

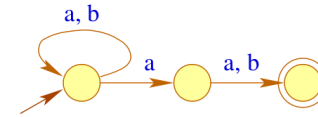
Title: Petter: Compilerbau (02.05.2019)

Date: Thu May 02 14:16:26 CEST 2019

Duration: 92:57 min

Pages: 31

The expected outcome:



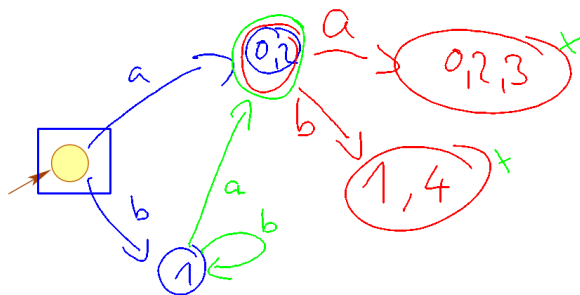
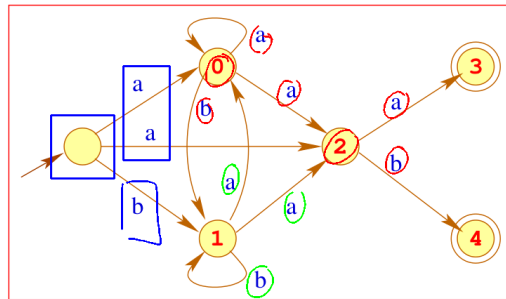
Remarks:

- 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

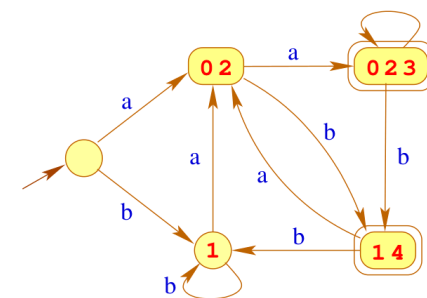
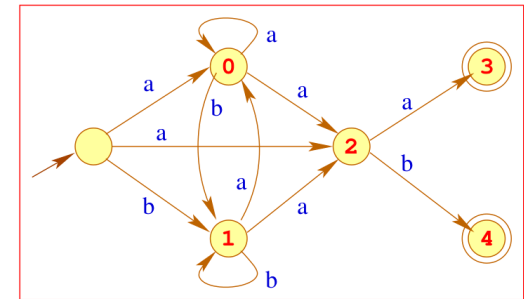
Powerset Construction

... for example:



Powerset Construction

... for example:



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

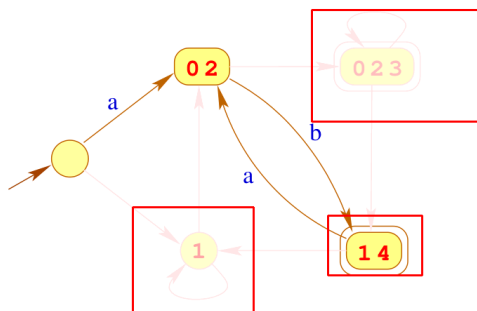
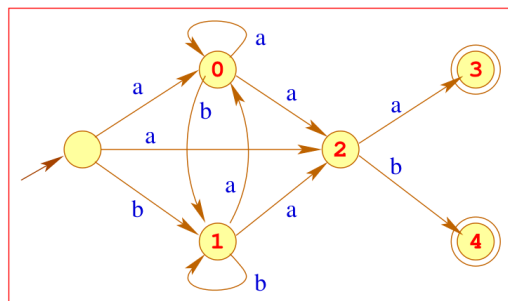
Transitions: $\delta_{\mathcal{P}}(Q', a) = \{q \in Q \mid \exists p \in Q' : (p, a, q) \in \delta\}$.

46 / 287

Powerset Construction

... for example:

a b a b



48 / 287

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

47 / 287

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

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

49 / 287

Scanner design

Input (simplified): a set of rules:

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

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 .

