

Syntaxanalyse: LR-Parser

BC George (FH Bielefeld)

Unless otherwise noted, this work is licensed under CC BY-SA 4.0.

Wiederholung

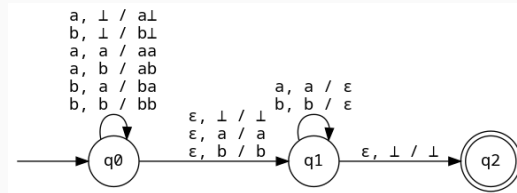


Abbildung 1: Ein PDA für $L = \{ww^R \mid w \in \{a, b\}^*\}$

- Baufeldbau von oben nach unten
- die Grammatik muss reduziert sein
- recursive-descent parser
- *First-* und *Follow*-Mengen bestimmen Wahl der Ableitungen
- tabellengesteuert
- nicht mehr rekursiv, sondern mit PDA

Motivation

LL ist nicht alles

Die Menge der *LL*-Sprachen ist eine echte Teilmenge der deterministisch kontextfreien Sprachen. Wir brauchen ein Verfahren, mit dem man alle deterministisch kontextfreien Sprachen parsen kann.

Bottom-Up-Analyse

Von unten nach oben

Bei LL -Sprachen muss man nach den ersten k Eingabezeichen entscheiden, welche Ableitung ganz oben im Baum als erste durchgeführt wird, also eine, die im Syntaxbaum ganz weit weg ist von den Terminalen, die die Entscheidung bestimmen. Es gibt deterministisch kontextfreie Sprachen, die nicht $LL(k)$ sind für irgendein k .

Bei der Bottom-Up-Analyse geht man den umgekehrten Weg. Der Parse Tree wird von unten nach oben aufgebaut, die Entscheidung, welche Produktion angewandt wird, erfolgt "näher" am Terminal. Mit Hilfe der Produktionen und der Vorschautoken werden die Ableitungen "rückwärts" angewandt und "Reduktionen" genannt.

Fehlermeldungen können näher am Programmtext erfolgen.

Baumaufbau von unten

zu parsen $id + id * id$

$$E' \rightarrow E$$

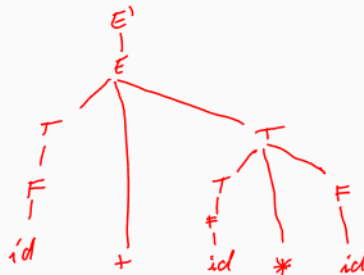
$$E \rightarrow E + T \mid F$$

$$F \rightarrow (E) \mid id$$

$$T \rightarrow T * F \mid F$$

Scanner $\rightarrow id + id * id$

↓
auf dem Stack packen?



Kann ein Stack helfen?

das ~~selbe~~ wie oben

Probleme damit?

Wir können nicht entscheiden, wann notisiert und wann weiter
eingesehen werden muss.



(„slippen“)

Konfliktfälle

Mehrdeutigkeiten = Konflikte beim Parsen

Es gibt Grammatiken, bei denen nicht aus dem Inhalt des Stacks und dem Eingabezeichen entschieden werden kann, wie fortgefahren wird, auch nicht, wenn man, wie auch schon im Fall LL , eine feste Zahl k von Vorschautoken berücksichtigt. Diese Grammatiken können mehrdeutig sein.

Folgen von falschen Entscheidungen:

- falscher Baum, falsche Bäume
- kein Baum

Mögliche Konflikte

- Reduce-Reduce-Konflikt: Es sind zwei oder mehr verschiedene Reduktionen möglich
- Shift-Reduce-Konflikt: Es kann nicht entschieden werden, ob eine Reduktion oder ein Shift durchgeführt werden soll.

Shiften bedeutet, das nächste Eingabesymbol miteinbeziehen.

