

Script generated by TTT

Title: Petter: Compilerbau (04.05.2015)

Date: Mon May 04 14:20:08 CEST 2015

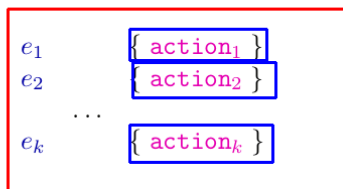
Duration: 95:37 min

Pages: 62

Chapter 5: Scanner design

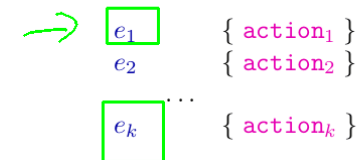
Scanner design

Input (simplified): a set of rules:



Scanner design

Input (simplified): a set of rules:



Output: a program,

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

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 / 82

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 ") ;

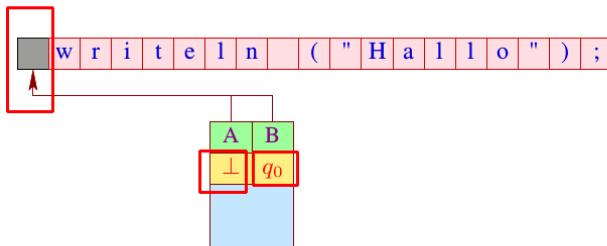


54 / 82

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 / 82

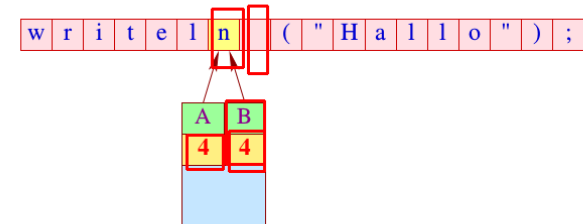
Implementation:

Idea (cont'd):

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

$$B := A; \quad A := \perp;$$

$$q_B := q_0; \quad q_A := \perp$$

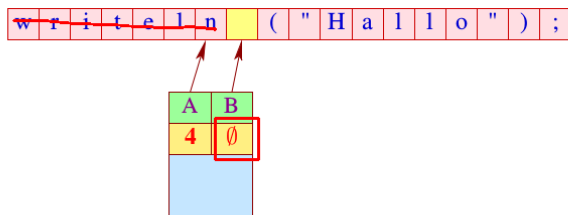


55 / 82

Implementation:

Idea (cont'd):

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

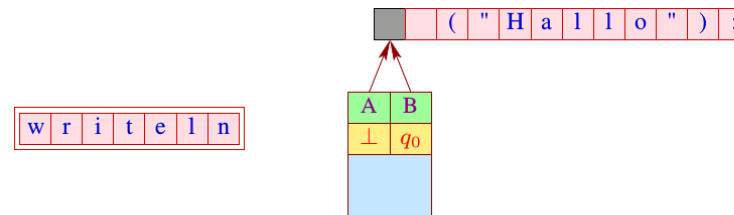
$$\begin{array}{l} B := A; \quad A := \perp; \\ q_B := q_0; \quad q_A := \perp \end{array}$$


55 / 82

Implementation:

Idea (cont'd):

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

$$\begin{array}{l} B := A; \quad A := \perp; \\ q_B := q_0; \quad q_A := \perp \end{array}$$


55 / 82

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 / 82

Input (generalized):

a set of rules:

$$\langle \text{state} \rangle \left\{ \begin{array}{l} e_1 \quad \{ \text{action}_1 \quad \text{yybegin}(\text{state}_1); \} \\ e_2 \quad \{ \text{action}_2 \quad \text{yybegin}(\text{state}_2); \} \\ \dots \\ e_k \quad \{ \text{action}_k \quad \text{yybegin}(\text{state}_k); \} \end{array} \right\}$$

