



Faculty of Science



An Introduction to Lexical and Syntax Analysis

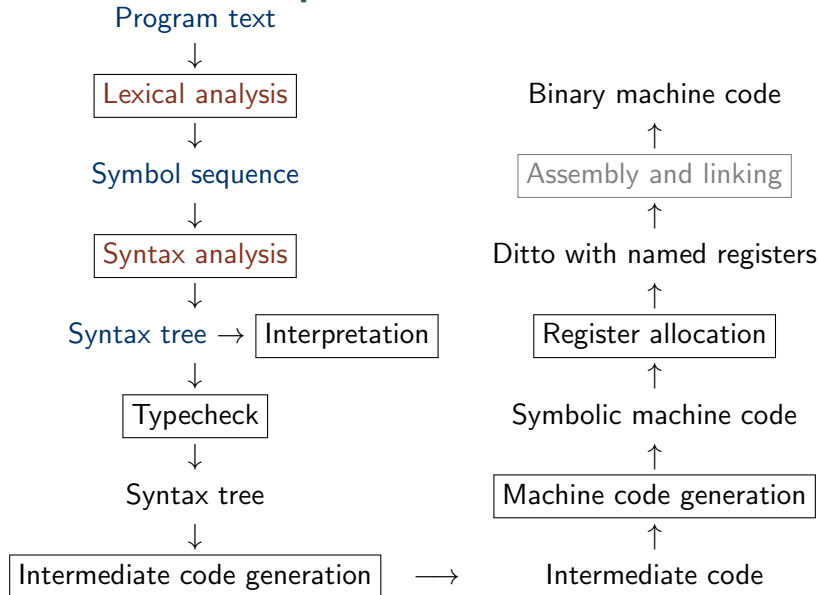
Slides by Jost Berthold and Cosmin E. Oancea
(With tweaks by Andrzej Filinski, andrzej@di.ku.dk)

Department of Computer Science (DIKU)
University of Copenhagen

April 2023 IPS Lecture Slides



Structure of a Compiler



1 Lexical Analysis; Regular Expressions

2 Syntax Analysis; Context-Free Grammars

Lexical Analysis

- **Lexical** (linguistics, NLP): relating to the words of the **vocabulary** of a language
 - (as opposed to **grammar**, i.e., correct construction of sentences)
 - “My mother **cooooookes** dinner not.”
- **Lexical Analyzer**, a.k.a. **lexer**, **scanner**, or **tokenizer**, splits the input program, seen as a stream of **characters**, into a sequence of **tokens**.
- **Tokens** (or **lexemes**) are the basic units of the (programming) language syntax, e.g., keywords, numbers, comments, parentheses, semicolon,

Compiler Phases

```
// My program
let result =
  let x = 10 :: 20 :: 0x30 :: [] in
    List.map (fun a -> 2 * 2 * a) x
```

- Input file also contains comments and meaningful formatting, which helps user only.
- Input file is read as a string, see below:

```
"// My program\nlet result =\n  let x = 10 :: 20 :: 0x30 :: [] in\n    List.map (fun a -> 2 * 2 * a) x"
```

Compiler Phases

```
// My program
let result =
  let x = 10 :: 20 :: 0x30 :: [] in
    List.map (fun a -> 2 * 2 * a) x
```

- Input file also contains comments and meaningful formatting, which helps user only.
- Input file is read as a string, see below:

```
"// My program\n let result =\n  let x = 10 :: 20 :: 0x30 :: [] in\n    List.map (fun a -> 2 * 2 * a) x"
```

Lexical Analysis: transforms a character stream to a token sequence.

