

# Topic: Syntactic Analysis

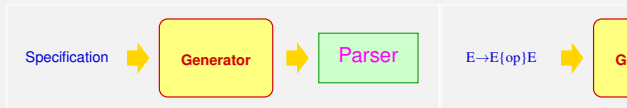
## Syntactic Analysis



- 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:



Specification of the hierarchical structure: contextfree grammars  
Generated implementation: Pushdown automata + X

## Syntactic Analysis

# 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**



## 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:

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

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

... a practical example:

```

S      → (stmt)
(stmt) → (if | (while | (exp);
(if)   → if ( (exp) ) (stmt) else (stmt)
(while) → while ( (exp) ) (stmt)
( exp ) → int | (lexp) | (lexp = (exp) | ...
(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  $j \geq 0$ ).

## 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   |      |     |

|                       |                   |                  |                   |                  |
|-----------------------|-------------------|------------------|-------------------|------------------|
| $E \rightarrow E+E^0$ | $E*E^1$           | $(E)^2$          | name <sup>3</sup> | int <sup>4</sup> |
| $E \rightarrow E+T^0$ | $T^1$             |                  |                   |                  |
| $T \rightarrow T*F^0$ | $F^1$             |                  |                   |                  |
| $F \rightarrow (E)^0$ | name <sup>1</sup> | int <sup>2</sup> |                   |                  |

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 → E + T
  → T + T
  → T * E + T
  → T * int + T
  → E * int + T
  → name * int + T
  → name * int + E
  → name * int + int
  
```

### Definition

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

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

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

## 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:

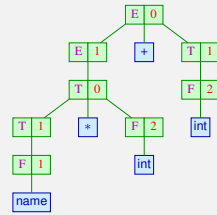
The order, in which disjunct fragments are rewritten is not relevant.

## Derivation Tree

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

... for example:

$E \rightarrow_0 E + T$   
 $\rightarrow_1 T + T$   
 $\rightarrow_0 T * E + T$   
 $\rightarrow_2 T * \text{int} + T$   
 $\rightarrow_1 E * \text{int} + T$   
 $\rightarrow_1 \text{name} * \text{int} + T$   
 $\rightarrow_1 \text{name} * \text{int} + E$   
 $\rightarrow_2 \text{name} * \text{int} + \text{int}$



A **derivation tree** for  $A \in N$ :  
 inner nodes: rule applications  
 root: rule application for  $A$   
 leaves: terminals or  $\epsilon$

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

11/55

## Special Derivations

### Attention:

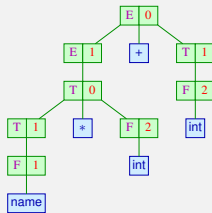
In contrast to arbitrary derivations, we find special ones, always rewriting the **leftmost** (or rather **rightmost**) occurrence of a nonterminal.

- 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

12/55

## Special Derivations

... for example:



**Leftmost derivation:**  
**Rightmost derivation:**  
**Reverse rightmost derivation:**

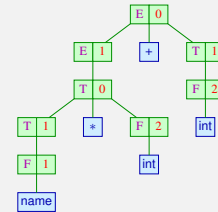
$(E, 0) (E, 1) (T, 0) (T, 1) (F, 1) (F, 2) (T, 1) (F, 2)$   
 $(E, 0) (T, 1) (F, 2) (E, 1) (T, 0) (F, 2) (T, 1) (F, 1)$   
 $(F, 1) (T, 1) (F, 2) (T, 0) (E, 1) (F, 2) (T, 1) (E, 0)$

13/55

## 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 .**

14/55

## Unique Grammars

### Definition:

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

... in our example:

|                         |                 |                |                 |                |
|-------------------------|-----------------|----------------|-----------------|----------------|
| $E \rightarrow E + E^0$ | $E * E^1$       | $(E)^2$        | $\text{name}^3$ | $\text{int}^4$ |
| $E \rightarrow E + T^0$ | $T^1$           |                |                 |                |
| $T \rightarrow T * F^0$ | $F^1$           |                |                 |                |
| $F \rightarrow (E)^0$   | $\text{name}^1$ | $\text{int}^2$ |                 |                |

The first one is ambiguous, the second one is unique

15/55

## 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.

16/55

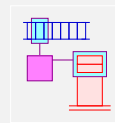
## Syntactic Analysis

### Chapter 2: Basics of Pushdown Automata

17/55

## 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.

18/55

### 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

19/55

### 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**.

20/55

... 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  |

$(0, aabbb) \vdash (11, aabbb)$   
 $\vdash (111, aabbb)$   
 $\vdash (1111, bbb)$   
 $\vdash (112, bb)$   
 $\vdash (12, b)$   
 $\vdash (2, \epsilon)$

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 on 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.

