

Compilerbau I

Zusammenfassung wichtiger Elemente der Einführung in den Compilerbau

von Alexander Köster
Student der Universität Siegen, CB1 2020
Letzte Aktualisierung: 25. Juli 2020

Dieser Kurs der theoretischen Informatik befasst sich mit den Grundlagen und Bauteilen eines Compilers.

Einzig für Lernzwecke erstellt.

Nicht geeignet als Klausurhilfe.

Das Erstellen einer eigenen Klausurhilfe führt zu einem besonders guten Lerneffekt.

Dieses Dokument sollte nur als Orientierung oder Vergleich dienen.

*Jeder sollte seine Klausurhilfe individuell auf seinen Lernstand und seine eigenen Probleme anpassen,
gut bekannte Dinge auslassen und schlecht merkbare Dinge hinzufügen.*

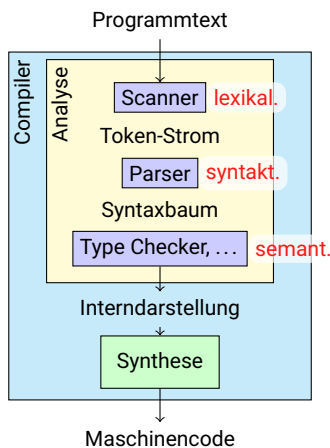
*Dieses Dokument beinhaltet viel mehr, als für eine tatsächliche Klausur durchschnittlich benötigt wird.
Für einen guten Ansatz, welche Inhalte man braucht, sollte sich an der Probeklausur orientiert werden.*

© study.woalk.de

Vervielfältigung ohne ausdrückliche Erlaubnis des Autors außerhalb der originalen Website untersagt.

0 Einführung

- Interpreter:** Keine Vorberechnung auf dem Programmtext erforderlich, Interpreter führt Programmtext „direkt“ aus.
 - Bsp.:** Shellskript (*sh*), Javascript (*js*)
 - Vorteil:** geringe Startup-Zeit
 - Nachteil:** Während der Ausführung werden die Bestandteile immer wieder analysiert \Rightarrow längere Laufzeit.
- Compiler:** Zwei Phasen:
 - Übersetzung des Programmtexts in Maschinencode
 - Ausführung des Maschinencodes
 - Vorteil:** Optimierungen sind im Übersetzungsschritt möglich, z. B. geschicktere Verwaltung der Variablen, d. h. die Ausführung des Programms wird effizienter.
 - Nachteil:** Übersetzung dauert selbst Zeit. Lohnt sich bei aufwendigen oder oft ausgeführten Programmen.



1 Analyse-Phase

1.1 Lexikalische Analyse

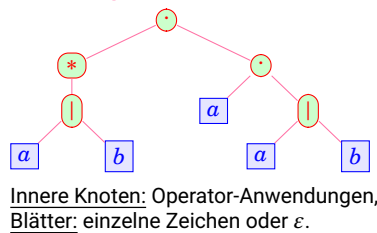
- Token** ist eine Folge von Zeichen, die zusammen eine Einheit bilden.
 - Token-Klassen:**
 - Namen (*identifier*) **Bsp.:** *xyz, pi, ...*
 - Konstanten (*constants*) **Bsp.:** *42, 3.14, "abc", ...*
 - Operatoren (*operators*) **Bsp.:** *+, -, ...*
 - reservierte Worte (*keywords*) **Bsp.:** *if, while, int, ...*
 - Sieber:** Vorverarbeitung klassifizierter Tokens
 - Wegwerfen irrelevanter Teile wie Leerzeichen, Kommentare, ...
 - Aussondern von **Pragmas**, d. h. Direktiven an den Compiler, die nicht Teil des Programms sind (**Bsp.:** *#include in C*)
 - Ersetzen der Tokens bestimmter Klassen durch ihre Bedeutung (Konstanten) bzw. Interndarstellung (Namen, die typischerweise zentral in einer **Symbol-Tabelle** verwaltet werden)
- Scanner und Sieber werden i. A. in einer Komponente zusammengefasst, indem man dem Scanner nach Erkennen eines Tokens eine Aktion gestattet.
- Scanner i. A. nicht von Hand programmiert, sondern aus Spezifikation autom. generiert.
 - Produktivität:** Die Komponente lässt sich schneller herstellen.
 - Korrektheit:** Die Komponente realisiert (beweisbar) die Spezifikation.
 - Effizienz:** Der Generator kann die Komponente mit den effizientesten Algorithmen ausstatten.

1.1.1 Grundlagen: Reguläre Ausdrücke

- Programmtext: endl. Alphabet Σ , **Bsp.:** ASCII
- Die Menge der Textabschnitte einer Token-Klasse ist i. A. eine **reguläre Sprache**.
- Reguläre Sprachen kann man mithilfe regulärer Ausdrücke (*reg. expressions*) spezifizieren.
- Def.:** Menge \mathcal{E}_Σ der (nicht-leeren) **regulären Ausdrücke** ist die kleinste Menge \mathcal{E} mit:
 - $\varepsilon \in \mathcal{E}$, wobei ε neues Symbol $\notin \Sigma$
 - $\Sigma \subseteq \mathcal{E}$
 - $e_1, e_2 \in \mathcal{E} \Rightarrow (e_1 | e_2), (e_1 \cdot e_2), e_1^* \in \mathcal{E}$
- ($|, \dots$ nennt man **Meta-Zeichen**.
- Um Klammern zu sparen: $* > \cdot > |$ und \cdot weglassen
- Reale Spezifikationsprachen bieten zusätzliche Konstrukte wie:
 - $e? \equiv (\varepsilon | e)$ und verzichten auf ε selbst.
 - $e^+ \equiv (e \cdot e^*)$
 - $[a - b] \equiv a | a_{i+1} | \dots | a_j$ wenn $a = a_i, b = a_j, \Sigma = \{a_1 \leq a_2 \leq \dots\}$ total geord.
- Def.:** Für den regulären Ausdruck $e \in \mathcal{E}_\Sigma$ ist die **spezifizierte Sprache** $\llbracket e \rrbracket \subseteq \Sigma^*$ ind. def. durch:
 - $\llbracket \varepsilon \rrbracket = \{\varepsilon\}$
 - $\llbracket a \rrbracket = \{a\}$
 - $\llbracket e^* \rrbracket = (\llbracket e \rrbracket)^*$, wobei für bel. Sprachen (Mengen) L gilt: $L^* := \{w_1 \dots w_k \mid k \geq 0, w_i \in L \forall i \in \underline{k}\}$
 - $\llbracket e_1 | e_2 \rrbracket = \llbracket e_1 \rrbracket \cup \llbracket e_2 \rrbracket$
 - $\llbracket e_1 \cdot e_2 \rrbracket = \llbracket e_1 \rrbracket \cdot \llbracket e_2 \rrbracket$, wobei für bel. Sprachen (Mengen) L_1, L_2 gilt: $L_1 \cdot L_2 := \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$

Bsp.: $\llbracket ab^*a \rrbracket = \{ab^n a \mid n \geq 0\}$

- Reguläre Ausdrücke als markierte geordnete Bäume: **Bsp.:** $(a|b)^* a(a|b)$



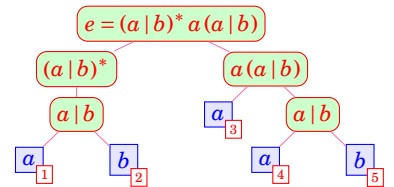
1.1.2 Grundlagen: Endliche Automaten

- Bsp.:**
- Knoten:** Zustände, **Kanten:** Übergänge, **Beschriftungen:** konsumierter Input
- ε -NFA** (nicht-deterministischer endl. Automat mit ε -Übergängen): $A = (Q, \Sigma, \delta, I, F)$ mit:
 - Q endl. Menge von Zuständen
 - Σ endl. Eingabe-Alphabet
 - $I \subseteq Q$ Menge der Anfangszustände
 - $F \subseteq Q$ Menge der akzept. Endzustände
 - $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times Q$ Übergangs-Relation
- Andere endl. Automaten:
 - Gibt es keine ε -Übergänge, ist A ein **NFA**.
 - Ist $\delta : Q \times \Sigma \rightarrow Q$ eine **Funktion** und $|I| = 1$, heißt A deterministisch (**DFA**).
- Berechnungen** sind Pfade im Graphen.
- Akzeptierende Berechnungen** führen von einem Zustand aus I zu einem Zustand aus F .
- Ein **akzeptiertes Wort** ist die Beschriftung einer akzeptierenden Berechnung. Die Menge aller akzeptierten Worte ist die **akzeptierte Sprache** vom endl. Automat A : $\mathcal{L}(A) = \{w \in \Sigma^* \mid \exists i \in I, f \in F : (i, w, f) \in \delta^*\}$

- Transitiver Abschluss** δ^* von δ ist die kleinste Menge δ^* mit, für alle $x \in \Sigma \cup \{\varepsilon\}, w \in \Sigma^*$:
 - $(p, \varepsilon, q) \in \delta^*$
 - Wenn $(p, x, p_1) \in \delta^* \wedge (p_1, w, q) \in \delta^* \Rightarrow (p, xw, q) \in \delta^*$
- δ^* beschreibt für je 2 Zustände, mit welchen Wörtern man von einem zum anderen kommt.
- Satz:** $\forall e \in \mathcal{E}_\Sigma$ kann (in lin. Zeit) ein ε -NFA konstruiert werden, der die Sprache $\llbracket e \rrbracket$ akzeptiert.