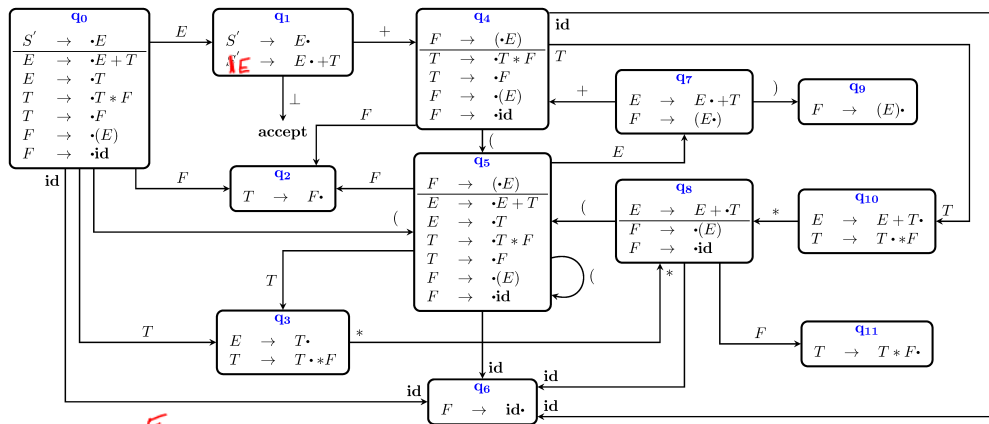
LR-Parsing

Da wollen wir hin

$S' \rightarrow E$ $F \rightarrow (E) | id$

$E \rightarrow E + T | T$

$T \rightarrow T * F | F$



Der Stack enthält Zustände, keine Terminals oder Nonterminals.

Der Top-of-Stack ist immer der aktuelle Zustand, am Anfang l_0 . Im Stack steht $l_0 \perp$.

Vorgehen: Im aktuellen Zustand nachschauen, ob das Eingabezeichen auf einem Pfeil steht.

- ja: Shiften, d. h. dem Pfeil folgen und den Zustand am Ende des Pfeils pushen. Dort weiter.
- nein: Reduzieren nach der Regel aus dem aktuellen Zustand mit dem Punkt hinten, d. h. so viele Zustände poppen, wie die Regel Elemente auf der rechten Seite hat. Der Zustand darunter wird aktuell, dem Pfeil mit dem zu reduzierenden Nonterminal der linken Seite der Regel folgen und pushen.

Am Schluss kann nur noch mit \perp akzeptiert werden.

Beispiel

Bsp: $id * id$ parsen

0) Zustand q_0

1) Scanner liefert $id \Rightarrow push$

2) id $*$

\Rightarrow reduzieren ($F \rightarrow id$)

$\underbrace{\quad}_{1 \text{ Element}}$

\Rightarrow $*$ poppen

wirte in q_0 , mit \neq weiter,

ergibt q_2

\vdots



Def.: Bei einer kontextfreien Grammatik G ist die *Rechtsableitung* von $\alpha \in (N \cup T)^*$ die Ableitung, die man erhält, wenn das am weitesten rechts stehende Nichtterminal in α abgeleitet wird. Man schreibt $\alpha \xRightarrow{*}_r \beta$.

Def.: Eine *Rechtssatzform* α einer Grammatik G ist ein Element aus $(N \cup T)^*$ mit $S \xRightarrow{*}_r \alpha$.

Def.: In dem Syntaxbaum von $S \xRightarrow{*}_r \alpha$ $A \Rightarrow_r \alpha \beta$ w einer kontextfreien Grammatik ist β ein *Handle* von der Produktion $A \rightarrow \beta$.

Bei der *LR*-Analyse eines Wortes w wird w von links nach rechts gelesen, dabei wird die Rechtsableitung von w in G von unten nach oben aufgebaut. Man spricht nicht nicht mehr von Ableitungen, sondern von Reduktionen.

Mehrdeutige Grammatiken können nicht *LR* sein.