**Definition: Deterministic Pushdown Automaton**

The pushdown automaton  $M$  is deterministic, if every configuration has maximally one successor configuration.

This is exactly the case if for distinct transitions  $(\gamma_1, x, \gamma_2), (\gamma'_1, x', \gamma'_2) \in \delta$  we can assume:

Is  $\gamma_1$  a suffix of  $\gamma'_1$ , then  $x \neq x' \wedge x \neq \epsilon \neq x'$  is valid.

... for example:

|    |   |    |
|----|---|----|
| 0  | a | 11 |
| 1  | a | 11 |
| 11 | b | 2  |
| 12 | b | 2  |

... this obviously holds

**Pushdown Automata**

**Theorem:**

For each context free grammar  $G = (N, T, P, S)$  a pushdown automaton  $M$  with  $\mathcal{L}(G) = \mathcal{L}(M)$  can be built.



The theorem is so important for us, that we take a look at two constructions for automata, motivated by both of the special derivations:

- $M_G^L$  to build **Leftmost derivations**
- $M_G^R$  to build **reverse Rightmost derivations**

Syntactic Analysis

Chapter 3:  
 Top-down Parsing

**Item Pushdown Automaton**

Construction: Item Pushdown Automaton  $M_G^L$

- Reconstruct a **Leftmost derivation**.
- Expand nonterminals using a rule.
- Verify successively, that the chosen rule matches the input.

⇒ The states are now **Items** (= rules with a bullet):

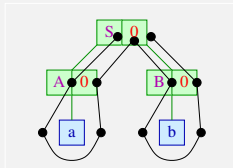
$$[A \rightarrow \alpha \bullet \beta], \quad A \rightarrow \alpha \beta \in P$$

The bullet marks the spot, how far the rule is already processed

**Item Pushdown Automaton – Example**

Our example:

$$S \rightarrow AB^0 \quad A \rightarrow a^0 \quad B \rightarrow b^0$$



**Item Pushdown Automaton – Example**

We add another rule  $S' \rightarrow S \S$  for initialising the construction:

Start state:  $[S' \rightarrow \bullet S \S]$   
 End state:  $[S' \rightarrow S \bullet \S]$   
 Transition relations:

|                                 |            |                                 |                              |
|---------------------------------|------------|---------------------------------|------------------------------|
| $[S' \rightarrow \bullet S \S]$ | $\epsilon$ | $[S' \rightarrow \bullet S \S]$ | $[S \rightarrow \bullet AB]$ |
| $[S \rightarrow \bullet AB]$    | $\epsilon$ | $[S \rightarrow \bullet AB]$    | $[A \rightarrow \bullet a]$  |
| $[A \rightarrow \bullet a]$     | $a$        | $[A \rightarrow a \bullet]$     |                              |
| $[S \rightarrow \bullet AB]$    | $\epsilon$ | $[S \rightarrow A \bullet B]$   |                              |
| $[S \rightarrow A \bullet B]$   | $\epsilon$ | $[S \rightarrow A \bullet B]$   | $[B \rightarrow \bullet b]$  |
| $[B \rightarrow \bullet b]$     | $b$        | $[B \rightarrow b \bullet]$     |                              |
| $[S \rightarrow A \bullet B]$   | $\epsilon$ | $[S \rightarrow AB \bullet]$    |                              |
| $[S' \rightarrow \bullet S \S]$ | $\epsilon$ | $[S \rightarrow AB \bullet]$    |                              |
|                                 |            | $[S' \rightarrow S \bullet \S]$ |                              |

**Item Pushdown Automaton**

The item pushdown automaton  $M_G^L$  has three kinds of transitions:

- Expansions:**  $([A \rightarrow \alpha \bullet B \beta], \epsilon, [A \rightarrow \alpha \bullet B \beta] [B \rightarrow \bullet \gamma])$  for  $A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P$
- Shifts:**  $([A \rightarrow \alpha \bullet a \beta], \alpha, [A \rightarrow \alpha a \bullet \beta])$  for  $A \rightarrow \alpha a \beta \in P$
- Reduces:**  $([A \rightarrow \alpha \bullet B \beta] [B \rightarrow \bullet \gamma], \epsilon, [A \rightarrow \alpha B \bullet \beta])$  for  $A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P$

Items of the form:  $[A \rightarrow \alpha \bullet]$  are also called **complete**  
 The item pushdown automaton shifts the bullet around the derivation tree ...

