

Compilerbau

0.1 Compiler = Programm zur Übersetzung Programmiersprache \rightarrow Maschinsprache.

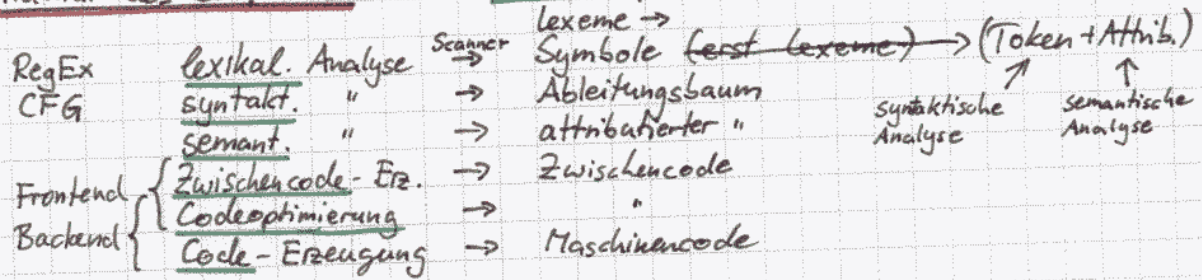
\rightarrow imperative Programme (Kontrollstrukturen, Variablen, Zuweisungen, Datenstrukturen)
(nicht: deklarative " , nebenläufige P., objektorientierte P.)

MS: von-Neumann-Rechner: Reduced (RISC), Complex (CISC)

0.2 Aspekte: Syntax: formaler hierarchischer Aufbau
Semantik: Bedeutung, Zustände
Pragmatik: Formulierung, natürliche Sprache, Maschinen-Abh.

äquivalent = semantisch äquivalent.

0.3 Struktur des Compilers n-Pass-Compiler: $n = \text{Läufe durch Programm}$



1. Lexikalische Analyse

Quellcode $\xrightarrow{\text{Scanner (RegEx) auf Einmal.}}$ Lexeme \rightarrow Symbole (Token + Attr.)

Symbolklassen: ~~Be~~

- Bezeichner	Var 1	- Zeiger in Symb. Tab.
- Zahlwörter	55	- Binärdarst. Zahl
- Schlüsselwörter	begin	- leer
- Einfache Symbole	+	
- Zusammenges. S.	::=	
- Leerzeichen	␣	
- Spezielle Symb.	! * # /	

Attribute:

Σ_0	Zeichensatz
$a \in \Sigma_0$	Zeichen
$P \in \Sigma_0^*$	Quellprogramm
$RegE(\Sigma)$	Regex. E
$[[]]$	Abb. Regex \rightarrow ^{matchende} Zeichenfolgen
$[[]]$	Automat
Q	Zustandsmenge
Σ	Alphabet
δ	Transitionsfunktion
$q_0 \in Q$	Anfangszustand
$F \subseteq Q$	Endzustandsmenge
Σ^E	$\Sigma \cup \{ \epsilon \}$
$T \in Q$	$\epsilon(T)$ ϵ -Hülle von T
$L([])$	durch [] erkannte Sprache

1.1. Scanner-Konstruktion: RegEx.

- $[[\cdot]]$ = alle Zeichenfolgen, die die RegEx "." matcht.
- alle Einzelzeichen sind Regex, sowie $(\alpha \vee \beta)$, $(\alpha \cdot \beta)$, (α^*)
- $[[\alpha \vee \beta]]$:= $[[\alpha]] \cup [[\beta]]$ // $[[\alpha \cdot \beta]]$:= $[[\alpha]] \cdot [[\beta]]$ // ...

$N = \langle Q, \Sigma, \delta, q_0, F \rangle \in NFA(\Sigma)$

1.1.1. Das einfache Matching-Problem

Wird Zeichenfolge durch RegEx gematcht?
Entscheide mit NFA (nondeterministic finite automata).