- Vor der Konstruktion des Automaten wird die Grammatik um eine neues Nonterminal S' und die neue Produktion $S' \rightarrow S$ erweitert. S' ist dann Startsymbol. ($S' \rightarrow S \perp$)
- Es wird ein Automat erstellt (s.o.)
- Es wird eine Parse Table aus dem Automaten erstellt, die den Parse-Vorgang steuert, mit Aktionsteil und Sprungteil.



≙ Pfeilen im Automaten

Im Stack stehen nur Zustandsnummern, am Anfang die Nummer des Startzustandes (+ Bottomzeichen, oft auch \$). Es ist nicht nötig, ^{terminalen} Symbole zu stacken.

- Lesen des obersten Stackelements ergibt Zustand q
- Lesen des nächsten Eingabezeichens ergibt Zeichen a
- Nachschlagen der Reaktion auf (q, a) in der Parse Table
- Durchführung der Reaktion

Mögliche “Actions” ohne Berücksichtigung von Vorschautoken

- Shift: Schiebe logisch das nächste Eingabesymbol auf den Stack (in Wirklichkeit Zustandsnummern)
- Reduce: (Identifiziere ein Handle oben auf dem Stack und ersetze es durch das Nichtterminal der dazugehörigen Produktion.) Das ist gleichbedeutend mit: Entferne so viele Zustände vom Stack wie die rechte Seite der zu reduzierenden Regel Elemente hat, und schreibe den Zustand, der im *Sprung-Teil* Goto-Teil für (q, a) steht, auf den Stack.
- Accept: Beende das Parsen erfolgreich
- Reagiere auf einen Syntaxfehler

0 Vorschautoken = LR(0)-Parsing

LR-Parsing ohne Vorschautoken

Wichtig: Das Handle, d. h. die rechte Seite einer zu reduzierenden Regel, erscheint oben auf dem Stack, nie weiter unten.

Je nach Anwendungsfall müssen beim Reduzieren von Handles weitere Aktionen ausgeführt werden: z. B. Syntaxbäume aufgebaut, Werte in Tabellen geschrieben werden, usw. Nicht alle rechten Seiten von Produktionen, die oben auf dem Stack stehen, sind auch Handles, manchmal muss nur geshiftet werden.

Bsp: Steht bei der Beispielgrammatik von ~~Folie 4~~ oben auf dem Stack ein T mit dem nächsten Eingabezeichen $*$, darf T nicht zu E reduziert werden.

Lösung: Der Parser merkt sich, wo er steht in noch nicht komplett reduzierten Regeln. Dazu benutzt er sogenannte *Items* oder $LR(0)$ – *Items*, auch *dotted Items* oder (*kanonische*) $LR(0)$ – *Elemente*.

⇒ Zustände

Items

Def.: Ein *Item* einer Grammatik G ist eine Produktion von G mit einem Punkt auf der rechten Seite der Regel vor, zwischen oder nach den Elementen.

Bsp.:

Zu der Produktion $A \rightarrow BC$ gehören die Items:

$[A \rightarrow \cdot BC]$

$[A \rightarrow BC \cdot]$

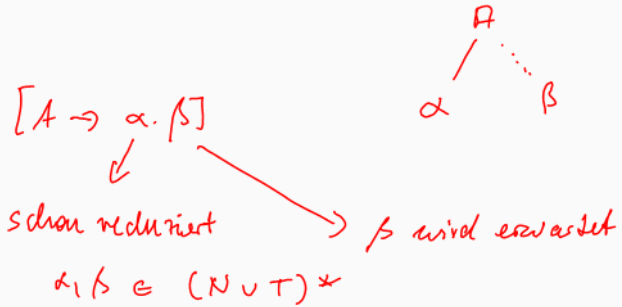
Item = Kern eines Zustands

$[A \rightarrow B \cdot C]$

$[A \rightarrow BC \cdot]$

Das zu $A \rightarrow \epsilon$ gehörende Item ist $[A \rightarrow \cdot]$

Was bedeuten die Items?



