

Kapitel 3

Syntaktische Analyse

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken
- 4 LR-, SLR-Grammatiken
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Syntaktische Analyse

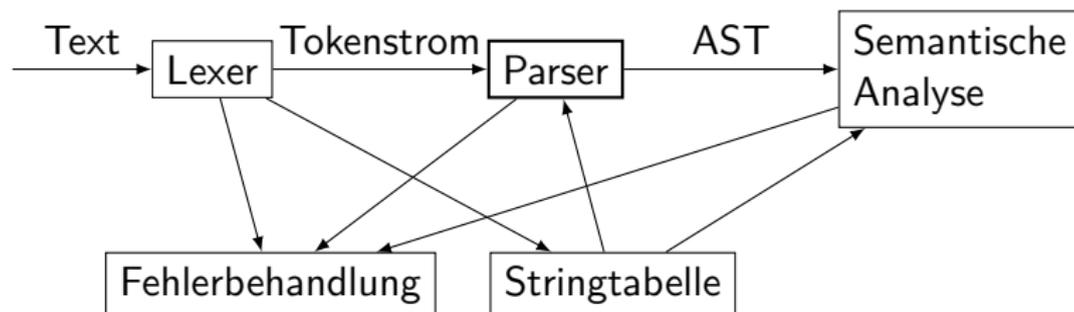
vorgegeben:

- Tokenstrom
- kontextfreie Grammatik (deterministisch?)

Aufgaben

- syntaktische Struktur bestimmen
- syntaktische Fehler melden, korrigieren (?)
- Ausgabe (immer): abstrakte Syntax (Rechts-/Linksableitung), Symbole (Bezeichner, Konstanten, usw.)

Eingliederung in den Übersetzer



ADT Parser

gelieferte Operationen:

- `parse()` : AST

benötigte Operationen:

Lexer/Tokenstrom:

- `next_token()` : Token

Fehlerbehandlung:

- `add_error(nr, pos)`

für Aufbau Strukturbaum:

- `production(nr), symbol(value)`

Aufgabe des Parsers, formal

Gegeben: Grammatik $G = (T, N, P, Z)$ mit
 T Alphabet, N Nichtterminale, P Produktionen, Z Zielsymbol

Gesucht: Entscheidung gehört Tokenstrom s zur Sprache $L(G)$,

- wenn ja, Produktionsfolge für Links-/Rechtsableitung
- wenn nein, Fehlerbehandlung zur Korrektur des Tokenstroms

Unterscheide konkrete Syntax G_k und abstrakte Syntax G_a :

Gesucht: Entscheidung, gehört Tokenstrom s zu $L(G_k)$,
wenn ja, Produktionsfolge für Links-/Rechtsableitung für G_a ?
Beziehung zwischen G_k und G_a ?

Annahmen für das Parsen

Syntax ist kontextfrei

- eigentlich ist sie kontext-sensitiv
- aber kontext-sensitive Grammatiken nicht in linearer Zeit parsbar (Kontextfreiheit ist selbsterfüllende Prophezeiung)
- der über die kontextfreie Grammatik hinausgehende Teil der Syntax heißt im Übersetzerbau statische Semantik

Syntax ist deterministisch kontextfrei

- keine wesentliche Einschränkung, da auch vom menschlichen Leser erwünscht

keine Rückkopplung zur lexikalischen Analyse

- sonst gäbe es mehrere Grundzustände des Lexers, gesteuert vom Parser

keine Rückkopplung semantische Analyse – syntaktische Analyse

- **typunabhängige Syntaktische Analyse**: Zustände des Parsers unabhängig von der Namens- und Typanalyse

Fragen

- Wie wird Sprache erkannt?
- Wie wird abstrakter Strukturbaum aufgebaut?
- Was geschieht bei Fehlern?

Historie, kf Grammatiken + Verarbeitung

1955	Definition und Klassifikation (Chomsky und Bar Hillel)
1957–1959	Kellerautomaten (Bauer&Samelson, sequentielle Formelübersetzung, 1959)
1961	formaler Zusammenhang kfG-Kellerautomat (Öttinger)
1958–1966	kfGs und BNF setzen sich für die Syntax von Programmiersprachen durch (Algol 58, Algol 60, ...)
1960–1972	Verfahren des rekursiven Abstiegs (Glennie) und dessen theoretische Fundierung als LL-Grammatiken (auch heute noch oft neu erfunden!)
1963–1969	deterministische kfGs: beschränkte Operatorpräzedenz, LR, SLR, LALR, ...
seit 1972	nichts wesentlich Neues außer Optimierung, Fehlerbehandlung

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken
- 4 LR-, SLR-Grammatiken
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Grundbegriffe

Kontextfreie Grammatik

Eine Grammatik $G = (T, N, P, Z)$ heißt kontextfrei, wenn für jede Produktion $A \rightarrow \alpha$ gilt: $A \in N$.

Sprache $L(G)$ einer Grammatik G

$L(G) = \{\omega \in T^* \mid Z \Rightarrow_G^* \omega\}$ ist die Menge aller in der Grammatik ableitbaren Wörter.

Linksableitung (\Rightarrow_L)

Es wird stets das linkeste Nichtterminal ersetzt.

Rechtsableitung (\Rightarrow_R)

Es wird stets das rechteste Nichtterminal ersetzt.

Beispiel

$G = (N, T, P, Z)$ mit

$$T = \{\text{bez}, +, *, (,)\}$$

$$N = \{Z, A, T, F\}$$

$$P = \{Z \rightarrow A, A \rightarrow T, T \rightarrow F, F \rightarrow \text{bez}, \\ A \rightarrow A + T, T \rightarrow T * F, F \rightarrow (A)\}$$

ist eine kontextfreie Grammatik.

Linksableitung für $a+(b+c)$:

$$A \Rightarrow A+T \Rightarrow T+T \Rightarrow F+T \Rightarrow a+T \Rightarrow a+(A) \Rightarrow \\ a+(A+T) \Rightarrow a+(T+T) \Rightarrow a+(F+T) \Rightarrow a+(b+T) \Rightarrow \\ a+(b+F) \Rightarrow a+(b+c)$$

Rechtsableitung für $a+(b+c)$:

$$A \Rightarrow A+T \Rightarrow A+F \Rightarrow A+(A) \Rightarrow A+(A+T) \Rightarrow \\ A+(A+F) \Rightarrow A+(A+c) \Rightarrow A+(T+c) \Rightarrow A+(F+c) \Rightarrow \\ A+(b+c) \Rightarrow T+(b+c) \Rightarrow F+(b+c) \Rightarrow a+(b+c)$$

Schreibweise der Produktionen

in der Theorie: $A \rightarrow x|y|\dots$, $A \in N$, $x, y \in V^*$, $V = T \cup N$

praktisch: Backus-Naur-Form (BNF)

- Nichtterminale in spitzen Klammern,
- Terminale als Symbole oder wie Nichtterminale
- ::= statt \rightarrow

$\langle \text{Ausdruck} \rangle ::= \langle \text{Term} \rangle \mid \langle \text{Ausdruck} \rangle + \langle \text{Term} \rangle$

Rechnereingabe: Erweiterte Backus-Naur-Form (EBNF)

- wie BNF, aber Bezeichner oft ohne spitze Klammern
- | (oder), . (Abschluß), () (Gruppierung), [] (optional), * (Wiederholung, auch 0-mal), + (Wiederholung, mindestens einmal) als Beschreibungssymbole
- Terminale durch Apostrophs o. ä. ausgezeichnet

$\text{Ausdruck} ::= \text{Term} ('+' \text{Term})^*$

Fortran-, Cobol-, Java-Beschreibung: Abarten von EBNF

Grammar Engineering (1/3)

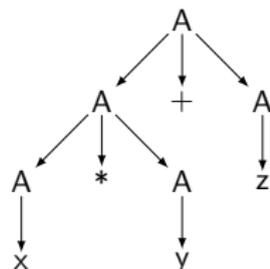
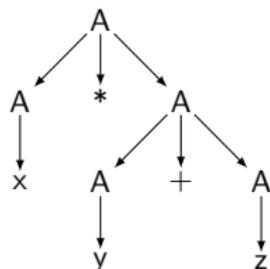
Forderungen:

- deterministische Grammatik: zu einer Eingabe existiert höchstens ein Syntaxbaum
- Operatorprioritäten: Grammatik erzeugt Syntaxbaum gemäß Prioritäten

Gegenbeispiel:

$$A \rightarrow A + A \mid A * A \mid \mathbf{bez} \mid (A)$$

$x*y+z$ hat 2 Syntaxbäume - sogar einen, der „Punkt vor Strich“ ignoriert.



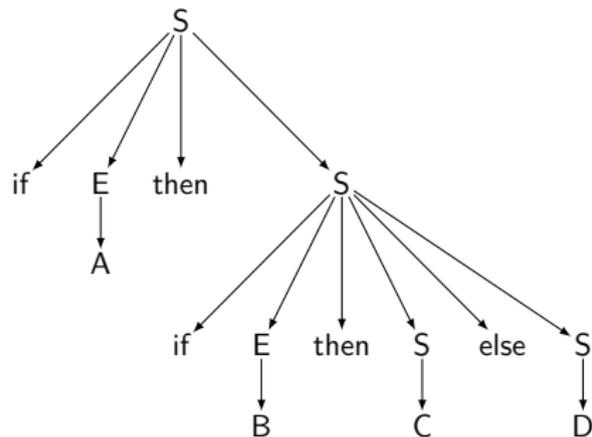
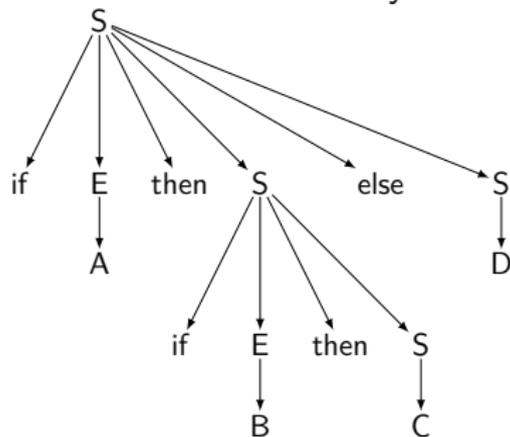
Grammar Engineering (2/3)

Weiteres Gegenbeispiel: „Dangling Else”

$$S \rightarrow \text{if } E \text{ then } S$$
$$| \text{if } E \text{ then } S \text{ else } S$$

if A then if B then C else D

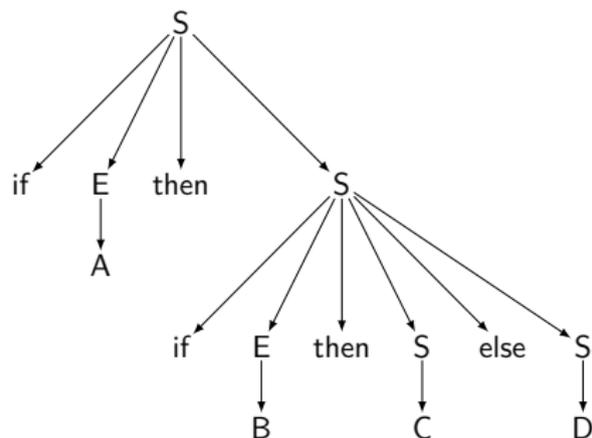
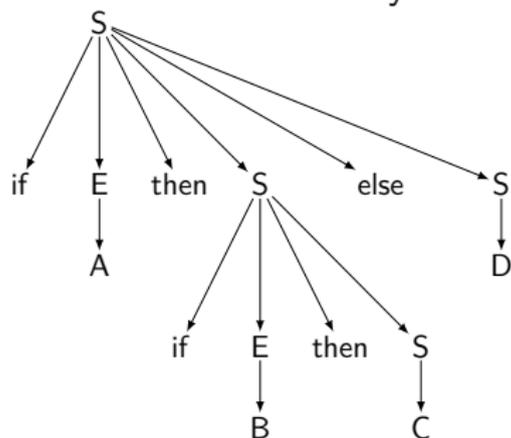
hat 2 verschiedene Syntaxbäume.



Grammar Engineering (3/3)

if A then if B then C else D

hat 2 verschiedene Syntaxbäume.



In der Praxis gehört ein **else** aber immer zum letzten **if**.

- Parsergeneratoren erkennen Mehrdeutigkeiten
- aber in Grammatiken mit hunderten Produktionen sind Mehrdeutigkeiten schwer zu beheben

Faustregeln

- ein Nichtterminal pro Prioritätsebene
- nicht zweimal dasselbe Nichtterminal auf der rechten Seite
- Links- oder Rechtsassoziativität von Operatoren wird durch links- bzw. rechtsrekursive Regeln ausgedrückt

Beispielgrammatik

Ausdrücke:

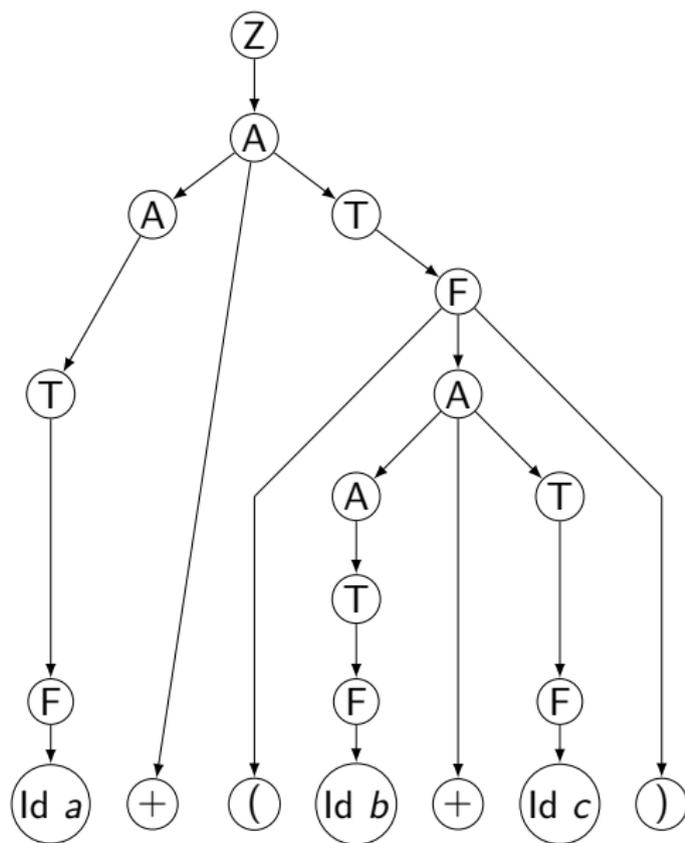
- (0) $Z \rightarrow A$
- (1) $A \rightarrow T$
- (2) $T \rightarrow F$
- (3) $F \rightarrow \text{bez}$
- (4) $A \rightarrow A + T$
- (5) $T \rightarrow T * F$
- (6) $F \rightarrow (A)$

EBNF:

- (0) $Z ::= A.$
- (1) $A ::= T ('+' T)^*.$
- (2) $T ::= F ('*' F)^*.$
- (3) $F ::= \text{bez} \mid '(' A ')'$.

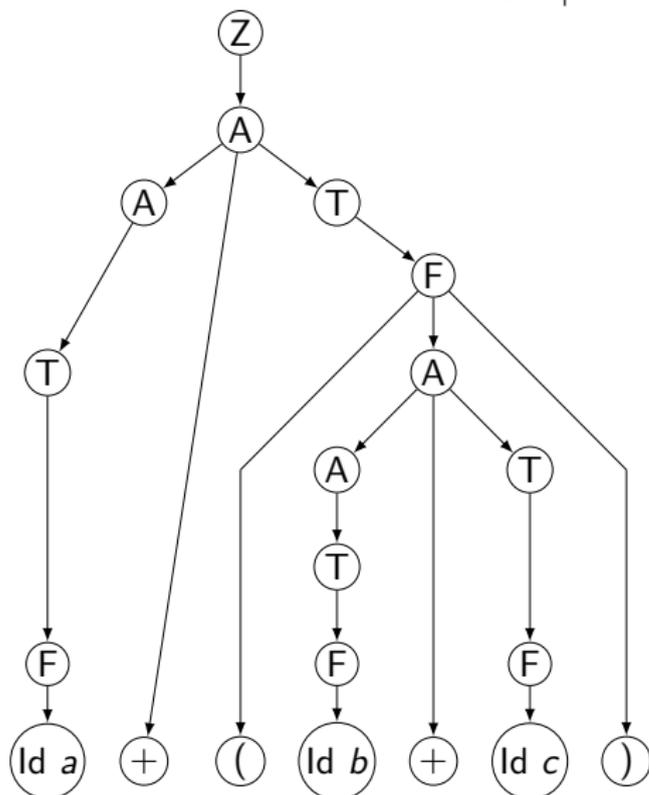
Ausdruck $a + (b + c)$

Konkreter Strukturbaum:

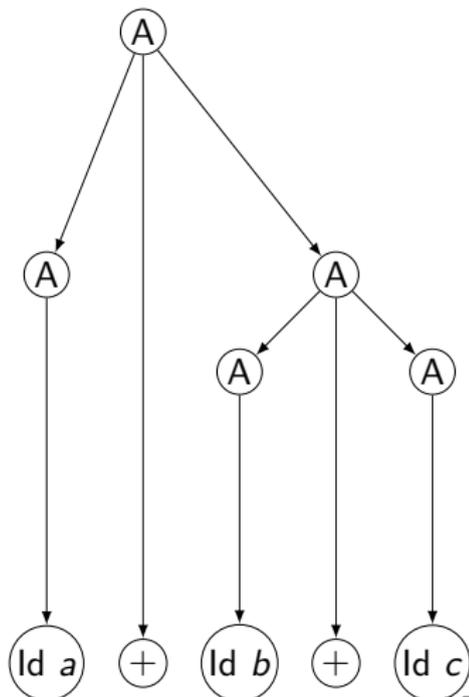


Konkrete und abstrakte Syntax

Prinzip der abstrakten Syntax: nur die für die Semantik wichtige Struktur behalten: $A \rightarrow A + A | A * A | \text{bez}$

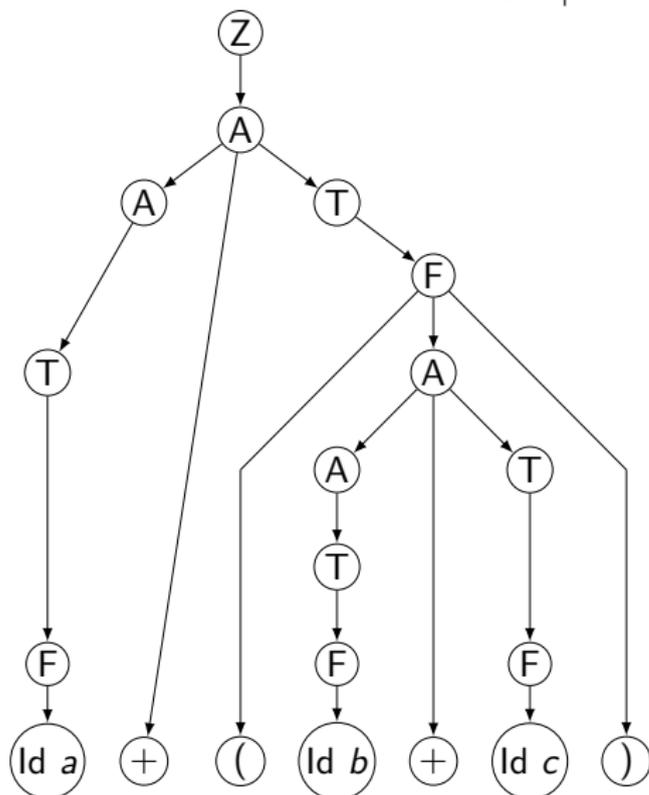


Syntaktische Analyse

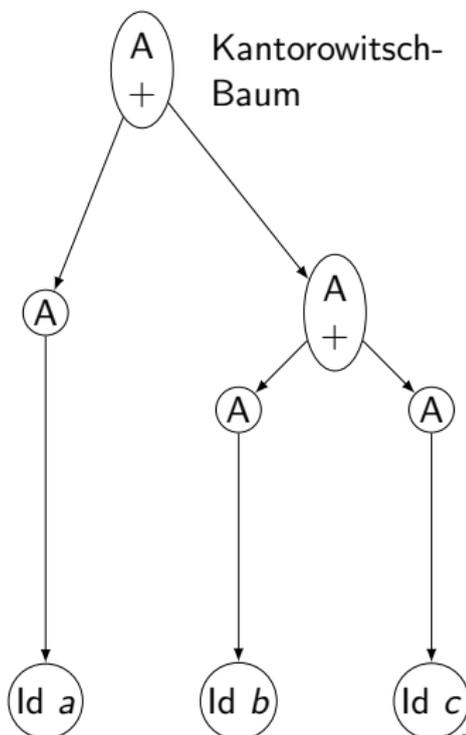


Konkrete und abstrakte Syntax

Prinzip der abstrakten Syntax: nur die für die Semantik wichtige Struktur behalten: $A \rightarrow A + A | A * A | \text{bez}$



Syntaktische Analyse



Übergang: Konkrete und abstrakte Syntax (1/2)

Konkrete Syntax G_k der zu übersetzenden Sprache (Datenstruktur: Tokenstrom)

- 1 explizite Strukturinformation (), begin end, etc.
- 2 Ketten- und Verteilerproduktionen $A \rightarrow B$ bzw.
 $A \rightarrow B \mid C \mid \dots$
- 3 Schlüsselwörter

Abstrakte Syntax G_a des Strukturbaums (Datenstruktur: Baum, AST)

- 1 Klammerung durch AST bereits eindeutig
- 2 Kettenproduktionen überflüssig, wenn keine semantische Bedeutung
- 3 Schlüsselwörter dienen dem eindeutigen Parsen, jetzt immer überflüssig, werden weggelassen.

Übergang: Konkrete und abstrakte Syntax (2/2)

- Abbildung konkrete auf abstrakte Syntax durch Parser (Verarbeitung von semantischen Aktionen), ggf. weitere Transformation während semantischer Analyse
- Produktionsnummer wird Knotentyp
- Operatoren als Attribute des Knotens für den Ausdruck

Abstrakte Syntax als abstrakte Algebra

Heute fasst man eine abstrakte Syntax als Signatur einer ordnungssortierten Termalgebra auf, und einen AST als Term gemäß dieser Signatur.

- Klassen der abstrakten Syntax entsprechen Sorten der Algebra
- innere Baumknoten entsprechen Operatoren (Funktionssymbolen) der Algebra inkl. Signatur.

Beispiel: abstrakte Syntax für Expressions und Statements

$$Stmt = IfStmt \mid IfElseStmt \mid WhileStmt \mid Assignment \mid Block \mid \dots$$
$$IfStmt :: Expr Stmt$$
$$IfElseStmt :: Expr Stmt Stmt$$
$$WhileStmt :: Expr Stmt$$
$$Block :: Decls StmtList$$
$$StmtList :: Stmt + \qquad Expr = Addop \mid MultOp \mid Var \mid \dots$$
$$Assignment :: Var Expr \qquad Addop :: Expr Expr$$
$$Var = \mathbf{Bez} \mid \dots \qquad Multop :: Expr Expr$$

Abstrakte Syntax als abstrakte Algebra

Entsprechende Bäume können auch als Terme dargestellt werden.
Schreibweise zur Konstruktion von Termen z.B.