- ϵ -Hülle: ohne Eingabe erreichbare Zustände.
- erweiterte Transitionsfunktion: Zustände, die durch Eingabewort erreichbar sind.

- DFA-Methode: RegEx Thompson \rightarrow NFA Potenzmengen-Konstruktion \rightarrow DFA.
 \rightarrow Entscheidung in $|w|+1$ Schritten.
- NFA-Methode: Durchlauf durch NFA mit Wort, auf allen möglichen Wegen.
 \rightarrow Platz: $O(|\alpha| + |w|)$, Zeit: $O(|\alpha| \cdot |w|)$ "DFA on the fly erzeugen"

1.1.2. Das erweiterte Matching-Problem : ~~matcht es auf irgendeine RegEx,~~

Problem: Zerlegung des Eingabeprogramms (wort) finden, eindeutig!

Lösung: Flm-Analyse: Längster Match, im Zweifel Automat höchste Priorität.

Berechnung:

1. Konstruiere einen DFA je RegEX.
2. Mache daraus Produktautomaten, und zwar so, dass jeder Endzustand die Klasse des niedrigsten matchenden Automaten hat.
(z.B.: Automat 3 und 5 würden matchen \rightarrow Endzustand hat Klasse 3.)

Produktiver Zustand = Zustand, aus dem Endzustand erreichbar ist.

Leseband mit Backtrack-Kopf und LookAhead-Kopf.

- Normal Mode: Match suchen. Bei Endzustand \rightarrow Backtrack-Mode.
(entlang produktiver z.) Bei nicht produktivem Zustand: lexerr.
- Backtrack Mode: längsten Match suchen. Bei Endzustand \rightarrow Backtrack-Fahne neu setzen.
(entlang produktiver z.) Bei nicht produktivem z. \rightarrow zurück zur Fahne.
- Eingabeende: Endzustand \rightarrow Token ausgesen.
Prod. Zustand \rightarrow lexerr
Prod. Zustand im Backtrack-Mode \rightarrow zurück zur Fahne.

(lwl²)

1.2. Praktische Aspekte der Scanner-Konstruktion

Lex-Spezifikation $\xrightarrow{\text{lex}}$ Scanner in C $\xrightarrow{\text{cc}}$ Ausführbarer Scanner

- Pascal: 2-Lookahead, Whitespace \rightarrow 1-Lookahead
- Attributberechnung: 2. Lauf: Zahlenumwandlungen, Bezeichner-Gleichheit prüfen, Schlüsselwörter erstmal ab, Bezeichner betrachten.



2. Syntaktische Analyse

... durch Parser. Erstelle Baumstruktur.

Syntaktische Einheiten: Variablen, Ausdrücke, Anweisungen, ...

Beschreibung durch CFG (BNF, EBNF, Syntaxdiagramme)

Erkennung durch Kellerautomat mit Ausgabe. \Rightarrow deterministisch!

(deterministischer Kellerautomat mit input-look-ahead), linearer Bedarf Platz + Zeit.)

• Top-Down-Analyse: Linksanalyse

• Bottom-Up-Analyse: gespiegelte Rechtsanalyse

2.1. Kontextfreie Grammatiken

$A, a, u, \alpha, A \rightarrow \alpha \quad S \xrightarrow{z} w$

Grammatik G genau dann eindeutig, wenn es genau einen Ableitungsbaum gibt,

l-Analyse von α : $S \xrightarrow{z} \alpha$. z = Folge von Regeln i . (r-Analyse analog.)

Syntaxanalyse von α : l/r-Analyse, oder SYNTAX ERROR.

Vorr: • G reduziert, d.h. • NT-Symbole auf rechter Seite erreichbar: $S \xRightarrow{*} \alpha A \beta$

• NT-Symbole " linker Seite produktiv: $A \xRightarrow{*} w$.