Konstruktion des ε -NFA für $e \in \mathcal{E}_\Sigma$:

- Nummeriere die Blätter. **Bsp.:**



Verwende $\text{num} : \mathcal{E}_\Sigma \rightarrow \mathcal{E}_{\mathbb{N} \times \Sigma}$, dazu die Hilfsfkt. $\text{num}' : \mathbb{N} \rightarrow (\mathcal{E}_\Sigma \rightarrow \mathcal{E}_{\mathbb{N} \times \Sigma})$ definiert für $i \in \mathbb{N}, x \in \Sigma, e_1, e_2 \in \mathcal{E}_\Sigma$ durch:

- $\text{num}'(i)(x) := [i, x]$
- $\text{num}'(i)(\varepsilon) := \varepsilon$
- $\text{num}'(i)(e_1 | e_2) := \text{num}'(i)(e_1) | \text{num}'(i + \ell(e_1))(e_2)$
- $\text{num}'(i)(e_1 \cdot e_2) := \text{num}'(i)(e_1) \cdot \text{num}'(i + \ell(e_1))(e_2)$
- $\text{num}'(i)(e_1^*) := \text{num}'(i)(e_1)^*$

wobei $\ell : \mathcal{E}_\Sigma \rightarrow \mathbb{N}$ die Anzahl Terminalsymbole ist, d. h. für $x \in \Sigma, e_1, e_2 \in \mathcal{E}_\Sigma$:

- $\ell(x) := 1$
- $\ell(\varepsilon) := 0$
- $\ell(e_1 | e_2) := \ell(e_1) + \ell(e_2)$
- $\ell(e_1 \cdot e_2) := \ell(e_1) + \ell(e_2)$
- $\ell(e_1^*) := \ell(e_1)$

Dann ist $\text{num} := \text{num}(1)$.

Bsp.: $\text{num}((a|b)^* a(a|b)) = ((1, a) | (2, b))^* (3, a) ((4, a) | (5, b))$

- Zustände:** $\bullet r, r \bullet$ für alle Knoten r von e

Anfangszustand: $\bullet e$, **Endzustand:** $e \bullet$

Übergangsrelation:

- Für Blätter $r = [i, x]: (\bullet r, x, r \bullet) \in \delta$
- Für Knoten der Form $r = r_1 | r_2$:
 - $(\bullet r, \varepsilon, \bullet r_1) \in \delta$
 - $(\bullet r, \varepsilon, \bullet r_2) \in \delta$
 - $(r_1 \bullet, \varepsilon, r \bullet) \in \delta$
 - $(r_2 \bullet, \varepsilon, r \bullet) \in \delta$
- Für Knoten der Form $r = r_1 \cdot r_2$:
 - $(\bullet r, \varepsilon, \bullet r_1) \in \delta$
 - $(r_1 \bullet, \varepsilon, \bullet r_2) \in \delta$
 - $(r_2 \bullet, \varepsilon, r \bullet) \in \delta$
- Für Knoten der Form $r = r_1^*$:
 - $(\bullet r, \varepsilon, r \bullet) \in \delta$
 - $(\bullet r, \varepsilon, \bullet r_1) \in \delta$
 - $(r_1 \bullet, \varepsilon, \bullet r_1) \in \delta$
 - $(r_1 \bullet, \varepsilon, r \bullet) \in \delta$
- Für Knoten der Form $r = r_1?$:
 - $(\bullet r, \varepsilon, r \bullet) \in \delta$
 - $(\bullet r, \varepsilon, \bullet r_1) \in \delta$
 - $(r_1 \bullet, \varepsilon, r \bullet) \in \delta$

• **Konstruktion des NFA ohne ϵ -Übergänge:**

Zunächst einige Definitionen:

- ▷ $\text{empty}[r] = \text{true} \Leftrightarrow \epsilon \in [r]$
 - Für Blätter $r = [i, x]$ ist $\text{empty}[r] := (x \equiv \epsilon)$
 - $\text{empty}[r_1 | r_2] := \text{empty}[r_1] \vee \text{empty}[r_2]$
 - $\text{empty}[r_1 \cdot r_2] := \text{empty}[r_1] \wedge \text{empty}[r_2]$
 - $\text{empty}[r_1^*] := \text{true}$
 - $\text{empty}[r_1?] := \text{true}$
- ▷ $\text{first}[r] := \{i \text{ in } r \mid (\bullet r, \epsilon, \bullet [i, x]) \in \delta^*, x \neq \epsilon\}$
 - Für Blätter $r = [i, x]$ ist $\text{first}[r] := \{i \mid x \neq \epsilon\}$
 - $\text{first}[r_1 | r_2] := \text{first}[r_1] \cup \text{first}[r_2]$
 - $\text{first}[r_1 \cdot r_2] := \begin{cases} \text{first}[r_1] \cup \text{first}[r_2] & \text{empty}[r_1] = \text{true} \\ \text{first}[r_1] & \text{empty}[r_1] = \text{false} \end{cases}$
 - $\text{first}[r_1^*] := \text{first}[r_1]$
 - $\text{first}[r_1?] := \text{first}[r_1]$
- ▷ $\text{next}[r] := \{i \in \mathbb{N} \mid (r \bullet, \epsilon, \bullet [i, x]) \in \delta^*, x \neq \epsilon\}$

Hier wird induktiv andersherum, also von oben, definiert:

 - Für die Wurzel e ist $\text{next}[e] := \emptyset$.
 - Ist der übergeordnete Knoten $r = \dots r_1 | r_2$, dann:
 $\text{next}[r_1] := \text{next}[r]$,
 $\text{next}[r_2] := \text{next}[r]$
 - $\dots r_1 \cdot r_2$, dann:
 $\text{next}[r_1] := \begin{cases} \text{first}[r_2] \cup \text{next}[r] & \text{empty}[r_2] = \text{true} \\ \text{first}[r_2] & \text{empty}[r_2] = \text{false} \end{cases}$,
 $\text{next}[r_2] := \text{next}[r]$
 - $\dots r_1^*$, dann:
 $\text{next}[r_1] := \text{first}[r_1] \cup \text{next}[r]$
 - $\dots r_1?$, dann:
 $\text{next}[r_1] := \text{next}[r]$
- ▷ $\text{last}[r] := \{i \text{ in } r \mid ([i, x] \bullet, \epsilon, r \bullet) \in \delta^*, x \neq \epsilon\}$
 - Für Blätter $r = [i, x]$ ist $\text{last}[r] := \{i \mid x \neq \epsilon\}$
 - $\text{last}[r_1 | r_2] := \text{last}[r_1] \cup \text{last}[r_2]$
 - $\text{last}[r_1 \cdot r_2] := \begin{cases} \text{last}[r_1] \cup \text{last}[r_2] & \text{empty}[r_2] = \text{true} \\ \text{last}[r_2] & \text{empty}[r_2] = \text{false} \end{cases}$
 - $\text{last}[r_1^*] := \text{last}[r_1]$
 - $\text{last}[r_1?] := \text{last}[r_1]$

Der **Berry-Sethi-/Glushkow-NFA** A_e ist dann:

Zustände: $\{q_0\} \cup \{i \in \mathbb{N} \mid i \text{ in Blatt } [i, x], x \neq \epsilon\}$
Startzustand: $I := \{q_0\}$

Endzustände: Falls $\text{empty}[e] = \text{false}$, dann $F := \text{last}[e]$; sonst $F := \{q_0\} \cup \text{last}[e]$.

(wobei \cup die disjunkte Vereinigung ist)

Übergänge: $\forall i, i' \in \mathbb{N} :$

- ▷ $(q_0, x, i) \in \delta$, falls $i \in \text{first}[e]$ und das Blatt i mit x beschriftet ist.
- ▷ $(i, y, i') \in \delta$, falls $i' \in \text{next}[[i, x]]$, das Blatt i mit x und i' mit y beschriftet ist.

• **Satz:** Zu jedem NFA $A = (Q, \Sigma, \delta, I, F)$ kann ein DFA $\mathcal{P}(A) = (\text{Pot}(Q), \Sigma, \delta_{\mathcal{P}}, \{i\}, F_{\mathcal{P}})$ konstruiert werden mit $\mathcal{L}(A) = \mathcal{L}(\mathcal{P}(A))$.

Teilmengen-/Potenzmengenkonstruktion:

Zustände: $\text{Pot}(Q)$ (alle Teilmengen von Q)

Startzustand: $i := I \in \text{Pot}(Q)$

Endzustände: $F_{\mathcal{P}} := \{Q' \subseteq Q \mid Q' \cap F \neq \emptyset\}$

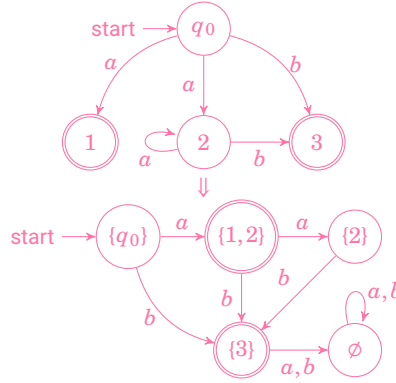
Übergangsfunktion:

$\delta_{\mathcal{P}}(Q', a) := \{q \in Q \mid \exists p \in Q' : (p, a, q) \in \delta\}$

Um die Anzahl Zustände mgl. klein zu halten,

beginnen wir mit $Q_{\mathcal{P}} = \{i\}$ und fügen weitere Zustände nur nach Bedarf hinzu.