**Item Pushdown Automaton**

**Discussion:**

- The **expansions** of a computation form a **leftmost derivation**
- Unfortunately, the expansions are chosen **nondeterministically**
- For proving correctness of the construction, we show that for every Item  $[A \rightarrow \alpha \bullet B \beta]$  the following holds:

$$([A \rightarrow \alpha \bullet B \beta], w) \vdash^* ([A \rightarrow \alpha B \bullet \beta], \epsilon) \quad \text{iff} \quad B \rightarrow^* w$$

- **LL-Parsing** is based on the item pushdown automaton and tries to make the expansions deterministic ...

### Item Pushdown Automaton

Example:  $S' \rightarrow S S \mid S \rightarrow \epsilon \mid a S b$

The transitions of the according Item Pushdown Automaton:

|   |   |            |   |
|---|---|------------|---|
| 0 | $[S' \rightarrow \bullet S \$]$                               | $\epsilon$ | $[S' \rightarrow \bullet S \$] [S \rightarrow \bullet]$       |
| 1 | $[S' \rightarrow \bullet S \$]$                               | $\epsilon$ | $[S' \rightarrow \bullet S \$] [S \rightarrow \bullet a S b]$ |
| 2 | $[S \rightarrow \bullet a S b]$                               | $a$        | $[S \rightarrow a \bullet S b]$                               |
| 3 | $[S \rightarrow a \bullet S b]$                               | $\epsilon$ | $[S \rightarrow a \bullet S b] [S \rightarrow \bullet]$       |
| 4 | $[S \rightarrow a \bullet S b]$                               | $\epsilon$ | $[S \rightarrow a \bullet S b] [S \rightarrow \bullet a S b]$ |
| 5 | $[S \rightarrow a \bullet S b] [S \rightarrow \bullet]$       | $\epsilon$ | $[S \rightarrow a S \bullet b]$                               |
| 6 | $[S \rightarrow a \bullet S b] [S \rightarrow a S b \bullet]$ | $\epsilon$ | $[S \rightarrow a S \bullet b]$                               |
| 7 | $[S \rightarrow a S \bullet b]$                               | $b$        | $[S \rightarrow a S b \bullet]$                               |
| 8 | $[S' \rightarrow \bullet S \$] [S \rightarrow \bullet]$       | $\epsilon$ | $[S' \rightarrow S \bullet \$]$                               |
| 9 | $[S' \rightarrow \bullet S \$] [S \rightarrow a S b \bullet]$ | $\epsilon$ | $[S' \rightarrow S \bullet \$]$                               |

Conflicts arise between the transitions (0,1) and (3,4), resp..

### Topdown Parsing

#### Problem:

Conflicts between the transitions prohibit an implementation of the item pushdown automaton as deterministic pushdown automaton.

#### Idea 1: GLL Parsing

For each conflict, we create a virtual copy of the complete configuration and continue computing in parallel.

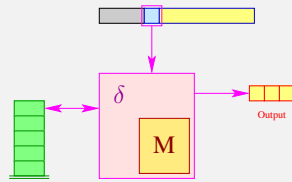
#### Idea 2: Recursive Descent & Backtracking

Depth-first search for an appropriate derivation.

#### Idea 3: Recursive Descent & Lookahead

Conflicts are resolved by considering a lookup of the next input symbols.

### Structure of the LL(1)-Parser:



- The parser accesses a frame of length 1 of the input;
- it corresponds to an item pushdown automaton, essentially;
- table  $M[q, w]$  contains the rule of choice.

