

Tutorium 12: Syntaktische Analyse

Paul Brinkmeier

09. Februar 2021

Tutorium Programmierparadigmen am KIT

Heutiges Programm

- Typinferenz
 - Herleitungsbäume
 - LET-Polymorphismus
- Übersicht Compilerbau
- Syntaktische Analyse

Typinferenz

Unifikationsalgorithmus: $\text{unify}(C) =$

```
if  $C == \emptyset$  then []  
else let  $\{\theta_l = \theta_r\} \cup C' = C$  in  
  if  $\theta_l == \theta_r$  then  $\text{unify}(C')$   
  else if  $\theta_l == Y$  and  $Y \notin FV(\theta_r)$  then  $\text{unify}([Y \dot{=} \theta_r] C') \circ [Y \dot{=} \theta_r]$   
  else if  $\theta_r == Y$  and  $Y \notin FV(\theta_l)$  then  $\text{unify}([Y \dot{=} \theta_l] C') \circ [Y \dot{=} \theta_l]$   
  else if  $\theta_l == f(\theta_l^1, \dots, \theta_l^n)$  and  $\theta_r == f(\theta_r^1, \dots, \theta_r^n)$   
    then  $\text{unify}(C' \cup \{\theta_l^1 = \theta_r^1, \dots, \theta_l^n = \theta_r^n\})$   
  else fail
```

$Y \in FV(\theta)$ **occur check**, verhindert zyklische Substitutionen

Korrektheitstheorem

$\text{unify}(C)$ terminiert und gibt *mgu* für C zurück, falls C unifizierbar, ansonsten **fail**.

Beweis: Siehe [Pie02]

Wiederholung: Unifikationsalgorithmus

$$C = \{X = a, Y = f(X), f(Z) = Y\}$$

- Problemstellung: Gleichungssystem C mit baumförmigen Termen
 - Bspw. Prolog-Terme, Typen
- Liefert: Unifikator σ_1 , der alle Gleichungen erfüllt
- Bspw. $\sigma_1 = [X \mapsto a, Y \mapsto f(X), Z \mapsto X]$

Wiederholung: Unifikationsalgorithmus

$$C = \{X = a, Y = f(X), f(Z) = Y\}$$

$$\begin{aligned}\text{unify}(C) &= \text{unify}(\{\underline{X = a}, Y = f(X), f(Z) = Y\}) \\ &= \text{unify}(\{Y = f(a), \underline{f(Z) = Y}\}) \circ [X \mapsto a] \\ &= \text{unify}(\{\underline{f(Z) = f(a)}\}) \circ [Y \mapsto f(Z)] \circ [X \mapsto a] \\ &= \text{unify}(\{\underline{Z = a}\}) \circ [Y \mapsto f(Z)] \circ [X \mapsto a] \\ &= [Z \mapsto a] \circ [Y \mapsto f(Z)] \circ [X \mapsto a] \\ &= [Z \mapsto a, Y \mapsto f(a), X \mapsto a]\end{aligned}$$

Klausuraufgabe WS16/18 A3 a) (6P.)

$$C_1 = \{\alpha_9 = \alpha_{10} \rightarrow \alpha_8, \alpha_9 = \alpha_4, \alpha_{10} = \text{bool}\}$$

$$C_2 = \{\alpha_{12} = \alpha_{13} \rightarrow \alpha_{11}, \alpha_{12} = \alpha_4, \alpha_{13} = \text{int}\}$$

- i. Geben Sie allgemeinste Unifikatoren σ_1 für C_1 und σ_2 für C_2 an.
- ii. Ist auch $C_1 \cup C_2$ unifizierbar?
- iii. Ist der Ausdruck

$\lambda a. \lambda f. f (a \text{ true}) (a \text{ 17})$

typisierbar? Begründen Sie ihre Antwort kurz.

$$C_1 = \{\alpha_9 = \alpha_{10} \rightarrow \alpha_8, \alpha_9 = \alpha_4, \alpha_{10} = \text{bool}\}$$

$$C_2 = \{\alpha_{12} = \alpha_{13} \rightarrow \alpha_{11}, \alpha_{12} = \alpha_4, \alpha_{13} = \text{int}\}$$

Geben Sie allgemeinste Unifikatoren σ_1 für C_1 und σ_2 für C_2 an.

$$\begin{aligned}\sigma_1 &= \text{unify}(\{\alpha_9 = \alpha_{10} \rightarrow \alpha_8, \alpha_9 = \alpha_4, \alpha_{10} = \text{bool}\}) \\ &= \dots = [\alpha_9 \dot{=} \text{bool} \rightarrow \alpha_8, \alpha_4 \dot{=} \text{bool} \rightarrow \alpha_8, \alpha_{10} \dot{=} \text{bool}]\end{aligned}$$

$$\begin{aligned}\sigma_2 &= \text{unify}(\{\alpha_{12} = \alpha_{13} \rightarrow \alpha_{11}, \alpha_{12} = \alpha_4, \alpha_{13} = \text{int}\}) \\ &= \dots = [\alpha_{12} \dot{=} \text{int} \rightarrow \alpha_{11}, \alpha_4 \dot{=} \text{int} \rightarrow \alpha_{11}, \alpha_{13} \dot{=} \text{int}]\end{aligned}$$

Klausuraufgabe WS16/18 A3 a) (6P.)

$$C_1 = \{\alpha_9 = \alpha_{10} \rightarrow \alpha_8, \alpha_9 = \alpha_4, \alpha_{10} = \text{bool}\}$$

$$C_2 = \{\alpha_{12} = \alpha_{13} \rightarrow \alpha_{11}, \alpha_{12} = \alpha_4, \alpha_{13} = \text{int}\}$$