Keywd_Let, Id "result", Equal, Keywd_Let, Id "x", Equal, Int 10, Op_Cons, Int 20, Op_Cons, Int 48, Op_Cons, LBracket, RBracket, Keywd_In, Id "List", Dot, Id "map", LParen, Keywd_fun Id "a", Arrow, Int 2, Multiply, Int 2, Multiply, Id "a", RParen, Id "x"

Compiler Phases

```
// My program
let result =
  let x = 10 :: 20 :: 0x30 :: [] in
    List.map (fun a -> 2 * 2 * a) x
```

- Input file also contains comments and meaningful formatting, which helps user only.
- Input file is read as a string, see below:

```
"// My program\n let result =\n  let x = 10 :: 20 :: 0x30 :: [] in\n    List.map (fun a -> 2 * 2 * a) x"
```

Lexical Analysis: transforms a character stream to a token sequence.

```
Keywd_Let, Id "result", Equal, Keywd_Let, Id "x", Equal, Int 10,
Op_Cons, Int 20, Op_Cons, Int 48, Op_Cons, LBracket, RBracket, Keywd_In,
Id "List", Dot, Id "map", LParen, Keywd_fun Id "a", Arrow, Int 2, Multiply,
Int 2, Multiply, Id "a", RParen, Id "x"
```

- **Tokens can be:** (fixed) vocabulary words, e.g., keywords (`let`), built-in operators (`*`, `::`), special symbols (`[`, `]`).
- **Identifiers** and **Number Literals** are **classes** of tokens, which are formed compositionally according to certain rules.

Formalism

Definition (Formal Languages)

Let Σ be an *alphabet*, i.e., a finite set of allowed characters.

- **A word** over Σ is a string of chars $w = a_1a_2 \dots a_n$, $a_i \in \Sigma$
 $n = 0$ is allowed and results in the empty string, denoted ϵ .
 Σ^* is the set of all words over Σ .
- **A language** L over Σ is a set of words over Σ , i.e., $L \subseteq \Sigma^*$.

Examples over the alphabet of lowercase Latin letters:

- Σ^* and \emptyset
- All C keywords: {auto, break, case, ..., while}
- $\{a^n b^n \mid n \geq 0\}$
- All palindromes: {kayak, racecar, mellem, retter, ...}
- $\{a^n b^n c^n \mid n \geq 0\}$

Formal languages encountered in compiler

Aim of compiler's front end: decide whether a program respects the object-language rules.

- **Lexical analysis:** decides whether the individual tokens are well formed. Requires the recognition of a simple language.
- **Syntactical analysis:** decides whether the composition of tokens is well formed: more complex language that checks compliance to grammar rules.
- **Type checking:** verifies that the program complies with (some of) the object language's typing rules. **Very complex (but still effectively decidable) language**

Language Examples: Number Literals in C++

- Integers in decimal format: 234, 0, or 8; but not 08 or iv
- Integers in hexadecimal format: 0X0123, 0xcafe; not 0X or 00Xa
- Floating point decimals: 0., .345, or 123.45; not 3, . or 0.1.2
- Scientific notation: 234E-45 or 0.E123 or .234e+45, ...

Language Examples: Number Literals in C++

- Integers in decimal format: 234, 0, or 8; but not 08 or iv
- Integers in hexadecimal format: 0X0123, 0xcafe; not 0X or 00Xa
- Floating point decimals: 0., .345, or 123.45; not 3, . or 0.1.2
- Scientific notation: 234E-45 or 0.E123 or .234e+45, ...
- A decimal integer is either 0 or a non-empty sequence of digits (0-9) that does not start with 0.
- A hexadecimal integer starts with 0x or 0X and is followed by one or more hexadecimal digits (0-9 or a-f or A-F).

