

Script generated by TTT

Title: Petter: Compilerbau (25.04.2016)

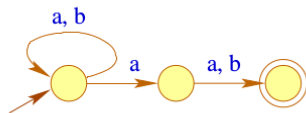
Date: Mon Apr 25 14:29:36 CEST 2016

Duration: 89:13 min

Pages: 53

Chapter 4: Turning NFAs deterministic

The expected outcome:



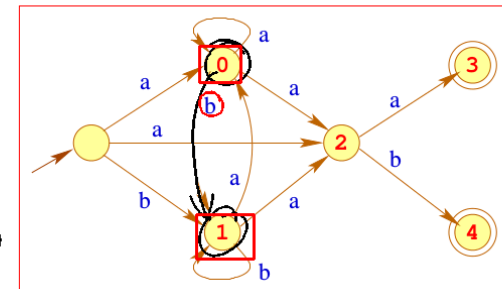
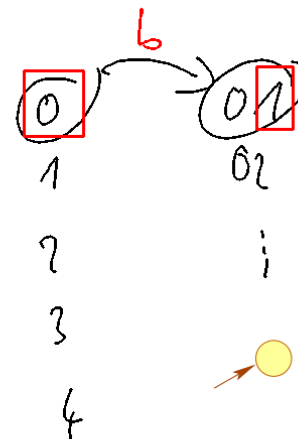
Remarks:

- ideal automaton would be even more compact (→ *Antimirov automata*, *Follow Automata*)
- but Berry-Sethi is rather directly constructed
- Anyway, we need a **deterministic** version

⇒ Powerset-Construction

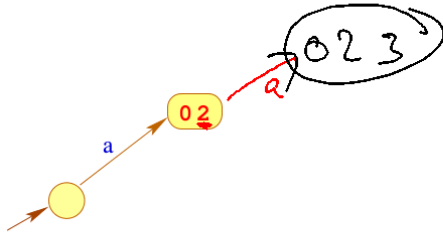
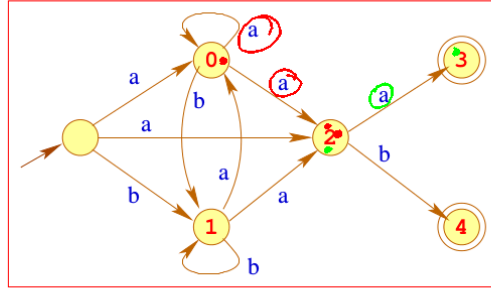
Powerset Construction

... for example:



Powerset Construction

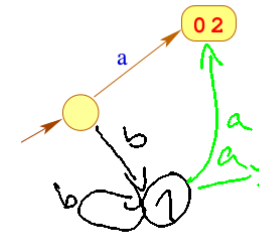
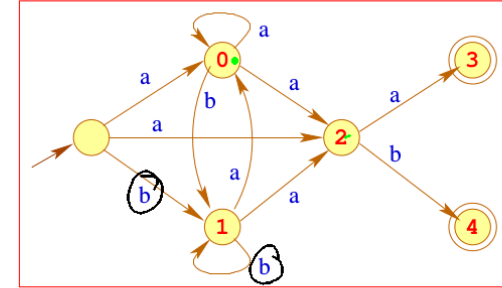
... for example:



46 / 282

Powerset Construction

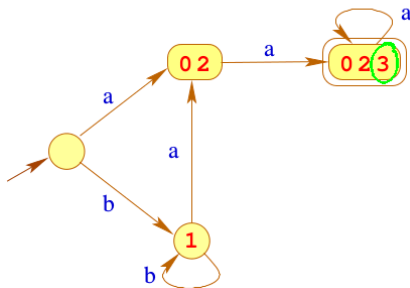
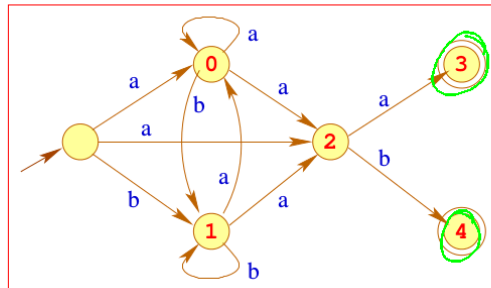
... for example:



46 / 282

Powerset Construction

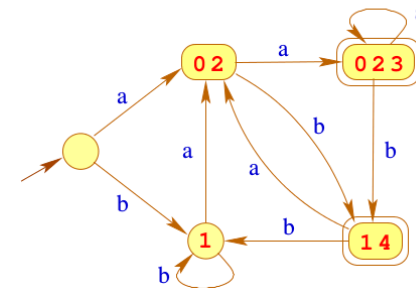
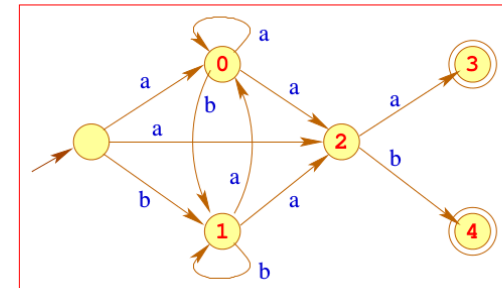
... for example:



46 / 282

Powerset Construction

... for example:



46 / 282

Powerset Construction

Theorem:

For every non-deterministic automaton $A = (Q, \Sigma, \delta, I, F)$ we can compute a deterministic automaton $\mathcal{P}(A)$ with

$$\mathcal{L}(A) = \mathcal{L}(\mathcal{P}(A))$$

47/282

Powerset Construction

Theorem:

For every non-deterministic automaton $A = (Q, \Sigma, \delta, I, F)$ we can compute a deterministic automaton $\mathcal{P}(A)$ with

$$\mathcal{L}(A) = \mathcal{L}(\mathcal{P}(A))$$

47/282

Construction:

States: Powersets of Q ;

Start state: I ;

Final states: $\{Q' \subseteq Q \mid Q' \cap F \neq \emptyset\}$;

Transitions: $\delta_{\mathcal{P}}[Q', a] = \{q \in Q \mid \exists p \in Q' : [p, a, q] \in \delta\}$.

47/282

Powerset Construction

Observation:

There are exponentially many powersets of Q

- Idea: Consider only **contributing** powersets. Starting with the set $Q_{\mathcal{P}} = \{I\}$ we only add further states **by need** ...
- i.e., whenever we can reach them from a state in $Q_{\mathcal{P}}$
- However, the resulting automaton can become enormously **huge** ... which is (sort of) not happening in **practice**

48/282

Powerset Construction

Observation:

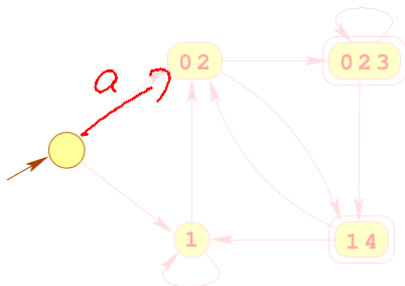
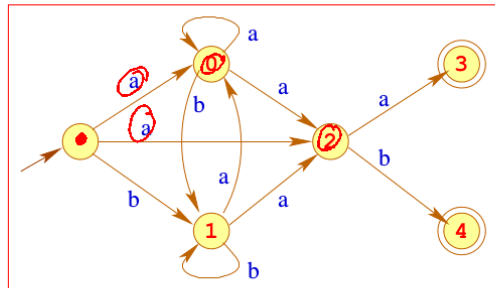
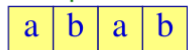
There are exponentially many powersets of Q

- Idea: Consider only **contributing** powersets. Starting with the set $Q_{\mathcal{P}} = \{I\}$ we only add further states **by need** ...
- i.e., whenever we can reach them from a state in $Q_{\mathcal{P}}$
- However, the resulting automaton can become enormously **huge** ... which is (sort of) not happening in **practice**
- Therefore, in tools like **grep** a regular expression's **DFA** is never created!
- Instead, only the sets, directly necessary for interpreting the input are generated **while processing the input**

48/282

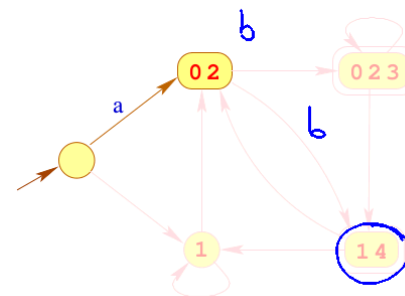
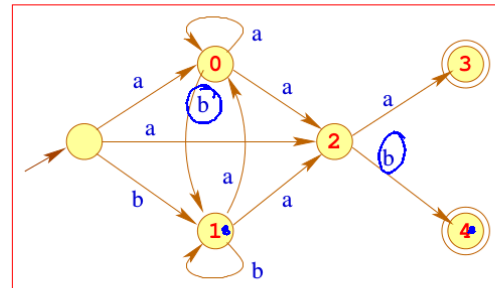
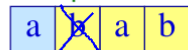
Powerset Construction