Ist auch $C_1 \cup C_2$ unifizierbar?

$$\sigma_1 = \dots = [\alpha_9 \dot{=} \text{bool} \rightarrow \alpha_8, \underline{\alpha_4 \dot{=} \text{bool} \rightarrow \alpha_8}, \alpha_{10} \dot{=} \text{bool}]$$

$$\sigma_2 = \dots = [\alpha_{12} \dot{=} \text{int} \rightarrow \alpha_{11}, \underline{\alpha_4 \dot{=} \text{int} \rightarrow \alpha_{11}}, \alpha_{13} \dot{=} \text{int}]$$

A: Nein, da die allgemeinsten Unifikatoren σ_1 und σ_2 einen Konflikt für α_4 enthalten: $\text{unify}(\{\text{bool} = \text{int}\}) = \text{fail}$

$$C_1 = \{\alpha_9 = \alpha_{10} \rightarrow \alpha_8, \alpha_9 = \alpha_4, \alpha_{10} = \text{bool}\}$$

$$C_2 = \{\alpha_{12} = \alpha_{13} \rightarrow \alpha_{11}, \alpha_{12} = \alpha_4, \alpha_{13} = \text{int}\}$$

Ist der Ausdruck

$$\lambda a. \lambda f. f (a \text{ true}) (a \text{ 17})$$

typisierbar? Begründen Sie ihre Antwort kurz.

A: Nein, da a mit zwei verschiedenen Typen verwendet wird.

Problemstellung bei Typinferenz: Zu einem gegebenen Term den passenden Typ finden.

- Struktur des Terms erkennen. Wo sind:
 - Lambdas?
 - Funktionsanwendungen?
 - Variablen/Konstanten?
- Entsprechenden Baum aufstellen.
- Typgleichungen finden.
- Gleichungssystem unifizieren.

Lambda-Terme bestehen aus

- Lambdas $\lambda p. b$
- Funktionsanwendungen $x y$
- Variablen x , Konstanten `true`, `17`

$\lambda a. \lambda f. f (a \text{ true})$

Struktur von Lambda-Termen

Lambda-Terme bestehen aus

- **Lambdas** $\lambda p. b$
- Funktionsanwendungen $x y$
- Variablen x , Konstanten `true`, `17`

$$\underbrace{\lambda a. \underbrace{\lambda f. f (a \text{ true})}_{\text{Lambda } p=f, b=f (a \text{ true})}}_{\text{Lambda } p=a, b=\lambda f. f (a \text{ true})}$$

Struktur von Lambda-Termen

Lambda-Terme bestehen aus

- Lambdas $\lambda p. b$
- **Funktionsanwendungen** $x y$
- Variablen x , Konstanten `true`, `17`

$$\lambda a. \lambda f. f \underbrace{(a \text{ true})}_{\substack{\text{Anwendung} \\ x=a, y=\text{true}}} \underbrace{\phantom{(a \text{ true})}}_{\substack{\text{Anwendung} \\ x=f, y=a \text{ true}}}$$

Struktur von Lambda-Termen

Lambda-Terme bestehen aus

- Lambdas $\lambda p. b$
- Funktionsanwendungen $x y$
- **Variablen** x , Konstanten `true`, `17`

$$\lambda a. \lambda f. \underbrace{f}_{\text{Variable}} \left(\underbrace{a}_{\text{Variable}} \text{ true} \right)$$

Struktur von Lambda-Termen

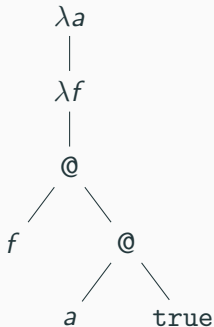
Lambda-Terme bestehen aus

- Lambdas $\lambda p. b$
- Funktionsanwendungen $x y$
- Variablen x , **Konstanten** `true`, `17`

$\lambda a. \lambda f. f (a \underbrace{\text{true}}_{\text{Konstante}})$

Lambda-Terme als Bäume

Wir können Lambda-Terme also als Bäume mit Lambda- und Anwendungsknoten und Variablen- und Konstantenblättern betrachten, um ihre Struktur zu untersuchen:



$\lambda a. \lambda f. f (a \text{ true})$

Cheatsheet: Typisierter Lambda-Kalkül

$$\frac{\Gamma, p : \pi \vdash b : \rho}{\Gamma \vdash \lambda p. b : \pi \rightarrow \rho} \text{ABS}$$

$$\frac{\Gamma \vdash f : \phi \rightarrow \alpha \quad \Gamma \vdash x : \phi}{\Gamma \vdash f x : \alpha} \text{APP}$$

$$\frac{\Gamma(t) = \tau}{\Gamma \vdash t : \tau} \text{VAR}$$

$$\frac{c \in \text{CONST}}{\Gamma \vdash c : \tau_c} \text{CONST}$$

- Typvariablen: τ, α, π, ρ
- Funktionstypen: $\tau_1 \rightarrow \tau_2$, rechtsassoziativ
- Typisierungsregeln sind eindeutig: Eine Regel pro Termform

Was bedeuten eigentlich \vdash , Γ und $:$?

$$\lambda a. \lambda f. f (a \text{ true})$$

Um zu einem solchen Term ein Typisierungsproblem zu beschreiben, notieren wir:

$$\Gamma \vdash \lambda a. \lambda f. f (a \text{ true}) : \tau$$

„Im Typkontext Γ hat der Term den Typen τ .“

- Γ : Enthält Typen für freie Variablen.
- $\dots \vdash \dots : \dots$ — Notation für Typisierungsproblem.