Beispiele

Addop(Bez(hinz), Bez(kunz))

IfStmt(Bez(test), Assignment(Bez(x), Addop(Bez(x), Bez(y))))

Block(Decls(...), StmtList(s₁, s₂, ..., s₄₂))

Assoziativitäten/Präzedenzen werden durch Termstruktur dargestellt

Achtung: Die abstrakte Syntax enthält keine semantischen Bedingungen z.B. „Typ einer **If**-Expression muss boolesch sein“

Semantische Aktionen

%Ausgabe

- Nach Erkennen des vorgehenden (Nicht-)Terminals ausgeführt
- Für AST: Konstruktor des entsprechenden Knotens im Ableitungsbaum für G_a aufrufen

&Ausgabe

- Wird ausgeführt, wenn Symbol erkannt aber noch nicht fortgeschaltet wurde
- Für AST: Konstruktoren werden gegebenenfalls Merkmale von Symbolen übergeben

Beachte:

- Semantische Aktionen basieren auf Seiteneffekt beim Parsen
- Symbole werden in der Reihenfolge abgenommen, in der sie in der Symbolfolge erscheinen

Ausgabe von Postfix- oder Präfixform

Ausgaberoutinen:

addop	gib aus '+'
mulop	gib aus '*'
bezeichner	gib aus bez
merke	merke bez
bez_aus	gib gemerkten bez aus, falls vorhanden

Postfixform, d.h. abstrakter Syntaxbaum als Rechtsableitung:

- 1 $A \rightarrow T (' + ' T \%addop)^*$
- 2 $T \rightarrow F (' * ' F \%mulop)^*$
- 3 $F \rightarrow \text{bez} \ \&\text{bezeichner} \mid '(' A ')'$

Präfixform, d.h. abstrakter Syntaxbaum als Linksableitung:

- 1 $A \rightarrow T (' + ' \%addop \%bez_aus T)^*$
- 2 $T \rightarrow F (' * ' \%mulop \%bez_aus F)^* \%bez_aus$
- 3 $F \rightarrow \text{bez} \ \&\text{merke} \mid '(' A ')'$

Beispiel Postfixform

$$1 \quad A \rightarrow T (' + ' T \%addop)^*$$

$$2 \quad T \rightarrow F (' * ' F \%mulop)^*$$

$$3 \quad F \rightarrow \text{bez \&bezeichner} | '(' A ')'$$

Ableitung für $x*y+z$

Ausgabe

A

$$\Rightarrow T + T$$

$$\Rightarrow F * F + T$$

$$\Rightarrow x * F + T$$

$$\Rightarrow x * y + T$$

$$\Rightarrow x * y + F$$

$$\Rightarrow x * y + z$$

x

xy*

xy*

xy*z+

Beispiel Präfixform

1 $A \rightarrow T ('+' \text{ \%addop } \%bez_aus T)^*$

2 $T \rightarrow F ('*' \text{ \%mulop } \%bez_aus F)^* \text{ \%bez_aus}$

3 $F \rightarrow \text{bez } \&\text{merke } | '(A)'$

Ableitung für $x*y+z$

Ausgabe

A

$$\Rightarrow T + T$$

$$\Rightarrow F * F + T$$

$$\Rightarrow x * F + T$$

$$\Rightarrow x * y + T$$

$$\Rightarrow x * y + F$$

$$\Rightarrow x * y + z$$

*x

*x+y

*x+y

*x+yz

Sonderfälle in abstrakter Syntax

- Bezeichner:
 - $A \rightarrow \text{bez}$ ist Kettenproduktion, soll aber wegen semantischer Analyse erhalten bleiben
- Klammern in Fortran:
 - Information eigentlich bereits in der Baumstruktur
 - **aber** Klammern sind bindend (kein Umordnen erlaubt)
 - sonst gilt eventuell Assoziativgesetz (Umordnen möglicherweise erlaubt)
 - müssen als Operator gespeichert werden
- Anweisungslisten in C:
 - sind Verteilersymbole
 - **aber** Strichpunkt-Operator legt Auswertungsfolge fest (auch ohne Datenabhängigkeiten), Code-Verschiebung verboten?

Abstrakte Syntax II

- abstrakte Syntax quellsprachenunabhängig?
 - Programmstruktur in semantischer Analyse aufgearbeitet, danach nur noch Prozeduren interessant
 - Prozeduraufrufe nur bezüglich Parameterübergabe unterschiedlich
 - Ablaufsteuerung identisch, eventuelle Ausnahme: Zählschleifen
 - Ausnahmebehandlung in allen modernen Sprachen identisch
 - Zuweisung, Ausdrucksoperatoren, usw.: identisch, manchmal vielleicht Ergänzungen erforderlich
- Konsequenz: weitere Verarbeitung (Transformation, Optimierung, Codegenerierung) weitgehend unabhängig von der Quellsprache
 - Systeme: UNCOL, ANDF, Dotnet
 - Dotnet kann als Postfixcodierung von UNCOL angesehen werden

Beseitigung von ε -Produktionen

Satz:

Für jede kfG G mit ε -Produktionen gibt es eine kfG G' ohne ε -Produktionen mit $L(G) \setminus \{\varepsilon\} = L(G')$ und umgekehrt.

Technik dazu: ε -Abschluß

Einsetzen von Ableitungen der Form $A \rightarrow \varepsilon$ in alle rechten Seiten der Form $X \rightarrow \alpha A \beta$, mit $\alpha, \beta \in (T \cup N)^*$

Beispiel zur Beseitigung von ε -Produktionen

$$(1) Z \rightarrow aS$$

$$(2) S \rightarrow aS \mid \varepsilon$$

Einsetzen von $S \rightarrow \varepsilon$ auf den rechten Seiten führt zu

$$(1) Z \rightarrow aS \mid a$$

$$(2) S \rightarrow aS \mid a$$

ohne ε -Produktionen.

Linksfaktorisierung

Produktionen $X \rightarrow Yb \mid Yc$ mit gleicher LS und gemeinsamem Anfang Y kann man nicht mit rekursivem Abstieg verarbeiten, wenn Länge $|y|$, $Y \Rightarrow^* y$, unbeschränkt, $|y| \geq 0$.

Lösung: den gemeinsamen Anfang ausklammern

Ersetze $X \rightarrow Yb \mid Yc$ durch $X \rightarrow YX'$, $X' \rightarrow b \mid c$

Analog kann man bei LR-Analyse rechtsfaktorisieren (seltener benötigt).

Beispiel zur Linksfaktorisierung

Die Produktionen

$S \rightarrow \text{if } E \text{ then } S \text{ endif}$

$S \rightarrow \text{if } E \text{ then } S \text{ else } S \text{ endif}$

haben gemeinsamen Anfang $\text{if } E \text{ then } S$.

Linksfaktorisierung ergibt:

$S \rightarrow \text{if } E \text{ then } S \ X$

$X \rightarrow \text{endif} \mid \text{else } S \text{ endif}$

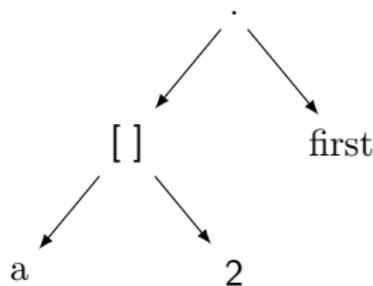
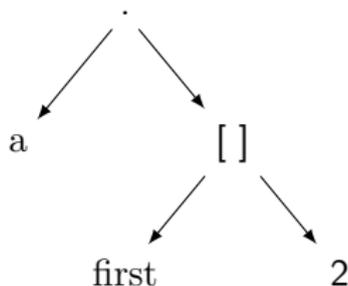
Typunabhängiges Parsen

Typunabhängiges Parsen

- Parsen ohne Kenntnis über Typen von Symbolen
- ist üblich, aber nicht immer ausreichend

Typabhängiges Parsen

- Bsp: ADA `a.first(2)`



Typabhängiges Parsen

Beispiel: Formate in FORTRAN

- `print(r 20, real_const)`
- `r 20` ist Format und muss anders behandelt werden, sonst `r` Bezeichner und `20` ganze Zahl

Parser umschaltbar, um Formate zu bearbeiten

- D.h., es gibt zwei verschiedene Parser
- Erst semantische Analyse erkennt Bezeichner `print`

Umschaltung also semantik- (oder typ-) gesteuert

Ähnliche Probleme in ABAP/4

Kellerautomaten

- Kellerautomat $A = (T, Q, R, q_0, F, S, s_0)$
 - T Eingabealphabet (Tokens)
 - Q Zustandsmenge
 - R Menge von Regeln $sqx \rightarrow s'q'x'$, $s, s' \in S^*$, $q, q' \in Q$, $x, x', x'' \in T^*$, $x = x''x'$
 - q_0 Anfangszustand
 - $F \subseteq Q$ Menge von Endzuständen
 - S Kelleralphabet
 - $s_0 \in S$ Anfangszeichen im Keller
- **Konfiguration:** $\underline{s}q\underline{x}$, \underline{s} vollständiger Kellerinhalt, \underline{x} restliche Eingabe
- Anfangskonfiguration: s_0q_0y , y vollständige Eingabe
- Regel $sqx \rightarrow s'q'x'$ anwendbar, wenn $\underline{s} = \underline{s}'s$, $\underline{x} = x\underline{x}'$
- Folgekonfiguration: $\underline{s}'s'q'x'x'$
- Halt bei Konfiguration sq , $q \in F$, Eingabe vollständig gelesen
- praktisch Endezeichen $\#$ erforderlich, Halt bei $sq\#$

Beispiel: Kellerautomat für Palindrome

Kellerautomat $A = (T, Q, R, q_0, F, S, s_0)$ mit

- $Q = \{q_0, q_1, q_2\}$
- $R = \{q_0t \rightarrow tq_0 \mid t \in T\} \cup \{q_0t \rightarrow q_1 \mid t \in T \cup \{\varepsilon\}\} \cup \{tq_1t \rightarrow q_1 \mid t \in T\} \cup \{s_0q_1\# \rightarrow q_2\#\}$
- $F = \{q_2\}$
- $S = T \cup \{s_0\}$

Abarbeitung von *otto*:

Keller	Zustand	Eingabe
s_0	q_0	<i>otto</i> #
s_0o	q_0	<i>tto</i> #
s_0ot	q_0	<i>to</i> #
s_0ot	q_1	<i>to</i> #
s_0o	q_1	<i>o</i> #
s_0	q_1	#
	q_2	#

Kontextfreie Grammatik und Kellerautomaten

Satz:

Für jede kontextfreie Grammatik G gibt es einen (nichtdeterministischen) Kellerautomaten A mit $L(A) = L(G)$.

⇒ das Akzeptionsproblem für kontextfreie Sprachen ist entscheidbar

Aber: Aufwand i.a. $\mathcal{O}(n^3)$

⇒ praktisch nur Teilklassen mit linearem Aufwand brauchbar, dazu Grammatik-Umformungen erforderlich

Aber: Sprachinklusion und Gleichheit nicht entscheidbar

⇒ keine eindeutige Normalform

Systematische Parserkonstruktion

- Es gibt weit mehr als 25 verschiedene Techniken zur Parserkonstruktion, vgl. Aho&Ullman, The Theory of Parsing and Compiling, 2 Bde, 1972
- Nur zwei Techniken, LL und LR, haben die Eigenschaften:
 - Der Parser liest die Quelle **einmal** von **links** nach rechts und baut dabei die **Links-** bzw. **Rechtsableitung** auf (daher die 2 Buchstaben).
 - Der Parser erkennt einen Fehler beim ersten Zeichen t , das nicht zu einem Satz der Sprache gehören kann. t heißt **parserdefinierte Fehlerstelle** (parser defined error): Wenn $x \notin L(G)$ und der Parser erkennt den Fehler beim Zeichen t , $x = x'tx''$, so gibt es einen Satz $y \in L(G)$ mit $y = x'y'$.
 - Alternative: Erkennen des Fehlers einige Zeichen später, keine syntaktische Fehlerlokalisierung möglich.

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken**
- 4 LR-, SLR-Grammatiken
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Herleitung der LL- und LR-Parser

- gegeben Grammatik $G = (T, N, P, Z)$, $V = T \cup N$, konstruiere **indeterministischen** Kellerautomat mit genau einem Zustand q , angesetzt auf Eingabe x

Für LL: (prädiktiv)

$tqt \rightarrow q, t \in T$

$Xq \rightarrow x_n \dots x_1 q,$

$X \rightarrow x_1 \dots x_n \in P$

Für LR: (reduzierend)

$qt \rightarrow tq, t \in T$

$x_1 \dots x_n q \rightarrow Xq,$

$X \rightarrow x_1 \dots x_n \in P$

- mache Kellerautomat deterministisch durch Hinzunahme Rechtskontext, also Vorhersage $Xqx' \rightarrow x_n \dots x_1 qx'$ bzw. Reduktion $x_1 \dots x_n qx' \rightarrow Xqx'$
 x' Anfang des unverarbeiteten Eingaberests
- **deterministisch machen geht nur für eingeschränkte Grammatikklassen**

Nichtdeterministische LL- und LR-Parser

Für LL: (prädiktiv)

Vergleich (compare):

$tqt \rightarrow q, t \in T$

Vorhersage (produce):

$Xq \rightarrow x_n \dots x_1 q,$

$X \rightarrow x_1 \dots x_n \in P$

Für LR: (reduzierend)

Schift (shift):

$qt \rightarrow tq, t \in T$

Reduktion (reduce):

$x_1 \dots x_n q \rightarrow Xq,$

$X \rightarrow x_1 \dots x_n \in P$

top-down Parser

vom Startsymbol zum Wort

bottom-up Parser

vom Wort zum Startsymbol

Anmerkung: Der Zustand q ist noch bedeutungslos, er wird später beim deterministisch Machen benötigt.

Textmengen

$k : x$

$k : x = x\#$ falls $x = x_1 \dots x_m \wedge m < k$

$k : x = x_1 \dots x_k$ falls $x = x_1 \dots x_m \wedge m \geq k$

$\text{Anf}_k(x) = \{u \mid \exists y \in T^* : x \Rightarrow^* y \wedge u = k : y\}$

in der Literatur auch $\text{First}_k(x)$ genannt

$A \Rightarrow_{R'} \alpha$ gdw. $A \Rightarrow_R \alpha \wedge \nexists B \in N : A \Rightarrow_R B\alpha \Rightarrow \alpha$

$\text{Anf}'_k(x) = \{u \in \text{Anf}_k(x) \mid \exists y \in T^* : x \Rightarrow_{R'} uy\}$

in der Literatur auch $\text{EFF}_k(x)$ genannt (EFFective First)

$\text{Folge}_k(x) = \{u \mid \exists m, y \in V^* : Z \Rightarrow^* mxy \wedge u \in \text{Anf}_k(y)\}$

in der Literatur auch $\text{Follow}_k(x)$ genannt

LL(k)-Grammatiken

Für $k \geq 1$ heißt eine kfG $G = (T, N, P, Z)$ eine LL(k)-Grammatik, wenn für alle Paare von Ableitungen

$$\begin{array}{ll} Z \Rightarrow_L \mu A \chi \Rightarrow \mu \nu \chi \Rightarrow^* \mu \gamma & \mu, \gamma \in T^*; \nu, \chi \in V^*, A \in N \\ Z \Rightarrow_L \mu A \chi' \Rightarrow \mu \omega \chi' \Rightarrow^* \mu \gamma' & \gamma' \in T^*; \omega, \chi' \in V^* \end{array}$$

gilt:

$$(k : \gamma = k : \gamma') \Rightarrow \nu = \omega$$

Also: Aus den nächsten k Zeichen kann unter Berücksichtigung des Kellerinhalts die nächste anzuwendende Produktion eindeutig vorhergesagt werden.

Die k Zeichen können aus der Produktion resultieren oder ganz oder teilweise dem Folgetext angehören, z.B. bei ε -Produktionen.

Beispiele von LL-Grammatiken

- $A \rightarrow TA', A' \rightarrow \varepsilon \mid '+' TA', T \rightarrow FT',$
 $T' \rightarrow \varepsilon \mid '*' FT', F \rightarrow \text{bez} \mid '(' A ')'$ ist LL(1).
- $Z \rightarrow aAab \mid bAbb, A \rightarrow \varepsilon \mid a$ ist LL(2), nicht LL(1):
Vorschau aa, ab, bb entscheidet unter Berücksichtigung der
Produktion für Z über die Produktion für A .
- $Z \rightarrow X, X \rightarrow Y \mid bYa, Y \rightarrow c \mid ca$ ist LL(3).
- $Z \rightarrow X, X \rightarrow Yc \mid Yd, Y \rightarrow a \mid bY$ ist für kein k LL(k);
aber Linksfaktorisieren macht daraus LL(1).
- Anweisungen, die mit Schlüsselwort `while`, `if`, `case`, usw.
beginnen, sind mit LL(1)-Technik vorhersagbar. Bei Beginn
mit Bezeichner sind Linksfaktorisierungen nötig.

Satz über Linksrekursion

Satz:

Eine linksrekursive kfG ist für kein k $LL(k)$.

Beweisidee:

Seien $A \rightarrow Ax$ und $A \rightarrow y$ linksrekursive bzw. terminierende Regeln. Jeder k -Anfang der terminierenden Regel ist auch k -Anfang der linksrekursiven Regel.

Elimination von Linksrekursion (1/2)

Satz:

Für jede kfG G mit linksrekursiven Produktionen gibt es eine kfG G' ohne Linksrekursion mit $L(G) = L(G')$.

Elimination von Linksrekursion (2/2)

Konstruktion:

- Nummeriere Nichtterminale beliebig X_1, \dots, X_n
- Für $i = 1, \dots, n$
 - Für $j = 1, \dots, i - 1$ ersetze $X_i \rightarrow X_j x$ durch $\{X_i \rightarrow y_j x \mid X_j \rightarrow y_j \in P\}$ (danach $i \leq j$, wenn $X_i \rightarrow X_j x \in P$)
 - Ersetze die Produktionsmengen $\{X_i \rightarrow X_j x\} \cup \{X_i \rightarrow z \mid z \neq X_j z'\}$ durch $\{Y_i \rightarrow x Y_i \mid X_i \rightarrow X_j x \in P\} \cup \{Y_i \rightarrow \varepsilon\} \cup \{X_i \rightarrow z Y_i \mid X_i \rightarrow z \in P \wedge z \neq X_j z'\}$ mit einem neuen Nichtterminal Y_i . (Nummerierung der Y_i mit $n + 1, n + 2, \dots$)
- **Ergebnis:** $i < j$, wenn $X_i \rightarrow X_j x \in P$

Beachte: in Schritt 2 Ersetzung durch $\{Y_i \rightarrow x, Y_i \rightarrow x Y_i \mid X_i \rightarrow X_j x \in P\} \cup \{X_i \rightarrow z, X_i \rightarrow z Y_i \mid X_i \rightarrow z \in P \wedge z \neq X_j z'\}$ ohne ε -Produktionen möglich, wenn x nicht mit X_j , $j \leq i$, beginnt.

Beispiel

- $A \rightarrow T \mid A + T$, $T \rightarrow F \mid T * F$, $F \rightarrow \text{bez} \mid (A)$ ist linksrekursiv
- Ersetzung: Schritt 1 leer, Schritt 2: $A \rightarrow T \mid A + T$ durch $A \rightarrow T A'$, $A' \rightarrow \varepsilon \mid + TA'$ ersetzen; $T \rightarrow F \mid T * F$ analog. Dies entspricht der EBNF $A \rightarrow T ('+' T)^*$, $T \rightarrow F ('*' F)^*$, $F \rightarrow \text{bez} \mid '(' A ')'$.
- Andere mögliche Ersetzung $A \rightarrow T \mid T A'$, $A' \rightarrow + T \mid TA'$
- **Vorsicht:** Die Ersetzung durch $A \rightarrow T \mid T + A$ ist semantisch unzulässig!
- Sie transformiert Links- in Rechtsassoziativität, verändert also die semantisch bedeutungsvolle Struktur.
- Beseitigung von Linksrekursion bei $LL(k)$ -Analyse nötig für alle Anweisungen, die mit $\langle \text{Bezeichner} \rangle \langle \text{Operator} \rangle$ anfangen können (Zuweisungen, Ausdrücke)

SLL(k)-Grammatiken

Für $k \geq 1$ heißt eine kfG $G = (T, N, P, Z)$ eine SLL(k)-Grammatik (**starke LL-Grammatik**), wenn für alle Paare von Ableitungen

$$\begin{array}{ll} Z \Rightarrow_L \mu A \chi \Rightarrow \mu \nu \chi \Rightarrow^* \mu \gamma & \mu, \gamma \in T^*; \nu, \chi \in V^*, A \in N \\ Z \Rightarrow_L \mu' A \chi' \Rightarrow \mu' \omega \chi' \Rightarrow^* \mu' \gamma' & \mu', \gamma' \in T^*; \omega, \chi' \in V^* \end{array}$$

gilt:

$$(k : \gamma = k : \gamma') \Rightarrow \nu = \omega$$

Also: Aus den nächsten k Zeichen kann **ohne** Berücksichtigung des Kellerinhalts die nächste anzuwendende Produktion eindeutig vorhergesagt werden.

SLL-Bedingung

Satz:

Eine Grammatik ist genau dann eine SLL(k)-Grammatik, wenn für alle Paare von Produktionen $A \rightarrow x \mid x'$, $x \neq x'$, die SLL(k)-Bedingung gilt:

$$\text{Anf}_k(x\text{Folge}_k(A)) \cap \text{Anf}_k(x'\text{Folge}_k(A)) = \emptyset$$

Beweis: trivial

- Also: SLL(k)-Eigenschaft durch Berechnung von Anf_k - und Folge_k -Mengen einfach nachzuprüfen.
- Wenn aus x, x' nur terminale Zeichenreihen mit mindestens k Zeichen ableitbar sind, trägt $\text{Folge}_k(A)$ nichts zum Ergebnis bei und kann entfallen.
- wichtiger Spezialfall: $k = 1$, $x \not\Rightarrow^* \varepsilon$, $x' \not\Rightarrow^* \varepsilon$. Dann muss

$$\text{Anf}_k(x) \cap \text{Anf}_k(x') = \emptyset$$

gelten. Falls $x \Rightarrow^* \varepsilon$, so muss außerdem gelten:

$$\text{Folge}_k(A) \cap \text{Anf}_k(x') = \emptyset$$

LL(1) und SLL(1)

Satz: Jede SLL(k)-Grammatik ist auch eine LL(k)-Grammatik.

Satz: Jede LL(1)-Grammatik ist eine SLL(1)-Grammatik.

Beweis (indirekt):

Angenommen, G ist LL(1), aber die SLL(1)-Bedingung ist nicht erfüllt. Dann gibt es Produktionen $A \rightarrow x \mid x'$, $x \neq x'$, und ein Zeichen

$$t \in \text{Anf}_1(x\text{Folge}_1(A)) \cap \text{Anf}_1(x'\text{Folge}_1(A)).$$

Fall $t \in \text{Anf}_1(x)$, $t \in \text{Anf}_1(x')$ verstößt gegen die LL(1)-Definition, da wegen $x = \nu$, $t \in \text{Anf}_1(\nu\chi)$ und $x' = \omega$, $t \in \text{Anf}_1(\omega\chi')$ gilt:
 $1 : \gamma = 1 : \gamma'$, jedoch $\nu \neq \omega$. Widerspruch.

Andere Fälle analog.

Satz nicht auf $k > 1$ verallgemeinerbar:

$Z \rightarrow aAab \mid bAbb$, $A \rightarrow \varepsilon$ | a ist LL(2), aber nicht SLL(2).

Konstruktion der LL(1)-Tabelle

$$LL[X, \mathbf{a}] = \{X \rightarrow X_1 \dots X_n \in P \mid \mathbf{a} \in \text{Anf}_1(X_1 \dots X_n \text{Folge}_1(X))\}$$

Es muss gelten $|LL[X, \mathbf{a}]| = 1$ für alle X, \mathbf{a} , sonst ist die Grammatik nicht LL(1).