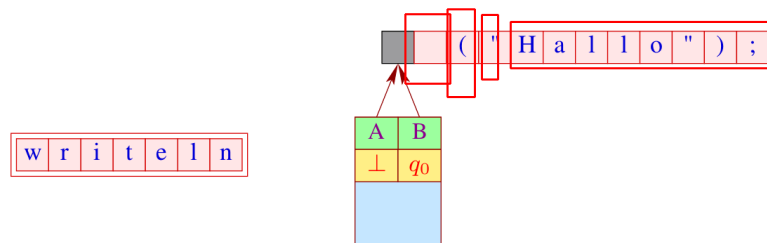
51 / 287

Implementation:

Idea (cont'd):

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

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



54 / 287

Implementation:

Idea:

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

52 / 287

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

55 / 287

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`.

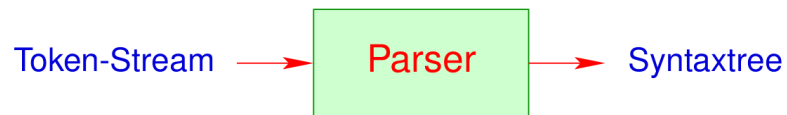
... for example:

```

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

```

Syntactic Analysis



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

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" { }
             }

```

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

Chapter 1: Basics of Contextfree Grammars

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

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...

Conventions

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

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

... 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

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

63 / 287

... a practical example:

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

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 $j \geq 0$).

$\langle \text{stmt} \rangle, \langle \text{stmt} \rangle \rightarrow \langle \text{rexpr} \rangle$

64 / 287

... a practical example:

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

64 / 287

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: $E \rightarrow E + T$

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

66 / 287

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:

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

Definition

The rewriting 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 $A \rightarrow \beta \in P$

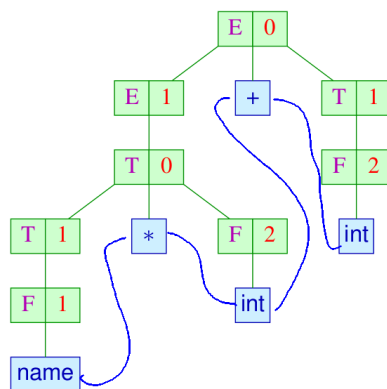
The **reflexive** and **transitive** closure of \rightarrow is denoted as: \rightarrow^*

66 / 287

Derivation Tree

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

... for example:

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


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

68 / 287

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:

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

67 / 287

Special Derivations

Attention:

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

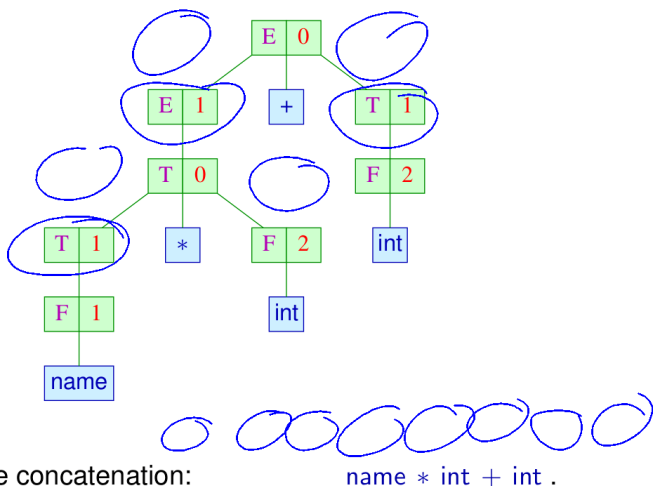
- These are called **leftmost** (or rather **rightmost**) derivations and are denoted with the index L (or R respectively).
- Leftmost (or rightmost) derivations correspond to a left-to-right (or right-to-left) **preorder**-DFS-traversal of the derivation tree.
- **Reverse** rightmost derivations correspond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

69 / 287

Unique Grammars

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

... for example:



71 / 287

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.

73 / 287

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

72 / 287

Syntactic Analysis

Chapter 2:

Basics of Pushdown Automata

74 / 287