Form  $v_1 \dots v_n$ , wobei  $v_1, \dots, v_n \in D_T, n \in \mathbb{N}$ . Die Menge aller Dyckwörter über dem Alphabet  $T$  wird mit  $W_T$  bezeichnet.

**Definition 2.1.1** Eine reguläre Grammatik ist eine Grammatik  $G = (N, T, S, P)$ , wobei gilt:

- $N$  ist die Menge der Nichtterminalzeichen.

- $T$  ist die Menge der Terminalzeichen mit  $T \cap N = \emptyset$ .
- $S$  ist das Startsymbol.
- $P$  ist die Menge der Produktionen der Form

$$\begin{aligned} X &\longrightarrow Yz \text{ und} \\ X &\longrightarrow z \end{aligned}$$

mit  $X, Y \in N, z \in T^*$ .

□

$L(G) = \{w \mid w \in T^* \wedge S \Rightarrow^* w\}$  ist die von der Grammatik  $G$  erzeugte reguläre Sprache.

**Definition 2.1.2** (Pumping Lemma für reguläre Sprachen)

Es sei  $L$  eine reguläre Sprache. Dann existiert eine natürliche Zahl  $n$ , so dass sich alle Wörter  $x \in L$  mit  $|x| \geq n$  in  $x = uvw$  zerlegen lassen, so dass folgende drei Bedingungen erfüllt sind:

1.  $|v| \geq 1$
2.  $|uv| \leq n$
3.  $uw^i w \in L, \forall i \in \mathbb{N}_0$ .

□

**Definition 2.1.3** Eine kontextfreie Grammatik ist eine Grammatik  $G = (N, T, S, P)$ , wobei gilt:

- $N$  ist die Menge der Nichtterminalzeichen.
- $T$  ist die Menge der Terminalzeichen  $N \cap T = \emptyset$ .
- $S$  ist das Startsymbol.
- $P$  ist die Menge der Produktionen der Form

$$X \longrightarrow \alpha$$

mit  $X \in N, \alpha \in (N \cup T)^*$ .

□

**Definition 2.1.4** (selbsteinbettende Nichtterminalzeichen)

Es sei  $G = (N, T, S, P)$  eine kontextfreie Grammatik. Ein Nichtterminalzeichen  $X \in N$  wird als selbsteinbettend bezeichnet, wenn eine Ableitung  $X \Rightarrow^* zXz_0, z, z_0 \in T^+$  existiert.

□

Die Ableitung eines Wortes  $w \in L(G)$ , wird bildlich durch einen Parsebaum beschrieben.

**Definition 2.1.5** (Parsebaum)

Es sei  $G = (N, T, S, P)$  eine kontextfreie Grammatik. Ein Parsebaum (Ableitungsbaum) ist ein Baum mit folgenden Eigenschaften:

- Die Wurzel wird mit dem Startsymbol  $S$  markiert.
- Jedes Blatt ist mit einem Zeichen  $a \in T$  oder  $\lambda$  gekennzeichnet.
- Jeder innere Knoten ist mit einem Nichtterminalzeichen  $X \in N$  gekennzeichnet.
- Es sei  $X \in N$  die Beschriftung eines inneren Knoten und  $X \rightarrow X_0 \dots X_n$  die ausgewählte Produktion, dann sind  $X_0, \dots, X_n \in (N \cup T)$  die Nachfolger des Knotens. Wenn  $X \rightarrow \lambda$  die gewählte Produktion ist, dann ist  $\lambda$  der einzige Nachfolger.

□

**Definition 2.1.6** Ein nichtdeterministischer endlicher Automat (NFA)  $\mathcal{R} = (S, A, s_0, \delta, F)$  ist ein System mit:

- $S$  ist die Menge der Zustände.
- $A$  ist die Menge der Eingabesymbole.
- $s_0$  ist der Startzustand.
- $F \subseteq S$  ist die Menge der Endzustände.
- $\delta : S \times A \rightarrow 2^S$  ist die Überföhrungsfunktion.

Wir nehmen an, die Überföhrungsfunktion eines NFA sei total. Falls für eine Situation kein Folgezustand definiert ist, dann bildet der NFA in die leere Menge ab.

□

Die Erweiterung der Funktion  $\delta$  für Wörter ist:  $\Delta : S \times A^* \rightarrow 2^S$ .

Eine Konfiguration  $(s, w)$  von  $\mathcal{R}$  ist ein Paar  $S \times A^*$ . Die Konfiguration  $(s, w)$  bedeutet, dass sich  $\mathcal{R}$  im Zustand  $s$  befindet und  $w$  die noch nicht gelesene Eingabe ist.

$L(\mathcal{R}) = \{w \in A^* \mid \Delta(s, w) \in F\}$  ist die von dem nichtdeterministischen endlichen Automat  $\mathcal{R}$  erkannte Sprache.

**Definition 2.1.7** Ein nichtdeterministischer Kellerautomat  $\mathcal{K} = (S, A, G, \Delta, s_0, \delta)$  ist ein System mit:

- $S$  ist die Menge der Zustände.

- $A$  ist die Menge der Eingabesymbole.
- $G$  ist die Menge der Kellersymbole.
- $\Delta \in G$  ist das Kellerendesymbol.
- $s_0$  ist der Startzustand.
- Die Überföhrungsfunktion  $\delta$  bildet von  $S \times (A \cup \{\lambda\}) \times G$  in die endlichen Teilmengen von  $S \times G^*$  ab.

Wir nehmen an, der Kellerautomat akzeptiert mit leerem Keller. □

Eine Konfiguration  $(s, w, z)$  von  $\mathcal{K}$  ist ein Tripel aus  $S \times A^* \times G^*$ , wobei sich  $\mathcal{K}$  im Zustand  $s$  befindet,  $w$  die noch nicht gelesene Eingabe und  $z$  der Kellerinhalt ist.

$L(\mathcal{K}) = \{w \in A^* \mid (s_0, w, \Delta) \vdash^* (s', \lambda, \lambda) \wedge s' \in S\}$  ist die von dem Kellerautomat  $\mathcal{K}$  erkannte Sprache.

**Definition 2.1.8** Ein deterministischer Kellerautomat  $\mathcal{K} = (S, A, G, \Delta, s_0, \delta, F)$  ist ein System mit:

- $S$  ist die Menge der Zustände.
- $A$  ist die Menge der Eingabesymbole.
- $G$  ist die Menge der Kellersymbole.
- $\Delta \in G$  ist das Kellerendesymbol.
- $s_0$  ist der Startzustand.
- $F \subseteq S$  ist die Menge der Endzustände.
- Für jedes  $s \in S$  und  $g \in G$  ist  $\delta(s, a, g)$  für alle  $a \in A$  nicht definiert, wenn  $\delta(s, \lambda, g)$  definiert ist.  
Für alle  $s \in S, g \in G$  und  $a \in (A \cup \{\lambda\})$  enthält  $\delta(s, a, g)$  höchstens ein Element.

□

**Definition 2.1.9** Es sei  $A$  ein Alphabet, das keines der Zeichen aus der Menge  $V = \{[, *, +, ?, \emptyset, (, )\}$  enthält. Reguläre Ausdröcke über einem Alphabet  $A$  werden wie folgt festgelegt:

- Alle  $p \in A$  sowie  $\emptyset$  sind reguläre Ausdröcke.
- Wenn  $p$  und  $q$  reguläre Ausdröcke sind, dann ist auch die Konkatenation  $(pq)$  ein regulärer Ausdruck.

- Wenn  $p$  und  $q$  reguläre Ausdrücke sind, dann ist auch die Vereinigung  $(p|q)$  ein regulärer Ausdruck.
- Wenn  $p$  ein regulärer Ausdruck ist, dann ist auch  $(p)^*$  ein regulärer Ausdruck.

□

Die Klammern können weggelassen werden, wenn sich dadurch die Semantik des regulären Ausdrucks nicht ändert.

Im Folgenden gilt:  $(p)^+ = (p)(p)^*$ ,  $(p)? = (p|\emptyset)$  und  $\lambda = \emptyset^*$ . Die Zeichen  $^*, +, ?, |$  werden als Operatoren bezeichnet.

Jeder reguläre Ausdruck  $r$  über einem Alphabet  $A$  legt eine Sprache  $L(r)$  folgendermaßen fest:

- Die durch  $\emptyset$  bezeichnete Sprache ist die leere Sprache.
- Die durch  $a \in A$  bezeichnete Sprache enthält nur das Wort  $a$ .
- Für reguläre Ausdrücke  $p$  und  $q$  über dem Alphabet  $A$  gilt:  $L(p|q) = L(p) \cup L(q)$ ,  $L(pq) = L(p)L(q)$ ,  $L(p^*) = (L(p))^*$ .

Es sei  $A$  ein Alphabet, dann ist  $Reg(A)$  die Menge der regulären Ausdrücke über diesem Alphabet.

Ein regulärer Ausdruck  $r$  über einem Alphabet  $A$  wird zu einem indizierten regulären Ausdruck  $r'$ , indem jedes Vorkommen eines Zeichens  $a \in A$  mit einem Index  $i \in \mathbb{N}$  versehen wird, so dass die Kombination aus dem Zeichen  $a$  und dem Index  $i$  eindeutig in dem regulären Ausdruck  $r'$  ist. Die Indizierung beginnt bei 1. Es sei  $r = uvuuwvww$ ,  $u, v, w \in A$  ein regulärer Ausdruck, dann ist der indizierte reguläre Ausdruck  $r' = u_1v_1u_2u_3w_1v_2w_2w_3w_4$ .

$\pi$  ist eine injektive Abbildung, die für den gegebenen regulären Ausdruck  $r$  einen indizierten regulären Ausdruck  $r'$  zurückgibt.

### Definition 2.1.10

Sei  $r$  ein regulärer Ausdruck über dem Alphabet  $A$  und sei  $V = \{*, |, (, ), \cdot\}$  die Menge der Operatoren. Die folgenden Abbildungen fassen den regulären Ausdruck als Wort über dem Alphabet  $A$  und den Operatoren  $V$  auf.

$$A' = \{a_i \mid a \in A, i \in \mathbb{N}\}$$

ist die Menge der indizierten Zeichen. Die folgenden Abbildungen

$$\gamma : (A \cup V)^* \times A \longrightarrow \mathbb{N}$$

$$\gamma(w, a) = |w|_a \text{ in } w$$

$$\psi : A \times \mathbb{N} \longrightarrow A'$$

$$\psi(a, i) = a_i$$

dienen dazu, die Anzahl des Vorkommens eines Buchstabens in einem Wort zu zählen und einem Buchstaben einen Index zuzuweisen.

Mit diesen Abbildungen können wir eine Abbildung  $\pi$  definieren, die zu einem regulären Ausdruck  $r$ , welcher als Wort aufgefasst wird, ein indiziertes Wort liefert, welches dann wieder als indizierter regulärer Ausdruck aufgefasst wird.

$$\begin{aligned} \pi : (A \cup V)^* &\longrightarrow (A' \cup V)^* \\ \pi(\epsilon) &= \epsilon \\ x \in (A \cup V)^*, a \in A : \pi(xa) &= \pi(x)\psi(a, \gamma(xa, a)) \\ x \in (A \cup V)^*, y \in V : \pi(xy) &= \pi(x)y \end{aligned}$$

□

**Beispiel 2.1.4** Indizierter regulärer Ausdruck

$$\begin{aligned} r &= (a|b)^*a(a|b) \\ \pi(r) = r' &= (a_1|b_1)^*a_2(a_3|b_2) \end{aligned}$$

□

Es sei  $\pi^{-1}$  die Abbildung, welche die Indizierung eines regulären Ausdrucks aufhebt:  $\pi^{-1}(r') = r$ . Ebenso hebt  $\pi^{-1}$  die Indizierung eines Wortes auf:  $\pi^{-1}(w') = w, w' \in L(\pi(r)), w \in L(r)$ .

Ein regulärer Ausdruck  $r$  heißt eindeutig, wenn gilt:

$$\forall w' \in L(\pi(r)) : \neg \exists w \in L(\pi(r)) : (\pi^{-1}(w) = \pi^{-1}(w') \wedge w \neq w').$$

## 2.2 k-lookahead deterministische reguläre Ausdrücke

Im Folgenden werden k-lookahead deterministische reguläre Ausdrücke definiert. Sie sind wichtig für die Untertypen von XML-artigen Grammatiken, die in dieser Arbeit untersucht werden. Eine XML-artige Grammatik  $G = (N, T, S, P)$  besitzt Produktionen der Form

$$X \longrightarrow xm\bar{x},$$

wobei  $m$  ein regulärer Ausdruck über Nichtterminalen ist und  $X \in N, x, \bar{x} \in T$ .

In einigen XML-artigen Grammatiktypen ist nicht jeder beliebige reguläre Ausdruck erlaubt. Eine allgemeine Einschränkung der regulären Ausdrücke ist, dass sie eindeutig sein müssen. Die Eigenschaft, dass jeder reguläre Ausdruck in einen eindeutigen regulären Ausdruck in maximal  $O(2^n)$  Zeit umgewandelt werden kann, ist allgemein bekannt, siehe [BEGO71]. Ein effektives Parsen verlangt Grammatiken, aus denen möglichst schnell ein XML Parser

konstruiert werden kann.

Eine präzisere Einschränkung der Inhaltsmodelle erfolgt durch die Einschränkung des  $k$ -Determinismus. Dieser hat zwei verschiedene Ausprägungen: Den  $k$ -lookahead Determinismus, welcher sich an die  $LL(1)$  Grammatiken im kontextfreien Bereich anlehnt und den  $k$ -block Determinismus [HW07]. In der Praxis sowie in dieser Arbeit wird nur der 1-lookahead Determinismus verwendet.

**Definition 2.2.11** Es sei  $A$  ein Alphabet und  $r$  ein regulärer Ausdruck über  $A$ .  $\sigma$  ist eine Projektion auf das erste Zeichen eines Wortes  $w \in A^*$ . Es sei  $w = u_1 \dots u_n$  mit  $u_1, \dots, u_n \in A$ , dann ist  $\sigma(w) = u_1$ . Es gilt  $\sigma(\lambda) = \lambda$ . Das Alphabet, das durch die Indizierung des regulären Ausdrucks  $r$  entsteht, wird mit  $A'$  bezeichnet.

Ein regulärer Ausdruck ist  $k$ -lookahead deterministisch,  $k \in \mathbb{N}$ , wenn gilt:

$$\begin{aligned} \forall u, v, w, x, y \in A'^* : (uxv \in L(\pi(r)) \wedge uyw \in L(\pi(r)) \wedge |x| = |y| = k) \\ \implies ((\pi^{-1}(x) \neq \pi^{-1}(y)) \vee (\pi^{-1}(x) = \pi^{-1}(y) \wedge \sigma(x) = \sigma(y))). \end{aligned}$$

□

Nicht jede reguläre Sprache kann durch einen  $k$ -lookahead deterministischen regulären Ausdruck beschrieben werden [HW07]. So ist die Sprache

$$L((a|b)^* a (a|b)^k), k \geq 1$$

nicht durch einen  $k$ -lookahead regulären Ausdruck erzeugbar.

## 2.3 Balancierte Grammatiken

Erweiterte kontextfreie Grammatiken sind Grammatiken, bei denen die rechte Seite der Produktionen aus einem regulären Ausdruck über Nichtterminalzeichen und Terminalzeichen besteht. Balancierte Grammatiken sind spezielle erweiterte kontextfreie Grammatiken. Die rechten Seiten der Produktionen bestehen aus einem Start- und einem Endtag und dazwischen befindet sich ein regulärer Ausdruck über Nichtterminalzeichen.

Die praktische Bedeutung dieser Grammatiken findet sich in den Beschreibungsmodellen von XML wieder. Wie in Kapitel 4 gezeigt wird, sind die balancierten Grammatiken die Obermenge der XML-artigen Grammatiken und wurden in [BB02b] definiert.

Ähnliche Grammatiken sind die Parenthesis Grammars, beschrieben von Knuth [Knu67] und McNaughton [McN67]. Sie sind eine Untermenge der balancierten Grammatiken und besitzen eine endliche Menge von Produktionen und genau ein Klammerpaar.

Eine weitere Verbindung kann zu den regulären Baumgrammatiken [CDGJ97] gezogen werden. In dieser Sichtweise werden die Wörter nicht mehr als Klammerpaar gesehen, sondern in eine Baumstruktur aufgelöst. Hierbei gibt es Ähnlichkeiten zu regulären Baumgrammatiken. Lee, Mani und Murata schlagen diese Variante zum Parsen vor [LMM00] [MLM01b].

Aus formalsprachlicher Sicht sind die balancierten Grammatiken sehr interessant, denn sie besitzen sehr viele positive Abschlusseigenschaften und viele Entscheidbarkeitsfragen sind für sie entscheidbar.

**Definition 2.3.12** Eine balancierte Grammatik ist ein 4-Tupel  $G = (N, T, S, P)$ , wobei gilt:

- $N$  ist eine endliche Menge von Nichtterminalsymbolen.
- $T$  ist eine endliche Menge von Terminalsymbolen mit

$$T = A \cup \bar{A} \cup B \text{ und } \bar{A} = \{\bar{a} \mid a \in A\}.$$

$A, \bar{A}$  und  $B$  sind disjunkt.

- $S$  ist die Menge der Startsymbole.
- $P$  ist die Menge von Produktionen der Form

$$X \longrightarrow am\bar{a},$$

mit  $m \in \text{Reg}(N \cup B)$ . Der Mittelteil  $m$  der Produktion wird als Inhaltsmodell bezeichnet. Für ein festes  $X \in N$  und ein  $a \in A$  gibt es genau eine Produktion

$$X \longrightarrow am_{X,a}\bar{a}. \quad \square$$

Das erste Terminalzeichen  $a \in A$  auf der rechten Seite einer Produktion wird als Farbe der Produktion bezeichnet.

Wenn  $B = \emptyset$  ist, dann ist die balancierte Grammatik „streng“ (engl. pure). Die entstehenden Wörter sind wohlgeformt und Dyckwörter. Wenn  $B \neq \emptyset$ , dann werden die entstehenden Wörter als Motzkinwörter [ST04] bezeichnet.

In dieser Arbeit werden ausschließlich „strenge“ balancierte Grammatiken behandelt. Der Zusatz „streng“ wird im Folgenden weggelassen.

**Definition 2.3.13** Eine Sprache ist eine balancierte Sprache, wenn sie durch eine balancierte Grammatik erzeugt werden kann. □

Es folgen ein paar Beispiele von balancierten Grammatiken und balancierten Sprachen:

**Beispiel 2.3.5** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit folgenden Produktionen:

$$\begin{aligned} S &\longrightarrow xY^*\bar{x} \\ Y &\longrightarrow y\bar{y} \end{aligned}$$

Die Grammatik ist balanciert und erzeugt die Sprache  $L = \{x(y\bar{y})^i\bar{x} \mid i \in \mathbb{N}_0\}$ .

□

**Beispiel 2.3.6** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit folgenden Produktionen:

$$\begin{aligned} S &\longrightarrow xY\bar{x} \\ Y &\longrightarrow y(S|A)\bar{y} \\ A &\longrightarrow a\bar{a} \end{aligned}$$

Die Grammatik ist balanciert und erzeugt durch ihren selbsteinbettenden Teil eine kontextfreie Sprache  $L = \{(xy)^n a\bar{a}(\bar{y}\bar{x})^n \mid n \geq 1\}$ .

□

**Beispiel 2.3.7**

$$L = \{x(y\bar{y})^n(z\bar{z})^n\bar{x} \mid n \geq 1\}$$

Behauptung: Die Sprache  $L$  ist nicht durch eine balancierte Grammatik erzeugbar.

Nehmen wir das Gegenteil an, dann muss es eine Ableitung

$$S \Longrightarrow xm\bar{x}, m \in \text{Reg}(N)$$

geben, aus welcher folgt

$$xm\bar{x} \Longrightarrow^* x(y\bar{y})^n(z\bar{z})^n\bar{x}.$$

Das Inhaltsmodell  $m$  hat die Gestalt

$$m = X_1X_2 \dots X_nY_1Y_2 \dots Y_n.$$

Jeder Dyckprime, ohne Beachtung der Operatoren, wird durch ein eigenes Nichtterminalzeichen erzeugt. Das Inhaltsmodell  $m$  ist somit  $X^nY^n$ . Es ist jedoch nicht durch einen regulären Ausdruck erzeugbar.

*Beweis.* Der Beweis folgt mittels Pumping Lemma. Die Pumping Lemma Konstante ist  $n$ . Wir wählen das Wort  $X^n Y^n$  und die Länge  $2n$ . Aufgrund der 1. Bedingung des Pumping Lemmas darf  $v$  nicht leer sein, aufgrund der 2. Bedingung kann es nur aus  $Y$  bestehen.

Die 3. Bedingung lässt zu, dass auch die Wörter  $X^{n-|v|} Y^n$  zulässig sind. Die Definition der Sprache erlaubt dies allerdings nicht. Der Ausdruck sowie das Inhaltsmodell sind somit nicht regulär. Die Sprache kann also nicht durch eine balancierte Grammatik erzeugt werden.  $\square$

**Beispiel 2.3.8** Es sei  $G = (N, T, S, P)$  eine kontextfreie Grammatik mit folgenden Produktionen:

$$\begin{aligned} S &\longrightarrow xY\bar{x} \\ Y &\longrightarrow y\bar{y}(Y|\lambda)z\bar{z} \end{aligned}$$

Die durch  $G$  erzeugte Sprache ist nur eine balancierte Sprache, wenn  $y = z$ . Dann würde die Sprache durch die balancierte Grammatik

$$\begin{aligned} S &\longrightarrow x(YY)^*\bar{x} \\ Y &\longrightarrow yy \end{aligned}$$

erzeugt werden.

Es würde andernfalls die Sprache  $L = \{x(y\bar{y})^n(z\bar{z})^n\bar{x} \mid n \geq 1\}$  erzeugt werden, die durch obigen Beweis ausgeschlossen wurde.  $\square$

Nun werden aus den balancierten Grammatiken, Grammatiken erzeugt, deren Produktionen indiziert sind.

### **Algorithmus:**

Es sei  $G = (\bar{N}, T, S', P)$  eine balancierte Grammatik. Aus  $G$  entsteht folgendermaßen eine indizierte balancierte Grammatik  $G'$ . Die Konstruktion der indizierten balancierten Grammatik erfolgt über die Konstruktion eines Baumes.

1. Das Startsymbol  $S$  ist der Wurzelknoten des Baumes. Wenn eine Grammatik mehrere Startsymbole besitzt, dann wird für jedes Startsymbol ein eigener Baum konstruiert.
2. Die Nachfolger von  $S$  im Baum sind alle Produktionen von  $S$ . Die rechten Seiten der Produktionen werden indiziert. Jedes Nichtterminalzeichen  $X$  in jedem Inhaltsmodell der Produktionen wird mit einem Index versehen, so dass die Kombination aus Nichtterminalzeichen und Index in der indizierten balancierten Grammatik eindeutig ist. Der Index beginnt bei 1. Somit entsteht aus jedem Inhaltsmodell  $m$  ein indiziertes Inhaltsmodell  $m'$ .

3. Danach wird für jedes indizierte Nichtterminalzeichen  $X_i$  und der Produktion  $X \rightarrow xm\bar{x}$  eine Produktion erzeugt und die Nichtterminalzeichen im Inhaltsmodell wieder indiziert:  $X_i \rightarrow xm''\bar{x}$ .
4. Nun gilt wieder, dass die Nachfolger eines indizierten Nichtterminalzeichens  $Y_i$  alle Produktionen von  $Y_i$  sind.
5. Die Schritte drei und vier werden solange wiederholt bis folgender Zustand eintritt: Es sei  $\psi$  eine Abbildung, die für ein indiziertes Nichtterminalzeichen  $Y_i$  die Indizierung löscht ( $Y$ ). Falls für ein Nichtterminalzeichen  $Y_i$  ein Nachfolger konstruiert werden soll und ein Vorgänger  $Y_j$  im Baum ( $Y_j \rightarrow ym\bar{y}$ ) existiert, so dass  $\psi(Y_i) = \psi(Y_j)$  gilt, dann wird eine Produktion  $Y_i \rightarrow ym\bar{y}$  konstruiert, für die keine weiteren Nachfolger konstruiert werden. Dies ist der Fall, wenn ein Nichtterminalzeichen selbsteinbettend auftritt.
6. Das Verfahren bricht ab, wenn keine Nachfolger mehr erzeugt werden müssen.

Die Produktionen der indizierten balancierten Grammatik werden aus den Knoten des Baums abgelesen. Daraus ergeben sich auch gleichzeitig die Nichtterminalzeichen, die Terminalzeichen und die Startsymbole.

**Definition 2.3.14**

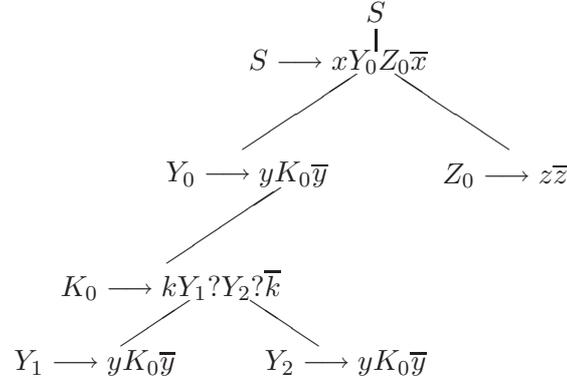
Eine indizierte balancierte Grammatik ist eine balancierte Grammatik, deren Produktionen durch den obigen Algorithmus indiziert wurden.

□

**Beispiel 2.3.9** Gegeben sei die balancierte Grammatik  $G = (N, T, \{S\}, P)$  mit den Produktionen:

$$\begin{array}{lcl}
 S & \longrightarrow & xYZ\bar{x} \\
 Y & \longrightarrow & yK\bar{y} \\
 K & \longrightarrow & kY?Y?\bar{k} \\
 Z & \longrightarrow & z\bar{z}
 \end{array}$$

Für die Grammatik  $G$  wird folgende indizierte Grammatik  $G'$  konstruiert:



□

## 2.4 Untertypen der balancierten Grammatiken

Im Folgenden werden eingeschränkte balancierte Grammatiken definiert. Die Einschränkungen beziehen sich vor allem darauf, dass die entstehenden Grammatiken zu den LL(1) Grammatiken gehören, wie im vierten Kapitel gezeigt wird. Die balancierten Grammatiken können Mehrdeutigkeiten in ihren Inhaltsmodellen enthalten. So stellt sich zu allererst die Frage, ob sie zu den deterministisch kontextfreien Grammatiken gehören oder nichtdeterministisch kontextfrei sind. Falls sie deterministisch kontextfrei sind, stellt sich weiter die Frage, mit welchem Zeitaufwand ein deterministischer Kellerautomat aus einer Grammatik konstruiert werden kann.

In den Arbeiten [LMM00] und [MLM01b] wurden XML Baumgrammatiken definiert. Nun werden auf ähnliche Weise Wortgrammatiken definiert.

**Definition 2.4.15** Es sei  $G = (N, T, S, P)$  mit  $T = A \cup \bar{A}$  eine balancierte Grammatik. Zwei Nichtterminale  $X$  und  $Y$  konkurrieren miteinander, wenn gilt:

$$\exists X, Y \in N \exists m_0, m_1 \in \text{Reg}(N) \exists \bar{x} \in \bar{A} :$$

$$\exists x \in A : (X \longrightarrow xm_0\bar{x} \in P \wedge Y \longrightarrow xm_1\bar{x} \in P \wedge X \neq Y)$$

□

Konkurrierende Nichtterminale erschweren das Parsen, da nicht mehr eindeutig anhand eines Terminals die Produktion hergeleitet werden kann. Somit müssen unter Umständen mehrere Produktionen auf ihre Gültigkeit hin überprüft werden, was zum Nichtdeterminismus und damit zu einer größeren Laufzeit führt.

**Beispiel 2.4.10** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, Y, Z\}$
- $T = \{x, \bar{x}, y, \bar{y}\}$
- $P = \left\{ \begin{array}{l} S \longrightarrow xY\bar{x}, \\ Y \longrightarrow yZ\bar{y}, \\ Z \longrightarrow y\bar{y} \end{array} \right\}$

Die Grammatik beinhaltet zwei Nichtterminale mit den gleichen Terminalen auf der rechten Seite der Produktion. Die Nichtterminale  $Y$  und  $Z$  konkurrieren miteinander.

□

**Definition 2.4.16** Eine XML Grammatik ist eine balancierte Grammatik ohne konkurrierende Nichtterminale und mit nur einem Startsymbol.

Für eine XML Grammatik  $G = (N, T, S, P)$  mit  $T = A \cup \bar{A}$  und dem Startsymbol  $S$  muss gelten:

Jede Produktion hat die Form:

$$X_a \longrightarrow am\bar{a}, \quad m \in \text{Reg}(N), \quad a \in A$$

und

$$\forall X, Y \in N \quad \forall m, m_0 \in \text{Reg}(N) \quad \forall \bar{x} \in \bar{A} :$$

$$\neg \exists x \in A : (X \longrightarrow xm\bar{x} \in P \wedge Y \longrightarrow xm_0\bar{x} \in P \wedge X \neq Y)$$

Zusätzlich besteht jedes Inhaltsmodell aus einem 1-lookahead deterministischen regulären Ausdruck.

□

Das bedeutet, es besteht eine bijektive Beziehung zwischen den Nichtterminalen  $N$  und den Terminalen  $A$ .

**Definition 2.4.17** Eine balancierte Sprache ist eine XML Sprache, wenn sie durch eine XML Grammatik erzeugt werden kann.

□

**Beispiel 2.4.11** Gegeben sei die Sprache

$$L = \{xyc\bar{c}y\bar{y}d\bar{d}y\bar{x}\}.$$

Sie ist keine XML Sprache. Nehmen wir an, sie wäre eine XML Sprache, dann müsste eine XML Grammatik existieren, die  $L$  erzeugt.

Um  $L$  zu beschreiben wird eine balancierte Grammatik  $G = (N, T, S, P)$  mit folgenden Produktionen benötigt:

$$P = \left\{ \begin{array}{l} S \longrightarrow xYZ\bar{x}, \\ Y \longrightarrow yC\bar{y}, \\ Z \longrightarrow yD\bar{y}, \\ C \longrightarrow c\bar{c}, \\ D \longrightarrow d\bar{d} \end{array} \right\}$$

Da die Sprache innerhalb der beiden  $y \dots \bar{y}$  Klammerpaare verschieden ist und deshalb zwei verschiedene Nichtterminalzeichen  $Y, Z$  mit den gleichen Terminalen auf der rechten Seite entstehen, gibt es zwei konkurrierende Nichtterminalzeichen.

Dies ist in einer XML Grammatik nicht erlaubt und die erzeugte Sprache stellt keine XML Sprache dar. □

**Definition 2.4.18** Eine einfache balancierte Grammatik ist eine balancierte Grammatik  $G = (N, T, S, P)$  mit folgenden Eigenschaften:

$\forall X, Y, Z \in N \forall z, \bar{z}, \bar{x} \in T \forall m, m_0, m_1 \in \text{Reg}(N), \alpha, \beta, \gamma, \omega \in N^* :$

$$\neg \exists x \in A : (X \longrightarrow xm\bar{x} \in P \wedge Y \longrightarrow xm_0\bar{x} \in P \wedge X \neq Y \wedge Z \longrightarrow zm_1\bar{z} \in P \wedge \alpha X \beta \in L(m_1) \wedge \gamma Y \omega \in L(m_1))$$

$S$  ist das einzige Startsymbol.

Zusätzlich besteht jedes Inhaltsmodell aus einem 1-lookahead deterministischen regulären Ausdruck. □

Es existieren also keine Produktionen, deren Inhaltsmodelle konkurrierende Nichtterminale enthalten und es existiert nur ein Startsymbol.

**Definition 2.4.19** Eine balancierte Sprache ist eine einfache balancierte Sprache, wenn sie durch eine einfache balancierte Grammatik erzeugt werden kann. □

**Beispiel 2.4.12** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, Y, Z\}$
- $T = \{x, \bar{x}, y, \bar{y}\}$
- $P = \left\{ \begin{array}{l} S \longrightarrow xYZ\bar{x}, \\ Y \longrightarrow y\bar{y}, \\ Z \longrightarrow y\bar{y} \end{array} \right\}$

Die Grammatik beinhaltet zwei konkurrierende Nichtterminale  $Y$  und  $Z$ , welche in einem Inhaltsmodell enthalten sind. Dies ist in einer einfachen balancierten Grammatik nicht erlaubt. □

**Definition 2.4.20** Eine konkurrenzverhindernde balancierte Grammatik ist eine balancierte Grammatik  $G = (N, T, S, P)$  mit folgenden Einschränkungen:

$\forall m, m_0, m_1 \in \text{Reg}(N) \forall X, Y, Z \in N \forall z, \bar{z}, \bar{x} \in T \forall U, V, W \in N^* :$

$$\neg \exists x \in A : (X \longrightarrow xm\bar{x} \in P \wedge Y \longrightarrow xm_0\bar{x} \in P \wedge X \neq Y \wedge \\ Z \longrightarrow zm_1\bar{z} \in P \wedge UXW \in L(m_1) \wedge UYV \in L(m_1))$$

$S$  ist das einzige Startsymbol.

Zusätzlich besteht jedes Inhaltsmodell aus einem 1-lookahead deterministischen regulären Ausdruck. □

**Beispiel 2.4.13** Es sei  $G = (N, T, S, P)$  eine balancierte Grammatik mit

- $N = \{X, Y, S\}$
- $T = \{x, \bar{x}, z, \bar{z}\}$
- $P = \{ \begin{array}{l} S \longrightarrow zX^*Y^*\bar{z} \\ X \longrightarrow x\bar{x} \\ Y \longrightarrow x\bar{x} \end{array} \}.$

Die Grammatik beinhaltet zwei Nichtterminale  $X$  und  $Y$  mit dem gleichen Terminal auf der rechten Seite der Produktion. Wir wählen  $U$ ,  $V$  und  $W$  leer. Somit kann sowohl  $UXW$ , als auch  $UYV$  entstehen. Dies ist in einer konkurrenzverhindernden balancierten Grammatik nicht erlaubt. □

## Kapitel 3

# Durchschnitt von regulären Sprachen und Dyckprimes

In der Praxis wurden die balancierten Grammatiken in verschiedene Untertypen unterteilt: XML Grammatiken, einfache balancierte Grammatiken und konkurrenzverhindernde balancierte Grammatiken. Die Einschränkungen werden, wie sich im nächsten Kapitel herausstellen wird, aus Gründen der besseren Parsebarkeit getroffen. Sie haben jedoch den Nachteil, dass bestimmte Sprachen nicht erzeugt werden können und einige Abschlusseigenschaften verloren gehen.

Es wird nun nach einem Untertyp gesucht, welcher einige Sprachen erzeugen kann, die die anderen Untertypen nicht erzeugen können und viele positive Abschlusseigenschaften besitzt.

Der Durchschnitt einer regulären Sprache mit der Menge  $D_T$  von Dyckprimes über einem Alphabet  $T$ , einer sogenannten RegDyck Sprache, scheint aus mehreren Gründen interessant zu sein. Er ist auf jeden Fall deterministisch kontextfrei und jede RegDyck Sprache ist auf einen regulären Ausdruck und eine Menge  $D_T$  von Dyckprimes zurückzuführen. Dabei wird ausgenutzt, dass der reguläre Ausdruck auch ein mehrdeutiger regulärer Ausdruck sein darf und für ihn ein NFA konstruiert werden kann, der in  $O(n)$  Zeit parst.

### 3.1 Konstruktion von RegDyck Automaten

Der Durchschnitt einer regulären Sprache und der Menge  $D_T$  von Dyckprimes muss in den deterministisch kontextfreien Sprachen liegen, da der Durchschnitt einer regulären Sprache und einer deterministisch kontextfreien Sprache deterministisch kontextfrei ist.

Es wird daher überprüft, wie der Durchschnitt dieser beiden Sprachen aussieht und ob er möglicherweise ein Untertyp der balancierten Sprachen ist. Zuerst wird ein RegDyck Automat aus einem beliebigen endlichen Automa-

ten und einem beliebigen Dyckprime Automaten konstruiert. Der RegDyck Automat hat die Besonderheit, dass seine Zustandsüberführungen zwar nicht-deterministisch sind, aber der Kellerinhalt immer eindeutig ist.

Es sei  $\mathcal{R} = (S_R, A_R, s_{0R}, \delta_R, F_R)$  ein nichtdeterministischer endlicher Automat:

- $S_R = \{s_{0R}, s_1, s_2, \dots, s_n\}$  ist die Menge der Zustände.
- $A_R = A \cup \bar{A} \cup B \cup \bar{B}$  ist das Eingabealphabet.
- $s_{0R}$  ist der Startzustand.
- $F_R \subseteq S_R$  ist die Menge der Endzustände.
- $\delta_R$  ist die Überföhrungsfunktion.

Es sei  $\mathcal{D} = (S_D, A_D, G_D, \Delta, s_{0D}, \delta_D, F_D)$  ein deterministischer Kellerautomat, der die Menge  $D_{A_D}$  von Dyckprimes erkennt:

- $S_D = \{s_{0D}, s_{1D}, s_{eD}\}$  ist die Menge der Zustände des Automaten.
- $A_D = A \cup \bar{A} \cup C \cup \bar{C}$  ist das Eingabealphabet.
- $G_D = \bar{A} \cup \bar{C} \cup \{\Delta\}$  ist das Kelleralphabet.
- $s_{0D}$  ist der Startzustand.
- $F_D = \{s_{0D}, s_{eD}\}$  ist die Menge der Endzustände.
- $\delta_D$  ist die Überföhrungsfunktion:

$$\begin{aligned} (s_{0D}, a, \Delta) &\longrightarrow (s_{1D}, \bar{a}\Delta) \quad \forall a \in (A \cup C) \\ (s_{1D}, a, \bar{z}) &\longrightarrow (s_{1D}, \bar{a}\bar{z}) \quad \forall a \in (A \cup C) \quad \forall \bar{z} \in (\bar{A} \cup \bar{C}) \\ (s_{1D}, \lambda, \Delta) &\longrightarrow (s_{eD}, \Delta) \\ (s_{1D}, \bar{a}, \bar{a}) &\longrightarrow (s_{1D}, \lambda) \quad \forall \bar{a} \in (\bar{A} \cup \bar{C}). \end{aligned}$$

Nun wird ein nichtdeterministischer Kellerautomat  $\mathcal{K} = (S_K, A_K, G_K, \Delta, s_{0K}, \delta_K)$  (der RegDyck Automat) konstruiert, der den Durchschnitt der Sprache des endlichen Automaten  $\mathcal{R}$  und der Sprache des deterministischen Kellerautomaten  $\mathcal{D}$  erkennt:

- $S_K = S_R \cup \{s_{eK}, s_{0K}, s_{1K}, s_E\}$  ist die Menge der Zustände.
- $A_K = A_R \cap A_D = A \cup \bar{A}$  ist das Eingabealphabet.
- $G_K = \bar{A} \cup \{\diamond, \Delta\}$  ist das Kelleralphabet.
- $s_{0K}$  ist der Startzustand.
- Es wird eine Startüberföhrung zur Initialisierung benötigt:

$$(s_{0K}, \lambda, \Delta) \longrightarrow (s_{1K}, \diamond \Delta).$$

Der RegDyck Automat kellert ein Zeichen  $\diamond$  ein. Danach arbeitet er wie der nichtdeterministische endliche Automat  $\mathcal{R}$  und der deterministische Kellerautomat  $\mathcal{D}$ . Zum Schluss löscht er das Zeichen  $\diamond$ , wenn sowohl  $\mathcal{R}$  als auch  $\mathcal{D}$  akzeptieren.

Für jede Überführung in  $\mathcal{R}$  der Form  $(s_{0R}, a) \longrightarrow s_j$  mit  $s_j \in S_R$  und  $a \in A$  wird eine Überführung

$$(s_{1K}, a, \diamond) \longrightarrow (s_j, \bar{a}\diamond)$$

erzeugt. Dieser Zustand  $s_{1K}$  wird nur einmal benutzt und ist danach nicht mehr in dem Wertebereich der Überführungsfunktion enthalten. In allen von  $\mathcal{K}$  übernommenen Zuständen des ursprünglichen Automaten  $\mathcal{R}$  wird, wenn das Kellersymbol  $\diamond$  an der obersten Stelle des Kellers steht, in einen Endzustand oder einen Fehlerzustand verwiesen. Das Zeichen  $\diamond$  zeigt damit an, dass ein Dyckprime abgearbeitet worden ist.