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

... for example:

$$\begin{array}{l} \langle \text{YYINITIAL} \rangle \quad \{ \text{"/"} \text{"/} \quad \{ \text{yybegin}(\text{COMMENT}); \} \\ \langle \text{COMMENT} \rangle \quad \{ \text{"} \text{*"} \text{"} \quad \{ \text{yybegin}(\text{YYINITIAL}); \} \\ \quad \cdot \mid \backslash \text{n} \quad \{ \} \\ \quad \} \end{array}$$

57 / 82

Remarks:

- “.” matches all characters different from “\n”.
- For every state we generate the scanner respectively.
- Method `yybegin (STATE);` switches between different scanners.
- Comments might be directly implemented as (admittedly overly complex) token-class.
- Scanner-states are especially handy for implementing **preprocessors**, expanding special fragments in regular programs.

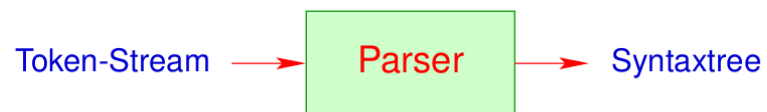
58 / 82

Topic:

Syntactic Analysis

59 / 82

Syntactic Analysis



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

60 / 82

Discussion:

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



61 / 82

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 / 82

Syntactic Analysis

Chapter 1: Basics of Contextfree Grammars

62 / 82

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 / 82

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...

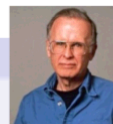
Definition: Context-Free Grammar

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**



Noam Chomsky



John Backus

63 / 82

Conventions

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

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

64 / 82

Conventions

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

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

... for example:

$$\begin{aligned} S &\rightarrow a S b \\ S &\rightarrow \epsilon \end{aligned}$$

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

64 / 82

Conventions

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

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

... for example:

$$\begin{aligned} S &\rightarrow a S b \\ S &\rightarrow \epsilon \end{aligned}$$

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 / 82

... a practical example:

$$\begin{aligned} S &\rightarrow \langle \text{stmt} \rangle \\ \langle \text{stmt} \rangle &\rightarrow \langle \text{if} \rangle \mid \langle \text{while} \rangle \mid \langle \text{rexpr} \rangle; \\ \langle \text{if} \rangle &\rightarrow \text{if} (\langle \text{rexpr} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ \langle \text{while} \rangle &\rightarrow \text{while} (\langle \text{rexpr} \rangle) \langle \text{stmt} \rangle \\ \langle \text{rexpr} \rangle &\rightarrow \text{int} \mid \langle \text{lexpr} \rangle \mid \langle \text{lexpr} \rangle = \langle \text{rexpr} \rangle \mid \dots \\ \langle \text{lexpr} \rangle &\rightarrow \text{name} \mid \dots \end{aligned}$$

65 / 82

... a practical example:

$S \rightarrow \langle \text{stmt} \rangle$
 $\langle \text{stmt} \rangle \rightarrow \langle \text{if} \rangle \mid \langle \text{while} \rangle \mid \langle \text{rexp} \rangle;$
 $\langle \text{if} \rangle \rightarrow \text{if} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$
 $\langle \text{while} \rangle \rightarrow \text{while} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle$
 $\langle \text{rexp} \rangle \rightarrow \text{int} \mid \langle \text{lexp} \rangle \mid \langle \text{lexp} \rangle = \langle \text{rexp} \rangle \mid \dots$
 $\langle \text{lexp} \rangle \rightarrow \text{name} \mid \dots$

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 $j \geq 0$).