$$\Gamma \vdash a + 42 : \text{int}$$

$$\text{CONST} = \{42\}, \tau_{42} = \text{int}$$

Damit die Aussage „ $a + 42$ hat in Γ den Typen `int`“ stimmt, müssen wir für Γ wählen:

$$\Gamma \vdash a + 42 : \text{int}$$

$$\text{CONST} = \{42\}, \tau_{42} = \text{int}$$

Damit die Aussage „ $a + 42$ hat in Γ den Typen `int`“ stimmt, müssen wir für Γ wählen:

- $\Gamma = a : \text{int}, + : \text{int} \rightarrow \text{int} \rightarrow \text{int}$

$$\Gamma \vdash a + 42 : \text{int}$$

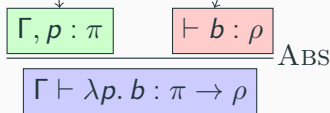
$$\text{CONST} = \{42\}, \tau_{42} = \text{int}$$

Damit die Aussage „ $a + 42$ hat in Γ den Typen int “ stimmt, müssen wir für Γ wählen:

- $\Gamma = a : \text{int}, + : \text{int} \rightarrow \text{int} \rightarrow \text{int}$
- Allgemeiner: $\Gamma = a : \alpha, + : \alpha \rightarrow \text{int} \rightarrow \text{int}$

Typisierungsregel für Lambdas

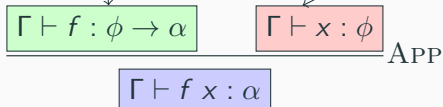
- „Unter Einfügung des Typs π von p in den Kontext...“
- „... ist b als Funktion von p typisierbar.“



- Daraus folgt:
- „ $\lambda p. b$ ist eine Funktion, die π s auf ρ s abbildet“

Typisierungsregel für Funktionsanwendungen

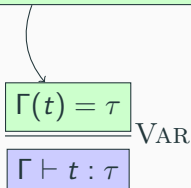
- „ f ist im Kontext Γ eine Funktion, die ϕ s auf α s abbildet.“
- „ x ist im Kontext Γ ein Term des Typs ϕ .“



- Daraus folgt:
- „ x eingesetzt in f ergibt einen Term des Typs α .“

Einfache Typisierungsregel für Variablen

- „Der Typkontext Γ enthält einen Typ τ für t .“



- Daraus folgt:
- „Variable t hat im Kontext Γ den Typ τ .“

$$x : \text{bool} \vdash \lambda f. f \ x : (\text{bool} \rightarrow \alpha) \rightarrow \alpha$$

„Unter der Annahme, dass x den Typ bool hat, hat $\lambda f. f \ x$ den Typ $(\text{bool} \rightarrow \alpha) \rightarrow \alpha$.“

Typisierung: Beispiel

$$\frac{\underline{x : \text{bool}}, \underline{f : \text{bool} \rightarrow \alpha} \vdash \underline{f\ x : \alpha}}{\underline{x : \text{bool}} \vdash \lambda \underline{f}. \underline{f\ x : (\text{bool} \rightarrow \alpha)} \rightarrow \underline{\alpha}}_{\text{ABS}}$$

Pattern-Matching: Der äußerste Term ist ein Lambda, also wenden wir die ABS-Regel an.

$$\underline{\Gamma} = \underline{x : \text{bool}}$$

$$\underline{p} = \underline{f}, \underline{b} = \underline{f\ x}$$

$$\underline{\pi} = \underline{\text{bool} \rightarrow \alpha}$$

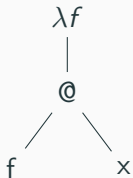
$$\underline{\rho} = \underline{\alpha}$$

$$\frac{\underline{\Gamma}, \underline{p : \pi} \vdash \underline{b : \rho}}{\underline{\Gamma} \vdash \lambda \underline{p}. \underline{b : \pi \rightarrow \rho}}_{\text{ABS}}$$

Typisierung: Beispiel

$$\frac{\frac{\Gamma(f) = \text{bool} \rightarrow \alpha}{\Gamma \vdash f : \text{bool} \rightarrow \alpha} \text{VAR} \quad \frac{\Gamma(x) = \text{bool}}{\Gamma \vdash x : \text{bool}} \text{VAR}}{\Gamma \vdash f x : \alpha} \text{APP}$$
$$\frac{\Gamma \vdash f x : \alpha}{\Gamma \vdash \lambda f. f x : (\text{bool} \rightarrow \alpha) \rightarrow \alpha} \text{ABS}$$

$$\Gamma = x : \text{bool}, f : \text{bool} \rightarrow \alpha$$



$\lambda f. f x$

Von Typisierungsregeln zu Typinferenz

Beim inferieren wird das Pattern-matching der Typen durch die Unifikation übernommen. Deswegen schreiben wir anstelle von konkreten Typen immer α_i und merken uns die Gleichungen für später:

$$\frac{\Gamma, p : \pi \vdash b : \rho}{\Gamma \vdash \lambda p. b : \pi \rightarrow \rho} \text{ABS} \quad \rightsquigarrow \quad \frac{\Gamma, p : \alpha_j \vdash b : \alpha_k}{\Gamma \vdash \lambda p. b : \alpha_i} \text{ABS} \quad \{\alpha_i = \alpha_j \rightarrow \alpha_k\}$$

Von Typisierungsregeln zu Typinferenz

Beim inferieren wird das Pattern-matching der Typen durch die Unifikation übernommen. Deswegen schreiben wir anstelle von konkreten Typen immer α_i und merken uns die Gleichungen für später:

$$\frac{\Gamma \vdash f : \phi \rightarrow \alpha \quad \Gamma \vdash x : \phi}{\Gamma \vdash f x : \alpha} \text{APP} \rightsquigarrow \frac{\Gamma \vdash f : \alpha_j \quad \Gamma \vdash x : \alpha_k}{\Gamma \vdash f x : \alpha_i \quad \{\alpha_j = \alpha_k \rightarrow \alpha_i\}} \text{APP}$$