Für jede Überführung in  $\mathcal{R}$  der Form  $(s_i, a) \longrightarrow s_j$  mit  $s_i, s_j \in S_R$  und  $a \in A$  werden Überführungen

$$(s_i, a, \bar{z}) \longrightarrow (s_j, \bar{a}\bar{z}) \quad \forall \bar{z} \in \bar{A}, \bar{a} \in \bar{A}$$

erzeugt. Es wird festgelegt, dass für ein Eingabezeichen  $a \in A$  der Automat  $\mathcal{K}$  definiert ist, egal welches Zeichen  $\bar{z} \in \bar{A}$  auf dem Keller liegt. Die einzige Bedingung ist, dass der endliche Automat  $\mathcal{R}$  eine Zustandsüberführung im momentanen Zustand für das Eingabezeichen  $a \in A$  besitzt.

Der Automat  $\mathcal{D}$  ist für ein Zeichen  $a \in A$  immer definiert, außer wenn er sich in dem Zustand  $s_{eD}$  oder  $s_E$  befindet. Er kellert das entsprechende Gegenstück aus der Menge  $\bar{A}$  ein. Da der Automat  $\mathcal{K}$  immer den gleichen Kellerinhalt wie  $\mathcal{D}$  haben soll, muss auch  $\mathcal{K}$  immer das entsprechende Gegenstück  $\bar{a} \in \bar{A}$  zu  $a \in A$  einkellern.

Für jede Überführung in  $\mathcal{R}$  der Form  $(s_i, \bar{a}) \longrightarrow s_j$  mit  $s_i, s_j \in S_R$  und  $\bar{a} \in \bar{A}$  wird eine Überführung

$$(s_i, \bar{a}, \bar{a}) \longrightarrow (s_j, \lambda)$$

erzeugt. Der Automat  $\mathcal{K}$  ist somit für ein Zeichen  $\bar{a} \in \bar{A}$  definiert, sofern der Automat  $\mathcal{R}$  eine Überführung aus dem Zustand  $s_i$  für dieses Zeichen besitzt. Zusätzlich muss sich noch der Kellerinhalt verändern. Für jedes Eingabezeichen  $\bar{a} \in \bar{A}$  löscht der Automat  $\mathcal{D}$  das Zeichen  $\bar{a}$

aus dem Keller, sofern  $\bar{a}$  an der obersten Stelle im Keller liegt. Diese Überführung ist hiermit auch im Automaten  $\mathcal{K}$  konstruiert.

Für jede Überführung in  $\mathcal{R}$  der Form  $(s_i, \bar{a}) \longrightarrow s_j$  mit  $s_i \in S_R, s_j \in S_R \setminus F_R$  und  $\bar{a} \in \bar{A}$  wird eine Überführung

$$(s_j, \lambda, \diamond) \longrightarrow (s_E, \diamond)$$

erzeugt. Der Zustand  $s_E$  ist ein Fehlerzustand. Dieser Zustand wird erreicht, wenn ein Dyckprime abgearbeitet ist, jedoch kein Endzustand des endlichen Automaten erreicht ist. Das bedeutet, dass zwar der Automat  $\mathcal{D}$  akzeptiert, nicht aber der Automat  $\mathcal{R}$ .

Es werden zusätzlich Überführungen für alle Endzustände  $s_f \in F_R$  erzeugt:

$$(s_f, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda).$$

Die Konstruktion beruht darauf, dass der Kellerautomat  $\mathcal{D}$  und der Kellerautomat  $\mathcal{K}$  zu jedem Zeitpunkt den gleichen Kellerinhalt haben. Weiter haben der Kellerautomat  $\mathcal{K}$  und der endliche Automat  $\mathcal{R}$  den gleichen Zustand nach dem Bearbeiten eines Eingabezeichens.

Der folgende Beweis teilt sich in die Betrachtung der Kellerinhalte von  $\mathcal{K}$  und  $\mathcal{D}$  und der Zustände von  $\mathcal{K}$  und  $\mathcal{R}$  auf.

**Satz 3.1.1** Seien  $\mathcal{R}, \mathcal{D}$  und  $\mathcal{K}$  die Automaten aus der obigen Konstruktion, dann erkennt  $\mathcal{K}$  den Durchschnitt von  $\mathcal{R}$  und  $\mathcal{D}$ .

*Beweis.* Der erste Abschnitt des Beweises befasst sich mit den Kellern der Automaten  $\mathcal{K}$  und  $\mathcal{D}$ .

Es sei  $\mathcal{D}$  in der Konfiguration  $(s_m, g, \Delta)$  und  $\mathcal{K}$  in der Konfiguration  $(s_i, g, \diamond \Delta)$ . Dann folgt nach  $|g|$  Schritten, dass  $\mathcal{D}$  in der Konfiguration  $(s_l, \lambda, \bar{v} \Delta)$  und  $\mathcal{K}$  in der Konfiguration  $(s_j, \lambda, \bar{v} \diamond \Delta)$  ist. Hierbei sind  $s_m, s_l \in S_D, s_i \in S_K \setminus \{s_{0K}\}, s_j \in S_K \setminus \{s_{0K}, s_{1K}\}, g \in (A \cup \bar{A})^*$  und  $\bar{v} \in \bar{A}^*$ .

Der Beweis folgt per Induktion über die Länge eines Teilwortes  $g'' \in (A \cup \bar{A})^*$  der Eingabe  $g = g''g' \in (A \cup \bar{A})^*$  mit  $g' \in (A \cup \bar{A})^*$ .

Sei  $g'' = \lambda$ . Dann existiert für den Kellerautomaten  $\mathcal{D}$  keine Überführung für  $s_m \in S_D \setminus \{s_{1D}\}$ .  $\mathcal{D}$  bleibt also in der ursprünglichen Konfiguration  $(s_m, g', \Delta)$ . Ist  $s_m = s_{1D}$ , so wird die Überführung

$$(s_{1D}, \lambda, \Delta) \longrightarrow (s_{eD}, \Delta)$$

angewendet.  $\mathcal{D}$  befindet sich dann in der Konfiguration  $(s_{eD}, g', \Delta)$ .

Der Kellerautomat  $\mathcal{K}$  verwendet für  $s_i \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}\}$  entweder die Überführung

$$(s_i, \lambda, \diamond) \longrightarrow (s_E, \diamond)$$

oder die Überführung

$$(s_i, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda).$$

$\mathcal{K}$  befindet sich dann in der Konfiguration  $(s_E, g', \diamond\Delta)$  oder  $(s_{eK}, g', \Delta)$ .

Der Fall  $s_i = s_{0K}$  kann nach Konstruktion nicht auftreten. In den Fällen  $s_i \in \{s_{1K}, s_E, s_{eK}\}$  ist der Automat für die Eingabe  $\lambda$  nicht definiert und bleibt in der Konfiguration  $(s_i, g', \diamond\Delta)$ .

Also steht in den Kellern der Automaten  $\mathcal{K}$  und  $\mathcal{D}$  (evtl. bis auf das Zeichen  $\diamond$  im Keller von  $\mathcal{K}$ ) dasselbe Wort.

Sei  $g'' = a$ . Dann verwendet  $\mathcal{D}$  für  $s_m = s_{0D}$  die Überführung

$$(s_{0D}, a, \Delta) \longrightarrow (s_{1D}, \bar{a}\Delta).$$

$\mathcal{D}$  befindet sich nun in der Konfiguration  $(s_{1D}, g', \bar{a}\Delta)$ .

Ist  $s_m = s_{1D}$ , so ist der Automat  $\mathcal{D}$  nicht definiert und bleibt in der Konfiguration  $(s_{1D}, a, \Delta)$ .

Ist der Automat  $\mathcal{K}$  in der Konfiguration  $(s_{1K}, ag', \diamond\Delta)$ , so verwendet er die Überführung

$$(s_{1K}, a, \diamond) \longrightarrow (s_w, \bar{a}\diamond)$$

und befindet sich anschließend in der Konfiguration  $(s_w, g', \bar{a}\diamond\Delta)$ , sofern  $\mathcal{R}$  auch für  $a$  definiert ist.

Ist  $s_i \neq s_{1K}$ , so ist der Automat  $\mathcal{K}$  nicht definiert und bleibt in der Konfiguration  $(s_i, ag', \diamond\Delta)$ .

Somit haben beide Automaten denselben Kellerinhalt (mit Ausnahme des Zeichens  $\diamond$  im Keller von  $\mathcal{K}$ ), da nur zwei Fälle eintreten können: Entweder ist  $\mathcal{D}$  im Zustand  $s_{0D}$  und  $\mathcal{K}$  im Zustand  $s_{1K}$  und beide Automaten sind definiert oder die Automaten befinden sich nicht in ihren Anfangszuständen und sind beide nicht definiert.

Sei  $g'' = \bar{a}$ . Dann ist für den Automaten  $\mathcal{D}$  in keinem Zustand eine Überführung mit leerem Keller für diese Eingabe definiert. Ebenso ist für den Automaten  $\mathcal{K}$  in keinem Zustand eine Überführung mit leerem Keller für eine solche Eingabe definiert. Folglich bleiben beide Automaten in ihrer Konfiguration und haben somit denselben Kellerinhalt.

Die Behauptung gelte nun für  $|g''| \leq n \in \mathbb{N}_0$ .

Es sei  $ag'$  oder  $\bar{a}g'$  der Rest eines Eingabewortes  $g = g''g'$  mit  $g'' = wa$  bzw.  $g'' = w\bar{a}$ ,  $|w| = n$ ,  $a \in A$ ,  $\bar{a} \in \bar{A}$ .

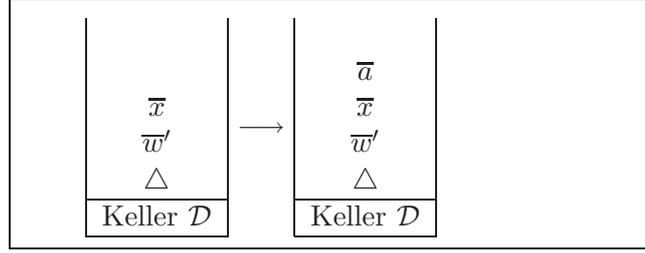


Abbildung 3.1: Kellerinhalt des Kellerautomaten  $\mathcal{D}$  nach Abarbeitung des Eingabezeichens  $a$ .

Im Folgenden sei  $\bar{w}' \in \bar{A}^*$ ,  $\bar{x} \in \bar{A}$  mit  $\bar{w} = \bar{x}\bar{w}'$ , wobei  $\bar{w}$  das nach Abarbeitung von  $w$  in den Kellern stehende Wort ist.

**1. Fall:**  $g' = wa$

Nach Induktionsvoraussetzung und Definition ist der Kellerautomat  $\mathcal{D}$  nach Abarbeitung von  $w$  in der Konfiguration  $(s_{1D}, ag', \bar{x}\bar{w}'\Delta)$  und der Automat  $\mathcal{K}$  in der Konfiguration  $(s_j, ag', \bar{x}\bar{w}' \diamond \Delta)$ , wobei  $s_j \in S_K \setminus \{s_{0K}, s_{1K}\}$  (nach Konstruktion). Der Kellerautomat  $\mathcal{D}$  benutzt die Überführung

$$(s_{1D}, a, \bar{x}) \longrightarrow (s_{1D}, \bar{a}\bar{x})$$

und wechselt in die Konfiguration

$$(s_{1D}, g', \bar{a}\bar{x}\bar{w}'\Delta). \text{ (Abbildung 3.1)}$$

Der Kellerautomat  $\mathcal{K}$  benutzt für  $s_j \in S_K \setminus \{s_{0K}, s_{1K}\}$  die Überführung

$$(s_j, a, \bar{x}) \longrightarrow (s_k, \bar{a}\bar{x})$$

und befindet sich danach in der Konfiguration

$$(s_k, g', \bar{a}\bar{x}\bar{w}' \diamond \Delta). \text{ (Abbildung 3.2)}$$

Die Fälle  $s_j \in \{s_{0K}, s_{1K}\}$  können nach Konstruktion des Automaten nicht auftreten.

Somit haben die beiden Automaten immer denselben Kellerinhalt (bis auf das Zeichen  $\diamond$  im Keller von  $\mathcal{K}$ ). Falls  $\mathcal{R}$  für eine Eingabe nicht definiert ist, dann gehört die Sprache aber auch nicht zum Durchschnitt der Sprachen von  $\mathcal{R}$  und  $\mathcal{D}$ .

Sind die Keller leer, so sind beide Kellerautomaten nicht definiert, da sie nach Konstruktion nicht in ihre Anfangszustände  $s_{0D}$  und  $s_{0K}$  bzw.  $s_{1K}$  zurückkehren können. Nur für diese Zustände sind Überführungen mit leerem Keller und nicht leerer Eingabe definiert. Beide Automaten bleiben in der ursprünglichen Konfiguration und haben somit denselben Kellerinhalt.

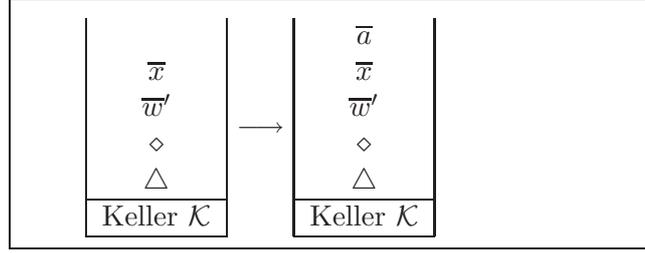


Abbildung 3.2: Kellerinhalt des Kellerautomaten  $\mathcal{K}$  nach Abarbeitung des Eingabezeichens  $a$ .

**2. Fall:**  $g'' = w\bar{a}$

Nach Induktionsvoraussetzung und Definition von  $\mathcal{D}$  wird nach Abarbeitung von  $w$  ein Eingabezeichen  $\bar{a}$  eingelesen und im Kellerautomaten  $\mathcal{D}$  besteht eine Konfiguration  $(s_{1D}, \bar{a}g', \bar{x}\bar{w}'\Delta)$ . Falls  $\bar{x} = \bar{a}$ , so verwendet der Kellerautomat  $\mathcal{D}$  folgende Überführung:

$$(s_{1D}, \bar{a}, \bar{a}) \longrightarrow (s_{1D}, \lambda)$$

Es folgt die Konfiguration

$$(s_{1D}, g', \bar{w}'\Delta). \quad (\text{Abbildung 3.3})$$

Ist  $\bar{x} \neq \bar{a}$ , so ist  $\mathcal{D}$  nicht definiert und bleibt in der Konfiguration

$$(s_{1D}, \bar{a}g', \bar{x}\bar{w}'\Delta).$$

Der Kellerautomat  $\mathcal{K}$  ist nach Induktionsvoraussetzung und Definition nach Abarbeitung von  $w$  in der Konfiguration  $(s_j, \bar{a}g', \bar{a}\bar{w}'\diamond\Delta)$  mit  $s_j \in S_K \setminus \{s_{0K}, s_{1K}\}$ . Falls  $\bar{x} = \bar{a}$ , benutzt er die Überführung, wenn  $\mathcal{R}$  auch dafür definiert ist

$$(s_j, \bar{a}, \bar{a}) \longrightarrow (s_k, \lambda)$$

und ist dann in der Konfiguration

$$(s_k, g', \bar{w}'\diamond\Delta). \quad (\text{Abbildung 3.4})$$

Ist  $\bar{x} \neq \bar{a}$ , so ist  $\mathcal{K}$  nicht definiert und bleibt in der Konfiguration

$$(s_j, \bar{a}g', \bar{x}\bar{w}'\diamond\Delta).$$

In beiden Fällen steht im Keller von  $\mathcal{K}$  dasselbe wie im Keller von  $\mathcal{D}$  (bis auf das Zeichen  $\diamond$  im Keller von  $\mathcal{K}$ ).

Die Fälle  $s_j \in \{s_{0K}, s_{1K}\}$  können nach Konstruktion des Automaten  $\mathcal{K}$  nicht auftreten.

Sind die Keller leer, so sind beide Kellerautomaten nicht definiert. Sie bleiben in der ursprünglichen Konfiguration und haben somit denselben Kellerinhalt.

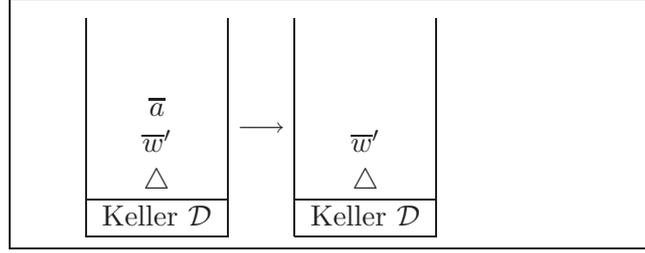


Abbildung 3.3: Kellerinhalt des Kellerautomaten  $\mathcal{D}$  nach Abarbeitung des Eingabezeichens  $\bar{a}$ .

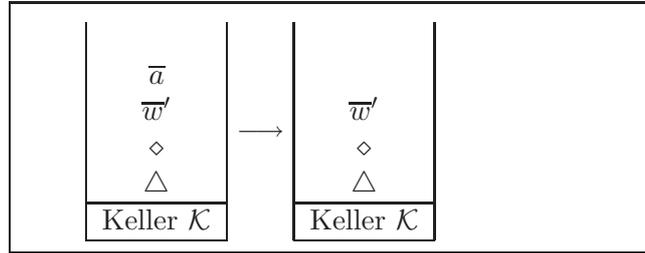


Abbildung 3.4: Kellerinhalt des Kellerautomaten  $\mathcal{K}$  nach Abarbeitung des Eingabezeichens  $\bar{a}$ .

Befindet sich der Automat  $\mathcal{K}$  in der Konfiguration  $(s_i, ag', \bar{w} \diamond \Delta)$ , dann wechselt er in die Konfiguration  $(s_j, g', \bar{a}\bar{w} \diamond \Delta)$ . Der Automat  $\mathcal{D}$  wechselt von der Konfiguration  $(s_{1D}, ag', \bar{w}\Delta)$  in  $(s_{1D}, g', \bar{a}\bar{w}\Delta)$ .

Ist  $\mathcal{K}$  in der Konfiguration  $(s_i, \bar{a}g', \bar{b}\bar{w}' \diamond \Delta)$  und  $\mathcal{D}$  in der Konfiguration  $(s_{1D}, \bar{a}g', \bar{b}\bar{w}'\Delta)$  mit  $\bar{b} \in \bar{A}, \bar{a} \neq \bar{b}$ , so gibt es für beide Automaten keine Überführung die angewendet werden kann. Somit sind beide Automaten in diesem Fall nicht definiert. Für  $\bar{b} = \bar{a}$  wechselt  $\mathcal{D}$  in die Konfiguration  $(s_{1D}, g', \bar{w}\Delta)$  und  $\mathcal{K}$  wechselt in die Konfiguration  $(s_j, g', \bar{w} \diamond \Delta)$ .

Der zweite Teil des Beweises folgt:

Es wird überprüft, ob sich der endliche Automat  $\mathcal{R}$  und der RegDyck Automat  $\mathcal{K}$  immer im gleichen Zustand befinden.

Es sei  $\mathcal{R}$  in der Konfiguration  $(s_t, g)$  und  $\mathcal{K}$  in der Konfiguration  $(s_j, g, \diamond \Delta)$ ,  $s_t \in S_R, s_j \in S_K \setminus \{s_{0K}\}$ . Dann folgt nach  $|g|$  Schritten, dass  $\mathcal{R}$  in der Konfiguration  $(S_i, \lambda), S_i \subseteq S_R$ , und  $\mathcal{K}$  in der Konfiguration  $(S_i, \lambda, \bar{v} \diamond \Delta)$  ist, da  $\mathcal{R}$  nichtdeterministisch sein kann. Hierbei ist  $\bar{v} \in \bar{A}^*$ .

Falls  $\mathcal{K}$  in der Konfiguration  $(s_{1K}, g, \diamond \Delta)$  ist, so ist  $\mathcal{R}$  in der Konfiguration  $(s_{0R}, g)$ . Da für den Zustand  $s_{1K}$  dieselben Überführungen konstruiert werden, wie sie für  $s_{0R}$  vorliegen, werden die beiden Zustände in dieselbe

Folgezustandsmenge  $S_i \subseteq S_R$  überführt. Im Fall  $s_j \neq s_{1K}$  gilt  $s_j = s_t$  oder  $s_j \in \{s_E, s_{eK}\}$ .

Der Beweis folgt per Induktion über die Länge eines Teilwortes  $g'' \in (A \cup \bar{A})^*$  der Eingabe  $g = g''g'$  mit  $g' \in (A \cup \bar{A})^*$ .

Sei  $g'' = \lambda$ . Der Automat  $\mathcal{R}$  ist für diese Eingabe nicht definiert und bleibt in der Konfiguration

$$(s_t, g').$$

Der Kellerautomat  $\mathcal{K}$  verwendet für  $s_j \notin \{s_{0K}, s_{1K}, s_E, s_{eK}\}$  entweder die Überführung

$$(s_j, \lambda, \diamond) \longrightarrow (s_E, \diamond)$$

oder die Überführung

$$(s_j, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda).$$

$\mathcal{K}$  befindet sich dann entweder in der Konfiguration  $(s_E, \lambda, \diamond)$  oder in der Konfiguration  $(s_{eK}, \lambda, \Delta)$ .  $\mathcal{R}$  und  $\mathcal{K}$  befinden sich nun in unterschiedlichen Zuständen. Dies ändert allerdings nichts daran, dass  $\mathcal{K}$  den Durchschnitt der Sprachen der Automaten  $\mathcal{R}$  und  $\mathcal{D}$  erkennt, da dieser Fall nur auftritt, wenn entweder die Eingabe gleich  $\lambda$  ist oder die Eingabe bereits abgearbeitet wurde. In beiden Fällen bleibt für  $\mathcal{R}$  nichts zu tun, allerdings muss  $\mathcal{K}$  noch überprüfen, ob es sich bei der Eingabe um einen Dyckprime handelt. Die Fälle  $s_j \in \{s_{eK}, s_E, s_{0K}\}$  können nach vorheriger Überlegung nicht eintreten. Der Fall  $s_j = s_{1K}$  wurde oben bereits betrachtet.

Die Behauptung gelte nun für  $|g''| \leq n \in \mathbb{N}_0$ .

$\mathcal{R}$  befindet sich nach der Abarbeitung eines Wortes  $w$ , wobei  $g = g''g'$  mit  $g'' = wa$  oder  $g'' = w\bar{a}$ ,  $|w| = n$ , in einer Zustandsmenge  $S_i \subseteq S_R$ . Ebenfalls befindet sich der Kellerautomat  $\mathcal{K}$  in der Zustandsmenge  $S_i$ . Die restliche Eingabe besteht aus einem Wort  $ag'$  oder  $\bar{a}g'$ .

Im Folgenden sei  $\bar{w}' \in \bar{A}^*$ ,  $\bar{x} \in \bar{A}$  mit  $\bar{w} = \bar{x}\bar{w}'$ , wobei  $\bar{w}$  das nach Abarbeitung von  $w$  im Keller von  $\mathcal{K}$  stehende Wort ist.

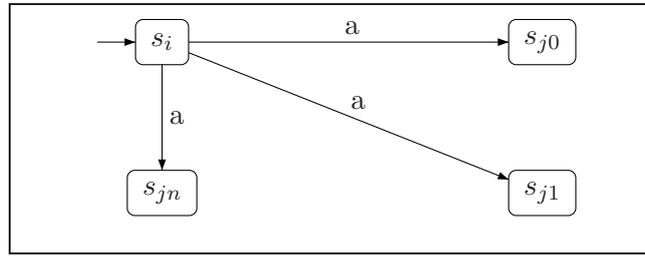
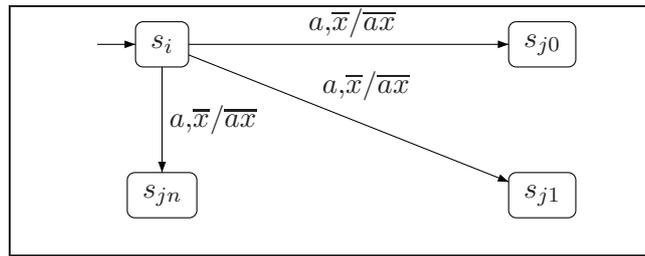
**1. Fall:**  $g'' = wa$

Die Konstruktion beschreibt für jede Überführung

$$(s_m, a) \longrightarrow s_k$$

in  $\mathcal{R}$  mit  $s_m \in S_i, s_k \in S_R$  eine Überführung

$$(s_m, a, \bar{x}) \longrightarrow (s_k, \bar{a}\bar{x})$$


 Abbildung 3.5: Automat  $\mathcal{R}$  im Zustand  $s_i$  mit der Eingabe  $a$ .

 Abbildung 3.6: Automat  $\mathcal{K}$  im Zustand  $s_i$  mit der Eingabe  $a$ . Die Beschriftung einer Kante mit  $a, \bar{x} / a\bar{x}$  bedeutet, dass, wenn der Automat ein Zeichen  $a$  einliest und das Zeichen  $\bar{x}$  auf dem Keller liegt, dann wechselt er in den Zustand zu dem die beschriftete Kante führt und legt  $a\bar{x}$  auf dem Keller ab.

in  $\mathcal{K}$ , welche in den gleichen Zustand überführt. Da dies für alle Keller-symbole  $\bar{x} \in G_K$  gilt, wird die Überführung unabhängig vom obersten Kellerzeichen ausgeführt. Gibt es in  $\mathcal{R}$  mehrere Überführungen der Form  $(s_m, a) \rightarrow s_k$ , so gibt es auch mehrere Überführungen in  $\mathcal{K}$ . Da  $\mathcal{R}$  nicht-deterministisch sein kann, befinden sich sowohl  $\mathcal{R}$  als auch  $\mathcal{K}$  in einer Zustandsmenge  $S_j$  (Abbildung 3.5 und 3.6).

## 2. Fall: $g'' = w\bar{a}$

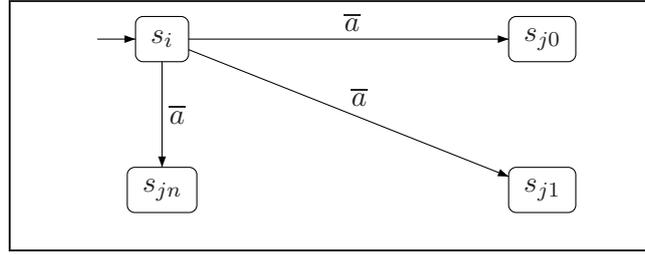
Für eine Überführung in  $\mathcal{R}$  der Form

$$(s_m, \bar{a}) \rightarrow s_k$$

mit  $s_m \in S_i, s_k \in S_R$  hat der Automaten  $\mathcal{K}$  die Überführung

$$(s_m, \bar{a}, \bar{a}) \rightarrow (s_k, \lambda).$$

Hierbei ist es für den Automaten  $\mathcal{K}$  nicht gleichgültig, welches Zeichen auf der obersten Position im Keller liegt. Es muss das passende Zeichen  $\bar{a}$  sein. Im Fall, dass es nicht das passende Zeichen ist, wäre der Automaten  $\mathcal{K}$  in einem anderen Zustand. Er würde nicht akzeptieren, genauso wie der Automaten  $\mathcal{D}$  nicht akzeptieren würde. Somit ist die erkannte Sprache wieder gleich. Ist  $\bar{a}$  an der obersten Stelle des Kellers, befinden sich beide Automaten  $\mathcal{K}$  und  $\mathcal{R}$  in der Zustandsmenge  $S_j$ .


 Abbildung 3.7: Automat  $\mathcal{R}$  im Zustand  $s_i$  mit der Eingabe  $\bar{a}$ .

Es sei  $(s_0, ag', \bar{w} \diamond \Delta)$  die Konfiguration von  $\mathcal{K}$  und  $(s_0, ag')$  die Konfiguration von  $\mathcal{R}$ . Beide Automaten werden in die Zustandsmenge  $S_j$  überführt, da in  $\mathcal{K}$  entsprechende Überführungen erzeugt worden sind. Der Kellerinhalt ist dabei beliebig.

Befinden sich die Automaten in den Konfigurationen  $(s_0, \bar{a}g', \bar{b}\bar{w}' \diamond \Delta)$  und  $(s_0, \bar{a}g')$ , verweist  $\mathcal{R}$  möglicherweise in einen Zustand und  $\mathcal{K}$  akzeptiert die Eingabe nicht, falls  $\bar{a} \neq \bar{b}$ .  $\mathcal{D}$  akzeptiert jedoch dann auch nicht, da  $\mathcal{K}$  und  $\mathcal{D}$  den gleichen Kellerinhalt haben.

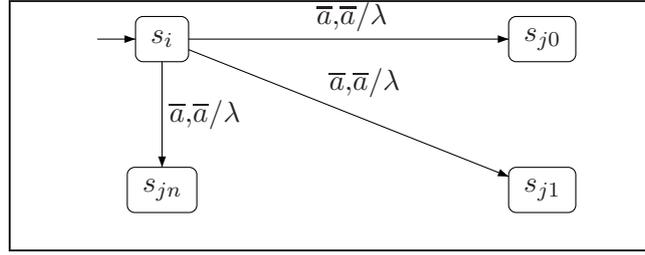
Der Automat  $\mathcal{K}$  akzeptiert eine Eingabe genau dann, wenn sein Keller leer ist, er also in der Konfiguration  $(s_{eK}, \lambda, \Delta)$  ist. Diese Konfiguration erreicht der Automat genau dann, wenn sowohl  $\mathcal{D}$  als auch  $\mathcal{R}$  die Eingabe akzeptieren, also der Kellerautomat  $\mathcal{D}$  einen leeren Keller hat und der endliche Automat  $\mathcal{R}$  sich in einem Endzustand befindet.

Befindet sich der Automat  $\mathcal{R}$  in einem Endzustand, aber der Keller von  $\mathcal{K}$  ist (bis auf das Zeichen  $\diamond$ ) nicht leer, so ist auch der Keller von  $\mathcal{D}$  nicht leer. Folglich wird weder von  $\mathcal{D}$  noch von  $\mathcal{K}$  die Eingabe akzeptiert. Die Eingabe liegt damit auch nicht im Durchschnitt.

Ist der Keller von  $\mathcal{D}$  leer, so ist der Keller von  $\mathcal{K}$  bis auf das Zeichen  $\diamond$  geleert, da sie immer denselben Kellerinhalt haben. Akzeptiert der Automat  $\mathcal{R}$  die Eingabe allerdings nicht, so befindet sich  $\mathcal{R}$  nicht in einem Endzustand, wenn er die Eingabe abgearbeitet hat. Hier existiert in  $\mathcal{K}$  keine Überführung, die das Kellersymbol  $\diamond$  im Keller löscht. Somit akzeptiert  $\mathcal{K}$  die Eingabe ebenfalls nicht. □

Nun wird noch die Gegenrichtung gezeigt. Es wird aus einem RegDyck Automaten  $\mathcal{K}$  ein Kellerautomat  $\mathcal{D}$  und ein endlicher Automat  $\mathcal{R}$  konstruiert. Der Durchschnitt ihrer Sprachen soll gleich der Sprache von  $\mathcal{K}$  sein. Dabei erkennt  $\mathcal{D}$  eine Menge  $D_T$  von Dyckprimes über einem Alphabet  $T$ .

Es sei  $\mathcal{K} = (S_K, A_K, G_K, \Delta, s_{0K}, \delta_K)$  ein RegDyck Automat mit:


 Abbildung 3.8: Automat  $\mathcal{K}$  im Zustand  $s_i$  mit der Eingabe  $\bar{a}$ .

- $S_K = \{s_{0K}, s_{1K}, s_E, s_{eK}, s_1, \dots, s_n\}$  ist die Menge der Zustände.
- $A_K = A \cup \bar{A}$  ist die Menge der Eingabesymbole.
- $G_K = \bar{A} \cup \{\Delta, \diamond\}$  ist die Menge der Kellerzeichen.
- $s_{0K}$  ist der Startzustand.
- $\delta_K$  ist die Überföhrungsfunktion.

Nun werden die beiden Automaten  $\mathcal{R}$  und  $\mathcal{D}$  konstruiert, deren Durchschnitt ihrer Sprachen gleich der Sprache des RegDyck Automaten  $\mathcal{K}$  ist.

Der nichtdeterministische endliche Automat  $\mathcal{R} = (S_R, A_R, s_{1K}, \delta_R, F_R)$  wird folgendermaßen festgelegt:

- $S_R = S_K \setminus \{s_E, s_{0K}, s_{eK}\}$  ist die Menge der Zustände.
- $A_R = A \cup \bar{A}$  ist das Eingabealphabet.
- $s_{1K}$  ist der Startzustand.
- $F_R = \{s_f \in S_K \mid \exists \delta_K(s_f, \lambda) = (s_{eK}, \lambda)\}$  ist die Menge der Endzustände.
- $\delta_R$ :  
Für jede Überföhrung in  $\delta_K$  der Form

$$(s_i, a, \bar{z}) \longrightarrow (s_j, \bar{a}\bar{z}), \quad a \in A, \bar{a}, \bar{z} \in \bar{A}, s_i, s_j \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}\},$$

wird eine Überföhrung

$$(s_i, a) \longrightarrow s_j$$

erzeugt.

Für jede Überföhrung der Form

$$(s_i, \bar{a}, \bar{a}) \longrightarrow (s_j, \lambda), \quad \bar{a} \in \bar{A}, s_i, s_j \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}\},$$

wird eine Überführung

$$(s_i, \bar{a}) \longrightarrow s_j$$

erzeugt.

Es wird außerdem für alle Überführungen der Form

$$(s_{1K}, a, \diamond) \longrightarrow (s_j, \bar{a}\diamond), a \in A, \bar{a} \in \bar{A}, s_j \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}\},$$

eine Überführung

$$(s_{1K}, a) \longrightarrow s_j$$

erzeugt.

Der deterministische Kellerautomat  $\mathcal{D} = (S_D, A_D, G_D, \Delta, s_{0D}, \delta_D, F_D)$ , der die Menge  $D_{A \cup \bar{A}}$  von Dyckprimes erkennt, ist folgendermaßen definiert:

- $S_D = \{s_{0D}, s_{1D}, s_{eD}\}$  ist die Menge der Zustände.
- $A_D = A \cup \bar{A}$  ist das Eingabealphabet.
- $G_D = \bar{A} \cup \{\Delta\}$  ist das Kelleralphabet.
- $s_{0D}$  ist der Startzustand.
- $F_D = \{s_{0D}, s_{eD}\}$  ist die Menge der Endzustände.
- $\delta_D$  ist definiert durch:

$$\begin{aligned} (s_{0D}, a, \Delta) &\longrightarrow (s_{1D}, \bar{a}\Delta) \quad \forall a \in A \\ (s_{1D}, a, \bar{z}) &\longrightarrow (s_{1D}, \bar{a}\bar{z}) \quad \forall a \in A \quad \forall \bar{z} \in \bar{A} \\ (s_{1D}, \lambda, \Delta) &\longrightarrow (s_{eD}, \Delta) \\ (s_{1D}, \bar{a}, \bar{a}) &\longrightarrow (s_{1D}, \lambda) \quad \forall \bar{a} \in \bar{A}. \end{aligned}$$

**Satz 3.1.2** Es sei  $\mathcal{K}$  der oben beschriebene RegDyck Automat über dem Alphabet  $A \cup \bar{A}$ , dann existiert ein Kellerautomat  $\mathcal{D}$ , welcher eine Dycksprache über dem Alphabet  $A \cup \bar{A}$  erkennt und ein endlicher Automat  $\mathcal{R}$ , deren Durchschnitt ihrer Sprachen gleich der Sprache von  $\mathcal{K}$  ist.

*Beweis.* Die Idee des Beweises ist, dass der Kellerinhalt des Kellerautomaten  $\mathcal{K}$  und des Kellerautomaten  $\mathcal{D}$  zu jeder Zeit gleich ist. Weiterhin befinden sich der Kellerautomat  $\mathcal{K}$  und der endliche Automat  $\mathcal{R}$  im gleichen Zustand.

Es sei  $g \in (A \cup \bar{A})^*$  die Eingabe.

Es sei  $\mathcal{D}$  in der Konfiguration  $(s_0, g, \Delta)$  und  $\mathcal{K}$  in der Konfiguration  $(s_k, g, \diamond\Delta)$ , wobei  $s_0 \in S_D, s_k \in S_K \setminus \{s_{0K}\}$ . Dann folgt nach  $|g|$  Schritten, dass  $\mathcal{D}$  in der Konfiguration  $(s_d, \lambda, \bar{v}\Delta)$  und  $\mathcal{K}$  in der Konfiguration  $(s_i, \lambda, \bar{v} \diamond \Delta)$  ist,

$$s_d \in S_D \setminus \{s_{0D}\}, s_i \in S_K \setminus \{s_{0K}, s_{1K}\}, \bar{v} \in \bar{A}^*.$$

Der Beweis folgt per Induktion über die Länge eines Teilwortes  $g'' \in (A \cup \bar{A})^*$  der Eingabe  $g = g''g'$  mit  $g' \in (A \cup \bar{A})^*$ .

Die Behauptung ist trivial für  $|g''| = 0$ , da in diesem Fall  $g'' = \lambda$  und  $\mathcal{D}, \mathcal{K}$  in der angegebenen Konfiguration bleiben bzw. in  $\mathcal{K}$  die Überführung

$$(s_k, \lambda, \diamond) \longrightarrow (s_E, \diamond)$$

oder die Überführung

$$(s_k, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda)$$

verwendet wird. Dabei ändert sich allerdings der Kellerinhalt nicht (bis auf das Zeichen  $\diamond$ , was nicht relevant ist). Somit haben  $\mathcal{D}$  und  $\mathcal{K}$  für die Eingabe  $\lambda$  den gleichen Kellerinhalt.

Die Behauptung gelte nun für  $|g''| \leq n \in \mathbb{N}_0$ .

Es sei  $ag'$  oder  $\bar{a}g'$  der Rest eines Eingabewortes  $g''g'$  mit  $g'' = wa$  bzw.  $g'' = w\bar{a}$ ,  $|w| = n$ .

**1. Fall:**  $g'' = wa$

Nach Induktionsvoraussetzung und Definition ist der Automat  $\mathcal{K}$  nach Abarbeitung von  $w$  in der Konfiguration  $(s_i, ag', \bar{xw}' \diamond \Delta)$  mit  $s_i \in S_K \setminus \{s_{0K}, s_{1K}\}$ . Nach Voraussetzung hat der Kellerautomat  $\mathcal{D}$  nach Abarbeitung von  $w$  den gleichen Kellerinhalt  $\bar{xw}'$  und den gleichen Rest  $ag'$  des Eingabewortes wie  $\mathcal{K}$ . Er ist nach Definition in der Konfiguration  $(s_{1D}, ag', \bar{xw}' \Delta)$ .

Der Automat  $\mathcal{K}$  benutzt für  $s_i \notin \{s_E, s_{eK}\}$  die Überführung

$$(s_i, a, \bar{x}) \longrightarrow (s_j, \bar{a}\bar{x}) \quad \forall \bar{x} \in G_K$$

und wechselt in die Konfiguration  $(s_j, g', \bar{a}\bar{xw}' \diamond \Delta)$ .

Der Fall  $s_i \in \{s_E, s_{eK}\}$  kann nicht auftreten, da die Eingabe bisher ungleich  $\lambda$  war und der Keller nicht leer ist. Nur in diesem Fall wechselt der Automat in einen der Zustände  $s_E$  oder  $s_{eK}$ .

Der Kellerautomat  $\mathcal{D}$  benutzt die Überführung

$$(s_{1D}, a, \bar{x}) \longrightarrow (s_{1D}, \bar{a}\bar{x}) \quad \forall \bar{x} \in G_D$$

und wechselt in den Zustand  $(s_{1D}, g', \bar{a}\bar{xw}' \Delta)$ .