2.2. TopDown-Analyse mit LL(k)-Grammatiken

Automat: • Ableitungsschritte (nicht deterministisch)

• Vergleichsschritte (deterministisch)

$(w, S, \epsilon) \vdash (\epsilon, \epsilon, z)$ (Eingabe, Keller, Ausgabe)

Ziel: Nichtdeterminismus durch k-Lookahead beseitigen

$\text{first}_k(\alpha)$: {die ersten k Symbole der möglichen Ableitungen}

Folgerungen: 1. $\text{first}_k(\alpha) \neq \emptyset$, da G reduziert.

2. ϵ in $\text{first}_k \Rightarrow k=0$ oder $\alpha \xRightarrow{*} \epsilon$.

3. $\alpha \xRightarrow{*} \beta \sim \text{first}_k(\beta) \subseteq \text{first}_k(\alpha)$

4. Wenn man von α nach x kommt, ist die first-Menge von x auch in der first-Menge von α enthalten. z.B.:

CFG Kontextfreie Grammatik

G Grammatik

Σ Eingabealphabet

\mathcal{K} Kelleralphabet

$1, \dots, p$ Ausgabealphabet

$\Sigma^* \times \mathcal{K} \times [p]^*$ Konfigurationsmenge.

$X \rightarrow i \mid \text{first}_3(\text{simon}) = \{\text{sim}\}$

$X \rightarrow a \mid \text{first}_3(\text{sXmon}) = \{\text{sim}, \text{sam}\}$

LL(k)-Grammatik: Lesen der Eingabe von links nach rechts.

• Linksableitungsschritt für $wA\alpha$ ist durch die nächsten auf w folgenden Symbole bestimmt.

• TD-Analyseautomat kann nun deterministisch mit k-Lookahead arbeiten.

Ziel: Bestimmung der A-Regel aus k-Look-ahead.

aus den ersten k Zeichen des abgeleiteten Teilwortes $S \xrightarrow{z} wA\alpha$ ergibt sich, aus welchem Satz es abgeleitet ist.

G ist LL(k)-Grammatik, wenn \forall Linksabl. gilt: $\text{first}_k(x) = \text{first}_k(y) \sim \beta = \gamma$

• Lemma: Die first_k -Mengen von 2 möglichen Regeln haben keine Schnittmenge.

• Brauche Rechtskontext (follow-Menge), falls first-Menge ϵ ist.