### Topdown Parsing

#### Idea:

- Emanate from the item pushdown automaton
- Consider the next input symbol to determine the appropriate rule for the next expansion
- A grammar is called LL(1) if a unique choice is always possible



#### Definition:

A reduced grammar is called LL(1), if for each two distinct rules  $A \rightarrow \alpha, A \rightarrow \alpha' \in P$  and each derivation  $S \xrightarrow{*} u A \beta$  with  $u \in T^*$  the following is valid:

$$\text{First}_1(\alpha \beta) \cap \text{First}_1(\alpha' \beta) = \emptyset$$

### Topdown Parsing

#### Example 1:

$S \rightarrow \text{if } ( E ) S \text{ else } S \mid \text{while } ( E ) S \mid E ;$   
 $E \rightarrow \text{id}$

#### Example 2:

$S \rightarrow \text{if } ( E ) S \text{ else } S \mid \text{if } ( E ) S \mid \text{while } ( E ) S \mid E ;$   
 $E \rightarrow \text{id}$

is LL(1), since  $\text{First}_1(E) = \{\text{id}\}$

... is not LL(k) for any  $k > 0$ .

### Lookahead Sets

#### Definition: First<sub>1</sub>-Sets

For a set  $L \subseteq T^*$  we define:

$$\text{First}_1(L) = \{\epsilon \mid \epsilon \in L\} \cup \{u \in T \mid \exists v \in T^* : uv \in L\}$$

Example:  $S \rightarrow \epsilon \mid a S b$

| $\text{First}_1(\{S\})$ |
|-------------------------|
| $\epsilon$              |
| $a b$                   |
| $a a b b$               |
| $a a a b b b$           |
| ...                     |

≡ the yield's prefix of length 1

### Lookahead Sets

#### Arithmetics:

$\text{First}_1(\_)$  is distributive with union and concatenation:

$$\begin{aligned} \text{First}_1(\emptyset) &= \emptyset \\ \text{First}_1(L_1 \cup L_2) &= \text{First}_1(L_1) \cup \text{First}_1(L_2) \\ \text{First}_1(L_1 \cdot L_2) &= \text{First}_1(\text{First}_1(L_1) \cdot \text{First}_1(L_2)) \\ &:= \text{First}_1(L_1) \odot_1 \text{First}_1(L_2) \end{aligned}$$

$\odot_1$  being 1-concatenation

#### Definition: 1-concatenation

Let  $L_1, L_2 \subseteq T \cup \{\epsilon\}$  with  $L_1 \neq \emptyset \neq L_2$ . Then:

$$L_1 \odot_1 L_2 = \begin{cases} L_1 & \text{if } \epsilon \notin L_1 \\ (L_1 \setminus \{\epsilon\}) \cup L_2 & \text{otherwise} \end{cases}$$

If all rules of  $G$  are productive, then all sets  $\text{First}_1(A)$  are non-empty.

### Lookahead Sets

For  $\alpha \in (N \cup T)^*$  we are interested in the set:

$$\text{First}_1(\alpha) = \text{First}_1(\{w \in T^* \mid \alpha \rightarrow^* w\})$$

#### Idea: Treat $\epsilon$ separately:

- Let  $\text{empty}(X) = \text{true}$  iff  $X \rightarrow^* \epsilon$ .
- $F_\epsilon(X_1 \dots X_m) = \bigcup_{i=1}^m F_\epsilon(X_i)$  if  $\neg \text{empty}(X_j) \wedge \bigwedge_{i=1}^{j-1} \text{empty}(X_i)$

We characterize the  $\epsilon$ -free  $\text{First}_1$ -sets with an inequality system:

$$\begin{aligned} F_\epsilon(a) &= \{a\} & \text{if } a \in T \\ F_\epsilon(A) &\supseteq F_\epsilon(X_j) & \text{if } A \rightarrow X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1}) \end{aligned}$$

### Lookahead Sets

for example...

$E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow ( E ) \mid \text{name} \mid \text{int}$

with  $\text{empty}(E) = \text{empty}(T) = \text{empty}(F) = \text{false}$

... we obtain:

$$\begin{aligned} F_\epsilon(S') &\supseteq F_\epsilon(E) & F_\epsilon(E) &\supseteq F_\epsilon(E) \\ F_\epsilon(E) &\supseteq F_\epsilon(T) & F_\epsilon(T) &\supseteq F_\epsilon(T) \\ F_\epsilon(T) &\supseteq F_\epsilon(F) & F_\epsilon(F) &\supseteq \{ (, \text{name}, \text{int}) \} \end{aligned}$$

### Fast Computation of Lookahead Sets

#### Observation:

- The form of each inequality of these systems is:

$$x \supseteq y \quad \text{resp.} \quad x \supseteq d$$

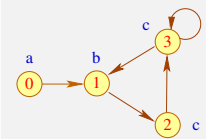
for variables  $x, y$  and  $d \in D$ .