Parsertabelle

Grammatik:

$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \varepsilon$$

$$F \rightarrow (E) \mid \mathbf{id}$$

Nichtterminal	Anf ₁	Folge ₁	Eingabesymbol					
			id	+	*	()	#
<i>E</i>	(, id), #	$E \rightarrow TE'$			$E \rightarrow TE'$		
<i>E'</i>	+, #), #		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
<i>T</i>	(, id	+, #	$T \rightarrow FT'$			$T \rightarrow FT'$		
<i>T'</i>	*, #	+, #		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
<i>F</i>	(, id	+, *, #	$F \rightarrow \mathbf{id}$			$F \rightarrow (E)$		

Abbildung: Parsertabelle

Parsertabelle

Grammatik:

$$S \rightarrow \mathbf{iEtSS'} \mid \mathbf{a}$$

$$S' \rightarrow \mathbf{eS} \mid \varepsilon$$

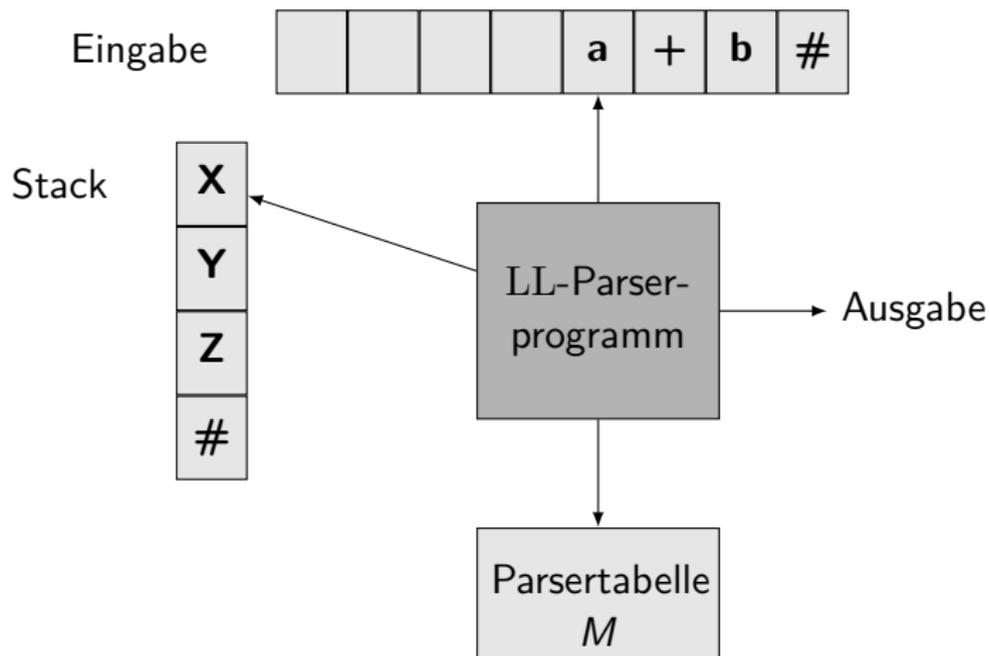
$$E \rightarrow \mathbf{b}$$

Parsertabelle:

Nichtterminal	Anf ₁	Folge ₁	Eingabesymbol					
			a	b	e	i	t	#
S	i, a	e, #	$S \rightarrow \mathbf{a}$			$S \rightarrow \mathbf{iEtSS'}$		
S'	e, #	e, #			$S' \rightarrow \mathbf{eS}$			$S' \rightarrow \varepsilon$
E	b	t		$E' \rightarrow \mathbf{b}$				

Nach Definition ist $LL[S', e] = \{S' \rightarrow eS, S' \rightarrow \varepsilon\}$ und somit die Grammatik nicht LL(1). Zur Auflösung des Konfliktes wird die zweite Produktion manuell aus dem Tabelleneintrag entfernt. Dadurch gehört ein **e** immer zum letzten **i**.

Modell eines tabellengesteuerten LL-Parsers



Verhalten eines LL-Parsers

Übereinstimmung	Stack	Eingabe	Aktion
	$E\#$	id+id*id#	
	$TE'\#$	id+id*id#	Ausgabe von $E \rightarrow TE'$
	$FT'E'\#$	id+id*id#	Ausgabe von $T \rightarrow FT'$
	id $T'E'\#$	id+id*id#	Ausgabe von $F \rightarrow \mathbf{id}$
id	$T'E'\#$	+id*id#	Übereinstimmung mit id
id	$E'\#$	+id*id#	Ausgabe von $T' \rightarrow \varepsilon$
id	$+TE'\#$	+id*id#	Ausgabe von $E' \rightarrow +TE'$
id+	$TE'\#$	id*id#	Übereinstimmung mit +
id+	$FT'E'\#$	id*id#	Ausgabe von $T \rightarrow FT'$
id+	id $T'E'\#$	id*id#	Ausgabe von $F \rightarrow \mathbf{id}$
id+id	$T'E'\#$	*id#	Übereinstimmung mit id
id+id	$*FT'E'\#$	*id#	Ausgabe von $T' \rightarrow *FT'$
id+id*	$FT'E'\#$	id#	Übereinstimmung mit *
id+id*	id $T'E'\#$	id#	Ausgabe von $F \rightarrow \mathbf{id}$
id+id*id	$T'E'\#$	#	Übereinstimmung mit id
id+id*id	$E'\#$	#	Ausgabe von $T' \rightarrow \varepsilon$
id+id*id	#	#	Ausgabe von $E' \rightarrow \varepsilon$

LL(1)-Parseralgorithmus

```
push('#'); push(Z); t = next_token();
while (t != '#') {
    if (stackEmpty()) { error("end of input expected"); }
    else if (top() ∈ T) {
        if (top() == t) {
            pop(); t = next_token();
        } else { error(top() + " expected"); pop(); }
    } else if (LL[top(), t] == ⊥) {
        error("illegal Symbol " + t); t = next_token();
    } else {
        (X → X1 ... Xn) = LL[top(), t];
        pop();
        for(i = n; i >= 1; --i)
            push(Xi);
    }
}
if (top() != '#')
    error("unexpected end of input");
```

LL(1)-Parser mit rekursivem Abstieg

- 1 Definiere Prozedur X für alle Nichtterminale X
- 2 Für alternative Produktionen $X \rightarrow X_1 \mid \dots \mid X_n$ sei Rumpf von X

```
switch t {  
  case Anf1(X1Folge1(X)) : Code für X1;  
  ...  
  case Anf1(XnFolge1(X)) : Code für Xn;  
  default : Fehler(...);  
}
```

- 3 Für rechte Seite $X_i = Y_1 \dots Y_m$ erzeuge:
 $C_1; \dots; C_m$; **return**;

Es gilt $C_i =$

- 1 **if** (t == Y_i) t = nächstesSymbol() **else** Fehler(...);
wenn $Y_i \in T$
- 2 $Y_i()$;
wenn $Y_i \in N$

Parser aus Grammatik in EBNF

Nichtterminal	X	$X()$;
Terminal	t	if (symbol == t) nextSymbol(); else Fehler();
Option	$[X]$	if (symbol \in Anf ₁ (X)) $X()$;
Iteration	X^+	do $X()$; while (symbol \in Anf ₁ (X));
	X^*	while (symbol \in Anf ₁ (X)) $X()$;
Liste	$X d$	$X()$; while (symbol \in Anf ₁ (d)) { $d()$; $X()$; }
semantische Aktion	$t\&Y$	if (symbol == t) { $Y()$; nextSymbol(); } else Fehler();
	$\%Z$	$Z()$;

Beispielgrammatik

Beispielgrammatik in EBNF-Notation zum Parser auf der nächsten Folie:

1 $Z ::= A$

2 $A ::= T \{ '+' T \}^*$

3 $T ::= F \{ '*' F \}^*$

4 $F ::= id \mid '(' A ')'$

Parser für Beispielgrammatik (1/2)

```
AST parse() { t = nextSymbol(); return Z(); }
```

```
AST Z() { return A(); }
```

```
AST A() {
```

```
    AST res = T(); // merke 1. Operand
```

```
    while (t == '+') {
```

```
        t = nextSymbol();
```

```
        AST res1 = new AST(plus);
```

```
        res1.left = res;
```

```
        res1.right = T();
```

```
        res = res1;
```

```
    }
```

```
    return res;
```

```
}
```

```
T() // analog A
```

Hinweis: Additionen/Multiplikationen werden im AST linksassoziativ interpretiert.

Parser für Beispielgrammatik (2/2)

```
AST F() {
    AST res = null;
    if (t == bez) {
        res = new AST(t); t=nextSymbol();
    }
    else if (t == '(') {
        t = nextSymbol(); res = A();
        if (t == ')')
            t = nextSymbol();
        else
            Fehlerbehandlung.Fehlereintrag(KlammerZuFehlt, t.pos);
    }
    else
        Fehlerbehandlung.Fehlereintrag(unzulässigesSymbol, t.pos);
    return res;
}
```

Praxis des rekursiven Abstiegs

- Einfügung von **semantischen Aktionen**:
Semantische Aktion formal wie Produktion $A \rightarrow \varepsilon$ behandeln, statt der Prozedur für ein Nichtterminal A die Ausgabeprozedur aufrufen.
- Rekursiver Abstieg baut Linksableitung auf.
Vorteil: beim Aufbau bereits erste **Berechnungen von semantischen Attributen möglich** (s. Kapitel "Semantische Analyse").
- **Problem**: Durch die Handprogrammierung können leicht während der Wartung syntaktische Eigenschaften eingeschleust werden, die die **Systematik der Syntax und die Unabhängigkeit Syntax-Semantik zerstören**. Negativbeispiel: ABAP 4
- Rekursiver Abstieg kann auch **tabellengesteuert** implementiert werden! Parser wird Interpretierer der Tabelle.
Vorteile: Vermeidung von Prozeduraufrufen, einfachere Fehlerbehandlung. Nachteil: nicht von Hand programmierbar.

Ziel: bei Prüfung der Anwendbarkeit von Regeln $sqx \rightarrow s'q'x'$
Kellerinhalt und Zustand sq mit **einem** Zustandssymbol codieren
(Prüfung mehrerer Einträge im Keller vermeiden)

Lösungsidee:

- bei LL und LR ist s rechte bzw. linke Seite einer Produktion
 $X \rightarrow x_1 \dots x_n$
- Übergänge $tqt \rightarrow q$ (bei LL) bzw. $qt \rightarrow tq$ (bei LR) sind nur zulässig, wenn in der Produktion ein Terminalzeichen t ansteht, $x_1 \dots x_n = x'tx'''$, wobei $x'' := tx'''$
- also: ersetze sqx durch **Situation** $[X \rightarrow x' \cdot x''; x]$, die durch den Punkt anzeigt, wie weit die Produktion abgearbeitet ist.
- Situationen $[X \rightarrow \cdot x''; x]$ oder $[X \rightarrow x' \cdot; x]$ sind erlaubt und notwendig.
- Verwende Situationen als Zustände **und** als Kellersymbole.
- Situationen heißen englisch *items*.

- 1 Initial $Q = \{q_0\}$ und $R = \emptyset$, mit $q_0 = [Z \rightarrow \cdot S; \#]$.
Anfangszustand und erster Kellerzustand q_0 .
Hinweis: Folge $_k(Z) = \{\#\}$.
- 2 Sei $q = [X \rightarrow \mu \cdot \nu; \Omega] \in Q$ und noch nicht betrachtet.
- 3 Wenn $\nu = \varepsilon$ setze $R := R \cup \{q\varepsilon \rightarrow \varepsilon\}$
Auskellern $q'q \rightarrow q'$ mit beliebigem q' .
- 4 Wenn $\nu = t\gamma$ mit $t \in T$ und $\gamma \in V^*$, setze
 $q' = [X \rightarrow \gamma t \cdot \gamma; \Omega]$. Setze $Q := Q \cup \{q'\}$ und
 $R := R \cup \{qt \rightarrow q'\}$.
- 5 Wenn $\nu = B\gamma$ mit $B \in N$ und $\gamma \in V^*$, setze
 $q' = [X \rightarrow \gamma B \cdot \gamma; \Omega]$ und
 $H = \{[B \rightarrow \cdot \beta_i; \text{Anf}_k(\gamma\Omega)] \mid B \rightarrow \beta_i \in P\}$.
Hinweis: $1 \leq i \leq m$, wenn es m Produktionen mit linker Seite B gibt. Setze $Q := Q \cup \{q'\} \cup H$ und
 $R := R \cup \{q\tau_i \rightarrow q'h_i\tau_i \mid h_i \in H, \tau_i \in \text{Anf}_k(\beta_i\gamma\Omega)\}$.
- 6 Wenn alle $q \in Q$ betrachtet wurden, stop. Sonst, gehe zu 2.

- 1 Initial $Q = \{q_0\}$ und $R = \emptyset$, mit $q_0 = [Z \rightarrow \cdot S; \#]$.
Anfangszustand und erster Kellerzustand q_0 .
Hinweis: $\text{Folge}_k(Z) = \{\#\}$.
- 2 Sei $q = [X \rightarrow \mu \cdot \nu; \Omega] \in Q$ und noch nicht betrachtet.
- 3 Wenn $\nu = \varepsilon$ setze $R := R \cup \{q\varepsilon \rightarrow \varepsilon\}$
Auskellern $q'q \rightarrow q'$ mit beliebigem q' .
- 4 Wenn $\nu = t\gamma$ mit $t \in T$ und $\gamma \in V^*$, setze
 $q' = [X \rightarrow \gamma t \cdot \gamma; \Omega]$. Setze $Q := Q \cup \{q'\}$ und
 $R := R \cup \{qt \rightarrow q'\}$.
- 5 Wenn $\nu = B\gamma$ mit $B \in N$ und $\gamma \in V^*$, setze
 $q' = [X \rightarrow \gamma B \cdot \gamma; \Omega]$ und
 $H = \{[B \rightarrow \cdot \beta_i; \text{Folge}_k(B)] \mid B \rightarrow \beta_i \in P\}$.
Hinweis: $1 \leq i \leq m$, wenn es m Produktionen mit linker Seite B gibt. Setze $Q := Q \cup \{q'\} \cup H$ und
 $R := R \cup \{q\tau_i \rightarrow q'h_i\tau_i \mid h_i \in H, \tau_i \in \text{Anf}_k(\beta_i \text{Folge}_k(B))\}$.
- 6 Wenn alle $q \in Q$ betrachtet wurden, stop. Sonst, gehe zu 2.

Einzig Regel 5 der LL(k) Konstruktion ändert sich:

5 Wenn $\nu = B\gamma$ mit $B \in N$ und $\gamma \in V^*$, setze

$q' = [X \rightarrow \gamma B \cdot \gamma; \Omega]$ und

$H = \{[B \rightarrow \cdot \beta_i; \text{Folge}_k(B)] \mid B \rightarrow \beta_i \in P\}$.

Hinweis: $1 \leq i \leq m$, wenn es m Produktionen mit linker Seite

B gibt. Setze $Q := Q \cup \{q'\} \cup H$ und

$R := R \cup \{q\tau_i \rightarrow q'h_i\tau_i \mid h_i \in H, \tau_i \in \text{Anf}_k(\beta_i \text{Folge}_k(B))\}$.

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken
- 4 LR-, SLR-Grammatiken**
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Nichtdeterministische LL- und LR-Parser

Für LL: (prädiktiv)

Vergleich (compare):

$tqt \rightarrow q, t \in T$

Vorhersage (produce):

$Xq \rightarrow x_n \dots x_1 q,$

$X \rightarrow x_1 \dots x_n \in P$

Für LR: (reduzierend)

Schift (shift):

$qt \rightarrow tq, t \in T$

Reduktion (reduce):

$x_1 \dots x_n q \rightarrow Xq,$

$X \rightarrow x_1 \dots x_n \in P$

top-down Parser

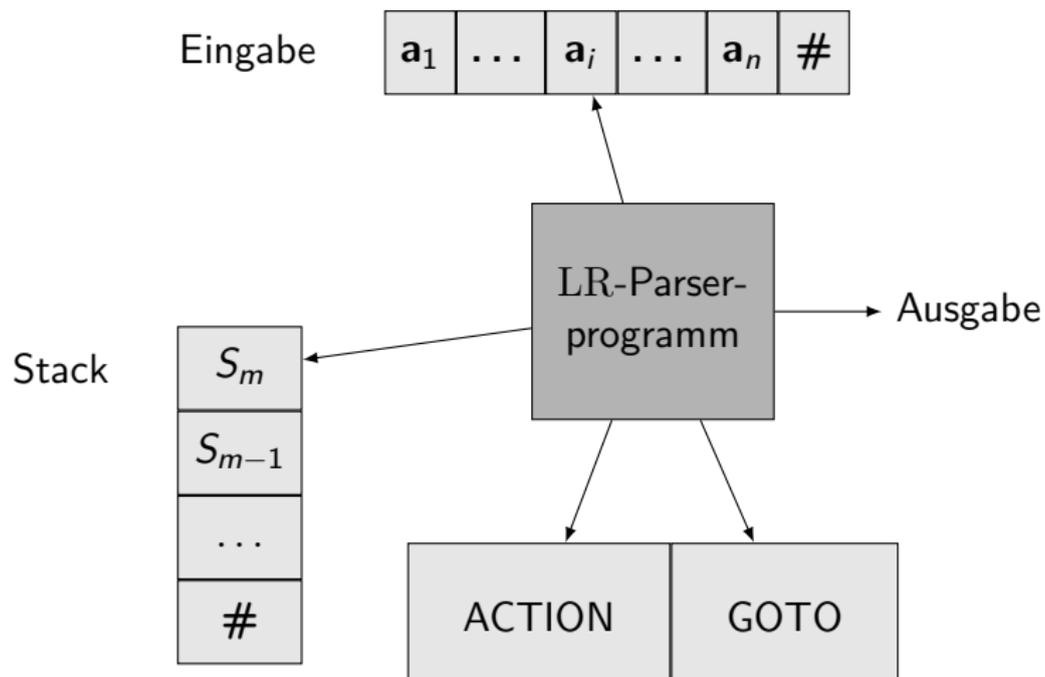
vom Startsymbol zum Wort

bottom-up Parser

vom Wort zum Startsymbol

Anmerkung: Der Zustand q ist noch bedeutungslos, er wird später beim deterministisch Machen benötigt.

Modell eines LR-Parsers



Beispiel: Grammatik G

$$E' \rightarrow E$$

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

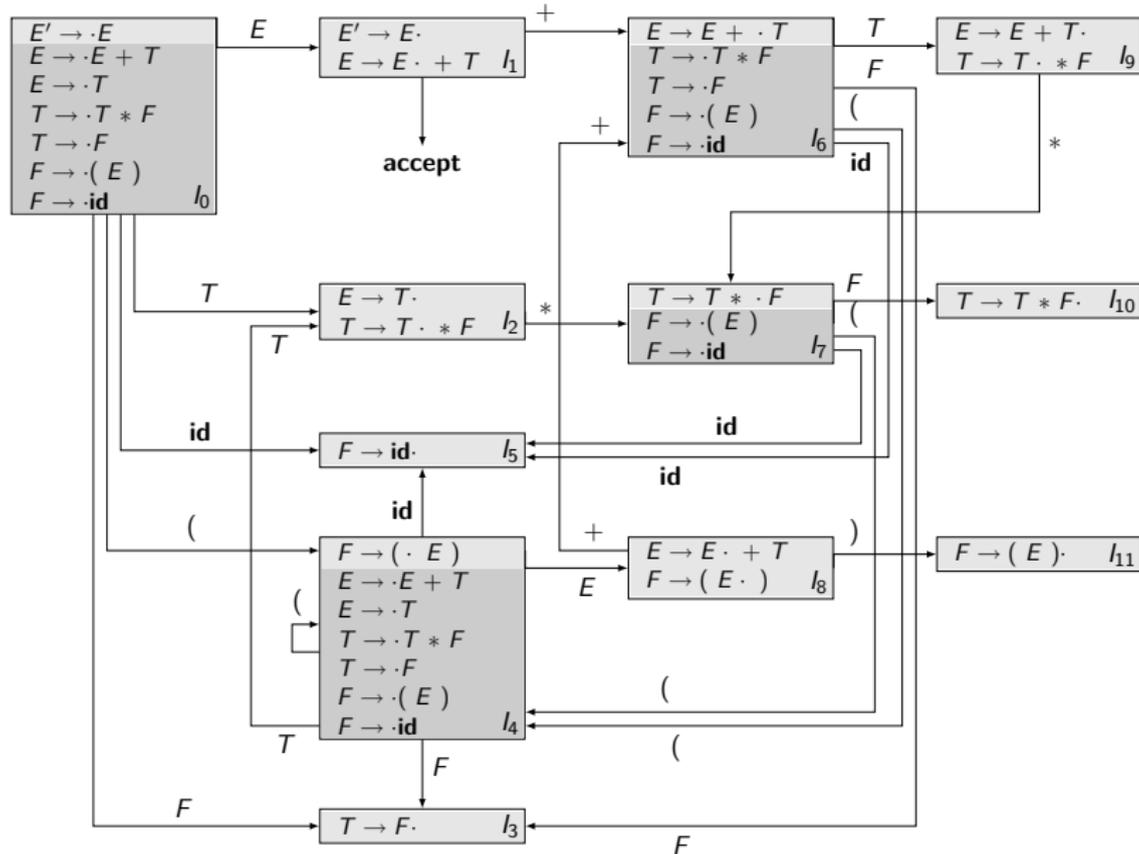
$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \mathbf{id}$$

Verhalten eines Shift-Reduce-Parsers

Stack	Eingabe	Aktion
#	id₁ * id₂ #	Verschieben (shift)
# id₁	* id₂ #	Reduzieren durch $F \rightarrow \mathbf{id}$
# F	* id₂ #	Reduzieren durch $T \rightarrow F$
# T	* id₂ #	Verschieben
# T *	id₂ #	Verschieben
# T * id₂	#	Reduzieren durch $F \rightarrow \mathbf{id}$
# T * F	#	Reduzieren durch $T \rightarrow T * F$
# T	#	Reduzieren durch $E \rightarrow T$
# E	#	Akzeptieren

LR(0)-Automat für Grammatik G



Parsertabelle für Ausdrucksgrammatik

Zustand	ACTION						GOTO		
	id	+	*	()	#	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r3	s7		r3	r3			
3		r5	r5		r5	r5			
4	s5			s4			8	2	3
5		r7	r7		r7	r7			
6	s5			s4				9	3
7	s5			s4					10
8		s6			s11				
9		r2	s7		r2	r2			
10		r4	r4		r4	r4			
11		r6	r6		r6	r6			

Syntaxanalyse von $id * id$

Stack	Symbole	Eingabe	Aktion
0	#	id * id #	Verschieben (shift)
0 5	# id	* id #	Reduzieren durch $F \rightarrow id$
0 3	# F	* id #	Reduzieren durch $T \rightarrow F$
0 2	# T	* id #	Verschieben
0 2 7	# $T *$	id #	Verschieben
0 2 7 5	# $T * id$	#	Reduzieren durch $F \rightarrow id$
0 2 7 10	# $T * F$	#	Reduzieren durch $T \rightarrow T * F$
0 2	# T	#	Reduzieren durch $E \rightarrow T$
0 1	# E	#	Akzeptieren

