

Script generated by TTT

Title: Simon: Compilerbau (29.04.2013)

Date: Mon Apr 29 14:17:08 CEST 2013

Duration: 94:24 min

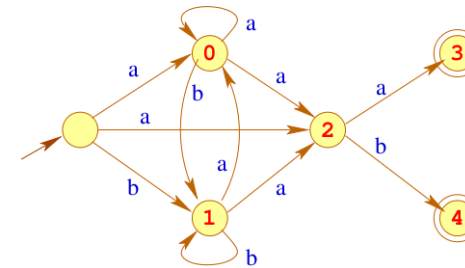
Pages: 63

Lexical Analysis

Chapter 4: Turning NFAs deterministic

Berry-Sethi Approach

... for example:



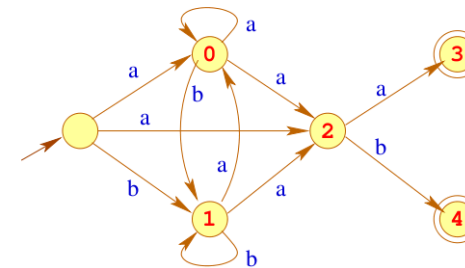
Remarks:

- This construction is known as **Berry-Sethi-** or **Glushkov-**construction.
- It is used for **XML** to define **Content Models**
- The result may not be, what we had in mind...

44 / 150

Berry-Sethi Approach

... for example:

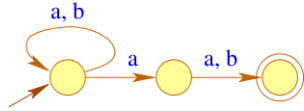


Remarks:

- This construction is known as **Berry-Sethi-** or **Glushkov-**construction.
- It is used for **XML** to define **Content Models**
- The result may not be, what we had in mind...

44 / 150

The expected outcome:



Remarks:

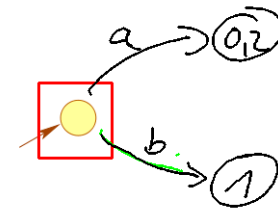
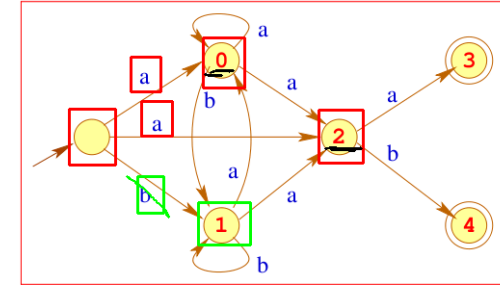
- ingoing edges do not necessarily have the same label here
- but Berry-Sethi is rather directly constructed
- Anyway, we need a **deterministic** technique

⇒ **Powerset-Construction**

46 / 150

Powerset Construction

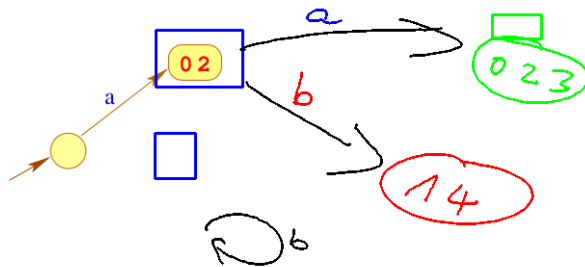
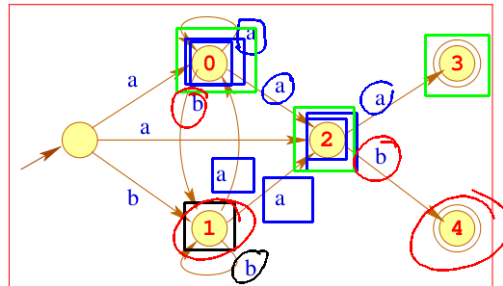
... for example:



47 / 150

Powerset Construction

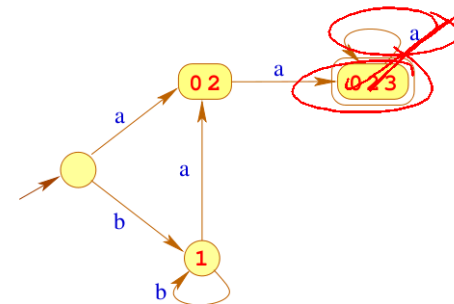
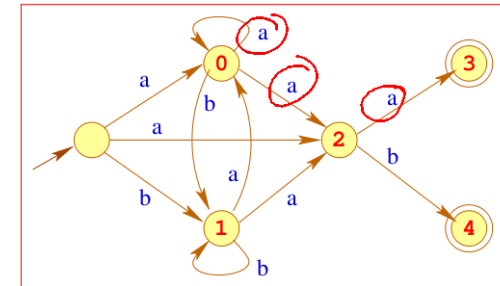
... for example:



47 / 150

Powerset Construction

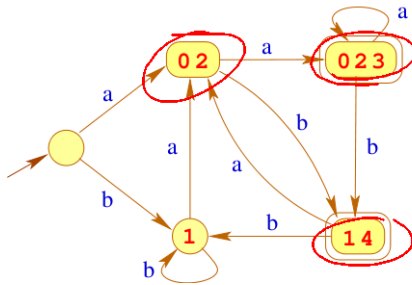
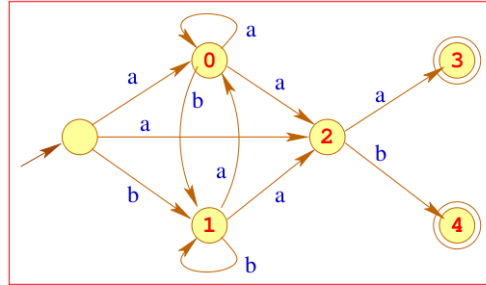
... for example:



47 / 150

Powerset Construction

... for example:



47 / 150

Powerset Construction

Theorem:

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

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

48 / 150

Powerset Construction

Theorem:

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

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

Construction:

States: Powersets of Q ;

Start state: I ;

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

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

48 / 150

Powerset Construction

Theorem:

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

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

48 / 150

Powerset Construction

Theorem:

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

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



48 / 150

Powerset Construction

Bummer!

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 $Q_{\mathcal{P}}$
- Even though, the resulting automaton can become enormously **huge**
... which is (sort of) not happening in **practice**

49 / 150

Powerset Construction

Bummer!

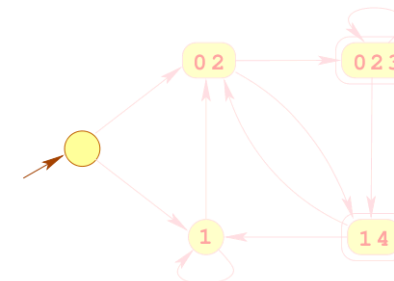
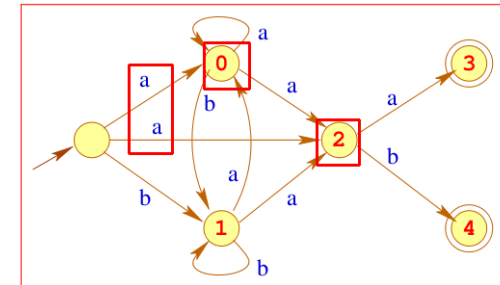
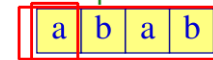
There are exponentially many powersets of Q

- **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 $Q_{\mathcal{P}}$
- Even though, 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**

49 / 150

Powerset Construction

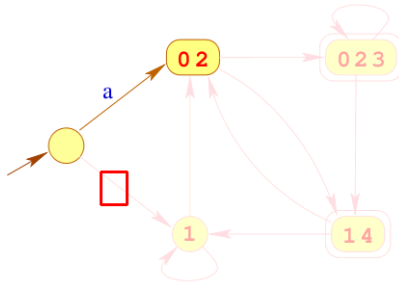
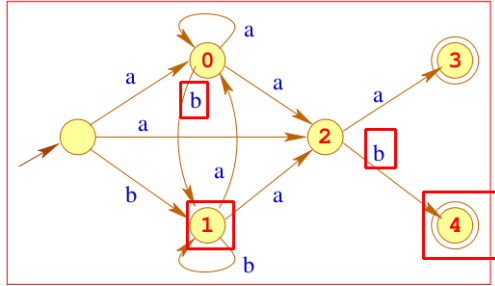
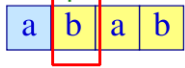
... for example:



50 / 150

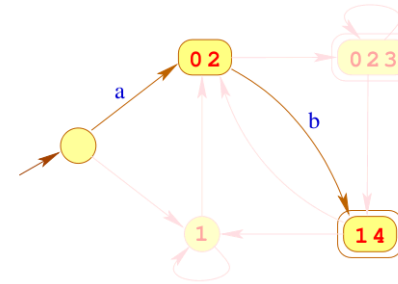
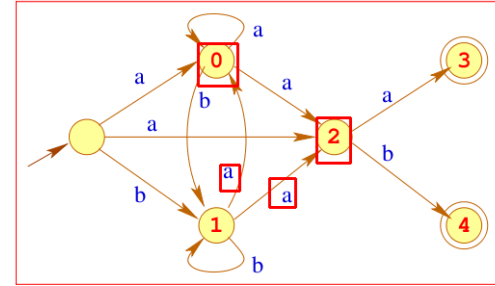
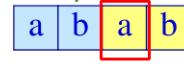
Powerset Construction