Language Examples: Number Literals in C++

- Integers in decimal format: 234, 0, or 8; but not 08 or iv
- Integers in hexadecimal format: 0X0123, 0xcafe; not 0X or 00Xa
- Floating point decimals: 0., .345, or 123.45; not 3, . or 0.1.2
- Scientific notation: 234E-45 or 0.E123 or .234e+45, ...
- A decimal integer is either 0 or a non-empty sequence of digits (0-9) that does not start with 0.
- A hexadecimal integer starts with 0x or 0X and is followed by one or more hexadecimal digits (0-9 or a-f or A-F).
- Floating-point csts have a “mantissa,” [...and] an “exponent,” [...]. The mantissa is as a sequence of digits followed by a period, followed by an optional sequence of digits[...]. The exponent, if present, specifies the magnitude[...] using e or E[...] followed by an optional sign (+ or -) and a sequence of digits. If an exponent is present, the trailing decimal point is unnecessary in whole numbers. <http://msdn.microsoft.com/en-us/library/tfh6f0w2.aspx>.

Not really a firm basis for automatic processing!

Regular Expressions

We need a formal, **compositional** (and intuitive) description of what tokens are, *and* automatic implementation of the token language.

Definition (Regular Expressions)

The set $RE(\Sigma)$ of regular expressions over alphabet Σ is defined:

- *Base Rules (Non Recursive):*
 - $\epsilon \in RE(\Sigma)$ describes the lang consisting of *only the empty string*.

Regular Expressions

We need a formal, **compositional** (and intuitive) description of what tokens are, *and* automatic implementation of the token language.

Definition (Regular Expressions)

The set $RE(\Sigma)$ of regular expressions over alphabet Σ is defined:

- *Base Rules (Non Recursive):*
 - $\epsilon \in RE(\Sigma)$ describes the lang consisting of *only the empty string*.
 - $a \in RE(\Sigma)$, for $a \in \Sigma$, describes the lang. of *one-letter word* a .
- *Recursive Rules: for every $r, s \in RE(\Sigma)$,*

Regular Expressions

We need a formal, **compositional** (and intuitive) description of what tokens are, *and* automatic implementation of the token language.

Definition (Regular Expressions)

The set $RE(\Sigma)$ of regular expressions over alphabet Σ is defined:

- *Base Rules (Non Recursive):*
 - $\epsilon \in RE(\Sigma)$ describes the lang consisting of *only the empty string*.
 - $a \in RE(\Sigma)$, for $a \in \Sigma$, describes the lang. of *one-letter word* a .
- *Recursive Rules: for every $r, s \in RE(\Sigma)$,*
 - $r \cdot s \in RE(\Sigma)$, *sequence/concatenation*: words in which first part is described by r , and second part by s .

Regular Expressions

We need a formal, **compositional** (and intuitive) description of what tokens are, *and* automatic implementation of the token language.

Definition (Regular Expressions)

The set $RE(\Sigma)$ of regular expressions over alphabet Σ is defined:

- *Base Rules (Non Recursive):*
 - $\epsilon \in RE(\Sigma)$ describes the lang consisting of **only the empty string**.
 - $a \in RE(\Sigma)$, for $a \in \Sigma$, describes the lang. of **one-letter word** a .
- *Recursive Rules: for every $r, s \in RE(\Sigma)$,*
 - $r \cdot s \in RE(\Sigma)$, **sequence/concatenation**: words in which first part is described by r , and second part by s .
 - $r \mid s \in RE(\Sigma)$, **alternative/union**: words described by r OR by s .

Regular Expressions

We need a formal, **compositional** (and intuitive) description of what tokens are, *and* automatic implementation of the token language.

Definition (Regular Expressions)

The set $RE(\Sigma)$ of regular expressions over alphabet Σ is defined:

- *Base Rules (Non Recursive):*
 - $\epsilon \in RE(\Sigma)$ describes the lang consisting of **only the empty string**.
 - $a \in RE(\Sigma)$, for $a \in \Sigma$, describes the lang. of **one-letter word** a .
- *Recursive Rules: for every $r, s \in RE(\Sigma)$,*
 - $r \cdot s \in RE(\Sigma)$, **sequence/concatenation**: words in which first part is described by r , and second part by s .
 - $r \mid s \in RE(\Sigma)$, **alternative/union**: words described by r OR by s .
 - $r^* \in RE(\Sigma)$, **repetition**: zero or more words described by r .