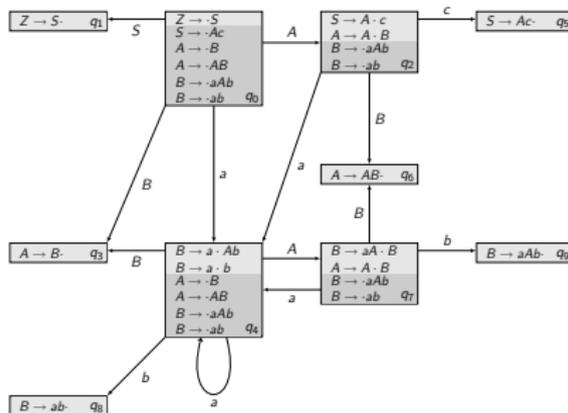
Verhalten eines LR-Parsers

Stack	Symbole	Eingabe	Aktion
0		id * id + id #	Verschieben
0 5	id	* id + id #	Reduzieren durch $F \rightarrow \mathbf{id}$
0 3	F	* id + id #	Reduzieren durch $T \rightarrow F$
0 2	T	* id + id #	Verschieben
0 2 7	T *	id + id #	Verschieben
0 2 7 5	T * id	+ id #	Reduzieren durch $F \rightarrow \mathbf{id}$
0 2 7 10	T * F	+ id #	Reduzieren durch $T \rightarrow T * F$
0 2	T	+ id #	Reduzieren durch $E \rightarrow T$
0 1	E	+ id #	Verschieben
0 1 6	E +	id #	Verschieben
0 1 6 5	E + id	#	Reduzieren durch $F \rightarrow \mathbf{id}$
0 1 6 3	E + F	#	Reduzieren durch $T \rightarrow F$
0 1 6 9	E + T	#	Reduzieren durch $E \rightarrow E + T$
0 1	E	#	Akzeptieren

Beispiel: Übergangstabelle eines LR(0)-Automaten

	a	b	c	#	A	B	S
0	4	-	-	-	2	3	1
1	-	-	-	#	-	-	-
2	4	-	5	-	-	6	-
3	+4	+4	+4	-	+4	+4	+4
4	4	8	-	-	7	3	-
5	+2	+2	+2	-	+2	+2	+2
6	+3	+3	+3	-	+3	+3	+3
7	4	9	-	-	-	6	-
8	+6	+6	+6	-	+6	+6	+6
9	+5	+5	+5	-	+5	+5	+5

- 1: $Z \rightarrow S$ 2: $S \rightarrow Ac$
 3: $A \rightarrow AB$ 4: $A \rightarrow B$
 5: $B \rightarrow aAb$ 6: $B \rightarrow ab$

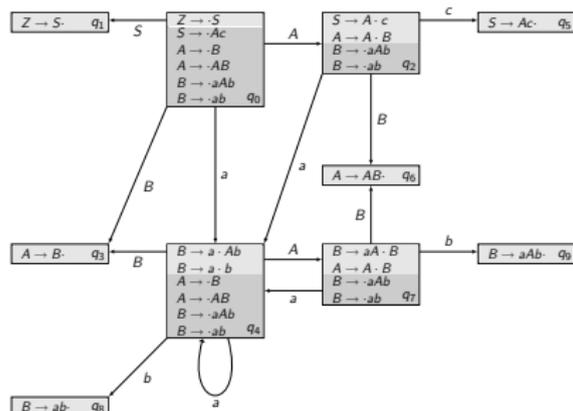


-	Fehler
+r	Reduziere mit Regel r
q	Schifte, neuer Zustand: q
#	HALT

Beispiel: Ablauf eines LR(0)-Automaten

Keller	Eingabe	Bemerkungen
0	aabbc#	schifte(a)
04	abbc#	schifte(a)
044	bbc#	schifte(b)
0448	bc#	reduziere($B \rightarrow ab$)
04	bc#	schifte(B)
043	bc#	reduziere($A \rightarrow B$)
04	bc#	schifte(A)
047	bc#	schifte(b)
0479	c#	reduziere($B \rightarrow aAb$)
0	c#	schifte(B)
03	c#	reduziere($A \rightarrow B$)
0	c#	schifte(A)
02	c#	schifte(c)
025	#	reduziere($S \rightarrow Ac$)
0	#	schifte(S)
01	#	HALT

- 1: $Z \rightarrow S$ 2: $S \rightarrow Ac$
 3: $A \rightarrow AB$ 4: $A \rightarrow B$
 5: $B \rightarrow aAb$ 6: $B \rightarrow ab$



rot: aktuelle Symbole bzw. zu reduzierender Kellerteil
 blau: Reduktionssymbol mit Schift im nächsten Schritt

Situationen (Items)

Ziel: bei Prüfung der Anwendbarkeit von Regeln $sqx \rightarrow s'q'x'$
Kellerinhalt und Zustand sq mit **inem** Zustandssymbol codieren
(Prüfung mehrerer Einträge im Keller vermeiden)

Lösungsidee:

- bei LL und LR ist s rechte bzw. linke Seite einer Produktion $X \rightarrow x_1 \dots x_n$
- Übergänge $tqt \rightarrow q$ (bei LL) bzw. $qt \rightarrow tq$ (bei LR) sind nur zulässig, wenn in der Produktion ein Terminalzeichen t ansteht, $x_1 \dots x_n = x'tx'''$, wobei $x'' := tx'''$
- also: ersetze sqx durch **Situation** $[X \rightarrow x' \cdot x''; x]$, die durch den Punkt anzeigt, wie weit die Produktion abgearbeitet ist.
- Situationen $[X \rightarrow \cdot x''; x]$ oder $[X \rightarrow x' \cdot; x]$ sind erlaubt und notwendig.
- Verwende Situationen als Zustände **und** als Kellersymbole.
- Situationen heißen englisch *items*.

Berechnung von CLOSURE

```
set<item> closure(item I)
{
  J = {I};
  do {
    changed = false;
    foreach((A → α · Bβ) ∈ J) {
      foreach((B → γ) ∈ G) {
        if((B → ·γ) ∉ J) {
          J = J ∪ {B → ·γ};
          changed = true;
        }
      }
    }
  } while (changed);

  return J;
}
```

Beispiel: closure

$$E' \rightarrow E$$

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

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

$$F \rightarrow (E) \mid \mathbf{id}$$

Berechnung von $\text{closure}(E' \rightarrow \cdot E)$ ergibt:

$$E' \rightarrow \cdot E$$

$$E \rightarrow \cdot E + T$$

$$E \rightarrow \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot \mathbf{id}$$

GOTO

$$\text{GOTO}(I, X) = \text{closure}(\{A \rightarrow \alpha X \cdot \beta \mid A \rightarrow \alpha \cdot X \beta \in I\})$$

Beispiel: Ist $I = \{[E' \rightarrow E \cdot], [E \rightarrow E \cdot + T]\}$, dann enthält $\text{GOTO}(I, +)$ folgende Situationen:

$$E \rightarrow E + \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot \text{id}$$

Berechnung der kanonischen LR(0)-Situationsmengen

```
void items(grammar G) {  
    C = { closure( {[S' → ·S]} ) };  
    do {  
        changed = false;  
        foreach(set<item> I ∈ C) {  
            foreach(grammar_symbol X) {  
                if (GOTO(I, X) ≠ ∅ && GOTO(I, X) ∉ C) {  
                    C = C ∪ { closure(GOTO(I, X)) };  
                    changed = true;  
                }  
            }  
        }  
    } while(changed);  
}
```

LR(1)-Parseralgorithmus

```
push(0);
a = next_token();
while (true) {
    if (ACTION[top(), a] == shift t) {
        push(t);
        a = next_token();
    } else if (ACTION[top(), a] == reduce  $A \rightarrow \beta$ ) {
        for (i = 0; i <  $|\beta|$ ; ++i)
            pop();
        push(GOTO[top(), A]);
    } else if (ACTION[top(), a] == accept) {
        break;
    } else {
        report_error(); break;
    }
}
```

SLR(k)-Grammatiken

SLR(k)-Grammatik (simple LR(k)):

Eine Grammatik heißt SLR(k), wenn sie LR(0) ist oder die SLR(k)-Übergangsfunktion bezüglich des charakteristischen LR(0)-Automaten C , $q \in Q$, $t \in V$ und $v \in T^{\leq k-1}$:

$$f(q, tv) = \begin{cases} \text{SHIFT} & \text{wenn } [X \rightarrow \mu \cdot t\gamma] \in q \wedge \\ & (v \in \text{Anf}'_{k-1}(\gamma \text{Folge}_{k-1}(X)) \vee k < 2) \\ \text{RED}(X \rightarrow x) & \text{wenn } [X \rightarrow x \cdot] \in q \wedge \\ & k : tv \in \text{Folge}_k(X) \\ \text{HALT} & \text{wenn } t = \# \wedge [Z \rightarrow S \cdot] \in q \\ \text{FEHLER} & \text{sonst} \end{cases}$$

eindeutig ist.

rot: Unterschied LR(k) und SLR(k)

SLR(1)-Übergangsfunktion am Beispiel

$$E' \rightarrow E$$

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

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

$$F \rightarrow (E) \mid \mathbf{id}$$

$$I_1 : E' \rightarrow E \cdot$$

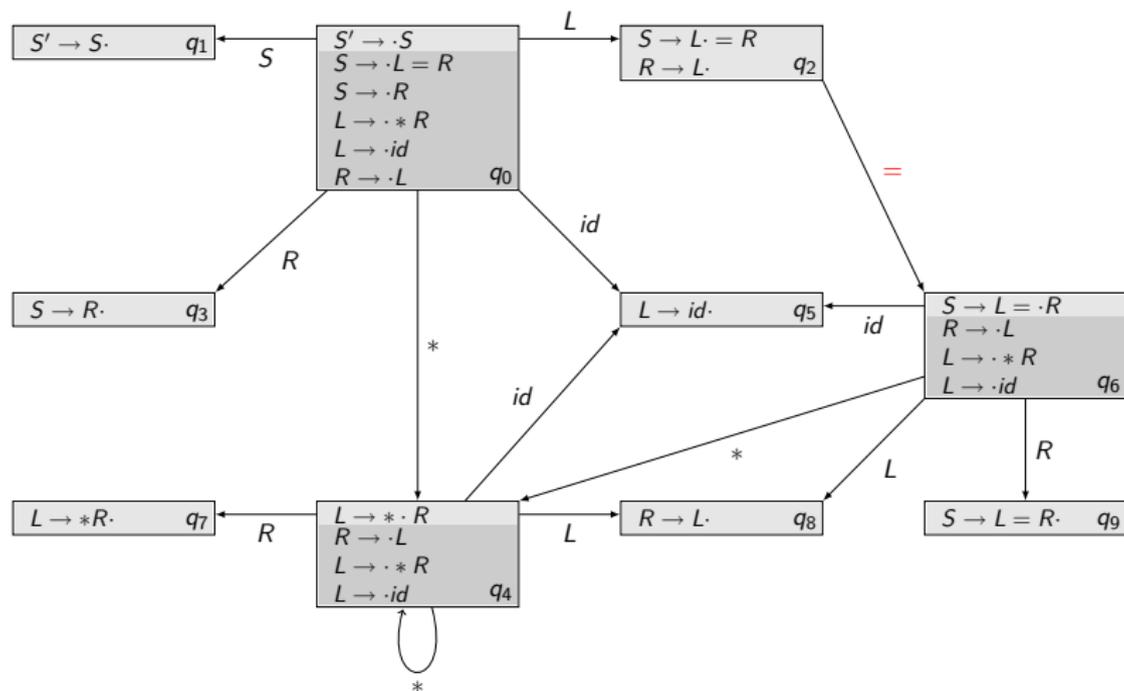
$$E \rightarrow E \cdot + T \mid T$$

Es gilt $\text{Folge}_1(E') = \{\#\}$. SLR(1)-Übergangsfunktion für I_1 ergibt:

$$f(I_1, +) = \text{SHIFT}$$

$$f(I_1, \#) = \text{HALT}$$

Beispiel für Konflikt bei SLR(1)

 $S \rightarrow L = R \mid R$
 $L \rightarrow * R \mid id$
 $R \rightarrow L$

 $\text{Folge}_1(S) = \{\#\}$
 $\text{Folge}_1(L) = \{=, \#\}$
 $\text{Folge}_1(R) = \{=, \#\}$

Motivation für LR/LALR (1/4)

LR-Parser erfassen alle deterministischen kf Grammatiken. Sie sind noch einigermaßen schnell berechenbar und es gibt ausgereifte Generatoren. Für den yacc-Generator produziert die Grammatik

```
%token IF ELSE THEN IDENT
```

```
%%
```

```
Statement: IfStatement | Expr ;
```

```
Expr: IDENT;
```

```
IfStatement: IF Expr THEN Statement |
```

```
IF Expr THEN Statement ELSE Statement ;
```

allerdings folgende Fehlermeldung

```
state 7 contains 1 shift /reduce conflict .
```

Motivation für LR/LALR (2/4)

Grammatik:

%token IF ELSE THEN IDENT

%%

Statement: IfStatement | Expr ;

Expr: IDENT;

IfStatement: IF Expr THEN Statement |

IF Expr THEN Statement ELSE Statement ;

Fehlermeldung:

state 7 contains 1 shift/reduce conflict.

...

state 7

IfStatement \rightarrow IF Expr THEN Statement . (rule 4)

IfStatement \rightarrow

IF Expr THEN Statement . ELSE Statement (rule 5)

ELSE shift, and go to state 8

ELSE [reduce using rule 4 (IfStatement)]

\$default reduce using rule 4 (IfStatement)

Motivation für LR/LALR (3/4)

Frage: Wie kann man die vorige Fehlermeldung verstehen?

- Zweck aller nachfolgenden Ausführungen ist nur darauf zu beziehen.
- Detailliertes Wissen über automatische Konstruktion vernachlässigbar
- Aber: Intuitive „Hand“konstruktion muss beherrscht werden.

Motivation für LR/LALR (4/4)

Lösungen:

- Sprache ändern: abschließendes „end“ einführen (zulässig?)
- Ausfaktorisieren:

```
%token IF ELSE THEN IDENT
```

```
%%
```

```
Statement: IfStatement | Expr ;
```

```
Expr: IDENT;
```

```
IfStatement:
```

```
    IF Expr THEN Statement |
```

```
    IF Expr THEN IfThenElseStat ELSE Statement ;
```

```
IfThenElseStat:
```

```
    IF Expr THEN IfThenElseStat ELSE IfThenElseStat |
```

```
    Expr ;
```

LR-Grammatiken

Ziel:

- alle deterministisch parsbaren kfG charakterisieren (LL(k) ist stark eingeschränkt)
- Rechtsableitung konstruieren

Endgültige Definition von D.E. Knuth 1966:

Eine kf-Grammatik heißt eine LR(k)-Grammatik, wenn für alle Paare von Ableitungen

$$\begin{array}{ll} Z \Rightarrow_R^* \mu A \omega \Rightarrow \mu \chi \omega & \mu \in V^*, \omega \in T^*, A \rightarrow \chi \in P \\ Z \Rightarrow_R^* \mu' B \omega' \Rightarrow \mu' \gamma \omega' & \mu' \in V^*, \omega' \in T^*, B \rightarrow \gamma \in P \end{array}$$

gilt:

$$(|\mu \chi| + k) : \mu \chi \omega = (|\mu \chi| + k) : \mu' \gamma \omega' \Rightarrow \mu = \mu', A = B, \chi = \gamma$$

Idee für deterministisch Machen des LR-Parsers

- Nimm alle verfügbare Information (Reduktionsklassen):
 - unendliche Vorschau auf die noch zu verarbeitende Eingabe
 - den gesamten Kellerinhalt

- Allerdings sind die **Reduktionsklassen nicht effektiv berechenbar**

Deshalb: Beschränkung auf k -Zeichen Vorschau und Erkenntnis, dass der zur jeweiligen Entscheidung nötige Kellerinhalt in einem bzw. beschränkt vielen Zuständen subsumierbar ist.

Dies führt zu k -Kellerklassen

- Mit Hilfe der k -Kellerklassen wird der charakteristische Automat definiert
- Aus dem charakteristischen Automaten ist die Übergangsfunktion des $LR(k)$ -Parsers ablesbar
- Aber: **Deterministisch machen geht nur für eingeschränkte Grammatikklassen**

Textmengen

$k : x$

$k : x = x\#$ falls $x = x_1 \dots x_m \wedge m < k$

$k : x = x_1 \dots x_k$ falls $x = x_1 \dots x_m \wedge m \geq k$

$\text{Anf}_k(x) = \{u \mid \exists y \in T^* : x \Rightarrow^* y \wedge u = k : y\}$

in der Literatur auch $\text{First}_k(x)$ genannt

$A \Rightarrow_{R'} \alpha$ gdw. $A \Rightarrow_R \alpha \wedge \nexists B \in N : A \Rightarrow_R B\alpha \Rightarrow \alpha$

$\text{Anf}'_k(x) = \{u \in \text{Anf}_k(x) \mid \exists y \in T^* : x \Rightarrow_{R'} uy\}$

in der Literatur auch $\text{EFF}_k(x)$ genannt (EFFective First)

$\text{Folge}_k(x) = \{u \mid \exists m, y \in V^* : Z \Rightarrow^* mxy \wedge u \in \text{Anf}_k(y)\}$

in der Literatur auch $\text{Follow}_k(x)$ genannt

Reduktionsklassen

Unter welchen Bedingungen soll der LR-Automat schiften / reduzieren?

maximal verfügbare Information:

Reduktionsklasse = $\{(Kellerinhalt, Eingaberest)\}$

R_0 für Schift und R_p für Produktionen $A_p \rightarrow y_p, p = 1, \dots, n$:

$$R_0 = \{(r'r, ss') \mid Z \Rightarrow_R^* r'As' \wedge A \Rightarrow_{R'} rs \wedge s \neq \varepsilon\}$$

$$R_p = \{(r'y_p, s) \mid Z \Rightarrow_R^* r'A_p s \wedge A_p \rightarrow y_p \in P\}$$

$R_i \cap R_j = \emptyset$ für $i \neq j$ bedeutet: Schiften bzw. Reduzieren mit Produktion p kann stets eindeutig entschieden werden. Jeder Satz besitzt eindeutige Rechtsableitung und eindeutigen Strukturbaum.

Grammatik ist eindeutig.

Leider: Eindeutigkeit von kf Grammatiken nicht entscheidbar, d.h.

$R_i \cap R_j = \emptyset$ nicht algorithmisch entscheidbar.

Problem ist der unbeschränkt lange Eingaberest.

Warum ist $\Rightarrow_{R'}$ nötig: Beispiel

$L(G) = \{cb, b\}$ und

$P = \{1: Z \rightarrow S, 2: S \rightarrow Ab, 3: A \rightarrow c, 4: A \rightarrow \varepsilon\}$

Es gilt $b \in \text{Anf}_1(S)$, aber $b \notin \text{Anf}'_1(S)$.

$R_0 = \{(\varepsilon, cb\#), (A, b\#), (\varepsilon, b\#)\},$

$R_1 = \{(S, \#)\},$

$R_2 = \{(Ab, \#)\},$

$R_3 = \{(c, b\#)\},$

$R_4 = \{(\varepsilon, b\#)\}$

Beispiel: Rechtsableitung $Z \Rightarrow_R S \Rightarrow_R Ab \Rightarrow_R b:$

ε auf A reduzieren, bevor das Zeichen b geschiftet wird.

Ohne Unterscheidung zwischen $\Rightarrow_{R'}$ und \Rightarrow_R : $(\varepsilon, b\#)$ gehört zu R_0 , also wird geschiftet aber Konfiguration $bq\#$ kann keiner Reduktionsklasse zugeordnet werden, d.h. **Sackgasse**

k -Kellerklassen

Idee: Beschränke den betrachteten Eingaberest auf $k \geq 0$ Zeichen:
 k -Kellerklasse K_p^k , $p = 0, \dots, n$ definiert durch

$$K_p^k = \{(r, k : s) \mid \exists (r, s) \in R_p\}$$

Vergleich mit LR-Definition: Mit

$$\begin{array}{ll} Z \Rightarrow_R^* \mu A \omega \Rightarrow \mu \chi \omega & \mu \in V^*, \omega \in T^*, A \rightarrow \chi \in P \\ Z \Rightarrow_R^* \mu' B \omega' \Rightarrow \mu' \gamma \omega' & \mu' \in V^*, \omega' \in T^*, B \rightarrow \gamma \in P \end{array}$$

gilt:

$$(|\mu \chi| + k) : \mu \chi \omega = (|\mu \chi| + k) : \mu' \gamma \omega' \Rightarrow \mu = \mu', A = B, \chi = \gamma$$

Satz: Eine Grammatik ist genau dann LR(k), wenn die k -Kellerklassen paarweise disjunkt sind. Eine Grammatik G ist genau dann deterministisch parsbar, wenn es ein $k \geq 0$ gibt, so dass G LR(k) ist.

Reguläre Erkenner für k -Kellerklassen

Satz(Büchi): Alle Kellerklassen K_p^k für $p = 0, \dots, n$ sind regulär.

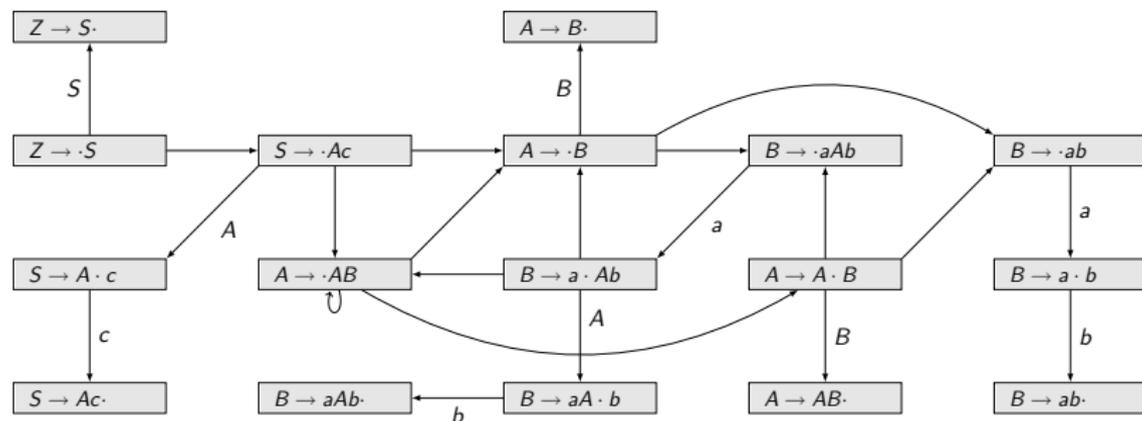
Nachweis durch Angabe einer rechtslinearen Grammatik G_p^k für jede Kellerklasse.