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: E

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

66 / 82

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 / 82

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 / 82

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 / 82

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:

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

67 / 82

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:

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

67 / 82

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:

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

67 / 82

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:

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

67 / 82

Definition

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

$$\alpha \rightarrow \alpha' \quad \text{iff} \quad \alpha = \alpha_1 \underline{A} \alpha_2 \wedge \alpha' = \alpha_1 \underline{\beta} \alpha_2 \quad \text{for an} \quad \underline{A} \rightarrow \underline{\beta} \in P$$

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:

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

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 / 82

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:

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

68 / 82

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:

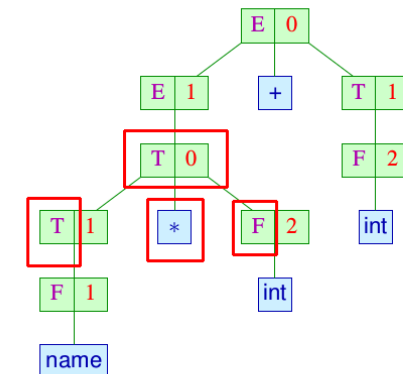
The order, in which disjunct fragments are rewritten is not relevant.

68 / 82

Derivation Tree

Derivations of a symbol are represented as **derivation trees**:

... for example:

$$\begin{aligned} \underline{E} &\rightarrow^0 \underline{E} + T \\ &\rightarrow^1 \underline{T} + T \\ &\rightarrow^0 \underline{T} * \underline{F} + T \\ &\rightarrow^2 \underline{T} * \underline{\text{int}} + T \\ &\rightarrow^1 \underline{F} * \underline{\text{int}} + T \\ &\rightarrow^1 \underline{\text{name}} * \underline{\text{int}} + T \\ &\rightarrow^1 \underline{\text{name}} * \underline{\text{int}} + \underline{F} \\ &\rightarrow^2 \underline{\text{name}} * \underline{\text{int}} + \underline{\text{int}} \end{aligned}$$


A **derivation tree** for $A \in N$:

inner nodes: rule applications

root: rule application for A

leaves: terminals or ϵ

The successors of (B, i) correspond to right hand sides of the rule

69 / 82

Special Derivations

Attention:

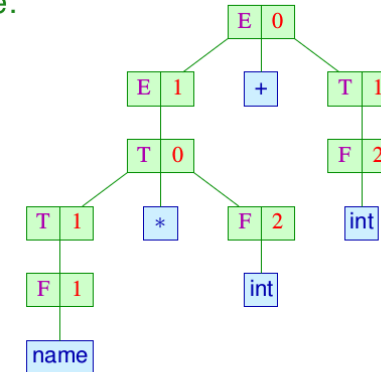
In contrast to arbitrary derivations, we find special ones, always rewriting the **leftmost** (or rather **rightmost**) occurrence of a nonterminal.

- These are called **leftmost** (or rather **rightmost**) derivations and are denoted with the index **L** (or **R** respectively).
- Leftmost (or rightmost) derivations correspond to a left-to-right (or right-to-left) **preorder**-DFS-traversal of the derivation tree.
- **Reverse rightmost** derivations correspond to a left-to-right **postorder**-DFS-traversal of the derivation tree

70 / 82

Special Derivations

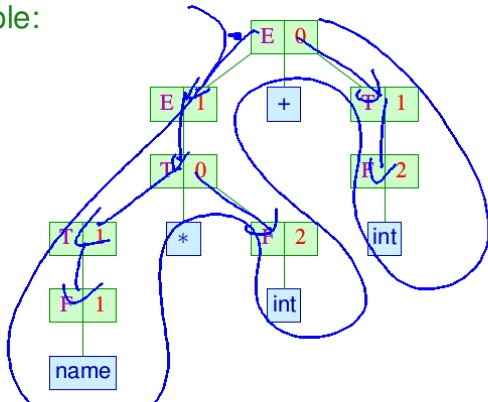
... for example:



71 / 82

Special Derivations

... for example:

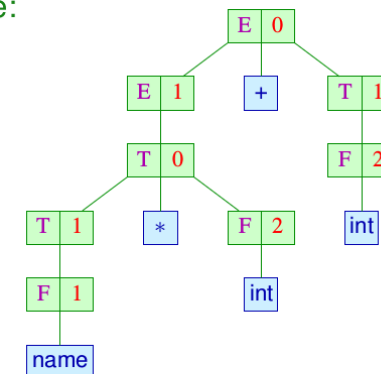


Leftmost derivation: $(E, 0) (E, 1) (T, 0) (T, 1) (F, 1) (F, 2) (T, 1) (F, 2)$

71 / 82

Special Derivations

... for example:



Leftmost derivation: $(E, 0) (E, 1) (T, 0) (T, 1) (F, 1) (F, 2) (T, 1) (F, 2)$

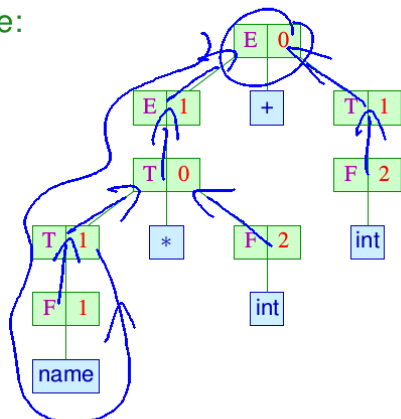
Rightmost derivation:

$(E, 0) (T, 1) (F, 2) (E, 1) (T, 0) (F, 2) (T, 1) (F, 1)$

71 / 82

Special Derivations

... for example:



Leftmost derivation: $(E, 0)(E, 1)(T, 0)(T, 1)(F, 1)(F, 2)(T, 1)(F, 2)$

Rightmost derivation:

$(E, 0)(T, 1)(F, 2)(E, 1)(T, 0)(F, 2)(T, 1)(F, 1)$

Reverse rightmost derivation:

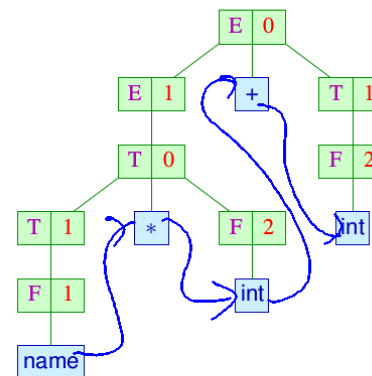
~~$(F, 1)(T, 1)(F, 2)(T, 0)(E, 1)(F, 2)(T, 1)(E, 0)$~~

71 / 82

Unique Grammars

The concatenation of leaves of a derivation tree t are often called yield(t).

... for example:



gives rise to the concatenation:

name * int + int.

72 / 82

Unique grammars

Definition:

Grammar G is called **unique**, if for every $w \in T^*$ there is maximally one derivation tree t of S with $\text{yield}(t) = w$.

... in our example:

$E \rightarrow E + E^0$	$E * E^1$	$(E)^2$	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		

The first one is ambiguous, the second one is unique

73 / 82

Conclusion:

- A derivation tree represents a possible hierarchical structure of a word.
- For programming languages, only those grammars with a unique structure are of interest.
- Derivation trees are one-to-one corresponding with leftmost derivations as well as (reverse) rightmost derivations.

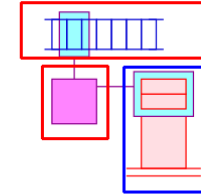
- **Leftmost derivations** correspond to a **top-down** reconstruction of the syntax tree.
- **Reverse rightmost derivations** correspond to a **bottom-up** reconstruction of the syntax tree.

74 / 82

Chapter 2: Basics of Pushdown Automata

Basics of Pushdown Automata

Languages, specified by context free grammars are accepted by **Pushdown Automata**:



The pushdown is used e.g. to verify correct nesting of braces.