Von Typisierungsregeln zu Typinferenz

Beim inferieren wird das Pattern-matching der Typen durch die Unifikation übernommen. Deswegen schreiben wir anstelle von konkreten Typen immer α_i und merken uns die Gleichungen für später:

$$\frac{\Gamma(t) = \tau}{\Gamma \vdash t : \tau} \text{VAR} \quad \rightsquigarrow \quad \frac{\Gamma(t) = \alpha_j}{\Gamma \vdash t : \alpha_i} \text{VAR} \quad \{\alpha_i = \alpha_j\}$$

Algorithmus zur Typinferenz

- Stelle Typherleitungsbaum auf
 - In jedem Schritt werden neue Typvariablen α_i angelegt
 - Statt die Typen direkt im Baum einzutragen, werden Gleichungen in einem Constraint-System eingetragen
- Unifiziere Constraint-System zu einem Unifikator
 - Robinson-Algorithmus, im Grunde wie bei Prolog
 - I.d.R.: Allgemeinster Unifikator (mgu)

$$\frac{\Gamma(t) = \alpha_j}{\Gamma \vdash t : \alpha_i} \text{VAR}$$

Constraint:
 $\{\alpha_i = \alpha_j\}$

$$\frac{\Gamma \vdash f : \alpha_j \quad \Gamma \vdash x : \alpha_k}{\Gamma \vdash f x : \alpha_i} \text{APP}$$

Constraint:
 $\{\alpha_j = \alpha_k \rightarrow \alpha_i\}$

$$\frac{\Gamma, p : \alpha_j \vdash b : \alpha_k}{\Gamma \vdash \lambda p. b : \alpha_i} \text{ABS}$$

Constraint:
 $\{\alpha_i = \alpha_j \rightarrow \alpha_k\}$

$$\vdash \lambda x. \lambda y. x : \alpha_1$$

Beispielhafte Aufgabenstellung: Finde den Typen α_1 .

$$\frac{\underline{x} : \alpha_2 \vdash \underline{\lambda y. x} : \alpha_3}{\vdash \lambda \underline{x}. \underline{\lambda y. x} : \alpha_1} \text{ABS}$$

Typgleichungen:

$$C = \{\underline{\alpha_1 = \alpha_2 \rightarrow \alpha_3}\}$$

$$\frac{\frac{\frac{x : \alpha_2, y : \alpha_4 \vdash x : \alpha_5}{\vdash \lambda y. x : \alpha_3} \text{ABS}}{\vdash \lambda x. \lambda y. x : \alpha_1} \text{ABS}}$$

Typgleichungen:

$$C = \{ \alpha_1 = \alpha_2 \rightarrow \alpha_3, \alpha_3 = \alpha_4 \rightarrow \alpha_5 \}$$

Herleitungsbaum: Beispiel

$$\frac{\frac{\frac{(x : \alpha_2, y : \alpha_4)(x) = \alpha_2}{\text{VAR}}}{x : \alpha_2, y : \alpha_4 \vdash \underline{x} : \alpha_5}{\text{ABS}}}{x : \alpha_2 \vdash \lambda y. x : \alpha_3}{\text{ABS}} \vdash \lambda x. \lambda y. x : \alpha_1$$

Typgleichungen:

$$C = \{\alpha_1 = \alpha_2 \rightarrow \alpha_3 \\ , \alpha_3 = \alpha_4 \rightarrow \alpha_5 \\ , \underline{\alpha_5 = \alpha_2}\}$$

$$\frac{\dots}{\vdash \lambda f. f (\lambda x. x) : \alpha_1} \text{ABS}$$

Findet den Typen α_1 . Teilpunkte gibt es für:

- Herleitungsbaum,
- Typgleichungsmenge C ,
- Unifikation per Robinsonalgorithmus.

Herleitungsbaum: Aufgabe

$$\frac{\frac{(f : \alpha_2)(f) = \alpha_2}{f : \alpha_2 \vdash f : \alpha_4} \text{VAR} \quad \frac{\frac{\Gamma(x) = \alpha_6}{\Gamma \vdash x : \alpha_7} \text{VAR} \quad \frac{}{f : \alpha_2 \vdash \lambda x. x : \alpha_5} \text{ABS}}{f : \alpha_2 \vdash f (\lambda x. x) : \alpha_3} \text{APP} \quad \frac{}{\vdash \lambda f. f (\lambda x. x) : \alpha_1} \text{ABS}$$

$$\Gamma = f : \alpha_2, x : \alpha_6$$

$$C = \{\alpha_1 = \alpha_2 \rightarrow \alpha_3, \alpha_4 = \alpha_5 \rightarrow \alpha_3, \\ \alpha_2 = \alpha_4, \\ \alpha_5 = \alpha_6 \rightarrow \alpha_7, \alpha_6 = \alpha_7\}$$

Let-Polymorphism

$$\lambda f. f f$$

- Diese Funktion verwendet f auf zwei Arten:
 - $\alpha \rightarrow \alpha$: Rechte Seite.
 - $(\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$: Linke Seite, nimmt f als Argument und gibt es zurück.

$$\lambda f. f f$$

- Diese Funktion verwendet f auf zwei Arten:
 - $\alpha \rightarrow \alpha$: Rechte Seite.
 - $(\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$: Linke Seite, nimmt f als Argument und gibt es zurück.
- Problem: $\alpha \rightarrow \alpha$ und $(\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$ sind nicht unifizierbar!
 - „occurs check“: α darf sich nicht selbst einsetzen.
- Idee: Bei jeder Verwendung eines polymorphen Typen erzeugen wir neue Typvariablen, um diese Beschränkung zu umgehen.

