

TDDD55 Compilers and interpreters

TDDE66 Compiler Construction



Formal Languages Part 2

Context Free Grammars

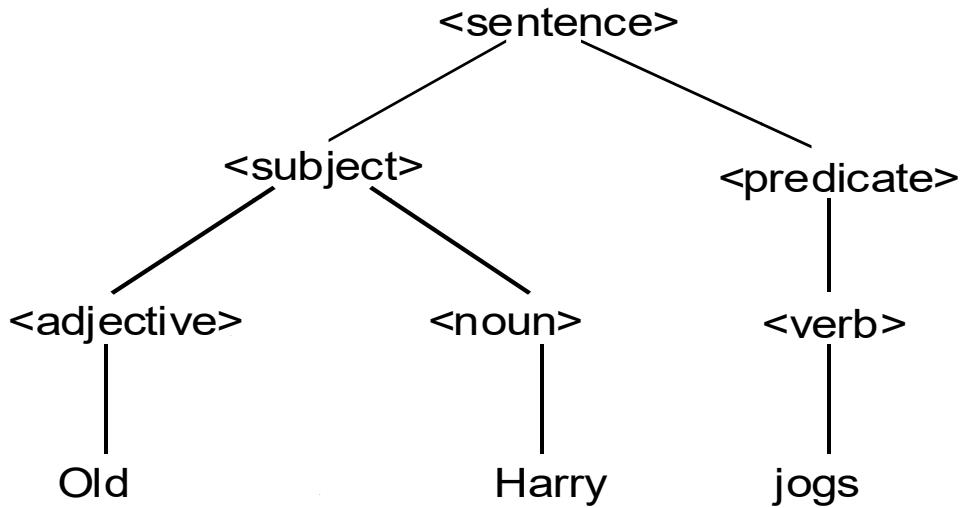
Context-Free Grammars

- Example: an English sentence:

Sentence: Old Harry jogs
 Constituents: subject predicate
 Word Class: adjective noun verb

- A grammar is used to describe the syntax.

BNF (Backus-Naur form) 1960 (meta-language to describe languages):



- <sentence> → <subject><predicate>
- <subject> → <adjective> <noun>
- <predicate> → <verb>
- <adjective> → old | big | strong | ...
- <noun> → Harry | brother | ...
- <verb> → jogs | snores | sleeps | ...

Grammars, cont.

- <sentence> is a *start symbol*.
- *Symbols to the left of “→” are called nonterminals.*
- Symbols not surrounded by “< >” are *terminals*.
- *Each line is a production.*

Symbol	Meaning
< ... >	syntactic classes
→	"consists of", "is" (also “::=”)
	"or"

A Grammar can be used to Produce or Derive Sentences

- Example: $\langle \text{sentence} \rangle \Rightarrow^* \text{Old Harry jogs}$
 - where $\langle \text{sentence} \rangle$ is the start symbol and \Rightarrow^* means derivation in zero or more steps.

Example Derivation:

$\langle \text{sentence} \rangle \Rightarrow \langle \text{subject} \rangle \langle \text{predicate} \rangle$
 $\Rightarrow \langle \text{adjective} \rangle \langle \text{noun} \rangle \langle \text{predicate} \rangle$
 $\Rightarrow \text{Old } \langle \text{noun} \rangle \langle \text{predicate} \rangle$
 $\Rightarrow \text{Old Harry } \langle \text{predicate} \rangle$
 $\Rightarrow \text{Old Harry } \langle \text{verb} \rangle$
 $\Rightarrow \text{Old Harry jogs}$

Definition: CFG (Context-free grammar)

- A CFG (Context-free grammar) is a quadruple (4 parts):

$$G = \langle N, \Sigma, P, S \rangle$$

where

N : Nonterminals.

Σ : terminal Symbols.

P : rules, Productions of the form
 $A \rightarrow a$ where $A \in N$ and
 $a \in (N \cup \Sigma)^*$

S : the Start symbol,
 a nonterminal, $S \in N$.

- (Sometimes $V = N \cup \Sigma$ is used, called the *vocabulary*.)