... for example:



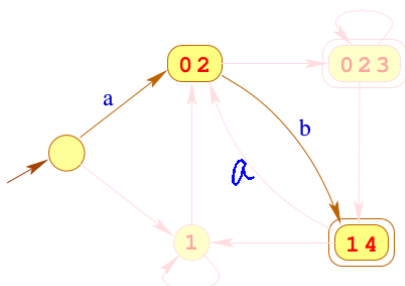
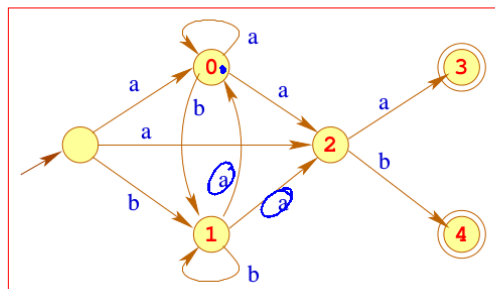
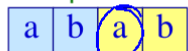
Powerset Construction

... for example:



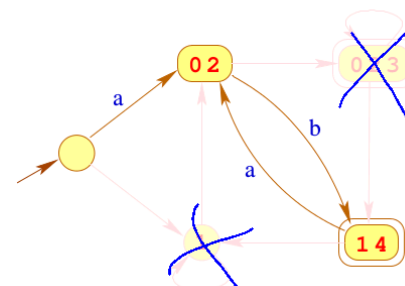
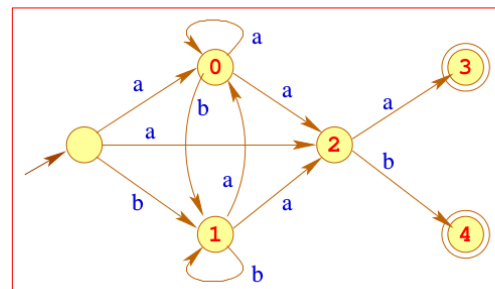
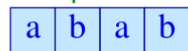
Powerset Construction

... for example:



Powerset Construction

... for example:



Remarks:

- For an input sequence of length n , maximally $O(n)$ sets are generated
- Once a set/edge of the DFA is generated, they are stored within a hash-table.
- Before generating a new transition, we check this table for already existing edges with the desired label.

50/282

Remarks:

- For an input sequence of length n , maximally $O(n)$ sets are generated
- Once a set/edge of the DFA is generated, they are stored within a hash-table.
- Before generating a new transition, we check this table for already existing edges with the desired label.

Summary:

Theorem:

For each regular expression e we can compute a deterministic automaton $A = \mathcal{P}(A_e)$ with

$$\mathcal{L}(A) = \llbracket e \rrbracket$$

50/282

Lexical Analysis

Chapter 5: Scanner design

51/282

Scanner design

Input (simplified):

a set of rules:

e_1 { action₁ }
 e_2 { action₂ }
...
 e_k { action_k }

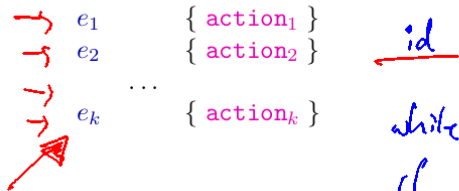
output ("codas" "dies")

52/282

Scanner design

if

Input (simplified): a set of rules:



Output: a program,

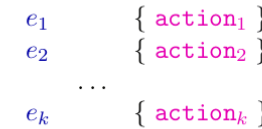
- ... reading a maximal prefix w from the input, that satisfies $e_1 \mid \dots \mid e_k$;
- ... determining the minimal i , such that $w \in [e_i]$;
- ... executing action_i for w .

52/282

Scanner design

while a

Input (simplified): a set of rules:



while
;
id

Output: a program,

- ... reading a maximal prefix w from the input, that satisfies $e_1 \mid \dots \mid e_k$;
- ... determining the minimal i , such that $w \in [e_i]$;
- ... executing action_i for w .