- Such systems are called pure unification problems
- Such problems can be solved in linear space/time.

for example:

$$D = 2^{\{a,b,c\}}$$

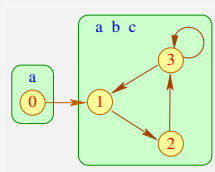
$$\begin{aligned} x_0 &\supseteq \{a\} \\ x_1 &\supseteq \{b\} & x_1 &\supseteq x_0 & x_1 &\supseteq x_3 \\ x_2 &\supseteq \{c\} & x_2 &\supseteq x_1 & & \\ x_3 &\supseteq \{c\} & x_3 &\supseteq x_2 & x_3 &\supseteq x_3 \end{aligned}$$



### Fast Computation of Lookahead Sets



Frank DeRemer & Tom Pennello



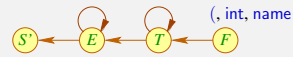
Proceeding:

- Create the **Variable Dependency Graph** for the inequality system.
- Within a **Strongly Connected Component** ( $\rightarrow$  Tarjan) all variables have the same value
- Is there no ingoing edge for an SCC, its value is computed via the smallest upper bound of all values within the SCC
- In case of ingoing edges, their values are also to be considered for the upper bound

### Fast Computation of Lookahead Sets

... for our example grammar:

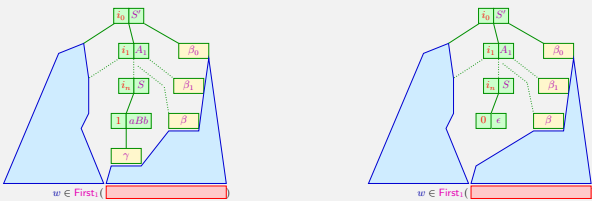
First<sub>1</sub> :



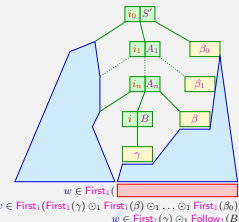
### Item Pushdown Automaton as LL(1)-Parser

context is relevant too:  $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

|                            |    |   |   |
|----------------------------|----|---|---|
| First <sub>1</sub> (input) | \$ | a | b |
| S                          | ?  | ? | ? |



### Item Pushdown Automaton as LL(1)-Parser



Inequality system for Follow<sub>1</sub>(B) = First<sub>1</sub>(β) ⊇<sub>1</sub> ... ⊇<sub>1</sub> First<sub>1</sub>(β<sub>0</sub>)

- Follow<sub>1</sub>(S) ⊇ { \$ }
- Follow<sub>1</sub>(B) ⊇ F<sub>ε</sub>(X<sub>j</sub>) if  $A \rightarrow \alpha B X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1})$
- Follow<sub>1</sub>(B) ⊇ Follow<sub>1</sub>(A) if  $A \rightarrow \alpha B X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_m)$

### Item Pushdown Automaton as LL(1)-Parser

Is  $G$  an LL(1)-grammar, we can index a lookahead-table with items and nonterminals:

#### LL(1)-Lookahead Table

We set  $M[B, w] = i$  with  $B \rightarrow \gamma^i$  if  $w \in \text{First}_1(\gamma) \cap_1 \text{Follow}_1(B)$

... for example:  $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

First<sub>1</sub>(S) = {ε, a}    Follow<sub>1</sub>(S) = {b, \$}

- S-rule 0 : First<sub>1</sub>(ε) ⊇<sub>1</sub> Follow<sub>1</sub>(S) = {b, \$}
- S-rule 1 : First<sub>1</sub>(aSb) ⊇<sub>1</sub> Follow<sub>1</sub>(S) = {a}

|   |    |   |   |
|---|----|---|---|
|   | \$ | a | b |
| S | 0  | 1 | 0 |

### Item Pushdown Automaton as LL(1)-Parser

For example:  $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

The transitions of the according Item Pushdown Automaton:

|   |                                 |                                 |                                 |                                 |
|---|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| 0 | $[S' \rightarrow \bullet S \$]$ | ε                               | $[S' \rightarrow \bullet S \$]$ | $[S \rightarrow \bullet]$       |
| 1 | $[S' \rightarrow \bullet S \$]$ | ε                               | $[S' \rightarrow \bullet S \$]$ | $[S \rightarrow \bullet a S b]$ |
| 2 | $[S \rightarrow \bullet a S b]$ | a                               | $[S \rightarrow \bullet a S b]$ |                                 |
| 3 | $[S \rightarrow \bullet a S b]$ | ε                               | $[S \rightarrow \bullet a S b]$ | $[S \rightarrow \bullet]$       |
| 4 | $[S \rightarrow \bullet a S b]$ | ε                               | $[S \rightarrow \bullet a S b]$ | $[S \rightarrow \bullet a S b]$ |
| 5 | $[S \rightarrow \bullet a S b]$ | $[S \rightarrow \bullet]$       | ε                               | $[S \rightarrow \bullet a S b]$ |
| 6 | $[S \rightarrow \bullet a S b]$ | $[S \rightarrow \bullet a S b]$ | ε                               | $[S \rightarrow \bullet a S b]$ |
| 7 | $[S \rightarrow \bullet a S b]$ | b                               | $[S \rightarrow \bullet a S b]$ |                                 |
| 8 | $[S' \rightarrow \bullet S \$]$ | $[S \rightarrow \bullet]$       | ε                               | $[S' \rightarrow \bullet S \$]$ |
| 9 | $[S' \rightarrow \bullet S \$]$ | $[S \rightarrow \bullet a S b]$ | ε                               | $[S' \rightarrow \bullet S \$]$ |

Lookahead table:

|   |    |   |   |
|---|----|---|---|
|   | \$ | a | b |
| S | 0  | 1 | 0 |

### Left Recursion

#### Attention:

Many grammars are not LL(k)!

A reason for that is:

#### Definition

Grammar  $G$  is called **left-recursive**, if

$$A \rightarrow^+ A \beta \quad \text{for an } A \in N, \beta \in (T \cup N)^*$$

Example:

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

... is left-recursive

### Left Recursion

#### Theorem:

Let a grammar  $G$  be reduced and left-recursive, then  $G$  is not LL(k) for any  $k$ .

#### Proof:

Let wlog.  $A \rightarrow A \beta \mid \alpha \in P$  and  $A$  be reachable from  $S$

Assumption:  $G$  is LL(k)

$$\Rightarrow \text{First}_k(\alpha \beta^n \gamma) \cap \text{First}_k(\alpha \beta^{n+1} \gamma) = \emptyset$$

- Case 1:  $\beta \rightarrow^+ \epsilon$  — Contradiction!!!
- Case 2:  $\beta \rightarrow^+ w \neq \epsilon \implies \text{First}_k(\alpha w^k \gamma) \cap \text{First}_k(\alpha w^{k+1} \gamma) \neq \emptyset$

### Right-Regular Context-Free Parsing

Recurring scheme in programming languages: Lists of sth...

$$S \rightarrow b \mid S a b$$

Alternative idea: **Regular Expressions**

$$S \rightarrow ( b a )^* b$$

#### Definition: Right-Regular Context-Free Grammar

A **right-regular context-free grammar (RR-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 rules with **regular expressions of symbols** as rhs,
- $S \in N$  the start symbol

Example: Arithmetic Expressions

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow T \mid ( + T )^* \\ T &\rightarrow F \mid ( * F )^* \\ F &\rightarrow ( E ) \mid \text{name} \mid \text{int} \end{aligned}$$

### Idea 1: Rewrite the rules from $G$ to $\langle G \rangle$ :

$$\begin{aligned} A &\rightarrow \langle \alpha \rangle & \text{if } A \rightarrow \alpha \in P \\ \langle \alpha \rangle &\rightarrow \alpha & \text{if } \alpha \in N \cup T \\ \langle \epsilon \rangle &\rightarrow \epsilon & \\ \langle \alpha^* \rangle &\rightarrow \epsilon \mid \langle \alpha \rangle \langle \alpha^* \rangle & \text{if } \alpha \in \text{RegExt}_{T,N} \\ \langle \alpha_1 \dots \alpha_n \rangle &\rightarrow \langle \alpha_1 \rangle \dots \langle \alpha_n \rangle & \text{if } \alpha_i \in \text{RegExt}_{T,N} \\ \langle \alpha_1 \mid \dots \mid \alpha_n \rangle &\rightarrow \langle \alpha_1 \rangle \mid \dots \mid \langle \alpha_n \rangle & \text{if } \alpha_i \in \text{RegExt}_{T,N} \end{aligned}$$