Somit haben  $\mathcal{D}$  und  $\mathcal{K}$  denselben Kellerinhalt nach Abarbeitung von  $g''$ .

**2. Fall:**  $g'' = w\bar{a}$

Nach Induktionsvoraussetzung und Definition befindet sich der Automat  $\mathcal{K}$  nach Abarbeitung von  $w$  in einer Konfiguration  $(s_i, \bar{a}g', \bar{xw}' \diamond \Delta)$  mit  $s_i \in S_K \setminus \{s_{0K}, s_{1K}\}$ . Der Kellerautomat  $\mathcal{D}$  befindet sich nach Induktionsvoraussetzung und Definition nach Abarbeitung von  $w$  in der Konfiguration

$(s_{1D}, \bar{a}g', \bar{x}\bar{w}' \triangle)$ .

Der Kellerautomat  $\mathcal{K}$  benutzt für  $s_i \notin \{s_E, s_{eK}\}$  und  $\bar{x} = \bar{a}$  die Überführung

$$(s_i, \bar{a}, \bar{a}) \longrightarrow (s_j, \lambda)$$

und befindet sich danach in der Konfiguration  $(s_j, g', \bar{w}' \triangle)$ .

Der Fall  $s_i \in \{s_E, s_{eK}\}$  kann nicht auftreten, da die Eingabe bisher ungleich  $\lambda$  war und der Keller nicht leer ist und nur in diesem Fall der Automat in einen der Zustände  $s_E$  oder  $s_{eK}$  wechselt.

Falls  $\bar{x} \neq \bar{a}$  ist der Automat  $\mathcal{K}$  nicht definiert und bleibt in der Konfiguration  $(s_i, \bar{a}g', \bar{x}\bar{w}' \triangle)$ .

Der Automat  $\mathcal{D}$  benutzt für  $\bar{x} = \bar{a}$  die Überführung

$$(s_{1D}, \bar{a}, \bar{a}) \longrightarrow (s_{1D}, \lambda)$$

und befindet sich danach in der Konfiguration  $(s_{0D}, g', \bar{w}' \triangle)$ .

Für  $\bar{x} \neq \bar{a}$  ist der Automat  $\mathcal{D}$  nicht definiert und bleibt in der Konfiguration  $(s_{1D}, \bar{a}g', \bar{x}\bar{w}' \triangle)$ .

Somit haben  $\mathcal{D}$  und  $\mathcal{K}$  denselben Kellerinhalt nach Abarbeitung von  $g''$ .

Befindet sich der Automat  $\mathcal{K}$  in der Konfiguration  $(s_0, ag', \bar{w} \triangle)$ , dann wechselt er in die Konfiguration  $(s_j, g', \bar{a}\bar{w} \triangle)$ . Der Automat  $\mathcal{D}$  wechselt von der Konfiguration  $(s_0, ag', \bar{w} \triangle)$  in  $(s_0, g', \bar{a}\bar{w} \triangle)$ .

Ist  $\mathcal{K}$  in der Konfiguration  $(s_0, \bar{a}g', \bar{b}\bar{w}' \triangle)$  und  $\mathcal{D}$  in der Konfiguration  $(s_0, \bar{a}g', \bar{b}\bar{w}' \triangle)$  gibt es für  $\bar{b} \neq \bar{a}$  in beiden Automaten keine Produktion die angewendet werden kann. Somit sind beide nicht definiert, sofern  $\bar{b} \neq \bar{a}$ . Sonst wechselt  $\mathcal{K}$  in die Konfiguration  $(s_j, g', \bar{w} \triangle)$  und  $\mathcal{D}$  in die Konfiguration  $(s_0, g', \bar{w}' \triangle)$ .

Der zweite Teil des Beweises befasst sich mit den Zuständen des RegDyck Automaten  $\mathcal{K}$  und den Zuständen des endlichen Automaten  $\mathcal{R}$ .

Es sei  $\mathcal{R}$  in der Konfiguration  $(s_r, g)$  und  $\mathcal{K}$  in der Konfiguration  $(s_k, g, \triangle)$ ,  $s_r \in S_R, s_k \in S_K \setminus \{s_{0K}\}$ . Dann folgt nach  $|g|$  Schritten, dass  $\mathcal{R}$  in der Konfiguration  $(S_i, \lambda), S_i \subseteq S_R$ , und  $\mathcal{K}$  in der Konfiguration  $(S_i, \lambda, \bar{v} \triangle)$  ist, da  $\mathcal{R}$  nichtdeterministisch sein kann. Hierbei ist  $\bar{v} \in \bar{A}^*$ .

Ist  $s_k = s_{1K}$ , so ist  $s_r = s_{0R}$ . Dann überführt sowohl  $\mathcal{K}$  als auch  $\mathcal{R}$  in die Menge  $S_i \subseteq S_R$ . Somit ist die Zustandsmenge nach der Überführung wieder gleich.

Im Fall  $s_k \neq s_{1K}$  gilt  $s_k = s_r$  oder  $s_k \in \{s_E, s_{eK}\}$ .

Der Beweis folgt per Induktion über die Länge eines Teilwortes  $g'' \in (A \cup \bar{A})^*$  der Eingabe  $g = g'', g' \in (A \cup \bar{A})^*$ .

Sei  $|g''| = 0$ , also  $g'' = \lambda$ . Die Überföhrungsfunktion des Automaten  $\mathcal{R}$  ist nicht definiert und der Automat bleibt in der angegebenen Konfiguration.

Falls  $s_k \in S_R$ , verwendet  $\mathcal{K}$  entweder die Überführung

$$(s_k, \lambda, \triangle) \longrightarrow (s_E, \triangle)$$

oder die Überführung

$$(s_k, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda).$$

Hier unterscheiden sich die Zustandsmengen von  $\mathcal{K}$  und von  $\mathcal{R}$ . Dies ändert allerdings nichts daran, dass  $\mathcal{K}$  den Durchschnitt der Sprachen von  $\mathcal{D}$  und  $\mathcal{R}$  erkennt. Denn in diesem Fall ist entweder ein Dyckprime abgearbeitet, der sowohl von  $\mathcal{R}$  als auch von  $\mathcal{D}$  akzeptiert wird, oder es ist ein Dyckprime, der nur von  $\mathcal{D}$  akzeptiert wird. Dann akzeptiert allerdings  $\mathcal{K}$  auch nicht, da  $\mathcal{R}$  nicht akzeptiert, also  $s_k$  kein Endzustand des Automaten  $\mathcal{R}$  ist.

Die Fälle  $s_k \in \{s_E, s_{eK}\}$  können nicht auftreten, da dafür bereits mindestens eine Überführung vollzogen worden sein muss.

Somit befinden sich beide Automaten in derselben Zustandsmenge, sofern die Eingabe von  $\mathcal{R}$  akzeptiert wird.

Die Behauptung gelte für  $|g''| \leq n \in \mathbb{N}_0$ .

Nach Induktionsvoraussetzung und Konstruktion befindet sich  $\mathcal{R}$  nach der Abarbeitung eines Wortes  $w$ , wobei  $g = g''g'$  mit  $g'' = wa$  oder  $g'' = w\bar{a}$ ,  $|w| = n$ , in einer Zustandsmenge  $S_i \subseteq S_R$ . Ebenfalls befindet sich der Kellerautomat  $\mathcal{K}$  in der Zustandsmenge  $S_i$  und hat den Kellerinhalt  $\overline{xw'}$ ,  $\bar{x} \in \bar{A}$ ,  $\bar{w}' \in \bar{A}^*$ . Die restliche Eingabe besteht aus einem Wort  $ag'$  oder  $\bar{a}g'$ .

**1. Fall:**  $g'' = wa$

Es sei  $(s_i, ag', \overline{xw'} \diamond \Delta)$  eine Konfiguration des Automaten  $\mathcal{K}$  und  $(s_i, ag')$  eine Konfiguration des Automaten  $\mathcal{R}$  mit  $s_i \in S_i$ .

Der Automat  $\mathcal{K}$  wendet die Überführung

$$(s_i, a, \bar{x}) \longrightarrow (s_j, \overline{a\bar{x}}), \bar{x} \in G_K$$

an. Die Überführung wird für ein beliebiges oberes Zeichen des Kellers benutzt. Danach befindet sich der Automat in der Konfiguration  $(s_j, g', \overline{axw'} \diamond \Delta)$ . Der endliche Automat wendet die Überführung

$$(s_i, a) \longrightarrow s_j$$

an. Nach Konstruktion von  $\mathcal{R}$  existiert für jedes Zeichen  $a$  eine solche Überführung, die in den gleichen Zustand führt wie  $\mathcal{K}$ . Somit befinden sich beide Automaten nach Abarbeiten des Zeichens  $a$  in dem gleichen Zustand  $s_j$ .

**2. Fall:**  $g'' = w\bar{a}$

Es sei  $(s_i, \bar{a}g', \overline{xw'} \diamond \Delta)$  eine Konfiguration des Automaten  $\mathcal{K}$  und  $(s_i, \bar{a}g')$  eine Konfiguration des Automaten  $\mathcal{R}$  mit  $s_i \in S_i$ .

Falls  $\bar{x} = \bar{a}$ , wendet der Automat  $\mathcal{K}$  die Überführung

$$(s_i, \bar{a}, \bar{a}) \longrightarrow (s_j, \lambda)$$

an und befindet sich danach in der Konfiguration  $(s_j, g', \bar{w}' \diamond \Delta)$ .

Sofern  $\bar{x} = \bar{a}$ , benutzt der Automat  $\mathcal{R}$  nach Konstruktion die Überführung

$$(s_i, \bar{a}) \longrightarrow s_j$$

und befindet sich in dem Zustand  $s_j$ .

Beide Automaten befinden sich in dem gleichen Zustand.

Gilt allerdings  $\bar{x} \neq \bar{a}$ , so akzeptiert der Automat  $\mathcal{K}$  nicht und der Automat  $\mathcal{R}$  wechselt in den nächsten Zustand. Automat  $\mathcal{D}$  akzeptiert allerdings auch nicht und somit wird diese Sprache weder von  $\mathcal{R}$  noch  $\mathcal{D}$  erkannt.

Befinden sich die Automaten in den Konfigurationen  $(s_0, \bar{a}g', \bar{b}\bar{w}' \diamond \Delta)$  und  $(s_0, \bar{a}g')$ , verweist  $\mathcal{R}$  möglicherweise in einen anderen Zustand und  $\mathcal{K}$  akzeptiert die Eingabe nicht. Falls  $\bar{b} \neq \bar{a}$ , akzeptiert  $\mathcal{D}$  ebenfalls nicht, da  $\mathcal{K}$  und  $\mathcal{D}$  den gleichen Kellerinhalt haben. Somit liegt die Eingabe nicht im Durchschnitt der Sprachen von  $\mathcal{R}$  und  $\mathcal{D}$ .  $\square$

## 3.2 RegDyck Grammatiken

RegDyck Grammatiken lassen sich den balancierten Grammatiken unterordnen (siehe Kapitel 4), da sie lediglich Einschränkungen gegenüber den balancierten Grammatiken besitzen. Wie im Folgenden gezeigt wird, ist die Sprachfamilie deren Sprachen durch RegDyck Grammatiken erzeugt werden, eine Teilmenge der Sprachfamilie deren Sprachen durch RegDyck Automaten erkannt werden.

Die Vorteile gegenüber anderen Sprachfamilien der XML-artigen Sprachen liegen in der Möglichkeit, mehrdeutige reguläre Ausdrücke und konkurrierende Nichtterminalzeichen in den Inhaltsmodellen der Grammatiken zuzulassen. Die Parser der RegDyck Grammatiken bestehen aus zwei Einheiten: dem Dyck-Automaten und dem endlichen Automaten, welche beide in  $O(n)$  Zeit parsen. Nachteile bestehen darin, dass zwei selbsteinbettende Nichtterminalzeichen  $X, Y$  mit gleichen Ableitungssequenzen nur bedingt erlaubt sind.

In der Praxis können Verbesserungen erreicht werden, da selbsteinbettende Nichtterminale, die gleiche Ableitungen erzeugen, nur sehr selten auftreten.

**Definition 3.2.21** Es sei  $G = (N, T, S, P)$  eine balancierte indizierte Grammatik. Eine Abbildung  $FA : N^* \longrightarrow \text{Reg}(T)$  liefert für eine Folge von Nichtterminalzeichen einen regulären Ausdruck von Terminalzeichen, wie folgt:

- Für jede Produktion  $X \longrightarrow xm\bar{x}$  aus  $P$  wird ein nichtdeterministischer endlicher Automat  $\mathcal{R}_X$  konstruiert, der die rechte Seite der Produktion erkennt.
- In den Überführungen sind Nichtterminalzeichen und Terminalzeichen der Grammatik enthalten. Nun werden alle Nichtterminalzeichen aus

den Überführungen folgendermaßen entfernt: Es sei  $(s_i, X) \longrightarrow (s_j)$  eine Überführung eines Automaten  $\mathcal{R}_Y$ , wobei  $s_i$  und  $s_j$  zwei Zustände des Automaten  $\mathcal{R}_Y$  sind und  $X$  ein Nichtterminalzeichen ist. Dann wird ein  $\lambda$ -Übergang von  $s_i$  zum Startzustand von  $\mathcal{R}_X$  und von den Endzuständen von  $\mathcal{R}_X$  zu  $s_j$  erzeugt. Die Überführung  $(s_i, X) \longrightarrow (s_j)$  wird in  $\mathcal{R}_Y$  entfernt. Zu beachten ist, dass dies auch durchgeführt wird, wenn ein Nichtterminalzeichen selbsteinbettend auftritt.

- Wenn  $FA(X)$ ,  $X \in N$ , erzeugt werden soll, wird aus dem nichtdeterministischen endlichen Automaten  $\mathcal{R}_X$  ein regulärer Ausdruck erzeugt.
- $FA(Z)$ , wobei  $Z$  eine Folge  $Z = Z_0Z_1 \dots Z_n$ ,  $n \in \mathbb{N}_0$ , von Nichtterminalzeichen ist, ist die Konkatenation  $FA(Z) = FA(Z_0)FA(Z_1) \dots FA(Z_n)$ .

Für die Definition von RegDyck Grammatiken wird die Eigenschaft benutzt, dass zwei Nichtterminalzeichen in Konflikt stehen. Um in Konflikt zu stehen, müssen die beiden Nichtterminalzeichen selbsteinbettend sein oder es muss die Möglichkeit bestehen, aus ihnen in einer Folge von Ableitungsschritten ein selbsteinbettendes Nichtterminalzeichen abzuleiten. Die selbsteinbettenden Nichtterminalzeichen müssen zusätzlich gleiche Ableitungssequenzen erzeugen. Es sei  $X$  ein selbsteinbettendes Nichtterminalzeichen und  $X \Longrightarrow^+ u_0u_1 \dots u_nXv_n \dots v_1v_0$  eine Ableitungsfolge, wobei  $u_i, v_i$  für  $i \in \{0, \dots, n\}$ ,  $u_iv_i \in D$ , Folgen von Terminalzeichen sind. Dann erkennt der RegDyck Automat unabhängig voneinander die Folgen  $u_i$  und  $v_i$  durch Schleifen. Da der RegDyck Automat weder in seinem Keller, noch in seinen Zuständen kodieren kann, wieviele Schleifen von  $u$ 's er durchlaufen hat, sondern nur welche Zeichen er noch erkennen muss, um einen Dyckprime zu erkennen, ist es unklar, wie oft er die Schleife von  $v$ 's durchlaufen muss.

Wenn nun ein anderes selbsteinbettendes Nichtterminalzeichen  $Y$  mit der Ableitungsfolge  $Y \Longrightarrow^+ u'_0u'_1 \dots u'_mYv'_m \dots v'_1v'_0$  existiert und es eine Kombination  $u_iv'_j \in D$ ,  $j \in \{0, \dots, m\}$ , gibt, kann es dazu kommen, dass der RegDyck Automat die durch die Grammatik beschriebene Sprache nicht korrekt erkennt, sondern zusätzlich noch andere Wörter.

Die Definition der RegDyck Grammatiken verhindert solche Konflikte von Nichtterminalzeichen.

Es gibt zwei Arten, dass zwei Nichtterminalzeichen in Konflikt stehen. Die erste Art kann auftreten, wenn zwei Nichtterminalzeichen in einem Inhaltsmodell enthalten sind. Die zweite Art kann auftreten, wenn ein Nichtterminalzeichen im Inhaltsmodell des anderen oder im Inhaltsmodell eines Nachfolgers des anderen enthalten ist. Dabei ist ein Nachfolger von  $X$  ein Nachfolger der Produktion von  $X$  im Baum der indizierten balancierten Grammatik. Die zweite Art muss genauer betrachtet werden. Es seien  $X$  und  $Y$  zwei Nichtterminalzeichen, wobei überprüft wird, ob  $X$  zu  $Y$  auf die zweite Art

in Konflikt steht. Dann muss die Produktion von  $Y$  ein Nachfolger von der Produktion von  $X$  im Baum der indizierten balancierten Grammatik sein. Nun müssen die Ableitungsfolgen untersucht werden, die  $X$  erzeugt, bis  $Y$  auftritt. In der Ableitungsfolge werden jedoch nur die Nachfolger untersucht, aus denen  $Y$  abgeleitet werden kann. Alle anderen Teile der Grammatik aus denen Dyckprimes erzeugt werden können, werden nicht untersucht. Deshalb werden in der Definition Produktionen verändert. Es werden nur die Ableitungsfolgen herausgefiltert, die zu einem Konflikt zwischen den beiden Nichtterminalzeichen führen könnten.

**Definition 3.2.22** Im Folgenden wird definiert, wann zwei Nichtterminalzeichen in Konflikt stehen. Es sei  $G = (N, T, S, P)$ ,  $T = A \cup \bar{A}$ , eine indizierte balancierte Grammatik und  $X$  und  $Y$  zwei Nichtterminalzeichen. Zwei Unterschiedliche Fälle sind zu betrachten:

1.  $Y$  steht in Konflikt zu  $X$ , wenn:  
 $\exists X, Y, Z \in N, \exists z \in A, \exists \bar{z} \in \bar{A}, \exists u, v, w \in N^* :$

$$(Z \longrightarrow zm\bar{z} \wedge uXwYv \in L(m) \wedge L(FA(X)FA(Y)) \cap D_T \neq \emptyset)$$

Der Konflikt wird aufgehoben, wenn ein Wort  $uCv \in L(m)$  existiert und  $L(G') = L(FA(X)FA(w)FA(Y)) \cap D_T$ ,  $G' = (N, T, C, P)$  gilt.

2. Ein Nichtterminalzeichen  $X$  steht auch in Konflikt zu  $Y$ , wenn gilt:  
 $V_Y$  ist die Menge der Nichtterminalzeichen, deren Produktionen Vorgänger und  $N_Y$  die Menge der Nichtterminalzeichen, deren Produktionen Nachfolger von der Produktion  $Y \longrightarrow ym_Y\bar{y}$  im Baum der indizierten balancierten Grammatik sind.

Für alle Produktionen  $Z \longrightarrow zm\bar{z}$ , wobei  $Z \in V_Y$ , wird das Inhaltsmodell gelöscht und ein neues Inhaltsmodell  $m_0$  erzeugt, welches durch die disjunkte Vereinigung der Elemente der Menge  $M_0 = \{V \in N \mid (\alpha V \beta \in L(m)) \wedge ((V = X) \vee (V \in N_X \wedge V \in V_Y)), \alpha, \beta \in N^*\}$  beschrieben wird.

Die Produktionen und Nichtterminalzeichen, die durch diese Konstruktion erzeugt werden, werden mit einem  $\#$  versehen.

$X$  steht im Konflikt zu  $Y$ , wenn gilt:

$$\exists Y \in N, X \in V_Y :$$

$$((L(FA(Y)FA(X\#)) \cap D_T \neq \emptyset) \vee (L(FA(X\#)FA(Y)) \cap D_T \neq \emptyset))$$

Die Konflikte werden aufgehoben, wenn  $L(FA(X)) \cap D_T = L(G')$ ,  $G' = (N, T, Y, P)$ .

□

Nun wird definiert, in welchem Fall ein Nichtterminalzeichen mit sich selbst in Konflikt steht.

**Definition 3.2.23** Es sei  $G = (N, T, S, P)$  eine balancierte Grammatik und  $X$  ein selbsteinbettendes Nichtterminalzeichen.  $X$  steht in Konflikt mit sich selbst, wenn Ableitungen

$$\begin{aligned} X &\Longrightarrow uXv \\ X &\Longrightarrow u'Xv' \end{aligned}$$

existieren, so dass  $uv' \in D_T$  und keine Ableitung

$$X \Longrightarrow uXv'$$

existiert. □

**Definition 3.2.24** Eine RegDyck Grammatik ist eine balancierte Grammatik  $G = (N, T, S, P)$  für die gilt:

- Es existieren keine Nichtterminalzeichen, die in Konflikt mit sich selbst stehen.
  - Es existieren keine Nichtterminalzeichen, die in Konflikt miteinander stehen.
- 

Es folgen nun ein paar Beispiele für RegDyck Grammatiken.

**Beispiel 3.2.14** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, Y, Z\}$
- $T = \{x, \bar{x}, y, \bar{y}, z, \bar{z}\}$
- $P = \{$ 
  - $S \longrightarrow xY\bar{x}$
  - $Y \longrightarrow y((ZS)|(SZ))^? \bar{y}$
  - $Z \longrightarrow z\bar{z}\}$

Das Nichtterminalzeichen  $S$  steht in Konflikt mit sich selbst, da die Ableitungen  $S \Longrightarrow^+ (xyz\bar{z})^n S(\bar{y}\bar{x})^n$  und  $S \Longrightarrow^+ (xy)^m S(z\bar{z}\bar{y}\bar{x})^m$  möglich sind, jedoch nicht  $S \Longrightarrow^+ (xyz\bar{z})^l S(z\bar{z}\bar{y}\bar{x})^l$  mit  $l, m, n \in \mathbb{N}$ .

Würde die Produktion von  $Y$  durch  $Y \longrightarrow y((ZS)|(SZ)|(ZSZ))^? \bar{y}$  ersetzt, wäre die Grammatik eine RegDyck Grammatik. □

**Beispiel 3.2.15** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, Y, Z\}$

- $T = \{x, \bar{x}, y, \bar{y}\}$
- $P = \{$ 
  - $S \longrightarrow xYZ\bar{x}$
  - $Y \longrightarrow yY?\bar{y}$
  - $Z \longrightarrow yZ?\bar{y}\}$

Die balancierte Grammatik  $G$  ist keine RegDyck Grammatik, weil in der Startproduktion zwei selbsteinbettende Nichtterminale  $Y$  und  $Z$  existieren. Beide stehen in Konflikt miteinander, da  $FA(Y) = (y)^+(\bar{y})^+$ ,  $FA(Z) = (y)^+(\bar{y})^+$  und  $L((y)^+(\bar{y})^+(y)^+(\bar{y})^+) \cap D_T \neq \emptyset$ .

Der Konflikt kann aufgehoben werden, wenn die Produktionsmenge folgendermaßen ersetzt wird:

$$\begin{aligned}
 S &\longrightarrow x((YZ)|A)\bar{x} \\
 A &\longrightarrow y(A|(YZ))\bar{y} \\
 Y &\longrightarrow yY?\bar{y} \\
 Z &\longrightarrow yZ?\bar{y}
 \end{aligned}$$

□

**Beispiel 3.2.16** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, X, Y, Z\}$
- $T = \{s, \bar{s}, y, \bar{y}\}$
- $P = \{$ 
  - $S \longrightarrow sXZ\bar{s}$
  - $X \longrightarrow yY\bar{y}$
  - $Y \longrightarrow yY?\bar{y}$
  - $Z \longrightarrow yZ?\bar{y}\}$

Die balancierte Grammatik  $G$  ist keine RegDyck Grammatik, da die Nichtterminalzeichen  $X, Z$  in Konflikt stehen, da  $FA(X) = y(y)^+(\bar{y})^+\bar{y}$  und  $FA(Z) = (y)^+(\bar{y})^+$  und  $L(y(y)^+(\bar{y})^+\bar{y}(y)^+(\bar{y})^+) \cap D_T \neq \emptyset$ . Nach gleichem Muster wie im vorherigen Beispiel könnte die Grammatik zu einer RegDyck Grammatik ergänzt werden.

□

**Beispiel 3.2.17** Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit:

- $N = \{S, Y, Z\}$
- $T = \{x, \bar{x}, z, \bar{z}\}$
- $P = \{$ 
  - $S \longrightarrow x(S|ZY Z)\bar{x}$
  - $Y \longrightarrow xY?\bar{x}$
  - $Z \longrightarrow z\bar{z}\}$

Die balancierte Grammatik  $G$  ist keine RegDyck Grammatik, da die Nichtterminalzeichen  $S, Y$  in Konflikt stehen. Es gilt  $FA(S^\#) = (x)^+(\bar{x})^+$  und  $FA(Y) = (x)^+(\bar{x})^+$  und  $L((x)^+(\bar{x})^+(x)^+(\bar{x})^+) \cap D_T \neq \emptyset$  gilt. Der Konflikt wird nicht verhindert, da in  $FA(S) \cap D_T$  das Wort  $w = xxxz\bar{x}\bar{x}\bar{x}z\bar{x}\bar{x}$  enthalten ist, jedoch nicht in  $L(G)$ .

□

### Konstruktion eines RegDyck Automaten aus einer RegDyck Grammatik

Es wird eine beliebige RegDyck Grammatik vorgegeben und für diese ein RegDyck Automat konstruiert.

Es sei  $G = (N, T, S, P), T = A \cup \bar{A}$  eine RegDyck Grammatik und  $L(G)$  die dazugehörige Sprache.

Alle Produktionen in  $P$  haben die Form

$$X \longrightarrow xm\bar{x},$$

wobei  $X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$ .

Die Grammatik  $G$  wird indiziert und zu  $G' = (N', T', S', P')$ .

Nun wird ein RegDyck Automat konstruiert, der genau die Sprache dieser indizierten Grammatik erkennt. Die Konstruktion eines erkennenden RegDyck Automaten  $\mathcal{K} = (S_K, A_K, G_K, s_{0K}, \delta_K)$  enthält die Konstruktion von RegDyck Automaten  $\mathcal{K}_X = (S_X, A_X, G_X, s_{0X}, \delta_X)$ . Für jedes Nichtterminalzeichen  $X \in N'$  wird ein erkennender Automat  $\mathcal{K}_X$  konstruiert. Dieser Automat erkennt die Sprache, die durch das Nichtterminalzeichen  $X$  erzeugt wird. Zum Schluss werden die einzelnen erkennenden Automaten gemäß der Grammatik miteinander verbunden, so dass die komplette Sprache der Grammatik  $G'$  durch den Automaten  $\mathcal{K}$  erkannt wird.

Für alle Nichtterminale  $X \in N'$  wird nun ein Teilautomat  $\mathcal{K}_X$  von  $\mathcal{K}$  konstruiert.

Wenn  $m = \lambda$ , dann müssen lediglich folgende Überführungen erzeugt werden:

$$\begin{aligned} (s_0, x, \bar{z}) &\longrightarrow (s_1, \bar{x}\bar{z}) \quad \forall \bar{z} \in G_X, \\ (s_1, \bar{x}, \bar{x}) &\longrightarrow (s_{X_e}, \lambda) \end{aligned}$$

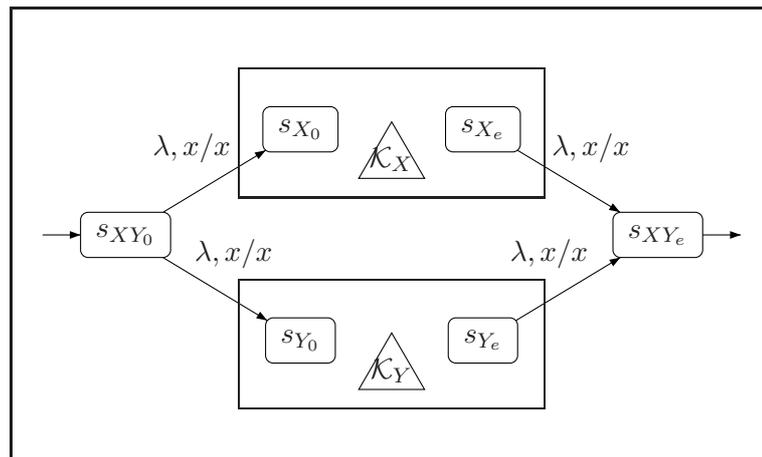
mit  $s_0, s_1, s_{X_e} \in S_X, x \in A_X, \bar{z}, \bar{x} \in G_X$ .

Wenn  $m \neq \lambda$ , dann muss noch beachtet werden, dass die Inhaltsmodelle aus regulären Ausdrücken bestehen und das Zusammensetzen der Teilautomaten erschwert wird.

Die Konstruktion erfolgt induktiv über den Aufbau der regulären Ausdrücke. In den folgenden Abbildungen ist  $s_{X_0}$  der Startzustand und  $s_{X_e}$  der Endzustand des Automaten  $\mathcal{K}_X$ . Der Rahmen um  $\mathcal{K}_X$  bedeutet, dass hier der Automat  $\mathcal{K}_X$  arbeitet (analoges gilt für  $\mathcal{K}_Y$ ).

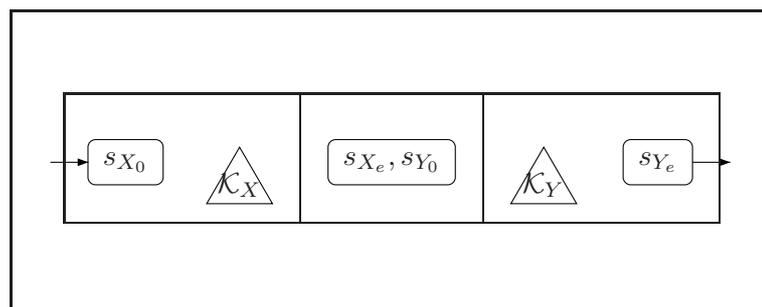
Es wird angenommen, dass  $\mathcal{K}_X$  und  $\mathcal{K}_Y$  RegDyck Automaten für die Produktionen der Nichtterminale  $X$  und  $Y$  sind.

1. Für den regulären Ausdruck  $X|Y$  wird  $\mathcal{K}_{X|Y}$  konstruiert.



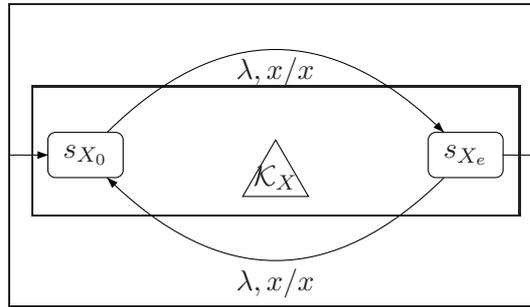
Hierbei wird über den gemeinsamen Startzustand  $s_{XY_0}$  von  $\mathcal{K}_X$  und  $\mathcal{K}_Y$  durch einen  $\lambda$ -Übergang in die Startzustände von  $\mathcal{K}_X$  und  $\mathcal{K}_Y$  verwiesen. Es existiert ein gemeinsamer Endzustand  $s_{XY_e}$ .

2. Für den regulären Ausdruck  $XY$  wird der zusammengesetzte RegDyck Automat  $\mathcal{K}_{XY}$  konstruiert.



Hierbei werden der Endzustand  $s_{X_e}$  und der Anfangszustand  $s_{Y_0}$  zu einem Zustand vereinigt.

3. Für den regulären Ausdruck  $X^*$  wird der RegDyck Automat  $\mathcal{K}_{X^*}$  konstruiert.



Hierbei wird ein  $\lambda$ -Übergang vom Startzustand  $s_{X_0}$  in den Endzustand  $s_{X_e}$  erzeugt und umgekehrt.

Ein Fall muss gesondert betrachtet werden. Es sei  $X$  ein Nichtterminalzeichen, für das ein RegDyck Automat  $\mathcal{K}_X$  konstruiert werden soll. Während der Konstruktion des RegDyck Automaten  $\mathcal{K}_X$  tritt das Nichtterminalzeichen  $X$  nochmal auf. Da die balancierte Grammatik indiziert ist, kann es nur innerhalb der Konstruktion des RegDyck Automaten  $\mathcal{K}_X$  auftreten und sonst nicht. Das Nichtterminalzeichen  $X$  ist somit ein selbsteinbettendes Nichtterminalzeichen. Der Automat, der für  $X$  nun ein weiteres mal konstruiert werden müsste, ist also derselbe Automat  $\mathcal{K}_X$ . Er wird jedoch nicht erneut konstruiert, weil die Konstruktion wegen der Selbsteinbettung wieder und wieder durchgeführt werden müsste. Für alle Konstruktionen (siehe oben Punkt 1 bis 3), die für das Inhaltsmodell durchgeführt werden müssen, muss lediglich der Startzustand und der Endzustand eines RegDyck Automaten bekannt sein. Beide sind von dem RegDyck Automaten  $\mathcal{K}_X$  bekannt. Somit kann dieser für die Konstruktion verwendet werden.

Es bleibt zu zeigen, dass alle  $\lambda$ -Übergänge ersetzt werden können, da RegDyck Automaten nur einen einzigen  $\lambda$ -Übergang vom Startzustand  $s_{0K}$  zum Initialisieren erlauben und  $\lambda$ -Übergänge zum Entfernen des Kellersymbols  $\diamond$  nach dem Abarbeiten eines Dyckprimes.

Es seien  $s_k, s_j \in S_X$  zwei Zustände eines Automaten  $\mathcal{K}_X$ .

Wenn eine Überführung

$$(s_k, \lambda, \bar{z}) \longrightarrow (s_j, \bar{z}), \bar{z} \in G_X$$

existiert, werden für alle Überführungen der Form

$$(s_j, a, \bar{z}) \longrightarrow (s_n, \bar{a}\bar{z}), a \in A_X, \bar{a}, \bar{z} \in G_X, s_n \in S_X$$

und

$$(s_j, \bar{a}, \bar{a}) \longrightarrow (s_n, \lambda)$$

Überführungen

$$(s_k, a, \bar{z}) \longrightarrow (s_n, \overline{a\bar{z}})$$

oder

$$(s_k, \bar{a}, \bar{a}) \longrightarrow (s_n, \lambda)$$

erzeugt.

Dies ist bei den RegDyck Automaten in jedem Fall erlaubt, da sie nichtdeterministische Zustandsüberföhrungsfunktionen besitzen dürfen.

Die Überföhrungen des Automaten  $\mathcal{K}$ , die dafür sorgen, dass die Dyckprimes erkannt werden, die die Grammatik erzeugt, sind durch die Konstruktion erzeugt. Nun müssen noch die Überföhrungen konstruiert werden, die dafür sorgen, dass nur gültige Dyckprimes erkannt werden.

Es seien

$$\begin{aligned} S &\longrightarrow xm_0\bar{x} \\ S &\longrightarrow ym_1\bar{y} \\ S &\longrightarrow wm_2\bar{w} \\ &\dots \end{aligned}$$

die Startproduktionen der Grammatik  $G'$ . Dann besitzt der Automat  $\mathcal{K}_S$  einen Startzustand  $s_i$ , für den es folgende Überföhrungen gibt:

$$\begin{aligned} (s_{S_0}, x, \bar{z}) &\longrightarrow (s_j, \overline{x\bar{z}}) \\ (s_{S_0}, y, \bar{z}) &\longrightarrow (s_k, \overline{y\bar{z}}) \\ (s_{S_0}, w, \bar{z}) &\longrightarrow (s_l, \overline{w\bar{z}}) \\ &\dots \end{aligned}$$

Es wird für  $\mathcal{K}$  die Startüberföhrung

$$(s_{0K}, \lambda, \Delta) \longrightarrow (s_{1K}, \diamond\Delta)$$

und die Überföhrungen

$$\begin{aligned} (s_{1K}, x, \diamond) &\longrightarrow (s_j, \overline{x\bar{\diamond}}) \\ (s_{1K}, y, \diamond) &\longrightarrow (s_k, \overline{y\bar{\diamond}}) \\ (s_{1K}, w, \diamond) &\longrightarrow (s_l, \overline{w\bar{\diamond}}) \\ &\dots \end{aligned}$$

erzeugt.

$s_{S_e}$  sei der Zustand der nach Abarbeitung eines Dyckprimes von  $\mathcal{K}_S$  erreicht wurde. Dafür wird in  $\mathcal{K}$  eine Überföhrung zum Leeren des Kellers definiert.

$$(s_{S_e}, \lambda, \diamond) \longrightarrow (s_{eK}, \lambda), s_{S_e} \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}\}$$

Für alle anderen Zustände  $s_i \in S_K \setminus \{s_{0K}, s_{1K}, s_E, s_{eK}, s_{S_e}\}$  werden Überföhrungen

$$(s_i, \lambda, \diamond) \longrightarrow (s_E, \diamond)$$

in  $\mathcal{K}$  erzeugt.

Damit ist der komplette RegDyck Automat  $\mathcal{K}$  konstruiert.

**Satz 3.2.3** Jede, von einer RegDyck Grammatik  $G$  erzeugte Sprache, kann durch einen RegDyck Automat  $\mathcal{K}$  erkannt werden.

*Beweis.* Es muss nur gezeigt werden, dass der RegDyck Automat genau die Ableitungen der Grammatik  $G$  erkennt, die Selbsteinbettung enthalten. Wenn keine Selbsteinbettung vorhanden ist, wird durch die Konstruktion des Automaten mittels der abgewandelten Thompson Methode klar, dass der RegDyck Automat  $\mathcal{K}_X$  genau die Sprache, die von der Grammatik  $G$  erzeugt wird, erkennt.

Es sei  $G = (N, T, S, P)$  eine indizierte RegDyck Grammatik und  $\mathcal{K} = (S_K, A_K, G_K, s_{0K}, \delta_K)$  der RegDyck Automat aus der obigen Konstruktion.

Es sei

$$S \Longrightarrow^+ uyw\alpha\bar{y}v$$

eine Ableitung der Grammatik  $G$ , wobei  $S \in N, u \in T^*, y, \bar{y} \in T, w \in D_T, \alpha \in \text{Reg}(N), v \in \text{Reg}(N \cup T)$ . Nun wird nach einem selbsteinbettenden Nichtterminalzeichen  $X$  aus der Menge  $M_\alpha = \{Y \mid Yq \in L(\alpha)\}$  abgeleitet,  $X \in M_\alpha$ .

$X$  erzeugt

$$X \Longrightarrow^* u_0u_1 \dots u_n X v_n \dots v_1v_0 \text{ mit}$$

$n \in \mathbb{N}, u_i v_i \in D_T, i \in \{1, \dots, n\}$ .

Der Automat  $\mathcal{K}_X$  erkennt einen selbsteinbettenden Teil mittels einer Schleife. In dem Teilautomat  $\mathcal{K}_X$ , der die Sprache des selbsteinbettenden Nichtterminalzeichen erkennt, wird, an der Stelle in dem das Nichtterminalzeichen selbsteinbettend auftritt, wieder in den Startzustand von  $\mathcal{K}_X$  verwiesen. So erkennt er für  $X$  die Sequenzen

$$u_0u_1 \dots u_n w_0 v_1 v_2 \dots v_j \text{ mit}$$

$j, k \in \mathbb{N}, w_0, u_k v_{n-k} \in D_T$ .

Jedoch wird durch den Keller des RegDyck Automaten gesichert, dass nur Dyckprimes erkannt werden. Das bedeutet, wenn der RegDyck Automat eine Schleife von  $u$  mehr durchlaufen hat wie von  $v$ , dann muss ein anderer Teilautomat existieren, der auch eine Sprache eines selbsteinbettenden Nichtterminalzeichen beschreibt und durch einen zusätzlichen Schleifendurchlauf den Keller leert.

Wir nehmen an, es gebe einen zweiten Teilautomaten  $\mathcal{K}_Y$  der dafür sorgt, dass der Keller geleert wird.

Dann muss in der Grammatik  $G$  ein selbsteinbettendes Nichtterminalzeichen  $Y$  existieren für das gilt:

$$d = u_0 u_1 \dots u_n X v_j \dots v_1 v_0 p_0 p_1 \dots p_l Y q_m \dots q_1 q_0 \text{ mit}$$

$j, l, m \in \mathbb{N}, d \in D_T$ .

