

(NB. Pages 21 - 39 are intended for those who need repeated study in formal languages)

Formal languages

Basic concepts for symbols, strings and languages:

Alphabet

A finite set of symbols.

Example:

- $\Sigma_b = \{0, 1\}$ binary alphabet
- $\Sigma_s = \{A, B, C, \dots, Z, \text{Å}, \text{Ä}, \text{Ö}\}$ Swedish characters
- $\Sigma_r = \{\text{WHILE}, \text{IF}, \text{BEGIN}, \dots\}$ reserved words

String

A finite sequence of symbols from an alphabet.

Example:

- 10011 from Σ_b
- KALLE from Σ_s
- WHILE DO BEGIN from Σ_r

Length of a string

Number of symbols in the string.

Example:

- x arbitrary string, $|x|$ length of the string x
- $|10011| = 5$ according to Σ_b
- $|\text{WHILE}| = 5$ according to Σ_s
- $|\text{WHILE}| = 1$ according to Σ_r

Empty string

The empty string is denoted ϵ , $|\epsilon| = 0$

Concatenation

Two strings x and y are joined together $x \bullet y = xy$

Example:

- $x = \text{AB}, y = \text{CDE}$ produce $x \bullet y = \text{ABCDE}$
- $|xy| = |x| + |y|$
- $xy \neq yx$ (not commutative)
- $\epsilon x = x\epsilon = x$

String exponentiation

- $x^0 = \epsilon$
- $x^1 = x$
- $x^2 = xx$
- $x^n = x \bullet x^{n-1}, n \geq 1$

Substrings: Prefix, suffix.

Example: $x = abc$

Prefix: Substring at the beginning.

Prefix of $x: abc$ (improper as the prefix = x), ab , a , ϵ

Suffix: Substring at the end.

Suffix of $x: abc$ (improper as the suffix = x), bc , c , ϵ

Languages

A finite or infinite set of strings which can be constructed from a special alphabet.

Alternatively: a subset of all the strings which can be constructed from an alphabet.

\emptyset = the empty language. NB! $\{\epsilon\} \neq \emptyset$.

Example: $\Sigma = \{0, 1\}$

- $L_1 = \{00, 01, 10, 11\}$ all strings of length 2
- $L_2 = \{1, 01, 11, 001, \dots, 111, \dots\}$
all strings which finish on 1
- $L_3 = \emptyset$ all strings of length 1 which finish on 01

Σ^* denotes the set of all strings which can be constructed from the alphabet.

$*$ = closure, Kleene closure.

$+$ = positive closure.

- e.g. $\Sigma = \{0, 1\}$
- $\Sigma^* = \{\epsilon, 0, 1, 00, 01, \dots, 111, 101, \dots\}$
- $\Sigma^+ = \Sigma^* - \{\epsilon\} = \{0, 1, 00, 01, \dots\}$

Operations on languages

Concatenation

L, M are languages.

$L \bullet M = LM = \{xy | x \in L \text{ and } y \in M\}$

$L\{\epsilon\} = \{\epsilon\}L = L$

$L\emptyset = \emptyset L = \emptyset$

Example: $L = \{ab, cd\}$ $M = \{uv, yz\}$

gives us

$LM = \{abuv, abyz, cduv, cdyz\}$

Exponents of languages

$$L^0 = \{\epsilon\}$$

$$L^1 = L$$

$$L^2 = L \cdot L$$

$$L^n = L \cdot L^{n-1}, n \geq 1$$

Union of languages

L, M are languages.
 $L \cup M = \{x \mid x \in L \text{ or } x \in M\}$
 Example: $L = \{ab, cd\}, M = \{uv, yz\}$
 gives us $L \cup M = \{ab, cd, uv, yz\}$

Closure

$$L^* = L^0 \cup L^1 \cup \dots \cup L^\infty$$

Positive closure

$$L^+ = L^1 \cup L^2 \cup \dots \cup L^\infty$$

$$LL^* = L^* - \{\epsilon\}, \text{ if } \epsilon \notin L$$

$$L^* = \{\epsilon\} \cup L^+$$

Example: $A = \{a, b\}$
 $A^* = \{\epsilon, a, b, aa, ab, ba, bb, \dots\}$
 = All possible sequences of a and b .
 A language over A is always a subset of A^* .

Regular expressions

Are used to describe simple languages, e.g. basic symbols, tokens.

Example: identifier = letter • (letter | digit)*

Regular expressions over an alphabet Σ denote a language (regular set).

Rules for constructing regular expressions

Σ is an alphabet, the regular expression r describes the language L_r , the regular expression s corresponds to the language L_s , etc.