Berechnung der $Closure_0$ von einer Menge I von Items

$[A \rightarrow \cdot B]$ $B \rightarrow CD \mid E$
 $B \rightarrow \cdot CD$
 $B \rightarrow \cdot E$

I ist ein Item

1. füge I zu $CLOSURE_0(I)$ hinzu
2. gibt es ein Item $[A \rightarrow \alpha \cdot B\beta]$ aus $CLOSURE_0(I)$ und eine Produktion $(B \rightarrow \gamma)$, füge $[B \rightarrow \cdot \gamma]$ zu $CLOSURE_0(I)$ hinzu



eigentliche Zustände

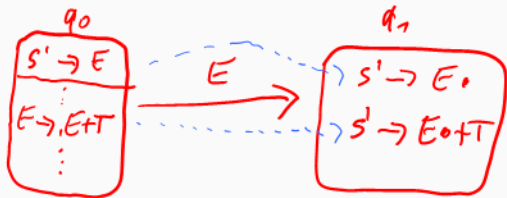
Berechnung der $GOTO_0$ -Sprungmarken

$\forall I, \forall X$

alter Zustand

$$GOTO_0(I, X) = CLOSURE_0(\{[A \rightarrow \alpha X \cdot \beta] \mid [A \rightarrow \alpha \cdot X\beta] \in I\})$$

für eine Itemmenge I und $X \in N \cup T, A \in N, \alpha, \beta \in (N \cup T)^*$.



← Das Item war vorher nicht in dem Zustand.

Konstruktion des $LR(0)$ - Automaten

1. Bilde die Hülle von $S' \rightarrow S$ und mache sie zum ersten Zustand.
2. Für jedes noch nicht betrachtete $\cdot X, X \in (N \cup T)$ in einem Zustand q des Automaten berechne $GOTO_0(q, X)$ und mache $GOTO_0(q, X)$ zu einem neuen Zustand r . Verbinde q mit einem Pfeil mit r und schreibe X an den Pfeil. Ist ein zu r identischer Zustand schon vorhanden, wird ~~q~~ ^{g} mit diesem verbunden und kein neuer erzeugt.

Konstruktion der Parse Table

closure

1. Erstelle eine leere Tabelle mit den Zuständen als Zeilenüberschriften. Für den Aktionstabellenteil überschreibe die Spalten mit den Terminalen, für den Sprungtabellenteil mit den Nonterminals.
2. Shift: Für jeden mit einem Terminal beschrifteten Pfeil aus einem Zustand erstelle in der Aktionstabelle die Aktion shift mit der Nummer des Zustands, auf den der Pfeil zeigt. Für Pfeile mit Nonterminals schreibe in die Sprungtabelle nur die Nummer des Folgezustands.
3. Schreibe beim Zustand $[S' \rightarrow S \cdot]$ ein *accept* bei dem Symbol \perp .
4. Für jedes Item mit $[A \rightarrow \beta \cdot]$ aus allen Zuständen schreibe für alle Terminals *reduce* und die Nummer der entsprechenden Grammatikregel in die Tabelle.

Und wenn in einer Zelle schon ein Eintrag ist?

Konflikt!