Typschemata und Instanziierung

- Idee: Bei jeder Verwendung eines polymorphen Typen erzeugen wir neue Typvariablen, um diese Beschränkung zu umgehen.
- Ein Typschema ist ein Typ, in dem manche Typvariablen allquantifiziert sind:

$$\phi = \forall \alpha_1. \dots \forall \alpha_n. \tau$$
$$\alpha_i \in FV(\tau)$$

- Typschemata kommen bei uns immer nur in Kontexten vor!
- Beispiele:
 - $\forall \alpha. \alpha \rightarrow \alpha$
 - $\forall \alpha. \alpha \rightarrow \beta \rightarrow \alpha$

- Ein Typschema spannt eine Menge von Typen auf, mit denen es instanziiert werden kann:

$$\forall \alpha. \alpha \rightarrow \alpha \succeq \text{int} \rightarrow \text{int}$$

$$\forall \alpha. \alpha \rightarrow \alpha \succeq \tau \rightarrow \tau$$

$$\forall \alpha. \alpha \rightarrow \alpha \not\succeq \tau \rightarrow \sigma$$

$$\forall \alpha. \alpha \rightarrow \alpha \not\succeq \tau \rightarrow \tau \rightarrow \tau$$

$$\forall \alpha. \alpha \rightarrow \alpha \succeq (\tau \rightarrow \tau) \rightarrow (\tau \rightarrow \tau)$$

Um Typschemata bei der Inferenz zu verwenden, müssen wir zunächst die Regel für Variablen anpassen:

$$\frac{\Gamma(x) = \phi \quad \phi \succeq_{\text{frische } \alpha_i} \tau}{\Gamma \vdash x : \alpha_j} \text{VAR}$$

Constraint: $\{\alpha_j = \tau\}$

- $\succeq_{\text{frische } \alpha_i}$ instanziiert ein Typschema mit α_i , die noch nicht im Baum vorkommen.
- Jetzt brauchen wir noch eine Möglichkeit, Typschemata zu erzeugen.

Mit einem LET-Term wird ein Typschema eingeführt:

$$\frac{\Gamma \vdash t_1 : \alpha_i \quad \Gamma' \vdash t_2 : \alpha_j}{\Gamma \vdash \text{let } x = t_1 \text{ in } t_2 : \alpha_k} \text{LET}$$

$$\sigma_{\text{let}} = \text{mgu}(C_{\text{let}})$$

$$\Gamma' = \sigma_{\text{let}}(\Gamma), x : \text{ta}(\sigma_{\text{let}}(\alpha_i), \sigma_{\text{let}}(\Gamma))$$

$$C'_{\text{let}} = \{\alpha_n = \sigma_{\text{let}}(\alpha_n) \mid \sigma_{\text{let}}(\alpha_n) \text{ ist definiert}\}$$

$$\text{Constraints: } C'_{\text{let}} \cup C_{\text{body}} \cup \{a_j = a_k\}$$

Beispiel: Let-Polymorphismus

$$\begin{array}{c}
 \frac{\dots}{\vdash \lambda x. x : \alpha_2} \text{ABS} \quad \frac{\frac{\frac{\Gamma'(f) = \forall \alpha_5. \alpha_5 \rightarrow \alpha_5}{\succeq \alpha_8 \rightarrow \alpha_8} \text{VAR}}{\Gamma' \vdash f : \alpha_6} \text{VAR} \quad \frac{\frac{\Gamma'(f) = \forall \alpha_5. \alpha_5 \rightarrow \alpha_5}{\succeq \alpha_9 \rightarrow \alpha_9} \text{VAR}}{\Gamma' \vdash f : \alpha_7} \text{VAR}}{\Gamma' \vdash f f : \alpha_3} \text{APP} \\
 \hline
 \vdash \text{let } f = \lambda x. x \text{ in } f f : \alpha_1 \text{LET}
 \end{array}$$

$$C_{\text{let}} = \{\alpha_2 = \alpha_4 \rightarrow \alpha_5, \alpha_4 \rightarrow \alpha_5\}$$

$$\sigma_{\text{let}} = [\alpha_2 \dot{\rightarrow} \alpha_5 \rightarrow \alpha_5, \alpha_4 \dot{\rightarrow} \alpha_5]$$

$$\Gamma' = x : \forall \alpha_5. \alpha_5 \rightarrow \alpha_5$$

$$C'_{\text{let}} = \{\alpha_2 = \alpha_5 \rightarrow \alpha_5, \alpha_4 = \alpha_5\}$$

$$C_{\text{body}} = \{\alpha_6 = \alpha_7 \rightarrow \alpha_3, \alpha_6 = \alpha_8 \rightarrow \alpha_8, \alpha_7 = \alpha_9 \rightarrow \alpha_9\}$$

$$C = C'_{\text{let}} \cup C_{\text{body}} \cup \{\alpha_3 = \alpha_1\}$$

Einführung in Compilerbau

- Ein bisschen...
 - Lexikalische Analyse (Tokenisierung)
 - Syntaktische Analyse (Parsen)
 - Semantische Analyse (Optimierung)
 - Codegenerierung