52/282

Implementation:

Idea:

- Create the DFA $\mathcal{P}(A_e) = (Q, \Sigma, \delta, q_0, F)$ for the expression $e = (e_1 \mid \dots \mid e_k)$;
- Define the sets:

$$F_1 = \{q \in F \mid q \cap \text{last}[e_1] \neq \emptyset\}$$

$$F_2 = \{q \in (F \setminus F_1) \mid q \cap \text{last}[e_2] \neq \emptyset\}$$
 ...

$$F_k = \{q \in (F \setminus (F_1 \cup \dots \cup F_{k-1})) \mid q \cap \text{last}[e_k] \neq \emptyset\}$$
- For input w we find: $\delta^*(q_0, w) \in F_i$ iff the scanner must execute action_i for w

53/282

Implementation:

Idea (cont'd):

- The scanner manages two pointers $\langle A, B \rangle$ and the related states $\langle q_A, q_B \rangle \dots$
- Pointer A points to the last position in the input, after which a state $q_A \in F$ was reached;
- Pointer B tracks the current position.

s t d o u t . w r i t e l n (" H a l l o ") ;

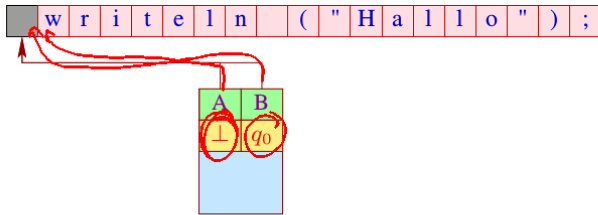


54/282

Implementation:

Idea (cont'd):

- The scanner manages two pointers $\langle A, B \rangle$ and the related states $\langle q_A, q_B \rangle \dots$
- Pointer A points to the last position in the input, after which a state $q_A \in F$ was reached;
- Pointer B tracks the current position.



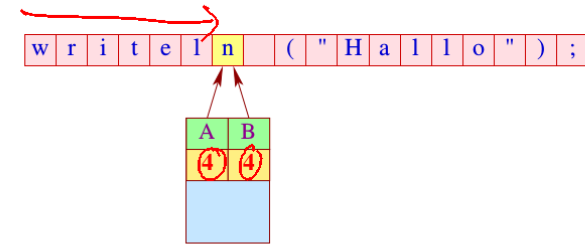
54 / 282

Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

$$\begin{aligned} B &:= A; & A &:= \perp; \\ q_B &:= q_0; & q_A &:= \perp \end{aligned}$$



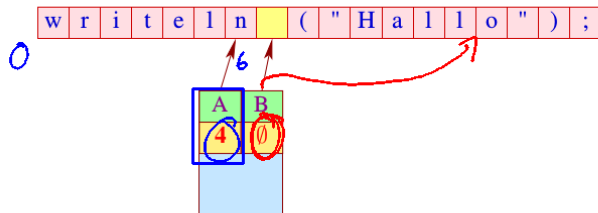
55 / 282

Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

$$\begin{aligned} B &:= A; & A &:= \perp; \\ q_B &:= q_0; & q_A &:= \perp \end{aligned}$$



55 / 282

Extension: States

- Now and then, it is handy to differentiate between particular **scanner states**.
- In different states, we want to recognize different token classes with different precedences.
- Depending on the consumed input, the scanner state can be changed

Example: Comments

Within a comment, identifiers, constants, comments, ... are ignored

56 / 282

Input (generalized): a set of rules:

```

<state> {
  e1 { action1 yybegin(state1); }
  e2 { action2 yybegin(state2); }
  ...
  ek { actionk yybegin(statek); }
}
    
```

- The statement `yybegin (statei);` resets the current state to `statei`.
- The start state is called (e.g. flex JFlex) `YYINITIAL`.

... for example:

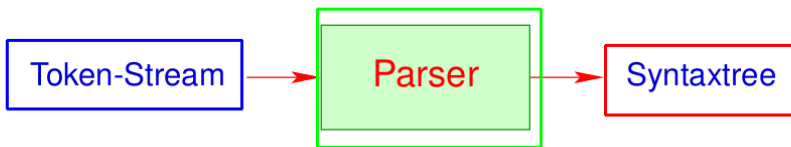
```

<YYINITIAL> { /*" { yybegin COMMENT; }
<COMMENT> { "*/" { yybegin YYINITIAL; }
             . | \n { }
    
```

Topic:

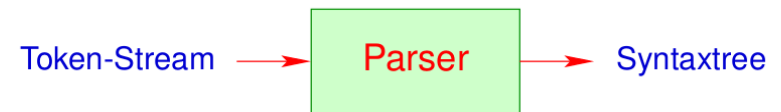
Syntactic Analysis

Syntactic Analysis



- Syntactic analysis tries to integrate Tokens into larger program units.