- Example:

- 1. $\langle \text{number} \rangle \rightarrow \langle \text{no} \rangle$
- 2. $\langle \text{no} \rangle \rightarrow \langle \text{no} \rangle \langle \text{digit} \rangle$
- 3. | $\langle \text{digit} \rangle$
- 4. $\langle \text{digit} \rangle \rightarrow 0|1|2|3|4|5|6|7|8|9$

- $N = \{ \langle \text{number} \rangle, \langle \text{no} \rangle, \langle \text{digit} \rangle \}$
- $\Sigma = \{ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$
- $S = \langle \text{number} \rangle$

Notational Conventions

$\alpha, \beta, \gamma \in V^*$	string of terminals and nonterminals
$A, B, C \in N$	nonterminals
$a, b, c \in \Sigma$	terminal symbols
$u, v, w, x, y, z \in \Sigma^*$	string of terminals

Derivations

□ Derivation

- $\alpha \Rightarrow \beta$ (pronounced “ α derives β ”)
 - Formally: $\gamma A \theta \Rightarrow \gamma \delta \theta$ if we have $A \rightarrow \delta$
 - Example: $\langle \text{number} \rangle \Rightarrow_{\text{rm}}^* \langle \text{no} \rangle \Rightarrow_{\text{rm}}^* \langle \text{no} \rangle \langle \text{digit} \rangle \Rightarrow_{\text{rm}}^* \langle \text{no} \rangle 2 \Rightarrow_{\text{rm}}^* \langle \text{digit} \rangle 2 \Rightarrow_{\text{rm}}^* 12$
 - $\langle \text{number} \rangle \Rightarrow \langle \text{no} \rangle$ direct derivation.
 - $\langle \text{number} \rangle \Rightarrow^* 12$ several derivations (zero or more).
 - $\langle \text{number} \rangle \Rightarrow^+ 12$ several derivations (one or more).
- Given $G = \langle N, \Sigma, P, S \rangle$ the language generated by G can be defined as $L(G)$:
- $$L(G) = \{ w \mid S \Rightarrow^+ w \text{ and } w \in \Sigma^* \}$$

Sentential form, Sentence

□ Sentential form

- A string α is a *sentential form* in G if
- $S \Rightarrow^* \alpha$ and $\alpha \in V^*$ (string of terminals and/or nonterminals)
- Example: <no> <digit> is a sentential form in $G(<\text{number}>)$. $<\text{no}>8$ is another sentential form

□ Sentence

- w is a *sentence* in G if $S \Rightarrow^+ w$ and $w \in \Sigma^*$.
- Example: 12 is a sentence in $G(<\text{number}>)$.

Left and Right Derivations

□ Left derivation

- \Rightarrow_{lm} means that we replace the *leftmost* nonterminal by some appropriate right side.

□ Left sentential form

- A sentential form which is part of a leftmost derivation.

□ Right derivation (canonical derivation)

- \Rightarrow_{rm} means that we replace the *rightmost* nonterminal by some appropriate right side.

□ Right sentential form

- A sentential form which is part of a rightmost derivation.

Rightmost Derivation, Handle

□ Reverse rightmost derivation

- $12 \leq_{rm}^* \langle \text{digit} \rangle 2 \leq_{rm}^* \langle \text{no} \rangle 2 \leq_{rm}^* \langle \text{no} \rangle \langle \text{digit} \rangle \leq_{rm}^* \langle \text{no} \rangle \leq_{rm}^* \langle \text{number} \rangle$

□ Handles

Consist of two parts:

- 1. A production $A \rightarrow \beta$
 - 2. A position
 - If $S \Rightarrow_{rm}^* \alpha A w \Rightarrow_{rm} \alpha \beta w$, the production $A \rightarrow \beta$ together with the position after α is a **handle** of $\alpha \beta w$.
1. $\langle \text{number} \rangle \rightarrow \langle \text{no} \rangle$
 2. $\langle \text{no} \rangle \rightarrow \langle \text{no} \rangle \langle \text{digit} \rangle$
 3. | $\langle \text{digit} \rangle$
 4. $\langle \text{digit} \rangle \rightarrow 0|1|2|3|4|5|6|7|8|9$

□ Example: The handle of $\langle \text{no} \rangle 2$ is the production $\langle \text{digit} \rangle \rightarrow 2$ and the position after $\langle \text{no} \rangle$ because:

- $\langle \text{number} \rangle \Rightarrow_{rm} \langle \text{no} \rangle \Rightarrow_{rm} \langle \text{no} \rangle \langle \text{digit} \rangle \Rightarrow_{rm} \langle \text{no} \rangle 2 \Rightarrow_{rm} \langle \text{digit} \rangle 2 \Rightarrow_{rm} 12$

□ Informally: a handle is what we *reduce* to what and where to get the previous sentential form in a rightmost derivation.