- Ein bisschen...
 - Lexikalische Analyse (Tokenisierung)
 - Syntaktische Analyse (Parsen)
 - Semantische Analyse (Optimierung)
 - Codegenerierung
- Klausur:
 - SLL(k)-Form beweisen
 - Rekursiven Abstiegsparser schreiben/vervollständigen
 - First/Follow-Mengen berechnen
 - Java-Bytecode

$$SGML \rightarrow < id > Children < / >$$
$$Children \rightarrow \epsilon \mid SGML \ Children$$
$$\{<id><id></><id></></>, <id></>, \dots\} \in G$$

$$SGML \rightarrow < id > Children < / >$$
$$Children \rightarrow \epsilon \mid SGML \ Children$$
$$\{<id><id></><id></></>, <id></>, \dots\} \in G$$

- Begründen Sie formal, dass die obige Grammatik nicht in SLL(1)-Form ist (3P.).

$$SGML \rightarrow < id > Children < / >$$
$$Children \rightarrow \epsilon \mid SGML \ Children$$
$$\{<id><id></><id></></>, <id></>, \dots\} \in G$$

- Begründen Sie formal, dass die obige Grammatik nicht in SLL(1)-Form ist (3P.).
- Entwickeln Sie für [eine linksfaktorierte Version der obigen Grammatik] einen rekursiven Abstiegsparser (16P.).

Übersetzen Sie folgenden Java-Programmausschnitt in Java-Bytecode (10P.):

```
if (((a < b) || !((a < c) || (c < b))) && !(c < 0)) {  
    c = b + a;  
}
```

Übersetzen Sie folgenden Java-Programmausschnitt in Java-Bytecode (10P.):

```
if (((a < b) || !((a < c) || (c < b))) && !(c < 0)) {  
    c = b + a;  
}
```

Hinweise:

- Codeerzeugung für bedingte Sprünge: Folien 447ff.
- Um eine Bedingung der Form !cond zu übersetzen, reicht es, cond zu übersetzen und die Sprungziele anzupassen.

Compiler

- Maschine(-nmodell) versteht i.d.R. eingeschränkten Instruktionssatz
 - Es gibt/gab zwar auch mal CISC-Maschinen, heute ist sind aber RISC(-ähnliche) Prozessoren am weitesten verbreitet
 - Gründe: RISC-Prozessoren sind wesentlich einfacher (= billiger) zu bauen
- Programme in Maschinensprache sind i.d.R. für Menschen nicht einfach zu Schreiben.

Compiler: Motivation

- Maschine(-nmodell) versteht i.d.R. eingeschränkten Instruktionssatz
 - Es gibt/gab zwar auch mal CISC-Maschinen, heute ist sind aber RISC(-ähnliche) Prozessoren am weitesten verbreitet
 - Gründe: RISC-Prozessoren sind wesentlich einfacher (= billiger) zu bauen
- Programme in Maschinensprache sind i.d.R. für Menschen nicht einfach zu Schreiben.
- Also: Erfinde einfacher zu Schreibende (\approx mächtigere) Sprache, die dann in die Sprache der Maschine übersetzt wird.
- Diesen Übersetzungsschritt sollte optimalerweise ein Programm erledigen, da wir sonst auch einfach direkt Maschinensprache-Programme schreiben können.

- Übersetzer für formale Sprachen nennt man Compiler
- Beispiele:
 - C, Haskell, Rust, Go \rightarrow X86
 - Java, Clojure, Kotlin \rightarrow Java-Bytecode
 - TypeScript \rightarrow JavaScript
 - Python \rightarrow Python-AST

- Übersetzer für formale Sprachen nennt man Compiler
- Beispiele:
 - C, Haskell, Rust, Go \rightarrow X86
 - Java, Clojure, Kotlin \rightarrow Java-Bytecode
 - TypeScript \rightarrow JavaScript
 - Python \rightarrow Python-AST
- Interpreter kann man auch als Compiler kategorisieren, sie zählen aber i.A. nicht dazu

- Übersetzer für formale Sprachen nennt man Compiler
- Beispiele:
 - C, Haskell, Rust, Go \rightarrow X86
 - Java, Clojure, Kotlin \rightarrow Java-Bytecode
 - TypeScript \rightarrow JavaScript
 - Python \rightarrow Python-AST
- Interpreter kann man auch als Compiler kategorisieren, sie zählen aber i.A. nicht dazu
- Single-pass vs. Multi-pass
 - Single-pass: Eingabe wird einmal gelesen, Ausgabe währenddessen erzeugt (ältere Compiler)
 - Multi-pass: Eingabe wird in Zwischenschritten in verschiedene Repräsentationen umgewandelt
 - Quellsprache, Tokens, AST, Zwischensprache, Zielsprache

Lexikalische Analyse

```
int x1 = 123;  
print("123");
```

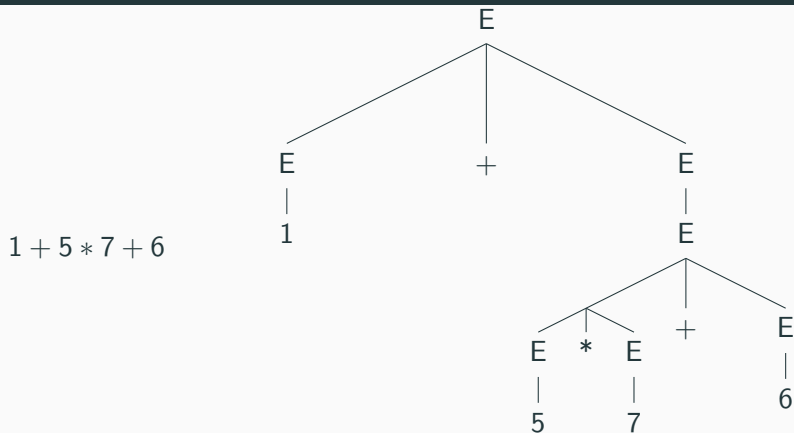
```
int, id[x1], assign,  
intlit[123], semi,  
id[print], lp,  
stringlit["123"], ...
```