Example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

Example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

Conventions:

- We do **not** differentiate between pushdown symbols and states
- The rightmost / upper pushdown symbol represents the state
- Every transition consumes / modifies the upper part of the pushdown

Definition: Pushdown Automaton

A pushdown automaton (PDA) is a tuple $M = (Q, T, \delta, q_0, F)$ with:



- Q a finite set of states;
- T an input alphabet;
- $q_0 \in Q$ the start state;
- $F \subseteq Q$ the set of final states and
- $\delta \subseteq Q^+ \times (T \cup \{\epsilon\}) \times Q^*$ a finite set of transitions

Definition: Pushdown Automaton

A pushdown automaton (PDA) is a tuple $M = (Q, T, \delta, q_0, F)$ with:



- Q a finite set of states;
- T an input alphabet;
- $q_0 \in Q$ the start state;
- $F \subseteq Q$ the set of final states and
- $\delta \subseteq Q^+ \times (T \cup \{\epsilon\}) \times Q^*$ a finite set of transitions

We define **computations** of pushdown automata with the help of transitions; a particular **computation state** (the current **configuration**) is a pair:

$$(\gamma, w) \in Q^* \times T^*$$

consisting of the **pushdown content** and the **remaining input**.

... for example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

$$(0, a a a b b b)$$

... for example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

$$(0, a a a b b b) \vdash (11, a a b b b)$$

... for example:

States: 0,1,2
Start state: 0
Final states: 0,2

0	a	11
1	a	11
11	b	2
12	b	2

(0, aaabb) ⊢ (11, aabb)
⊢ (111, abbb)

... for example:

States: 0,1,2
Start state: 0
Final states: 0,2

0	a	11
1	a	11
11	b	2
12	b	2

(0, aaabb) ⊢ (11, aabb)
⊢ (111, abbb)
⊢ (1111, bbb)

... for example:

States: 0,1,2
Start state: 0
Final states: 0,2

0	a	11
1	a	11
11	b	2
12	b	2

(0, aaabb) ⊢ (11, aabb)
⊢ (111, abbb)
⊢ (1111, bbb)
⊢ (112, bb)

... for example:

$\{a^n b^n \mid n \geq 0\}$

States: 0,1,2
Start state: 0
Final states: 0,2

0	a	11
1	a	11
11	b	2
12	b	2

~~ξ~~ ξ

(0, aaabb) ⊢ (11, aabb)
⊢ (111, abbb)
⊢ (1111, bbb)
⊢ (112, bb)
⊢ (12, b)
⊢ (2, ε)

A computation step is characterized by the relation $\vdash \subseteq (Q^* \times T^*)^2$ with

$$(\alpha\gamma, \boxed{w}) \vdash (\alpha\gamma', \boxed{w}) \text{ for } (\gamma, \boxed{w}, \gamma') \in \delta$$

80 / 82

A computation step is characterized by the relation $\vdash \subseteq (Q^* \times T^*)^2$ with

$$(\alpha\gamma, xw) \vdash (\alpha\gamma', w) \text{ for } (\gamma, x, \gamma') \in \delta$$

Remarks:

- The relation \vdash depends of the pushdown automaton M
- The reflexive and transitive closure of \vdash is denoted by \vdash^*
- Then, the language accepted by M is

$$\mathcal{L}(M) = \{w \in T^* \mid \exists f \in F : \boxed{q_0}, \boxed{w} \vdash^* \boxed{f}, \boxed{\epsilon}\}$$

80 / 82

A computation step is characterized by the relation $\vdash \subseteq (Q^* \times T^*)^2$ with

$$(\alpha\gamma, xw) \vdash (\alpha\gamma', w) \text{ for } (\gamma, x, \gamma') \in \delta$$

Remarks:

- The relation \vdash depends of the pushdown automaton M
- The reflexive and transitive closure of \vdash is denoted by \vdash^*
- Then, the language accepted by M is

$$\mathcal{L}(M) = \{w \in T^* \mid \exists f \in F : (q_0, w) \vdash^* (f, \epsilon)\}$$

We accept with a **final state** together with **empty input**.

80 / 82