Es gibt zwei Möglichkeiten, wie die Ableitung auftreten können:

1.  $d' = \alpha \beta u_0 u_1 \dots u_n w_0 v_1 v_2 \dots v_j \bar{\beta} \gamma u'_0 u'_1 \dots u'_m w'_0 v'_1 v'_2 \dots v'_i \bar{\gamma} \bar{\alpha}$   
mit  $\alpha, \bar{\alpha}, \beta, \bar{\beta}, \gamma, \bar{\gamma} \in T^*, d' \in D_T, n \neq j, m \neq i$ . der

2.  $d' = \alpha u'_0 u'_1 \dots u'_m \beta u_0 u_1 \dots u_n w_0 v_1 v_2 \dots v_j \bar{\beta} v'_1 v'_2 \dots v'_i \bar{\alpha}$

Im ersten Fall würde jedoch  $Y$  zu  $X$  in Konflikt stehen (nach der ersten Variante der Konflikte), weil  $FA(X)FA(Y) \cap D_T \neq \emptyset$ . Im zweiten Fall würde  $Y$  zu  $X$  im Konflikt stehen (nach der zweiten Variante der Konflikte), da  $FA(X)FA(Y^\#) \cap D_T \neq \emptyset$  oder  $FA(Y^\#)FA(X) \cap D_T \neq \emptyset$ . Oder  $X$  stände im Konflikt mit sich selbst. Es sind jedoch keine Nichtterminalzeichen, die im Konflikt stehen erlaubt und auch keine Nichtterminalzeichen, die mit sich selbst im Konflikt stehen. Dies führt zum Widerspruch der Annahme.

Es zeigt sich, dass der Automat  $\mathcal{K}$  genau die Sprache der Grammatik  $G$  erkennt.

□

Es wurde nur gezeigt, dass jede RegDyck Grammatik durch einen RegDyck Automaten erkannt wird. Die umgekehrte Richtung bleibt zu zeigen.

Für die Untersuchung der XML-artigen Grammatiken ist es unwichtig, ob aus jedem RegDyck Automaten eine RegDyck Grammatik erzeugt werden kann. Die Vorteile, dass mehrdeutige Inhaltsmodelle benutzt werden können und dass konkurrierende Nichtterminalzeichen in den Grammatiken vorhanden sein dürfen bleiben bestehen. Interessant wäre die Frage, ob für jede, durch einen RegDyck Automaten beschriebene Sprache, eine balancierte Grammatik existiert, die die gleiche Sprache erzeugt. Wenn dies nicht der Fall wäre, könnten durch den RegDyck Automaten noch andere nicht balancierte Mengen von Dyckprimes beschrieben werden.

# Kapitel 4

## Einordnung der Modelle

In dem folgenden Kapitel werden die verschiedenen XML-artigen Sprachfamilien eingeordnet. So wird betrachtet, ob es eine Hierarchie zwischen den einzelnen Sprachfamilien gibt und in welchem Bereich der Chomsky Hierarchie die XML-artigen Sprachfamilien einzuordnen sind. Die Ergebnisse sollen für Validierungsalgorithmen benutzt werden. Es gilt zu überprüfen, wie lange benötigt wird, um einen Parser zu konstruieren und wie lange der Parser benötigt, um XML Daten auf ihre Struktur hin zu überprüfen.

### 4.1 XML Ausdrücke

In Beispielen werden XML-artige Sprachen oft durch Ausdrücke angegeben, die den regulären Ausdrücken ähneln.

**Beispiel 4.1.18**  $L(a^n(b\bar{b})^*c\bar{c}a^n) = \{a^n(b\bar{b})^i c\bar{c}a^n \mid i, n \in \mathbb{N}\}$

Sie werden als XML Ausdrücke bezeichnet und in diesem Unterkapitel definiert. Es wird überprüft, ob mittels dieser Ausdrücke alle Sprachen, die durch eine balancierte Grammatik erzeugt werden können, auch durch einen XML Ausdruck dargestellt werden können.

Die entstehenden Ausdrücke werden XML Ausdrücke genannt.

**Definition 4.1.25** Es sei  $T = A \cup \bar{A}$  ein Alphabet, das keines der Zeichen  $V = \{ |, *, +, ?, \emptyset, (, ) \}$  enthält. XML Ausdrücke über  $T$  werden wie folgt festgelegt:

- $x\bar{x}, x \in A, \bar{x} \in \bar{A}$  und  $\emptyset$  sind Unterausdrücke.
- $xu_0u_1 \dots u_i\bar{x}, i \in \mathbb{N}$  ist ein Unterausdruck, wenn  $u_0, u_1, \dots, u_i, i \in \mathbb{N}$  Unterausdrücke sind.
- Wenn  $r$  und  $s$  Unterausdrücke sind, dann ist auch die Konkatenation  $(rs)$  ein Unterausdruck.

- Wenn  $r$  und  $s$  Unterausdrücke sind, dann ist auch  $(r|s)$  ein Unterausdruck.
- Wenn  $r$  ein Unterausdruck ist, dann ist auch  $(r)^*$  ein Unterausdruck.
- $x\bar{x}$ ,  $x \in A$ ,  $\bar{x} \in \bar{A}$  und  $\emptyset$  sind XML Ausdrücke.
- $xu_0u_1 \dots u_i\bar{x}$  ist ein XML Ausdruck, wobei  $u_0u_1 \dots u_i$  Unterausdrücke sind.
- Es gibt eine besondere Art von Ausdrücken, die als rekursive Ausdrücke bezeichnet werden. Sie können sowohl als XML Ausdrücke, als auch als Unterausdrücke benutzt werden.

Es sei  $q$  ein Ausdruck der Form  $xv$ , wobei  $x$  ein Zeichen aus  $A$  und  $v$  eine beliebige Sequenz von Unterausdrücken über  $T$  und Zeichen  $a \in A$  ist. Es sei  $\psi$  eine Abbildung, die aus einem Ausdruck  $q$  alle Unterausdrücke löscht.  $q$  und eine dazugehörige Sequenz  $\bar{q}$  können mit einer Variablen versehen werden ( $q^n, \bar{q}^n, n \in \mathbb{N}$ ), wenn gilt:  $\psi(q)\psi(\bar{q}) \in D_T$ , wobei  $D_T$  die Menge der Dyckprimes über  $T = A \cup \bar{A}$  ist.

Der Teil  $w$  zwischen  $q$  und  $\bar{q}$  ( $q^n w \bar{q}^n$ ) ist ein beliebiger Unterausdruck.

□

Die Klammern können weggelassen werden, wenn die Semantik des Ausdrucks dadurch erhalten bleibt.

Es sei  $r$  ein Unterausdruck, dann gilt:  $r^+ = rr^*$  und  $r? = (r|\emptyset^*)$ . Für  $\emptyset^*$  kann auch  $\lambda$  geschrieben werden.

Jeder XML Ausdruck  $p$  über einem Alphabet  $A$  legt eine Sprache  $L(p)$  folgendermaßen fest:

- Die durch  $\emptyset$  bezeichnete Sprache ist die leere Sprache.
- Die durch  $a\bar{a}$ ,  $a \in A$ ,  $\bar{a} \in \bar{A}$  bezeichnete Sprache enthält nur das Wort  $a\bar{a}$ .
- Für Unterausdrücke  $p$  und  $q$  über  $A$  gilt:  $L(p|q) = L(p) \cup L(q)$ ,  $L(pq) = L(p)L(q)$ ,  $L(p^*) = (L(p))^*$ .
- Es sei  $u$  ein Unterausdruck, dann ist die durch den Ausdruck  $au\bar{a}$  beschriebene Sprache  $aL(u)\bar{a}$ .
- Die von  $v^n u \bar{v}^n$ ,  $v, \bar{v} \in (A \cup \bar{A})^+$  erzeugte Sprache ist  $L(v)^n L(u) L(\bar{v})^n$  mit  $n \in \mathbb{N}$ .

Der Ausdruck  $(a^m \bar{a}^m)^*$  bedeutet hierbei, dass die Sprache  $a^m \bar{a}^m a^i \bar{a}^i a^j \bar{a}^j \dots$  erzeugt werden kann.

**Satz 4.1.4** Die Sprache jedes XML Ausdrucks kann durch eine balancierte Grammatik erzeugt werden.

*Beweis.* Es sei  $p = xu\bar{x}$ ,  $x \in A$ ,  $\bar{x} \in \bar{A}$  ein XML Ausdruck, wobei  $u$  ein Unterausdruck ist. Es wird eine balancierte Grammatik  $G = (N, T, S, P)$  konstruiert. Für den XML Ausdruck  $p$  wird eine Produktion

$$X \longrightarrow xm_u\bar{x}, X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$$

erzeugt. Wenn der Unterausdruck  $u = \lambda$  ist, dann ist auch das Inhaltsmodell der Produktion  $m_u = \lambda$ . Die Konstruktion des Inhaltsmodells  $m_u$  erfolgt über den induktiven Aufbau der XML Ausdrücke. Durch ein wiederholtes Anwenden der folgenden Schritte wird die balancierte Grammatik  $G$  vervollständigt:

- Es sei  $u = (u_0)^*$  der Unterausdruck, dann wird  $m_u = (m_{u_0})^*$  erzeugt.
- Es sei  $u = (u_0|u_1)$  der Unterausdruck, dann wird  $m_u = (m_{u_0}|m_{u_1})$  erzeugt.
- Es sei  $u = (u_0u_1)$  der Unterausdruck, dann wird  $m_u = (m_{u_0}m_{u_1})$  erzeugt.
- Es sei  $u = yu_0\bar{y}$ ,  $y \in A$ ,  $\bar{y} \in \bar{A}$  ein Unterausdruck, dann wird  $m_u = Y$  und eine Produktion

$$Y \longrightarrow ym_{u_0}\bar{y}$$

erzeugt. Wenn der Unterausdruck  $u_0 = \lambda$  ist, dann ist auch das Inhaltsmodell der Produktion  $m_{u_0} = \lambda$ . Wenn  $u_0 \neq \lambda$ , dann wird aus  $u_0$  das Inhaltsmodell  $m_{u_0}$  erzeugt.

- Es sei  $u$  ein XML Ausdruck oder ein Unterausdruck der Form  $(yvx)^n w(\bar{xv}\bar{y})^n$ , wobei  $x$  und  $y$  Zeichen aus  $A$  sind und  $v$  eine Sequenz von Unterausdrücken und Zeichen aus  $A$ . So werden wie oben beschrieben für  $y$ ,  $w$  und  $v$  Nichtterminalzeichen und Produktionen erzeugt. Es sei  $Y$  das Nichtterminalzeichen, das durch die Produktion,  $Y \longrightarrow ym\bar{y}$ ,  $u$  erzeugt und  $\alpha$  ein regulärer Ausdruck über Nichtterminalzeichen, der  $w$  erzeugt, dann wird für  $x$  die Produktion

$$X \longrightarrow xY|\alpha\bar{x}$$

konstruiert.

□

**Satz 4.1.5** Es kann nicht jede Sprache, die durch eine balancierte Grammatik erzeugt wird, durch einen XML Ausdruck erzeugt werden.

*Beweis.* Es können zum Beispiel keine selbsteinbettenden Nichtterminale behandelt werden, die mit einem Stern \* versehen sind. Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit der einzigen Produktion:

$$S \longrightarrow xS^*\bar{x}.$$

Angenommen, die Sprache der balancierten Grammatik könnte durch einen XML Ausdruck erzeugt werden. Wird  $L(G)$  als Baum betrachtet, dann kann ein Knoten des Baums beliebig viele Nachfolger und beliebig viele nachfolgende Ebenen besitzen.

Ein XML Ausdruck kann beliebig viele Nachfolger durch den Kleene Stern ausdrücken. Jedoch folgt aus der Definition der balancierten Grammatik, dass in der beschriebenen Sprache  $L(G)$  auf jeder neuen Ebene wieder beliebig viele Knoten mit beliebig vielen Nachfolgern erzeugt werden können. Dies muss beliebig oft auf beliebig vielen Ebenen möglich sein. Das ist in einem XML Ausdruck nicht möglich, da sich der XML Ausdruck nicht wieder selbst aufrufen kann. Es können nur beliebig viele, aber fest viele Ebenen erzeugt werden, woraus sich ein Widerspruch zur Annahme ergibt. Somit können nicht alle Sprachen, die durch eine balancierte Grammatik erzeugt werden, auch durch einen XML Ausdruck dargestellt werden.  $\square$

## 4.2 Einordnung in die Chomsky Hierarchie

Die in Kapitel 2 definierten XML-artigen Grammatiken werden nun in die Chomsky Hierarchie eingeordnet. Darüber hinaus wird eine Hierarchie der XML-artigen Sprachfamilien nachgewiesen und schließlich wird überprüft, ob die Sprachen eine Teilfamilie der nichtdeterministisch kontextfreien Sprachen oder der deterministisch kontextfreien Sprachen sind. Dies wird eine entscheidende Rolle für die Parser spielen, da deterministisch kontextfreie Sprachen in linearer Zeit geparkt werden können.

### 4.2.1 XML-artige Sprachen und reguläre Sprachen

Es ist intuitiv klar, dass die balancierten Grammatiken kontextfreie Sprachen beschreiben können und somit nicht zu den regulären Sprachen gehören.

**Satz 4.2.6** Die balancierten Sprachen sind keine Teilmenge der regulären Sprachen.

*Beweis.* Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit einer einzigen Produktion

$$S \longrightarrow aS^*\bar{a}.$$

Intuitiv wird klar, dass die erzeugte Sprache ( $L = \{a^m \bar{a}^m \mid m \in \mathbb{N}\}$ ) nicht durch eine reguläre Grammatik erzeugt werden kann.

Dies wird mittels Pumping Lemma bewiesen. Die Pumping Lemma Konstante sei  $n$ . Wir wählen das Wort  $a^n \bar{a}^n$  und die Länge  $2n$ . Aufgrund der 1. Bedingung des Pumping Lemmas darf  $v$  nicht leer sein, aufgrund der 2. Bedingung kann es nur aus  $a$ 's bestehen.

Die 3. Bedingung erlaubt, dass auch die Wörter  $a^{n-|v|} \bar{a}^n$  zulässig sind. Dies ist ein Widerspruch zur Definition der Sprache. Der Ausdruck ist nicht regulär. Die balancierten Sprachen sind daher keine Teilmenge der regulären Sprachen.  $\square$

**Satz 4.2.7** Die balancierten Sprachen sind keine Obermenge der regulären Sprachen.

*Beweis.* Es sei  $L$  eine reguläre Sprache und  $a \in L$  ein Terminalzeichen.  $a$  kann nicht durch eine balancierte Grammatik erzeugt werden, da zu jedem  $a$  ein Gegenstück  $\bar{a}$  erzeugt werden muss. Die Sprache  $L$  kann somit nicht durch eine balancierte Grammatik erzeugt werden.  $\square$

In [BB02a] wird allerdings gezeigt, dass sehr viele balancierte Sprachen zu den regulären Sprachen gehören. Diese Untermenge der balancierten Sprachen wird dort auch untersucht.

**Korollar 4.2.1** Die XML Sprachen, die einfachen balancierten Sprachen und die konkurrenzverhindernden Sprachen sind keine Obermengen der regulären Sprachen.

Da die XML Sprachen, die einfachen balancierten Sprachen und die konkurrenzverhindernden balancierten Sprachen die Beispielsprache aus Satz 4.2.7 beinhalten, gehören diese auch nicht zu den regulären Sprachen.

Das nächste Unterkapitel zeigt, dass die Grammatiken zu den kontextfreien Grammatiken gehören.

## 4.2.2 Auflösung der Struktur von balancierten Grammatiken

In diesem Abschnitt soll gezeigt werden, dass die balancierten Grammatiken eine Unterklasse der kontextfreien Grammatiken darstellen. Dies ist nicht sofort ersichtlich, da die balancierten Grammatiken Produktionen enthalten, deren rechte Seite einen regulären Ausdruck über Nichtterminalzeichen enthält. Es kann somit unendlich viele kontextfreie Produktionen geben.

Nun werden die rechten Seiten der Produktionen in eine endliche Form gebracht. Aus der Zuordnung zu den kontextfreien Grammatiken folgt, dass die balancierten Grammatiken durch Kellerautomaten erkannt werden können.

**Satz 4.2.8** Die balancierten Sprachen gehören zu den kontextfreien Sprachen.

*Beweis.* Gegeben sei eine balancierte Grammatik  $G = (N, T, S, P)$ .  $P$  sind Produktionen der Form  $X \rightarrow xm\bar{x}$ . Da  $m$  aus einem regulären Ausdruck aus Nichtterminalen besteht, kann das Inhaltsmodell aus einer unbeschränkten Folge von Nichtterminalen bestehen. Diese unbeschränkten Inhaltsmodelle werden aufgelöst und damit gezeigt, dass die definierte Sprache kontextfrei ist.

Die Produktionen werden aufgelöst, so dass keine Symbole der Menge  $\{*, (, ), |\}$  enthalten sind. Die folgenden Schritte werden wiederholt, bis alle rechten Seiten der Produktionen endlich sind. Es gilt  $X \in N, u, v, w, z \in \text{Reg}(N), x, \bar{x} \in T$ :

1. Es werden alle Gruppierungen aufgelöst. Es sei  $X \rightarrow xu(w)v\bar{x}$  eine Produktion, dann wird eine neue Produktion  $A \rightarrow w, A \in N$  eingeführt und die ursprüngliche Produktion durch  $X \rightarrow xAv\bar{x}$  ersetzt.
2. Es sei  $X \rightarrow xuw^*v\bar{x}$  eine Produktion, dann werden zwei Produktionen  $A \rightarrow wA, A \in N$  und  $A \rightarrow \lambda$  erzeugt. Zusätzlich wird  $X \rightarrow xAv\bar{x}$  erzeugt und die ursprüngliche Produktion gelöscht.
3. Es sei  $X \rightarrow xu(v|w)z\bar{x}$  eine Produktion, dann werden die Produktionen  $A \rightarrow v, A \in N$  und  $A \rightarrow w$  erzeugt. Zusätzlich wird  $X \rightarrow xAz\bar{x}$  eingeführt und die ursprüngliche Produktion gelöscht.

Nach der Umformung existieren nur noch endlich viele kontextfreie Produktionen. Damit ist gezeigt, dass jede balancierte Grammatik eine kontextfreie Sprache erzeugt.  $\square$

Zur Anschauung sei folgendes Beispiel verwendet:

**Beispiel 4.2.19**

Gegeben sei eine balancierte Grammatik:

$G = (\{S, C, D, E\}, \{a, \bar{a}, c, \bar{c}, d, \bar{d}, e, \bar{e}\}, \{S\}, P)$  mit den Produktionen

$$\begin{aligned} S &\rightarrow a(CD)^*E\bar{a} \\ C &\rightarrow c(S|\lambda)\bar{c} \\ D &\rightarrow d\bar{d} \\ E &\rightarrow e\bar{e}. \end{aligned}$$

Es entsteht folgende Grammatik  $G_1 = (\{S, S_1, S_2, A_1, A_2, C, D, E_1, E\}, \{a, \bar{a}, c, \bar{c}, d, \bar{d}, e, \bar{e}\}, \{S_1\}, P_1)$ :

1. Schritt	2. Schritt	3. Schritt
		$S \rightarrow aA_2E_1\bar{a}$
	$S \rightarrow aA_2E_1\bar{a}$	$A_2 \rightarrow A_1A_2$
	$A_2 \rightarrow A_1A_2$	$A_2 \rightarrow \lambda$
$S \rightarrow aA_1^*E^*\bar{a}$	$A_2 \rightarrow \lambda$	$A_1 \rightarrow CD$
$A_1 \rightarrow CD$	$A_1 \rightarrow CD$	$E_1 \rightarrow EE_1$
$E \rightarrow e\bar{e}$	$E_1 \rightarrow EE_1$	$E_1 \rightarrow \lambda$
$C \rightarrow cS_1\bar{c}$	$E_1 \rightarrow \lambda$	$E \rightarrow e\bar{e}$
$S_1 \rightarrow S \lambda$	$E \rightarrow e\bar{e}$	$C \rightarrow cS_1\bar{c}$
$D \rightarrow d\bar{d}$	$C \rightarrow cS_1\bar{c}$	$S_1 \rightarrow S_2$
	$S_1 \rightarrow S \lambda$	$S_2 \rightarrow \lambda$
	$D \rightarrow d\bar{d}$	$S_2 \rightarrow S$
		$D \rightarrow d\bar{d}$

Im ersten Schritt werden alle Klammern aufgelöst. Im zweiten Schritt werden die Sterne ersetzt. Im dritten Schritt wird das Zeichen | ersetzt.

Nun muss eine genauere Einordnung vorgenommen werden. So wird überprüft, ob die balancierten Sprachen in den deterministisch kontextfreien Sprachen enthalten sind und welche Rolle die anderen Grammatiktypen der XML-artigen Grammatiken spielen.

*kontextfreie Sprachen*

⋮

*balancierte Sprachen* ... *reguläre Sprachen*

**Bild 4.2.1:** Erste Einordnung in die Chomsky Hierarchie.

### 4.2.3 Balancierte Grammatiken und die Chomsky Hierarchie

Ein erstes Ziel ist es, aus einer beliebigen balancierten Grammatik einen deterministischen Kellerautomaten zu konstruieren.

Ein wichtiger Schritt ist die Konstruktion eines erkennenden Kellerautomaten für jedes Nichtterminalzeichen. Da die Inhaltsmodelle aus regulären Ausdrücken bestehen, muss gezeigt werden, wie die Sprachen zweier Automaten, die eine balancierte Sprache erkennen, vereinigt, konkateniert, iteriert und gruppiert werden. Zuletzt müssen noch selbsteinbettende Nichtterminalzeichen behandelt werden. Diese werden vorerst in der Konstruktion nicht berücksichtigt. Zum Schluss werden die Überführungen, in denen selbsteinbettende Nichtterminalzeichen enthalten sind, ersetzt.

Die Konstruktion verläuft in zwei Schritten. Zuerst wird aus der indizierten

balancierten Grammatik ein nichtdeterministischer Kellerautomat konstruiert. Danach wird der nichtdeterministische Kellerautomat in einen deterministischen Kellerautomaten umgewandelt.

In der Konstruktion wird die Eigenschaft der balancierten Sprachen genutzt, worin zu jedem Zeichen  $a \in A$  aus der Eingabe genau ein Zeichen  $\bar{a} \in \bar{A}$  existiert. Die Zuordnung von Zeichen  $a \in A$  und Zeichen  $\bar{a} \in \bar{A}$  wird spezifiziert. Es sei  $G' = (N, T, S', P)$  eine indizierte balancierte Grammatik und

$$X \longrightarrow xm\bar{x}, X \in N, x, \bar{x} \in T, m \in \text{Reg}(N)$$

eine Produktion aus  $P$ . Im erkennenden nichtdeterministischen Automaten wird für das gelesene Zeichen  $x$ , das Zeichen  $(\alpha, X, x)$  in den Keller gelegt.  $\alpha$  ist dabei die Menge der Startzustände der Kellerautomaten, welche die Sprachen von Nichtterminalzeichen beschreiben, die das Inhaltsmodell  $m$  beinhaltet, plus der Zustände, aus denen das Zeichen  $\bar{x}$  dieser Produktion erkannt wird. Diese Menge ist eindeutig. Da die Grammatik indiziert ist, kann ein Startzustand nur dann in zwei verschiedenen Kellerzeichen enthalten sein, wenn ein Nichtterminalzeichen selbst einbettend auftritt. Dann ist jedoch das Inhaltsmodell gleich und somit ist eindeutig, welcher Zustand zu welchem Kellerzeichen gehört. Somit ist im Kellerzeichen enthalten, welche Produktion noch abzuarbeiten ist und welche Automaten dafür in Frage kommen, das Inhaltsmodell abzuarbeiten.

Im deterministischen Kellerautomaten muss es nicht eindeutig sein, welches Kellerzeichen eingekellert wird (bei Mehrdeutigkeit). In diesem Fall wird ein neues Symbol, welches die Menge aller möglichen Kellersymbole enthält, in den Keller gelegt.

**Satz 4.2.9** Balancierte Sprachen werden durch deterministische Kellerautomaten erkannt.

*Beweis.* Es sei  $G' = (N, T, S', P)$  eine indizierte balancierte Grammatik. Die Grammatik muss indiziert sein, da die Nichtterminalzeichen in verschiedenen Produktionen vorkommen und für jedes Nichtterminal ein erkennender Kellerautomat konstruiert wird. Wenn die balancierte Grammatik nicht indiziert würde, würde für ein Nichtterminalzeichen  $X$ , welches in zwei verschiedenen Inhaltsmodellen auftritt, nur ein Kellerautomat erzeugt.

Zuerst wird aus der Grammatik  $G'$  ein nichtdeterministischer Kellerautomat  $\mathcal{D}' = (S', T', G', \Delta, s'_0, \delta', F')$  konstruiert und dieser wird in einen deterministischen Kellerautomaten  $\mathcal{D} = (S, T, G, \Delta, s_0, \delta, F)$  umgewandelt. Vorerst wird für jede Produktion ein nichtdeterministischer Kellerautomat, der die Sprache der Produktion erkennt, konstruiert. Es sei

$$X \longrightarrow xm\bar{x}, m \in \text{Reg}(N), X \in N, x, \bar{x} \in T$$

eine Produktion aus der Grammatik  $G'$ .

Im Folgenden sind  $s_{(i, \mathcal{X})}$  Zustände des Automaten  $\mathcal{D}'_X$ .  $a, x, y, \bar{a}, \bar{x}, \bar{y}$  sind

Terminalzeichen.  $\overline{(\alpha, X, x)}$  sind Kellerzeichen, die eingekellert werden, wenn eine Überführung für das Zeichen  $x$  der Produktion für  $X$  benutzt wird.  $\alpha$  ist eine Menge von Zuständen, die später definiert wird.

Zuerst sei der Fall betrachtet, dass  $m = \lambda$  ist. Dann ist der erkennende Teilautomat  $\mathcal{D}'_X$  von  $\mathcal{D}'$  durch folgende Überführungen definiert:

$$\begin{aligned} (s_{(0,\mathcal{X})}, x, \Delta) &\longrightarrow (s_{(1,\mathcal{X})}, \overline{(\{s_{(1,\mathcal{X})}\}, X, x)\Delta}) \\ (s_{(1,\mathcal{X})}, \bar{x}, \overline{(\{s_{(1,\mathcal{X})}\}, X, x)}) &\longrightarrow (s_{(e,(X,x))}, \lambda) \end{aligned}$$

Es ist hierbei wichtig zu beachten, dass für ein Zeichen  $x \in T$  ein Zeichen in den Keller gelegt wird, welches eindeutig die Produktion identifiziert.

Im Fall, dass eine Produktion  $X \longrightarrow xm\bar{x}$  mit dem Inhaltsmodell  $m \neq \lambda$  existiert, muss das Inhaltsmodell mitbehandelt werden. Das Inhaltsmodell besteht aus einem regulären Ausdruck über Nichtterminalzeichen. Es wird beschrieben, wie die Automaten der Nichtterminalzeichen zusammengesetzt werden und dabei insgesamt ein nichtdeterministischer Kellerautomat entsteht. Die Konstruktion erfolgt induktiv über den Aufbau der regulären Ausdrücke. Dabei genügen allgemein die Konkatenation zweier regulärer Ausdrücke, die Vereinigung, die Gruppierung und Vervielfältigung mit  $*$ . Zuletzt müssen noch selbsteinbettende Nichtterminale behandelt werden.

Die Gruppierung besagt, dass der eingeklammerte Teilausdruck und somit sein erkennender nichtdeterministischer Kellerautomat als ein Ausdruck oder ein Automat anzusehen sind. Hier muss also nur beachtet werden, dass bei Vereinigung und Vervielfältigung der komplette Automat verwendet wird.

Für die Produktion  $X \longrightarrow xm\bar{x}$ ,  $m \neq \lambda$  wird der Automat  $\mathcal{D}'_X$  mit folgenden Überführungen erzeugt:

$$\begin{aligned} (s_{(0,\mathcal{X})}, x, \Delta) &\longrightarrow (s_{(1,\mathcal{X})}, \overline{(\{s_{(2,\mathcal{X})}\} \cup \beta, X, x)\Delta}) \\ (s_{(1,\mathcal{X})}, m, \overline{(\{s_{(2,\mathcal{X})}\} \cup \beta, X, x)}) &\longrightarrow (s_{(2,\mathcal{X})}, (\{s_{(2,\mathcal{X})}\} \cup \beta, X, x)\Delta) \\ (s_{(2,\mathcal{X})}, \bar{x}, (\{s_{(2,\mathcal{X})}\} \cup \beta, X, x)\Delta) &\longrightarrow (s_{(e,(X,x))}, \lambda) \end{aligned}$$

Es sei  $V$  die Menge der Nichtterminalzeichen, die in  $m$  enthalten sind. Dann ist  $\beta$  die Menge der Startzustände der Kellerautomaten, die die Sprachen der Nichtterminalzeichen aus  $V$  erkennen. Die Startzustände werden im Kellerzeichen kodiert, damit bei der Umwandlung in einen deterministischen Kellerautomaten für jeden Startzustand eines Teilautomaten von  $\mathcal{D}'$  klar ist, für welches Kellerzeichen aus der Kombination mehrerer Kellerzeichen, eine Überführung definiert wird.

In dem Kellerautomat  $\mathcal{D}'_X$  gibt es eine Überführung, in der das Inhaltsmodell  $m$  enthalten ist. Diese Überführung wird noch ersetzt, wenn der erkennende Kellerautomat konstruiert ist.

Selbsteinbettende Nichtterminale werden in den Überführungen erst ganz zum Schluss behandelt.

Wenn der nichtdeterministische Kellerautomat, der die Sprache von  $m$  erkennt, eingefügt wird, muss an dem einzufügenden Automaten noch etwas verändert werden.

Es sei  $\mathcal{M}$  der Kellerautomat, der die Sprache von  $m$  erkennt. Es sei

$$(s_{(0,\mathcal{M})}, y, \Delta) \longrightarrow (s_{(1,\mathcal{M})}, \overline{(\gamma, Y, y)}\Delta)$$

eine Startüberführung des Automaten, welche sich noch auf das Kellerendsymbol bezieht. Jedoch ist der Automat  $\mathcal{M}$  jetzt Teil eines „großen“ Automaten, nämlich des Automaten  $\mathcal{D}_X$ , in den sie eingefügt worden sind.

Die Überführung wird nun zu

$$(s_{(0,\mathcal{M})}, y, \overline{(\alpha, X, x)}) \longrightarrow (s_{(1,\mathcal{M})}, \overline{(\gamma, Y, y)}(\alpha, X, x))$$

geändert.

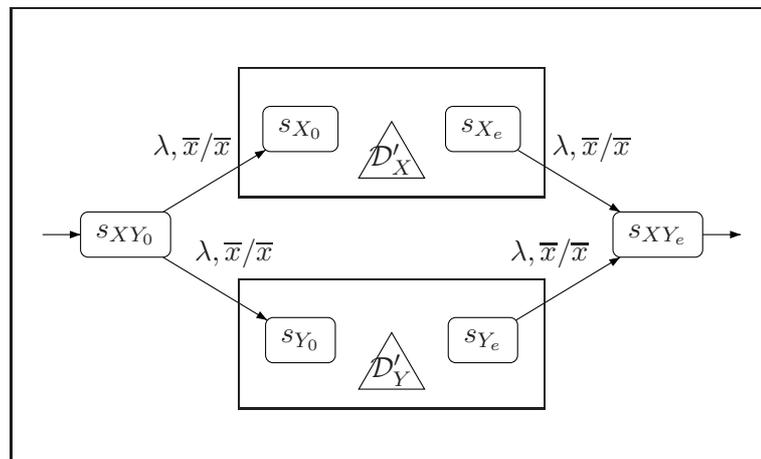
Hier ist auch der Unterschied zu der Konstruktion eines RegDyck Automaten aus einer RegDyck Grammatik erkennbar. In einem RegDyck Automaten ist die Überführungsfunktion des Automaten für ein Zeichen  $a \in A$  für jedes Kellersymbol definiert. Hinzu kommt, dass bei der Konstruktion aus der balancierten Grammatik die Kellerzeichen eindeutig einer Produktion zuzuordnen sind und die Startzustände der Kellerautomaten des Inhaltsmodells enthalten sind.

Die folgende Konstruktion ähnelt der Konstruktion aus Kapitel 3.

Hier wird die Konstruktion der nichtdeterministischen Kellerautomaten über den induktiven Aufbau der regulären Ausdrücke für die Inhaltsmodelle gezeigt.

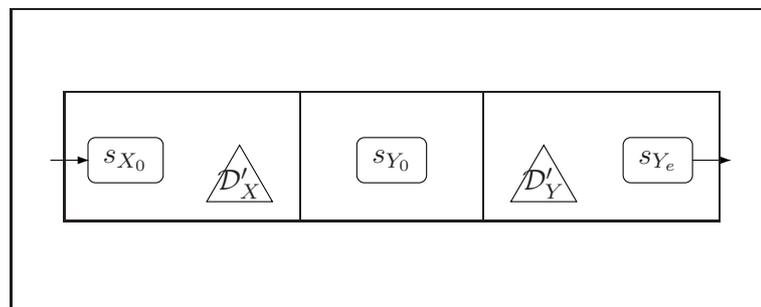
Zuerst wird die Vereinigung beschrieben.

1. Für den regulären Ausdruck  $X|Y$  wird  $\mathcal{D}'_{X|Y}$  konstruiert.



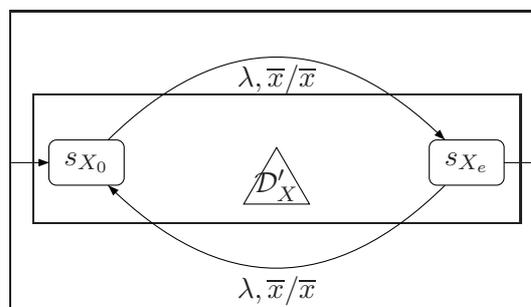
Hierbei wird über den gemeinsamen Startzustand  $s_{XY_0}$  von  $\mathcal{D}'_{X|Y}$  durch einen  $\lambda$ -Übergang in die Startzustände von  $\mathcal{D}'_X$  und  $\mathcal{D}'_Y$  verwiesen. Es existiert ein gemeinsamer Endzustand  $s_{XY_e}$ .  $\bar{x}$  sind dabei Kellerzeichen der Form  $(\alpha, X, x)$ .

2. Für den regulären Ausdruck  $XY$  wird der zusammengesetzte Kellerautomat  $\mathcal{D}'_{XY}$  konstruiert.



Hierbei wird der Endzustand  $s_{X_e}$  des Kellerautomaten  $\mathcal{D}'_X$  zum Startzustand  $s_{Y_0}$  des Automaten  $\mathcal{D}'_Y$ .

3. Für den regulären Ausdruck  $X^*$  wird der Kellerautomat  $\mathcal{D}'_{X^*}$  konstruiert.



Hierbei wird ein  $\lambda$ -Übergang vom Startzustand  $s_{X_0}$  in den Endzustand  $s_{X_e}$  erzeugt und umgekehrt.

Ein Fall muss gesondert betrachtet werden. Es sei  $X$  ein Nichtterminalzeichen, für das ein Kellerautomat  $\mathcal{D}'_X$  konstruiert werden soll. Während der Konstruktion des Kellerautomaten  $\mathcal{D}'_X$  tritt das Nichtterminalzeichen  $X$  nochmal auf. Da die balancierte Grammatik indiziert ist, kann es nur innerhalb der Konstruktion des Kellerautomaten  $\mathcal{D}'_X$  auftreten und sonst nicht. Das Nichtterminalzeichen  $X$  ist somit ein selbsteinbettendes Nichtterminalzeichen. Der Automat, der für  $X$  nun ein weiteres mal konstruiert werden müsste, ist also der selbe Automat  $\mathcal{K}_X$ . Er wird jedoch nicht erneut konstruiert, da die Konstruktion wegen der Selbsteinbettung wieder und wieder durchgeführt werden müsste. Für alle Konstruktionen (siehe oben Punkt 1 bis 3), die für das Inhaltsmodell durchgeführt werden müssen, muss lediglich der Startzustand und der Endzustand eines Automaten bekannt sein. Beide sind von dem Kellerautomaten  $\mathcal{D}'_X$  bekannt. Somit kann dieser für die Konstruktion verwendet werden.

Es bleibt zu zeigen, dass alle  $\lambda$ -Übergänge ersetzt werden können, um den nichtdeterministischen Kellerautomaten einfacher in einem deterministischen Kellerautomaten umzuwandeln.

Es seien  $s_k, s_j \in S_X$  zwei Zustände eines Automaten  $\mathcal{D}'_X$ . Wenn eine Überführung

$$(s_k, \lambda, \bar{z}) \longrightarrow (s_j, \bar{z}), \bar{z} \in G_X,$$

existiert, werden für alle Überführungen der Form

$$(s_j, a, \bar{z}) \longrightarrow (s_n, \bar{a}\bar{z}), a \in A_X, \bar{a}, \bar{z} \in G_X, s_n \in S_X,$$

und

$$(s_j, \bar{a}, \bar{a}) \longrightarrow (s_n, \lambda)$$

Überführungen

$$(s_k, a, \bar{z}) \longrightarrow (s_n, \bar{a}\bar{z})$$

oder

$$(s_k, \bar{a}, \bar{a}) \longrightarrow (s_n, \lambda)$$

erzeugt.

Für jede Produktion der Form

$$X \longrightarrow xm\bar{x}, X \in N, x, \bar{x} \in T, m \in \text{Reg}(N)$$

wurde ein erkennender nichtdeterministischer Kellerautomat konstruiert. Das Inhaltsmodell wurde aus den einzelnen Automaten der Nichtterminalzeichen nach einer Abwandlung der Thompson Methode zusammengesetzt. Wenn keine Selbsteinbettung auftritt, ist klar, dass  $L(G) = L(\mathcal{D}')$ . Bei Selbsteinbettung

$$X \Rightarrow^+ uXv$$

erkennt der Teilautomat  $\mathcal{D}'_X$  von Automaten  $\mathcal{D}'$ , die Sequenzen  $u$ , bzw.  $v$  durch Schleifen, indem er wieder in sich selbst verweist.

Wir nehmen an, es gilt  $L(G) \neq L(\mathcal{D})$ . Für Grammatiken, die keine selbst-einbettenden Nichtterminalzeichen besitzen, ist aufgrund der Konstruktion klar, dass die Sprachen, die sie erzeugen, gleich sind. Damit  $L(G) \neq L(\mathcal{D})$  gilt, müssten selbst-einbettende Nichtterminalzeichen in der Grammatik enthalten sein. Wenn der Kellerautomat  $\mathcal{D}'$  gleich viele Sequenzen von  $u$  und  $v$  durchläuft, erkennt er die Sprache, die das selbst-einbettende Nichtterminalzeichen  $X$  erzeugt. Wenn dies nicht der Fall ist, muss ein anderer Teilautomat für ein anderes Nichtterminalzeichen Sequenzen erzeugen, die  $u$  oder  $v$  zu einem Dyckprime ergänzen, da der Keller geleert werden muss.

Es muss also ein zweites selbst-einbettendes Nichtterminalzeichen  $Y$  geben mit

$$Y \Rightarrow^+ u'Yv'$$

und gelten  $u'v \in D_T$  oder  $uv' \in D_T$ .  $\mathcal{D}'$  kellert für jedes Zeichen  $a \in A$  ein Zeichen ein, das eindeutig die Produktion, aus der es entstanden ist, identifiziert. Weiter ist nicht mehr jedes Zeichen  $a \in A$  für alle Kellerzeichen definiert, sondern nur noch, wenn das Kellerzeichen auf dem Keller liegt, welches in der Produktion steht, in der  $a$  enthalten ist.