- Lexikalische Analyse (Tokenisierung) verarbeitet eine Zeichensequenz in eine Liste von Tokens.
- Tokens sind Zeichengruppen, denen eine Semantik innewohnt:
 - `int` — Typ einer Ganzzahl
 - `=` — Zuweisungsoperator
 - `x1` — Variablen- oder Methodenname
 - `123` — Literal einer Ganzzahl
 - `"123"` — String-Literal
 - etc.
- Lösbar mit regulären Ausdrücken, Automaten

- Syntaktische Analyse stellt die unterliegende Struktur der bisher linear gelesenen Eingabe fest:
 - Blockstruktur von Programmen
 - Baumstruktur von HTML-Dateien
 - Header, Inhalt-Struktur von Mails
 - Verschachtelte mathematische Ausdrücke
- Syntaktische Analyse ist das größte Compiler-Thema in PP.

- Syntaktische Analyse stellt die unterliegende Struktur der bisher linear gelesenen Eingabe fest:
 - Blockstruktur von Programmen
 - Baumstruktur von HTML-Dateien
 - Header, Inhalt-Struktur von Mails
 - Verschachtelte mathematische Ausdrücke
- Syntaktische Analyse ist das größte Compiler-Thema in PP.
- Übliche Vorgehensweise (in PP):
 - Grammatik G erfinden
 - ggf. G in einfache Form G' bringen
 - rekursiven Abstiegsparser für G' implementieren
- Alternativ: Parser-Kombinatoren, Yacc, etc.

Beispiel: Mathematische Ausdrücke



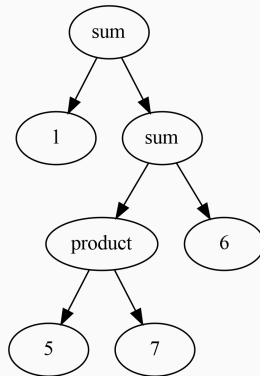
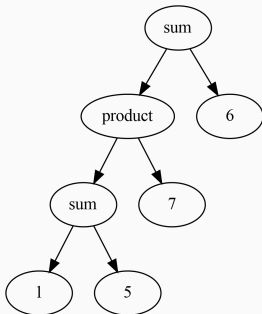
- Zu beachten: Punkt-vor-Strich (Präzedenz), Klammerung, etc.
- Nicht mehr mit regulären Ausdrücken lösbar
- „Offensichtliche“ Grammatik oft nicht einfach zu Parsen

Beispiel: Mathematische Ausdrücke

$$E \rightarrow \text{num} \mid (E) \mid E + E \mid E * E$$

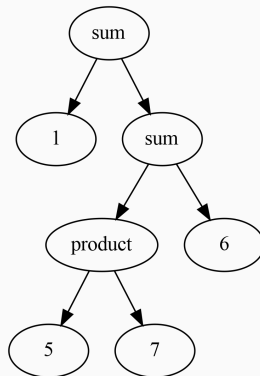
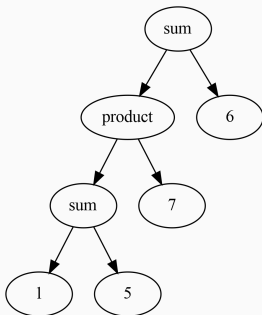
Beispiel: Mathematische Ausdrücke

$$E \rightarrow \text{num} \mid (E) \mid E + E \mid E * E$$



Beispiel: Mathematische Ausdrücke

$$E \rightarrow \text{num} \mid (E) \mid E + E \mid E * E$$



- Ableitungsbaum nicht eindeutig \leadsto schlecht
- Ableitungsbaum garantiert nicht Punkt-vor-Strich \leadsto schlecht

Wie zeichnen sich „gute“ Grammatiken aus?

Präzedenz, Linksfaktorisierung

Wie zeichnen sich „gute“ Grammatiken aus?

Operatorpräzedenz schon in Grammatik definiert:

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

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

$$F \rightarrow \text{num} \mid (E)$$

Präzedenz, Linksfaktorisierung

Wie zeichnen sich „gute“ Grammatiken aus?

Operatorpräzedenz schon in Grammatik definiert:

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

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

$$F \rightarrow \text{num} \mid (E)$$

Vermeidung von Linksrekursion (Linksfaktorisierung):

$$E \rightarrow T \ EList$$

$$EList \rightarrow \epsilon \mid + T \ EList \mid - T \ EList$$

$$T \rightarrow F \ TList$$

$$TList \rightarrow \epsilon \mid * F \ TList \mid / F \ TList$$

$$F \rightarrow \text{num} \mid (E)$$

First-/Followmenge, Indizmenge

$$EList \rightarrow \epsilon \mid + \ T \ EList \mid - \ T \ EList$$

Wie können wir bspw. bei *EList* entscheiden, welche Produktion anzuwenden ist?

First-/Followmenge, Indizmenge

$$EList \rightarrow \epsilon \mid + \ T \ EList \mid - \ T \ EList$$

Wie können wir bspw. bei $EList$ entscheiden, welche Produktion anzuwenden ist?

- \leadsto definiere Indizmenge $IM_k(A \rightarrow \alpha) = \text{First}_k(\alpha \text{Follow}_k(A))$
- Wenn nächste k Token in $IM_k(EList \rightarrow \phi) \leadsto$ weiter mit ϕ

First-/Followmenge, Indizmenge

$$EList \rightarrow \epsilon \mid + \ T \ EList \mid - \ T \ EList$$

Wie können wir bspw. bei $EList$ entscheiden, welche Produktion anzuwenden ist?