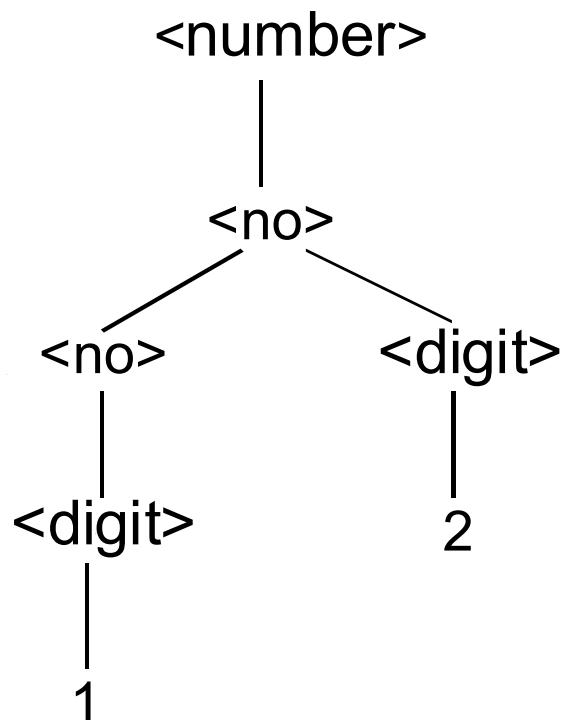
Reduction

□ Reduction of a grammar rule

In reverse right derivation, find a **right side** in some rule according to the grammar in the given right sentential form and **replace** it with the corresponding **left side**, i.e., nonterminal

Parse trees (derivation trees)

- ❑ A parse tree can correspond to several different derivations.



Parse tree for 12

Example Grammar:

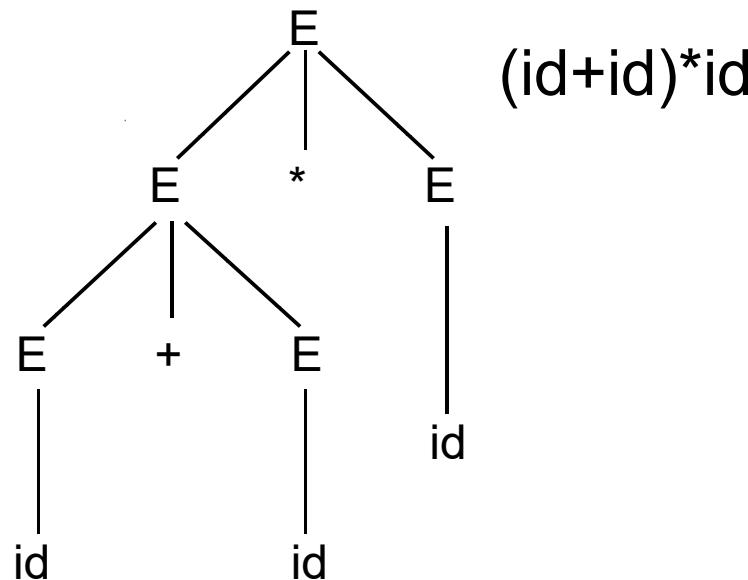
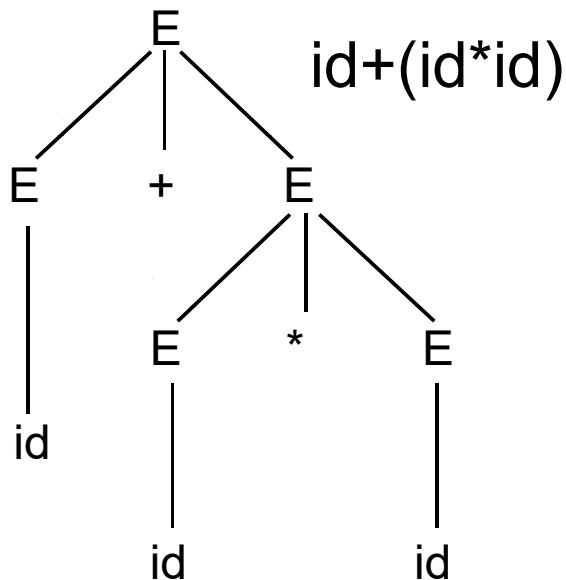
1. $\langle \text{number} \rangle \rightarrow \langle \text{no} \rangle$
2. $\langle \text{no} \rangle \rightarrow \langle \text{no} \rangle \langle \text{digit} \rangle$
3. | $\langle \text{digit} \rangle$
4. $\langle \text{digit} \rangle \rightarrow 0|1|2|3|4|5|6|7|8|9$

Ambiguous Grammars

- A grammar G is *ambiguous* if a sentence in G has several different parse trees.

e.g.
$$\begin{array}{lcl} E & \rightarrow & E + E \\ & | & E * E \\ & | & E \uparrow E \\ & | & id \end{array}$$

- $id + id * id$ has two different parse trees.