Der Kellerautomat kann diese Sequenz nicht erkennen, da im Kellerzeichen die Information enthalten ist, dass jedes Zeichen aus  $u$  von der Produktion  $X$  entstanden ist und jedes Zeichen  $u'$  von dem Nichtterminalzeichen  $Y$ . Wenn  $\mathcal{D}'$  nicht gleich viele Schleifen von  $u$  durchläuft wie von  $v$ , dann kann  $\mathcal{D}'$  seinen Keller nicht mehr leeren. Deshalb müssen gleich viele Schleifen von  $u$  und  $v$  durchlaufen werden. Daraus folgt ein Widerspruch zu der Annahme.

Der nichtdeterministische Kellerautomat  $\mathcal{D}' = (S', T', G', \Delta, s'_0, \delta', F')$  ist nun erzeugt. Es fehlt noch die Umwandlung in einen deterministischen Kellerautomaten. Wichtig zu erkennen ist, dass der Kellerautomat  $\mathcal{D}'$  zwar nichtdeterministisch ist, jedoch für ein Eingabezeichen klar ist, wenn ein Zeichen  $a \in A$  in der Eingabe steht, wird ein Zeichen eingekellert. Darüber hinaus wird, wenn ein Zeichen  $\bar{a} \in \bar{A}$  in der Eingabe steht, ein Zeichen aus dem Keller entfernt.

Nun werden zwei verschiedene Konstruktionen durchgeführt. Es wird die Potenzmenge der Zustände  $2^{S'}$  gebildet und ähnlich wie bei endlichen Automaten die Umwandlung der Zustandsübergänge deterministisch gemacht.

Weiter wird, wenn in mehrere Zustände verwiesen wird, eine Kombination aus allen Kellerzeichen, die im nichtdeterministischen Fall in den Keller gelegt würden, in den Keller gelegt. Also wird auch die Potenzmenge der Kellerzeichen gebildet  $2^{G'}$ . Der deterministische Kellerautomat  $\mathcal{D} = (S, T, G, \Delta, s_0, \delta, F)$  wird folgermaßen aus  $\mathcal{D}'$  konstruiert:

- $S = 2^{S'} \cup \{m_{S''} \mid S'' \in 2^{S'}\}$  ist die Menge der Zustände.
- $T = T'$  ist die Menge der Eingabesymbole.
- $G = 2^{G'}$  ist die Menge der Kellersymbole.
- $s_0 = s'_0$  ist der Startzustand.
- $F$  ist die Menge der Endzustände. Es sei

$$X \longrightarrow xm\bar{x}, m \in \text{Reg}(N), X \in N, x, \bar{x} \in T$$

eine Produktion aus der Grammatik  $G'$ , wobei  $X$  ein Startzeichen der Grammatik ist. Dann ist

$$s_{(e,(X,x))}$$

ein Endzustand des erkennenden Automaten.

Alle Zustände  $S'', s_{(e,(X,x))} \in S'', S'' \in S$  sind Endzustände.

- $\delta$  ist die Überföhrungsfunktion. Im Folgenden gilt:  $R, R' \in S, a, \bar{a} \in T, C, C' \in G$ .

1. Es werden Überföhrungen

$$\begin{aligned} \delta(R, a, C) &= (R', C'C) \\ R' &= \{s' \mid s \in R \wedge \overline{(\alpha, X, x)} \in C \wedge s \in \alpha \wedge \\ &\quad \delta(s, a, \overline{(\alpha, X, x)}) = (s', g'), s' \in S', g' \in G'\} \\ C' &= \{g' \mid s \in R \wedge \overline{(\alpha, X, x)} \in C \wedge s \in \alpha \wedge \\ &\quad \delta(s, a, \overline{(\alpha, X, x)}) = (s', g'), s' \in S', g' \in G'\} \end{aligned}$$

erzeugt.

Die Überföhrungen definieren das Verhalten des deterministischen Kellerautomaten für Eingabezeichen  $a \in A$ . Der Kellerautomat  $\mathcal{D}$  reagiert darauf, indem er in die Menge der Zustände verweist, in die der nichtdeterministische Kellerautomat  $\mathcal{D}'$  für dieses Eingabezeichen verweist. Zusätzlich werden alle möglichen Kellerzeichen, die der Automat  $\mathcal{D}'$  in den Keller legt, als Vereinigung in den Keller gelegt.

2. Es werden Überföhrungen

$$\begin{aligned} \delta(R, \bar{a}, C) &= (m_{R'}, \lambda) \\ R' &= \{s' \mid s \in R \wedge \overline{(\alpha, X, a)} \in C \wedge s \in \alpha \wedge \\ &\quad \delta(s, \bar{a}, \overline{(\alpha, X, a)}) = (s', \lambda), s' \in S'\} \end{aligned}$$

erzeugt.

Die Überführungen definieren das Verhalten des deterministischen Kellerautomaten  $\mathcal{D}$  für Eingabezeichen  $\bar{a} \in \bar{A}$ . Der Kellerautomat  $\mathcal{D}$  überprüft dabei, ob ein Kellerzeichen, in dessen Tupel  $\overline{(\alpha, X, a)}$  das Zeichen  $a$  enthalten ist, an oberster Stelle im Keller liegt. Wenn dies der Fall ist, verweist er in einen Zwischenzustand, in dem alle Zustände, in die der nichtdeterministische Kellerautomat  $\mathcal{D}'$  verweist, enthalten sind.

3. Es werden Überführungen

$$\begin{aligned} \delta(m_{R'}, \lambda, C) &= (R', C') \\ C' &= \{\overline{(\alpha, X, x)} \mid s \in R' \wedge s \in \alpha \wedge \overline{(\alpha, X, x)} \in C\} \end{aligned}$$

erzeugt.

Diese Überführungen definieren das Verhalten des deterministischen Kellerautomaten  $\mathcal{D}$  für die Zwischenzustände  $m_{T'}$ . In diese Zustände wird verwiesen, wenn der Automat  $\mathcal{D}$  ein Zeichen aus dem Keller entfernt. Durch die Zwischenzustände wird das oberste Kellerzeichen aktualisiert. In einem Kellerzeichen des deterministischen Kellerautomaten sind alle Kellerzeichen vereinigt, die der nichtdeterministische Kellerautomat für ein einziges Eingabezeichen in den Keller legt. Der deterministische Kellerautomat durchläuft alle Wege des nichtdeterministischen Kellerautomaten gleichzeitig. Wenn der deterministische Kellerautomat erkennt, dass ein Weg falsch war, wird beim deterministischen Kellerautomat  $\mathcal{D}$  durch die Zwischenzustände das Kellerzeichen aktualisiert. Somit kodiert der Automat  $\mathcal{D}$ , dass die Kellerzeichen nur durch die bestimmte Menge aus den ursprünglichen Produktionen entstanden sind.

Es ist zu beachten, dass sowohl die Keller des nichtdeterministischen Kellerautomaten, als auch der Keller des deterministischen Kellerautomaten immer die gleiche Höhe haben, da für ein Zeichen  $a \in A$  genau ein Zeichen in den Keller geschrieben wird und für ein Zeichen  $\bar{a}$  immer ein Zeichen aus dem Keller entfernt wird.

Es sei  $\pi_i, i \in \{1, 2, 3\}$  die Projektion auf die  $i$ -te Komponente einer Konfiguration. Es seien  $\mathcal{D} = (S, T, G, \Delta, s'_0, \delta, F)$  und  $\mathcal{D}' = (S', T', G', \Delta, s'_0, \delta', F')$  die Kellerautomaten aus der obigen Konstruktion.

Die Gleichheit  $L(\mathcal{D}) = L(\mathcal{D}')$  wird induktiv über die Länge der Eingabe  $w \in T^*$  gezeigt. Zunächst wird behauptet, dass

$$\Delta(s'_0, w, \Delta) = (s, \lambda, \overline{zw''}) \Leftrightarrow \Delta'(s'_0, w, \Delta) = (s, \lambda, \overline{zw'}) , \bar{z} \in \bar{G}, w', w'' \in \bar{G}^* .$$

Es wird die Betrachtungsweise gewählt, dass ein nichtdeterministischer Kellerautomat unter Umständen mehrere Zustandsübergänge gleichzeitig durchläuft. Dabei ist im nichtdeterministischen Kellerautomat der Kellerinhalt  $\overline{zw}$  pro Zeichen  $\overline{a} \in \overline{zw}$  die Vereinigung aus allen Kellerzeichen der möglichen Wege, die der Automat durchläuft.

Die Behauptung ist trivial für  $|w|$ , da in diesem Fall  $w = \lambda$  und  $\Delta(s'_0, \lambda, \Delta) = (s'_0, \lambda, \Delta) = \Delta'(s'_0, w, \Delta)$ . Es sei nun angenommen, die Behauptung gelte für  $|w| \geq i$ . Für das Wort  $wa$  der Länge  $i + 1$  und  $a \in A$  gilt

$$\Delta(s'_0, wa, \Delta) = \delta(\pi_1(\Delta(s'_0, w, \Delta)), a, \pi_3(\Delta(s'_0, w, \Delta))).$$

Aus der Annahme folgt, dass

$$\Delta(s'_0, w, \Delta) = (s, \lambda, \overline{zw''}) \Leftrightarrow \Delta'(s'_0, w, \Delta) = (s, \lambda, \overline{zw''}).$$

Aus der Definition folgt weiter, dass

$$\begin{aligned} \delta(s, a, \overline{z}) &= (s', \lambda, \overline{yzw''}) \Leftrightarrow \delta'(s, a, \overline{z}) = \\ &= \left( \bigcup_{s'' \in s} \pi_1(\delta'(s'', a, \overline{z''})), \lambda, \bigcup_{s'' \in s} \pi_3(\delta'(s'', a, \overline{z''})) \right) = (s', \lambda, \overline{yzw''}), \overline{z''} \in \overline{z}. \end{aligned}$$

Also gilt auch

$$\Delta(s'_0, wa, \Delta) = (s, \lambda, \overline{yzw''}) \Leftrightarrow \Delta'(s'_0, wa, \Delta) = (s, \lambda, \overline{yzw''}).$$

Es sei nun angenommen, die Behauptung gelte für  $|w| \geq i$ . Für das Wort  $w\overline{a}$  der Länge  $i + 1$  und  $\overline{a} \in \overline{A}$  gilt

$$\Delta(s'_0, w\overline{a}, \overline{ay_1w'}) = \delta(\pi_1(\Delta(s'_0, w, \Delta)), \overline{a}, \pi_3(\Delta(s'_0, w, \Delta))).$$

Aus der Annahme folgt, dass

$$\Delta(s'_0, w, \Delta) = (s, \lambda, \overline{ay_1w''}) \Leftrightarrow \Delta'(s'_0, w, \Delta) = (s, \lambda, \overline{ay_1w''}).$$

Aus der Definition folgt weiter, dass

$$\begin{aligned} \delta(s, \overline{a}, \overline{a}) &= (m_{s'}, \lambda, \overline{y_1z'w'}) \longrightarrow \delta(m_{s'}, \lambda, \overline{y_1}) = (s', \lambda, \overline{y_1w'}) \\ &\Leftrightarrow \delta'(s, \overline{a}, \overline{a}) = \left( \bigcup_{s'' \in s} \pi_1(\delta'(s'', \overline{a}, \overline{a})), \lambda, \bigcup_{s'' \in s} \pi_3(\delta'(s'', \overline{a}, \overline{a})) \right). \end{aligned}$$

Hierbei ist erneut zu beachten, dass in den Kellerzeichen die Zustände eindeutig den Kellerzeichen zugeordnet sind.

Also gilt auch

$$\Delta(s'_0, w\overline{a}, \Delta) = (s, \lambda, \overline{y_1w''}) \Leftrightarrow \Delta'(s'_0, w\overline{a}, \Delta) = (s, \lambda, \overline{y_1w''}).$$

Somit ist die Gleichheit  $L(\mathcal{D}) = L(\mathcal{D}')$  gezeigt. □

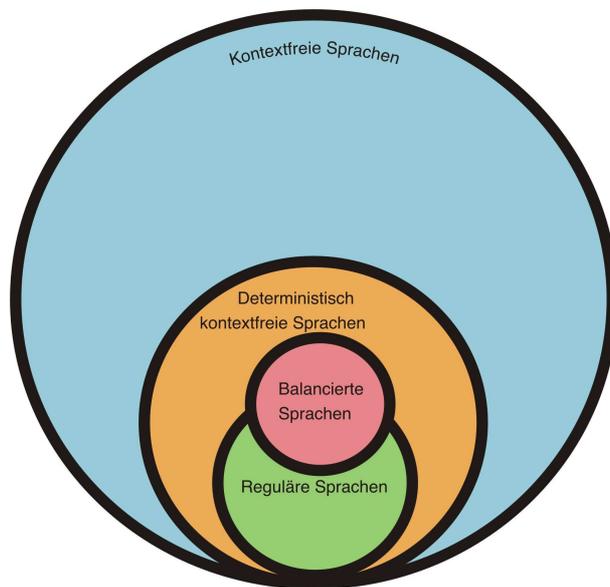


Abbildung 4.1: Genaue Einordnung in die Chomsky Hierarchie.

**Satz 4.2.10** XML Grammatiken gehören zu den  $LL(k)$  Grammatiken.

*Beweis.* Es wird gezeigt, dass dieser Grammatiktyp durch einen Kellerautomaten mit einem  $1$ -lookahead erkannt werden kann.

Gegeben sei eine XML Grammatik  $G = (N, T, S, P)$ . Für  $k$  wird  $1$  gewählt.

Abgesehen von der Einschränkung, dass jede Produktion in  $P$  die Form  $X \rightarrow xm\bar{x}$ ,  $X \in N$ ,  $m \in Reg(N)$ ,  $x, \bar{x} \in T$ , haben muss, sind drei Eigenschaften der XML Grammatiken entscheidend:

1. Es darf keine zwei Produktionen geben, die mit dem selben Terminalzeichen auf der rechten Seite der Produktion beginnen. Das bedeutet es kann eindeutig, anhand des nächsten Eingabesymbols, entschieden werden, welche Produktion benutzt werden muss.
2. Eine weitere Eigenschaft ist, dass die Inhaltsmodelle aller Produktionen aus einem  $1$ -lookahead deterministischem regulären Ausdruck bestehen müssen.
3. Zusätzlich ist jede Produktion der Grammatik  $G$  in doppelter Greibach Normalform.

Die Menge der Terminalzeichen  $a$ , die direkt hinter dem Nichtterminalzeichen  $A$  in einer Ableitung auftreten können  $S \Rightarrow^* \alpha A a \beta$ ,  $\alpha, \beta \in (N \cup T)^*$ ,  $S \in N$ ,  $a \in T$ , werden mit  $FOLLOW(A)$  bezeichnet.

Demgegenüber stehen die Eigenschaften von  $LL(1)$  Grammatiken [ASU99a] [ASU99b]:

Es seien  $A \rightarrow \alpha|\beta$  mit  $\alpha, \beta \in (T \cup N)^*$  zwei unterschiedliche Produktionen in  $P$ :

1. Für kein Terminalzeichen  $a$  können Zeichenketten aus  $\alpha$  und  $\beta$  abgeleitet werden, die beide mit  $a$  beginnen.
2. Nur einer der beiden Ausdrücke  $\alpha$  oder  $\beta$  kann zur leeren Zeichenkette abgeleitet werden.
3. Wenn  $\beta \Rightarrow^* \lambda$ , dann darf nicht  $\alpha$  in eine Zeichenkette beginnend mit  $a \in FOLLOW(A)$  ableiten.

Nach der Umwandlung der Grammatik  $G$  nach Satz 4.2.8, wonach alle rechten Seiten der Produktionen aus einer endlichen Menge bestehen, gibt es für ein Nichtterminal  $X \in N$  folgende Möglichkeiten von Produktionen:

- $X \rightarrow xA_1A_2 \dots A_n\bar{x}$
- $X \rightarrow A_1A_2 \dots A_n$
- $X \rightarrow xA_1A_2 \dots A_n\bar{x}X_i$  und  $X_i \rightarrow \lambda$
- $X \rightarrow xA_1A_2 \dots A_n\bar{x}$  und  $X_i \rightarrow \lambda$ .

für  $x \in A, \bar{x} \in \bar{A}, A_1, A_2, \dots, A_n \in N$ .

Für den ersten Typ von Produktionen  $X \rightarrow xA_1A_2 \dots A_n\bar{x}$  ist klar, welche Produktion  $X$  gewählt wird, da für dieses Nichtterminal  $X$  nur genau diese eine Produktion existiert. Für den zweiten Typ von Produktionen  $X \rightarrow A_1A_2 \dots A_n$  gilt selbiges.

Der dritte und vierte Typ  $X \rightarrow xA_1A_2 \dots A_n\bar{x}X$  und  $X \rightarrow \lambda$  bzw.  $X \rightarrow xA_1A_2 \dots A_n\bar{x}$  und  $X \rightarrow \lambda$  gehorcht der ersten und der zweiten Regel für  $LL(1)$  Grammatiken, da die erste Variante von  $X$  auf jeden Fall mit dem Terminalzeichen  $x$  beginnt und die zweite Variante ein  $\lambda$  ist. Die dritte Regel wird durch die Forderung eingehalten, dass der reguläre Ausdruck 1-deterministisch sein muss und dass es keine zwei Nichtterminale geben darf, deren Produktionen mit einem gleichen Terminalzeichen  $x$  beginnen. Innerhalb der Produktion  $Y \rightarrow yX_1X_2 \dots X_n\bar{y}$  ist eindeutig, dass nach  $X_{i-1} X_i$  das einzige Nichtterminal mit einem  $x \in A$  auf der rechten Seite der Produktion ist. Somit ergibt sich, dass  $x \notin FOLLOW(X_i)$ , falls  $X_i \Rightarrow^* \lambda$ .  $\square$

**Satz 4.2.11** Die einfachen balancierten Grammatiken gehören zu den  $LL(k)$  Grammatiken.

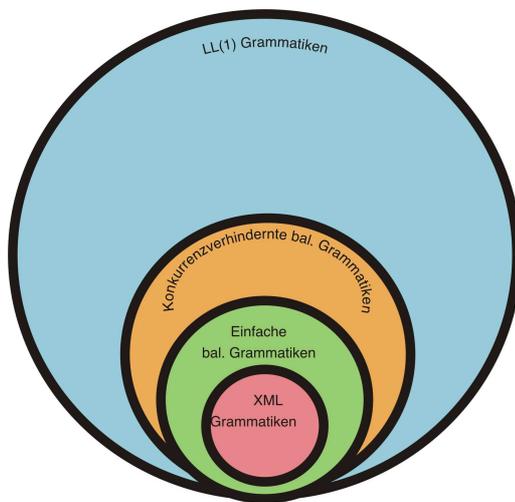


Abbildung 4.2: XML-artige Grammatiken und  $LL(1)$  Grammatiken.

*Beweis.* Da bei den einfachen balancierten Grammatiken keine konkurrierenden Nichtterminale innerhalb einer Produktion auftreten dürfen, ist auch hierbei für ein Terminalzeichen  $x$  innerhalb einer Produktion eindeutig klar, von welchem Nichtterminalzeichen dieses erzeugt wird.  $\square$

**Satz 4.2.12** Die konkurrenzverhindernden balancierten Grammatiken gehören zu den  $LL(k)$  Grammatiken.

*Beweis.* Die Bedingung der konkurrenzverhindernden balancierten Grammatiken war, dass für Nichtterminalsequenzen  $U, V, W$  und zwei konkurrierende Nichtterminale  $A, B \in N$  nicht  $UAV$  und  $UBW$  innerhalb einer Produktion erzeugt werden darf.

Es sei  $x$  ein Terminalzeichen, welches sowohl von  $A$  als auch von  $B$  erzeugt wird. Da nach einer Abarbeitung von einer beliebigen Sequenz  $U$  nicht  $A$  und  $B$  folgen dürfen und für den regulären Ausdruck des Inhaltsmodells die 1-lookahead deterministische Eigenschaft gilt, ist eindeutig klar, welche Produktion für welches Nichtterminalzeichen angewendet werden muss.  $\square$

**Satz 4.2.13** Jede Sprache  $L$ , die von einem RegDyck Automaten erkannt wird, ist deterministisch kontextfrei.

*Beweis.* Der Durchschnitt einer regulären Sprache und einer deterministisch kontextfreien Sprache ist deterministisch kontextfrei.  $\square$

### 4.3 Hierarchie der XML-artigen Grammatiken

In diesem Abschnitt wird eine Hierarchie der XML-artigen Sprachfamilien gezeigt.

Es ist klar, dass die XML Sprachen, die einfachen balancierten Sprachen, die konkurrenzverhindernden balancierten Sprachen und die RegDyck Sprachen Teilmengen der balancierten Sprachen sind. Die Inhaltsmodelle der Produktionen der Untertypengrammatiken sind mit Einschränkungen in Bezug zu den Produktionen in balancierten Grammatiken belegt. Es wird untersucht, ob die Sprachen der Untertypen echte Teilmengen der balancierten Sprachen sind.

**Satz 4.3.14** Es sei  $\mathcal{L}$  die Familie der balancierten Sprachen und  $\mathcal{L}_1$  sei die Familie der konkurrenzverhindernden Sprachen, dann gilt  $\mathcal{L}_1 \subset \mathcal{L}$ .

*Beweis.* Es sei

$$L((x(a(c\bar{c})?a)^*(a(d\bar{d})\bar{a})^*\bar{x}))$$

die Sprache, die von folgenden Produktionen einer balancierten Grammatik  $G = (N, T, \{S\}, P)$  erzeugt wird:

$$P = \left\{ \begin{array}{l} S \longrightarrow xA^*A_1^*\bar{x} \\ A \longrightarrow aC?a \\ A_1 \longrightarrow aD\bar{a} \\ C \longrightarrow c\bar{c} \\ D \longrightarrow d\bar{d} \end{array} \right\}$$

Es lässt sich keine Grammatik für die Sprache finden, welche in der Produktion für  $X$  keine zwei konkurrierenden Nichtterminale enthält, da die beiden Dycksprachen

$$\{(a(c\bar{c})?a)^*\} \text{ und } \{(a(d\bar{d})\bar{a})^*\}$$

zwischen  $x$  und  $\bar{x}$  erzeugt werden.

Um die beiden Dycksprachen zu erzeugen, werden mindestens zwei unterschiedliche Nichtterminale benötigt, die beide mit einem Stern gekennzeichnet sind.

Angenommen, die Grammatik enthält nur ein Nichtterminalzeichen, dann wird davon die Sprache  $\{x(a(c\bar{c})?|(d\bar{d})\bar{a})^*(a(c\bar{c})?|(d\bar{d})\bar{a})^*\bar{x}\}$  erzeugt. Diese Sprache ist jedoch nicht die Ausgangssprache.

Die beiden konkurrierenden Nichtterminale, die die Dycksprachen erzeugen, stehen hintereinander in der Produktion  $X \longrightarrow xA^*A_1^*\bar{x}$ . Die Nichtterminale können beide auftreten oder komplett weggelassen werden.

Es folgt, dass diese Grammatik keine konkurrenzverhindernde balancierte Grammatik ist, da dadurch Sequenzen  $UAW$  und  $UA_1V$ , mit  $U, V, W \in N^*$ , erzeugt werden können. Solche Inhaltsmodelle sind in den konkurrenzverhindernden balancierten Grammatiken nicht erlaubt, somit ist die Sprache keine konkurrenzverhindernde balancierte Sprache. Es folgt, dass die Familie der konkurrenzverhindernden balancierten Sprachen eine echte Untermenge der Familie der balancierten Sprachen ist.  $\square$

**Satz 4.3.15** Es sei  $\mathcal{L}$  die Familie der konkurrenzverhindernden balancierten Sprachen und  $\mathcal{L}_1$  die Familie der einfachen balancierten Sprachen, dann gilt  $\mathcal{L}_1 \subset \mathcal{L}$ .

*Beweis.* Es sei

$$L((x(a(c\bar{c})?\bar{a})(a(d\bar{d})\bar{a})^*\bar{x}))$$

die Sprache, die von folgender konkurrenzverhindernden Grammatik  $G = (N, T, S, P)$  erzeugt wird:

$$P = \{ \begin{array}{l} S \longrightarrow xAA_1^*\bar{x} \\ A \longrightarrow aC?\bar{a} \\ A_1 \longrightarrow aD\bar{a} \\ C \longrightarrow c\bar{c} \\ D \longrightarrow d\bar{d} \end{array} \}$$

Es ist wieder notwendig zwei konkurrierende Nichtterminale in einer Grammatik zu erzeugen.

Konkurrierende Nichtterminale innerhalb eines Inhaltsmodelles sind in den einfachen balancierten Grammatiken nicht erlaubt. Somit ist die Sprache keine einfache balancierte Sprache und es folgt, dass die Familie der einfachen balancierten Sprachen eine echte Untermenge der Familie der konkurrenzverhindernden balancierten Sprachen ist.  $\square$

**Satz 4.3.16** Es sei  $\mathcal{L}$  die Familie der einfachen balancierten Sprachen und  $\mathcal{L}_1$  die Familie der XML Sprachen, dann gilt  $\mathcal{L}_1 \subset \mathcal{L}$ .

*Beweis.* Es sei

$$L((x(a(c\bar{c})?\bar{a})(d(a(e\bar{e})\bar{a})\bar{d})^*\bar{x}))$$

die Sprache die von folgender Grammatik  $G = (N, T, S, P)$  erzeugt wird:

$$P = \{ \begin{array}{l} S \longrightarrow xAD^*\bar{x} \\ A \longrightarrow aC?\bar{a} \\ C \longrightarrow c\bar{c} \\ D \longrightarrow dA_1\bar{d} \\ A_1 \longrightarrow aE\bar{a} \\ E \longrightarrow e\bar{e} \end{array} \}$$

Es müssen wieder zwei konkurrierende Nichtterminale erzeugt werden. Konkurrierende Nichtterminale sind in den XML Grammatiken jedoch nicht erlaubt, daher kann die angegebene Sprache nicht durch eine XML Grammatik erzeugt werden.

Es folgt, dass die Familie der XML Sprachen eine echte Untermenge der Familie der einfachen balancierten Sprachen ist.  $\square$

**Satz 4.3.17** Es sei  $\mathcal{L}$  die Familie der balancierten Sprachen und  $\mathcal{L}_1$  die Familie der RegDyck Sprachen, dann gilt  $\mathcal{L}_1 \subset \mathcal{L}$ .

*Beweis.* Es sei

$$L((x(a^n(c\bar{c})?a^n)^*(a^m(d\bar{d})\bar{a}^m)*\bar{x})), m, n \in \mathbb{N}$$

die Sprache, die von folgender balancierter Grammatik  $G = (N, T, \{S\}, P)$  erzeugt wird:

$$P = \{ \begin{array}{l} S \longrightarrow xA^*A_1^*\bar{x} \\ A \longrightarrow a(A|C?)\bar{a} \\ A_1 \longrightarrow a(A_1|D)\bar{a} \\ C \longrightarrow c\bar{c} \\ D \longrightarrow d\bar{d} \end{array} \}$$

Die Definition der RegDyck Sprachen erlaubt allerdings nicht die Wörter  $\{(xa^n(c\bar{c})?a^n a^m(d\bar{d})\bar{a}^m\bar{x}) \mid m, n \in \mathbb{N}\}$  in einer Sprache, ohne dass dabei auch die Wörter

$$\{(xa^k a^n(c\bar{c})?a^n a^m(d\bar{d})\bar{a}^m\bar{a}^k\bar{x}) \mid m, n, k \in \mathbb{N}\}$$

in der Sprache enthalten sind.

Es folgt somit, dass es balancierte Sprachen gibt, die keine RegDyck Sprachen sind.  $\square$

**Korollar 4.3.2** Die in Satz 4.3.17 angegebene Sprache ist auch keine konkurrenzverhindernde balancierte Sprache, da die Sprache lediglich eine Abwandlung der Sprache aus Satz 4.3.14 ist. Es gibt somit balancierte Sprachen, die keine konkurrenzverhindernde balancierte Sprachen und keine RegDyck Sprachen sind.

**Satz 4.3.18** Es sei  $\mathcal{L}_1$  die Familie der konkurrenzverhindernden balancierten Sprachen und  $\mathcal{L}$  die Familie der RegDyck Sprachen, dann gilt  $\mathcal{L}_1 \not\subset \mathcal{L}$ .

*Beweis.* Es sei

$$L((x(a^n(c\bar{c})?a^n)(a^m(d\bar{d})\bar{a}^m)*\bar{x})), m, n \in \mathbb{N}$$

die Sprache, die von folgender balancierter Grammatik  $G = (N, T, S, P)$  erzeugt wird:

$$P = \{ \begin{array}{l} S \longrightarrow xAA_1^*\bar{x} \\ A \longrightarrow a(A|C?)\bar{a} \\ A_1 \longrightarrow a(A_1|D)\bar{a} \\ C \longrightarrow c\bar{c} \\ D \longrightarrow d\bar{d} \end{array} \}$$

Die Definition der RegDyck Sprachen erlauben jedoch nicht die Wörter  $\{(xa^n(c\bar{c})?\bar{a}^na^m(d\bar{d})\bar{a}^m\bar{x}) \mid m, n \in \mathbb{N}\}$  in einer Sprache, ohne dass dabei auch die Wörter

$$\{(xa^k a^n(c\bar{c})?\bar{a}^na^m(d\bar{d})\bar{a}^m\bar{a}^k\bar{x}) \mid m, n, k \in \mathbb{N}\}$$

in der Sprache enthalten sind.

Es folgt somit, dass es konkurrenzverhindernde balancierte Sprachen gibt, die keine RegDyck Sprachen sind.  $\square$

**Korollar 4.3.3** Die in Satz 4.3.18 angegebene Sprache ist auch keine einfache balancierte Sprache, da die Sprache lediglich eine Abwandlung der Sprache aus Satz 4.3.15 ist. Es gibt also balancierte Sprachen, die keine einfachen balancierten Sprachen und keine RegDyck Sprachen sind.

**Satz 4.3.19** Es sei  $\mathcal{L}_1$  die Familie der RegDyck Sprachen und  $\mathcal{L}$  die Familie der konkurrenzverhindernden Sprachen, dann gilt  $\mathcal{L}_1 \not\subseteq \mathcal{L}$ .

*Beweis.* Es sei  $G = (N, T, \{S\}, P)$ ,

$$P = \{ \begin{array}{l} S \longrightarrow xA^*A_1^*\bar{x} \\ A \longrightarrow a\bar{a} \\ A_1 \longrightarrow ad\bar{d}\bar{a} \end{array} \},$$

eine RegDyck Grammatik, die die Sprache

$$L((x(a\bar{a})^n(ad\bar{d}\bar{a})^m\bar{x})), m, n \in \mathbb{N}$$

erzeugt.

Es existiert keine konkurrenzverhindernde balancierte Grammatik, die diese Sprache erzeugen kann, da im Inhaltsmodell der zu konstruierenden Startproduktion zwei konkurrierende Nichtterminale erzeugt werden müssten und es für diese sowohl eine Sequenz  $UAW, U, W \in N^*$ , als auch eine Sequenz  $UA_1V, V \in N^*$  geben kann.  $\square$

**Korollar 4.3.4** Es existieren RegDyck Sprachen, die keine einfachen balancierten Sprachen und keine XML Sprachen sind. Also sind die RegDyck Sprachen keine Teilmenge dieser Sprachen.

Aus dem Beweis von Satz 4.3.19 wird klar, dass es erstens auch einfache balancierte Sprachen gibt, die keine RegDyck Sprachen und XML Sprachen sind und dass es XML Sprachen gibt, die keine RegDyck Sprachen sind.

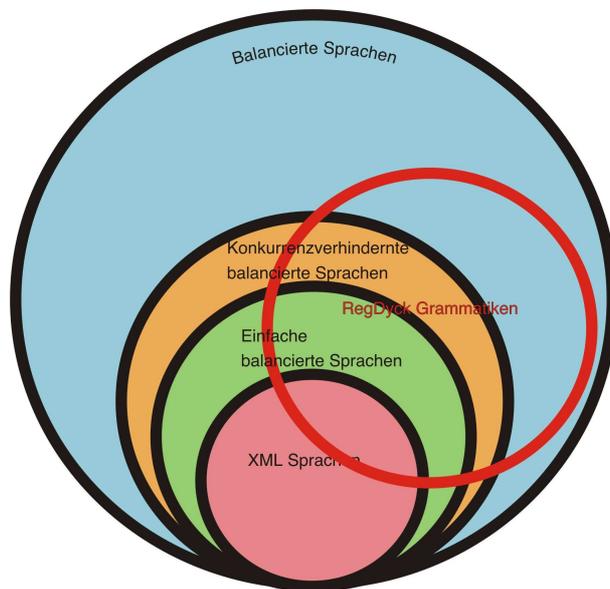


Abbildung 4.3: Hierarchie der XML-artigen Grammatiken.

Zweitens gibt es auch RegDyckSprachen, die keine einfachen balancierten Sprachen und keine XML Sprachen sind.

## 4.4 Parser Generierung

Dieses Unterkapitel beschäftigt sich mit den Parsergeneratoren und den Parsern von XML-artigen Grammatiken. In der Praxis werden XML Dokumente oft zusammen mit ihren Beschreibungsmodellen verschickt. Zur Validierung muss zuerst das Beschreibungsmodell gelesen und daraus ein geeigneter Parser konstruiert werden.

Es stellt sich die Frage, wie lange es dauert, um einen Parser zu erzeugen und wie lange der Parser braucht, um ein XML Dokument zu validieren.

### 4.4.1 Sprachen und Methoden

Kapitel 4.2 hat gezeigt, dass die XML Grammatiken, die einfachen balancierten Grammatiken und die konkurrenzverhindernden balancierten Grammatiken zu den  $LL(1)$  Grammatiken gehören.

Für diese Art von Grammatiken gibt es Untersuchungen [SSS88] für die Laufzeiten der Parsergeneratoren und der Parser.

Sprache	max. Laufz. Parsergenerator	max. Laufz. Parser
XML Grammatiken	$O(n)$	$O(n)$
einfache balancierte Grammatiken	$O(n)$	$O(n)$
konk. bal. Grammatiken	$O(n)$	$O(n)$
RegDyck Grammatiken	?	?
balancierte Grammatiken	?	?

Für die balancierten Grammatiken und die RegDyck Grammatiken müssen die Laufzeiten noch untersucht werden. Zusätzlich wird noch eine weitere Variante für Parser von XML Grammatiken, einfachen balancierten Grammatiken und konkurrenzverhindernden balancierten Grammatiken vorgestellt und untersucht. Teile der neuen Variante können in der Praxis für weiterführende Anwendungen benutzt werden, wodurch sich ein Zeitersparnis ergibt.

### 4.4.2 Konstruktion von codeterministischen Grammatiken aus balancierten Grammatiken

In [BB02a] wird gezeigt, dass aus einer beliebigen balancierten Grammatik eine codeterministische balancierte Grammatik erzeugt werden kann.

Mit Hilfe der codeterministischen balancierten Grammatik soll ein spezieller Parsingalgorithmus verwendet werden.

Eine balancierte Grammatik  $(N, T, S, P)$  wird als codeterministisch bezeichnet, falls gilt, wenn  $X \rightarrow m$  und  $X' \rightarrow m$ , dann ist  $X = X'$ ,  $X, X' \in N, m \in Reg(N \cup T)$ .

Im Folgenden wird die Konstruktion in Bezug auf ihre Komplexität analysiert. Die Ergebnisse sollen später zum Vergleich der Grammatikmodelle

genutzt werden.

Die Größe eines endlichen Automaten wird durch die Anzahl der Zustände multipliziert mit der Anzahl der Eingabesymbole  $|S| \cdot |A|$  festgelegt. Die Größe einer Grammatik ist durch die Länge aller Produktionen gegeben  $\sum_{p \in P} |p|$ .

Um eine balancierte Grammatik in eine codeterministische balancierte Grammatik zu transformieren sind die folgenden Schritte notwendig.

- Es sei  $G = (N, T, S_0, P)$  eine balancierte Grammatik. Ausgangspunkt sind die rechten Seiten der Produktionen. Für ein Nichtterminalzeichen  $X$  werden alle rechten Seiten der Produktionen für  $X$  zu einem regulären Ausdruck  $R_X$  zusammengefasst. Es sei  $m$  die Größe der Grammatik. Es sei  $n$  die Länge des regulären Ausdrucks der Produktion für  $X$ .  $\mathcal{A}_X = (S, A, s_0, \delta, F)$  sei ein deterministischer endlicher Automat, der  $R_X$  erkennt. Die Umwandlung eines regulären Ausdrucks in einen nichtdeterministischen endlichen Automaten nach der Thompson Methode verlangt maximal vier Zustände für das Konstruieren eines erkennenden Automaten pro Zeichen. Bei der Länge  $n$  sind das maximal  $4 \cdot n$  Zustände. Dadurch, dass die Potenzmenge der Zustände bei der Umwandlung in einen deterministischen endlichen Automaten konstruiert wird, kann die Umwandlung hierbei im schlimmsten Fall  $O(2^{4n})$  Zustände zur Folge haben. Die für das Konstruieren notwendige Zeit ist auch  $O(2^{4n})$ .
- Für jeden Automat  $\mathcal{A}_X$  wird ein Potenzautomat  $\mathcal{A}'_X = (S', A', s'_0, \delta', F')$  gebildet, welcher folgendermaßen definiert ist:
  1.  $S' = 2^{|S|}$  ist die Menge der Zustände.
  2.  $A' = 2^{|A|} \setminus \emptyset$  ist die Menge der Eingabesymbole.
  3.  $s'_0 = \{s_0\}$  ist der Startzustand.
  4.  $\delta'$ , für  $s' \in S'$  und  $a' \in A'$

$$\delta(s', a') = \bigcup_{s \in s', a \in a'} \delta(s, a)$$

5.  $F' = \{s' \in S' \mid s' \cap F \neq \emptyset\}$

Für den Potenzautomat ergibt sich die Größe von

$$2^{2^{4n}} \cdot (2^{|T|} - 1) \in O(2^{2^{4n}}).$$

Nun wird für alle Teilmengen von Nichtterminalzeichen  $\alpha \in 2^N, \alpha = \{X_1, \dots, X_i\}$  folgende Konstruktion durchgeführt:

1.  $D_\alpha = L(\mathcal{A}'_{X_1}) \cap L(\mathcal{A}'_{X_2}) \cap \dots \cap L(\mathcal{A}'_{X_i})$

2. Für alle  $\beta \in 2^{(N \setminus \alpha)} - \{\emptyset\}, \beta = \{Y_1, \dots, Y_j\}$  wird die Sprache  $P_{\alpha, \beta} = D_\alpha \cap L(\mathcal{A}'_{y_1}) \cap L(\mathcal{A}'_{y_2}) \cap \dots \cap L(\mathcal{A}'_{y_i})$  gebildet.
3.  $G_\alpha = \bigcup_\beta P_{\alpha, \beta}$
4.  $S_\alpha = D_\alpha - G_\alpha$

Die Menge  $D_\alpha$  beschreibt für eine Menge von Nichtterminalzeichen den Durchschnitt der erzeugten Wörter aus den Produktionen.

Die Produktionen der entstehenden codeterministischen balancierten Grammatik sind:

$$\alpha \longrightarrow S_\alpha, \text{ für alle } \alpha \in 2^N.$$

Nun sei angenommen die Produktionen der balancierten Grammatik  $G$  haben alle die Länge  $n$ . Dann existieren  $m/n$  Nichtterminale. Die Automaten  $\mathcal{A}'_{x_i}, i \in \mathbb{N}$  haben somit eine maximale Größe von  $O(2^{2^{4n}})$ . Es seien  $S_{\mathcal{X}}$  und  $S_{\mathcal{Y}}$  die Zustandsmengen der Automaten  $\mathcal{X}, \mathcal{Y}$ . Die Konstruktion des Durchschnitts dieser Automaten benötigt  $O(S_{\mathcal{X}} \cdot S_{\mathcal{Y}})$  Zeit. Da alle Teilmengen konstruiert werden müssen und für diese jeweils der Durchschnitt berechnet werden muss, ergibt sich folgende Zeit:

$$\sum_{k=1}^{\lceil \frac{m}{n} \rceil} 2^{2^{4n \cdot (k-1) \cdot \lceil \frac{m}{k} \rceil}} \in \Omega\left(\frac{m}{n} \cdot 2^{2^{4n \cdot \lceil \frac{m}{n} \rceil}}\right), \lceil \frac{m}{n} \rceil \geq 3$$

Die Laufzeit des Parsergenerators ist sehr groß. Wenn eine balancierte Grammatik aus einer einzigen Produktion besteht, erreicht der Parser eine Laufzeit von  $O(2^{2^m})$ .