Die Beispielgrammatik G1

$$(0) S' \rightarrow S$$

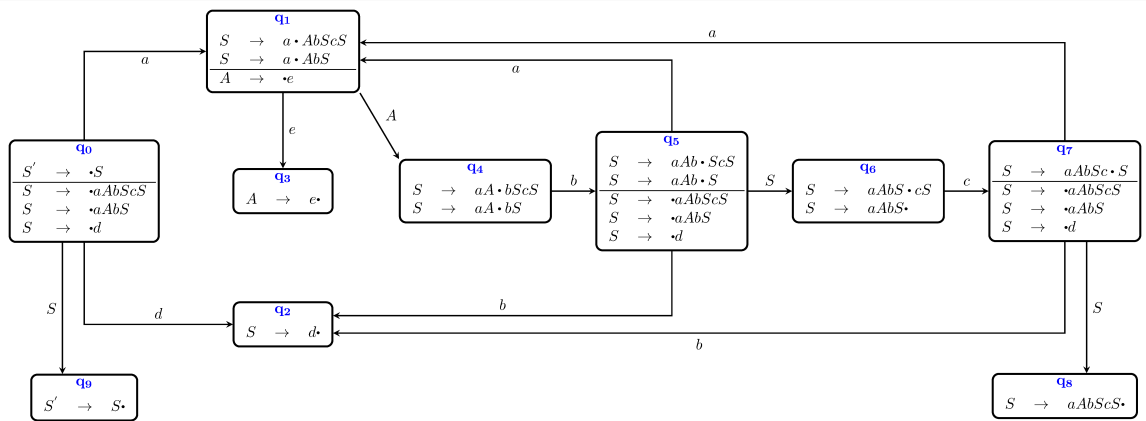
$$(1) S \rightarrow aAbScS$$

$$(2) S \rightarrow aAbS$$

$$(3) S \rightarrow d$$

$$(4) A \rightarrow e$$

Der LR(0)-Automat zu G1



Die LR(0)-Parsertabelle zu G1

Zustand	Aktion						Sprung	
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	\perp	<i>S</i>	<i>A</i>
q_0	s1			s2			q_9	
q_1					s3			q_4
q_2	r3	r3	r3	r3	r3	r3		
q_3	r4	r4	r4	r4	r4	r4		
q_4		s5						
q_5	s1	s2		s2			q_6	
q_6	r2	r2	s7,r2	r2	r2	r2		
q_7	s1	s2					q_8	
q_8	r1	r1	r1	r1	r1	r1		
q_9						acc		

s = shift

r = reduce

Regelnummer

Zustandsnummer

1 Vorschautoken = LR(1)-Parsing

Ist eine Grammatik nicht LR(0), kann sie vielleicht mit einem Vorschautoken geparkt werden. Hier gibt es drei Verfahren:

- SLR(1)-Parsing
- (kanonisches) LR(1)-Parsing
- LALR(1)-Parsing

SLR

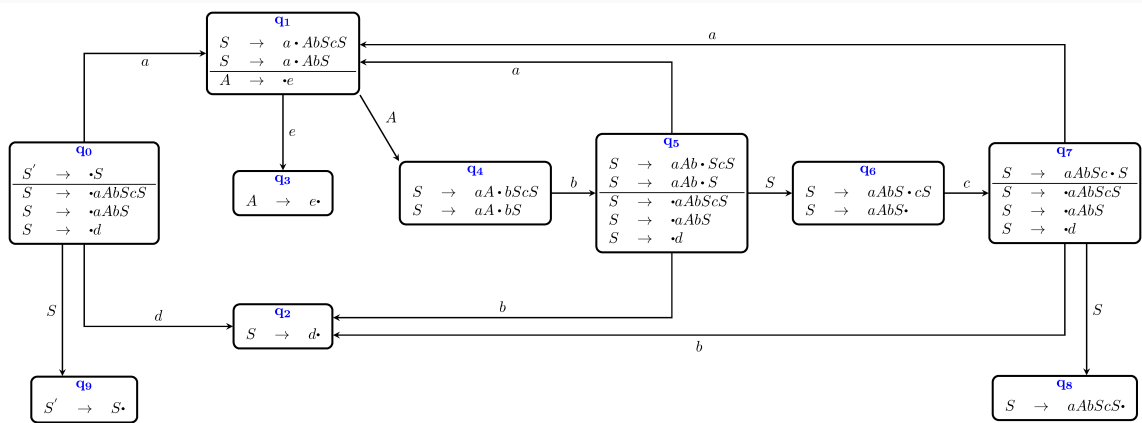
Simple LR(1) = (SLR)-Parsing

$A \rightarrow \beta$ wird nur reduziert, wenn das Vorschautoken in der *FOLLOW*-Menge von A ist.

⇒ Es ändert sich nur die Parse Table:

Bei allen LR(0)-Items in der Tabelle, die einen Punkt am Ende der rechten Seite stehen haben, trage in der Aktionstabelle beim zugehörigen Zustand die Reduktion mittels der zugehörigen Regel bei allen Terminals ein, die in der FOLLOW-Menge des Nonterminals auf der linken Seite der Regel enthalten sind.