Bsp.:



• Viele Tools wie `grep` bauen den DFA nur während der Abarbeitung auf, mit nur genau den Zuständen, wie sie für die jeweilige Eingabe mit Länge n nötig sind. $\Rightarrow \mathcal{O}(n)$

• \Rightarrow **Satz:** $\forall e \in \mathcal{E}_{\Sigma}$ kann ein DFA $A = \mathcal{P}(A_e)$ konstruiert werden mit $\mathcal{L}(A) = \llbracket e \rrbracket$.

1.1.3 Design eines Scanners

• **Eingabe:** eine Menge von Regeln:

- e_1 {action₁}
- e_2 {action₂}
- \vdots
- e_k {action_k}

mit $k \in \mathbb{N}$, $e_i \in \mathcal{E}_{\Sigma}$ und $\epsilon \notin \llbracket e_i \rrbracket \forall i \in \mathbb{k}$.

Ausgabe: ein Programm, das

- ▷ von der Eingabe ein *maximales Präfix* w liest, das $e_1 | \dots | e_k$ erfüllt,
- ▷ das minimale i ermittelt mit $w \in \llbracket e_i \rrbracket$,
- ▷ für w dann action_i ausführt.

• **Konstruktion:**

1. Konstruiere den DFA $\mathcal{P}(A_e)$ (Pot.-A. des Berry-Sethi-NFAs) zu $e := e_1 | e_2 | \dots | e_k$.
2. Definiere die Mengen:
 - ▷ $F_1 := \{q \in F \mid q \cap \text{last}[e_1] \neq \emptyset\}$
 - ▷ $F_2 := \{q \in (F \setminus F_1) \mid q \cap \text{last}[e_2] \neq \emptyset\}$
 - \vdots
 - ▷ $F_k := \{q \in (F \setminus (F_1 \cup \dots \cup F_{k-1})) \mid q \cap \text{last}[e_k] \neq \emptyset\}$

Für Eingabe w gilt: $\delta_{\mathcal{P}}^*(q_0, w) \in F_i \Leftrightarrow$ der Scanner für w action_i ausführen soll.

• **Tokenizing-/Reps' Maximal-Munch-Algorithmus**

INPUT: String $w \in \Sigma^+$

$s := 1; k := 1$

FORALL $q \in Q_{\mathcal{P}}, i \in |w| + 1$ **DO**

fehlversuch $[q, i] := \text{false}$

WHILE **true** **DO**

$q := q_0; p := \perp; S := \{(q, k)\}$

WHILE $k \leq |w|$ **AND** $\delta(q, w[k]) \neq \emptyset$ **DO**

IF **fehlversuch** $[q, k] = \text{true}$ **THEN**

BREAK

$q := \delta(q, w[k]); k := k + 1$

IF $q \in F$ **THEN**

$p = q; j := k; S := \emptyset$

ELSE $S := S \cup \{(q, k)\}$

ENDWHILE

IF $p = \perp$ **THEN** **RETURN** (*failure*)

FORALL $(q, k) \in S$ **DO**

fehlversuch $[q, k] := \text{true}$

 let i such that $p \in F_i$

 write($w[s, j - 1]$)

 action_i()

$s := j; k := j$

IF $j = |w| + 1$ **THEN** **RETURN**

ENDWHILE

Der **Maximal-Munch-Algorithmus** benötigt nur $\mathcal{O}(|Q| \cdot |w|)$, Tokenizing braucht $\mathcal{O}(|w|^2)$.

• In unterschiedlichen **Scannerzuständen** können verschiedene Token-Klassen erkannt werden, **Bsp.:** **Kommentare.**

Eingabe: eine Menge von Regeln:

```
(state) {
  e1 {action1 yybegin(state1)}
  e2 {action2 yybegin(state2)}
  :      :      :
  ek {actionk yybegin(statek)}
}
```

Bsp.:

```
<YYINITIAL> {
  /* {yybegin(COMMENT)}
}
<COMMENT> {
  /* {yybegin(YYINITIAL)}
  .\n { }
}
```

1.1.4 Implementierung von DFAs

• Die **Klassifizierungsfunktion** $r : Q \rightarrow \mathbb{N}$ ist für $\mathcal{P}(A_e)$ definiert durch $r(q) := \begin{cases} 0 & q \notin F \\ i & q \in F_i \end{cases}$

• **Def.** Äquivalenzrel. auf Zuständen: $p \equiv_r q \Leftrightarrow \forall w \in \Sigma^* : r(\delta^*(p, w)) = r(\delta^*(q, w))$

Konstruktion des minimierten DFA für $\mathcal{L}(A)$:

Zustände: $[q]_r, \forall q \in Q$

Anfangszustand: $[q_0]_r$

Klassifizierung: $r([q]_r) := r(q)$

Übergangsfunktion: $\delta([p]_r, a) := [\delta(p, a)]_r$

Dazu: Konstruktion der Äquiv.-klassen:

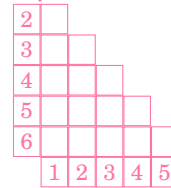
1. $\bar{Q} := \{r^{-1}(i) \mid i \in \mathbb{N}, r^{-1}(i) \neq \emptyset\}$
(Partition von Q in die Fasern von r , **dabei ist die Faser** $r^{-1}(i) := \{q \in Q \mid r(q) = i\}$, d. h. Zustände gleicher Klassifizierung.)
2. Teste $\forall \bar{q} \in \bar{Q} \forall p_1, p_2 \in \bar{q} \forall a \in \Sigma$, ob $r(\delta(p_1, a)) \neq r(\delta(p_2, a))$. In dem Fall muss \bar{q} entsprechend aufgeteilt werden.

Naive Implementierung: $\mathcal{O}(n^2)$, Optimierung auf $\mathcal{O}(n \cdot \log n)$ möglich.

• **Konstruktion des minimierten DFA, alt. Alg.:**

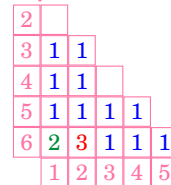
1. Starte Tabelle aller Zustandspaare (treppenförmig, horiz. letzten Zustand weglassen, vert. 1. Zustand weglassen)

Bsp.:



2. Markiere $\{q, q'\}$, falls:
 - ▷ Erkennungsäq. (GTI): $q \in F, q' \notin F$
 - ▷ r -Äquivalenz: $r(q) \neq r(q')$
3. \forall unmarkierten Paare: teste $\forall a \in \Sigma$, ob $\{\delta(q, a), \delta(q', a)\}$ bereits markiert. Ja \Rightarrow markiere $\{q, q'\}$.

Bsp.:

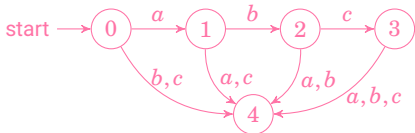


1. Schritt
2. Schritt
3. Schritt

4. Wdh., bis keine Änd. der Tabelle mehr.
5. \forall jetzt noch unmark. Paare gilt: $q \equiv q'$
6. Verschmelze je äquivalente Zustands-paare, vereinige Übergänge der beiden.

- Speichern der Übergangsfunktion als Tabelle indiziert mit $Q \times \Sigma$.

Bsp.:



δ	0	1	2	3	4
a	1	4	4	4	4
b	4	2	4	4	4
c	4	4	3	4	4

- Reduktion der Tabellengröße
 - durch Zusammenfassung von Zeichen in Zeichenklassen, Bsp.: $[a-zA-Z_]$
 - häufigsten Wert nicht speichern & durch „default“ ersetzen, Bsp.:

δ	0	1	2	3	4
a	1				
b		2			
c			3		

- Übereinanderlegen der Zeilen, Bsp.:

δ	0	1	2
A	1	2	3
valid	a	b	c

Passen die Zeilen nicht direkt übereinander, schieben wir diese nach rechts, bis sie passen, und merken uns die Verschiebung als displacement. Je weniger zusätzliche Spalten nötig sind, desto besser. Bsp.:

δ	0	1	2	3	4
a	1		2	1	
b	3	2		0	

↓

δ	0	1	2	3	4	5
a			1		2	1
b	3	2		0		

displacement(a) := 2,
displacement(b) = 0

1.2 Syntaktische Analyse

- Zusammenfassen von Tokens mit Parser zu größeren Programmeinheiten, Bsp.:
 - Ausdrücke
 - Statements
 - bedingte Verzweigungen (if)
 - Schleifen (while, for), ...
- Auch Parser werden i. A. nicht von Hand progr., sondern aus Spezifikation generiert.
- Spezifikation: kontextfreie Grammatiken mit Terminalalphabet = Menge der Token-Klassen
Generierte Impl.: Kellerautomaten + X

1.2.1 Grundlagen: Kontextfreie Grammatiken

- Def.: **kontextfr. Grammatik:** $G = (N, T, P, S)$
 - N : Menge der Nichtterminale
 - T : endl. Menge der Terminale
 - P : Menge der Produktionen/Regeln der Form $A \rightarrow \alpha$ mit $A \in N, \alpha \in (N \cup T)^*$
 - $S \in N$ Startsymbol

Konvention: In Beispielen ist Spezifikation von N und T i. A. implizit anhand ggb. Regeln P .

Bsp.: $S \rightarrow aSb$
 $S \rightarrow \epsilon$

spezifizierte Sprache: $\{a^n b^n \mid n \geq 0\}$

- Sammele & nummeriere \forall Nichtterminal A die rechten Seiten: (A, j) ist j -te Regel von A .
Bsp. KfGI: $E \rightarrow E + T^0 \mid T^1$
 $T \rightarrow T * F^0 \mid F^1$
 $F \rightarrow (E)^0 \mid name^1 \mid int^2$
- Grammatiken sind **Wortersetzungssysteme**. Folge v. Ersetzungsschritten heißt **Ableitung**.