Rewriting to Unambiguous Grammar

- Rewrite the grammar to make it unambiguous:
 - +, * are to have the right priority and
 - +, * are to be left associative while
 - ↑ is to be right associative.
- Example: $a+b+c+d$ is interpreted as $(a+b)+c+d$, using the rewritten grammar:

$$\begin{array}{ll} E \rightarrow & E + T \quad (\text{left associative}) \\ | & T \\ T \rightarrow & T * F \quad (\text{left associative}) \\ | & F \\ F \rightarrow & P \uparrow F \quad (\text{right associative}) \\ | & P \\ P \rightarrow & id \end{array}$$

Small Parse Tree Exercise

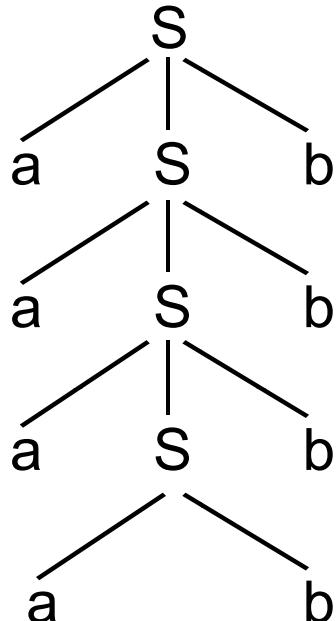
See 00-LectureExercises

Example Palindrome Grammars

- Palindrome: a string that is symmetrical around its center
- Example: The following grammar generates
 $\{ a^n b^n \mid n \geq 1 \}$.

$$S \rightarrow a S b \\ | \\ a b$$

Example
parse tree:



Another example:
Grammar describing binary
palindromes of *odd* lengths ≥ 1 :

$$S \rightarrow 0 S 0 \\ | \\ 1 S 1 \\ | \\ 0 \\ | \\ 1$$

Example derived strings: 0, 1, 000,
010, 111

Binary Palindrome Exercise

See 00-LectureExercises

Example

Excerpt from a Pascal Grammar

```
goal>      → <progdecl> .
<progdecl> → <prog_hedr> ; <block>
<prog_hedr> → program <idname>
                ( <idname_list> )
                | program <idname>
<block>      → <decls> begin
                <stat_list>
                end
<decls>      → <labels> <consts>
                <types> <vars> <procs>
<labels>      → label <label_decl> ;
                |
                ε
<label_decl> → <label_decl> , <labelid>
                |
                <labelid>
<labelid>     → <int>
                |
                <id>
<consts>      → const <const_decls>
                |
                ε
```

```
<const_decls>   → <const_decls> <const_decl_c>
                  |
                  <const_decl_c>
<const_decl_c>  → <const_decl> ;
<const_decl>    → <idname> = <const>
<types>         → type <type_decls>
                  |
                  ε
<type_decls>   → <type_decls> <type_decl_c>
                  |
                  <type_decl_c>
<type_decl_c>  → <type_decl> ;
<type_decl>    → <idname> = <type>
<vars>          → var <var_decls>
                  |
                  ε
<var_decls>    → <var_decls> <var_decl_c>
                  |
                  <var_decl_c>
<var_decl_c>   → <var_decl> ;
<var_decl>     → <id_list> : <type>
<procs>         → <proc_decls>
                  |
                  ε
<proc_decls>   → <proc_decls> <proc>
                  |
                  <proc>
```

Excerpt from a Pascal Grammar, Cont.

```
<proc>      → procedure <phead_c> forward ;
             | procedure <phead_c> <block> ;
             | function <fhead_c> forward ;
             | function <fhead_c> <block> ;

<fhead_c>    → <fhead> ;

<fhead>      → <idname> <params> : <type_id>

<phead_c>    → <phead> ;

<phead>      → <idname> <params>
             | ε

<params>      → ( <param_list> )
             | ε

<param>       → var <par_decl>
             | <par_decl>
             | ε

<par_decl>    → <id_list> : <type_id>

<param_list>  → <param_list> ; <param>
             | <param>

<id_list>     → <id_list> , <id>
             | <id>

...
```