Der SLR-Automat der Grammatik G1:



Die SLR-Parsertabelle der Grammatik G1

Zustand	Aktion						Sprung	
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	\perp	<i>S</i>	<i>A</i>
q_0	s1			s2			q_9	
q_1					s3			q_4
q_2			r3			r3		
q_3		r4						
q_4		s5						
q_5	s1	s2					q_6	
q_6			s7,r2			r2		
q_7	s1	s2					q_8	
q_8			r1			r1		
q_9						acc		

Kanonische LR(1)-Syntaxanalyse

Mehr geht nicht: Kanonische LR(1)-Syntaxanalyse = LR-Analyse

Beim SLR-Verfahren wird nach $A \rightarrow \beta$ reduziert, wenn das Vorschautoken in $\text{Follow}(A)$ liegt. Dabei kann es vorkommen, dass das Vorschautoken ein Element davon ist, aber genau bei dieser Regel kann es nicht dem A folgen. Es wird also falsch reduziert, und es entstehen zu viele Einträge in der Tabelle (Konflikte!).

Jetzt werden nicht Follow-Mengen von Nichtterminalen, sondern LOOKAHEAD-Mengen von Produktionen berechnet.

$S \rightarrow A \alpha B$
 $A \rightarrow b B | c$
 $B \rightarrow A$

$\text{Follow}(B) = \{ \alpha \dots \}$
 $S \Rightarrow A \alpha B \Rightarrow b B \alpha B \dots$

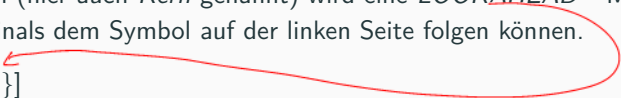
$\begin{array}{l} B \\ A \alpha B \end{array} \left\{ \begin{array}{l} B \rightarrow A \text{ hätte nicht} \\ \text{angewandt} \\ \text{werden dürfen.} \end{array} \right.$
 $\Rightarrow \text{Follow}(A) \text{ als Kriterium}$
 nicht nicht

\Rightarrow Wir benötigen lookahead von Produktionen

Die LR(1)-Items

Zu jedem LR(0)-Item (hier auch *Kern* genannt) wird eine *LOOKAHEAD* - Menge L hinzugefügt, die angibt, welche Terminals dem Symbol auf der linken Seite folgen können.

z. B. $[S' \rightarrow \cdot S, \{\perp\}]$



1. füge I zu $CLOSURE_1(I)$ hinzu
2. gibt es ein LR(1) - Item $[A \rightarrow \alpha \cdot B\beta, L]$ aus $CLOSURE_1(I)$ und eine Produktion $(B \rightarrow \gamma)$, füge $[B \rightarrow \cdot\gamma, FIRST(\beta L)]$ zu $CLOSURE_1(I)$ hinzu $\setminus (\alpha, \beta \text{ dürfen } \epsilon \text{ sein})$.

$GOTO_1(I, X)$ = eine Produktion

$CLOSURE_1(\{[A \rightarrow \alpha X \cdot \beta, L] \mid [A \rightarrow \alpha \cdot X \beta, L] \in I\})$

für eine Itemmenge I und $X \in N \cup T, A \in N, \alpha, \beta \in (N \cup T)^*$.

schon da

Der LR(1)-Automat

Der Automat wird analog zum LR(0)-Automaten erstellt mit dem Startzustand

$$[S' \rightarrow \cdot S, \{\perp\}]$$

Die Tabelle unterscheidet sich nur bei der Reduktion von der LR(0)-Tabelle:

Reduktionsoperationen werden in den Spalten der Terminals eingetragen, die in der LOOKAHEAD-Menge der entsprechenden Regel enthalten sind.