Bsp.: bzgl. Grammatik KfGI:

- $E \rightarrow E + T$ mit Regel 0
- $\rightarrow T + T$ mit Regel 1
- $\rightarrow T * F + T$ mit Regel 0
- $\rightarrow T * int + T$ mit Regel 2
- $\rightarrow F * int + T$ mit Regel 1
- $\rightarrow name * int + T$ mit Regel 1
- $\rightarrow name * int + F$ mit Regel 1
- $\rightarrow name * int + int$ mit Regel 2

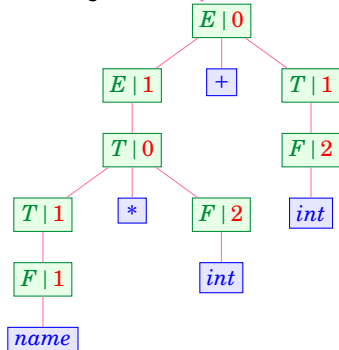
Formal: \rightarrow ist Relation auf $(N \cup T)^*$, wobei $\alpha \rightarrow \alpha' \Leftrightarrow \alpha = \alpha_1 A \alpha_2 \wedge \alpha' = \alpha_1 \beta \alpha_2$ für ein $A \rightarrow \beta \in P, A \in N, \alpha, \alpha', \alpha_1, \alpha_2, \beta \in (N \cup T)$.

Reflexiv-transitiver Abschluss von \rightarrow sei \rightarrow^* .

Bemerkungen:

- \rightarrow hängt von Grammatik ab.
- In jedem Ableitungsschritt kann Stelle & Regel zum Ersetzen ausgesucht werden.
- Von Grammatik G spezifizierte Sprache ist $\mathcal{L}(G) := \{w \in T^* \mid S \rightarrow^* w\}$.

• **Ableitungsbaum:** Bsp.: $t :=$



innere Knoten: Regel-Anwendungen,
Wurzel: Regel-Anwendung für A ,
Blätter: Terminale oder ϵ
Kinder eines Knotens (B, i) entsprechen der rechten Seite der i -ten Regel von B .

- Links-/Rechts-Ableitungen** sind solche, bei denen man stets das linke bzw. rechte Vorkommen eines Nichtterminals ersetzt.
 - entsprechen links-rechts- bzw. rechts-links-preorder-DFS-Durchlauf durch den Ableitungsbaum.
 - Reverse Rechts-Abl.** entsprechen links-rechts-postorder-DFS-Durchlauf.

DFS (depth-first search): Tiefensuche

- Konkatenation der Blätter eines Ableitungsbaums t heißt yield(t).
Bsp.: yield(t) = name * int + int
- Die Grammatik G heißt **eindeutig**, falls es $\forall w \in T^*$ max. einen Ableitungsbaum t von S gibt mit yield(t) = w .

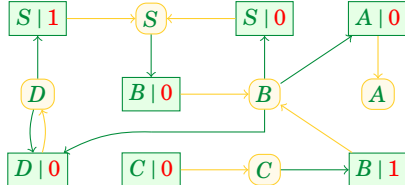
1.2.2 Überflüssige Nichtterminale & Regeln

- Def.: $A \in N$ heißt...
 - produktiv**, falls $A \rightarrow^* w$ für ein $w \in T^*$.
 - erreichbar**, falls $S \rightarrow^* \alpha A \beta$ für geeignete $\alpha, \beta \in (T \cup N)^*$.

Bsp.: $S \rightarrow aBB \mid bD$
 $A \rightarrow Bc$
 $B \rightarrow Sd \mid C$
 $C \rightarrow a$
 $D \rightarrow BD$

Produktiv: S, A, B, C ; Erreichbar: S, B, C, D

• Für Produktivität: **And-Or-Graph**. Bsp.:



And-Knoten: ■ (nummerierte) Regeln

Or-Knoten: • Nichtterminale

Kanten: \rightarrow ($(B, i), B$) für alle Regeln (B, i) ,
 \nearrow ($A, (B, i)$) falls $(B, i) \equiv B \rightarrow \alpha_1 A \alpha_2$

Algorithmus für Produktivität:

```
Subset(N) result := ∅ (d.h. result ⊆ N)
int count[P]
Subset(P) rhs[N]
FORALL A ∈ N DO rhs[A] := ∅
FORALL (A, i) ∈ P DO
  count[(A, i)] := 0
  init(A, i)
  ↳ gehe rechte Seite der Produktion (A, i) durch. Für jedes versch. Nichtt. X darin:
  count[(A, i)] += und rhs[X].add((A, i)).
```

ENDFORALL

```
Subset(P) W := {r | count[r] = 0}
WHILE W ≠ ∅ DO
```

Sei $(A, i) \in W$ beliebig.
Entferne (A, i) aus W .

IF $A \notin result$ THEN

result := result \cup {A}

FORALL $r \in rhs[A]$ DO

count[r] --

IF count[r] = 0 DO $W := W \cup \{r\}$

ENDIF

ENDFORALL

ENDIF

ENDWHILE

Alle prod. Nichtterminale sind nun in result.

Verwende Tabelle: Bsp.:

	$S \rightarrow aBB$...	$C \rightarrow a$	$D \rightarrow BD$	W	result
$C \rightarrow a$	1	...	0	2	$C \rightarrow a$	\emptyset
$B \rightarrow C$	1	...	0	2	$B \rightarrow C$	$\{C\}$
...
$B \rightarrow Sd$	0	...	0	1	\emptyset	$\{S, A, B, C\}$

Spalten: Alle Regeln, W , result.

Zeilen: nacheinander Regeln aus W

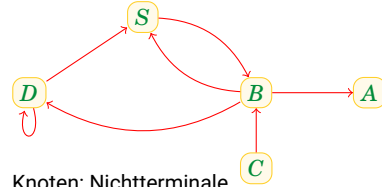
Zellen: count

Eigenschaften des Algorithmus:

- Initialisierung erfordert lineare Laufzeit.
- Jede Regel wird maximal einmal in W eingefügt.
- Jede Regel in W hat auf der linken Seite ein produktives Nichtterminal stehen.
- Jedes A wird maximal einmal in result eingefügt.
- ⇒ Laufzeit linear in der Größe der Grammatik.
- Um den Test „ $A \in result$ “ einfach zu machen, repräsentiert man die Menge result durch ein Array.
- W wie auch die Mengen rhs[A] werden dagegen als Listen repräsentiert.

Der Algorithmus funktioniert auch, um kleinste Lösungen von Boole'schen Ungleichungssystemen zu bestimmen.

- Leerheitsproblem:** $\mathcal{L}(G) \neq \emptyset \Leftrightarrow S$ produktiv
- Für Erreichbarkeit: **Abhängigkeitsgraph**. Bsp.:



Knoten: Nichtterminale

Kanten: (A, B) , falls $B \rightarrow \alpha_1 A \alpha_2 \in P$

A ist erreichbar, falls es im Abhängigkeitsgraphen einen Pfad von A nach S gibt.

- Erreichbarkeit kann mithilfe von DFS in linearer Zeit berechnet werden.
- Damit ist Berechnung von produktiven & erreichbaren Nichtterminalen linear.

• Def.: Grammatik G heißt **reduziert**, wenn alle Nichtterminale von G produktiv & erreichbar.

• Satz: Zu jeder kontextfreien Grammatik G mit $\mathcal{L}(G) \neq \emptyset$ kann in linearer Zeit eine reduzierte Gr. G' konstruiert werden mit $\mathcal{L}(G') = \mathcal{L}(G)$.

Konstruktion der reduzierten Grammatik:

1. Berechne die Teilmenge $N_1 \subseteq N$ aller produktiven Nichtterminale von G .
Da $\mathcal{L}(G) \neq \emptyset$, ist insbesondere $S \in N_1$.
2. Entferne unproduktive Regeln:
 $P_1 := \{A \rightarrow \alpha \in P \mid A \in N_1 \wedge \alpha \in (N_1 \cup T)^*\}$
3. Berechne die Teilmenge $N_2 \subseteq N_1$ aller produktiven & erreichbaren Nichtt. von G .
Da $\mathcal{L}(G) \neq \emptyset$, ist insbesondere $S \in N_2$.
4. Entferne unerreichbare Regeln:
 $P_2 := \{A \rightarrow \alpha \in P \mid A \in N_2 \wedge \alpha \in (N_2 \cup T)^*\}$

Ergebnis: $G' := (N_2, T, P_2, S)$

1.2.3 Grundlagen: Kellerautomaten

- **Def.: Kellerautomat (PDA, push-down automaton):** $M = (Q, T, \delta, q_0, F)$ mit:
 - ▷ Q endl. Menge von Zuständen (bzw. Kellersymbolen),
 - ▷ T Eingabealphabet,
 - ▷ $q_0 \in Q$ Anfangszustand,
 - ▷ $F \subseteq Q$ Menge der Endzustände,
 - ▷ $\delta \subseteq Q^+ \times (T \cup \{\epsilon\}) \times Q^*$ endl. Menge von Übergängen (gelesener Kellerinhalt, gel. Eingabe, hinzuzufügender Kellerinhalt)
- **Def.: Konfiguration** des PDA: $(\gamma, w) \in Q^* \times T^*$
Ein **Berechnungsschritt** des PDA wird durch die Relation $\vdash \subseteq (Q^* \times T^*)^2$ beschrieben, wobei $(\alpha\gamma, xw) \vdash (\alpha\gamma', w)$ für $(\gamma, x, \gamma') \in \delta$.
 - ▷ Relation \vdash hängt vom PDA ab.
 - ▷ Reflexiv-transitive Hülle bezeichnen wir mit \vdash^* .
- Die von M akzeptierte Sprache ist:
 $\mathcal{L}(M) = \{w \in T^* \mid \exists f \in F : (q_0, w) \vdash^* (f, \epsilon)\}$
- PDA M heißt **deterministisch**, falls \forall paarw. versch. Übergänge $(\gamma_1, x, \gamma_2), (\gamma_1', x', \gamma_2') \in \delta$ gilt: Ist γ_1 Suffix von $\gamma_1' \Rightarrow x \neq x' \wedge x \neq \epsilon \neq x'$.
- **Satz:** \forall kontextfr. Grammatik $G = (N, T, P, S)$ kann PDA konstr. werden mit $\mathcal{L}(G) = \mathcal{L}(M)$.