- We may use parentheses (...) for grouping regular expressions.
- We will often omit the explicit \cdot in sequences.
- Sequence groups tighter than alternative: $a|bc^* = a|(b(c^*))$.

Demonstrating Regular-Expression Combinators

$r \cdot s$ Assume the languages of regular expressions r and s are $L(r) = \{ "a", "b" \}$ and $L(s) = \{ "c", "de" \}$, respectively.

Then $L(r \cdot s) =$

Demonstrating Regular-Expression Combinators

$r \cdot s$ Assume the languages of regular expressions r and s are $L(r) = \{ "a", "b" \}$ and $L(s) = \{ "c", "de" \}$, respectively.

Then $L(r \cdot s) = \{ "ac", "ade", "bc", "bde" \}$.

When matching keywords, `if` is the concatenation of two regular expressions: `i` and `f`.

r^* Assume the language of regular expression r is $L(r) = \{ "a", "bb" \}$.

Then

$L(r^*) =$

Demonstrating Regular-Expression Combinators

$r \cdot s$ Assume the languages of regular expressions r and s are $L(r) = \{ "a", "b" \}$ and $L(s) = \{ "c", "de" \}$, respectively.

Then $L(r \cdot s) = \{ "ac", "ade", "bc", "bde" \}$.

When matching keywords, `if` is the concatenation of two regular expressions: `i` and `f`.

r^* Assume the language of regular expression r is $L(r) = \{ "a", "bb" \}$.

Then

$L(r^*) = \{ "", "a", "bb", "aa", "abb", "bba", "bbbb", "aaa", \dots \}$.

Examples: Integers and Variable Names in C++

- Integers in decimal format: 234, 0, 8; but not 08 or iv
- Integers in hexadecimal format: 0X0123, 0xcafe; not 0X, 00Xa
- A variable name consists of letters, digits and underscore, and it must begin with a letter or underscore.

Examples: Integers and Variable Names in C++

- Integers in decimal format: 234, 0, 8; but not 08 or iv
- Integers in hexadecimal format: 0X0123, 0xcafe; not 0X, 00Xa
- A variable name consists of letters, digits and underscore, and it must begin with a letter or underscore.

- Integers in decimal format:

$(1|2|\dots|9)(0|1|2|\dots|9)^* \mid 0$

Shorthand via “character range” $[-]$: $[1-9][0-9]^* \mid 0$

- Integers in hexadecimal format:

$0(x|X)[0-9a-fA-F][0-9a-fA-F]^*$

Shorthand via “at least one” $(+)$: $0(x|X)[0-9a-fA-F]^+$

- Variable names: $[a-zA-Z_][a-zA-Z_0-9]^*$

Useful Abbreviations for Regular Expressions

- **Character Sets:** $[a_1 a_2 \dots a_n] := (a_1 \mid a_2 \mid \dots \mid a_n)$, i.e., one of $a_1, a_2, \dots, a_n \in \Sigma$.
- **Negated Character Sets:** $[\neg a_1 a_2 \dots a_n]$ describes any $a \in \Sigma \setminus \{a_1, a_2, \dots, a_n\}$.
- **Character Ranges:** $[a_1 - a_n] := (a_1 \mid a_2 \mid \dots \mid a_n)$, where $\{a_i\}$ is ordered, i.e., one character in the range a_1 through a_n .
- **Optional Parts:** $r? := (r \mid \epsilon)$ for $r \in RE(\Sigma)$, optionally a string described by r .
- **Repeated Parts:** $r^+ := (r r^*)$ for $r \in RE(\Sigma)$, *at least one* string described by r (but possibly more).