Man kann daraus einen nichtdeterministischen endlichen Automaten C ableiten, der die regulär erkennbaren Wortanfänge der durch die (original) Grammatik G gegebenen Sprache akzeptiert. Dieser hat:

- Situationen als Zustände $\in Q$
- $[Z \rightarrow \cdot S; \#]$ als Startzustand (Z nichtrekursives Start-NT von G)
- die Übergangsfunktion f
 - $f([X \rightarrow x \cdot vy; \omega], v) = [X \rightarrow xv \cdot y; \omega]$ mit $v \in V$
 - $f([X \rightarrow x \cdot By; \omega], \varepsilon) = [B \rightarrow \cdot b; \tau]$ mit $B \rightarrow b \in P$,
 $\tau \in \text{Anf}_k(\gamma\omega)$
- Endzustände sind belanglos für das weitere Vorgehen

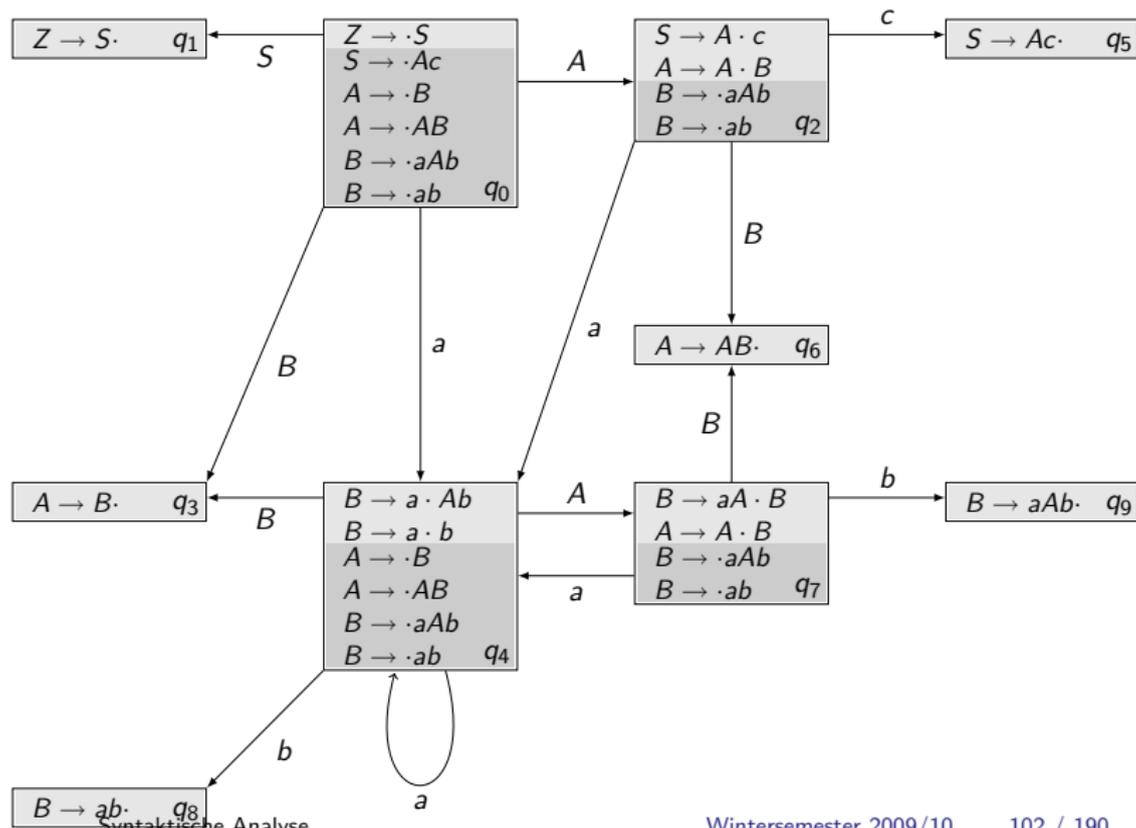
Nichtdeterministischer, endlicher Automat C

$Z \rightarrow S, S \rightarrow Ac, A \rightarrow AB, A \rightarrow B, B \rightarrow aAb, B \rightarrow ab$



Charakteristischer Automat für $k = 0$

$Z \rightarrow S, S \rightarrow Ac, A \rightarrow AB, A \rightarrow B, B \rightarrow aAb, B \rightarrow ab$



Hilfskonstruktionen mit dem charakteristischen Automaten

Der nichtdeterministische endliche Automat C kann mit den bekannten Verfahren deterministisch gemacht werden. Ergebnis:

charakteristischer Automat

- Definiere $\text{basis}(q, t, \omega)$ als Menge von Situationen die mit Einlesen des (Nicht-)Terminals t aus Situationen von q ausgehend erreichbar sind:

$$\text{basis}(q, t, \omega) = \{[X \rightarrow \mu t \cdot \gamma; \omega] \mid [X \rightarrow \mu \cdot t\gamma; \omega] \in q\}$$

- Definiere $\text{next}(q, t, \omega)$ als die transitive Hülle von $\text{basis}(q, t, \omega)$

$$\text{next}(q, t, \omega) = H(\text{basis}(q, t, \omega)), \text{ wenn } \text{basis}(q, t, \omega) \neq \emptyset$$

$$\text{mit } H(M) = M \cup \{[B \rightarrow \cdot \beta; \tau] \mid \exists [X \rightarrow \mu \cdot B\gamma; \xi] \in H(M) : B \rightarrow \beta \in P \wedge \tau \in \text{Anf}_k(\gamma\xi)\}$$

Hinweis: $H(M)$ erweitert Situationen (oder Situationsmengen) auf die Zustände des charakteristischen Automaten.

Übergangsfunktion des LR(k)-Automaten

Für einen charakteristischen Automaten C , $q \in Q$, $t \in V$ und $v \in T^{\leq k-1}$:

$$f(q, tv) = \begin{cases} \text{next}(q, t, \omega) & \text{wenn } [X \rightarrow \mu \cdot t\gamma; \omega] \in q \wedge \\ & (v \in \text{Anf}'_{k-1}(\gamma\omega) \vee k < 2) \\ \text{RED}(X \rightarrow x) & \text{wenn } [X \rightarrow x \cdot; k : tv] \in q \\ \text{HALT} & \text{wenn } t = \# \wedge [Z \rightarrow S \cdot; \#] \in q \\ \text{FEHLER} & \text{sonst} \end{cases}$$

- Zustand q heißt **inadäquat**, wenn Übergangsfunktion $f(q, tv)$ für irgendein tv nicht eindeutig die Fälle eins, zwei und drei unterscheiden kann. Nur zwei Möglichkeiten existieren:
 - Schift-Reduktionskonflikt: schiften und reduzieren möglich
 - Reduktions-Reduktionskonflikt: Reduktion mit zwei verschiedenen Produktionen möglich
- Eine kfG G ist genau dann LR(k)-Grammatik, wenn der charakteristische Automat keine inadäquaten Zustände besitzt.

Ablauf des LR(k)-Automaten

- Startzustand des Automaten ist: $H(\{[Z \rightarrow \cdot S; \#]\})$
- In jedem Schritt wird die Übergangsfunktion f angewandt:
 - Beim Schiften wird statt des konkreten Symbols ein Zustand des charakteristischen Automaten auf den Keller gelegt. Nur durch Schiften von Terminalen wird die Eingabe verkürzt; Schiften von Nichtterminalen läßt die Eingabe unverändert.
 - RED($X \rightarrow x$) bedeutet Reduzieren mit Produktion $X \rightarrow x$, d.h.:
 - Löschen von $|x|$ Kellereinträgen.
 - Oberster Kellereintrag wird aktueller Zustand.
 - Die Übergangsfunktion $f(q, Xtv)$ wird angewandt. Dies führt sofort zu einem Schift des Nichtterminals X und kann auch in einem Schritt mit dem Kellerlöschen behandelt werden.
 - Bei HALT wird die Eingabe akzeptiert.
 - Bei FEHLER ist die Eingabe nicht akzeptabel.

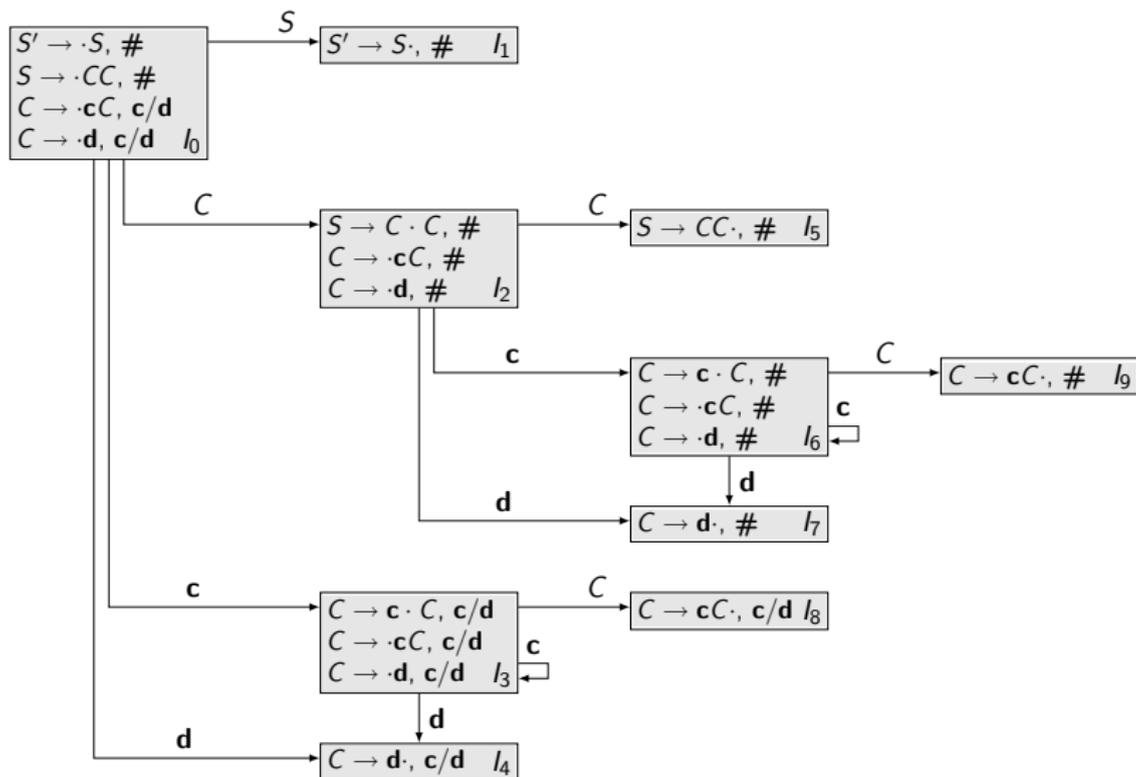
Beispiel: Grammatik G'

$$S' \rightarrow S$$

$$S \rightarrow C C$$

$$C \rightarrow \mathbf{c} C \mid \mathbf{d}$$

GOTO-Graph für die Grammatik G'



Beispiel: Ein LR(2)-Kellerautomat

1: $Z \rightarrow X$, 2: $X \rightarrow Y$, 3: $X \rightarrow bYa$, 4: $Y \rightarrow c$, 5: $Y \rightarrow ca$

Keller	Eingabe	Bemerkungen
0	b caa#	schifte(<i>b</i>)
03	c aa#	schifte(<i>c</i>)
031	a a#	schifte(<i>a</i>)
0312	a #	reduziere($Y \rightarrow ca$)
03	a #	schifte(<i>Y</i>)
035	a #	schifte(<i>a</i>)
0358	#	reduziere($X \rightarrow bYa$)
0	#	schifte(<i>X</i>)
01	#	HALT

rot: aktuelle Symbole bzw. zu
reduzierender Kellerteil

grün: Vorschau

blau: Reduktionssymbol mit
Schift im nächsten Schritt

	bc	c#	ca	a#	aa	#	X#	Y#	Ya
0	3	7	7	-	-	-	9	8	-
1	-	-	-	+4	2	-	-	-	-
2	-	-	-	+5	-	-	-	-	-
3	-	-	1	-	-	-	-	-	4
4	-	-	-	5	-	-	-	-	-
5	-	-	-	-	-	+3	-	-	-
6	-	-	-	-	-	+5	-	-	-
7	-	-	-	6	-	+4	-	-	-
8	-	-	-	-	-	+2	-	-	-
9	-	-	-	-	-	#	-	-	-

-	Fehler
+r	Reduziere mit Regel r
q	Schifte, neuer Zustand: q
#	HALT

LR-Grammatiken: Probleme

Probleme:

- Wegen Unterscheidung verschiedener Rechtskontexte hat der LR(k)-Automat bereits für $k = 1$ sehr viele Zustände (verglichen mit LR(0))
 - Beispiel: für Ausdrucksgrammatik mit „+“ 17 statt 9 für $k = 0$.
- LR(1)-Automat nur automatisch generierbar
- sehr große Tabellen
- Test der Eigenschaft nur durch Konstruktion möglich

In der Praxis Beschränkung auf Unterklassen von LR(1) mit LR(0)-Zustandsmenge.

SLR(k)-Grammatiken (1/2)

SLR(k)-Grammatik (simple LR(k)):

Eine Grammatik heißt SLR(k), wenn sie LR(0) ist oder die modifizierte Übergangsfunktion bezüglich des charakteristischen LR(0)-Automaten C , $q \in Q$, $t \in V$ und $v \in T^{\leq k-1}$:

$$f(q, tv) = \begin{cases} \text{next}(q, t, \omega) & \text{wenn } [X \rightarrow \mu \cdot t\gamma] \in q \wedge \\ & (v \in \text{Anf}'_{k-1}(\gamma \text{Folge}_{k-1}(X)) \vee k < 2) \\ \text{RED}(X \rightarrow x) & \text{wenn } [X \rightarrow x \cdot] \in q \wedge \\ & k : tv \in \text{Folge}_k(X) \\ \text{HALT} & \text{wenn } t = \# \wedge [Z \rightarrow S \cdot] \in q \\ \text{FEHLER} & \text{sonst} \end{cases}$$

keine inadäquaten Zustände liefert.

rot: Unterschied LR(k) und SLR(k)

SLR(k)-Grammatiken (2/2)

- **Achtung:** $\text{next}(q, t, \varepsilon)$ operiert auf einem charakteristischen LR(0)-Automaten, also ohne Rechtskontexte.
- **Problem:** SLR(k) genügt zwar für Ausdrucksgrammatiken, berücksichtigt aber Linkskontext zu wenig (da die Vorausschau ohne Kenntnis der schon verarbeiteten mithin links stehenden Zeichen bestimmt wird).

LALR(k)-Grammatiken

Sei $\text{kern}(q) = \{[X \rightarrow \mu \cdot \gamma] \mid [X \rightarrow \mu \cdot \gamma; \Omega_i] \in q\}$.

Eine Grammatik heißt **LALR(k) (look ahead LR(k))**, wenn es keine inadäquaten Zustände gibt, falls man im LR(k)-Automaten alle Zustände q, q' mit $\text{kern}(q) = \text{kern}(q')$ zusammenlegt.

- **Satz:** Jeder SLR(k)- oder LALR(k)-Automat hat die gleiche Anzahl von Zuständen wie der LR(0)-Automat zur gleichen Grammatik.
- Der Unterschied der parsbaren Sprachen zwischen LALR(k) und LR(k) ist praktisch unerheblich.
- Alle verbreiteten LR-Parsergeneratoren (yacc, pgs, lalr, bison, ...) konstruieren LALR(1)-Automaten.

LALR(k)-Übergangsfunktion

Für einen charakteristischen Automaten C' in dem alle Zustände mit gleichem Kern zusammengelegt sind, $q \in Q$, $t \in V$ und $v \in T^{\leq k-1}$:

$$f(q, tv) = \begin{cases} \text{next}(q, t, \omega) & \text{wenn } [X \rightarrow \mu \cdot t\gamma; \omega] \in q \wedge \\ & (v \in \text{Anf}'_{k-1}(\gamma\omega) \vee k < 2) \\ \text{RED}(X \rightarrow x) & \text{wenn } [X \rightarrow x.; k : tv] \in q \\ \text{HALT} & \text{wenn } t = \# \wedge [Z \rightarrow S.; \#] \in q \\ \text{FEHLER} & \text{sonst} \end{cases}$$

erhalten wir die LALR(k)-Übergangsfunktion.

- **rot:** Unterschied LR(k) und LALR(k)
Einziger Unterschied: Zusammenlegen der Zustände „modulo Kern“
- Im Unterschied zu SLR(k) benutzt LALR(k) einen „echten“ Rechtskontext von $X \rightarrow x$, der schärfer trennt als Folge $_k(X)$.

nicht SLR, aber LALR

1: $Z \rightarrow A$

2: $A \rightarrow aBb$

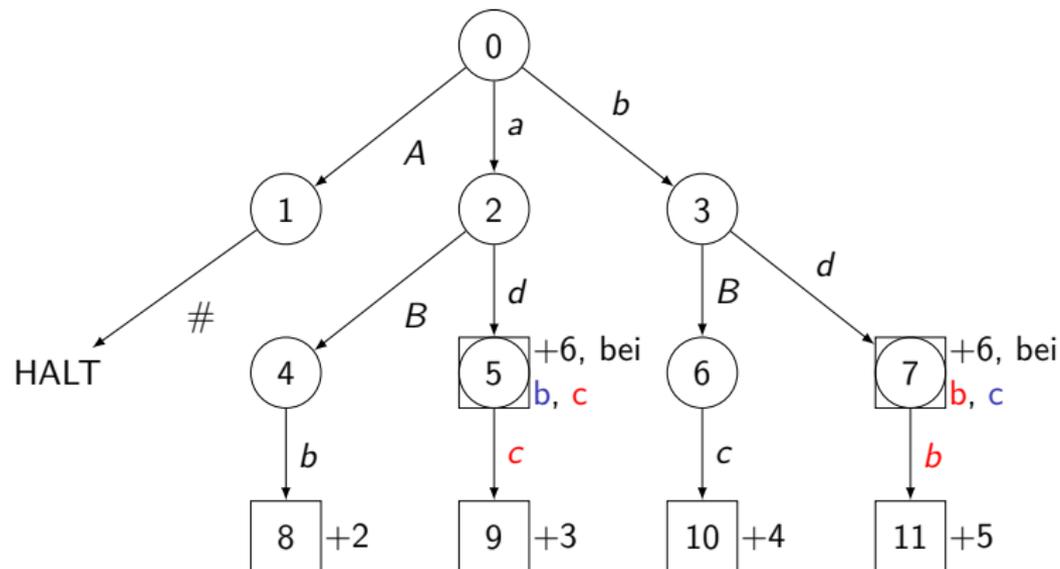
3: $A \rightarrow adc$

4: $A \rightarrow bBc$

5: $A \rightarrow bdb$

6: $B \rightarrow d$

$\text{Folge}_1(B) = \{b, c\}$



praktisches Beispiel (Algo 60): Unterscheide ;id: und [id:

Behebung inadäquater Zustände

- Parsergenerator:
 - Schift-Reduzier-Konflikt: Schiften wird bevorzugt
 - Reduzier-Reduzier-Konflikt: Zuerst spezifizierte Reduktion wird häufig bevorzugt. **Meist fehlerbehaftet.**
- Maßnahmen:
 - Automatische Konfliktauflösung auf Korrektheit prüfen!
 - „Faktorisieren“ gemeinsamer Produktionsteile (siehe if-then-else)
 - Verwendung von Präzedenzen (z.B. bei bison)
 - Vergrößern der erkannten Sprache, nachher Einschränken mit semantischer Analyse
 - An Beispielen lernen
 - Mächtigeren Generator (höheres k , LR statt LALR) verwenden, allerdings in der Regel nicht hilfreich