• $\text{fi}(\beta) = \text{first}_k(\beta)$

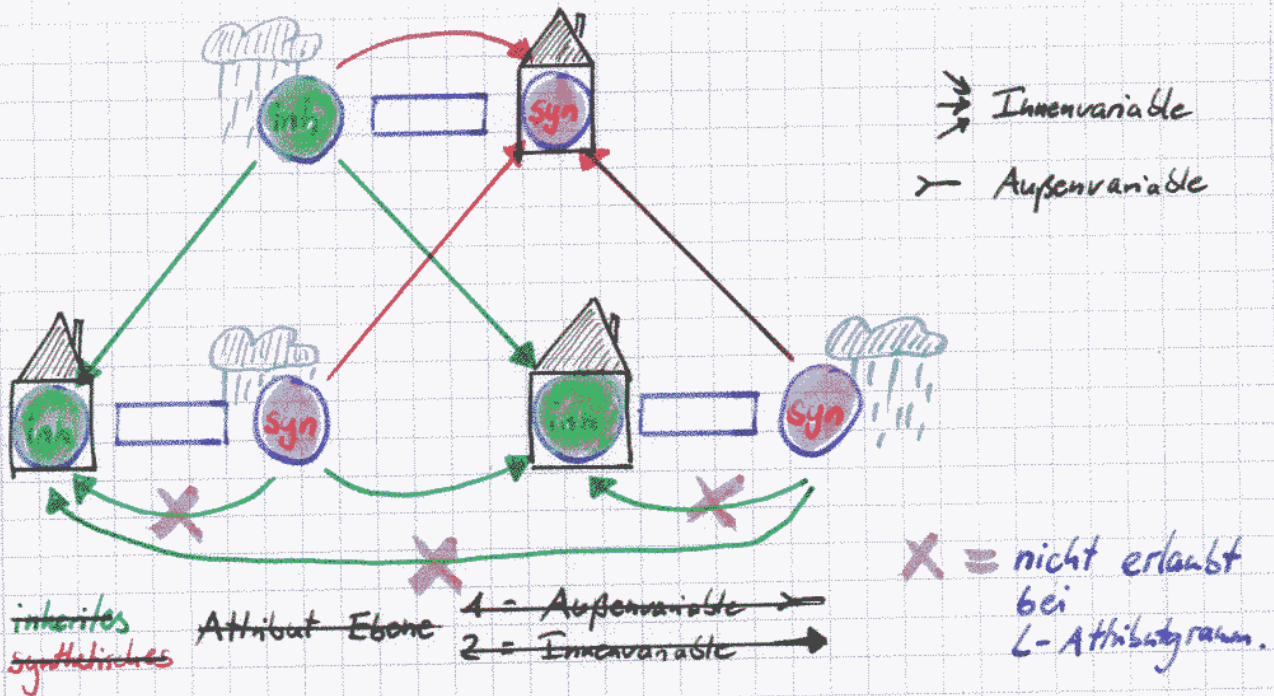
• $\text{fo}(A) = \text{follow}_k(A)$

• $\text{la}(A \rightarrow \beta) = \text{fi}(\beta \text{fo}(A))$

Berechnung fi-Mengen: auf was können die NT-Symbole der rechten Seiten abgeleitet werden?

Berechnung fo-Mengen: was kann hinter den NT-Symbolen der rechten Seite kommen?

Berechnung la-Mengen: fi-Menge der rechten Seite. Wenn ϵ dabei, dann fo-Menge der linken Seite dazunehmen.

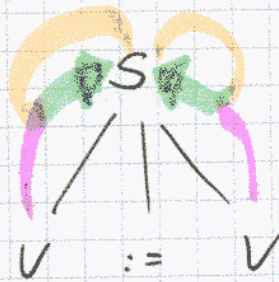
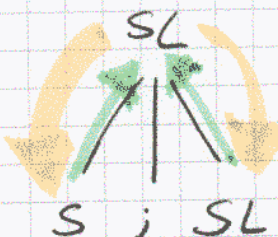
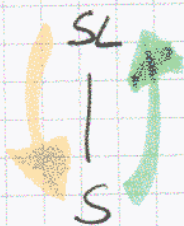
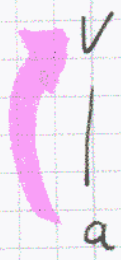
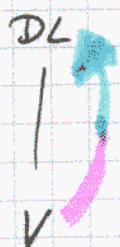
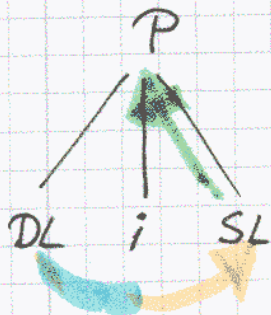


- inherite/synthetische Attribute sind disjunkt.
- Innen-(Außen-) Variablen sind disjunkt.

	inherit	synthetisch
Ebene 1	A	I
Ebene 2	I	A

Zirkularitäten beim Verkleben.