Syntactic Analysis



- Syntactic analysis tries to integrate Tokens into larger program units.
- Such units may possibly be:
 - Expressions;
 - Statements;
 - Conditional branches;
 - loops; ...

Discussion:

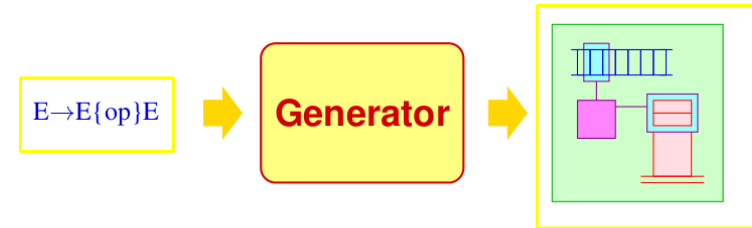
In general, parsers are not developed by hand, but **generated** from a specification:



61 / 282

Discussion:

In general, parsers are not developed by hand, but **generated** from a specification:



Specification of the hierarchical structure: contextfree grammars

Generated implementation: Pushdown automata + X

61 / 282

Syntactic Analysis

Chapter 1: Basics of Contextfree Grammars

62 / 282

Basics: Context-free Grammars

- Programs of programming languages can have arbitrary numbers of tokens, but only finitely many **Token-classes**.
- This is why we choose the set of **Token-classes** to be the finite alphabet of **terminals T** .
- The nested structure of program components can be described elegantly via **context-free** grammars...

63 / 282

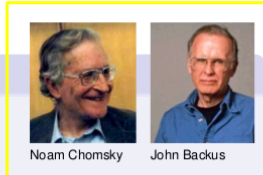
Basics: Context-free Grammars

- Programs of programming languages can have arbitrary numbers of tokens, but only finitely many **Token-classes**.
- This is why we choose the set of **Token-classes** to be the finite alphabet of terminals T .
- The nested structure of program components can be described elegantly via **context-free** grammars...

Definition: Context-Free Grammar

A **context-free grammar (CFG)** is a 4-tuple $G = (N, T, P, S)$ with:

- N the set of **nonterminals**,
- T the set of **terminals**,
- P the set of **productions or rules**, and
- $S \in N$ the **start symbol**



63/282

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \quad \text{with} \quad A \in N, \quad \alpha \in (N \cup T)^*$$

64/282

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \quad \text{with} \quad A \in N, \quad \alpha \in (N \cup T)^*$$

... for example:

$$\begin{array}{l} S \rightarrow aSb \\ S \rightarrow \epsilon \end{array}$$

Specified language: $\{a^n b^n \mid n \geq 0\}$

64/282

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \quad \text{with} \quad A \in N, \quad \alpha \in (N \cup T)^*$$

... for example:

$$\begin{array}{l} S \rightarrow aSb \\ S \rightarrow \epsilon \end{array}$$

Specified language: $\{a^n b^n \mid n \geq 0\}$

Conventions:

In examples, we specify nonterminals and terminals in general **implicitly**:

- nonterminals are: $A, B, C, \dots, \langle \text{exp} \rangle, \langle \text{stmt} \rangle, \dots;$
- terminals are: $a, b, c, \dots, \text{int}, \text{name}, \dots;$

64/282

... a practical example:

```

S      → ⟨stmt⟩
⟨stmt⟩ → ⟨if⟩ | ⟨while⟩ | ⟨rexp⟩;
⟨if⟩   → if ( ⟨rexp⟩ ) ⟨stmt⟩ else ⟨stmt⟩
⟨while⟩ → while ( ⟨rexp⟩ ) ⟨stmt⟩
⟨rexp⟩ → int | ⟨lexp⟩ | ⟨lexp⟩ = ⟨rexp⟩ | ...
⟨lexp⟩ → name | ...
    
```

... a practical example:

```

S      → ⟨stmt⟩
⟨stmt⟩ → ⟨if⟩0 | ⟨while⟩1 | ⟨rexp⟩2;
⟨if⟩   → if ( ⟨rexp⟩ ) ⟨stmt⟩ else ⟨stmt⟩
⟨while⟩ → while ( ⟨rexp⟩ ) ⟨stmt⟩
⟨rexp⟩ → int | ⟨lexp⟩ | ⟨lexp⟩ = ⟨rexp⟩ | ...
⟨lexp⟩ → name | ...
    
```

More conventions:

- For every nonterminal, we collect the right hand sides of rules and list them together.
- The j -th rule for A can be identified via the pair (A, j) (with $j \geq 0$).