Achtung: Es gibt keine allgemeingültigen Verfahren

1: $Z \rightarrow A$

2: $A \rightarrow A + T$

3: $A \rightarrow T$

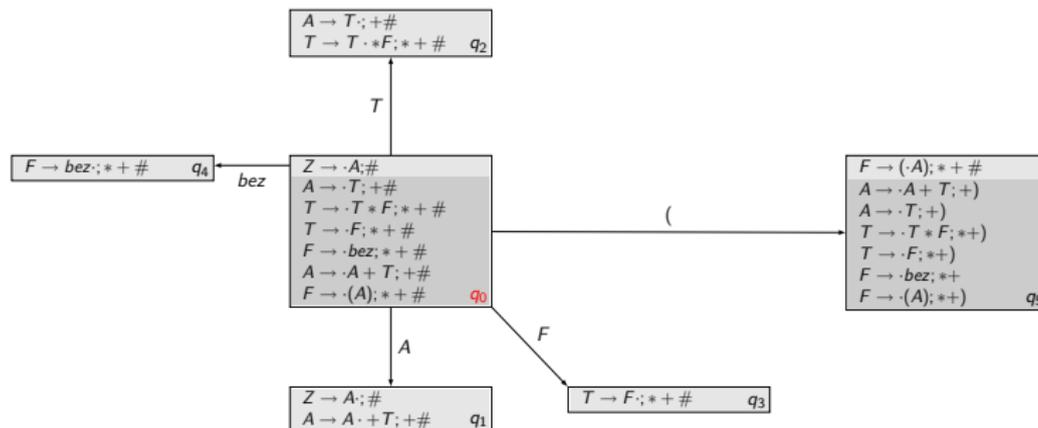
4: $T \rightarrow T * F$

5: $T \rightarrow F$

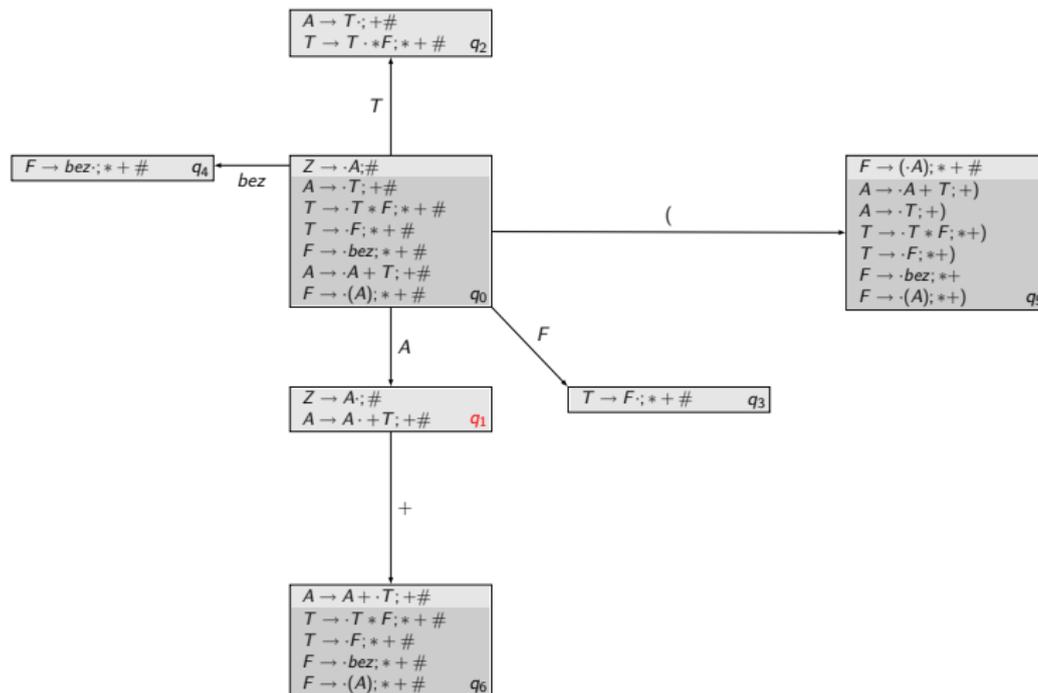
6: $F \rightarrow bez$

7: $F \rightarrow (A)$

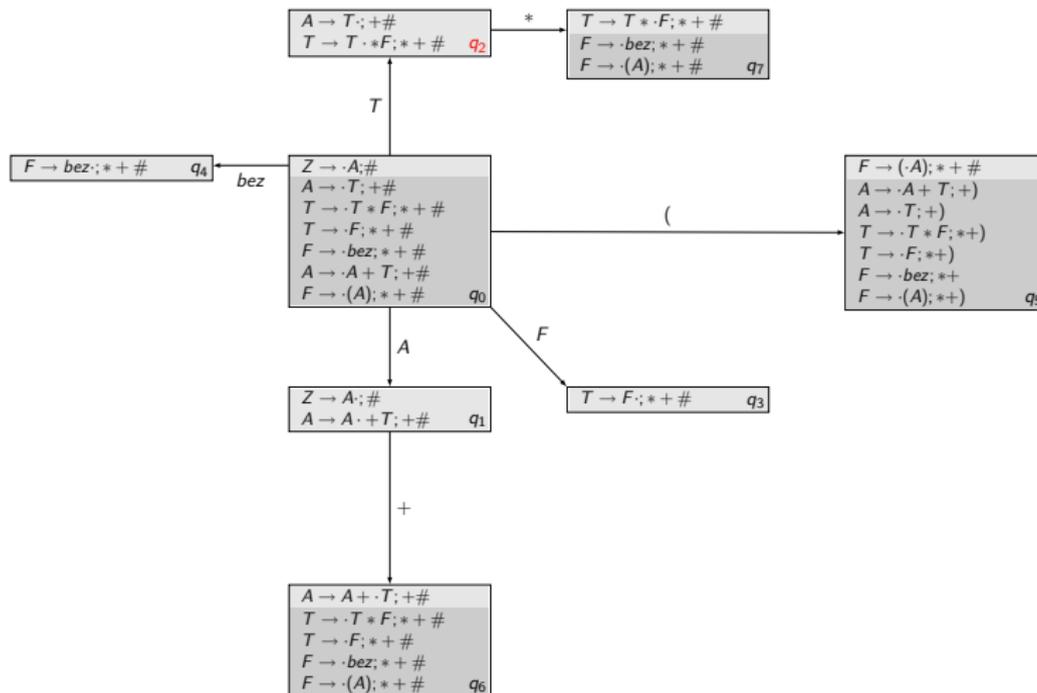
LALR(1)-Konstruktion am Beispiel



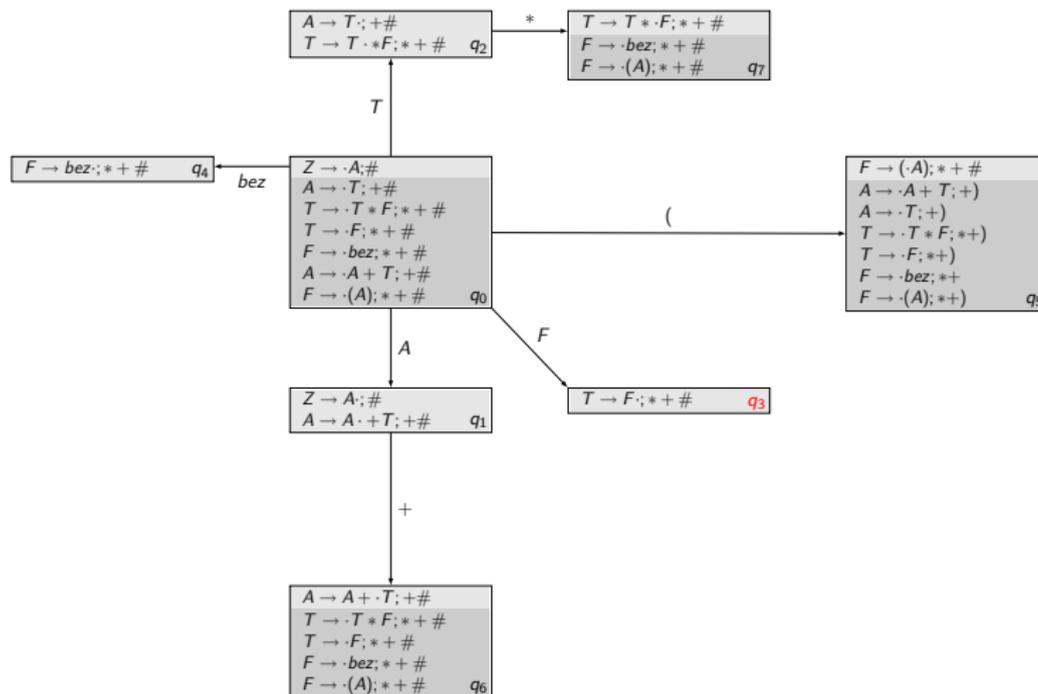
LALR(1)-Konstruktion am Beispiel



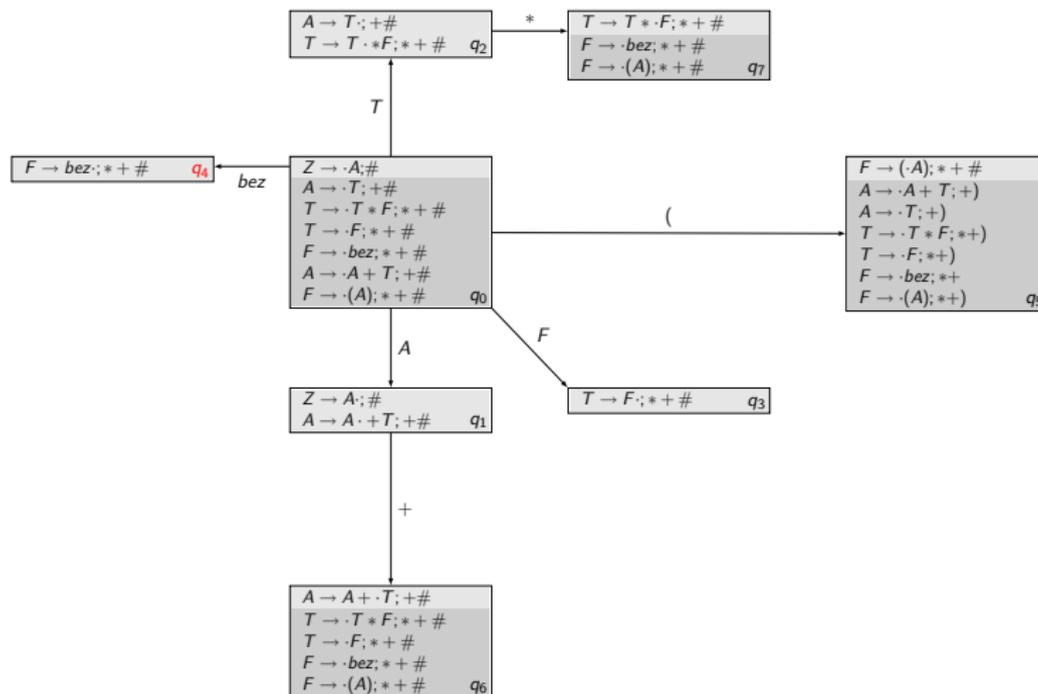
LALR(1)-Konstruktion am Beispiel



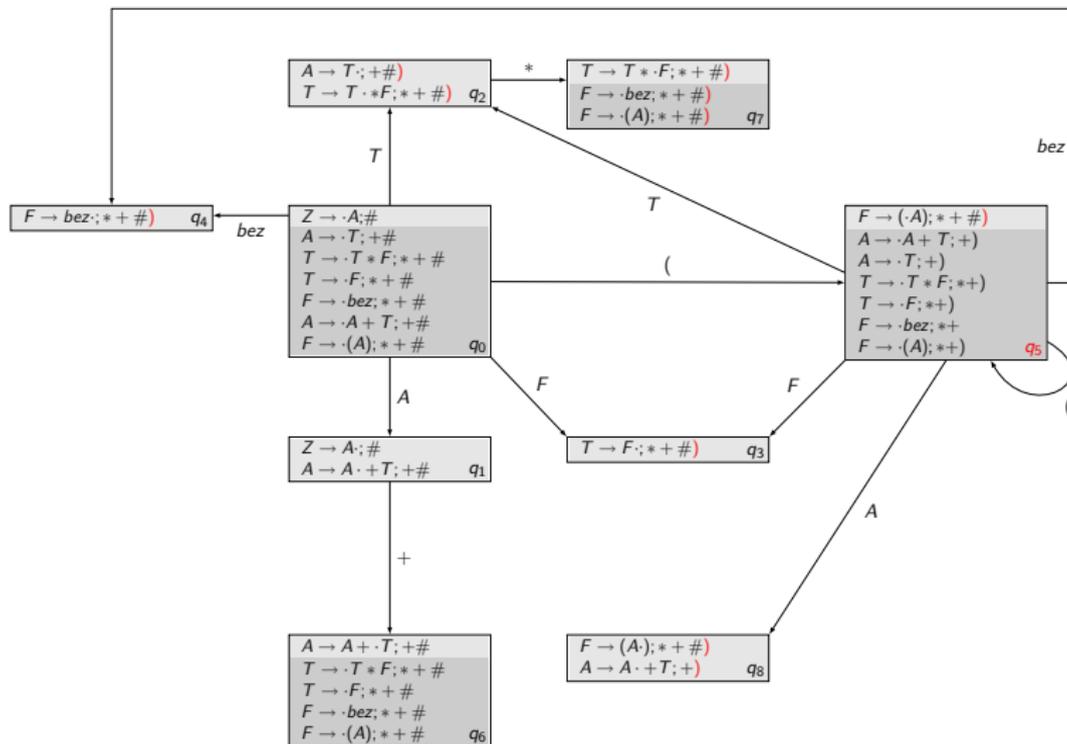
LALR(1)-Konstruktion am Beispiel



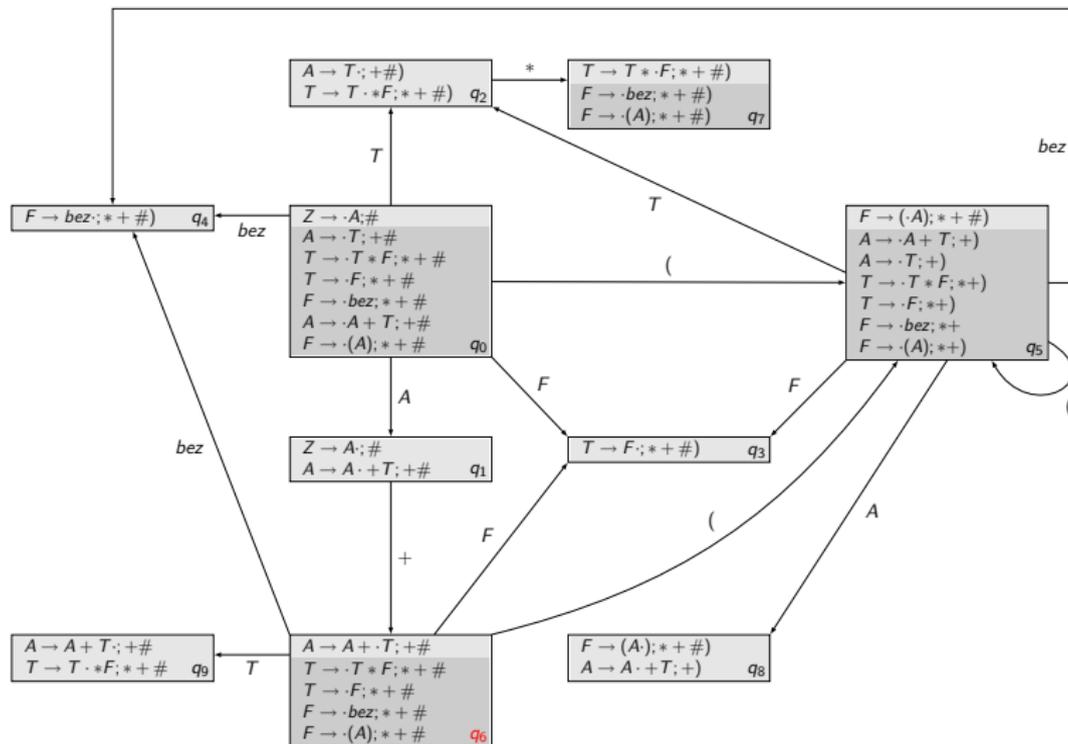
LALR(1)-Konstruktion am Beispiel



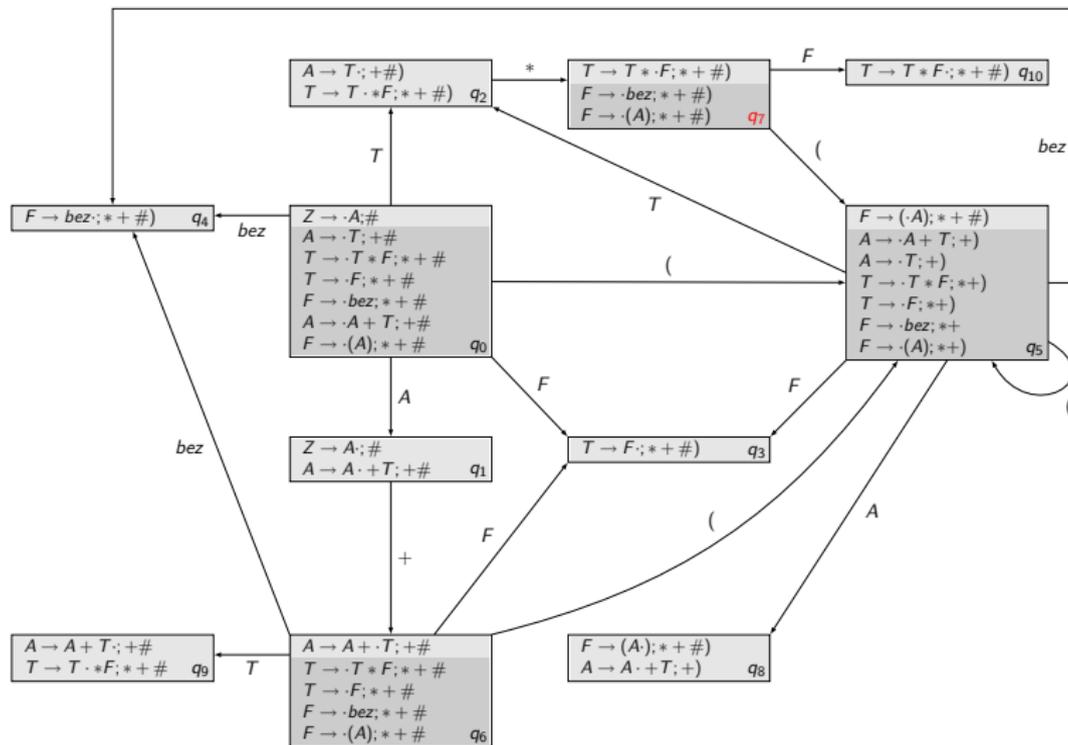
LALR(1)-Konstruktion am Beispiel



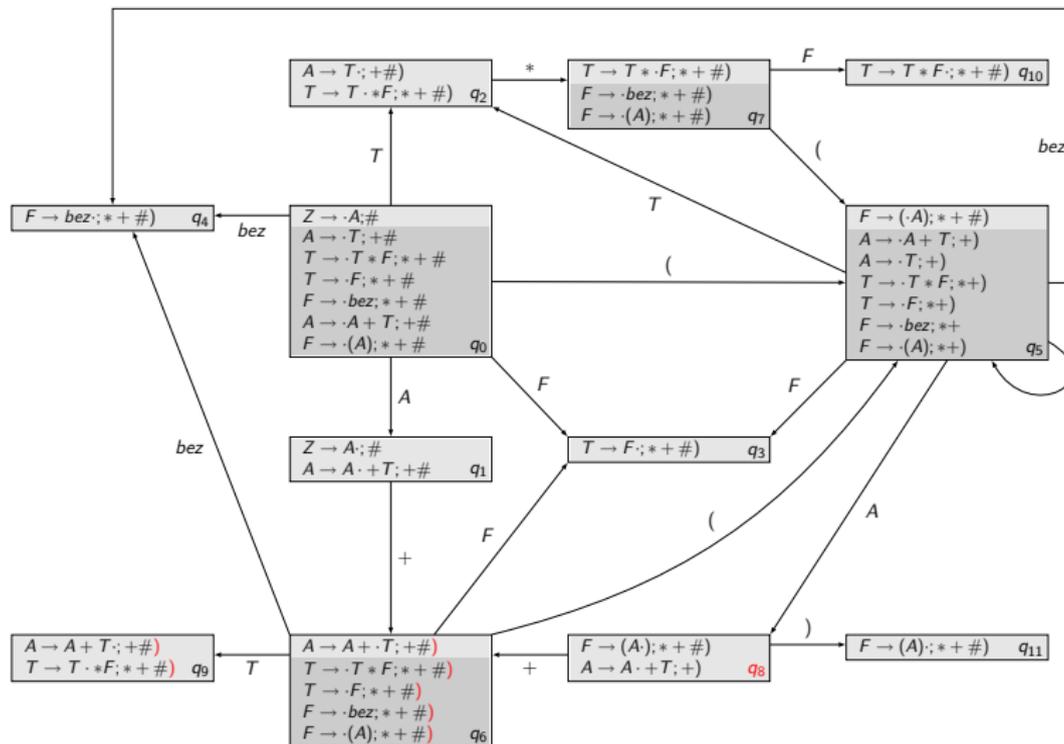
LALR(1)-Konstruktion am Beispiel



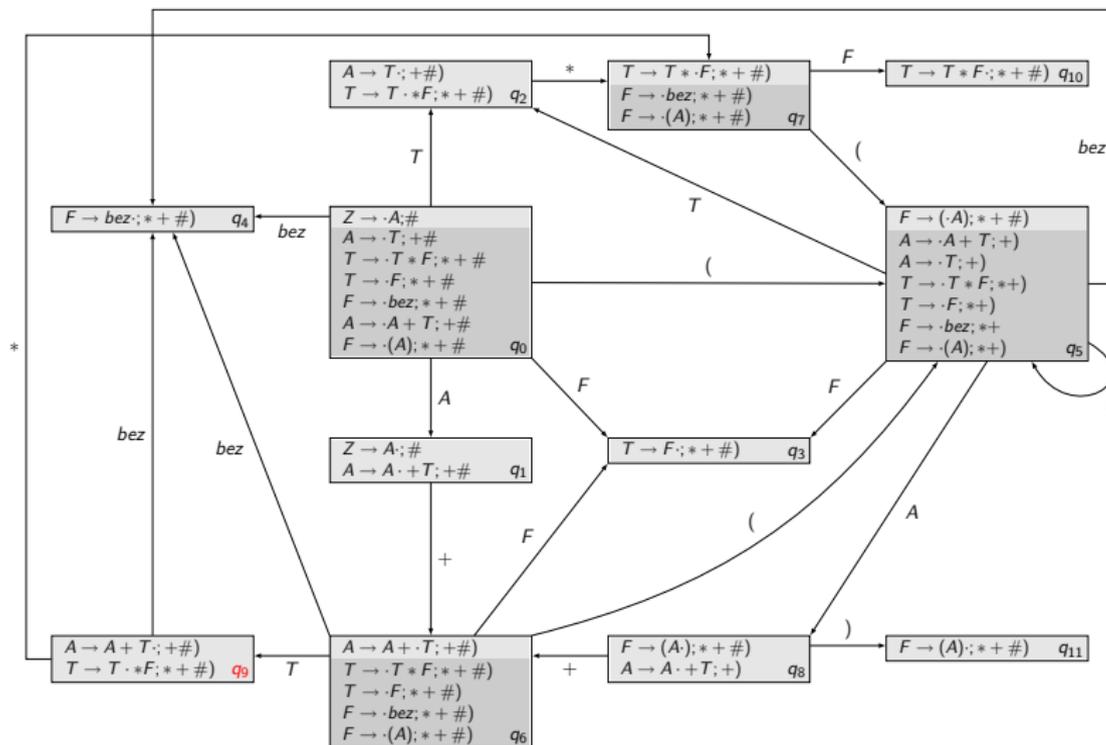
LALR(1)-Konstruktion am Beispiel



LALR(1)-Konstruktion am Beispiel



LALR(1)-Konstruktion am Beispiel



- q_0 $f(q_0, A) = q_1, f(q_0, T) = q_2, f(q_0, F) = q_3,$
 $f(q_0, bez) = q_4, f(q_0, () = q_5$
- q_1 $f(q_1, \#) = \text{HALT}, f(q_1, +) = q_6$
- q_2 $f(q_2, +\#)) = \text{Reduziere}(A \rightarrow T) f(q_2, *) = q_7$
- q_3 $f(q_3, * + \#)) = \text{Reduziere}(T \rightarrow F)$
- q_4 $f(q_4, * + \#)) = \text{Reduziere}(F \rightarrow bez)$
- q_5 $f(q_5, A) = q_8, f(q_5, T) = q_2, f(q_5, F) = q_3,$
 $f(q_5, bez) = q_4, f(q_5, () = q_5$
- q_6 $f(q_6, T) = q_9, f(q_6, F) = q_3, f(q_6, bez) = q_4,$
 $f(q_6, () = q_5$
- q_7 $f(q_7, F) = q_{10}, f(q_7, bez) = q_4, f(q_7, () = q_5$
- q_8 $f(q_8,)) = q_{11}, f(q_8, +) = q_6$
- q_9 $f(q_9, +\#)) = \text{Reduziere}(A \rightarrow A + T), f(q_9, *) = q_7$
- q_{10} $f(q_{10}, * + \#)) = \text{Reduziere}(T \rightarrow T * F)$
- q_{11} $f(q_{11}, * + \#)) = \text{Reduziere}(F \rightarrow (A))$

	bez	()	+	*	#	A	T	F
0	4	5	-	-	-	-	1	2	3
1	-	-	-	6	-	*			
2	-	-	+3	+3	7	+3			
3	+5	+5	+5	+5	+5	+5			
4	+6	+6	+6	+6	+6	+6			
5	4	5	-	-	-	-	8	2	3
6	4	5	-	-	-	-		9	3
7	4	5	-	-	-	-			10
8	-	-	11	6	-	-			
9	-	-	+2	+2	7	+2			
10	+4	+4	+4	+4	+4	+4			
11	+7	+7	+7	+7	+7	+7			

- Fehler

i Schiffe, neuer Zustand: i

+r Reduziere mit Regel r

„ “ Don't care

Fakten zu LR

- Automat $qt \rightarrow tq, t \in T, x_1 \dots x_n q \rightarrow Xq, X \rightarrow x_1 \dots x_n \in P$ deterministisch machen
 - dazu Reduktionsklassen bestimmen, Kellerklassen ableiten
 - Kellerklassen sind regulär, Produktionen mit Situationen als Nichtterminale
- charakteristischen Automaten herleiten:
nicht-deterministischen Automaten in deterministischen überführen (Hüllenbildung)
- LR(k)-Automat zu groß: SLR(k)-Klasse (Vorschau: Folge $_k(A)$) ungenügend
- LALR(k): Zustände des LR(k)-Automaten mit gleichem Kern verschmelzen
- SLR- und LALR-Automat hat gleiche Zustandsmenge wie LR(0)-Automat

Fragen zu LR

- Wie lautet die Übergangsfunktion?
- Wie bestimmt man Zustände und Übergänge des LALR(1)-Automaten von Hand?
- Wie korrigiert man Grammatik, wenn es inadäquate Zustände gibt?

Generatoren für die syntaktische Analyse

- Parsergeneratoren ermöglichen eine kompakte Spezifikation syntaktischen Regeln und Aktionen beim erkennen solcher Regeln (eine Domain Specific Language).
- Eingabe: Ein Tokenstrom
- Ausgabe: Bei erkannten Regeln lassen sich beliebige Aktionen durchführen, die in der Beschreibungssprache spezifiziert werden. Mögliche Anbindungen an einen Compiler:
 - Aufbau des AST als Datenstruktur.
 - unmittelbare semantische Analyse und Codeausgabe (Single-Pass-Compiler).
 - Erzeugen von Zwischensprachen.

Generatoren Kriterien

Heute existieren sehr viele unterschiedliche Parsergeneratoren. Die Liste auf

http://wikipedia.org/wiki/Comparison_of_parser_generators

hat zur Zeit 101 Einträge!