Regular expression r	Language L_r
ϵ	$\{\epsilon\}$
$a, a \in \Sigma$	$\{a\}^\dagger$
union: $(s) \mid (t)$	$L_s \cup L_t$
concatenation: $(s) \cdot (t)$	$L_s \cdot L_t$
repetition: $(s)^*$	L_s^*
repetition: $(s)^+$	L_s^+

\dagger Each symbol in the alphabet Σ is a regular expression which denotes $\{a\}$.
 $*$ = repetition, zero or more times.
 $+$ = repetition, one or more times.

Priorities	
Highest	$*, +$
	\cdot
Lowest	$ $

- Example: $\Sigma = \{a, b\}$
- $r=a \quad L_r=\{a\}$
 - $r=a^* \quad L_r=\{\epsilon, a, aa, aaa, \dots\} = \{a\}^*$
 - $r=a \mid b \quad L_r=\{a, b\} = \{a\} \cup \{b\}$
 - $r=(a \mid b)^*$
 $L_r=\{a, b\}^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$
 - $r=(a^* b^*)^*$
 $L_r=\{a, b\}^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$
 - $r=a \mid ba^*$
 $L_r=\{a, b, ba, baa, baaa, \dots\} = \{a \text{ or } ba^i \mid i \geq 0\}$

NB! $\{a^n b^n \mid n \geq 0\}$ can not be described with regular expressions.

$r=a^* b^*$ gives us $L_r=\{a^i b^j \mid i, j \geq 0\}$ does not work.

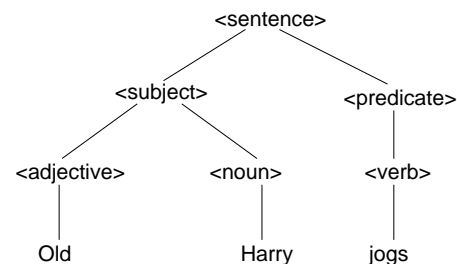
$r=(ab)^*$ gives us $L_r=\{(ab)^i \mid i \geq 0\} = \{\epsilon, ab, abab, \dots\}$ does not work.

Regular expressions can not "count" (have no memory).

Context-free grammars

Example: an English sentence

Sentence:	Old	Harry	jogs
Constituent:	subject		predicate
Word class:	adjective	noun	verb



A grammar is used to describe the syntax. BNF (Backus-Naur form) 1960 (metalanguage to describe languages):

- $\langle \text{sentence} \rangle \rightarrow \langle \text{subject} \rangle \langle \text{predicate} \rangle$
- $\langle \text{subject} \rangle \rightarrow \langle \text{adjective} \rangle \langle \text{noun} \rangle$
- $\langle \text{predicate} \rangle \rightarrow \langle \text{verb} \rangle$
- $\langle \text{adjective} \rangle \rightarrow \text{old} \mid \text{big} \mid \text{strong} \mid \dots$
- $\langle \text{noun} \rangle \rightarrow \text{Harry} \mid \text{brother} \mid \dots$
- $\langle \text{verb} \rangle \rightarrow \text{jogs} \mid \text{snores} \mid \text{sleeps} \mid \dots$

- <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"

The grammar can be used to produce or derive sentences.

Example: <sentence> $\xrightarrow{*}$ Old Harry jogs
where <sentence> is the start symbol and " $\xrightarrow{*}$ " means derivation in zero or more steps.

e.g. Derivation

<sentence> \Rightarrow <subject> <predicate>
 \Rightarrow <adjective> <noun> <predicate>
 \Rightarrow Old <noun> <predicate>
 \Rightarrow Old Harry <predicate>
 \Rightarrow Old Harry <verb>
 \Rightarrow Old Harry jogs

Definition:

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

- G = < N, Σ , P, S > where
- N : nonterminals.
- Σ : terminal symbols.
- P : rules, productions of the form
A → a where A ∈ N and a ∈ (N ∪ Σ)*
- S : the start symbol, a nonterminal, S ∈ N.

(Sometimes V = N ∪ Σ is used, called the *vocabulary*.)