Konstruktion 1: Shift-Reduce-Parser

- ▷ Eingabe sukzessive auf Keller schieben.
- ▷ Liegt oben auf dem Keller eine vollständige rechte Seite (*Handle*), wird dieses durch die zugehörige linke Seite ersetzt (reduziert).

$M_G^{(1)} = (Q, T, \delta, q_0, \{f\})$ mit

- ▷ $Q := T \cup N \cup \{q_0, f\}$
- ▷ $\delta := \{(q, x, qx) \mid q \in Q, x \in T\} \cup \{(q\alpha, \epsilon, qA) \mid q \in Q, A \rightarrow \alpha \in P\} \cup \{(q_0S, \epsilon, f)\}$

Eigenschaften:

- ▷ Folge der Reduktionen entspricht einer reversen Rechts-Abl. für Eingabe.
- ▷ Zur Korrektheit:
Es gilt $(q, w) \vdash^* (qA, \epsilon) \Leftrightarrow A \rightarrow^* w$
- ▷ $M_G^{(1)}$ i. A. nicht-deterministisch.
Deterministisches Verfahren: identifizierbare Reduktionsstellen \rightarrow **LR-Parsing**

Konstruktion 2: Item-Kellerautomat

- ▷ Rekonstruiere eine Linksableitung.
- ▷ Expand. Nichtterminale mithilfe Regel.
- ▷ Verifiziere sukzessive, dass die gewählte Regel mit der Eingabe übereinstimmt. \Rightarrow die Zustände sind jetzt *Items*.
- ▷ Ein *Item* ist eine Regel mit Punkt:
 $[A \rightarrow \alpha \bullet \beta], A \rightarrow \alpha \beta \in P$
Der Punkt gibt an, wie weit die Regel bereits abgearbeitet wurde.

$M_G^{(2)} = (Q, T, \delta, q_0, \{f\})$

- ▷ $Q := \{\text{alle Items}\} \cup \{[S' \rightarrow \bullet S], [S' \rightarrow S \bullet]\}$
mit neuem Terminalsymbol S' .

- ▷ $q_0 := [S' \rightarrow \bullet S]$
- ▷ $f := [S' \rightarrow S \bullet]$
- ▷ Übergänge:
Expansionen:
 $([A \rightarrow \alpha \bullet B\beta], \epsilon, [A \rightarrow \alpha \bullet B\beta] [B \rightarrow \bullet \gamma])$
für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$,
Shifts:
 $([A \rightarrow \alpha \bullet a\beta], a, [A \rightarrow \alpha a \bullet \beta])$
für $A \rightarrow \alpha a\beta \in P$,
Reduce:
 $([A \rightarrow \alpha \bullet B\beta] [B \rightarrow \gamma \bullet], \epsilon, [A \rightarrow \alpha B \bullet \beta])$
für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$,

Items der Form $[A \rightarrow \alpha \bullet]$ heißen **vollständig**.

Eigenschaften:

- ▷ Expansionen bilden Linksableitung.
- ▷ Expansionen nicht-deterministisch.
Durch Vorausschau deterministisch machen \rightarrow **LL-Parsing**.
- ▷ Zur Korrektheit:
Es gilt \forall Item $[A \rightarrow \alpha \bullet B\beta]$:
 $([A \rightarrow \alpha \bullet B\beta], w) \vdash^* ([A \rightarrow \alpha B \bullet \beta], \epsilon) \Leftrightarrow B \rightarrow^* w$

1.2.4 Vorausschau-Mengen

- **Def.:** Für Menge $L \subseteq T^*$ sei $\text{First}_k(L) := \{u \in L \mid |u| < k\} \cup \{u \in T^k \mid \exists v \in T^* : uv \in L\}$ die **Menge der Präfixe** der Länge k .
- $\text{First}_k(_)$ ist verträglich mit Vereinigung und Konkatination:
▷ $\text{First}_k(\emptyset) = \emptyset$
▷ $\text{First}_k(L_1 \cup L_2) = \text{First}_k(L_1) \cup \text{First}_k(L_2)$
▷ $\text{First}_k(L_1 \cdot L_2) = \text{First}_k(\text{First}_k(L_1) \cdot \text{First}_k(L_2)) = \text{First}_k(L_1) \circ \text{First}_k(L_2)$
- $T^{\leq k} := \bigcup_{0 \leq i \leq k} T^i, \mathbb{D}_k := 2^{T^{\leq k}}$
- $\circ : \mathbb{D}_k \times \mathbb{D}_k \rightarrow \mathbb{D}_k$ distributiv in jedem Arg.:
▷ $L \circ \emptyset = \emptyset, \emptyset \circ L = \emptyset$
▷ $L \circ (L_1 \cup L_2) = (L \circ L_1) \cup (L \circ L_2)$
▷ $(L_1 \cup L_2) \circ L = (L_1 \circ L) \cup (L_2 \circ L)$
- Für $\alpha \in (N \cup T)^*$ sei
 $\text{First}_k(\alpha) := \text{First}_k(\{w \in T^* \mid \alpha \rightarrow^* w\})$
- Für $k \geq 1$ gilt:
▷ $\text{First}_k(x) = \{x\}$ für $x \in T \cup \{\epsilon\}$
▷ $\text{First}_k(\alpha_1\alpha_2) = \text{First}_k(\alpha_1) \circ \text{First}_k(\alpha_2)$

- **Satz:** Die Mengen $\text{First}_k(A)$ für $A \in N$ sind die kleinste Lösung des Ungleichungssystems $\text{First}_k(A) \supseteq \text{First}_k(X_1) \circ \dots \circ \text{First}_k(X_m)$ für alle $A \rightarrow X_1 \dots X_m \in P$.

Bsp.: bzgl. Grammatik KfGI:

- ▷ $\text{First}_2(E) \supseteq \text{First}_2(E+T) \supseteq \text{First}_2(E) \circ (+) \circ \text{First}_2(T)$
- ▷ $\text{First}_2(T) \supseteq \text{First}_2(T * F) \supseteq \text{First}_2(T) \circ (*) \circ \text{First}_2(F)$
- ▷ $\text{First}_2(F) \supseteq \text{First}_2(E) \supseteq (\circ \text{First}_2(E) \circ \circ)$
- ▷ $\text{First}_2(E) \supseteq \text{First}_2(T)$
- ▷ $\text{First}_2(T) \supseteq \text{First}_2(F)$
- ▷ $\text{First}_2(F) \supseteq \{\text{name int}\}$

Berechnung: Fixpunktiteration (s. 1.2.5), denn:

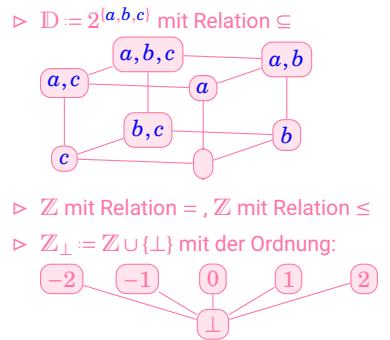
- ▷ Die Menge \mathbb{D}_k der möglichen Werte für $\text{First}_k(A)$ ist vollständiger Verband.
- ▷ Operatoren auf den rechten Seiten sind monoton, d. h. hier: verträglich mit \subseteq .

1.2.5 Exkurs: Vollständige Verbände

Dieses Kapitel ist nicht klausurrelevant in SS'20.

- **Def.:** Eine Menge \mathbb{D} mit Relation $\sqsubseteq \subseteq \mathbb{D} \times \mathbb{D}$ ist ein **Verband (lattice)**, falls $\forall a, b, c \in \mathbb{D}$:
▷ $a \sqsubseteq a$ (reflexiv)
▷ $a \sqsubseteq b \wedge b \sqsubseteq a \Rightarrow a = b$ (antisymmetrisch)
▷ $a \sqsubseteq b \wedge b \sqsubseteq c \Rightarrow a \sqsubseteq c$ (transitiv)

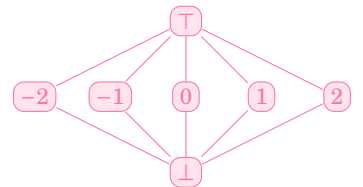
Bsp.:



- **Def.:** $d \in \mathbb{D}$ heißt **obere Schranke** für $X \subseteq \mathbb{D}$, falls $x \sqsubseteq d \forall x \in X$.
- **Def.:** d heißt **kleinste obere Schranke (lowest upper bound, lub)** von $X \subseteq \mathbb{D}$, falls:
▷ d eine obere Schranke für X ist
▷ $d \sqsubseteq y$ für jede obere Schranke y für X
Schreibe $d =: \sup X =: \bigsqcup X$.
- **Def.:** Ein **vollständiger Verband (complete lattice, cl)** \mathbb{D} ist eine Halbordnung \sqsubseteq , in der jede Teilmenge $X \subseteq \mathbb{D}$ eine *lub* $\bigsqcup X \in \mathbb{D}$ besitzt.