Bei der Auswahl eines Generators hilft es sich an den folgenden Kriterien zu orientieren:

- Parsingtechnik: LALR(1), LL(k), GLR, Packrat
- Programmiersprache(n) in denen der Parser erzeugt werden kann
- Anbindung an Werkzeuge zur lexikalischen Analyse und attributierter Grammatiken
- Einbettung in Entwicklungsumgebungen

Generatoren (eine Auswahl)

- *Yacc*: Erzeugen LALR(1)-Parser, bekannt aus dem Unix-/C-Umfeld.
- *Bison*: Erweiterte und Verbesserte Open-Source Variante von *Yacc*.
- Weitere LALR(1)-Parsergeneratoren: *SableCC*, *cup*
- Einer bekannter Generator im Java-Umfeld ist *antlr*. Dieser erzeugt trotz seines Namens LL(k) Generatoren. Er kann auch Parser in anderen Sprachen als Java erzeugen.
- Weitere LL(k)-Generatoren: *JavaCC* (Java), *Coco/R* (Java, C#, weitere), *happy* (haskell), *ocamlyacc* (ocaml)

Parser für XML

```
<?xml version='1.0'?>  
<!-- my personal books -->  
<books>  
  <book name='Goedel, Escher, Bach' />  
</books>
```

input → *prolog element | element*
prolog → **<? xml version = string ? >**
element → *simple_tag | start_end_tags*
start_end_tags → **< name attributes > content < / name >**
simple_tag → **< name attributes / >**
attributes → ϵ | *attribute attributes*
attribute → **name = string**
content → ϵ | *element content*

Parser für XML

```
<?xml version='1.0'?>
```

```
<!-- my personal books -->
```

```
<books>
```

```
  <book name='Goedel, Escher, Bach' />
```

```
</books>
```

input → *prolog element* | *element*

prolog → <? **xml version = string** ? >

element → *simple_tag* | *start_end_tags*

start_end_tags → < **name attributes** > *content* < / **name** >

simple_tag → < **name attributes** / >

attributes → ϵ | *attribute attributes*

attribute → **name = string**

content → ϵ | *element content*

Der erste Teil einer Bison Spezifikation besteht aus Deklarationen:

- Der Typ eines Tokens auf dem Stack wird als C-union mit %union angegeben.
- Die Tokenbezeichner werden mit Hilfe von %token festgelegt
- Vor einem Token kann dessen Typ in spitzen Klammern angegeben werden.
- Nach dem Token kann eine optionale Umschreibung als String angegeben werden.
- Alternativ kann man auch %left und %right angeben um eine Präzedenz zur Auflösung von Konflikten bei mehrdeutigen Grammatiken anzugeben.

Deklarationen: Beispiel

```
/* token value definitions */  
%union {  
    const char *string;  
}
```

```
/* token definitions */  
%token <string> T_STRING  
%token <string> T_NAME  
%token T_XML T_VERSION  
%token T_LESS_QUESTIONMARK "<?"  
%token T_LESS_SLASH "</"  
%token T_SLASH_GREATER "/>"  
%token T_QUESTIONMASK_GREATER "?>"
```

Regelspezifikation

Eine grammatische Regel der Form:

$$\text{Kopf} \rightarrow \text{Rumpf}_1 \mid \text{Rumpf}_2 \mid \dots \mid \text{Rumpf}_n$$

wird folgendermaßen spezifiziert:

```
kopf : rumpf1 { /* semantische Aktion 1 */ }  
      | rumpf2 { /* semantische Aktion 2 */ }  
      ...  
      | rumpfn { /* semantische Aktion n */ }  
;
```

In einem Rumpf wird eine Folge von Terminalen und Nichtterminalen angegeben. Nichtterminale werden durch ihren Bezeichner angegeben. Terminale durch Bezeichner, 'x' bei einzelnen Zeichen, oder "xx" (falls xx als Umschreibung eines Tokens definiert wurde).

Beispiel: Statements

statement : if_statement | while_statement | for_statement
| expression_statement | compound_statement
;

if_statement : T_IF '(' expression ')' statement
| T_IF '(' expression ')' statement T_ELSE statement
;

while_statement: T_WHILE '(' expression ')' statement ;

for_statement: T_FOR '(' expression ';' expression ';' expression ')' statement ;

expression_statement: expression ';' ;

compound_statement: '{' statements '}' ;

statements: /* empty */ | statement statements ;

Beispiel: Expressions

```
expression: expression '+' expression
| expression '-' expression
| expression '*' expression
| expression '/' expression
| call_expression
| '(' expression ')'
| T_IDENT
| T_LITERAL
;
```

```
call_expression: expression '(' arguments ')';
```

```
arguments: /* empty */ | expression | multi_arguments ;
```

```
multi_arguments: expression
| expression ',' multi_arguments ;
```

Beachte: Mehrdeutigkeiten werden hier durch Spezifikation von Tokenpräzedenzen aufgelöst.

Beispiel: XML Parser

input: prolog element | element ;

prolog: "<?" T_XML T_VERSION '=' T_STRING ">?" ;

element: start_end_tags | simple_tag ;

start_end_tags: '<' T_NAME attributes '>'
content
"</" T_NAME '>'
;

attributes: /* empty */ | attributes attribute ;

attribute: T_NAME '=' T_STRING ;

simple_tag: '<' T_NAME attributes ">" ;

content: /* empty */ | content element ;

```
// Definition of a Java class
jClassDefinition[int mods]
  returns [JClassDeclaration self = null]
{
  String superClass = null;
  CClassType[] interfaces = CClassType.EMPTY;
  CParseClassContext context = new CParseClassContext();
  TokenReference sourceRef = buildTokenReference();
  JavadocComment javadoc = getJavadocComment();
  JavaStyleComment[] comments = getStatementComment();
}
:
...
```

```
...  
"class" ident:IDENT  
superClass = jSuperClassClause[]  
interfaces = jImplementsClause[]  
jClassBlock[context] // the body of the class  
{  
    self = new JClassDeclaration(sourceRef,  
    mods, ident.getText(),  
    superClass, interfaces,  
    context.getFields(),  
    context.getMethods(),  
    context.getInnerClasses(),  
    context.getBody(),  
    javadoc, comments);  
};
```

/ 19.8.1 Production from §8.1: Class Declaration */*

class_declaration:

modifiers CLASS_TK identifier **super** interfaces

{ create_class (\$1, \$3, \$4, \$5); }

class_body

{;}

| CLASS_TK identifier **super** interfaces

{ create_class (0, \$2, \$3, \$4); }

...

```
...
class_body
  {;}
| modifiers CLASS_TK error
  { yyerror ("Missing class name"); RECOVER; }
| CLASS_TK error
  { yyerror ("Missing class name"); RECOVER; }
| CLASS_TK identifier error
  { if (!ctxp->class_err)
    yyerror ("{' expected");
    DRECOVER(class1);
  }
| modifiers CLASS_TK identifier error
  { if (!ctxp->class_err)
    yyerror ("{' expected"); RECOVER;}
;
...
```

...

super:

```
{ $$ = NULL; }  
| EXTENDS_TK class_type { $$ = $2; }  
| EXTENDS_TK class_type error  
  { yyerror ("'{' expected"); ctxp->class_err=1; }  
| EXTENDS_TK error  
  {  
    yyerror ("Missing super class name");  
    ctxp->class_err=1;  
  }  
;
```

Tabellenoptimierung

- LR(0)-Reduktionszustand: Zustand, in dem auf jeden Fall reduziert wird (Kontext unerheblich)
- LR(0)-Reduktionszustände können beseitigt werden, indem im vorigen Zustand bereits reduziert wird (Schift-Reduktionszustand)
- Kettenproduktionen eliminieren
- echte und unechte (don't care) Fehlerübergänge unterscheiden.
- Unecht: Übergang kann nie erreicht werden, z.B. alle leeren Übergänge mit Nichtterminalen
- Fehlerübergänge ausfaktorisieren in Fehlermatrix F:
 $f(q, t) = \text{if } F[q, t] \text{ then Fehler else Eintrag_in_Übergangsmatrix}$
- Übergangsmatrix komprimieren: leere Übergänge berücksichtigen
- Übergangsmatrix weiter komprimieren (siehe nächste Folie)

Methoden zur weiteren Kompression der Übergangsmatrix:

- Graphenfärben: Unverträglichkeitskanten zwischen ungleichen/unverträglichen Tabelleneinträgen. Gleichgefärbte Zeilen/Spalten(-Teile) können zusammengelegt werden.
- Index-Zugriffs-Methoden: Umsortieren und Kombinieren der Einträge, so dass mit konstant vielen Indizierungen der Eintrag findbar ist.
- Listen-Suche: Suche nach dem richtigen Eintrag via Schlüssel und variabler Vergleichszahl.
- Hinweise:
 - Verwendbar bei beliebigen dünn/gleichförmig besetzten Tabellen
 - Historisch: Bedingung für Anwendbarkeit ($|Tabelle| < \text{Hauptspeicher}$) Heute: Geschwindigkeit ($|Tabelle| < \text{Cachegröße}$)

Kettenproduktionen



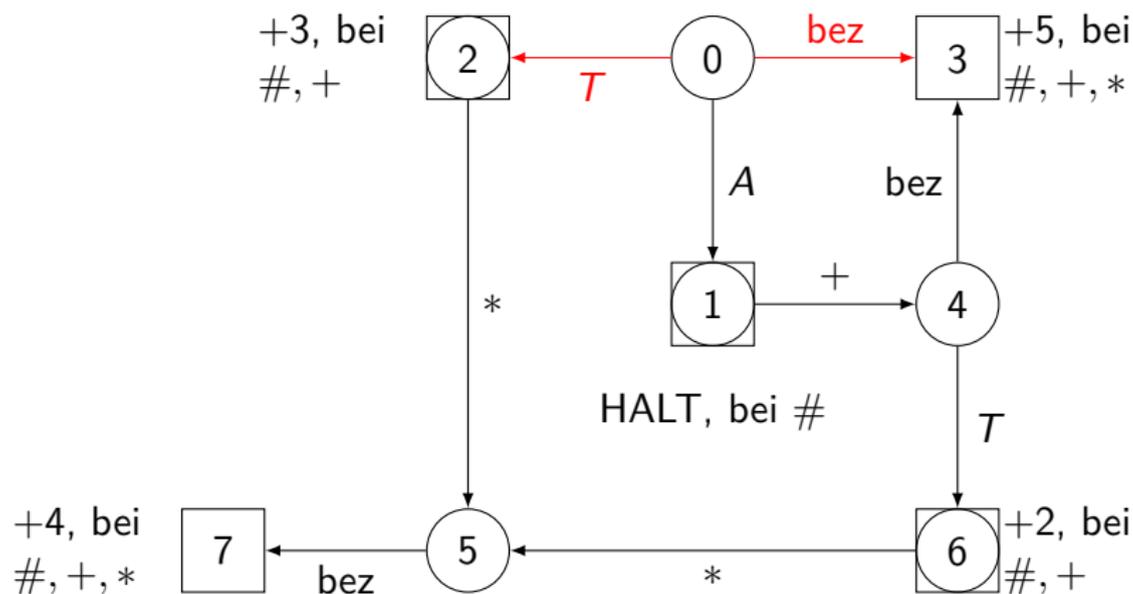
1: $Z \rightarrow A$

2: $A \rightarrow A + T$

3: $A \rightarrow T$

4: $T \rightarrow T * \text{bez}$

5: $T \rightarrow \text{bez}$



Kettenproduktionen Eliminieren



1: $Z \rightarrow A$

2: $A \rightarrow A + T$

3: $A \rightarrow T$

4: $T \rightarrow T * \text{bez}$

5: $T \rightarrow \text{bez}$

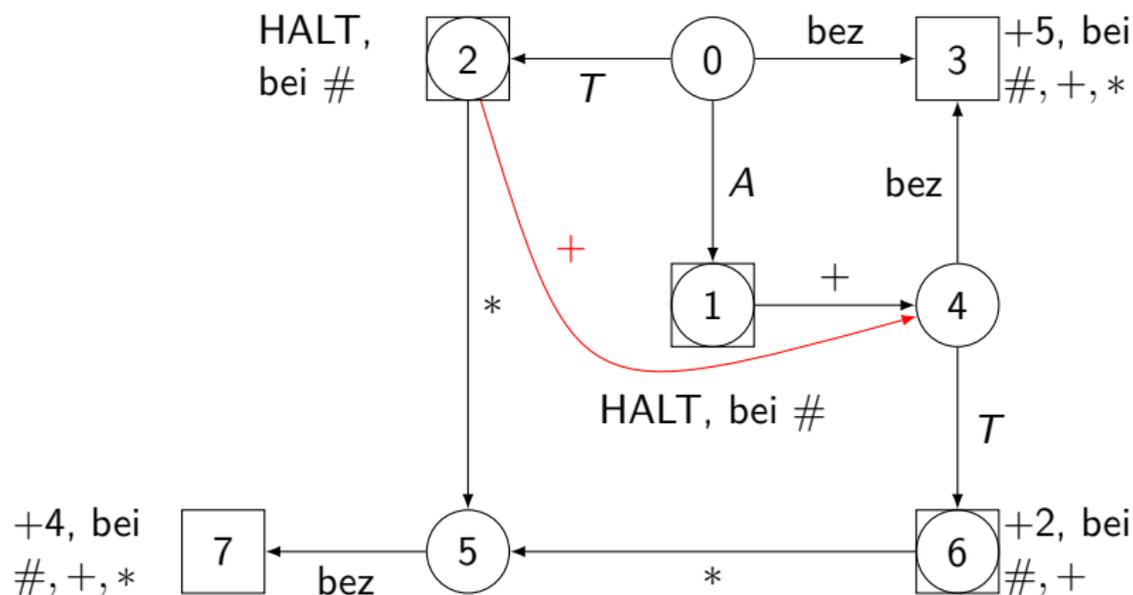


Tabelle mit Schiff-Reduktionen, ohne Kettenproduktionen



	bez	()	+	*	#	A	T	F
0	-6	3	-	-	-	-	1	2	2
1			-	4		*			
2	-	-	-	4	5	*			
3	-6	3	-	-	-	-	6	7	7
4	-6	3	-	-	-	-		8	8
5	-6	3	-	-	-	-			-4
6			-7	4		-			
7	-	-	-7	4	5	-			
8	-	-	+2	+2	5	+2			

- Fehler

i Schiffe, neuer Zustand: i

„ “ Don't care

+r Reduziere mit Regel r

-r Schiffe, reduziere dann mit Regel r

Fehlermatrix und komprimierte Übergangsmatrix



	bez	()	+	*	#
0	f	f	t	t	t	t
1			t	f		f
2	t	t	t	f	f	f
3	f	f	t	t	t	t
4	f	f	t	t	t	t
5	f	f	t	t	t	t
6			f	f		t
7	t	t	f	f	f	t
8	t	t	f	f	f	f

	A	T,F
0,1,2	1	2
3	6	7
4		8
5,6,7,8		-4

	bez	()	+	*	#
0,1,2,3	-6	3	-7	4	5	*
4,5,6,7						
8			+2	+2	5	+2

Tabellengröße nach Elimination von Ketten (Beispiel ADA 83):

- 95 Terminale, 252 Nichtterminale
- 540 Zustände (kodiert in 2 Bytes)
- Tabellengröße: 374.760 Bytes

Reduktion des Platzbedarfes durch:

- Einfache Tabellenkomprimierung
- Elimination der Fehlerübergänge bei Nichtterminalen

Tabellengröße nach diesen Reduktionen 22.584 Bytes (ca. 6%)

	Bison	yacc	PGS	Lalr	Ell
Grammatik	LALR(1)	LALR(1)	LALR(1)	LALR(1)	LL(1)
Spezifiziert in	BNF	BNF	EBNF	EBNF	EBNF
Geschwindigkeit in [10 ³ Symbol / Sekunde]	8.93	15.94	17.32	34.94	54.64
Geschwindigkeit in [10 ³ Zeilen / Minute]	150	270	290	580	910
Tabellengröße in [bytes] (komprimiert)	7724	9968	9832	9620	-
Parsergröße in [bytes]	10900	12200	14140	16492	18048

Eingabe: Modula-2 Code.

Hardware: PCS Cadmus mit MC68020 Prozessor (16.7 MHz).

Komplexität Parsen in $\mathcal{O}(n)$

Für alle Klassen gilt:

- Jeder Ableitungsschritt benötigt, da niemals ein Rücksetzen notwendig ist, konstanten Aufwand $\mathcal{O}(1)$.
- Parsen benötigt Aufwand $\mathcal{O}(n)$ wobei n die Anzahl der Ableitungsschritte ist.

Beweis über Grad der Knoten und mögliche Höhe des Baumes
(Kettenproduktionen und Epsilonproduktionen berücksichtigen)

n : Anzahl der Produktionen, k : Länge der Vorausschau

Grammatiktyp	Parsergenerierung	Test der Grammatik
LL(1)	$\mathcal{O}(n^2)$	$\mathcal{O}(n^2)$
SLL(k)	$\mathcal{O}(n^{k+1})$	$\mathcal{O}(n^{k+1})$
LL(k)	$\mathcal{O}(2^{n^{k+1}+(k+1)\log n})$	$\mathcal{O}(n^{2k})$
SLR(1)	$\mathcal{O}(2^{n+\log n})$	$\mathcal{O}(n^2)$
SLR(k)	$\mathcal{O}(2^{n+k\log n})$	$\mathcal{O}(n^{k+2})$
LR(k)	$\mathcal{O}(2^{n^{k+1}+k\log n})$	$\mathcal{O}(n^{2(k+1)})$

Sätze über kontextfreie Grammatiken

Satz 1: Für jede $LR(k)$ -Grammatik G mit $k > 1$ gibt es eine $LR(1)$ -Grammatik G' mit $L(G) = L(G')$.

Beweis durch Rechtsfaktorisierung.

Satz 2: Jede $LL(k)$ -Grammatik ist auch $LR(k)$.

Satz 3: Es gibt $LR(k)$ -Grammatiken, die für kein k' $LL(k')$ sind.

Satz 4: Es ist entscheidbar, ob es für eine gegebene $LR(k)$ -Grammatik G ein k' gibt, so dass G $LL(k')$ ist.

Satz 5: Es ist unentscheidbar, ob für eine Sprache L eine Grammatik G existiert, so dass G $LL(1)$ ist.

Satz 6: Es ist unentscheidbar, ob es für eine Sprache L eine Grammatik G gibt, so dass G $LL(k)$ oder $LR(k)$ ist.

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken
- 4 LR-, SLR-Grammatiken
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Fehlerbehandlung (1/2)

Ziel: Analyse soll formal korrekte Ergebnisse liefern (oder wegen zu vieler Fehler abbrechen), um möglichst viele Fehler zu finden.

- Gefahr von Folgefehlern wegen unpassender Korrektur unvermeidbar
- Unterscheide **Fehler** (= Ursache) und **Fehlersymptom** (= beobachtbare Wirkung)

Reaktionen:

- **Bericht:** Immer notwendig
- **Reparatur**, wenn Fehlerursache feststellbar und korrigierbar
- **Wiederaufsetzen:** Internen Zustand konsistent machen und weitere Fehler suchen
- **Abbruch** bei zu vielen Fehlern oder bei zu hohem Ressourcenverbrauch (Tabellenüberlauf, ...)

Hinweis: Dieser Abschnitt behandelt Fehlermeldungen aller Übersetzerphasen, nicht nur syntaktische!

Fehlerbehandlung (2/2)

Fehler: Nicht (vom Programmierer) intendiertes Programm

Fehlersymptom:

- Sichtbare Auswirkung des Fehlers: Verletzung der Sprachdefinition
- parserdefinierte Fehlerstelle.

Diagnose:

- Versuch, den Fehler auf der Grundlage des Symptoms zu erkennen.
- Entsprechende Fehlermeldung, Reparatur oder Wiederaufsetzen.

Beispiel

Zuweisung:	$x := (a + b * c;$
Fehler (vermutlich):) nach b fehlt.
Fehlersymptom:) fehlt vor ;

Position des Fehlersymptoms ist nicht die Fehlerstelle!

Fehlerklassen:

- **Anomalien** (verdächtig, aber nicht gefährlich)
 - **Notiz**, z.B. keine Standardkonstruktion
 - **Kommentar**, z.B. unpassender Programmierstil
 - **Warnung** bei möglichen Fehler, z.B. unbenutzte Variable
- **Fehler**
 - **einfacher Fehler**: reparierbar, Code kann erzeugt werden
 - **fataler Fehler**: kein Code erzeugbar, nur Wiederaufsetzen möglich
 - **Abbruchfehler**: Übersetzer gibt auf (z.B. wegen Ressourcenbeschränkung)

- Ausgabe von
 - Position (Datei, Zeile, Spalte),
 - Fehlerklasse,
 - Meldung
- Intern: Meldung kodiert durch ganze Zahl
- Meldungstexte in getrennter Datei zur Anpassung der Sprache
- Meldungen fallen nicht in der Reihenfolge des Eingabetexts an:
 - Meldungen aufsammeln
 - möglichst nach Position sortiert ausgeben

Unterscheidung nach Übersetzerphase

- 1 Symbolfehler:** unzulässiges Eingabezeichen, Abschluß Textkonstante oder Kommentar fehlt (falls erkennbar), verfrühtes Eingabeende: kein Endzustand des Symbolentschlüsselungsautomaten erreicht.
Reparatur: falsche Zeichen oder vorhandenen Symbolanfang ignorieren
- 2 Syntaktischer Fehler:** Satz gehört nicht zu der durch den Parser definierten Obermenge der Sprache.
- 3 Semantischer Fehler:** Fehler in der statischen Semantik.
- 4 Semantischer Fehler,** der erst in der Optimierung entdeckt wird, z.B. Verletzung Indexgrenzen bei Reihungsindizierung.
- 5 Ressourcenbeschränkung** verletzt (stets Abbruch), in allen Phasen möglich.

Parserdefinierte Fehlerstelle

parserdefinierte Fehlerstelle t :

$$\forall r \in \Sigma^* : str \notin L \text{ und } \exists s' : ss' \in L.$$

LL-, SLL-, LR-, SLR- und LALR-Parser finden die parserdefinierte Fehlerstelle.

Andere Parser, z.B. für Präzedenzgrammatiken, finden sie nicht.

Angenommen t parserdefinierte Fehlerstelle, d.h. $str \notin L$,
 $\exists s' : ss' \in L$.

Sei s im Keller und $q = [P \rightarrow P_{anf} \cdot P_{rest}; \Omega]$

Terminal t ist parserdefinierte Fehlerstelle gdw:

$$\nexists k : ts \in \text{Anf}_k(P_{rest}\Omega)$$

Minimale Anzahl von Operationen

- Einsetzen
- Streichen
- Ersetzen (= Streichen + Einsetzen)

ausführen

Tatsächliche Korrektur

- Panischer Modus: Wiederaufsetzen an Anweisungs- oder Vereinbarungsende
- Systematische Fortsetzung an parserdefinierter Fehlerstelle $str \notin L$: frühest möglichen Wiederaufsetzpunkt finden
- Totalkorrektur:
 - Eingabe $str = s_1x_1 \dots s_nx_n$, $s_1 = s$, n minimal, so aufteilen, dass alle s_i Ausschnitte einer korrekten Eingaben sind, die durch eventuell unbrauchbare Texte x_i verknüpft sind
 - die x_i durch korrekte Texte y_i so ersetzen, dass $s_1y_1 \dots s_ny_n$ korrekt ist
 - **Nachteil:** quadratischer Aufwand, praktisch nicht eingesetzt
- Die Zeichen, an denen wieder aufgesetzt werden kann, bilden die Ankermenge.

zahlreiche weitere Verfahren bekannt

Panischer Modus

alle Symbole bis Anweisungs- oder Vereinbarungsende streichen
Keller soweit abbauen, dass Folgesymbol ; end, }, ... akzeptiert
wird

Fehlermeldung ausgeben und Analyse fortsetzen

Vorteil:

- einfache Implementierung

Nachteile:

- keine Analyse des Anweisungsrests
- Schwierigkeiten mit korrektem Abschluß von Klammerungen
if ... then, then ... else, usw.

Fehlerbehandlung für LL-Parser

Panikmodus: Überlese Tokens bis nächstes Token in synchronisierender Menge (Ankermenge). Ankermengen werden beim rekursiven Abstieg mitübergeben bzw. stehen zu einem Nichtterminal auf dem Stack.

Berechnung der Ankermenge für Nichtterminal A : $Ank(A)$

- 1 erste Näherung: $FOLLOW(A)$. Aber das reicht nicht: z.B. fehlendes Semikolon würde zu Folgefehlern führen.
- 2 Erweiterung von $Ank(A)$ um solche Tokens, die „übergeordnete“ Strukturen beginnen. Beispiel: *statement* ist übergeordnet zu *expression*: $FIRST(statement)$ wird zu $Ank(expression)$ hinzugefügt.
- 3 wenn $\varepsilon \in L(A)$, wird an Fehlerstelle $A \rightarrow \varepsilon$ expandiert und normal fortgefahren.
- 4 wenn ein Topstack-Terminal nicht kommt, wird es gepoppt nebst Meldung „Terminal xx expected“.