Dazu kommt noch die Umwandlung der deterministischen endlichen Automaten in einen regulären Ausdruck. Es sei  $k$  die Anzahl der Zustände eines endlichen Automaten, dann ist die Laufzeit zur Konstruktion eines regulären Ausdrucks  $O(k^3)$  [HU79]. Dies muss jedoch nicht mehr in die Laufzeitangabe einfließen, da es durch die  $O$ -Notation verschwinden würde.

Es wird noch alleine die Größe der codeterministischen balancierten Grammatik betrachtet. Es werden alle Zwischenberechnungen (die Konstruktion eines deterministischen Automaten aus dem regulären Ausdruck der Nichtterminalzeichen und die Potenzautomaten) außer Acht gelassen und nur die Menge der entstehenden Nichtterminalzeichen der codeterministischen balancierten Grammatik mit der Menge der Nichtterminalzeichen in der ursprünglichen Grammatik verglichen. Die codeterministische balancierte Grammatik könnte maximal  $2^{|N|}$  Nichtterminalzeichen besitzen, da die Potenzmenge der Nichtterminalzeichen gebildet wird.

Für die Praxis ergibt die Umwandlung einer balancierten Grammatik in eine codeterministische Grammatik keinen Sinn. Jedoch können alle Sprachen, die von einer balancierten Grammatik erzeugt werden, auch durch eine codeterministische balancierte Grammatik erzeugt werden. So können

Beschreibungsmodelle mit codeterministischen balancierten Grammatiken verlangt werden, ohne dabei die Beschreibungsmächtigkeit einzuschränken.

#### 4.4.3 Konstruktion eines deterministischen Kellerautomaten aus einer balancierten Grammatik

Die Konstruktion eines deterministischen Kellerautomaten  $\mathcal{D} = (S, A, G, \Delta, s_0, \delta, F)$  aus einer balancierten Grammatik  $G = (N, T, S, P)$  wird in zwei Phasen unterteilt. Die erste Phase besteht aus der Konstruktion eines nicht-deterministischen Kellerautomaten  $\mathcal{D}' = (S', A', G', \Delta', s'_0, \delta', F')$  aus der balancierten Grammatik  $G$ . Die zweite Phase besteht aus der Konstruktion des deterministischen Kellerautomaten  $\mathcal{D}$  aus dem nichtdeterministischen Kellerautomaten  $\mathcal{D}'$ .

Die Größe der balancierten Grammatik ist wie im letzten Unterkapitel durch die Länge aller Produktionen gegeben  $\sum_{p \in P} |p|$ . Die Größe von  $G$  sei  $n$ . Zur Konstruktion des nichtdeterministischen Kellerautomaten  $\mathcal{D}'$  wurde die abgewandelte Thompson Methode benutzt. Wie schon im letzten Unterkapitel beschrieben, erzeugt diese Methode pro Zeichen höchstens 4 Zustände. Bei der abgewandelten Methode werden lediglich die Kelleroperationen hinzugefügt. Dies erzeugt keinen zusätzlichen Aufwand. Es ergibt sich eine Größe von:

$$|S'| = 4 \cdot n$$

Die Menge der Eingabezeichen ermittelt sich durch:

$$|A'| = \frac{|T|}{2}$$

Die maximale Anzahl der Kellersymbole ist:

$$|G'| = \frac{n}{2}$$

In der Konstruktion vom deterministischen Kellerautomaten  $\mathcal{D}$  aus dem nichtdeterministischen Kellerautomaten  $\mathcal{D}'$  wird die Potenzmenge der Zustände und die Potenzmenge der Kellersymbole gebildet. Für den deterministischen Kellerautomaten ergibt sich dadurch eine Größe von:

$$S \cdot A \cdot G = 2^{4n} \cdot \frac{|T|}{2} \cdot 2^{\frac{n}{2}}$$

In der  $O$ -Notation ist die Laufzeit gleich der Größe des Automaten, da nur die Kellerautomaten erzeugt werden müssen und keine zusätzlichen Konstruktionen durchgeführt werden. Es ergibt sich eine Laufzeit von:

$$O\left(2^{5n} \cdot \frac{|T|}{2}\right) = O(2^{5n})$$

#### 4.4.4 Konstruktion eines Parsebaums für XML Grammatiken

In [ML02] wurden spezielle Bäume für das Transformieren von XML Dokumenten in relationale Modelle verwendet. Diese Art von Bäumen wird hier erweitert und die Struktur genau definiert. Es wird beschrieben, wie die Bäume für das Parsen verwendet werden können. Dadurch kann Zeit eingespart werden, weil die Konstruktion dieser Bäume, sowohl für die Parsergenerierung, als auch für die Transformation der XML Dokumente genutzt werden kann.

Der folgende Parser kann für XML Grammatiken, einfache balancierte Grammatiken und konkurrenzverhindernde balancierte Grammatiken verwendet werden.

Beim Überprüfen der Gültigkeit eines Dokumentes ergibt sich das Problem, dass die Beschreibungsmodelle von Dokument zu Dokument variieren. Es kann kein allgemeiner Algorithmus verwendet werden, der eine feste Syntax überprüft, wie beispielsweise bei der Syntaxüberprüfung einer Programmiersprache.

Es wird ein Algorithmus benötigt, der ein Beschreibungsmodell, zum Beispiel eine DTD, einliest und ein Programm erzeugt, welches überprüft, ob das Dokument gültig ist.

Folgende DTD soll als Beispiel dienen:

**Beispiel 4.4.20** Die Beispiel DTD *Person* ist folgendermaßen definiert:

```
<!ELEMENT Person (Vorname*, Nachname, Adresse)>
<!ELEMENT Vorname (#PCDATA)>
<!ELEMENT Nachname (#PCDATA)>
<!ELEMENT Adresse (Ort, Strasse?)>
<!ELEMENT Ort(#PCDATA)>
<!ELEMENT Strasse (#PCDATA)>
```

□

Eine DTD ist von der Beschreibungsmächtigkeit durch eine XML Grammatik ausdrückbar und allgemein auch ähnlich aufgebaut wie eine kontextfreie Grammatik.

**Beispiel 4.4.21** Die Überführung der DTD aus dem letzten Beispiel ergibt eine XML Grammatik mit folgenden Produktionen:

```

PERSON  → <Person> VORN* NACHN ADR</Person>
VORN    → <Vorname></Vorname>
NACHN   → <Nachname></Nachname>
ADR     → <Adresse> ORT STRASSE? </Adresse>
ORT     → <Ort></Ort>
STRASSE → <Strasse></Strasse>

```

□

Es wird ein Parsebaum aufgebaut und das Dokument wird in Bezug auf den Baum überprüft. Der Baum muss variabel für Wiederholungen sein und die Möglichkeit haben optionale Produktionen zu berücksichtigen.

Der Validierungsalgorithmus muss beim Wurzelement beginnen und dann eine beliebige Anzahl von Unterelementen anhängen. Wiederholungen und Optionen werden durch eine Kennzeichnung des Knotens behandelt.

Folgende Schritte sind durchzuführen:

- Die Nichtterminalzeichen innerhalb der Produktionen werden indiziert.
- Für jedes Nichtterminalzeichen wird ein Knoten erzeugt. In jedem Knoten muss der Name des Elementes eingetragen werden.
- Für die Wiederholungszahl werden zwei Felder benötigt (Feld *Optional* und Feld *Wiederholung* jeweils mit *True* und *False* belegbar). Wenn das Nichtterminalzeichen im Eltern-Element des Elementes mit einem Fragezeichen „?“ versehen ist, wird das Feld *Optional* mit *True* und das Feld *Wiederholung* mit *False* belegt. Für ein Pluszeichen „+“ wird das Feld *Wiederholung* mit *True* und das Feld *Optional* mit *False* gekennzeichnet. Bei einem Stern „\*“ wird beides mit *True* belegt.
- Ein Feld *Wert* muss mit Werten belegbar sein (hier wird der Name des Elementes gespeichert).
- Zuletzt muss die Möglichkeit bestehen, auf verschiedene Elemente zu verweisen. Dazu wird eine Liste *next* aus Zeigern erzeugt.

**Implementierung des Baums** Ein Algorithmus liest zuerst das Beschreibungssystem ein und baut danach den Baum auf. Der Baum, der für die Grammatik aufgebaut werden muss, wird mit Hilfe von Zeigern implementiert.

Die Elemente bestehen aus einem Feld *Wert* vom Typ *String*, einem Feld *Optional* vom Typ *Boolean*, einem Feld *Wiederholung* vom Typ *Boolean* und einem Feld mit den Verweisen. Die Verweise bestehen aus einer Liste von Zeigern. Diese Liste verweist auf alle Nachfolger des Knotens.

#### Beispiel 4.4.22

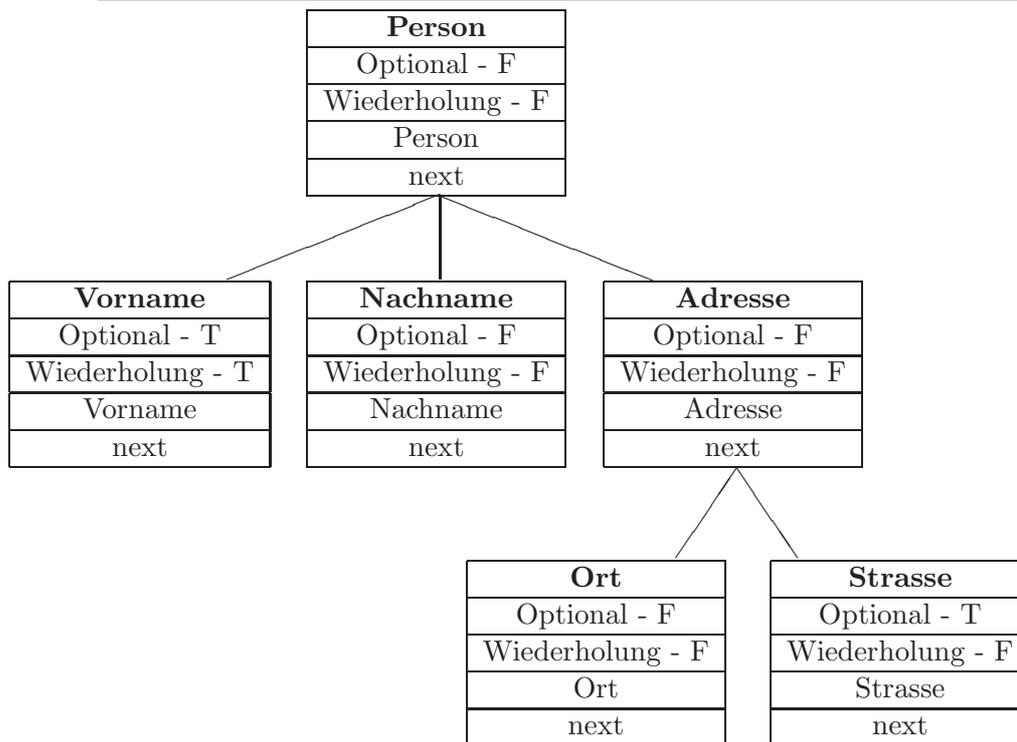


Abbildung 4.4: Beispiel: Parsebaum zur Validierung eines XML Dokuments.

Typ der Knoten

TYPE

node = ^nodeelement;

list = ^liste;

liste = Record

  elementlist : node;

  next : list;

end;

nodeelement = Record

  Wert : String;

  Optional : bool;

  Wiederholung : bool;

  nodes : list;

end;

Die Inhaltsmodelle der Grammatiken bestehen aus regulären Ausdrücken von Nichtterminalzeichen. Es sei hier nochmal darauf hingewiesen, dass sowohl die XML Grammatiken, als auch die einfachen balancierten Grammatiken und die konkurrenzverhindernden balancierten Grammatiken aus 1-lookahead deterministischen regulären Ausdrücken bestehen. Die Nichtterminalzeichen sind die Wurzeln der Bäume.

Um eine genaue Abfolge der Nachfolgebäume abzubilden, wird eine Liste erzeugt, die nur eine Abarbeitung der Nachfolgebäume in der Reihenfolge des regulären Ausdrucks zulässt.

Die Operatoren des regulären Ausdrucks werden in den Knoten untergebracht. Die Elementdeklaration ermöglicht es, verschiedene Varianten von Unterelementen in der Elementdeklaration anzugeben und diese zu gruppieren:

**Beispiel 4.4.23** `<!ELEMENT((Name | Persnr Lieblingsessen))>`

Dieses Element kann wahlweise einen Namen oder eine Personalnummer als Unterelement besitzen. Für jede Gruppe wird ein Knoten erzeugt, dessen *Wert* leer ist. Als Nachfolgebäume werden die Elemente der Gruppe angehängt. Bei Optionen wird den Listenelementen eine Signatur „optional“ gegeben. In der Gültigkeitsanalyse muss entschieden werden, welcher Unterbaum der Prüfung unterzogen werden soll. Auch bei Wiederholungen ganzer Gruppen ist im Knoten der Gruppe das Feld *Wiederholung* mit *true* belegt.

Nun wird für die Generierung des Baums die Komplexität bestimmt. Dabei soll, in Abhängigkeit von der Größe der Produktionen, eine Angabe gemacht werden, wie oft die einzelnen Zeichen durchlaufen werden.

Jedes Nichtterminalzeichen in den Produktionen wird indiziert. Anschließend wird für jedes Nichtterminalzeichen sowie für jede Gruppe ein Knoten erzeugt. Daraus ergibt sich eine Laufzeit von:

$$O(n).$$

**Probleme der Überprüfung** Bei dieser Methode können zwei verschiedene Arten von Problemen auftreten:

- Mehrdeutigkeiten und 1-lookahead Determinismus
- Ringschlüsse innerhalb einer DTD

**Mehrdeutigkeiten und 1-lookahead Determinismus** Die Beschreibungsmodelle ermöglichen es im allgemeinen Fall, dass mehrdeutige Grammatiken definiert werden können. Die Mehrdeutigkeiten können bei Alternativen von Unterelementen auftreten.

**Beispiel 4.4.24** Ein Element mit einem mehrdeutigen Inhaltsmodell:

`ELEMENT((a* b) | b)`

□

Das Problem besteht darin, den Baum zu finden, der aufgebaut werden soll. Wenn das Symbol  $b$  eingelesen wird, ist nicht klar, ob das  $a$  weggelassen worden ist (was durch die Kennzeichnung des Sterns „\*“ richtig sein kann) oder ob der Baum mit dem Unterelement  $b$  aufgebaut werden soll.

Die Mehrdeutigkeiten müssen in der Definition ausgeschlossen werden, damit es zu keinen Fehlern kommt. Dies ist bei XML Grammatiken, einfachen balancierten Grammatiken und konkurrenzverhindernden balancierten Grammatiken der Fall.

Desweiteren muss die 1-lookahead Determinismus Eigenschaft eingehalten werden. Es existiert zum Beispiel ein Beschreibungsmodell

$$\text{ELEMENT}((a^* b) \mid (b c)),$$

welches zwar nicht mehrdeutig ist, in dem es jedoch trotzdem nicht für jedes Eingabezeichen eindeutig ist, welcher Teil des regulären Ausdrucks dieses Zeichen erzeugt.

Auch dies wird von den oben genannten Grammatiken eingeschränkt.

**Ringschlüsse innerhalb einer DTD** Bei Ringschlüssen (Selbsteinbettung) innerhalb eines Baums verweisen Knoten auf Vorgängerknoten. Dieses Problem wurde bisher noch nicht behandelt. Eine Möglichkeit besteht darin, sich alle Vorgängerverweise zu merken und im Anschluss zu bearbeiten. Dies kann mit einem Stapel implementiert werden, auf welchen immer die letzte Zeigerposition gelegt wird. Nach der Beendigung eines Baums wird der Stapel überprüft und gegebenenfalls an die Stelle des Zeigers zurückgesprungen.

#### 4.4.5 Konstruktion eines regulären Ausdrucks und einer Dyck Grammatik aus einer RegDyck Grammatik

Es sei  $G = (N, T, S, P)$ ,  $T = A \cup \bar{A}$  eine balancierte Grammatik. Die Grammatik  $G$  wird indiziert und zu  $G' = (N', T', S', P')$ .

Die Konstruktion eines erkennenden RegDyck Automaten  $\mathcal{K}$  enthält die Konstruktion von Automaten  $\mathcal{K}_X$ ,  $X \in N'$ , die einzelne Teile der Grammatik erkennen. Diese Automaten werden miteinander verbunden.

Die Automaten für die Inhaltsmodelle werden durch eine Methode, die der Thompson Methode für reguläre Ausdrücke ähnlich ist, zusammengesetzt. Die Automaten für selbsteinbettende Sprachteile werden über zwei Schleifen konstruiert.

Der entstehende Automat ist ein RegDyck Automat. Aus diesem wird ein endlicher Automat konstruiert, indem nur die Zustandsüberführungen übernommen werden.

Nun wird noch der Dyck Automat erzeugt:

$\mathcal{D} = (S_D, A_D, G_D, \Delta, s_{0D}, \delta_D, F_D)$  ist ein deterministischer Kellerautomat, der die Menge  $D_{A \cup \bar{A}}$  von Dyckprimen erkennt:

- $S_D = \{s_{0D}, s_{1D}, s_{eD}\}$  ist die Menge der Zustände des Automaten.
- $A_D = A \cup \bar{A}$  ist das Eingabealphabet.
- $G_D = \bar{A} \cup \{\Delta\}$  ist das Kelleralphabet.
- $s_{0D}$  ist der Startzustand.
- $F_D = \{s_{0D}, s_{eD}\}$  ist die Menge der Endzustände.
- $\delta_D$  ist die Überföhrungsfunktion:

$$\begin{aligned} (s_{0D}, a, \Delta) &\longrightarrow (s_{1D}, \bar{a}\Delta) \quad \forall a \in A \\ (s_{1D}, a, \bar{z}) &\longrightarrow (s_{1D}, \bar{a}z) \quad \forall a \in A \quad \forall z \in \bar{A} \\ (s_{1D}, \lambda, \Delta) &\longrightarrow (s_{eD}, \Delta) \\ (s_{1D}, \bar{a}, \bar{a}) &\longrightarrow (s_{1D}, \lambda) \quad \forall \bar{a} \in \bar{A}. \end{aligned}$$

Die Konstruktion eines endlichen Automaten aus einer Grammatik mit Hilfe der Erweiterung der Thompson Methode, ist in

$$O(n)$$

Zeit möglich. Im Vergleich zu der originalen Thompson Methode werden hierbei jedoch nur die Kelleroperationen zusätzlich berücksichtigt. Diese stehen aber, wie schon in Kapitel 3 erwähnt, anhand der Eingabesymbole fest und die Konstruktion benötigt keine zusätzliche Zeit. Wichtig hierbei ist die Tatsache, dass nichtdeterministische endliche Automaten konstruiert werden können.

Die Konstruktion des Dyck Automaten benötigt in Bezug auf die Größe der Grammatik konstante Zeit  $O(1)$ , da die Größe des Dyck Automaten von der Größe des Alphabets abhängig ist. Daraus ergibt sich eine Laufzeit von

$$O(n) + O(1) = O(n).$$

Diese Methode erscheint vor diesem Hintergrund geeignet für den Parsingvorgang.

## 4.5 Parsing Algorithmen

Im letzten Unterkapitel wurde die Konstruktion von Parsergeneratoren beschrieben, zusätzlich wurden sie auf ihre Komplexität analysiert. In diesem Kapitel werden zwei Parsing Algorithmen beschrieben. Der erste Parsing Algorithmus überprüft codeterministische balancierte Grammatiken, der zweite nutzt die Parsebäume, die im letzten Unterkapitel vorgestellt worden sind, zur Syntaxüberprüfung. Die anderen Parser, die im letzten Kapitel konstruiert worden sind, müssen nicht mehr analysiert werden, da ihre Laufzeiten bereits bekannt sind.

### 4.5.1 Parsen mittels codeterministischer Grammatiken

Es sei  $G = (N, T, S, P)$  eine codeterministische balancierte Grammatik.

Es wird für jede Produktion  $p \in P$  ein nichtdeterministischer endlicher Automat  $\mathcal{R}_p$  konstruiert, der die Sprache der rechten Seiten der Produktion erkennt. Aus allen Automaten  $\mathcal{R}_p$  wird ein nichtdeterministischer endlicher Automat  $\mathcal{A}$  konstruiert, welcher alle Produktionen erkennt.

Aus den Endzuständen ist dabei ersichtlich, welche Produktion der Automat erkannt hat. Die codeterministischen balancierten Grammatiken haben die Eigenschaften, dass jede rechte Seite einer Produktion eindeutig ist.

Aufgrund der speziellen Syntax der Sprache (öffnende und schließende Klammern) können sofort die verschiedenen Ebenen der Produktionen erkannt werden. Dies wird im folgenden Algorithmus ausgenutzt.

Der Algorithmus benutzt eine Tabelle. Die Tabelle besteht aus zwei verschiedenen Spalten (*Ebene* und *Zeichenkette*).

Die eingelesenen Zeichen beim Parsingvorgang werden auf verschiedenen Ebenen aufgeschrieben.

Die Ebenen der weiter unten liegenden Zeichen sind die Nachfolgebene und kennzeichnen die Zeichen des Inhaltsmodells.

Für die Überprüfung der Syntax wird folgender Algorithmus verwendet:

- Lese die Zeichen der Eingabe nacheinander von vorne nach hinten ein.
- Bei einem Zeichen  $a \in A$  schreibe dieses Zeichen auf die nächste Ebene.
- Bei einem Zeichen  $\bar{a} \in \bar{A}$  schreibe das Zeichen auf die aktuelle Ebene und überprüfe die Zeichenkette mit dem Automaten  $\mathcal{A}$ . Das zurückgegebene Nichtterminal schreibe in die Zeichenkette auf die um eine höhere Ebene. Dann lösche den Inhalt der niedrigsten Ebene.

Wenn die Zeichenkette abgearbeitet ist und das Startsymbol erreicht wurde, war die Überprüfung erfolgreich.

Bei einem der folgenden drei Fälle wird abgebrochen und die Analyse war nicht erfolgreich:

1.  $\mathcal{A}$  liefert kein Ergebnis.
2. Die Zeichenkette ist abgearbeitet und das Startsymbol ist nicht erreicht.
3. Es wird ein Zeichen  $\bar{a} \in \bar{A}$  eingelesen und es gibt kein passendes Zeichen  $a \in A$  auf dieser Ebene.

Jedes Zeichen wird nur einmal eingelesen und verarbeitet. Das bedeutet eine Komplexität in Bezug auf die Eingabe von  $O(n)$ . Jedoch werden in den Zwischenschritten Zeichenketten erzeugt, die durch den nichtdeterministischen endlichen Automaten  $\mathcal{A}$  erkannt werden müssen. Es sei  $n$  die Länge

der Eingabe, dann können  $n/2$  Nichtterminale entstehen, die der Automat zusätzlich zu der Eingabe verarbeiten muss. Damit ergibt sich eine Komplexität von

$$O(n) + O(n/2) = O(n).$$

Dazu kommt die Zeit, die der Automat  $\mathcal{A}$  benötigt, um die Zeichenketten zu verarbeiten. Insgesamt ergibt sich dadurch:

$$2 \cdot (O(n) + O(n/2)) = O(n).$$

**Beispiel 4.5.25** Es sei  $G = (N, T, \{S\}, P)$  eine codeterministische balancierte Grammatik:

- $N = \{S, B, C, D, E\}$
- $T = A \cup \bar{A}, A = \{a, b, c, d, e\}, \bar{A} = \{\bar{a}, \bar{b}, \bar{c}, \bar{d}, \bar{e}\}$
- $P = \{$

$$\begin{array}{lcl} S & \longrightarrow & aDC^*B^+\bar{a} \\ B & \longrightarrow & b\bar{b} \\ C & \longrightarrow & cE^+\bar{c} \\ D & \longrightarrow & d\bar{d} \\ E & \longrightarrow & e\bar{e} \end{array}$$

Der nichtdeterministische endliche Automat für die Grammatik sieht folgendermaßen aus:

- $S' = \{s_0, s_a, s_D, s_C, s_B, s_{\bar{a}}, s_b, s_{\bar{b}}, s_c, s_E, s_{\bar{c}}, s_d, s_{\bar{d}}, s_e, s_{\bar{e}}\}$
- $A' = \{a, b, c, d, e, \bar{a}, \bar{b}, \bar{c}, \bar{d}, \bar{e}, S, B, C, D, E\}$
- $B' = \{S, B, C, D, E\}$
- $s_0$
- $F = \{s_{\bar{a}}, s_{\bar{b}}, s_{\bar{c}}, s_{\bar{d}}, s_{\bar{e}}\}$
- Die Überführungen des Automaten sind folgende

$$\begin{array}{ll}
 (s_0, a) & \longrightarrow (s_a, \lambda) \\
 (s_0, b) & \longrightarrow (s_b, \lambda) \\
 (s_0, c) & \longrightarrow (s_c, \lambda) \\
 (s_0, d) & \longrightarrow (s_d, \lambda) \\
 (s_0, e) & \longrightarrow (s_e, \lambda) \\
 (s_a, D) & \longrightarrow (s_D, \lambda) \\
 (s_D, C) & \longrightarrow (s_C, \lambda) \\
 (s_D, B) & \longrightarrow (s_B, \lambda) \\
 (s_C, C) & \longrightarrow (s_C, \lambda) \\
 (s_C, B) & \longrightarrow (s_B, \lambda) \\
 (s_B, B) & \longrightarrow (s_B, \lambda) \\
 (s_B, \bar{a}) & \longrightarrow (s_{\bar{a}}, S) \\
 (s_b, \bar{b}) & \longrightarrow (s_{\bar{b}}, B) \\
 (s_c, E) & \longrightarrow (s_E, \lambda) \\
 (s_E, E) & \longrightarrow (s_E, \lambda) \\
 (s_E, \bar{c}) & \longrightarrow (s_{\bar{c}}, C) \\
 (s_d, \bar{d}) & \longrightarrow (s_{\bar{d}}, D) \\
 (s_e, \bar{e}) & \longrightarrow (s_{\bar{e}}, E)
 \end{array}$$

Es sei die Zeichenkette  $\underline{a}\bar{d}\bar{d}\bar{c}\bar{e}\bar{e}\bar{c}\bar{c}\bar{e}\bar{e}\bar{e}\bar{e}\bar{e}\bar{c}\bar{b}\bar{b}\bar{a}$  gegeben.

**Ebene Zeichenkette**

0.

Im ersten Schritt wird das erste Zeichen  $a \in A$  eingelesen:

$$\underline{a}\bar{d}\bar{d}\bar{c}\bar{e}\bar{e}\bar{c}\bar{c}\bar{e}\bar{e}\bar{e}\bar{e}\bar{c}\bar{b}\bar{b}\bar{a}$$

**Ebene Zeichenkette**

0. a

nach dem ersten Schritt ergibt sich folgende Konfiguration:

$$\underline{a}\bar{d}\bar{d}\bar{c}\bar{e}\bar{e}\bar{c}\bar{c}\bar{e}\bar{e}\bar{e}\bar{e}\bar{c}\bar{b}\bar{b}\bar{a}$$

nach dem 5. Schritt ergibt sich folgende Konfiguration:

**Ebene Zeichenkette**

0. a D  
 1. c  
 2. e

$$\underline{a}\bar{d}\bar{d}\bar{c}\bar{e}\bar{e}\bar{c}\bar{c}\bar{e}\bar{e}\bar{e}\bar{e}\bar{c}\bar{b}\bar{b}\bar{a}$$

Nachdem das Zeichen  $\bar{e} \in \bar{A}$  eingelesen wurde, wird die Zeichenkette auf der zweiten Stufe von  $\mathcal{A}$  ausgewertet und es wird das Nichtterminal  $E$  zurückgegeben. Das Nichtterminal wird auf der Ebene eins eingetragen. Die Zeichen auf der Ebene zwei werden gelöscht.

Ebene	Zeichenkette
0.	a D
1.	c E

Nach dem 17. Schritt ergibt sich folgende Konfiguration:

Ebene	Zeichenkette
0.	a

$$add\bar{c}e\bar{e}c\bar{c}e\bar{e}e\bar{e}e\bar{c}bb\bar{a}$$

Das Zeichen  $\bar{a} \in \bar{A}$  wird eingelesen. Das Nichtterminal  $S$  wird ausgegeben. Es ist das Startsymbol.

Die Zeichenkette ist fertig bearbeitet.

Die Analyse war erfolgreich.

#### 4.5.2 Parsen mittels Parsebäumen

Im Folgenden wird der Parsingalgorithmus für die Parsebäume der XML Dokumente beschrieben. Das XML Dokument wird schrittweise durchlaufen. Dabei ist jedes Element ein Zeichen der dazugehörigen kontextfreien Grammatik. Es wird vor dem ersten Knoten des Parsebaums begonnen.

**Algorithmus** Der Baum wird in Preorder durchlaufen.

1. Das erste Zeichen  $x$  ( $x \in A$ , Starttag) wird eingelesen. Falls  $x \in \bar{A}$ , dann ist das Dokument fehlerhaft.
2. Konnte kein nächstes Zeichen eingelesen werden, dann wird überprüft, ob im Baum ebenfalls das Ende erreicht ist und alle Knoten behandelt sind oder ob alle weiteren Knoten mit *Optional* mit *True* belegt sind. Stimmt dies, ist die Gültigkeitsanalyse erfolgreich, andernfalls nicht.
3. Es wird zum nächsten Knoten des Baums gesprungen. (Wird in eine nichtadressierte Speicherstelle verwiesen, wird ein Fehler gemeldet.)
4. Das Zeichen wird mit dem Wert des Knotens des Baums verglichen. Ist der Vergleich erfolgreich, dann wird mit dem nächsten Punkt des Algorithmus fortgefahren. Falls der Vergleich nicht korrekt ist, wird überprüft, ob *Option* auf *True* gesetzt ist. Wenn dies der Fall ist, wird zum nächsten Knoten in der Liste gesprungen. Ist dies nicht der Fall, wird ein Fehler ausgelöst.
5. Das nächste Zeichen  $x$  wird eingelesen. Falls  $x = \bar{z} \in \bar{A}$  ist, wird überprüft, ob der aktuelle Knoten mit dem Wert  $z \in A$  belegt ist. Ist dies der Fall, dann ist dieser Knoten abgearbeitet und es wird zu dem Vorgängerknoten gesprungen und mit Punkt 5 fortgefahren. Sonst wird ein Fehler ausgelöst. Falls  $x \in A$  ist, wird mit Punkt 6 fortgefahren.

6. Falls *Wiederholen* im Knoten auf *True* steht, wird wieder mit dem aktuellen Knoten verglichen. Ist der Vergleich korrekt, wird mit Punkt 5 fortgefahren. Wenn der Vergleich nicht korrekt ist oder *Wiederholen* auf *False* steht, wird mit Punkt 2 fortgefahren.

Der Algorithmus liest jedes Zeichen einmal ein. Die Operationen im Baum benötigen konstante Zeit. Es folgt eine Laufzeit von:

$$O(n).$$

Für ein Dokument, welches der DTD aus Beispiel 4.4.23 unterliegt, muss der Parsebaum der Abbildung 4.4, wie im Algorithmus aus Kapitel 4.4.2 beschrieben, durchlaufen werden.

### Beispiel 4.5.26

```
<Person>
  <Vorname> Klaus </Vorname>
  <Vorname> Ernst </Vorname>
  <Nachname> Müller </Nachname>
  <Adresse>
    <Ort> Wetzlar </Ort>
  </Adresse>
</Person>
```

Das erste Zeichen (Elementname, Person) des Dokumentes wird gelesen. Dieser Name wird mit dem Wert des ersten Knotens des Baums verglichen. Sie stimmen überein.

Das nächste Zeichen (Elementname, Vorname) wird eingelesen. Wiederholen steht nicht auf *True*, deswegen wird mit Punkt 2 gearbeitet und zum nächsten Knoten verzweigt. Das Zeichen wird mit dem Knoten Vorname verglichen. Der Vergleich ist erfolgreich.

Das nächste Zeichen (Elementname, Vorname) wird eingelesen. Wiederholen ist im Knoten Vorname mit *True* markiert, deswegen wird der Elementname wieder mit dem Knoten Vorname verglichen. Der Vergleich ist korrekt, also wird mit Punkt 4 fortgefahren.

Das nächste Zeichen (Elementname, Nachname) wird eingelesen. Es wird wieder mit Vorname verglichen, da Wiederholen auf *True* steht. Der Vergleich ist nicht erfolgreich, deswegen wird mit Punkt 2 fortgefahren und in den nächsten Knoten verzweigt. Dieser wird mit dem Elementnamen verglichen. Der Vergleich ist erfolgreich.

Das nächste Zeichen (Elementname, Adresse) wird eingelesen. Wiederholen steht für das Element Nachname nicht auf *True*, deswegen wird der nächste Knoten angezeigt. Der Vergleich stimmt.

Das nächste Zeichen (Elementname, Ort) wird eingelesen. Der nächste Knoten wird eingelesen, da Wiederholen in Adresse auf *False* steht. Der Vergleich stimmt wieder.

Das nächste Zeichen wird eingelesen, allerdings ist das Ende schon erreicht. Es wird nachgeschaut, ob alle noch offenen Knoten im Baum optional sind. Dies ist der Fall, also stimmt die Analyse.

Der Algorithmus wird an die Syntaxanalyse angebunden, damit das Dokument nicht noch mal abgearbeitet werden muss.

Durch die Ergebnisse der letzten Unterkapitel entsteht die folgende Tabelle.

<b>Sprache</b>	<b>max. Laufz. Parsergenerator</b>	<b>max. Laufz. Parser</b>
XML Grammatiken	$O(n)$	$O(n)$
einfache balancierte Grammatiken	$O(n)$	$O(n)$
konk. bal. Grammatiken	$O(n)$	$O(n)$
RegDyck Grammatiken	$O(n)$	$O(n)$
balancierte Grammatiken durch codet. bal. Grammatiken	$\Omega\left(\frac{m}{n} \cdot 2^{24n \lceil \frac{m}{n} \rceil}\right)$	$O(n)$
balancierte Grammatiken det. Kellerautomat	$O(2^{5n})$	$O(n)$

## Kapitel 5

# Eigenschaften von XML-artigen Sprachen

In dem folgenden Kapitel werden Abschlusseigenschaften der XML-artigen Sprachen untersucht. Darüber hinaus werden Entscheidbarkeitsfragen der balancierten Sprachen geklärt.

Aus praktischer Motivation sind Kenntnisse über Abschlusseigenschaften sehr wichtig. Zum Beispiel können Werkzeuge zum Verbinden oder zum Vereinigen zweier Grammatiken konstruiert werden.

Aus theoretischer Sicht sind die Abschlusseigenschaften und Entscheidbarkeitsfragen Teil der Analyse und geben Aussagen über die Mächtigkeit der Erzeugungsmechanismen.

Die booleschen Operationen Komplement, Vereinigung und Durchschnitt wurden in [BB02a] für die balancierten Sprachen gezeigt. In [BB02b] wurden die booleschen Eigenschaften Durchschnitt, Vereinigung und Komplement der XML Sprachen betrachtet. Das Komplement bezieht sich im Gegensatz zur allgemeinen Definition auf die Menge der Dyckprimes.

Im Folgenden sollen die booleschen Eigenschaften für die einfachen balancierten Sprachen, konkurrenzverhindernden balancierten Sprachen und RegDyck Sprachen untersucht werden.

Konkatenation, Iteration und Substitution werden für alle Sprachtypen in einer erweiterten Form betrachtet, da es in den Wörtern der XML-artigen Sprachen einen festen Start- und Endtag geben muss und z. B. die Konkatenation der Sprachen  $L = \{x\bar{x}\}$  und  $L' = \{y\bar{y}\}$

$$L^\# = \{x\bar{x}y\bar{y}\}$$

ergeben würde. Zusätzlich werden die Operationen Spiegelung und Homomorphismus untersucht.

In dem zweiten Teil des Kapitels werden Entscheidbarkeitsfragen der Spra-

chen geklärt. Für die Praxis ist die Frage interessant, ob eine gegebene balancierte Grammatik eine Grammatik eines Untertyps ist. Damit könnte vorab entschieden werden, welche Art von Parser konstruiert werden muss. So muss nicht für jede balancierte Grammatik eine zeitaufwändige Konstruktion vollzogen werden.

Darüber hinaus wird untersucht, ob entschieden werden kann, ob eine balancierte Grammatik mehrdeutig ist und ob eine balancierte Sprache eine Teilmenge einer anderen ist.

## 5.1 Abschlusseigenschaften

Die XML-artigen Grammatiken wurden in fünf unterschiedliche Typen unterteilt: balancierte Grammatiken, XML Grammatiken, einfache balancierte Grammatiken, konkurrenzverhindernde balancierte Grammatiken und Reg-Dyck Grammatiken.

In den folgenden Unterkapiteln werden die Sprachen der einzelnen Grammatiktypen auf ihre Abschlusseigenschaften hin untersucht.

### 5.1.1 Vereinigung

**Satz 5.1.20** Die balancierten Sprachen sind unter Vereinigung abgeschlossen [BB02a].

**Satz 5.1.21** Die XML Sprachen sind nicht unter Vereinigung abgeschlossen [BB02b].

**Satz 5.1.22** Die einfachen balancierten Sprachen sind nicht unter Vereinigung abgeschlossen.

*Beweis.*  $L_0$  und  $L_1$  sind zwei einfache balancierte Sprachen, die aus den einzigen Wörtern

$$xcd\bar{d}\bar{c}\bar{x} \text{ und} \\ xc\bar{c}\bar{e}\bar{e}\bar{x}$$

bestehen.

Angenommen, die einfachen balancierten Sprachen wären unter Vereinigung abgeschlossen, dann müssten die Sprachen  $L_0$  und  $L_1$  vereinigt wieder eine einfache balancierte Sprache ergeben.

Die Grammatik der Vereinigung  $G = (N, T, S, P)$  der Sprachen muss in

der Startproduktion  $S \rightarrow xm\bar{x}$  im Inhaltsmodell mindestens zwei Nichtterminale  $C_0$  und  $C_1$ , die  $cd\bar{d}\bar{c}$  und  $c\bar{c}$  erzeugen und ein Zeichen  $E$ , welches  $e\bar{e}$  erzeugt ( $S \rightarrow x(C_0|(C_1E))\bar{x}$ ), enthalten.

Wird nur ein Nichtterminal  $C$  benutzt ( $S \rightarrow xC|CE\bar{x}$ ), dann muss es sowohl  $cd\bar{d}\bar{c}$  als auch  $c\bar{c}$  erzeugen. Dies führt dazu, dass auch das Wort

$$xcd\bar{d}\bar{c}e\bar{e}\bar{x}$$

erzeugt wird. Dieses konnte von den ursprünglichen Grammatiken nicht erzeugt werden.

Wenn zwei Nichtterminale zur Definition von  $cd\bar{d}\bar{c}$  und  $c\bar{c}$  verwendet werden, z.B.

$$\begin{aligned} C_0 &\rightarrow cD\bar{c} \\ C_1 &\rightarrow c\bar{c}, \end{aligned}$$

dann muss das Inhaltsmodell  $m = (C_0|(C_1E))$  sein.

Diese Konstruktion ist allerdings in einfachen balancierten Grammatiken nicht erlaubt.

Es dürfen keine zwei Nichtterminale mit dem gleichen Terminal auf der rechten Seite in einer Produktion stehen. Dies führt zum Widerspruch zur Annahme. Die einfachen balancierten Sprachen sind somit nicht unter Vereinigung abgeschlossen.  $\square$

**Satz 5.1.23** Die konkurrenzverhindernden balancierten Sprachen sind nicht unter Vereinigung abgeschlossen.

*Beweis.*  $L_0$  und  $L_1$  sind zwei konkurrenzverhindernde balancierte Sprachen, die aus den Wörtern

$$\begin{aligned} &x(cd\bar{d}\bar{c})^*\bar{x} \text{ und} \\ &x(c\bar{c})^*e\bar{e}\bar{x} \end{aligned}$$

bestehen.