Bsp.:

- ▷ $\mathbb{D} := 2^{\{a,b,c\}}$ ist ein *cl*.
- ▷ \mathbb{Z} mit $=$ oder \leq sowie \mathbb{Z}_\perp sind kein *cl*!
- ▷ $\mathbb{Z}_\perp := \mathbb{Z} \cup \{\perp, \top\}$ ist der **flache Verband**:



- **Satz:** In jedem *cl* \mathbb{D} besitzt jedes $X \subseteq \mathbb{D}$ eine **größte untere Schranke** $\bigsqcap X := \inf X \in \mathbb{D}$.
Und zwar gilt: $\bigsqcap X = \bigsqcup U_X$, wobei $U_X := \{u \in \mathbb{D} \mid \forall x \in X : u \sqsubseteq x\} \subseteq \mathbb{D}$ die Menge der unteren Schranken von X .
- **Def.:** Abbildung $f : \mathbb{D}_1 \rightarrow \mathbb{D}_2$ heißt **monoton**, falls $\forall a, b \in \mathbb{D}_1 : a \sqsubseteq_1 b \Rightarrow f(a) \sqsubseteq_2 f(b)$.
Bsp.: $\mathbb{D}_1 = \mathbb{D}_2 = \mathbb{Z}$ mit Ordnung \leq . Dann sind $\text{inc } x := x + 1$ und $\text{dec } x := x - 1$ monoton; $\text{inv } x := -x$ ist *nicht* monoton.

Ungleichungssystem über Verbänden:

GESUCHT: möglichst kleine Lösung für $x_i \supseteq f_i(x_1, \dots, x_n), i \in n = \{1, 2, \dots, n\}$, wobei $f_i : \mathbb{D}^n \rightarrow \mathbb{D}$ monoton $\forall i \in n$.

Lösung: Fixpunktiteration

Sei $F : \mathbb{D}^n \rightarrow \mathbb{D}^n, F(x_1, \dots, x_n) = (y_1, \dots, y_n)$ mit $y_i := f_i(x_1, \dots, x_n) \forall i \in n$.

Sind alle f_i monoton, so ist auch F monoton. Konstruiere $\perp, F\perp, F^2\perp, F^3\perp, \dots$

Gilt $F^k\perp = F^{k+1}\perp$, so ist eine Lösung gefunden. Ist \mathbb{D} endlich, so gibt es k immer.

Bsp. FII: $\mathbb{D} := 2^{\{a,b,c\}}, \sqsubseteq = \subseteq$
 $x_1 \supseteq \{a\} \cup x_3, x_2 \supseteq x_3 \cap \{a, b\}, x_3 \supseteq x_1 \cup \{c\}$
Iteration:

	0	1	2	3	4
x_1	\emptyset	$\{a\}$	$\{a, c\}$	$\{a, c\}$	$\{a, c\}$
x_2	\emptyset	\emptyset	\emptyset	$\{a\}$	$\{a\}$
x_3	\emptyset	$\{c\}$	$\{a, c\}$	$\{a, c\}$	$\{a, c\}$

Noch zu klären: *kleinste Lösung?*

- **Def.:** Eine Funktion $f : \mathbb{D} \rightarrow \mathbb{D}$ heißt **stetig**, falls \forall aufsteigende Kette $d_0 \sqsubseteq \dots \sqsubseteq d_m \sqsubseteq \dots$ gilt: $f(\bigsqcup_{m \geq 0} d_m) = \bigsqcup_{m \geq 0} (f(d_m))$.
- Jede stetige Funktion ist auch monoton.
- **Def.:** Eine Halbordnung $(\mathbb{D}, \sqsubseteq)$ heißt **vollständig (complete partial order, CPO)**, falls alle aufsteigenden Ketten ein *lub* haben.
- Jeder vollständige Verband ist eine *CPO*.
- **Satz (Kleene):** In einer vollständigen Halbordnung \mathbb{D} hat jede stetige Funktion $f : \mathbb{D} \rightarrow \mathbb{D}$

einen kleinsten Fixpunkt $d_0 = \bigcup_{k \geq 0} f^k \perp$.

- Für $P := \{x \in \mathbb{D} \mid x \sqsupseteq f(x)\}$ ist der kleinste Fixpunkt $d_0 = \min P$ (d. h. untere Schr. & in P).
- Die **kleinste Lösung eines Ungleichungssystems** ist der kleinste Fixpunkt von F .
- Round-Robin-Iteration (RR):**

Benutze bei der Iteration nicht die Werte der letzten Iteration, sondern die jeweils aktuellen.

Bsp.: Bezüglich Ungleichungssystem aus FI1:

	0	1	2	3
x_1	\emptyset	$\{a\}$	$\{a, c\}$	$\{a, c\}$
x_2	\emptyset	\emptyset	$\{a\}$	$\{a\}$
x_3	\emptyset	$\{a, c\}$	$\{a, c\}$	$\{a, c\}$

Algorithmus Round-Robin-Iteration:

FORALL $i \in \underline{n}$ DO $x_i := \perp$ ENDFORALL
REPEAT

```

finished := true
FORALL  $i \in \underline{n}$  DO
    new :=  $f_i(x_1, \dots, x_n)$ 
    IF  $\neg(x_i \sqsupseteq new)$ 
        finished := false
         $x_i := x_i \sqcup new$ 
    ENDFORALL
UNTIL finished
  
```

- Länge h der längsten echt aufsteigenden Kette nennen wir **Höhe** von \mathbb{D} .
- Bsp.:** Für First_k ist die Höhe des Verbands exponentiell in k .
⇒ man beschränkt sich i. A. auf kleine k .
- Anzahl Runden in RR-Iteration ist beschränkt durch $\mathcal{O}(n \cdot h)$, wobei n die Anzahl Variablen.
- Effizienz von RR-Iteration hängt von Anordnung der Variablen ab.

1.2.6 Top-down Parsing

- Def.:** Eine reduzierte kontextfr. Gr. heißt $\text{LL}(k)$ („**Left-to-right parsing, Leftmost derivation**“, Vorausschau der Länge k “), falls für verschiedene Regeln $A \rightarrow \alpha, A \rightarrow \alpha' \in P$ und jede Linksableitung $S \xrightarrow{*}_L uA\beta$ mit $u \in T^*$ gilt: $\text{First}_k(\alpha\beta) \cap \text{First}_k(\alpha'\beta) = \emptyset$.
- Bsp. KfG2:** $S \rightarrow \text{if}(E)S \text{ else } S \mid E$
 $E \rightarrow \text{id}$
- ist $\text{LL}(1)$.
- Bsp. KfG3:** $S \rightarrow \text{if}(E)S \text{ else } S \mid \text{if}(E)S \mid \text{while}(E)S \mid E$
 $E \rightarrow \text{id}$
- ist nicht $\text{LL}(k)$ für jedes $k > 0$, denn: Betrachte Regeln $(S, 0), (S, 1)$, d. h. $\alpha = \text{if}(E)S \text{ else } S$ und $\alpha' = \text{if}(E)S$, und die Linksableitung $S \xrightarrow{*}_L \text{if}(id)S$, d. h. $\beta = \epsilon$. Es gilt offensichtlich: $\text{First}_k(\alpha\beta) \cap \text{First}_k(\alpha'\beta) \supseteq \{\text{if}\} \neq \emptyset \forall k > 0$.

- LL(k)-Parser:** sieht Fenster der Länge k der Eingabe; realisiert i. W. Item-Kellerautomaten, tabelliert nächste zu wählende Regeln.

Bsp.: bzgl. KfG2:

Zustände: Items

Tabelle:	if	while	id
$[\dots \rightarrow \dots \bullet S \dots]$	0	1	2
$[\dots \rightarrow \dots \bullet E \dots]$	-	-	0

- Menge der möglichen nächsten k Zeichen ist: $\text{First}_k(\alpha\beta) = \text{First}_k(\alpha) \circ \text{First}_k(\beta)$, wobei: α die rechte Seite der passenden Regel, β ein möglicher rechter Kontext von A . $\text{First}_k(\beta)$ dynamisch akkumulieren:
- Ein **erweiterter Item** ist ein Paar $[A \rightarrow \alpha \bullet \gamma, L]$ für $A \rightarrow \alpha\gamma \in P, L \subseteq T^{\leq k}$. L wird als Vorausschau-Menge zum Repräsentieren von $\text{First}_k(\beta)$ für den rechten Kontext β von A verwendet.
- Erweiterter Item-Kellerautomat:**
Zustände: erweiterte Items
Anfangszustand: $[S' \rightarrow \bullet S, \{\epsilon\}]$

Endzustand: $[S' \rightarrow S \bullet, \{\epsilon\}]$

Übergänge:

Expansionen:

$([A \rightarrow \alpha \bullet B\beta, L], \epsilon,$
 $[A \rightarrow \alpha \bullet B\beta, L] [B \rightarrow \bullet \gamma, \text{First}_k(\beta) \circ L])$
für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$,

Shifts:

$([A \rightarrow \alpha \bullet a\beta, L], a, [A \rightarrow \alpha a \bullet \beta, L])$
für $A \rightarrow \alpha a\beta \in P$,

Reduce:

$([A \rightarrow \alpha \bullet B\beta, L] [B \rightarrow \bullet \gamma, L'], \epsilon,$
 $[A \rightarrow \alpha B \bullet \beta, L])$ für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$

- Vorausschau-Tabelle:**

$M([A \rightarrow \alpha \bullet B\beta, L], w) :=$
 $\{i \mid (B, i) = (B \rightarrow \gamma),$
 $w \in \text{First}_k(\gamma) \circ \text{First}_k(\beta) \circ L\}$

- Satz:** reduzierte kontextfr. Gr. G ist $\text{LL}(k)$ ⇔ \forall Eingabewort zu jedem Zeitpunkt in der Berechnung des erw. Item-Kellerautomaten $|M([A \rightarrow \alpha \bullet B\beta, L], w)| \leq 1$ gilt. Hierbei ist $[A \rightarrow \alpha \bullet B\beta, L]$ das aktuelle oberste Kellersymbol und w besteht aus den nächsten k Zeichen der Eingabe (bzw. $k' < k$ Zeichen, falls Rest der Eingabe nur noch k' lang).