... for example:



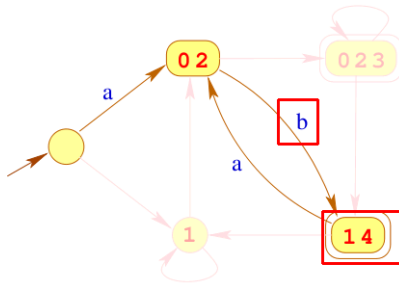
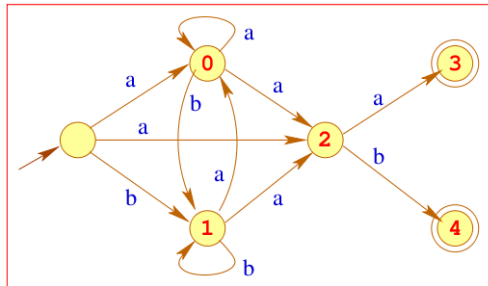
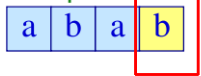
Powerset Construction

... for example:



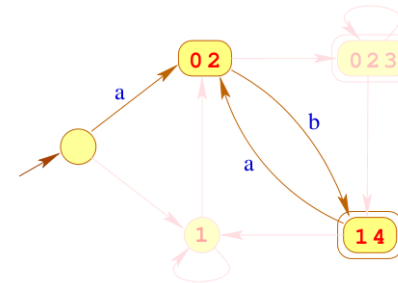
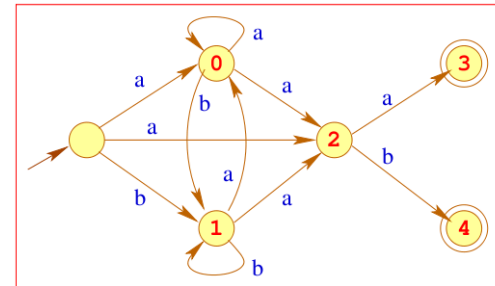
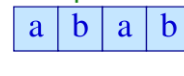
Powerset Construction

... for example:



Powerset Construction

... for example:



Remarks:

- For an input sequence of length n , maximally $\mathcal{O}(n)$ 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.

51 / 150

Remarks:

- For an input sequence of length n , maximally $\mathcal{O}(n)$ 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

$$\mathcal{L}(A) = [e]$$

51 / 150

Lexical Analysis

Chapter 5: Scanner design

52 / 150

Scanner design

Input (simplified):

a set of rules:

e_1	{ action ₁ }
e_2	{ action ₂ }
...	
e_k	{ action _k }

Special information

53 / 150

Scanner design

Input (simplified): a set of rules:

$$\begin{array}{ll} e_1 & \{ \text{action}_1 \} \\ e_2 & \{ \text{action}_2 \} \\ \dots & \\ e_k & \{ \text{action}_k \} \end{array}$$

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 $w \in \llbracket e_i \rrbracket$;
- ... executing action_i for w .

53 / 150

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:

$$\begin{aligned} F_1 &= \{ q \in F \mid q \cap \text{last}[e_1] \neq \emptyset \} \\ F_2 &= \{ q \in (F \setminus F_1) \mid q \cap \text{last}[e_2] \neq \emptyset \} \\ &\dots \\ F_k &= \{ q \in (F \setminus (F_1 \cup \dots \cup F_{k-1})) \mid q \cap \text{last}[e_k] \neq \emptyset \} \end{aligned}$$
- For input w we find: $\delta^*(q_0, w) \in F_i$ iff the scanner must execute action_i for w

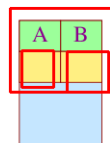
54 / 150

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 ") ;

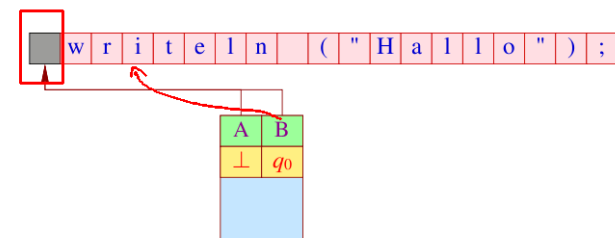


55 / 150

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.