Example:

- 1.<number> → <no>
- 2.<no> → <no> <digit>
3. | <digit>
- 4.<digit> → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- N = { <number>, <no>, <digit> }
- Σ = { 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 }
- S = <number>

A grammar G denotes the language L(G),
L(G) is generated by the grammar G.

$$L(\text{<number>}) = \{ w \mid w \in \text{digit}^+ \}$$

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

Derivation

$\alpha \Rightarrow \beta$ (pronounced "α derives β")

Formally: $\gamma A \theta \Rightarrow \gamma \delta \theta$ if we have $A \rightarrow \delta$

Example: <number> $\xRightarrow{=}$ <no> $\xRightarrow{=}$ <no> <digit> $\xRightarrow{=}$ <no> 2 $\xRightarrow{=}$ <digit> 2 $\xRightarrow{=}$ 12

- <number> \Rightarrow <no> direct derivation.
- <number> $\xrightarrow{*}$ 12 several derivations (zero or more).
- <number> $\xrightarrow{+}$ 12 several derivations (one or more).

Given G = < N, Σ , P, S > the language generated by G can be defined as L(G):

$$L(G) = \{ w \mid S \xrightarrow{*} w \text{ och } w \in \Sigma^+ \}$$

Sentential form

A string α is a *sentential form* in G if

$$S \xrightarrow{*} \alpha \text{ and } \alpha \in V^* \text{ (string of terminals and nonterminals.)}$$

Example: <no> <digit> is a sentential form in G(<number>).

Sentence

w is a *sentence* in G if $S \xrightarrow{*} w$ and $w \in \Sigma^*$.

Example: 12 is a sentence in G(<number>).

Left derivation

\xRightarrow{L} 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)

\xRightarrow{R} 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.

Reverse rightmost derivation

$12 \xrightarrow{\text{red}} \langle \text{digit} \rangle 2 \xrightarrow{\text{red}} \langle \text{no} \rangle 2 \xrightarrow{\text{red}} \langle \text{no} \rangle \langle \text{digit} \rangle$
 $\xrightarrow{\text{red}} \langle \text{no} \rangle \xrightarrow{\text{red}} \langle \text{number} \rangle$

Handles

Consist of two parts:

1. A production $A \rightarrow \beta$
2. A position

If $S \xrightarrow{\text{red}} \alpha Aw \xrightarrow{\text{red}} \alpha \beta w$, the production is $A \rightarrow \beta$ and the position after a is a handle of $\alpha \beta w$.

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:

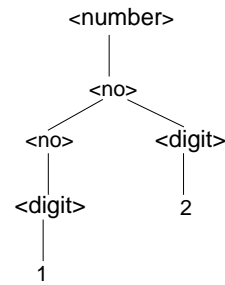
$\langle \text{number} \rangle \xrightarrow{\text{red}} \langle \text{no} \rangle \xrightarrow{\text{red}} \langle \text{no} \rangle \langle \text{digit} \rangle \xrightarrow{\text{red}} \langle \text{no} \rangle 2 \xrightarrow{\text{red}} \langle \text{digit} \rangle 2 \xrightarrow{\text{red}} 12$

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

Reduction

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)



Parse tree för 12

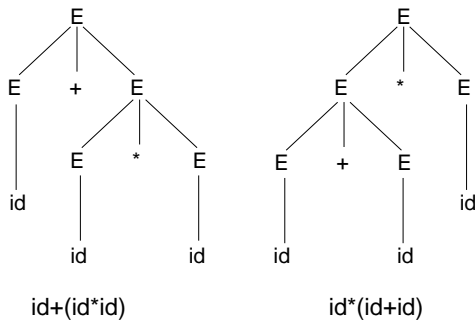
A parse tree can correspond to several different derivations.

Ambiguous grammars

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

e.g. $E \rightarrow E + E$
 $\quad | E * E$
 $\quad | E \uparrow E$
 $\quad | id$

id+id*id has two different parse trees:



Rewrite the grammar to make it unambiguous:

- $+$, $*$ are to have the right priority and
- $+$, $*$ are to be left associative while
- \uparrow is to be right associative.

Example: $a+b+c+d$ is interpreted as $(a+b)+c+d$.

$E \rightarrow E + T$ (left associative)
 $\quad | T$
 $T \rightarrow T * F$ (left associative)
 $\quad | F$
 $F \rightarrow P \uparrow F$ (right associative)
 $\quad | P$
 $P \rightarrow id$

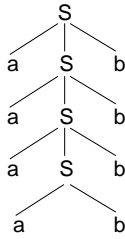
Home assignment:

Draw the parse trees for:

1. $id * id * id$
2. $id \uparrow id \uparrow id$

Example: The following grammar generates
{ $a^n b^n \mid n \geq 1$ }.

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



e.g. $S \rightarrow 0 S 0$
| $1 S 1$
| 0
| 1

describes binary palindromes of odd length.

Home assignment:
Write a CFG for binary palindromes of all lengths ≥ 1 ,
i.e. ϵ is not included.

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>
    
```

```

<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>
...
    
```