Die Beispielgrammatik G2

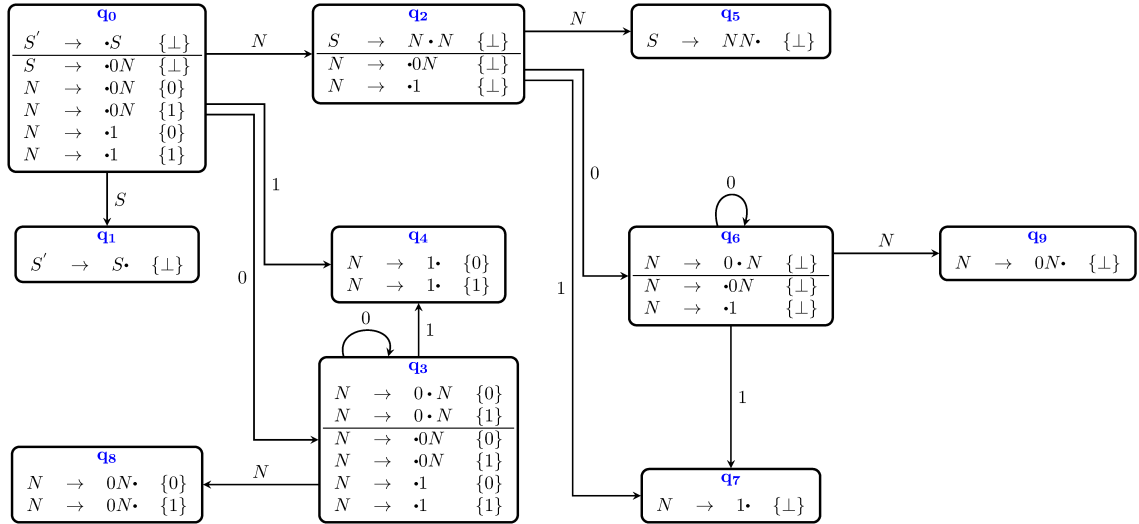
$$(0) S' \rightarrow S$$

$$(1) S \rightarrow NN$$

$$(2) N \rightarrow 0N$$

$$(3) N \rightarrow 1$$

Der LR(1)-Automat der Grammatik G2



Die LR(1)-Parsertabelle der Grammatik G2

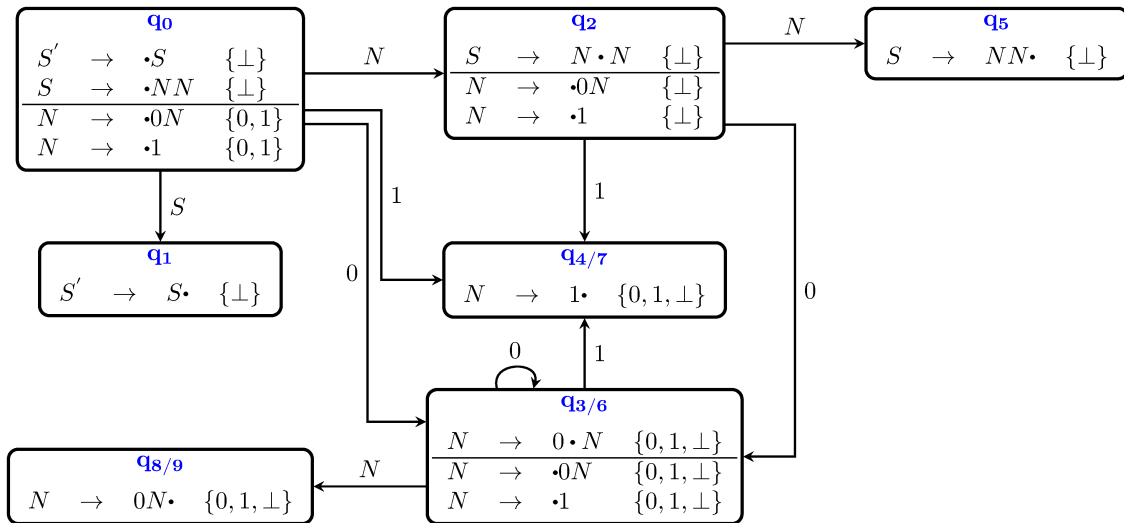
Zustand	Aktion			Sprung	
	0	1	\perp	S	N
q_0	s3	s4		q_1	q_2
q_1			acc		
q_2	s6	s7			q_5
q_3	s3	s4			q_8
q_4	r3	r3			
q_5			r1		
q_6	s6	s7			q_9
q_7			r3		
q_8	r2	r2			
q_9			r2		