- erw. Item-Kellerautomat mit k -Vorausschautabelle erlaubt deterministische Rekonstruktion einer Linksableitung für $\text{LL}(k)$ -Grammatik.

Bsp. KfG4: $S \rightarrow \epsilon \mid aSb$

Vorausschau-Tabelle:

	ϵ	a	b
$[S' \rightarrow \bullet S, \{\epsilon\}]$	0	1	-
$[S \rightarrow a \bullet Sb, \{\epsilon\}]$	-	1	0
$[S \rightarrow a \bullet Sb, \{b\}]$	-	1	0

- Anzahl der Vorausschau-Mengen L kann sehr groß sein.
- Hängt auszuwählende Regel nicht von den Erweiterungen der Items ab, kann man den Item-Kellerautomat auch ohne Erw. benutzen. Hängt auszuw. Regel nur von der aktuellen Vorausschau w ab, heißt G auch **stark LL(k)**.
- Bsp.:** Die Grammatik aus KfG4 ist anhand der Vorausschau-Tabelle offensichtlich **stark LL(k)**.
- Def.:** $\text{Follow}_k(A) := \bigcup \{ \text{First}_k(\beta) \mid S \xrightarrow{*}_L uA\beta \}$
- Def.:** Reduzierte kontextfr. Grammatik G heißt **stark LL(k)**, falls für je zwei verschiedene Regeln $A \rightarrow \alpha, A \rightarrow \alpha' \in P$ gilt: $\text{First}_k(\alpha) \circ \text{Follow}_k(A) \cap \text{First}_k(\alpha') \circ \text{Follow}_k(A) = \emptyset$
- Bsp.:** KfG4 ist tatsächlich **stark LL(k)**, denn:
 $\text{Follow}_1(S) = \{\epsilon, b\} \Rightarrow$
 $\text{First}_1(\epsilon) \circ \text{Follow}_1(S) = \{\epsilon\} \circ \{\epsilon, b\} = \{\epsilon, b\}$
 $\text{First}_1(aSb) \circ \text{Follow}_1(S) = \{a\} \circ \{\epsilon, b\} = \{a\}$

- Ist G **stark LL(k)**, können wir die Vorausschautabelle mit Nichtterminalen (statt erw. Items) indizieren:

$M[B, w] := \begin{cases} i & (B, i) = (B \rightarrow \gamma) \wedge \\ & w \in \text{First}_k(\gamma) \circ \text{Follow}_k(B) \\ - & \text{keine solche Regel existiert} \end{cases}$

Bsp.: bzgl. KfG4:

	ϵ	a	b
S	0	1	0

- Satz:**
▷ Jede **stark-LL(k)**-Gr. ist auch $\text{LL}(k)$.
Gilt i. A. nicht andersherum! Ausnahme:
▷ Jede $\text{LL}(1)$ -Gr. ist bereits **stark LL(1)**.
- Zu jeder $\text{LL}(k)$ -Grammatik kann eine äquivalente **starke LL(k)**-Gr. konstruiert werden.
- Berechnung von $\text{Follow}_k(B)$:**
Stelle Ungleichungssystem auf:
 $\text{Follow}_k(S) \supseteq \{\epsilon\}$ (gilt immer)
 $\text{Follow}_k(B) \supseteq$
 $\text{First}_k(X_1) \circ \text{First}_k(X_m) \circ \text{Follow}_k(A)$
für $A \rightarrow \alpha B X_1 \dots X_m \in P$

Kleinste Lösung ist $\text{Follow}_k(B)$.

- Größe der auftretenden Mengen steigt mit k rapide. In praktischen Systemen wird darum meist nur der Fall $k = 1$ implementiert.

1.2.7 Schnelle Berechnung von Vorausschau-Mengen

- Im Fall $k = 1$ lassen sich First_1 und Follow_1 besonders effizient berechnen.
- Seien $L_1, L_2 \subseteq T \cup \{\epsilon\}$ mit $L_1 \neq \emptyset \neq L_2$. Dann ist:
$$L_1 \circ L_2 = \begin{cases} L_1 & \epsilon \notin L_1 \\ (L_1 \setminus \{\epsilon\}) \cup L_2 & \text{sonst} \end{cases}$$
- Ist G reduziert, so sind $\forall A : \text{First}_1(A) \neq \emptyset$.
- Def.:** $\text{empty}(X) = \text{true} \Leftrightarrow X \xrightarrow{*} \epsilon$
- Def.:** die ϵ -freien **First₁-Mengen** seien:
 $F_\epsilon(a) := \{a\}$ für $a \in T$,
 $F_\epsilon(A) := \text{First}_1(A) \setminus \{\epsilon\}$ für $A \in N$.

- Berechnung:** Ungleichungssystem für $F_\epsilon(A)$:

$F_\epsilon(A) \supseteq F_\epsilon(X_j)$
falls $A \rightarrow X_1 \dots X_m \in P, j \in \underline{m}$,
 $\text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1})$.

Bsp.: bzgl. Grammatik aus KfG1:

$\text{empty}(E) = \text{empty}(T) = \text{empty}(F) = \text{false}$

Ungleichungssystem:

- ▷ $F_\epsilon(S') \supseteq F_\epsilon(E)$
- ▷ $F_\epsilon(E) \supseteq F_\epsilon(E)$
- ▷ $F_\epsilon(E) \supseteq F_\epsilon(T)$
- ▷ $F_\epsilon(T) \supseteq F_\epsilon(T)$
- ▷ $F_\epsilon(T) \supseteq F_\epsilon(F)$
- ▷ $F_\epsilon(F) \supseteq \{(\text{name}, \text{int})\}$

- Ungleichungssystem zu $\text{Follow}_1(A)$:**

$\text{Follow}_1(S) \supseteq \{\epsilon\}$
 $\text{Follow}_1(B) \supseteq F_\epsilon(X_j)$
falls $A \rightarrow \alpha B X_1 \dots X_m \in P$,
 $\text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1})$
 $\text{Follow}_1(B) \supseteq \text{Follow}_1(A)$
falls $A \rightarrow \alpha B X_1 \dots X_m \in P$,
 $\text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_m)$

Bsp.: bzgl. Grammatik aus KfG1:

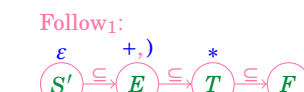
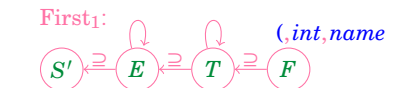
- ▷ $\text{Follow}_1(S') \supseteq \{\epsilon\}$
- ▷ $\text{Follow}_1(E) \supseteq \text{Follow}_1(S')$
- ▷ $\text{Follow}_1(E) \supseteq \{+, \cdot\}$
- ▷ $\text{Follow}_1(T) \supseteq \{*\}$
- ▷ $\text{Follow}_1(T) \supseteq \text{Follow}_1(E)$
- ▷ $\text{Follow}_1(F) \supseteq \text{Follow}_1(T)$

- Die Form dieser Ungleichungssysteme ist:
 $x \supseteq y$ bzw. $x \supseteq d$
für Variablen x, y und $d \in \mathbb{D}$. Solche Systeme heißen **reine Vereinigungsprobleme** und sind lösbar mit linearem Aufwand.

Berechnung:

- ▷ Konstruiere den **Variablen-Abhängigkeitsgraph** zum Ungleichungssystem.

Bsp.: bzgl. KfG1 (s. o. für deren Ungl.-S.):



- ▷ Innerhalb einer **starken Zusammenhangskomponente (SZK)**, d. h. einem maximalen Teilgraph, in dem man von jedem zu jedem Knoten kommen kann, haben alle Variablen den gleichen Wert.
- ▷ Hat eine SZK keine eingehenden Kanten, erhält man ihren Wert, indem man alle Werte innerhalb der SZK vereinigt.
- ▷ Gibt es eingehende Kanten, muss man zusätzlich die Werte an deren Startkno-

ten hinzufügen.

1.2.8 Bottom-up-Analyse

- Viele Grammatiken sind nicht LL(k)!
Wie parsen wir also diese?
- Def.:** Grammatik G heißt **links-rekursiv**, falls $A \rightarrow^+ A\beta$ für ein $A \in N$, $\beta \in (T \cup N)^*$.