$(\langle \text{stmt} \rangle, 3) \stackrel{1}{=} \langle \text{stmt} \rangle \rightarrow \langle \text{while} \rangle$

Pair of grammars:

E	\rightarrow	$E+E$	$ $	$E * E$	$ $	(E)	$ $	name	$ $	int
E	\rightarrow	$E+T$	$ $	T						
T	\rightarrow	$T * F$	$ $	F						
F	\rightarrow	(E)	$ $	name	$ $	int				

Both grammars describe the same language

Pair of grammars:

E	\rightarrow	$E+E^0$	$ $	$E * E^1$	$ $	$(E)^2$	$ $	name ³	$ $	int ⁴
E	\rightarrow	$E+T^0$	$ $	T^1						
T	\rightarrow	$T * F^0$	$ $	F^1						
F	\rightarrow	$(E)^0$	$ $	name ¹	$ $	int ²				

Both grammars describe the same language

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example: \underline{E}

67/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example: $\underline{E} \rightarrow \underline{E} + T$

67/282

Pair of grammars:

$E \rightarrow E+E$	$E * E$	(E)	name	int
$E \rightarrow E+T$	T			
$T \rightarrow T * F$	F			
$F \rightarrow (E)$	name	int		

Both grammars describe **the same language**

66/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example: $\underline{E} \rightarrow \underline{E} + T$
 $\rightarrow \underline{T} + T$

67/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example:

$$\begin{aligned} E &\rightarrow E + T \\ &\rightarrow T + T \\ &\rightarrow T * F + T \end{aligned}$$

67/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example:

$$\begin{aligned} E &\rightarrow E + T \\ &\rightarrow T + T \\ &\rightarrow T * F + T \\ &\rightarrow T * \text{int} + T \end{aligned}$$

67/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example:

$$\begin{aligned} E &\rightarrow E + T \\ &\rightarrow T + T \\ &\rightarrow T * F + T \\ &\rightarrow T * \text{int} + T \\ &\rightarrow F * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + F \\ &\rightarrow \text{name} * \text{int} + \text{int} \end{aligned}$$

*(Handwritten annotations: red arrows and boxes highlight the derivation steps, and a red rule $T \rightarrow T * F$ is written to the right.)*

Definition

The derivation relation \rightarrow is a relation on words over $N \cup T$, with

$\alpha \rightarrow \alpha'$ iff $\alpha = \alpha_1 A \alpha_2 \wedge \alpha' = \alpha_1 \beta \alpha_2$ for an $A \rightarrow \beta \in P$

67/282

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example:

$$\begin{aligned} E &\rightarrow E + T \\ &\rightarrow T + T \\ &\rightarrow T * F + T \\ &\rightarrow T * \text{int} + T \\ &\rightarrow F * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + F \\ &\rightarrow \text{name} * \text{int} + \text{int} \end{aligned}$$

Definition

The derivation relation \rightarrow is a relation on words over $N \cup T$, with

$\alpha \rightarrow \alpha'$ iff $\alpha = \alpha_1 A \alpha_2 \wedge \alpha' = \alpha_1 \beta \alpha_2$ for an $A \rightarrow \beta \in P$

The **reflexive** and **transitive** closure of \rightarrow is denoted as: \rightarrow^*

67/282

Derivation

Remarks:

- The relation \rightarrow depends on the grammar
- In each step of a derivation, we may choose:
 - * a spot, determining **where** we will rewrite.
 - * a rule, determining **how** we will rewrite.
- The language, specified by G is:

$$\mathcal{L}(G) = \{w \in T^* \mid S \rightarrow^* w\}$$

Derivation

Remarks:

- The relation \rightarrow depends on the grammar
- In each step of a derivation, we may choose:
 - * a spot, determining **where** we will rewrite.
 - * a rule, determining **how** we will rewrite.
- The language, specified by G is:

$$\mathcal{L}(G) = \{w \in T^* \mid S \rightarrow^* w\}$$

Attention:

The order, in which disjunct fragments are rewritten is not relevant.