Lookahead-LR = LALR

LALR(1)

Zusammenfassung aller LR(1)-Zustände, die sich nur in den LOOKAHEAD-Mengen unterscheiden
Parsergeneratoren generieren oft direkt aus einem LR(0)- einen LALR(1)-Zustands- Übergangsgraphen durch Hinzufügen der LOOKAHEAD-Mengen.

Der LALR-Automat der Grammatik G2



Die LALR-Parsertabelle der Grammatik G2

Zustand	Aktion			Sprung	
	0	1	\perp	S	N
q_0	s(3/6)	s(4/7)		q_1	q_2
q_1			acc		
q_2	s(3/6)	s(4/7)			q_5
$q_{3/6}$	s(3/6)	s(4/7)			$q_{8/9}$
$q_{4/7}$	r3	r3	r3		
q_5			r1		
$q_{8/9}$	r2	r2	r2		

keine Konflikte

Neuer Zustand aus q_5 und q_6

$k \geq 2$ **Vorschautoken**

Zu jeder LR(k)-Sprache gibt es eine LR(1)-Grammatik.

Mehrdeutige Grammatiken

Es gibt auch Auswege

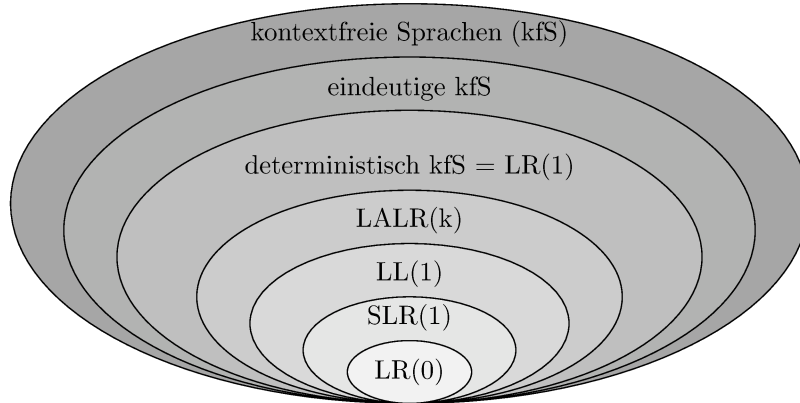
Mehrdeutige Grammatiken sind oft leichter zu lesen und kleiner als die Grammatiken, die man erhält, wenn man die Mehrdeutigkeit auflöst, sofern möglich. Also die Grammatik mehrdeutig lassen!

Folgendes kann trotzdem helfen:

- Angabe von Vorrangregeln
- Angabe von Assoziativität $a^{b^c} = a^{(b^c)}$
- Voreinstellung des Parsergenerators: z. B. Shiften bei Shift-Reduce-Konflikten
- Voreinstellung des Parsergenerators: z. B. Reduzieren nach der Regel, die in der Grammatik zuerst kommt bei Reduce-Reduce-Konflikten

Hierarchie der kontextfreien Sprachen

Hierarchie der kontextfreien Sprachen



Wrap-Up

- mit Bottom-Up-Parsing LR(1) kann man alle deterministisch kontextfreien Sprachen parsen
- ein Vorschautoken genügt
- LR(0)-, SLR- und LALR- Parsing sind vereinfachte Verfahren für Teilmengen der LR-Sprachen



Unless otherwise noted, this work is licensed under CC BY-SA 4.0.

Exceptions

- TODO (what, where, license)