Properties of Regular Expression Combinators

- $|$ is associative: $(r|s)|t = r|(s|t) = r|s|t$
- $|$ is commutative: $s|t = t|s$
- $|$ is idempotent: $s|s = s$
- Also, by definition: $s? = s|\epsilon$
- \cdot is associative: $(rs)t = r(st) = rst$
- ϵ is neutral element for \cdot : $s\epsilon = \epsilon s = s$
- \cdot distributes over $|$: $r(s|t) = rs|rt$, and $(r|s)t = rt|st$.
- $*$ is idempotent: $(s^*)^* = s^*$.
- Also, $s^*s^* = s^*$, and $ss^* = s^+ = s^*s$ by definition!

In all of these, $=$ means “describes the same language as”.

Working with regular expressions

Matching problem:

- Given RE r and word w , does $w \in L(r)$?
- Not obvious how to decide (efficiently), e.g., does $\text{"abbabab"} \in L((abb|ab)^*(b|ba)^*)$?
- Naive **backtracking** algorithm takes time $O(2^{|r| \cdot |w|})$.
- But with some cleverness can get it down to $O(|r| \cdot |w|)$, or even $O(|w|)$, using an $O(2^{|r|})$ -sized lookup table.
 - Will see how, later in the course!

Tokenization problem:

- Given REs r_1, \dots, r_n (each describing a valid token form), and word w (entire input string), do there exist i_1, \dots, i_m , such that $w \in L(r_{i_1} \cdots r_{i_m})$ (input divided into tokens)?
- When multiple tokenizations possible:
 - At each point, prefer **longest possible** prefix of remaining input.
 - At each point, prefer **lowest-numbered** RE matching that prefix.
- Will see how algorithms and data structures for matching generalize to efficient tokenization. Tool: `(fs)lex`.

1 Lexical Analysis; Regular Expressions

2 Syntax Analysis; Context-Free Grammars

Syntax Analysis (Parsing)

Relates to the correct construction of sentences, i.e., grammar.

- 1 Checks that grammar is respected, otherwise **syntax error**, and
- 2 Arranges tokens into a **syntax tree** reflecting the text structure: leaves are tokens, which if read from left to right results in the original text!

mother
cooks
dinner
My.

syntax error

My dinner
cooks
mother.

semantic error

Essential tool and theory used are *Context-Free Grammars*:
a notation suitable for human understanding that can be transformed
into an efficient implementation.

Context-Free Grammar (CFG) Definition

- 1 a set of *terminals* Σ – the language alphabet, e.g., the set of tokens produced by lexer. (Convention: use lowercase letters.)
- 2 a set of *non-terminals* N , denoting sets of recursively defined strings. (Convention: use uppercase letters.)
- 3 a *start symbol* $S \in N$, denoting the lang defined by the grammar.
- 4 a set P of productions of form $Y \rightarrow X_1 \dots X_n$, where $Y \in N$ is a (single) non-terminal, and $X_i \in (\Sigma \cup N)$, $\forall i$ can be a terminal or non-terminal. Each production describes some of the strings of the corresponding non-terminal Y .

May **abbreviate** productions $Y \rightarrow \vec{X}_1, Y \rightarrow \vec{X}_2$ as $Y \rightarrow \vec{X}_1 \mid \vec{X}_2$.

G: $S \rightarrow aS$

$S \rightarrow \epsilon$

Context-Free Grammar (CFG) Definition

- 1 a set of *terminals* Σ – the language alphabet, e.g., the set of tokens produced by lexer. (Convention: use lowercase letters.)
- 2 a set of *non-terminals* N , denoting sets of recursively defined strings. (Convention: use uppercase letters.)
- 3 a *start symbol* $S \in N$, denoting the lang defined by the grammar.
- 4 a set P of productions of form $Y \rightarrow X_1 \dots X_n$, where $Y \in N$ is a (single) non-terminal, and $X_i \in (\Sigma \cup N)$, $\forall i$ can be a terminal or non-terminal. Each production describes some of the strings of the corresponding non-terminal Y .

May **abbreviate** productions $Y \rightarrow \vec{X}_1, Y \rightarrow \vec{X}_2$ as $Y \rightarrow \vec{X}_1 \mid \vec{X}_2$.