C (gcc) – Panischer Modus

```
int main ( ) {  
    int j, i ;  
    if (i<j) {{ /*FEHLER*/  
        i = j ;  
    } /*FEHLER–SYMPTOM 1: } erwartet*/  
    else {  
        if ( i != j ){  
            j: = i ; /*FEHLER 2 bleibt unerkannt*/  
        } /*Ende if */  
    } /*Ende main*/  
    return 0 ; /*FEHLER–SYMPTOM 2: return unerwartet*/  
}
```

test.c: In function 'main':

test.c: 6: parse error before 'else '

test.c: At top level:

test.c: 11: parse error before 'return '

Systematische Korrektur (Röhrich)

gegeben parserdefinierter Fehlerstelle $str \notin L$

gesucht: Fortsetzung $sy \in L$:

- Bestimme eine **Ankermenge** $D = \{d \in \Sigma \mid sy = ss'ds'' \in L\}$
- Suche ein $d \in D$, so dass $r = r'dr''$ und $|r'|$ minimal ist.
- Ersetze $tr'dr''$ durch $s'dr''$. $\frac{str=str'dr''}{ss'dr''}$
- gib Fehlermeldung aus und setze Analyse fort
- **Vorteile:**
 - nahe bei Minimalkorrektur, einfache Fehlermeldung
 - fast vollautomatisch erzeugbar
 - terminiert, da Eingabe um d verkürzt.
- **Nachteile:**
 - Korrektur nicht zwangsläufig korrekt
 - Schwierigkeiten bei Listenkonstruktionen (Anweisungs-, Bezeichner-, Vereinbarungslisten, usw.)
 - bei rekursivem Abstieg Vorbereitung notwendig
 - bei LR-Analyse Adaption des Generators notwendig

ALGOL 60 - Systematische Korrektur (Röhrich)

```
BEGIN
```

```
  INTEGER ARRAY A,B(1...5 1...10);
```

^

```
***Error in front of 1 <,> inserted
```

```
  INTEGER I,J,K,L;
```

```
  UP:I+J>K+L*4 THEN GO L1 ELSE K IS 2;
```

^

^

^

```
***Error in front of <+> <:=> inserted
```

```
***Error in front of THEN <;> inserted
```

```
  THEN deleted
```

```
***Error in front of <L1> TO inserted
```

```
  ELSE <K> IS <2> deleted
```

ADA (adac) - Fehlerbehandlung

```
procedure Main is
  I;J : Integer;
begin
  if I > J then then /* FEHLER 1: I in Ankermenge,
                    then gestrichen */

    I:=J
  else
    if I /= J then
      J = ; I /* FEHLER 2: = unerwartet,
              ; in Ankermenge, = gestrichen */
    end if /* SYMPTOM: ; erwartet, end in Ankermenge,
            ; eingefügt */
  end if;
end Main;
```

ADA (adac) - Fehlerausgabe

```
error.ada...
```

```
4, 19: Error syntax error
```

```
4, 19: Information expected tokens:
```

```
IDENTIFIER CHAR_STRING NULL PRAGMA CASE RETURN
```

```
FOR BEGIN EXIT GOTO DELAY ABORT RAISE REQUEUE IF
```

```
WHILE LOOP DECLARE ACCEPT SELECT <<
```

```
5, 3: Information      restart point
```

```
8, 7: Error           syntax error
```

```
8, 7: Information     expected tokens: . ; ( , : : =
```

```
8, 9: Information     restart point
```

```
8, 6: Error           syntax error
```

```
9, 6: Information     expected tokens: . ; ( , : : =
```

```
9, 6: Repair          token inserted : ;
```

Durchführung

- 1 Zeichne in jedem Zustand eine Situation aus, deren weitere Verfolgung zur Generierung der Fortsetzung führt. Die Fortsetzung muss terminieren. Ankermenge sind alle Symbole in dieser Fortsetzung.
- 2 Streiche in der Eingabe alle Zeichen bis zum ersten Zeichen in der Ankermenge.
- 3 Dieses Zeichen wird während der Verarbeitung der generierten Fortsetzung in einem Zustand q' akzeptiert und es gibt einen Übergang vom Fehlerzustand q nach q'
- 4 Setze die Zeichen ein, die den Automaten vom Zustand q nach q' bringen.

LL Durchführung

- Zeichne für jedes Nichtterminal eine Produktion aus, die in der Fortsetzungserzeugung vorhergesagt wird. Die Produktion darf nicht rekursiv sein!
- Bestimme für die laufenden und jede vorhergesagte Produktion die noch fehlenden Symbole. Diese bilden die Ankermenge. Die noch laufenden Produktionen sind aus dem Keller ersichtlich.
- Schwierigkeit bei rekursivem Abstieg: Keller nur nach Prozedurrückkehr zugänglich.
 - Abhilfe: Ankermenge bereits während der normalen Syntaktische Analyse aufbauen und in getrenntem Keller (Datenstruktur) ablegen oder als Argument bei Prozeduraufruf mitgeben.

LL Beispiel Automat

$$Z \rightarrow A, A \rightarrow FA', A' \rightarrow \varepsilon \mid + FA', F \rightarrow i \mid (A)$$

$q_0 : [Z \rightarrow \cdot A]$	$q_1 : [Z \rightarrow A \cdot]$	$*q_0 i \rightarrow q_1 q_2 i, q_0 (\rightarrow q_1 q_2 (,$
$q_2 : [Z \rightarrow \cdot FA']$	$q_3 : [A \rightarrow F \cdot A']$	$*q_1 \rightarrow \varepsilon,$
$q_4 : [F \rightarrow \cdot i]$	$q_5 : [F \rightarrow \cdot (A)]$	$*q_2 i \rightarrow q_3 q_4 i, q_2 (\rightarrow q_3 q_5 (,$
$q_6 : [A \rightarrow FA' \cdot]$	$q_7 : [A' \rightarrow \cdot \varepsilon]$	$*q_3 \# \rightarrow q_6 q_7 \#, q_3) \rightarrow q_6 q_7),$
$q_8 : [A' \rightarrow \cdot + FA']$		$q_3 \rightarrow q_6 q_8 +, *q_4 i \rightarrow q_9,$
$q_9 : [F \rightarrow i \cdot]$		$*q_5 (\rightarrow q_{10}, *q_6 \rightarrow \varepsilon, *q_7 \rightarrow \varepsilon,$
$q_{10} : [F \rightarrow (\cdot A)]$		$*q_8 + \rightarrow q_{11}, *q_9 \rightarrow \varepsilon,$
$q_{11} : [A' \rightarrow + \cdot FA']$		$*q_{10} i \rightarrow q_{12} q_2 i, q_{10} (\rightarrow q_{12} q_2 (,$
$q_{12} : [F \rightarrow (A \cdot)]$		$*q_{11} i \rightarrow q_{13} q_4 i, q_{11} (\rightarrow q_{13} q_5 (,$
$q_{13} : [A' \rightarrow + F \cdot A']$		$*q_{12}) \rightarrow q_{14},$
$q_{14} : [F \rightarrow (A) \cdot]$		$*q_{13} \# \rightarrow q_{15} q_7 \#,$
$q_{15} : [A' \rightarrow + FA' \cdot]$		$q_{13}) \rightarrow q_{15} q_7), q_{13} + \rightarrow q_{15} q_{15} +,$
		$*q_{14} \rightarrow \varepsilon, *q_{15} \rightarrow \varepsilon$

* zeichnet die vorherzusagenden Produktionen bzw. Zustandsübergänge aus

LL Beispiel $i + \#$

Parse bis
zum Fehler

$q_0 i + \#$
 $q_1 q_2 i + \#$
 $q_1 q_3 q_4 i + \#$
 $q_1 q_3 q_9 + \#$
 $q_1 q_3 + \#$
 $q_1 q_6 q_8 + \#$
 $q_1 q_6 q_{11} \#$

Fortsetzung finden

$q_1 q_6 q_{11} D = \{i(\{)\}$
 $q_1 q_6 q_{13} q_4$
 $q_1 q_6 q_{13} q_9$
 $q_1 q_6 q_{13} D = \{i(\#)+\}$
 $q_1 q_6 q_{15} q_7$
 $q_1 q_6 q_{15}$
 $q_1 q_6$
 q_1

Wiederaufsetzen

$q_1 q_6 q_{11} \#$
 $q_1 q_6 q_{13} q_4 \#$
 $q_1 q_6 q_{13} q_9 \#$
 $q_1 q_6 q_{13} \#$

i generiert mit
 $q_4 i \rightarrow q_9$
weiter normal

Problem am Beispiel (1/2)

Grammatik:

$$1 \quad Z \rightarrow I \mid A$$

$$2 \quad I \rightarrow \text{if } E \text{ then } Z \text{ end}$$

$$3 \quad A \rightarrow E := E$$

$$4 \quad E \rightarrow id \ E_{rest}$$

$$5 \quad E_{rest} \rightarrow \varepsilon \mid = id$$

SLL-Analyse des fehlerhaften Satzes:

if a:=b then ... end

mit parserdefiniertem Fehler := landet in Situation

[I \rightarrow if E · then Z end; Ω]

Beachte: $:= \in \text{Folge}_1(E)$

Problem am Beispiel (2/2)

Grammatik:

$$1 \quad Z \rightarrow I \mid A$$

$$2 \quad I \rightarrow \text{if } B \text{ then } Z \text{ end}$$

$$3 \quad A \rightarrow E := E$$

$$4 \quad \text{redB} \rightarrow E$$

$$5 \quad E \rightarrow id E_{rest}$$

$$6 \quad E_{rest} \rightarrow \varepsilon \mid = id$$

Parse von **if** $a:=b$ then ... end

mit parserdefiniertem Fehler $:=$ landet in Situation

$[I \rightarrow \text{if } E \cdot \text{ then } Z \text{ end}; \Omega]$

Beachte: $:= \notin \text{Folge}_1(B)$

Fehlerbehandlung für LR-Parser

Beispiel für Grammatik:

$$E \rightarrow E + E \mid E * E \mid (E) \mid \mathbf{id}$$

Erzeuge 4 Aktionen zur Fehlerbehebung:

e1 **id** erwartet aber nicht gefunden.

Aktion: Shift 3; Nachricht „Operand fehlt“.

e2 **)** gefunden ohne vorherige **(**.

Aktion: Überspringe Token; Nachricht „) ohne Gegenstück“.

e3 Operator erwartet aber **id** oder **)** gefunden.

Aktion: Shift 4; Nachricht „Operator fehlt“.

e4 **#** gefunden aber noch Klammern geöffnet.

Aktion: Nachricht „) fehlt“.

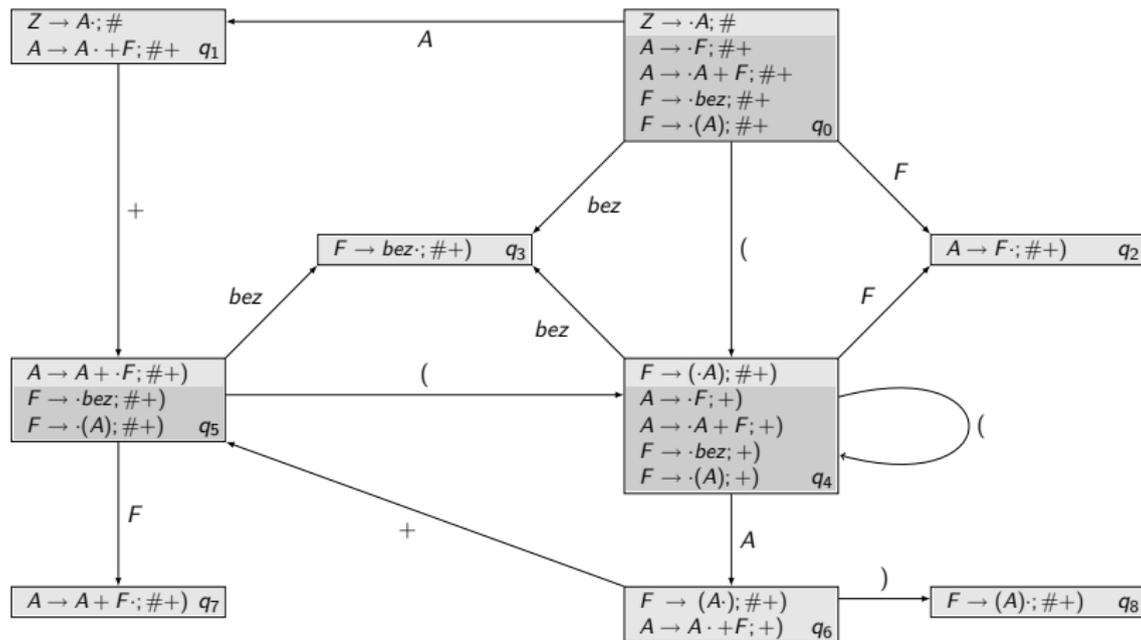
LR-Parsertabelle mit Fehlerrountinen

Zustand	ACTION						GOTO
	id	+	*	()	#	<i>E</i>
0	s3	e1	e1	s2	e2	e1	1
1	e3	s4	s5	e3	e2	acc.	
2	s3	e1	e1	s2	e2	e1	6
3	r4	r4	r4	r4	r4	r4	
4	s3	e1	e1	s2	e2	e1	7
5	s3	e1	e1	s2	e2	e1	8
6	e3	s4	s5	e3	s9	e4	
7	r1	r1	s5	r1	r1	r1	
8	r2	r2	r2	r2	r2	r2	
9	r3	r3	r3	r3	r3	r3	

LR Durchführung

- Zeichne für jedes Nichtterminal eine nichtrekursive Produktion aus
- Zustände sind nicht Mengen, sondern geordnete Listen von Situationen. Die Situationen für ausgezeichnete Produktionen kommen jeweils vor allen anderen Situationen mit gleicher linker Seite.
- Erzeugung der Fortsetzung: Benutze in jedem Zustand jeweils die erste Situation und ihren Übergang.
- Die ausgezeichnete Situation gehört zur Basis des Zustandes.
- Verfahren terminiert wegen Nichtrekursivität der ausgezeichneten Produktionen
- Sonderbehandlung von Trennzeichen für Listen. Sie werden initial in entsprechende Ankermenge übernommen.

LR Beispiel-Automat



Übergänge

q	bez	()	+	#	A	F
0	-4	4	-	-	-	1	-3
1	-	-	-	5	*1		
2	-4	4	-	-	-	6	-3
4	-4	4	-	-	+4	6	-3
5	-4	4	-	-	+4		-2
6	-	-	-5	5	-		

Grau hinterlegt: Ausgezeichnete Übergänge

LR Beispiel $i +) i \#$

Parse bis
zum Fehler

$q_0 i +) i \#$
 $q_0 q_1 +) i \#$
 $q_0 q_1 q_5 i \#$

Fortsetzung finden

$q_0 q_1 q_5 D = \{i(\}$
 $q_0 q_1 D = \{i(+\# \}$

Wiederaufsetzen

$q_0 q_1 q_5 i \#$ wobei $)$ gelöscht
weiter normal

Fehlerbehandlung in der semantischen Analyse

für jede Eigenschaft oder Wertetyp einen ausgezeichneten Wert *unbekannt* definieren

- Operationen mit unbekannt liefern wieder unbekannt
- Fehlermeldung nur, wenn unbekannt als Ergebnis, nicht als Operand erscheint.

Quellen fataler Fehler:

- fehlende Vereinbarung: Bezeichner erhält Typ unbekannt (oder Typ durch Typinferenz aus Anwendung erschließen)
- unzulässige Operation: Ergebnistyp ist unbekannt
- inkompatible Operanden oder sonstige Verletzung von Konsistenzbedingungen: Fehler melden, Ergebnistyp unbekannt, mit Analyse fortfahren

Kapitel 3: Syntaktische Analyse

- 1 Einbettung
- 2 Theoretische Grundlage: Kontextfreie Grammatiken
- 3 LL- und SLL-Grammatiken
- 4 LR-, SLR-Grammatiken
 - LALR-Konstruktion
 - Parsergeneratoren
 - Bison
 - Optimierungen und Komplexität
- 5 Fehlerbehandlung
- 6 Earley Parser

Earley Parser

- Algorithmus zum Parsen kontextfreier Grammatiken mit dynamischer Programmierung.
- Verarbeitet mehrdeutige Grammatiken und erzeugt alle möglichen Parsebäume eines Satzes.
- Laufzeit liegt im allgemeinen in $O(n^3)$; in $O(n^2)$ für eindeutige Grammatiken; in $O(n)$ für fast alle LR(k) Grammatiken.
- Häufiger Einsatz in der Computerlinguistik.

Algorithmus

- Erweitere die Grammatik um eine Regel $S' \rightarrow S$.
- Für jede Eingabeposition wird Menge S_i von Situationen $[X \rightarrow \alpha \cdot \beta, j]$ berechnet. j gibt die Eingabeposition an. Initiale Menge: $\{[S' \rightarrow \cdot S, 0]\}$.
- Für jede Eingabeposition i wende 3 Schritte iterativ bis zum Fixpunkt an:
 - **Prediction:** Für jede Situation $[X \rightarrow \alpha \cdot A \beta, j]$ in S_i füge eine neue Situation $[A \rightarrow \cdot \gamma, i]$ für jede Produktion der Grammatik mit A auf der linken Seite ein (transitive Hülle bilden).
 - **Completion:** Für jede Situation $[A \rightarrow \gamma \cdot, j]$ in S_i füge für jede Situation $[X \rightarrow \alpha \cdot A \beta, k]$ in S_j eine neue Situation $[X \rightarrow \alpha A \cdot \beta, k]$ in S_i ein.
 - **Scanning:** Sei \mathbf{a} das nächste Token. Füge für jede Situation $[X \rightarrow \alpha \cdot \mathbf{a} \beta, j]$ in S_i eine Situation $[X \rightarrow \alpha \mathbf{a} \cdot \beta, j]$ in S_{i+1} ein.

Beispiel

Grammatik:

$$S \rightarrow S \text{ und } S \mid S \text{ oder } S$$
$$S \rightarrow \text{blau} \mid \text{gestreift} \mid \text{glatt} \mid \text{teuer}$$

Eingabe: **blau und gestreift oder glatt und teuer**

Beispiel

Eingabe: **blau** und gestreift oder glatt und teuer # (Pos. 0)

Situationen S_0 :

$[S' \rightarrow \cdot S, 0]$

Iteration:

$[S \rightarrow \cdot S \text{ und } S, 0]$

$[S \rightarrow \cdot S \text{ oder } S, 0]$

$[S \rightarrow \cdot \text{blau}, 0]$

$[S \rightarrow \cdot \text{gestreift}, 0]$

$[S \rightarrow \cdot \text{glatt}, 0]$

$[S \rightarrow \cdot \text{teuer}, 0]$

Beispiel

Eingabe: **blau und gestreift oder glatt und teuer** # (Pos. 1)

Situationen S_1 :

$[S \rightarrow \mathbf{blau}\cdot, 0]$

Iteration:

$[S' \rightarrow S\cdot, 0]$

$[S \rightarrow S \cdot \mathbf{und} S, 0]$

$[S \rightarrow S \cdot \mathbf{oder} S, 0]$

Beispiel

Eingabe: **blau und gestreift** oder **glatt und teuer** # (Pos. 2)

Situationen S_2 :

$[S \rightarrow S \text{ und } \cdot S, 0]$

Iteration:

$[S \rightarrow \cdot S \text{ und } S, 2]$

$[S \rightarrow \cdot S \text{ oder } S, 2]$

$[S \rightarrow \cdot \text{blau}, 2]$

$[S \rightarrow \cdot \text{gestreift}, 2]$

$[S \rightarrow \cdot \text{glatt}, 2]$

$[S \rightarrow \cdot \text{teuer}, 2]$

Beispiel

Eingabe: **blau und gestreift oder glatt und teuer** # (Pos. 3)

Situationen S_3 :

$[S \rightarrow \text{gestreift}\cdot, 2]$

Iteration:

$[S \rightarrow S \text{ und } S\cdot, 0]$

$[S' \rightarrow S\cdot, 0]$

$[S \rightarrow S \cdot \text{und } S, 0]$

$[S \rightarrow S \cdot \text{oder } S, 0]$

$[S \rightarrow S \cdot \text{und } S, 2]$

$[S \rightarrow S \cdot \text{oder } S, 2]$

Beispiel

Eingabe: blau und gestreift oder **glatt** und teuer # (Pos. 4)

Situationen S_4 :

$[S \rightarrow S \text{ oder} \cdot S, 0]$

$[S \rightarrow S \text{ oder} \cdot S, 2]$

Iteration:

$[S \rightarrow \cdot S \text{ und } S, 4]$

$[S \rightarrow \cdot S \text{ oder } S, 4]$

$[S \rightarrow \cdot \text{blau}, 4]$

$[S \rightarrow \cdot \text{gestreift}, 4]$

$[S \rightarrow \cdot \text{glatt}, 4]$

$[S \rightarrow \cdot \text{teuer}, 4]$

Beispiel

Eingabe: blau und gestreift oder glatt **und** teuer # (Pos. 5)

Situationen S_5 :

$[S \rightarrow \text{glatt}\cdot, 4]$

Iteration:

$[S \rightarrow S \text{ oder } S\cdot, 0]$

$[S \rightarrow S \text{ oder } S\cdot, 2]$

$[S \rightarrow S \cdot \text{und } S, 4]$

$[S \rightarrow S \cdot \text{oder } S, 4]$

$[S' \rightarrow S\cdot, 0]$

$[S \rightarrow S \cdot \text{und } S, 0]$

$[S \rightarrow S \cdot \text{oder } S, 0]$

$[S \rightarrow S \text{ und } S\cdot, 0]$

$[S \rightarrow S \cdot \text{und } S, 2]$

$[S \rightarrow S \cdot \text{oder } S, 2]$

Beispiel

Eingabe: **blau und gestreift oder glatt und teuer** # (Pos. 6)

Situationen S_6 :

$[S \rightarrow S \text{ und } \cdot S, 4]$

$[S \rightarrow S \text{ und } \cdot S, 0]$

$[S \rightarrow S \text{ und } \cdot S, 2]$

Iteration:

$[S \rightarrow \cdot S \text{ und } S, 6]$

$[S \rightarrow \cdot S \text{ oder } S, 6]$

$[S \rightarrow \cdot \text{blau}, 6]$

$[S \rightarrow \cdot \text{gestreift}, 6]$

$[S \rightarrow \cdot \text{glatt}, 6]$

$[S \rightarrow \cdot \text{teuer}, 6]$

Beispiel

Eingabe: **blau und gestreift oder glatt und teuer** # (Pos. 7)

Situationen S_7 :

$[S \rightarrow \text{teuer}\cdot, 6]$

Iteration:

$[S \rightarrow S \cdot \text{und } S, 6]$

$[S \rightarrow S \cdot \text{oder } S, 6]$

$[S \rightarrow S \text{ und } S\cdot, 4]$

$[S \rightarrow S \text{ und } S\cdot, 0]$

$[S \rightarrow S \text{ und } S\cdot, 2]$

$[S \rightarrow S \cdot \text{und } S, 4]$

$[S \rightarrow S \cdot \text{oder } S, 4]$

$[S' \rightarrow S\cdot, 0]$

$[S \rightarrow S \cdot \text{und } S, 0]$

$[S \rightarrow S \cdot \text{oder } S, 0]$

$[S \rightarrow S \cdot \text{und } S, 2]$

$[S \rightarrow S \cdot \text{oder } S, 2]$