Angenommen, die konkurrenzverhindernden balancierten Sprachen wären unter Vereinigung abgeschlossen, dann müssten die Sprachen  $L_0$  und  $L_1$  vereinigt wieder eine konkurrenzverhindernde balancierte Sprache ergeben.

Die Grammatik der Vereinigung  $G = (N, T, S, P)$  muss mindestens zwei Nichtterminale  $C_0$  und  $C_1$  und ein Nichtterminalzeichen  $E$  in der Startproduktion  $S \rightarrow xm\bar{x}$  im Inhaltsmodell enthalten ( $S \rightarrow x(C_0^*|(C_1^*E))\bar{x}$ ).

Wird nur ein Nichtterminal  $C$  benutzt ( $S \rightarrow x(C^*|(C^*E))\bar{x}$ ), dann muss es sowohl  $cd\bar{d}\bar{c}$  als auch  $c\bar{c}$  erzeugen. Dies führt dazu, dass auch das Wort

$$xcd\bar{d}\bar{c}e\bar{e}\bar{x}$$

erzeugt wird. Dieses konnte von den ursprünglichen Grammatiken nicht erzeugt werden.

Wenn zwei Nichtterminale zur Definition der Wörter  $\{(c\bar{d}\bar{d}\bar{c})^*\}$  und  $\{(c\bar{c})^*\}$  verwendet werden, z.B.

$$\begin{aligned} C_0 &\longrightarrow cD\bar{c} \\ C_1 &\longrightarrow c\bar{c}, \end{aligned}$$

dann muss das Inhaltsmodell  $m = (C_0^* | (C_1^* E))$  sein. Diese Konstruktion ist allerdings in konkurrenzverhindernden balancierten Grammatiken nicht erlaubt, da das Inhaltsmodell  $m$  die Sequenzen  $UC_0W$  und  $UC_1V$  mit  $U, V, W \in N^*$  erzeugt. Dies führt zum Widerspruch zur Annahme. Somit sind die konkurrenzverhindernden balancierten Sprachen nicht unter Vereinigung abgeschlossen.  $\square$

**Satz 5.1.24** Die RegDyck Sprachen sind unter Vereinigung abgeschlossen.

*Beweis.*  $L_0$  und  $L_1$  sind zwei RegDyck Sprachen, die durch die Grammatiken  $G_0 = (N, T, S, P)$  und  $G_1 = (N_1, T_1, S_1, P_1)$  erzeugt werden.

In  $P$  liegt die Startproduktion

$$S_0 \longrightarrow xm_0\bar{x}$$

und in  $P_1$  liegt die Startproduktion

$$S_1 \longrightarrow xm_1\bar{x}.$$

Dann wird eine Menge von Startzeichen erzeugt, in der beide Startzeichen der ursprünglichen Grammatiken enthalten sind.

Die Vereinigung von zwei RegDyck Sprachen ist folgendermaßen zu konstruieren:

Es seien  $L_0$  und  $L_1$  die beiden Sprachen, die vereinigt werden und  $G_0$  und  $G_1$  die erzeugenden Grammatiken.

Ohne Beschränkung der Allgemeinheit gelte  $N_0 \cap N_1 = \emptyset$ . (Sonst erfolgt eine Umbenennung in einer der beiden Mengen, so dass diese Forderung erfüllt ist.) Die Vereinigung der beiden Sprachen wird durch eine neue Grammatik  $G' = (N', T', S', P')$  erzeugt:

- $N' = N_0 \cup N_1$  sind die Nichtterminalzeichen.
- $T' = T_0 \cup T_1$  sind die Terminalzeichen.
- $S_0 \cup S_1$  sind die Startsymbole.
- $P' = P_0 \cup P_1$  sind die Produktionen.

$\square$

### 5.1.2 Durchschnitt

Der Abschluss unter Durchschnitt wurde für die balancierten Sprachen und die XML Sprachen schon in [BB02a] gezeigt. Die Beweisidee für die RegDyck Sprachen kann auch auf die einfachen balancierten Sprachen und die konkurrenzverhindernden balancierten Sprachen übertragen werden. In dem Beweis für die XML Sprachen wird jedoch von Grammatiken mit beschränkten Inhaltsmodellen ausgegangen, obwohl die Inhaltsmodelle der XML Grammatiken aus regulären Ausdrücken bestehen dürfen. Der Abschluss der balancierten Sprachen ist eine Folgerung aus den Ergebnissen, dass die balancierten Sprachen unter Vereinigung und Komplement abgeschlossen sind.

Nun wird der Beweis für die RegDyck Grammatiken durchgeführt. Er kann jedoch auf alle anderen XML-artigen Sprachfamilien übertragen werden.

**Satz 5.1.25** Die RegDyck Sprachen sind unter Durchschnitt abgeschlossen.

*Beweis.*  $L$  und  $L_1$  sind zwei RegDyck Sprachen, die durch die Grammatiken  $G = (N, T, S, P)$  und  $G' = (N', T', S', P')$  erzeugt werden.

Die Grammatik  $G^\# = (N^\#, T^\#, S^\#, P^\#)$ , die den Durchschnitt erzeugt, wird Schritt für Schritt aus den Produktionen der Grammatiken  $G$  und  $G'$  erzeugt.

Wenn

$$X \longrightarrow x\bar{x}$$

eine Produktion in  $G$  ist und

$$X' \longrightarrow x\bar{x}$$

eine Produktion in  $G'$  ist, dann wird in  $G^\#$  eine Produktion

$$(X, X') \longrightarrow x\bar{x}$$

erzeugt. Wenn

$$X \longrightarrow xm\bar{x}$$

eine Produktion in  $G$  ist und

$$X' \longrightarrow xm'\bar{x}$$

eine Produktion in  $G'$  ist, dann wird in  $G^\#$  eine Produktion

$$(X, X') \longrightarrow x(m \cap m')\bar{x}$$

erzeugt. Der Durchschnitt des Inhaltsmodells  $m \cap m'$  wird im Folgenden behandelt.

Es sei  $P_Y, Y \in N$ , die Menge aller Nichtterminalzeichen aus der Grammatik  $G'$ , die mit  $Y$  konkurrieren

$$P_Y = \{Y' \mid Y \longrightarrow ym\bar{y} \in P \wedge Y' \longrightarrow ym'\bar{y} \in P'\},$$

mit  $Y' \in N'$ ,  $y, \bar{y} \in T'$ ,  $m \in \text{Reg}(N)$ ,  $m' \in \text{Reg}(N')$ . Analog gilt dies für die Menge  $P_{Y'}$ ,  $Y' \in P'$ .

Nun werden alle Vorkommen von Nichtterminal  $Z$  in  $m$  durch die Vereinigung

$$\bigcup_{Z' \in P_Z} (Z, Z')$$

ersetzt und die Vorkommen von Nichtterminalzeichen  $Z'$  in  $m'$  durch

$$\bigcup_{Z \in P_{Z'}} (Z, Z')$$

ersetzt und es wird der Durchschnitt  $m \cap m'$  gebildet. Da  $m$  und  $m'$  regulär sind, kann dies durchgeführt werden. Der Durchschnitt zweier balancierter Sprachen ist somit effektiv erzeugbar. □

**Satz 5.1.26** Die balancierten Sprachen sind unter Durchschnitt abgeschlossen [BB02a].

**Satz 5.1.27** Die XML Sprachen sind unter Durchschnitt abgeschlossen [BB02b].

**Satz 5.1.28** Die einfachen balancierten Sprachen sind unter Durchschnitt abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der RegDyck Sprachen. □

**Satz 5.1.29** Die konkurrenzverhindernden balancierten Sprachen sind unter Durchschnitt abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der RegDyck Sprachen. □

### 5.1.3 Komplement

Die Komplementbildung muss für die balancierten Sprachen in einer erweiterten Form betrachtet werden, da das Komplement jeder balancierten Sprache Wörter enthält, die keine Dyckprimes sind und somit auch nicht mehr durch eine balancierte Grammatik erzeugt werden können.

**Definition 5.1.26** Das erweiterte Komplement einer balancierten Sprache  $L$  über einem Alphabet  $T$  ist die Sprache  $D_T \setminus L$ , wobei  $D_T$  die Menge der Dyckprimes über dem Alphabet  $T$  ist.

**Satz 5.1.30** Die balancierten Sprachen sind unter erweitertem Komplement

abgeschlossen [BB02a].

**Satz 5.1.31** Die XML Sprachen sind nicht unter erweitertem Komplement abgeschlossen [BB02b].

**Satz 5.1.32** Die einfachen balancierten Sprachen sind nicht unter erweitertem Komplement abgeschlossen.

*Beweis.* Angenommen, die einfachen balancierten Sprachen wären unter erweitertem Komplement abgeschlossen, dann müsste das erweiterte Komplement einer Sprache  $L$ , die aus dem einzigen Wort

$$xc\overline{c\overline{c\overline{c\overline{x}}}}$$

besteht, durch eine einfache balancierte Grammatik erzeugbar sein. Das bedeutet, im erweiterten Komplement würden unter anderem alle Wörter  $xd_c d_{c_0} \overline{x}$  liegen. Dabei ist  $d_c$  die Menge aller Dyckprimes, die mit  $c$  beginnen und  $d_{c_0}$  die Menge aller Dyckprimes, die mit  $c$  beginnen, ohne das Wort  $c\overline{c\overline{c}}$ . Die Grammatik müsste somit eine Produktion  $X \rightarrow xCC_1\overline{x}$  mit mindestens zwei verschiedenen Nichtterminalzeichen  $C, C_1$  besitzen, da zwei unterschiedliche Sprachen erzeugt werden müssten. Dieses wären zwei konkurrierende Nichtterminalzeichen, da beide Nichtterminalzeichen Sprachen erzeugen, die mit  $c$  beginnen. Dies ist in einer einfachen balancierten Sprache im Inhaltsmodell nicht erlaubt, woraus sich der Widerspruch ergibt.  $\square$

**Satz 5.1.33** Die konkurrenzverhindernden balancierten Sprachen sind nicht unter erweitertem Komplement abgeschlossen.

*Beweis.* Angenommen, die konkurrenzverhindernden balancierten Sprachen wären unter erweitertem Komplement abgeschlossen, dann müsste das erweiterte Komplement einer Sprache  $L$ , die aus dem einzigen Wort

$$xc\overline{c\overline{c\overline{c\overline{x}}}}$$

besteht, durch eine konkurrenzverhindernde balancierte Grammatik erzeugbar sein. Das bedeutet, im erweiterten Komplement würden alle Wörter  $xd_c \overline{x}$  und  $xd_{c_0} \overline{x}$  liegen. Dabei ist  $d_c$  die Menge aller Dyckprimes, die mit  $c$  beginnen und  $d_{c_0}$  die Menge aller Dyckprimes, die mit  $c$  beginnen, ohne das Wort  $c\overline{c\overline{c}}$ . Die Grammatik müsste somit eine Produktion  $X \rightarrow x(C|C_1)\overline{x}$  mit mindestens zwei verschiedenen Nichtterminalzeichen  $C, C_1$  besitzen, da zwei unterschiedliche Sprachen erzeugt werden müssten. Dieses wären zwei konkurrierende Nichtterminalzeichen, da beide Nichtterminalzeichen Sprachen erzeugen, die mit  $c$  beginnen und aus dem Inhaltsmodell können Wörter der Form  $UCW$  und  $UC_1V$  abgeleitet werden, wobei  $U, W, V$  Folgen von Nichtterminalzeichen sind, die auch leer sein können.

Dies ist in einer konkurrenzverhindernden balancierten Sprache im Inhaltsmodell nicht erlaubt, woraus sich der Widerspruch ergibt.  $\square$

**Satz 5.1.34** Die RegDyck Sprachen sind nicht unter erweitertem Komplement abgeschlossen.

*Beweis.* Es sei  $L = \{x^m x^k \bar{x}^k x^n \bar{x}^n \bar{x}^m \mid k, m, n \in \mathbb{N}\}$  eine RegDyck Sprache. Angenommen, die RegDyck Sprachen wären unter Komplement abgeschlossen, dann muss eine RegDyck Grammatik existieren, die das erweiterte Komplement der Sprache  $L$  erzeugt.

In dem erweitertem Komplement von  $L$  sind auch die Wörter der Form  $xx^n \bar{x}^n x^m \bar{x}^m \bar{x}$ ,  $n, m \in \mathbb{N}$  enthalten. Um diese Wörter zu erzeugen, muss eine Startproduktion der Form  $S \rightarrow xXX\bar{x}$  existieren. Es muss genau zwei Nichtterminalzeichen geben, die selbsteinbettend sind und jeweils  $x^i \bar{x}^i$  erzeugen. Dies führt dazu, dass  $FA(X) = x^* \bar{x}^*$ . Das bedeutet,  $X$  steht zu  $X$  im Konflikt. Im Durchschnitt  $(x^* \bar{x}^*)(x^* \bar{x}^*) \cap D_{\{x, \bar{x}\}} \neq \emptyset$  liegen auch die Worte  $xx^n \bar{x}^n x^m \bar{x}^m \bar{x}$ ,  $n, m \in \mathbb{N}$ . Diese sind jedoch in  $L$  enthalten und dürfen somit nicht im Komplement vorkommen. Das bedeutet, der Konflikt kann nicht aufgehoben werden. Es ergibt sich ein Widerspruch zur Annahme, da das erweiterte Komplement der Sprache  $L$  nicht durch eine RegDyck Grammatik erzeugt werden kann.

Es folgt, dass die RegDyck Sprachen nicht unter erweitertem Komplement abgeschlossen sind.  $\square$

Die balancierten Sprachen sind unter allen booleschen Eigenschaften abgeschlossen. Für die Praxis erscheint jedoch vor allem wichtig zu sein, ob die Untertypen der balancierten Sprachen abgeschlossen sind. Dabei zeigt sich, dass die RegDyck Sprachen unter den Operationen Vereinigung und Durchschnitt abgeschlossen sind. Die in der Praxis verwendeten Sprachen (XML Sprachen, einfache balancierte Sprachen) und die konkurrenzverhindernden balancierten Sprachen sind lediglich unter Durchschnitt abgeschlossen.

#### 5.1.4 Iteration

Es sei  $L'$  eine balancierte Sprache mit den Wörtern  $\{x^n \bar{x}^n \mid n \geq 1\}$ . Wird diese Sprache iteriert, so entsteht die Sprache

$$L = \{x^{k_1} \bar{x}^{k_1} x^{k_2} \bar{x}^{k_2} \dots x^{k_p} \bar{x}^{k_p} \mid k_i \geq 1, 1 \leq i \leq p\} \cup \{\lambda\}.$$

Die Sprache ist nach Definition nicht balanciert, da es ein äußeres schließendes Tagpaar geben muss.

**Definition 5.1.27** Die erweiterte Iteration einer balancierten Sprache  $L$  wird im Folgenden definiert.

Sei  $G = (N, T, S, P)$  eine balancierte Grammatik, die  $L$  erzeugt. Die iterierte

Sprache  $L^*$  wird in ein neues Tagpaar  $s', \bar{s}' \notin T$  eingeschlossen. Es wird eine neue Grammatik  $G' = (N', T', S', P')$  erzeugt mit

- $N' = N \cup \{S'\}$
- $T' = T \cup \{s', \bar{s}'\}$
- $S' \notin N$
- $P' = P \cup \{S' \longrightarrow s'(\bigcup_{S' \in S} S')^* \bar{s}'\}$ .

$G'$  ist offensichtlich wieder eine balancierte Grammatik. □

**Satz 5.1.35** Die balancierten Sprachen sind unter der erweiterten Iteration abgeschlossen.

*Beweis.* Aus der Definition der erweiterten Iteration folgt die Behauptung. □

**Satz 5.1.36** Die XML Sprachen sind unter der erweiterten Iteration abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen. □

**Satz 5.1.37** Die einfachen balancierten Sprachen sind unter der erweiterten Iteration abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen. □

**Satz 5.1.38** Die konkurrenzverhindernde balancierte Sprachen sind unter der erweiterten Iteration abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen. □

**Satz 5.1.39** Die RegDyck Sprachen sind nicht unter der erweiterten Iteration abgeschlossen.

*Beweis.* Es sei  $L$  die folgende RegDyck Sprache

$$L = \{x^n \bar{x}^n \mid n \geq 1\}.$$

In der Iteration sind z.B. die Wörter

$$K = \{s' x^k \bar{x}^k x^l \bar{x}^l \bar{s}' \mid k \geq 1, l \geq 1\}$$

enthalten.

Angenommen, die RegDyck Sprachen wären unter Komplement abgeschlossen, dann müsste es eine Startproduktion

$$S' \longrightarrow s'S^*\bar{s}'$$

geben, wobei es mindestens ein Nichtterminalzeichen  $S$  gibt, welches mit einem Stern versehen ist, welches die Sprache  $L$  erzeugt. Jedoch lässt sich aus dem Inhaltsmodell von  $S'$  ein Wort  $SS$  ableiten. Da  $S$  selbsteinbettend ist,  $FA(S) = x^+\bar{x}^+$  und  $FA(S)FA(S) \cap D \neq \emptyset$ , steht  $S$  zu  $S$  im Konflikt. Dies ist ein Widerspruch zur Definition der RegDyck Sprachen.  $\square$

### 5.1.5 Konkatenation

Es seien  $L$  und  $L_0$  zwei balancierte Sprachen, welche konkateniert werden sollen und  $G = (N, T, S, P)$  und  $G_0 = (N_0, T_0, S_0, P_0)$  die erzeugenden Grammatiken.  $L$  erzeuge eine beliebige Zeichenkette der Form  $x \dots \bar{x}$  und  $L_0$  erzeuge eine beliebige Zeichenkette der Form  $x_0 \dots \bar{x}_0$ , wobei  $x, \bar{x} \in T$  und  $x_0, \bar{x}_0 \in T_0$ .

Die Konkatenation würde eine Zeichenkette der Form

$$x \dots \bar{x}x_0 \dots \bar{x}_0 \text{ oder} \\ x_0 \dots \bar{x}_0x \dots \bar{x}$$

ergeben. Diese Form ist nach Definition nicht erlaubt, da jede Produktion von einem Tagpaar umgeben sein muss.

Es sei auch hier eine abgewandelte Form der Konkatenation vorgestellt.

**Definition 5.1.28** Die erweiterte Konkatenation von zwei balancierten Sprachen  $L$  und  $L_0$  wird wie folgt definiert:

Sei  $G = (N, T, S, P)$  eine balancierte Grammatik, die  $L$  erzeugt und  $G_0 = (N_0, T_0, S_0, P_0)$  eine balancierte Grammatik, die  $L_0$  erzeugt. Die konkatenierte Sprache  $LL_0$  wird in ein neues Tagpaar  $s', \bar{s}' \notin (T \cup T_0)$  eingeschlossen. Es wird eine neue Grammatik  $G' = (N', T', S', P')$  erzeugt mit

- $N' = N \cup N_0 \cup \{S'\}$
- $T' = T \cup T_0 \cup \{s', \bar{s}'\}$
- $S' \notin (N \cup N_0)$
- $P' = P \cup P_0 \cup \{S' \longrightarrow s'(\bigcup_{S' \in S} S')(\bigcup_{S'' \in S_0} S'')\bar{s}'\}$ .

Analog wird eine neue Grammatik  $G' = (N', T', S', P')$  für die konkatenierte Sprache  $L_0L$  erzeugt:

- $N' = N \cup N_0 \cup \{S'\}$
- $T' = T \cup T_0 \cup \{s', \bar{s}'\}$
- $S' \notin (N \cup N_0)$

- $P' = P \cup P_0 \cup \{S' \longrightarrow s'(\bigcup_{S'' \in S_0} S'')(\bigcup_{S' \in S} S')\bar{s}'\}$

$G'$  ist in beiden Fällen offensichtlich wieder eine balancierte Grammatik.  $\square$

**Satz 5.1.40** Die balancierten Grammatiken sind unter der erweiterten Konkatenation abgeschlossen.

*Beweis.* Die Behauptung folgt aus der Definition der erweiterten Konkatenation.  $\square$

**Satz 5.1.41** Die XML Sprachen sind nicht unter der erweiterten Konkatenation abgeschlossen.

*Beweis.* Wir nehmen an, die XML Sprachen wären unter erweiterter Konkatenation abgeschlossen.

Es seien  $L$  und  $L_0$  zwei XML Sprachen, die konkateniert werden sollen.  $G$  und  $G_0$  seien die beiden erzeugenden XML Grammatiken.

Die Grammatik  $G_0$  kann konkurrierende Nichtterminale in Bezug zu den Nichtterminalen der Grammatik  $G$  enthalten. Für eine XML Grammatik  $G' = (N', T', S', P')$ ,  $T' = (A' \cup \bar{A}')$  muss gelten:

$\forall X, Y \in N' \forall m, m_0 \in \text{Reg}(N')$

$$\neg \exists x \in A' : (X \longrightarrow xm\bar{x} \in P' \wedge Y \longrightarrow xm_0\bar{x} \in P' \wedge X \neq Y)$$

Dies ist für die Konkatenation zweier XML Sprachen nicht sichergestellt, da es in beiden Grammatiken Produktionen mit den gleichen Terminalzeichen geben kann.

$G$  bestehe aus der einzigen Produktion

$$S \longrightarrow s\bar{s}$$

und  $G_0$  bestehe aus den beiden Produktionen

$$S_0 \longrightarrow sX\bar{s} \text{ und} \\ X \longrightarrow x\bar{x}.$$

Die Grammatik  $G'$  der Konkatenation muss die Startproduktion

$$S' \longrightarrow s'SS_0\bar{s}'$$

beinhalten. Genau wie bei der Vereinigung ist es nicht möglich, die Sprache, die  $S$  erzeugt und die Sprache, die  $S_0$  erzeugt, durch nur ein Nichtterminalzeichen auszudrücken, um die Konkurrenz der Nichtterminalzeichen  $S$  und  $S_0$  zu verhindern.

In XML Grammatiken sind keine konkurrierenden Nichtterminalzeichen erlaubt, somit sind die XML Sprachen nicht unter Konkatenation abgeschlossen.  $\square$

Es seien  $L$  und  $L_0$  zwei XML Sprachen und  $G$  und  $G_0$  die erzeugenden Grammatiken. Es wird angenommen, dass in  $G$  und  $G_0$  nur Terminalzeichen existieren, die auch in den Produktionen benutzt werden. Eine Konkatenation der Sprachen  $L$  und  $L_0$  kann auf jeden Fall durchgeführt werden, wenn

$$T \cap T_0 = \emptyset$$

gilt.

**Satz 5.1.42** Die einfachen balancierten Sprachen sind nicht unter der erweiterten Konkatenation abgeschlossen.

*Beweis.* Wir nehmen an, die einfachen balancierten Sprachen wären unter Konkatenation abgeschlossen.

Es seien  $L$  und  $L_0$  zwei einfache balancierte Sprachen, die konkateniert werden sollen.  $G$  und  $G_0$  seien die beiden erzeugenden einfachen balancierten Grammatiken.

Die Grammatik  $G_0$  kann konkurrierende Nichtterminale in Bezug zu den Nichtterminalen der Grammatik  $G$  enthalten. Für einfache balancierte Grammatiken  $G' = (N', T', S', P')$ ,  $T' = (A' \cup \bar{A}')$  muss jedoch gelten:

$$\forall X, Y, Z \in N' \forall z, \bar{z} \in T' \forall m, m_0, m_1 \in \text{Reg}(N') \forall \alpha, \beta, \gamma \in N^*$$

$$\neg \exists x \in A' : (X \longrightarrow xm\bar{x} \in P' \wedge Y \longrightarrow xm_0\bar{x} \in P' \wedge X \neq Y \wedge Z \longrightarrow zm_1\bar{z} \in P' \wedge \alpha X \beta Y \gamma \in L(m_1))$$

Dies ist für die Konkatenation zweier einfachen balancierten Sprachen nicht sichergestellt, da eine neue Produktion erzeugt werden muss, in der beide Startsymbole enthalten sind. Im Fall, dass die Startproduktionen miteinander konkurrieren, verstößt die neue Produktion gegen die Definition der einfachen balancierten Grammatiken.

Der Beweis erfolgt über die gleichen Grammatiken wie bei den XML Sprachen.

$G$  bestehe aus der einzigen Produktion

$$S \longrightarrow s\bar{s}$$

und  $G_0$  bestehe aus den beiden Produktionen

$$S_0 \longrightarrow sX\bar{s} \text{ und} \\ X \longrightarrow x\bar{x}.$$

Die Grammatik  $G'$  der Konkatenation muss die Startproduktion

$$S' \longrightarrow s'SS_0\bar{s}'$$

beinhalten. Genau wie bei der Vereinigung ist es nicht möglich, die Sprache, die  $S$  erzeugt und die Sprache, die  $S_0$  erzeugt, durch nur ein Nichtterminalzeichen auszudrücken, um die Konkurrenz der Nichtterminalzeichen  $S$  und  $S_0$  zu verhindern.

In einfachen balancierten Grammatiken sind keine konkurrierenden Nichtterminalzeichen innerhalb einer Produktion erlaubt, somit sind die XML Sprachen nicht unter Konkatenation abgeschlossen.  $\square$

Es seien  $L$  und  $L_0$  zwei einfache balancierte Sprachen und  $G$  und  $G_0$  die erzeugenden Grammatiken. Eine Konkatenation der Sprachen  $L$  und  $L_0$  kann auf jeden Fall durchgeführt werden, wenn

$$\forall m \in \text{Reg}(N) \forall m_0 \in \text{Reg}(N_0) \\ \neg \exists x \in A : (S \longrightarrow xm\bar{x} \in P \wedge S_0 \longrightarrow xm_0\bar{x} \in P_0).$$

**Satz 5.1.43** Die konkurrenzverhindernden balancierten Sprachen sind unter der erweiterten Konkatenation abgeschlossen.

*Beweis.* Es seien  $L$  und  $L_0$  zwei konkurrenzverhindernde balancierte Sprachen, die konkateniert werden sollen.  $G = (N, T, S, P)$  und  $G_0 = (N_0, T_0, S_0, P_0)$  seien die beiden erzeugenden konkurrenzverhindernden balancierten Grammatiken.

Die Grammatik  $G_0$  kann konkurrierende Nichtterminale in Bezug zu den Nichtterminalen der Grammatik  $G$  enthalten. Für eine konkurrenzverhindernde balancierte Grammatik  $G' = (N', T', S', P')$ ,  $T' = (A' \cup \bar{A}')$ , muss jedoch gelten:

$$\forall m, m_0, m_1 \in \text{Reg}(N') \forall X, Y, Z \in N' \forall z \in A' \forall U, V, W \in N^* \\ \neg \exists x \in A' : (X \longrightarrow xm\bar{x} \in P' \wedge Y \longrightarrow xm_0\bar{x} \in P' \wedge X \neq Y \wedge \\ Z \longrightarrow zm_1\bar{z} \in P' \wedge UXW \in L(m_1) \wedge UYV \in L(m_1))$$

Angenommen, die konkurrenzverhindernden balancierten Sprachen wären nicht unter Konkatenation abgeschlossen, dann müsste es eine Sprache geben, die nur durch eine Grammatik beschrieben werden kann, dessen Inhaltsmodell  $m_1$  mindestens zwei konkurrierende Nichtterminalzeichen  $X$  und  $Y$  enthält, für die gilt:  $UXW \in L(m_1)$  und  $UYV \in L(m_1)$ . Die Grammatiken  $G$  und  $G_0$  enthalten keine Produktionen für die dieses gilt, da sie konkurrenzverhindernd balanciert sind. Es bleibt die neue Startproduktion  $S' \longrightarrow s'SS_0\bar{s}'$ , bzw.  $S' \longrightarrow s'S_0S\bar{s}'$  zu untersuchen.

Da bei der Konkatenation eine eindeutige Reihenfolge  $LL_0$  oder  $L_0L$  entsteht, kann in der Startproduktion  $S'$  der Sprache der Konkatenation  $L'$  nur genau eine Sequenz der Startsymbole entstehen, entweder  $SS_0$  oder  $S_0S$  und somit also nicht  $USW$  und  $US_0V$ . Dies führt zum Widerspruch zur Annahme. Die konkurrenzverhindernden balancierten Grammatiken sind somit unter Konkatenation abgeschlossen.  $\square$

**Satz 5.1.44** Die RegDyck Sprachen sind nicht unter der erweiterten Konkatenation abgeschlossen.

*Beweis.* Angenommen, die RegDyck Sprachen wären unter erweiterter Konkatenation abgeschlossen. Es seien  $L$  und  $L_0$  zwei RegDyck Sprachen, die durch die RegDyck Grammatiken  $G = (N, T, S, P)$  und  $G_0 = (N_0, T_0, S_0, P_0)$  erzeugt werden, wobei

$$\begin{aligned} L &= \{x^n \bar{x}^n \mid n \geq 1\} \text{ und} \\ L_0 &= \{x^m \bar{x}^m \mid m \geq 1\} \\ LL_0 &= L' = \{sx^n \bar{x}^n x^m \bar{x}^m \bar{s} \mid n \geq 1, m \geq 1\} \end{aligned}$$

ergeben.

Die Konkatenation muss durch eine balancierte Grammatik mit der Startproduktion  $S \rightarrow sXX\bar{s}$  erzeugt werden, wobei die Bezeichnung der Nichtterminalzeichen unwichtig ist. Die beiden Nichtterminalzeichen müssen selbst-einbettend sein oder in ein selbsteinbettendes Nichtterminalzeichen ableiten. Es muss gelten:  $FA(X) = (x)^+(\bar{x})^+$ . Daraus folgt, dass  $FA(X)FA(X) \cap D_{(T \cup T_0)} \neq \emptyset$ . Die balancierte Grammatik, die die Konkatenation der beiden Sprachen  $L$  und  $L_0$  beschreibt, kann somit keine RegDyck Grammatik sein. Dies führt zu einem Widerspruch zur Annahme.

Die konkurrenzverhindernden balancierten Sprachen sind nicht unter erweiterter Konkatenation abgeschlossen. □

Die Abschlusseigenschaften „erweiterte Iteration“ und „erweiterte Konkatenation“ sind in der Praxis sehr wichtig. Die balancierten Sprachen und die konkurrenzverhindernden balancierten Sprachen sind unter beiden Operationen abgeschlossen. Die XML Sprachen und einfachen balancierten Sprachen sind unter Iteration abgeschlossen. Wie in Kapitel 4 gezeigt wurde, benötigt ein Parsergenerator gleich viel Zeit, um einen Parser für XML Sprachen zu konstruieren wie für einen Parser für konkurrenzverhindernde balancierte Sprachen. Jedoch zeigt sich aus den Abschlusseigenschaften, dass es sinnvoll ist, konkurrenzverhindernde balancierte Sprachen in der Praxis zu verwenden. Darüber hinaus zeigte sich in Kapitel 4, dass die konkurrenzverhindernden balancierten Sprachen ausdrucksmächtiger sind, als die XML Sprachen. Die RegDyck Sprachen sind weder unter erweiterter Iteration noch unter erweiterter Konkatenation abgeschlossen, was die Praxisrelevanz nicht mindert, weil sie unter Vereinigung und Durchschnitt abgeschlossen sind.

### 5.1.6 Spiegelung

Die Spiegelung wird ähnlich, wie in der Konstruktion aus [HU79] für kontextfreie Grammatiken, durchgeführt.

**Satz 5.1.45** Die balancierten Sprachen sind unter Spiegelung abgeschlossen.

*Beweis.* Es sei  $G = (N, T, S, P), T = (A \cup \bar{A})$  eine balancierte Grammatik. Hierbei muss erlaubt werden, dass Endtags vor Starttags stehen dürfen, da jede Produktion

$$X \longrightarrow xm\bar{x} \in P$$

umgedreht wird zu

$$X \longrightarrow \bar{x}mx.$$

Die gespiegelte Grammatik sei  $G' = (N', T', S', P')$ .

Das Vertauschen von Start- und Endtag erfordert zusätzlich

$$A' = \bar{A}, \bar{A}' = A.$$

In einem nächsten Schritt werden alle Inhaltsmodelle gespiegelt. Es sei

$$X \longrightarrow xm\bar{x} \in P$$

eine Produktion in  $G$  mit dem Inhaltsmodell  $m$ . Dann sei das Inhaltsmodell der Produktion der gespiegelten Sprache  $m^R$  und die Produktion

$$X \longrightarrow \bar{x}mx \in P'.$$

Die Spiegelung des Inhaltsmodells ist möglich, da die regulären Sprachen unter Spiegelung abgeschlossen sind.  $\square$

**Satz 5.1.46** Die XML Sprachen sind unter Spiegelung abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.47** Die einfachen balancierten Sprachen sind unter Spiegelung abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.48** Die konkurrenzverhindernden balancierten Sprachen sind unter Spiegelung abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.49** Die RegDyck Sprachen sind unter Spiegelung abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

### 5.1.7 Substitution

Es sei  $G = (N, T, S, P)$  eine balancierte Grammatik, welche eine Sprache  $L$  mit den Wörtern  $x \dots \bar{x}$  erzeugt, wobei  $x, \bar{x} \in T$ . Für alle Terminalzeichen  $y \in T$  der Sprache  $L$  wird eine neue Sprache  $L_y$  festgelegt. Vorerst wird ein Wort  $w$  der Sprache  $L$  abgeleitet. Danach wird für jedes Zeichen  $y \in T$  des Wortes  $w$  ein Wort der Sprache  $L_y$  erzeugt.

Es sei  $a_0 \dots \bar{a}_0$  ein Wort der Sprache  $L_x$  und  $a_1 \dots \bar{a}_1$  ein Wort der Sprache  $L_{\bar{x}}$ . Dann würde das resultierende Wort die Form

$$a_0 \dots \bar{a}_0 \dots a_1 \dots \bar{a}_1$$

haben.

Dies verstößt allerdings gegen die Definition der balancierten Sprachen.

Es wird ein äußeres schließendes Tagpaar benötigt, so dass

$$s a_0 \dots \bar{a}_0 \dots a_1 \dots \bar{a}_1 \bar{s}$$

erzeugt wird.

Es sei auch hier eine abgewandelte Form der Substitution vorgestellt.

**Definition 5.1.29** Die erweiterte Substitution einer balancierten Sprache  $L$  wird wie folgt definiert:

Sei  $G = (N, T, S, P)$  eine balancierte Grammatik, die  $L$  erzeugt.

Für jedes Zeichen  $a \in T$  wird eine neue balancierte Sprache  $L_a$  definiert. Es wird ein Wort aus  $L$  erzeugt und jedes Vorkommen von  $a$  durch ein Wort aus  $L_a$  ersetzt. Es sei  $w$  das entstehende Wort aus der Substitution, dann wird das Wort durch ein disjunktes Paar aus Start- und Endtag  $s, \bar{s}$  neu eingeklammert,  $sw\bar{s}$ .

**Satz 5.1.50** Die balancierten Sprachen sind nicht unter der erweiterten Substitution abgeschlossen.

*Beweis.* Es sei  $G = (N, T, \{S\}, P)$  eine balancierte Grammatik mit der einzigen Produktion

$$S \longrightarrow xS\bar{x}.$$

Es kann ein Wort

$$xxxx\bar{xxxx}$$

erzeugt werden.

Im Allgemeinen würde die Sprache  $\{x^n \bar{x}^n \mid n \geq 1\}$  erzeugt werden.

Alle Zeichen  $x$  werden durch das Wort  $b\bar{b}$  substituiert. Alle Zeichen  $\bar{x}$  werden durch das Wort  $c\bar{c}$  substituiert. Das entstehende Wort

$$\bar{b}\bar{b}\bar{b}\bar{b}\bar{b}\bar{b}\bar{c}\bar{c}\bar{c}\bar{c}\bar{c}\bar{c}\bar{c}$$

wird von einem Tagpaar  $s, \bar{s}$  umschlossen und es entsteht

$$s\bar{b}\bar{b}\bar{b}\bar{b}\bar{b}\bar{b}\bar{c}\bar{c}\bar{c}\bar{c}\bar{c}\bar{c}\bar{s}.$$

Im Allgemeinen ergibt sich die Sprache  $\{s(\bar{b}\bar{b})^n(\bar{c}\bar{c})^n\bar{s} \mid n \geq 1\}$ .

Angenommen, die balancierten Sprachen wären unter Substitution abgeschlossen, dann muss es eine Startproduktion in der substituierten Grammatik geben:

$$S \longrightarrow sm\bar{s}, m \in \text{Reg}(N).$$

Die Startproduktion muss folgendes Inhaltsmodell

$$m = B^n C^n$$

enthalten, da beide Dyckprimes  $\bar{b}\bar{b}$  und  $\bar{c}\bar{c}$  durch zwei verschiedene Nicht-terminalzeichen erzeugt werden und beide mit  $n$  potenziert sind. Das Inhaltsmodell muss eine kontextfreie Sprache beschreiben. Die Inhaltsmodelle dürfen aber nur reguläre Ausdrücke enthalten. Dies ist ein Widerspruch zur Annahme.

Es ergibt sich, dass die balancierten Sprachen nicht unter der erweiterten Substitution abgeschlossen sind.  $\square$

**Satz 5.1.51** Die XML Sprachen sind nicht unter der erweiterten Substitution abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.52** Die einfachen balancierten Sprachen sind nicht unter der erweiterten Substitution abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.53** Die konkurrenzverhindernden balancierten Sprachen sind nicht unter der erweiterten Substitution abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.54** Die RegDyck Sprachen sind nicht unter der erweiterten Substitution abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

### 5.1.8 Homomorphismus

**Satz 5.1.55** Die balancierten Sprachen sind nicht unter Homomorphismus abgeschlossen.

*Beweis.* Es sei  $G = (N, T, \{S\}, P), T = (A \cup \bar{A})$  eine balancierte Grammatik mit der einzigen Produktion

$$S \longrightarrow x\bar{x}.$$

Der Homomorphismus  $h : T \longrightarrow \{x\}$  ist definiert durch

$$\begin{aligned} h(x) &\longrightarrow x \\ h(\bar{x}) &\longrightarrow x. \end{aligned}$$

Die Sprache von  $G$  ist  $L(G) = \{x\bar{x}\}$ . Wenn der Homomorphismus angewendet wird, entsteht die Sprache  $\{xx\}$ . Diese Sprache kann nicht durch eine balancierte Grammatik erzeugt werden, da jede Produktion der balancierten Grammatiken aus mindestens einem Zeichen  $y \in A$  und aus einem Zeichen  $\bar{y} \in \bar{A}$  bestehen muss.  $\square$

**Satz 5.1.56** Die XML Sprachen sind nicht unter Homomorphismus abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.57** Die einfachen balancierten Sprachen sind nicht unter Homomorphismus abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.58** Die konkurrenzverhindernden balancierten Sprachen sind nicht unter Homomorphismus abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

**Satz 5.1.59** Die RegDyck Sprachen sind nicht unter Homomorphismus abgeschlossen.

*Beweis.* Der Beweis ist analog zu dem der balancierten Sprachen.  $\square$

Es hat sich gezeigt, dass alle XML-artigen Sprachen unter Spiegelung abgeschlossen sind, wenn es zulässig ist, dass Starttags und Endtags vertauscht werden. Unter Homomorphismus und Substitution ist keine XML-artige Sprache abgeschlossen.

Zusammenfassend kann gesagt werden, dass aus Sicht der Abschlusseigenschaften die balancierten Sprachen in der Praxis bevorzugt werden sollten. Wenn diese aus Gründen großer Laufzeiten bei der Erstellung eines Parsers nicht gewählt werden können, sollten entweder die konkurrenzverhindernden

balancierten Sprachen oder die RegDyck Sprachen gewählt werden. Sollen die XML Dokumente, mit denen gearbeitet wird, vereinigt werden, dann müssen RegDyck Grammatiken verwendet werden. Werden die Dokumente konkateniert oder iteriert, sind die konkurrenzverhindernden balancierten Sprachen auszuwählen.