G: $S \rightarrow aS$

$S \rightarrow \epsilon$

regular-expression

language a^*

G: $S \rightarrow aSb$

$S \rightarrow \epsilon$

Context-Free Grammar (CFG) Definition

- 1 a set of *terminals* Σ – the language alphabet, e.g., the set of tokens produced by lexer. (Convention: use lowercase letters.)
- 2 a set of *non-terminals* N , denoting sets of recursively defined strings. (Convention: use uppercase letters.)
- 3 a *start symbol* $S \in N$, denoting the lang defined by the grammar.
- 4 a set P of productions of form $Y \rightarrow X_1 \dots X_n$, where $Y \in N$ is a (single) non-terminal, and $X_i \in (\Sigma \cup N)$, $\forall i$ can be a terminal or non-terminal. Each production describes some of the strings of the corresponding non-terminal Y .

May **abbreviate** productions $Y \rightarrow \vec{X}_1, Y \rightarrow \vec{X}_2$ as $Y \rightarrow \vec{X}_1 \mid \vec{X}_2$.

G: $S \rightarrow aS$

$S \rightarrow \epsilon$

regular-expression

language a^*

G: $S \rightarrow aSb$

$S \rightarrow \epsilon$

describes language

$\{a^n b^n, \forall n \geq 0\}$

G: $S \rightarrow aSa \mid bSb \mid \dots$

$S \rightarrow a \mid b \mid \dots \mid \epsilon$

Context-Free Grammar (CFG) Definition

- 1 a set of *terminals* Σ – the language alphabet, e.g., the set of tokens produced by lexer. (Convention: use lowercase letters.)
- 2 a set of *non-terminals* N , denoting sets of recursively defined strings. (Convention: use uppercase letters.)
- 3 a *start symbol* $S \in N$, denoting the lang defined by the grammar.
- 4 a set P of productions of form $Y \rightarrow X_1 \dots X_n$, where $Y \in N$ is a (single) non-terminal, and $X_i \in (\Sigma \cup N)$, $\forall i$ can be a terminal or non-terminal. Each production describes some of the strings of the corresponding non-terminal Y .

May **abbreviate** productions $Y \rightarrow \vec{X}_1, Y \rightarrow \vec{X}_2$ as $Y \rightarrow \vec{X}_1 \mid \vec{X}_2$.

G: $S \rightarrow aS$

$S \rightarrow \epsilon$

regular-expression

language a^*

G: $S \rightarrow aSb$

$S \rightarrow \epsilon$

describes language

$\{a^n b^n, \forall n \geq 0\}$

G: $S \rightarrow aSa \mid bSb \mid \dots$

$S \rightarrow a \mid b \mid \dots \mid \epsilon$

describes palindromes,

e.g., *abba*, *babab*.

The latter two languages cannot be described with regular expressions.

Example: Deriving Words

Nonterminals recursively refer to themselves or each other
(cannot do that with regular expressions):

$$G: S \rightarrow aSB \quad (1)$$

$$S \rightarrow \epsilon \quad (2) \quad G: S \rightarrow aSB \mid \epsilon \quad S = \{a \cdot x \cdot y \mid x \in S, y \in B\} \cup \{\epsilon\}$$

$$B \rightarrow Bb \quad (3) \quad B \rightarrow Bb \mid b \quad B = \{x \cdot b \mid x \in B\} \cup \{b\}$$

$$B \rightarrow b \quad (4)$$

Words of the language can be constructed by

- starting with the start symbol S , and
- successively replacing nonterminals with right-hand sides.