... and generate the according LL(k)-Parser  $M_{\langle G \rangle}^T$

Example: Arithmetic Expressions cont'd

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow T \mid ( + T )^* ( T ( + T )^* ) \\ T &\rightarrow F \mid ( * F )^* ( F ( * F )^* ) \\ F &\rightarrow ( E ) \mid \text{name} \mid \text{int} \\ \langle T ( + T )^* \rangle &\rightarrow T \langle ( + T )^* \rangle \\ \langle ( + T )^* \rangle &\rightarrow \epsilon \mid \langle ( + T ) \rangle \langle ( + T )^* \rangle \\ \langle + T \rangle &\rightarrow + T \\ \langle F ( * F )^* \rangle &\rightarrow F \langle ( * F )^* \rangle \\ \langle ( * F )^* \rangle &\rightarrow \epsilon \mid \langle * F \rangle \langle ( * F )^* \rangle \\ \langle * F \rangle &\rightarrow * F \end{aligned}$$



### Definition:

An *RR-CFG*  $G$  is called *RLL(1)*, if the corresponding CFG  $(G)$  is an *LL(1)* grammar.

### Discussion

- directly yields the table driven parser  $M_{(G)}^L$  for *RLL(1)* grammars
- however: mapping regular expressions to recursive productions unnecessarily strains the stack  
→ instead directly construct automaton in the style of Berry-Sethi

### Idea 2: Recursive Descent RLL Parsers:

*Recursive descent* RLL(1)-parsers are an alternative to table-driven parsers; apart from the usual function `scan()`, we generate a program frame with the lookahead function `expect()` and the main parsing method `parse()`:

```

int next;
void expect(Set E){
  if ({ε,next} ∩ E = ∅){
    cerr << "Expected" << E << "found" << next;
    exit(0);
  }
  return;
}
void parse(){
  next = scan();
  expect(First1(S));
  S();
  expect({EOF});
}

```

### Idea 2: Recursive Descent RLL Parsers:

For each  $A \rightarrow \alpha \in P$ , we introduce:

```

void A(){
  generate(α)
}

```

with the meta-program *generate* being defined by structural decomposition of  $\alpha$ :

```

generate(r1 ... rk) = generate(r1)
                    expect(First1(r2));
                    generate(r2)
                    :
                    expect(First1(rk));
                    generate(rk)
generate(ε)         = ;
generate(a)         = next = scan();
generate(A)         = A();

```

### Idea 2: Recursive Descent RLL Parsers:

```

generate(r*)       = while (next ∈ Fε(r)) {
                    generate(r)
                    }
generate(r1 | ... | rk) = switch(next) {
                    labels(First1(r1)) generate(r1) break;
                    :
                    labels(First1(rk)) generate(rk) break;
                    }
labels({α1, ..., αm}) = label(α1); ... label(αm);
label(α)               = case α
label(ε)               = default

```

### Topdown-Parsing

#### Discussion

- A practical implementation of an *RLL(1)*-parser via *recursive descent* is a straight-forward idea
- However, **only a subset** of the deterministic contextfree languages can be parsed this way.
- As soon as *First1()* sets are not disjoint any more,
  - **Solution 1:** For many accessibly written grammars, the alternation between right hand sides happens too early. Keeping the common prefixes of right hand sides joined and introducing a new production for the actual diverging sentence forms often helps.
  - **Solution 2:** Introduce *ranked* grammars, and decide conflicting lookahead always in favour of the higher ranked alternative  
→ relation to *LL* parsing not so clear any more  
→ not so clear for “*\**” operator how to decide
  - **Solution 3:** Going from *LL(1)* to *LL(k)*  
The size of the occurring sets is rapidly increasing with larger *k*  
**Unfortunately**, even *LL(k)* parsers are not sufficient to accept all deterministic contextfree languages. (regular lookahead → *LL(\*)*)
- In practical systems, this often motivates the implementation of *k* = 1 only ...