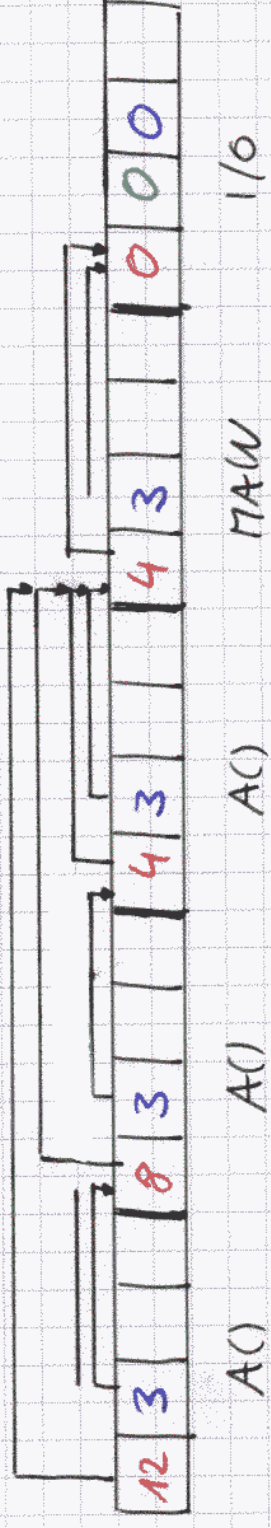
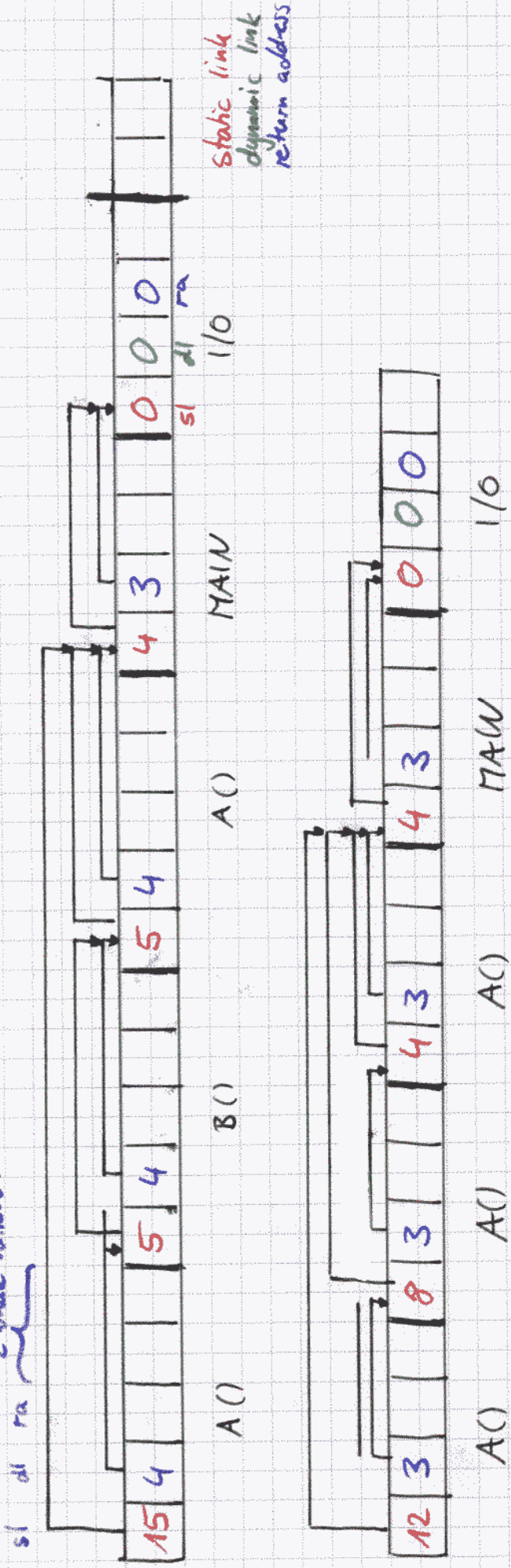
$D(A, \mathcal{E}) =$ Abhängigkeiten in der Wurzel A (für Baum \mathcal{E})



- DECL
- ENV
- DV

● V

2 lokale Variablen



Symbol - Tabelle

Tisch → const, 5 ^(Wert)

Stuhl → const, 7

Papier → var, 2 ^(Level) 5 ^(Offset)

Stift → var 1 3

	St.	Pa.	
--	-----	-----	--

Haus → proc ^(Adresse) 1002 ^(Level) 2 ^(Anzahl Variablen) 5 ^{(Offset) i.d. Proz.}

Boot → proc 1005 1 3