Deriving $aaabbbb$ (each step replaces \underline{Y} on LHS with $\overline{X_1 \dots X_n}$):

Example: Deriving Words

Nonterminals recursively refer to themselves or each other
(cannot do that with regular expressions):

$$G: S \rightarrow aSB \quad (1)$$

$$S \rightarrow \epsilon \quad (2) \quad G: S \rightarrow aSB \mid \epsilon \quad S = \{a \cdot x \cdot y \mid x \in S, y \in B\} \cup \{\epsilon\}$$

$$B \rightarrow Bb \quad (3) \quad B \rightarrow Bb \mid b \quad B = \{x \cdot b \mid x \in B\} \cup \{b\}$$

$$B \rightarrow b \quad (4)$$

Words of the language can be constructed by

- starting with the start symbol S , and
- successively replacing nonterminals with right-hand sides.

Deriving $aaabbbb$ (each step replaces \underline{Y} on LHS with $\overline{X_1 \dots X_n}$):

$$\underline{S} \Rightarrow^1 \overline{a\underline{S}B} \Rightarrow^1 \overline{aa\underline{S}B}B \Rightarrow^4 \overline{aa\underline{S}b}B \Rightarrow^1 \overline{aaa\underline{S}B}bB$$

Example: Deriving Words

Nonterminals recursively refer to themselves or each other
(cannot do that with regular expressions):

$$G: S \rightarrow aSB \quad (1)$$

$$S \rightarrow \epsilon \quad (2) \quad G: S \rightarrow aSB \mid \epsilon \quad S = \{a \cdot x \cdot y \mid x \in S, y \in B\} \cup \{\epsilon\}$$

$$B \rightarrow Bb \quad (3) \quad B \rightarrow Bb \mid b \quad B = \{x \cdot b \mid x \in B\} \cup \{b\}$$

$$B \rightarrow b \quad (4)$$

Words of the language can be constructed by

- starting with the start symbol S , and
- successively replacing nonterminals with right-hand sides.

Deriving $aaabbbb$ (each step replaces \underline{Y} on LHS with $\overline{X_1 \dots X_n}$):

$$\begin{aligned} \underline{S} &\Rightarrow^1 \overline{aSB} \Rightarrow^1 \overline{aaSB}B \Rightarrow^4 \overline{aaSb}B \Rightarrow^1 \overline{aaaSB}bB \\ &\Rightarrow^2 \overline{aaaB}bbB \Rightarrow^3 \overline{aaaBbb}B \Rightarrow^4 \overline{aaaBbbb} \Rightarrow^4 aaabbbb \end{aligned}$$

Definition: Derivation Relation

Let $G = (\Sigma, N, S, P)$ be a grammar.

The derivation relation \Rightarrow on $(\Sigma \cup N)^*$ is defined as:

- For a nonterminal $X \in N$ and a production $(X \rightarrow \beta) \in P$,
 $\alpha_1 X \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2$, for all $\alpha_1, \alpha_2 \in (\Sigma \cup N)^*$
- Describes one derivation step using one of the productions.
- Each step may be optionally annotated with the grammar-rule number.

G: $S \rightarrow aSB$ (1)

$S \rightarrow \epsilon$ (2) $S \Rightarrow^1 aSB \Rightarrow^1 aaSBB \Rightarrow^2 aa BB$

$B \rightarrow Bb$ (3)

$B \rightarrow b$ (4)

Definition: Derivation Relation

Let $G = (\Sigma, N, S, P)$ be a grammar.

The derivation relation \Rightarrow on $(\Sigma \cup N)^*$ is defined as:

- For a nonterminal $X \in N$ and a production $(X \rightarrow \beta) \in P$,
 $\alpha_1 X \alpha_2 \Rightarrow \alpha_1 \beta \alpha_2$, for all $\alpha_1, \alpha_2 \in (\Sigma \cup N)^*$
- Describes one derivation step using one of the productions.
- Each step may be optionally annotated with the grammar-rule number.

G: $S \rightarrow aSB$ (1)

$S \rightarrow \epsilon$ (2)

$B \rightarrow Bb$ (3)

$B \rightarrow b$ (4)

$S \Rightarrow^1 aSB \Rightarrow^1 aaSBB \Rightarrow^2 aa BB$
 $\Rightarrow^3 aaBbB \Rightarrow^4 aabbB \Rightarrow^4 aabbbb.$

- Here we have used **leftmost derivation**, i.e., always expanded the leftmost terminal first. Could also use **right-most derivation**.
- $aaabbbb$ and $aabbb \in L(G)$.