- \leadsto definiere Indizmenge $IM_k(A \rightarrow \alpha) = \text{First}_k(\alpha \text{Follow}_k(A))$
- Wenn nächste k Token in $IM_k(EList \rightarrow \phi) \leadsto$ weiter mit ϕ
- $IM_1(EList \rightarrow \epsilon) = \text{First}_1(\epsilon \text{Follow}_1(EList)) = \{), \#\}$
- $IM_1(EList \rightarrow + \ T \ EList) = \text{First}_1(+ \ T \ EList \text{Follow}_1(EList)) = \{+\}$
- $IM_1(EList \rightarrow - \ T \ EList) = \text{First}_1(- \ T \ EList \text{Follow}_1(EList)) = \{-\}$

First-/Followmenge, Indizmenge

$$EList \rightarrow \epsilon \mid + \ T \ EList \mid - \ T \ EList$$

Wie können wir bspw. bei $EList$ entscheiden, welche Produktion anzuwenden ist?

- \leadsto definiere Indizmenge $IM_k(A \rightarrow \alpha) = \text{First}_k(\alpha \text{Follow}_k(A))$
- Wenn nächste k Token in $IM_k(EList \rightarrow \phi) \leadsto$ weiter mit ϕ
- $IM_1(EList \rightarrow \epsilon) = \text{First}_1(\epsilon \text{Follow}_1(EList)) = \{\}, \#\}$
- $IM_1(EList \rightarrow + \ T \ EList) = \text{First}_1(+ \ T \ EList \text{Follow}_1(EList)) = \{+\}$
- $IM_1(EList \rightarrow - \ T \ EList) = \text{First}_1(- \ T \ EList \text{Follow}_1(EList)) = \{-\}$
- $\text{First}_k(A)$: Menge an möglichen ersten k Token in A
- $\text{Follow}_k(A)$: Menge an möglichen ersten k Token nach A

Grammatik ist in SLL(k)-Form

$$:\Leftrightarrow \forall A \rightarrow \alpha, A \rightarrow \beta \in P : IM_k(A \rightarrow \alpha) \cap IM_k(A \rightarrow \beta) = \emptyset$$

- SLL(k): Bei jedem Nichtterminal muss die zu wählende Produktion an den nächsten k Token wählbar sein.
- Nichtterminale mit nur einer Produktion sind hier irrelevant
- Schwierig daran: Follow-Mengen berechnen

Grammatik ist in SLL(k)-Form

$$:\Leftrightarrow \forall A \rightarrow \alpha, A \rightarrow \beta \in P : IM_k(A \rightarrow \alpha) \cap IM_k(A \rightarrow \beta) = \emptyset$$

- SLL(k): Bei jedem Nichtterminal muss die zu wählende Produktion an den nächsten k Token wählbar sein.
- Nichtterminale mit nur einer Produktion sind hier irrelevant
- Schwierig daran: Follow-Mengen berechnen

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

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

$$F \rightarrow \text{num} \mid (E)$$

- Begründet formal, dass obige Grammatik nicht SLL(1).
- Berechnet $\text{Follow}_1(N)$ für $N \in \{E, T, F\}$.

$$E \rightarrow T \text{ EList}$$

$$\text{EList} \rightarrow \epsilon \mid + T \text{ EList} \mid - T \text{ EList}$$

$$T \rightarrow F \text{ TList}$$

$$\text{TList} \rightarrow \epsilon \mid * F \text{ TList} \mid / F \text{ TList}$$

$$F \rightarrow \text{num} \mid (E)$$

- Yay, unsere Grammatik hat jetzt SLL(1)-Form!
- Aber was bringt das?

Rekursive Abstiegsparser

$$E \rightarrow T \text{ EList}$$

$$\text{EList} \rightarrow \epsilon \mid + T \text{ EList} \mid - T \text{ EList}$$

$$T \rightarrow F \text{ TList}$$

$$\text{TList} \rightarrow \epsilon \mid * F \text{ TList} \mid / F \text{ TList}$$

$$F \rightarrow \text{num} \mid (E)$$

- Yay, unsere Grammatik hat jetzt SLL(1)-Form!
- Aber was bringt das?
- G ist jetzt einfach ausprogrammierbar:
 - 1 Methode per Nichtterminal: `parseE()`, `parseEList()`, ...
 - `Token[k]`-Instanzattribut für k langen Lookahead
 - `expect(TokenType)`-Methode, um Token zu verarbeiten

Rekursive Abstiegsparser

$$E \rightarrow T \text{ EList}$$
$$\text{EList} \rightarrow \epsilon \mid + T \text{ EList} \mid - T \text{ EList}$$
$$T \rightarrow F \text{ TList}$$
$$\text{TList} \rightarrow \epsilon \mid * F \text{ TList} \mid / F \text{ TList}$$
$$F \rightarrow \text{num} \mid (E)$$

- Yay, unsere Grammatik hat jetzt SLL(1)-Form!
- Aber was bringt das?
- G ist jetzt einfach ausprogrammierbar:
 - 1 Methode per Nichtterminal: `parseE()`, `parseEList()`, ...
 - `Token[k]`-Instanzattribut für k langen Lookahead
 - `expect(TokenType)`-Methode, um Token zu verarbeiten
- Vervollständigt `demos/java/exrparser/ExprParser.java!`

- PP beschäftigt sich (bis auf Typinferenz) nur kurz mit semantischer Analyse
- Hier geht es um Optimierungen, Typchecks, etc.
- \leadsto weiterführende (Master-)Vorlesungen am IPD

Ende

- Di, 16.02.: Letztes Tut
- Mi, 17.02.: Letzte Vorlesung
- Themen: Mehr Syntaxanalyse (Beispiele), Codeerzeugung