## Zusammenfassung der Abschlusseigenschaften

Abschlusseigenschaften	reg. Spr.	XMLS	EBL	KVBL	RegDyck Spr.	BL	det. kon. Spr.	kontextf. Spr.
Vereinigung	+	-	-	-	+	+	-	+
Durchschnitt	+	+	+	+	+	+	-	-
erw. Komplement	+	-	-	-	-	+	+	-
Spiegelung	+	+	+	+	+	+	-	+
erw. Konkatenation	+	-	-	+	-	+	-	+
erw. Iteration	+	+	+	+	-	+	-	+
erw. Substitution	+	-	-	-	-	-	-	+
Homomorphismus	+	-	-	-	-	-	-	+

**reg. Spr.** - reguläre Sprachen, **XMLS** - XML Sprachen, **EBL** - einfache balancierte Sprachen, **KVBL** - konkurrenzverhindernde balancierte Sprachen, **RegDyck Spr.** - RegDyck Sprachen, **BL** - balancierte Sprachen, **det. kon. Spr.** - deterministisch kontextfreie Sprachen, **kontextf. Spr.** - kontextfreie Sprachen

### 5.1.9 Beispiele für Abschlusseigenschaften von XML-artigen Grammatiken

#### Erweiterte Konkatenation von balancierten Sprachen

$$G_1 = (N_1 = \{S_1, A\}, T_1 = \{x, \bar{x}, a, \bar{a}\}, \{S_1\}, P_1 = \{S_1 \longrightarrow xA^*\bar{x}, A \longrightarrow a\bar{a}\})$$

$$G_2 = (N_2 = \{S_2, B\}, T_2 = \{x, \bar{x}, b, \bar{b}\}, \{S_2\}, P_2 = \{S_2 \longrightarrow xB^*\bar{x}, B \longrightarrow b\bar{b}\})$$

$$G_{1,2} = (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, S, P_1 \cup P_2 \cup \{S \longrightarrow sS_1S_2\bar{s}\})$$

$G_{1,2}$  ist die Grammatik, die die erweiterte Konkatenation von  $G_1$  und  $G_2$  erzeugt. Sie ist wieder eine balancierte Grammatik. Durch sie kann z.B. das Wort

$$sxa\bar{a}a\bar{a}\bar{x}\bar{x}\bar{s}$$

erzeugt werden. □

#### Erweiterte Konkatenation von konkurrenzverhindernden balancierten Sprachen

$$G_1 = (N_1 = \{S_1, A, B\}, T_1 = \{x, \bar{x}, a, \bar{a}, b, \bar{b}\}, S_1, P_1 = \{S_1 \longrightarrow xA^*B\bar{x}, A \longrightarrow a\bar{a}, B \longrightarrow b\bar{b}\})$$

$$G_2 = (N_2 = \{S_2, D\}, T_2 = \{x, \bar{x}, d, \bar{d}\}, S_2, P_2 = \{S_2 \longrightarrow xD^*\bar{x}, D \longrightarrow d\bar{d}\}) :$$

$$G_{1,2} = (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, S, P_1 \cup P_2 \cup \{S \longrightarrow sS_1S_2\bar{s}\})$$

Die erstellte Grammatik  $G_{1,2}$  erzeugt die erweiterte Konkatenation von  $G_1$  und  $G_2$ . Sie ist wieder eine konkurrenzverhindernde balancierte Grammatik, da im Inhaltsmodell der Produktion  $S$  keine konkurrierende Nichtterminale enthalten sind.

Diese Grammatik erzeugt z.B. das Wort

$$sxa\bar{a}a\bar{a}b\bar{b}\bar{x}d\bar{d}\bar{x}\bar{s}.$$

□

#### Erweiterte Konkatenation von einfachen balancierten Sprachen

$$G_1 = (N_1 = \{S_1, A\}, T_1 = \{x, \bar{x}, a, \bar{a}\}, \{S_1\}, P_1 = \{S_1 \longrightarrow xAA^*\bar{x}, A \longrightarrow a\bar{a}\})$$

$$G_2 = (N_2 = \{S_2, B\}, T_2 = \{x, \bar{x}, a, \bar{a}\}, \{S_2\}, P_2 = \{S_2 \longrightarrow xB^*\bar{x}, B \longrightarrow a\bar{a}\})$$

$$G_{1,2} = (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, S, P_1 \cup P_2 \cup \{S \longrightarrow sS_1S_2\bar{s}\})$$

$G_{1,2}$  erzeugt die erweiterte Konkatenation der Grammatiken  $G_1$  und  $G_2$ . Allerdings ist diese Grammatik keine einfache balancierte Grammatik, da die Variablen  $S_1$  und  $S_2$  konkurrieren und somit die Bedingung für eine einfache balancierte Grammatik nicht erfüllt wird.

$$G_3 = (N_3 = \{S_3, A\}, T_3 = \{x, \bar{x}, a\bar{a}\}, S_3, P_3 = \{S_3 \longrightarrow xA^? \bar{x}, A \longrightarrow a\bar{a}\})$$

$$G_4 = (N_4 = \{S_4, B\}, T_4 = \{y, \bar{y}, b\bar{b}\}, S_4, P_4 = \{S_4 \longrightarrow yB^* \bar{y}, B \longrightarrow b\bar{b}\})$$

$$G_{3,4} = (N_3 \cup N_4 \cup \{S\}, T_3 \cup T_4, S, P_3 \cup P_4 \cup S \longrightarrow sS_3S_4\bar{s})$$

$G_{3,4}$  ist die Konkatenation der Grammatiken  $G_3$  und  $G_4$ . Sie ist eine einfache balancierte Grammatik, da  $S_3$  und  $S_4$  im Inhaltsmodell nicht konkurrieren.  $\square$

### Erweiterte Iteration einer balancierten Sprache

$$G_1 = (N_1 = \{S_1, A, B\}, T_1 = \{x, \bar{x}, a, \bar{a}, b, \bar{b}\}, \{S_1\}, P_1 = \{S_1 \longrightarrow xA^*B^+ \bar{x}, A \longrightarrow a\bar{a}, B \longrightarrow b\bar{b}\})$$

$$G = (N_1 \cup \{S\}, T_1, S, P_1 \cup \{S \longrightarrow sS_1^* \bar{s}\})$$

$G$  erzeugt die iterierte Sprache von  $G_1$ . Sie ist wieder eine balancierte Grammatik. Sie erzeugt z.B. das Wort

$$sxa\bar{a}a\bar{a}b\bar{b}\bar{x}xa\bar{a}a\bar{a}b\bar{b}\bar{x}s.$$

$\square$

### Erweiterte Iteration einer XML Sprache

$$G_1 = (N_1 = \{S_1, A\}, T_1 = \{x, \bar{x}, a, \bar{a}\}, S_1, P_1 = \{S_1 \longrightarrow xA^? \bar{x}, A \longrightarrow a\bar{a}\})$$

$$G = (N_1 \cup \{S\}, T_1, S, P_1 \cup \{S \longrightarrow sX^* \bar{s}\})$$

$G$  erzeugt die iterierte Sprache von  $G_1$ . Sie ist wieder eine XML Grammatik, da sie keine Variablen enthält, die konkurrieren. Sie erzeugt z.B. das Wort

$$sxa\bar{a}\bar{x}xa\bar{a}\bar{x}s.$$

$\square$

## 5.2 Entscheidbarkeitsprobleme

In dem folgenden Unterkapitel soll beantwortet werden, ob es entscheidbar ist, ob eine gegebene balancierte Grammatik eine Grammatik eines Untertyps der balancierten Sprachen ist. In Kapitel 4 wurde diskutiert, dass die Parsergeneratoren unter Umständen sehr lange Zeit benötigen, um einen Parser für balancierte Sprachen zu konstruieren. Deshalb könnte in der Praxis ein Algorithmus vorgeschaltet werden, der entscheidet, ob eine gegebene balancierte Grammatik eine Grammatik eines Untertyps ist.

Des Weiteren wird überprüft, ob entschieden werden kann, dass eine balancierte Grammatik mehrdeutig ist und ob eine balancierte Sprache eine Teilmenge einer anderen balancierten Sprache ist.

### 5.2.1 Balancierte Grammatiken und ihre Untertypen

Eine XML Grammatik zeichnet aus, dass jedes Nichtterminalzeichen  $X \in N$ , zu dem eine Produktion

$$X \longrightarrow xm\bar{x}$$

mit  $x \in A$  und  $m \in \text{Reg}(N)$  existiert, eindeutig durch das Terminalzeichen  $x$  identifizierbar ist. D.h. es existiert keine Produktion

$$Y \longrightarrow xm'\bar{x}$$

mit  $X \neq Y, Y \in N$  und  $m' \in \text{Reg}(N)$ .

Um zu entscheiden, ob eine gegebene balancierte Grammatik eine XML Grammatik ist, muss genau diese Eigenschaft überprüft werden.

Die Idee des Beweises ist, dass alle Produktionen der gegebenen balancierten Grammatik darauf untersucht werden, ob mindestens zwei verschiedene Nichtterminalzeichen durch dasselbe Terminalzeichen identifiziert werden.

**Definition 5.2.30** Sei  $G = (N, T, S, P), T = A \cup \bar{A}$  eine balancierte Grammatik. Dann liefert die Abbildung

$$\psi : P \longrightarrow N \times A$$

zu jeder Produktion das erzeugende Nichtterminal- und das identifizierende Terminalzeichen als Paar. Die Abbildung ist einfach zu implementieren. Sie liefert für eine Produktion der Grammatik  $G$  das Nichtterminalzeichen, was auf der linken Seite der Produktion steht und das erste Terminalzeichen (Starttag) auf der rechten Seite der Produktion.

Die Abbildung

$$\kappa : P \longrightarrow \text{Reg}(N)$$

ordnet einer Produktion ihr Inhaltsmodell zu (pro Terminalzeichen darf höchstens ein Inhaltsmodell existieren). Die Abbildung  $\kappa$  ist implementierbar. Sie entfernt von der rechten Seite einer Produktion lediglich den Start- und den Endtag  $(x, \bar{x})$  und gibt den Rest zurück.

Sei  $p \in P$  eine Produktion. Da  $G$  eine balancierte Grammatik ist, hat  $p$  die Form  $X \rightarrow xm\bar{x}$  mit  $X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$ . Dann ist  $\psi(p) = (X, x)$  und  $\kappa(p) = m$ . Hierbei ist zu beachten, dass  $m = \lambda$  als Inhaltsmodell zugelassen ist.

**Satz 5.2.60** Es ist entscheidbar, ob eine balancierte Grammatik eine XML Grammatik ist.

*Beweis.* Sei  $G = (N, T, S, P)$  eine gegebene balancierte Grammatik mit Produktionen der Form

$$X \rightarrow xm\bar{x}$$

mit  $X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$ .

Sei

$$K_x = \{X \in N \mid \psi(p) = (X, x), p \in P\}$$

die Menge, die jedem  $x \in A$  die Nichtterminale zuordnet, die durch  $x$  identifiziert werden. Die Menge  $K_x$  lässt sich konstruieren, da die Abbildung  $\psi$ , wie weiter oben beschrieben, implementierbar ist und die Mengen  $A, N$  und  $P$  endlich sind. Nun bleibt zu überprüfen, ob jede Menge  $K_x$  entweder kein, genau ein oder mehrere Elemente enthält. Falls alle Mengen kein oder genau ein Element enthalten, so ist die gegebene balancierte Grammatik eine XML Grammatik, andernfalls nicht.  $\square$

Eine einfache balancierte Grammatik zeichnet aus, dass im Inhaltsmodell einer Produktion nur nicht konkurrierende Nichtterminale auftreten. Um zu entscheiden, ob eine gegebene balancierte Grammatik eine einfache balancierte Grammatik ist, muss diese Eigenschaft nachgewiesen werden. Es wird überprüft, ob die Inhaltsmodelle aller Produktionen der gegebenen balancierten Grammatik jeweils zwei oder mehr konkurrierende Nichtterminale enthalten.

**Satz 5.2.61** Es ist entscheidbar, ob eine balancierte Grammatik eine einfache balancierte Grammatik ist.

*Beweis.* Sei  $G = (N, T, S, P), T = A \cup \bar{A}$  eine gegebene balancierte Grammatik. Die Produktionen haben die Form

$$X \rightarrow xm\bar{x}$$

mit  $X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$ . Sei

$$K_Y = \{X \in N \setminus \{Y\} \mid \psi(p) = (X, x), \psi(p') = (Y, x), p, p' \in P, x \in A\}$$

die Menge der mit  $Y \in N$  konkurrierenden Nichtterminale. Die Menge  $K_Y$  ist ebenfalls konstruierbar, da die Abbildung  $\psi$  implementierbar ist und die Mengen  $A, N$  und  $P$  endlich sind.

Sei

$$\pi_i : \text{Reg}(N) \longrightarrow N$$

die Projektion auf die  $i$ -te Komponente des übergebenen Ausdrucks. [Die erste Stelle hat hier den Index 1.] Bei der Indizierung und der Länge eines Ausdrucks werden die Operatoren  $+, *, |$  und  $?$  ignoriert. Die Länge  $n_m$  eines Ausdrucks  $m$  ist also die Anzahl der in  $m$  auftretenden Zeichen  $X \in N$ .

Sei

$$K_{Y,p} = \{X \in N \mid \pi_i(\kappa(p)) = X, \pi_j(\kappa(p)) = Y, X \neq Y, i, j \in \{1, \dots, n_{\kappa(p)}\}\}$$

die Menge aller Nichtterminale  $X$  des Inhaltsmodells der Produktion  $p$ , in dem sowohl  $X$  als auch das fest gewählte  $Y$  vorkommen. Die Abbildung  $\kappa$  und die Projektion  $\pi$  sind implementierbar. Die entstehende Menge  $K_{Y,p}$  ist endlich, da  $\kappa$  einen regulären Ausdruck und  $\pi$  eines der Zeichen des regulären Ausdrucks zurückgibt.

Bilde nun den Durchschnitt

$$K_Y \cap K_{Y,p}$$

für alle  $Y \in N$  und für alle  $p \in P$ . Falls alle Durchschnitte leer sind, so ist die gegebene balancierte Grammatik eine einfache balancierte Grammatik, andernfalls nicht.  $\square$

Eine konkurrenzverhindernde balancierte Grammatik zeichnet aus, dass für zwei konkurrierende Nichtterminale  $X$  und  $Y$ , die in demselben Inhaltsmodell einer Produktion auftreten, für alle Kombinationen  $U, V, W \in N^*$  gilt: Aus einem Inhaltsmodell  $m$  folgt

- nicht sowohl  $UXV \in L(m)$
- als auch  $UYW \in L(m)$ .

Um zu entscheiden, ob eine gegebene balancierte Grammatik eine konkurrenzverhindernde balancierte Grammatik ist, muss genau diese Eigenschaft überprüft werden.

Die Idee des Beweises ist, dass für alle Inhaltsmodelle zunächst geprüft wird, ob zwei konkurrierende Nichtterminale auftreten. Falls dies der Fall ist, wird noch überprüft, ob bei gleicher Anfangssequenz im Inhaltsmodell durch die konkurrierenden Nichtterminale unterschiedliche Worte abgeleitet werden können.

**Satz 5.2.62** Es ist entscheidbar, ob eine balancierte Grammatik eine konkurrenzverhindernde balancierte Grammatik ist.

*Beweis.* Sei  $G = (N, T, S, P), T = A \cup \bar{A}$  eine gegebene balancierte Grammatik. Produktionen haben die Form

$$X \longrightarrow xm\bar{x}$$

mit  $X \in N, x \in A, \bar{x} \in \bar{A}, m \in \text{Reg}(N)$ .

Sei

$$K_Y = \{X \in N \setminus \{Y\} \mid \psi(p) = (X, x), \psi(p') = (Y, x), p, p' \in P, x \in A\}$$

die Menge der mit  $Y \in N$  konkurrierenden Nichtterminale.

Es sei  $W_{p,Y}$  die Menge aller Teilworte aus  $L(r), r = \kappa(p), r \in \text{Reg}(N)$ , die mit  $Y \in N$  enden.

$$W_{p,Y} = \{u \mid u = \alpha \wedge \alpha Y \beta \in L(r) \wedge r = \kappa(p), \alpha, \beta \in N^*\}$$

Es werden nun die Durchschnitte gebildet

$$W_{p,Y} \cap W_{p,Z}$$

für alle  $Y \in N$  mit  $Z \in K_Y$  und für alle Produktionen  $p \in P$ .

Sind alle Durchschnitte leer (enthalten sie also nicht einmal das leere Wort  $\lambda$ ), so ist die gegebene balancierte Grammatik eine konkurrenzverhindernde balancierte Grammatik, andernfalls nicht.  $\square$

### 5.2.2 Inklusion

**Satz 5.2.63** Seien  $L_1$  und  $L_2$  zwei balancierte Sprachen. Es ist entscheidbar, ob  $L_1 \subseteq L_2$ .

*Beweis.* Seien  $G_1$  und  $G_2$  zwei balancierte Grammatiken und  $L_1$  bzw.  $L_2$  die von ihnen erzeugten Sprachen. Die balancierten Sprachen sind deterministisch kontextfrei und somit ist es entscheidbar, ob zwei balancierte Sprachen äquivalent sind. Zusätzlich sind die balancierten Sprachen unter Durchschnitt abgeschlossen.

Bilde nun den Schnitt  $L$  der beiden erzeugten Sprachen

$$L = L_1 \cap L_2.$$

Da die Familie der balancierten Sprachen unter Schnittbildung effektiv abgeschlossen ist, existiert eine balancierte Grammatik  $G$ , die  $L$  erzeugt. Ist  $L$  äquivalent zu  $L_1$ , so gilt  $L_1 \subseteq L_2$ . Ist  $L$  äquivalent zu  $L_2$ , so ist  $L_2 \subseteq L_1$ . Wenn  $L$  zu keiner der beiden Sprachen äquivalent ist, so sind die Sprachen  $L_1$  und  $L_2$  nicht ineinander enthalten.  $\square$

### 5.2.3 Mehrdeutigkeit

**Definition 5.2.31** Eine balancierte Grammatik  $G$  ist mehrdeutig, wenn es für die indizierte balancierte Grammatik  $G'$  ein Wort in  $L(G)$  gibt, wofür die Grammatik  $G'$  mehr als einen Ableitungsbaum besitzt.

**Satz 5.2.64** Mit balancierten Grammatiken können mehrdeutige Grammatiken beschrieben werden.

**Mehrdeutigkeiten** Es gibt verschiedene Ursachen der Mehrdeutigkeit.

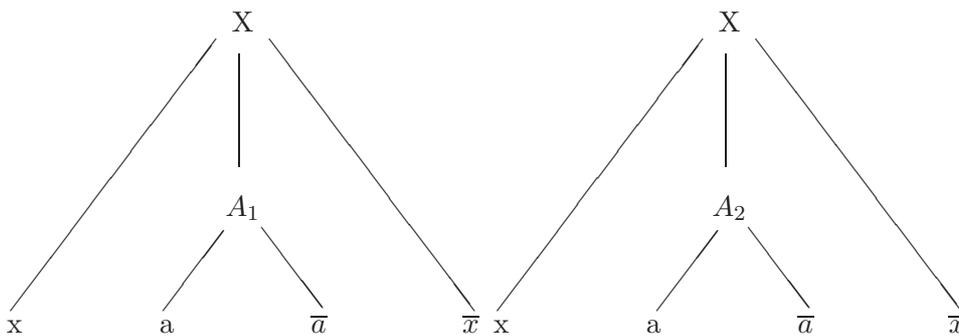
**Beispiel 5.2.27** Zum einen kann die Mehrdeutigkeit durch ein Inhaltsmodell, welches aus einem mehrdeutigen regulären Ausdruck besteht, hervorgerufen werden. Es sei  $G = (N, T, \{X\}, P)$  eine balancierte Grammatik mit folgenden Produktionen:

$$\begin{aligned} X &\longrightarrow x(A \mid B)^* A^* \bar{x} \\ A &\longrightarrow a\bar{a} \\ B &\longrightarrow b\bar{b}. \end{aligned}$$

Die Produktion von  $X$  wird nun indiziert.

$$X \longrightarrow x_1(A_1 \mid B_1)^* A_2^* \bar{x}_1$$

Für das Wort  $xa\bar{a}\bar{x}$  lassen sich zwei verschiedene Ableitungen erzeugen:  $x_1a_1\bar{a}_1\bar{x}_1, x_1a_2\bar{a}_2\bar{x}_1$ .



**Beispiel 5.2.28** Die zweite Form von Mehrdeutigkeiten wird durch eine spezielle Eigenschaft der balancierten Grammatiken hervorgerufen. Balancierte Grammatiken erlauben, im Gegensatz zu ihren Untertypen, jegliche Art von Produktionen der Form

$$X \longrightarrow am\bar{a}.$$

Sie erlauben somit auch eine beliebige Konkatenation von konkurrierenden Nichtterminalen.

Durch den Satz von [BB02a], dass aus einer balancierten Sprache eindeutig eine balancierte Grammatik erzeugt werden kann, wird die Sprache

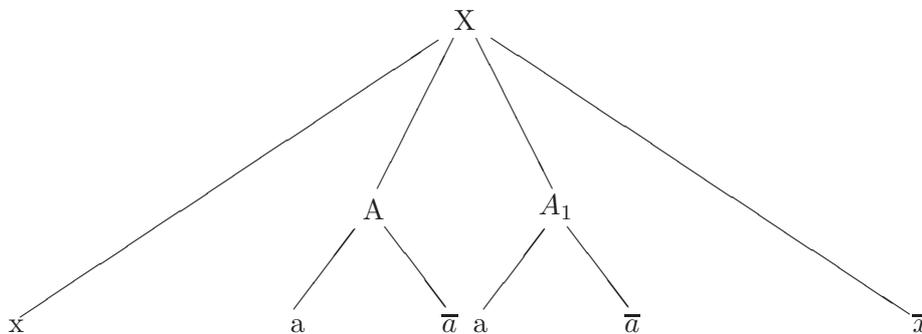
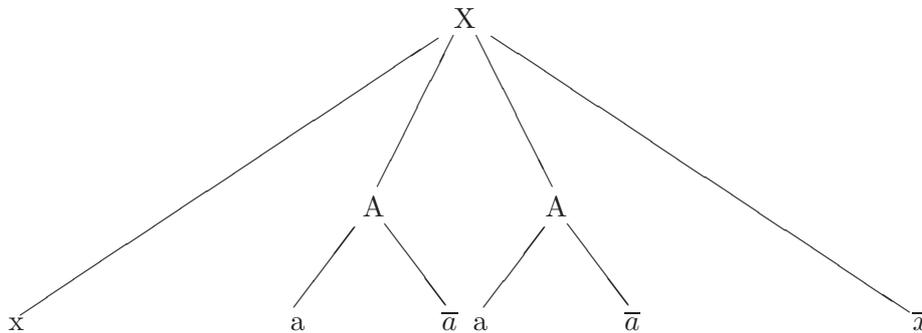
$$(x(a(c\bar{c})?a)^*(a(d\bar{d})?a)^*\bar{x})$$

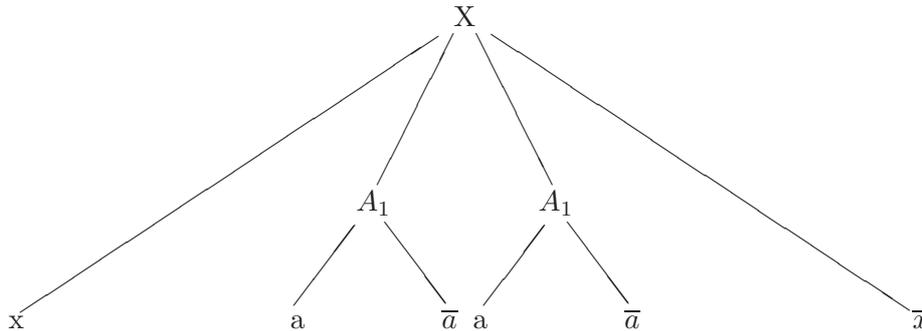
durch die balancierte Grammatik  $G = (N, T, \{X\}, P)$

$$P = \{ \begin{array}{l} X \longrightarrow xA^*A_1^*\bar{x} \\ A \longrightarrow aC?a \\ A_1 \longrightarrow aD?a \\ C \longrightarrow c\bar{c} \\ D \longrightarrow d\bar{d} \end{array} \}$$

beschrieben.

Für ein Wort  $w = xa\bar{a}a\bar{x}$  gibt es folgende drei Ableitungsbäume:





Die Grammatik ist mehrdeutig.

Mehrdeutigkeiten führen zu Problemen beim Parsen, da für ein gegebenes Wort nicht klar ist, welcher Parsebaum aufgebaut werden muss. Das kann zu hohen Laufzeiten führen.

Eine wichtige Eigenschaft einer Sprache in Bezug auf Mehrdeutigkeit ist es, ob eine Sprache inhärent mehrdeutig ist. Das bedeutet, dass die Sprache nur durch mehrdeutige Grammatiken erzeugt werden kann. Für die Einschränkung, dass die Grammatik eine balancierte Grammatik sein muss, zeigte das letzte Beispiel, dass es Sprachen gibt, für die nur mehrdeutige balancierte Grammatiken gefunden werden können.

**Satz 5.2.65** Es ist entscheidbar, ob eine balancierte Grammatik eindeutig ist.

*Beweis.* Sei  $G = (N, T, S, P)T = A \cup \bar{A}$  eine gegebene balancierte Grammatik. Elemente aus  $P$  haben die Form

$$X \longrightarrow xm\bar{x}$$

mit  $X \in N, x \in A, \bar{x} \in \bar{A}$  und  $m \in \text{Reg}(N)$ .

Um zu entscheiden, ob eine balancierte Grammatik eindeutig ist, muss geprüft werden, ob zu jedem beliebigen Wort aus der von ihr erzeugten Sprache nur genau ein Ableitungsbaum existiert.

Betrachte hierfür zunächst die Inhaltsmodelle aller Produktionen. Sie sind reguläre Ausdrücke. Für die regulären Ausdrücke ist es entscheidbar, ob sie eindeutige oder mehrdeutige reguläre Ausdrücke sind [BEGO71]. Falls ein Inhaltsmodell mehrdeutig ist, so folgt direkt, dass die Grammatik mehrdeutig ist, wie aus dem Beispiel 5.2.27 trivial zu folgern ist.

Seien also die Inhaltsmodelle alle eindeutige reguläre Ausdrücke.

Betrachte nun die Sprachen, die von jedem einzelnen Nichtterminalzeichen erzeugt werden. Sei  $X \in N$ . Die von  $X$  erzeugte Sprache ist

$$L_X = \{w \in (A \cup \bar{A})^* \mid X \Rightarrow^* w\}.$$

Bilde die Schnitte

$$L_X \cap L_Y$$

für alle  $X, Y \in N, X \neq Y$ . Wenn für zwei Nichtterminale der Schnitt leer ist, können sie nicht dasselbe Wort erzeugen. Also führen sie in Kombination miteinander im Inhaltsmodell einer Produktion nicht zur Mehrdeutigkeit. Ist der Schnitt nicht leer, so kann aus diesen zwei Nichtterminalen dasselbe Wort erzeugt werden.

Sei

$$N_X = \{Y \in N \mid L_X \cap L_Y \neq \emptyset, X \neq Y\}$$

die Menge der Nichtterminale, die in Kombination mit  $X$  zu Mehrdeutigkeit führen können. Es muss also noch überprüft werden, ob zu zwei Nichtterminalen  $X \in N, Y \in N_X$  und einem Inhaltsmodell  $m$  Sequenzen  $U, U', V, V' \in \text{Reg}(N)$  existieren, mit Ableitungen

$$m \Rightarrow^* UXV \Rightarrow^* w$$

und

$$m \Rightarrow^* U'YV' \Rightarrow^* w$$

für ein  $w \in (A \cup \bar{A})^*$ . Da  $X$  und  $Y$  dasselbe Wort ableiten können, muss

$$U \longrightarrow u, U' \longrightarrow u, V \longrightarrow v, V' \longrightarrow v$$

gelten für  $u, v \in (A \cup \bar{A})^*$ . Dann gibt es für das Wort  $w$  mehrere Ableitungen und die Grammatik ist mehrdeutig.

Sei

$$A_{m,X} = \{U \in \text{Reg}(N) \mid m = UXW, W \in \text{Reg}(N)\}$$

die Menge aller Nichtterminalkombinationen, die im Inhaltsmodell  $m$  vor einem gegebenen  $X \in N$  auftreten können.

Sei

$$E_{m,X} = \{W \in \text{Reg}(N) \mid m = UXW, U \in \text{Reg}(N)\}$$

die Menge aller Nichtterminalkombinationen, die im Inhaltsmodell  $m$  nach einem  $X \in N$  auftreten können.

Sei  $G_U$  die balancierte Grammatik zu  $U \in \text{Reg}(N)$  mit  $P_U = \{X \longrightarrow xU\bar{x}\}$ .

Sei  $U \in A_{m,X}$  und  $U' \in A_{m,Y}$  und  $G_U$  und  $G_{U'}$  die balancierten Grammatiken wie oben beschrieben. Bilde nun die Schnitte

$$L(G_U) \cap L(G_{U'}).$$

Sind alle Durchschnitte leer, so existieren keine Nichtterminalkombinationen, die vor  $X$  bzw.  $Y$  in den Inhaltsmodellen auftreten und gleiche Worte erzeugen können. Damit kann keine Mehrdeutigkeit auftreten.

Sei ein Schnitt nicht leer, es gelte also

$$L(G_U) \cap L(G_{U'}) \neq \emptyset$$

für ein  $U \in A_{m,X}$  und ein  $U' \in A_{m,Y}$ .  
Betrachte nun die Durchschnitte

$$L(G_W) \cap L(G_{W'})$$

für alle  $W \in E_{m,X}$  und  $W' \in E_{m,Y}$  und deren Grammatiken  $G_W$  bzw.  $G_{W'}$ . Sind diese Schnitte leer, so kann keine Mehrdeutigkeit auftreten. Ist mindestens einer der Schnitte nicht leer, so ist die Grammatik mehrdeutig, da dann für ein Wort aus  $L(G)$  mehrere Ableitungen existieren.  $\square$

Auch in den Entscheidbarkeitsfragen zeigt sich, dass die balancierten Sprachen sehr nutzbar für die Praxis sind. Es kann entschieden werden, ob eine balancierte Grammatik eine Grammatik eines Untertyps ist, ob eine balancierte Sprache eine Teilmenge einer anderen Sprache ist und ob eine balancierte Grammatik mehrdeutig ist.

Konkrete Anwendungen könnten sein, dass vor der Generierung eines Parsers für ein Beschreibungsmodell entschieden wird, von welchem Typ die Grammatik ist und danach entschieden wird, welche Art von Parser konstruiert wird.

### 5.2.4 Beispiele zu Mehrdeutigkeit

#### Beispiel 1:

Das Inhaltsmodell ist ein mehrdeutiger regulärer Ausdruck.  
Seien

$$X \longrightarrow xA^*A^*\bar{x}, A \longrightarrow a\bar{a}$$

Produktionen einer balancierten Grammatik  $G = (N, T, \{X\}, P)$ .

Das Wort  $xa\bar{a}a\bar{a}\bar{x}$  kann durch diese Produktionen erzeugt werden. Die Kombination  $a\bar{a}a\bar{a}$  kann sowohl aus  $A^*$  als auch aus  $A^*A^*$  abgeleitet werden. Sei  $xA_1^*A_2^*\bar{x}$  der indizierte Ausdruck der rechten Seite der Startproduktion. Somit existieren die Ableitungen  $X \longrightarrow xA_1A_1\bar{x} \longrightarrow xa_1\bar{a}_1a_1\bar{a}_1\bar{x}$  und  $X \longrightarrow xA_1A_2\bar{x} \longrightarrow xa_1\bar{a}_1a_2\bar{a}_2\bar{x}$  für das Wort  $xa\bar{a}a\bar{a}\bar{x}$ .  $\square$

#### Beispiel 2:

Die Inhaltsmodelle sind eindeutige reguläre Ausdrücke  
Seien

$$\begin{aligned} S &\longrightarrow sXY\bar{s}, X \longrightarrow xA^*B^*\bar{x}, Y \longrightarrow xA^*D^*\bar{x}, \\ A &\longrightarrow a\bar{a}, B \longrightarrow b\bar{b}, D \longrightarrow d\bar{d} \end{aligned}$$

Produktionen einer balancierten Grammatik  $G = (N, T, \{S\}, P)$ .

Das Wort  $s\bar{x}s$  kann durch diese Produktionen erzeugt werden. Es existieren die Ableitungen  $X \rightarrow \bar{x}$  und  $Y \rightarrow \bar{x}$  für das Teilwort  $\bar{x}$ .

□

## Zusammenfassung der Entscheidungsfragen

Entscheidbarkeit		reguläre Sprachen	balancierte Sprachen	det. kon. Spr.	kontextf. Spr.
	Wortproblem	+	+	+	+
	Äquivalenz	+	+	+	-
	Leerheit	+	+	+	+
	Unendlichkeit	+	+	+	+
	Inklusion	+	+	-	-
	Mehrdeutigkeit	+	+	+	-
	XML S.		+		
	einfache bal. S.		+		
	konkurrenzverhindernde bal. S.		+		
RegDyck G.		+			

**det. kon. Spr.** - deterministisch kontextfreie Sprachen, **kontextf. Spr.** - kontextfreie Sprachen, **einfache bal. S.** - einfache balancierte Sprachen, **konkurrenzverhindernde bal. S.** - konkurrenzverhindernde balancierte Sprachen

# Literaturverzeichnis

- [ASU99a] AHO, A. ; SETHI, R. ; ULLMANN, J.: *Compilerbau Teil 1*. 2. Auflage. München : Oldenbourg Wissenschaftsverlag, 1999
- [ASU99b] AHO, A. ; SETHI, R. ; ULLMANN, J.: *Compilerbau Teil 2*. 2. Auflage. München : Oldenbourg Wissenschaftsverlag, 1999
- [BB00] BERSTEL, J. ; BOASSON, L.: XML grammars. In: *Mathematical Foundations of Computer Science, LNCS* Bd. 1893, Springer, 2000, S. 182–191
- [BB02a] BERSTEL, J. ; BOASSON, L.: Balanced Grammars and Their Languages. In: *Formal and Natural Computing, LNCS* Bd. 2300, Springer, 2002, S. 3–25
- [BB02b] BERSTEL, J. ; BOASSON, L.: Formal properties of XML grammars and languages. *Acta Informatica* 38 (2002), S. 649–671
- [BEGO71] BOOK, R. ; EVEN, S. ; GREIBACH, S. ; OTT, G.: Ambiguity in Graphs and Expressions. *IEEE Transactions on Computers* 66 (1971), S. 451–472
- [BK93a] BRÜGGEMANN-KLEIN, A.: *Formal models in document processing*, Faculty of Mathematics at the University of Freiburg, Dissertation, 1993
- [BK93b] BRÜGGEMANN-KLEIN, A.: Unambiguity of extended regular expressions in SGML document grammars / Technical University of Munich. 1993. – Forschungsbericht
- [BKW92] BRÜGGEMANN-KLEIN, A. ; WOOD, D.: Deterministic regular languages. In: *9th Annual Symposium on Theoretical Aspects of Computer Science, LNCS* Bd. 577, Springer, 1992, S. 173–184
- [BKW98] BRÜGGEMANN-KLEIN, A. ; WOOD, D.: One-Unambiguous Regular Languages. *Information and Computation* 142 (1998), S. 182–206

- [Bos05] <http://www.javaworld.com/javaworld/javaone00/j1-00-bosak.html>
- [CDGJ97] COMON, H. ; DAUCHET, M. ; GILLERON, R. ; JACQUEMARD, F.: *Tree Automata Techniques and Applications*. Forschungsbericht, 1997
- [Claa] <http://www.thaiopensource.com/trex>
- [Clab] <http://www.oasis.open.org/committees/relax-ng/tutorial.html>
- [DFS99] DEUTSCH, A. ; FERNANDEZ, M. ; SUCIU, D.: Storing semistructured data with STORED. In: *Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data*, 1999, S. 431–442
- [DOM07] <http://www.w3.org/DOM/>
- [DTD07] <http://www.w3.org/TR/html401/sgml/dtd.html>
- [FK99] FLORESCU, D. ; KOSSMANN, D.: Storing and Querying XML Data using RDBMS. In: *IEEE Data Bulletin, VLDB Journal* 14 (1999), Nr. 1, S. 30–49
- [GH67] GINSBURG, S. ; HARRISON, M.: Bracketed context-free Languages. *Journal of Computer and System Sciences* 1 (1967), S. 1–23
- [GP00] GOLDFARB, C. ; PRESCOD, P.: *The XML Handbook*. 2. Auflage. London : Prentice Hall, 2000
- [Har99] HAROLD, E.: *XML Bible*. IDG Books Worldwide, 1999
- [HU79] HOPCROFT, J. ; ULLMAN, J.: *Introduction to Automata Theory, Language, and Computation*. München : Addison-Wesley, 1979
- [HW07] HAN, Yo-Sub ; WOOD, D.: Generalizations of One-deterministic Regular Languages. In: *Conference on Language and Automata Theory and Applications (LATA '07), Universitat Rovira I Virgili*, 2007, S. 273–285
- [Knu67] KNUTH, D.: A characterization of parenthesis languages. *Information Control* 11 (1967), S. 269–289
- [LC00] LEE, D. ; CHU, W. W.: Constraints-Preserving Transformation from XML Document Type Definition to Relational Schema. In: *International Conference on Conceptual Modeling / the Entity Relationship Approach, LNCS Bd. 1920/2000*, Springer, 2000, S. 323–338

- [LMCC02] LEE, D. ; MANI, M. ; CHIU, F. ; CHU, W. W.: NeT and CoT: Translating Relational Schemas to XML Schemas using Semantic Constraints. In: *ACM International Conference on Information and Knowledge Management*, 2002
- [LMM00] LEE, D. ; MANI, M. ; MURATA, M.: “Reasoning about XML Schema Languages using Formal Language Theory” / IBM Almaden Research Center. 2000. – Forschungsbericht
- [McN67] MCNAUGHTON, R.: Paranthesis grammars. *Journal of the ACM* 14 (1967), S. 490–500
- [Meg] <http://www.megginson.com/SAX/index.html>
- [ML02] MANI, M. ; LEE, D.: XML to Relational Conversion using Theory of Regular Tree Grammars. In: *Proceedings of the VLDB 2002 Workshop EEXTT and CAiSE 2002 Workshop DTWeb on Efficiency and Effectiveness of XML Tools and Techniques and Data Integration over the Web-Revised Papers, LNCS, Springer* Bd. 2590, 2002, S. 81–93
- [MLM01a] MANI, M. ; LEE, D. ; MUNTZ, R. R.: Semantic Data Modeling Using XML Schemas. In: *LNCS, Springer* Bd. 2224, 2001, S. 149–163
- [MLM01b] MURATA, M. ; LEE, D. ; MANI, M.: Taxonomy of XML Schema Languages using Formal Language Theory / IBM Almaden Research Center. 2001. – Forschungsbericht
- [SAX07] <http://www.saxproject.org/>
- [See03] SEEMANN, M.: *Native XML Datenbanken im Einsatz*. Paderborn : Software and Support Verlag, 2003
- [SSB<sup>+</sup>01] SHANMUGASUNDARAM, J. ; SHEKITA, E. ; BARR, R. ; CAREY, M. ; LINDSAY, B. ; PIRAHESH, H. ; REINWALD, B.: Efficiently publishing relational data as XML documents. *VLDB Journal: Very Large Data Bases* 10 (2001), S. 133–154
- [SSS88] SIPPU, S. ; SOISALON-SOININEN, E.: *Parsing Theory*. Bd. 1. Berlin : Springer, 1988
- [ST04] SAPOUNAKIS, A. ; TSIKOURAS, P.: On k-colored Motzkin words. *Journal of Integer Sequences* 4 (2004)
- [STZ<sup>+</sup>99] SHANMUGASUNDARAM, J. ; TUFTE, K. ; ZHANG, C. ; DEWITT, D. J. ; NAUGHTON, J. F.: Relational Databases for Querying

---

XML Documents: Limitations and Opportunities. In: *Proceedings of 25th International Conference on Very Large Data Bases, VLDB Journal*, 1999, S. 302-314

[WWW] <http://www.w3.org/TR/REC-xml>

[XML05] <http://www.w3.org/TR/2004/REC-xml-20040204/>

[XML07] <http://www.w3.org/XML/Schema>