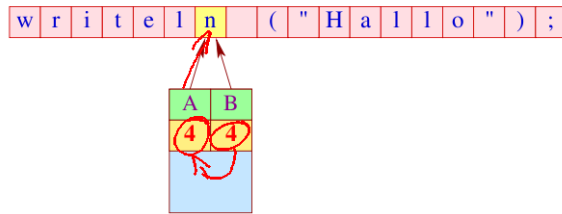
55 / 150

Implementation:

Idea (cont'd):

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

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



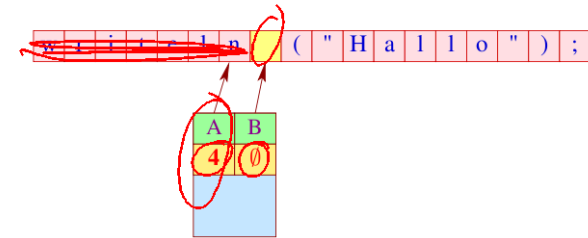
56 / 150

Implementation:

Idea (cont'd):

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

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



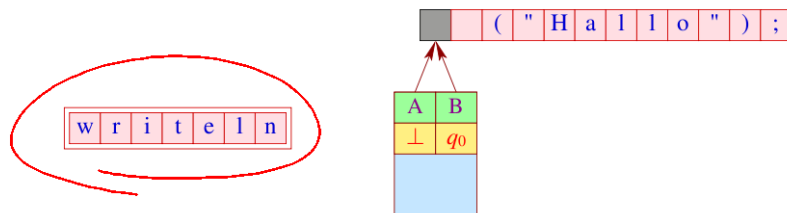
56 / 150

Implementation:

Idea (cont'd):

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

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



56 / 150

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

/ class = +;*

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

57 / 150

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 { }
            }
```

58 / 150

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.

59 / 150

Topic:

Syntactic Analysis

Syntactic Analysis



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

60 / 150

61 / 150

Syntactic Analysis



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

61 / 150

Discussion:

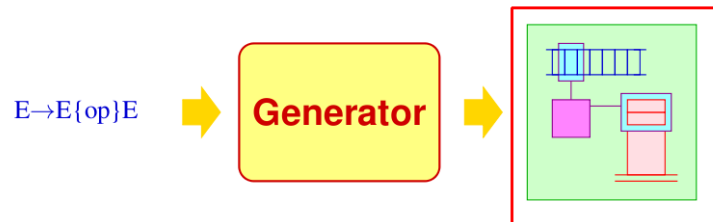
In general, parsers are not developed by hand, but **generated** from a specification:



62 / 150

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

62 / 150

Syntactic Analysis

Chapter 1: Basics of contextfree Grammars

63 / 150

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:

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



Noam Chomsky John Backus

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

64 / 150

Conventions

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

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

65 / 150

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:

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned}$$

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

65 / 150

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:

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned}$$

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: $a, b, c, \dots, \text{int, name}, \dots;$

65 / 150

... further examples:

```

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

66 / 150

... further examples:

```

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

Further 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 $j \geq 0$).

66 / 150

further 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

67 / 150

further 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

67 / 150

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: $\underline{E} \rightarrow \underline{E} + \underline{T}$

68 / 150

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: $\underline{E} \rightarrow \underline{E} + \underline{T}$
 $\rightarrow \underline{T} + \underline{T}$

68 / 150

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: $\underline{E} \rightarrow \underline{E} + \underline{T}$
 $\rightarrow \underline{T} + \underline{T}$
 $\rightarrow \underline{T} * \underline{F} + \underline{T}$
 $\rightarrow \underline{T} * \text{int} + \underline{T}$

68 / 150

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: $\underline{E} \rightarrow \underline{E} + \underline{T}$
 $\rightarrow \underline{T} + \underline{T}$
 $\rightarrow \underline{T} * \underline{F} + \underline{T}$
 $\rightarrow \underline{T} * \text{int} + \underline{T}$
 $\rightarrow \underline{F} * \text{int} + \underline{T}$
 $\rightarrow \text{name} * \text{int} + \underline{T}$
 $\rightarrow \text{name} * \text{int} + \underline{F}$
 $\rightarrow \text{name} * \text{int} + \text{int}$

Definition

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

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

68 / 150

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

A derivation \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^*

68 / 150

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:

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

69 / 150

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:

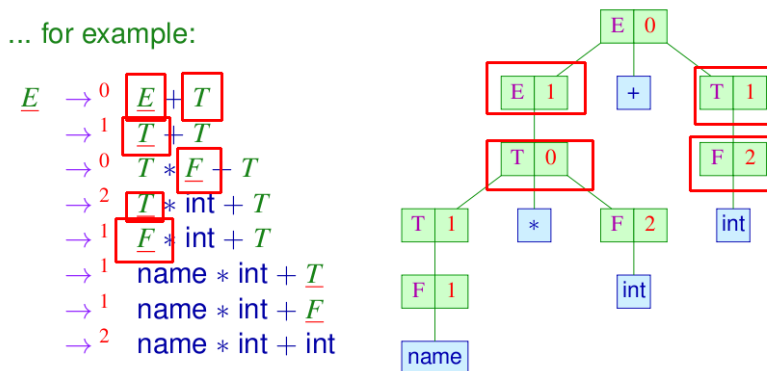
The order, in which disjunct fragments are rewritten is not relevant.

69 / 150

Derivation tree

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

... for example:



A derivation tree for $A \in N$:

- inner nodes:** rule applications
- root:** rule application for A
- leaves:** terminals or ϵ

70 / 150

Special Derivations

Attention:

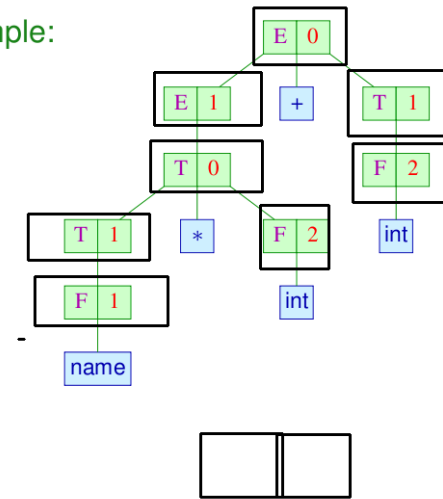
In contrast to arbitrary derivations, we find special ones, always rewriting the **leftmost** (or rather **rightmost**) occurrence of a nonterminal.

- 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

71 / 150

Special Derivations

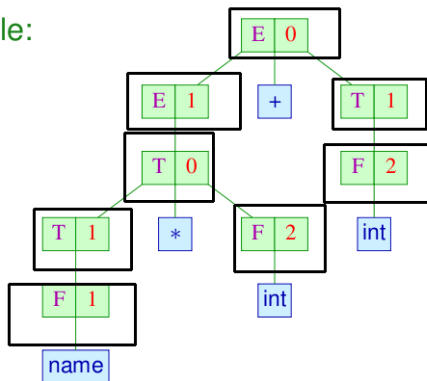
... for example:



72 / 150

Special Derivations

... for example:



Leftmost derivation:

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

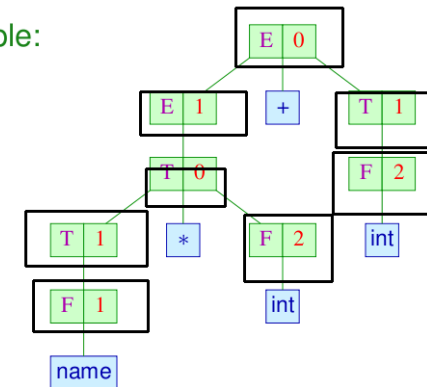
Rightmost derivation:

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

72 / 150

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:

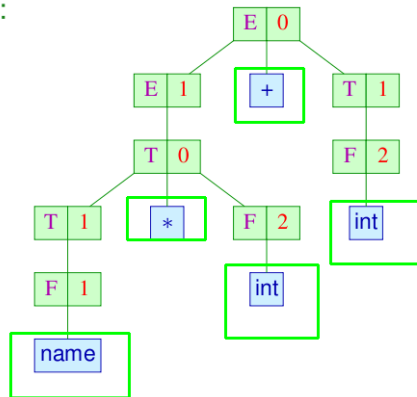
$(F, 1) (T, 1) (F, 2) (T, 0) (E, 1) (F, 2) (T, 1) (E, 0)$

72 / 150

Unique grammars

The concatenation of leaves of a derivation tree t are often called $\text{yield}(t)$.

... for example:



gives rise to the concatenation:

`name * int + int .`

73 / 150

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

74 / 150

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.

75 / 150

Syntactic Analysis

Chapter 2: Basics of pushdown automata

76 / 150