Bsp.: KfG1 ist links-rekursiv.

- Satz:** Ist G reduziert & links-rekursiv, dann ist \bar{G} nicht LL(k) $\forall k$.

- $\alpha\gamma$ ist **zuverlässiges Präfix** für das vollst. Item $[B \rightarrow \gamma\bullet]$, wenn folgende Berechnung von $M_G^{(1)}$ existiert:

$$(q_0\alpha\gamma, v) \vdash (q_0\alpha B, v) \vdash^* (q_0S, \epsilon)$$

Dies ist der Fall genau dann, wenn $S \xrightarrow{*}_R \alpha Bv$.

- Das Item $[B \rightarrow \gamma\bullet\beta]$ heißt **gültig** für $a' \Leftrightarrow S \xrightarrow{*}_R \alpha Bv$ mit $a' = \alpha\gamma$.

- Berechnung der Menge zuverlässiger Präfixe** durch NFA, den **charakterist. Automat** $c(G)$:

Zustände: Items

$$\text{Anfangszustand: } [S' \rightarrow \bullet S]$$

$$\text{Endzustände: } \{[B \rightarrow \gamma\bullet] \mid B \rightarrow \gamma \in P\}$$

Übergänge:

- $\{[A \rightarrow \alpha \bullet X\beta], X, [A \rightarrow \alpha X \bullet \beta]\}$
für $X \in (N \cup T)$, $A \rightarrow \alpha X\beta \in P$
- $\{[A \rightarrow \alpha \bullet B\beta], \epsilon, [B \rightarrow \bullet \gamma]\}$
für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$

- Den **kanonischen LR(0)-Automaten** $LR(G)$ erhält man aus $c(G)$, indem man ϵ -Übergänge entfernt & Teilmengenkonstruktion anwendet.

Alternative Konstruktion von LR(G):

Hilfsfunktion $\delta_\epsilon^*(q) := q \cup$

$$\{[B \rightarrow \bullet \gamma] \mid \exists [A \rightarrow \alpha \bullet B'\beta'] \in q,$$

$$\beta \in (N \cup T)^* : B' \xrightarrow{*} B\beta\}$$

Dann, DFA:

Zustände: Mengen von Items

$$\text{Anfangszustand: } \delta_\epsilon^*\{[S' \rightarrow \bullet S]\}$$

$$\text{Endzustände: } \{q \mid \exists A \rightarrow \alpha \in P : [A \rightarrow \alpha\bullet] \in q\}$$

$$\text{Übergänge: } \delta(q, X) := \delta_\epsilon^*\{[A \rightarrow \alpha X \bullet \beta] \mid [A \rightarrow \alpha \bullet X\beta] \in q\}$$

↓ Alles folgende nicht mehr klausurrelevant in SS'20.

- LR(0)-Parser:** Sei $LR(G) = (Q, T, \delta, q_0, F)$. Dann ist der LR(0)-Parser:

Zustände: $Q \cup \{f\}$ (wobei f neu)

Anfangszustand: q_0 Endzustand: f

Übergänge:

Shift: (q, a, pq) , falls $q = \delta(p, a) \neq \emptyset$

Reduce: $(pq_1 \dots q_m, \epsilon, pq)$,

falls $[A \rightarrow X_1 \dots X_m \bullet] \in q_m$, $q = \delta(p, A)$

Finish: (q_0p, ϵ, f) , falls $[S' \rightarrow S \bullet] \in p$

- Die akzeptierenden Berechnungen des LR(0)-Parsers stehen in 1:1-Beziehung zu denen des Shift-Reduce-Parsers $M_G^{(1)}$.

- Akzeptierte Sprache ist $\mathcal{L}(G)$.

- Die Folge der Reduktionen einer akzept. Berechnung für ein Wort $w \in T$ liefert eine reverse Rechts-Abl. von G für w .

- Achtung:** Nur deterministisch, wenn jeder Endzustand des kanon. LR(0)-Automaten $LR(G)$ keine Konflikte enthält. Sonst i. A. nicht-deterministisch!

- 2 Gründe für Nicht-Determinismus v. LR(0)-P.:

- ▷ **Reduce-Reduce-Konflikt:**

$$[A \rightarrow \gamma\bullet], [A' \rightarrow \gamma'\bullet] \in q$$

mit $A \neq A' \vee \gamma \neq \gamma'$.

- ▷ **Shift-Reduce-Konflikt:**

$$[A \rightarrow \gamma\bullet], [A' \rightarrow \alpha \bullet \alpha\beta] \in q$$

mit $a \in T$, $q \in Q$.

Solche Zustände q heißen **ungeeignet**.

- Satz:** Die reduzierte Grammatik G ist LR(0) \Leftrightarrow der kanonische LR(0)-Automat $LR(G)$ keine ungeeigneten Zustände enthält.

- Def.:** Die reduzierte kontextfreie Grammatik \bar{G} heißt **LR(k)-Grammatik**, falls für Wörter $w, x \in T^*$ mit $\text{First}_k(w) = \text{First}_k(x)$ aus

$$S \xrightarrow{*}_R \alpha Aw \rightarrow \alpha \beta w \quad \text{und}$$

$$S \xrightarrow{*}_R \alpha' A' w' \rightarrow \alpha \beta x$$

folgt: $\alpha = \alpha' \wedge A = A' \wedge w' = x$

- Ein **LR(k)-Item** ist ein Paar $[B \rightarrow \alpha \bullet \beta, x]$ mit $x \in \text{Follow}_k(B)$.

Dieses Item ist gültig für $\gamma\alpha$, falls $S \xrightarrow{*}_R \gamma Bw$ mit $\{x\} = \text{First}_k(w)$.

- Berechnung** der gült. LR(k)-Items für zuverl. Präf.: **charakterist. LR(k)-Automat** $c(G, k)$:

Zustände: LR(k)-Items

$$\text{Anfangszustand: } [S' \rightarrow \bullet S, \epsilon]$$

$$\text{Endzustände: } \{[B \rightarrow \gamma\bullet, x] \mid B \rightarrow \gamma \in P, x \in \text{Follow}_k(B)\}$$

Übergänge:

- $\{[A \rightarrow \alpha \bullet \alpha X\beta, x], X, [A \rightarrow \alpha X \bullet \beta, x]\}$
für $X \in (N \cup T)$, $A \rightarrow \alpha X\beta \in P$
- $\{[A \rightarrow \alpha \bullet B\beta, x], \epsilon, [B \rightarrow \bullet \gamma, x']\}$
für $A \rightarrow \alpha B\beta, B \rightarrow \gamma \in P$,
 $x' \in \text{First}_k(\beta) \circ \{x\}$

- Den **kanon. LR(k)-Automaten** $LR(G, k)$ erhält man aus $c(G, k)$, indem man ϵ -Übergänge entfernt & Teilmengenkonstruktion anwendet.

Alternative Konstruktion von LR(G, k):

Hilfsfunktion $\delta_\epsilon^*(q) := q \cup$

$$\{[B \rightarrow \bullet \gamma, x] \mid \exists [A \rightarrow \alpha \bullet B'\beta', x'] \in q,$$

$$\beta \in (N \cup T)^* : B' \xrightarrow{*} B\beta \wedge$$

$$x \in \text{First}_k(\beta\beta') \circ \{x'\}\}$$

Zustände: Mengen von LR(k)-Items

$$\text{Anfangszustand: } \delta_\epsilon^*\{[S' \rightarrow \bullet S, \epsilon]\}$$

$$\text{Endzustände: } \{q \mid \exists A \rightarrow \alpha \in P : [A \rightarrow \alpha\bullet, x] \in q\}$$

$$\text{Übergänge: } \delta(q, X) := \delta_\epsilon^*\{[A \rightarrow \alpha X \bullet \beta, x] \mid [A \rightarrow \alpha \bullet X\beta, x] \in q\}$$

- 2 mögliche Konflikte:

- ▷ **Reduce-Reduce-Konflikt:**

$$[A \rightarrow \gamma\bullet, x], [A' \rightarrow \gamma'\bullet, x] \in q$$

mit $A \neq A' \vee \gamma \neq \gamma'$

- ▷ **Shift-Reduce-Konflikt:**

$$[A \rightarrow \gamma\bullet, x], [A' \rightarrow \alpha \bullet \alpha\beta, y] \in q$$

mit $a \in T$ und $x \in \{a\} \circ \text{First}_k(\beta) \circ \{y\}$
für einen Zustand $q \in Q$.

Solche Zustände heißen **LR(k)-ungeeignet**.

- Satz:** Eine reduzierte kontextfr. Gr. G ist LR(k) \Leftrightarrow der kanon. LR(k)-Automat $LR(G, k)$ keine LR(k)-ungeeigneten Zustände besitzt.

- i. A. hat $LR(G, k)$ sehr viel mehr Zustände als $LR(G) = LR(G, 0)$.

Man betrachtet darum i. A. Teilklassen von LR(k)-Grammatiken, bei denen man nur $LR(G)$ benutzt.

- LR(k)-Parser:**

- ▷ **goto-Tabelle** kodiert die Zustandsübergänge: $\text{goto}[q, X] := \delta(q, X) \in Q$

- ▷ **action-Tabelle** beschreibt für jeden Zustand q und mögliche Vorausschau w die erforderliche Aktion. Mögliche Aktionen: $\text{shift}, \text{reduce } (A \rightarrow \gamma), \text{error}$

- Konfliktlösung:**

- ▷ **simple-LR(k):**

- Reduce-Reduce-Konflikt:

⋮