Transitive Derivation Relation Definition

Let $G = (\Sigma, N, S, P)$ be a grammar and \Rightarrow its derivation relation.

The transitive derivation relation \Rightarrow^* is defined as:

- $\alpha \Rightarrow^* \alpha$, for $\alpha \in (\Sigma \cup N)^*$, derived in 0 steps,
- for $\alpha, \beta \in (\Sigma \cup N)^*$, $\alpha \Rightarrow^* \beta$ iff there exists $\gamma \in (\Sigma \cup N)^*$ such that $\alpha \Rightarrow \gamma$, and $\gamma \Rightarrow^* \beta$, i.e., derived in at least one step.

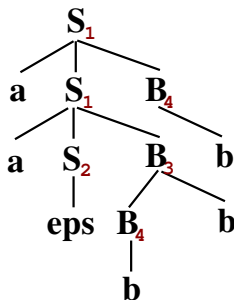
The Language of a Grammar consists of all the words that can be obtained via the transitive derivation relation:

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}.$$

For example $aaabbbb$ and $aabbb \in L(G)$,
because $S \Rightarrow^* aaabbbb$ and $S \Rightarrow^* aabbb$.

Syntax Trees

$$\begin{aligned}
 G: S &\rightarrow aSB \quad (1) \\
 S &\rightarrow \epsilon \quad (2) \\
 B &\rightarrow Bb \quad (3) \\
 B &\rightarrow b \quad (4)
 \end{aligned}$$



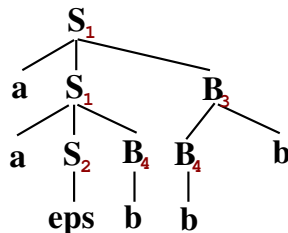
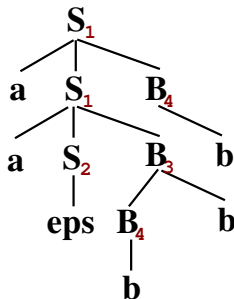
Syntax trees describe the “structure” of the derivation (independent of the order in which nonterminals have been chosen to be derived).

Leftmost derivation always derives the leftmost nonterminal first, and corresponds to a *depth-first, left-to-right, preorder tree traversal*:

$$\underline{S} \Rightarrow^1 \underline{aSB} \Rightarrow^1 \underline{aaSB} \Rightarrow^2 \underline{aaBB} \Rightarrow^3 \underline{aaBbB} \Rightarrow^4 \underline{aabbB} \Rightarrow^4 \underline{aabb}.$$

Syntax Trees & Ambiguous Grammars

$$\begin{aligned}
 G: S &\rightarrow aSB \quad (1) \\
 S &\rightarrow \epsilon \quad (2) \\
 B &\rightarrow Bb \quad (3) \\
 B &\rightarrow b \quad (4)
 \end{aligned}$$



Syntax trees describe the “structure” of the derivation (independent of the order in which nonterminals have been chosen to be derived).

The grammar is said to be **ambiguous** if there exists a word that can be derived in two ways, corresponding to different syntax trees.

$$\begin{aligned}
 \underline{S} &\Rightarrow^1 \underline{aSB} \Rightarrow^1 \underline{aaSB} \Rightarrow^2 \underline{aaBB} \Rightarrow^3 \underline{aaBb} \Rightarrow^4 \underline{aabbB} \Rightarrow^4 \underline{aabb\bar{b}}. \\
 \underline{S} &\Rightarrow^1 \underline{aSB} \Rightarrow^1 \underline{aaSB} \Rightarrow^2 \underline{aaBB} \Rightarrow^4 \underline{aabB} \Rightarrow^3 \underline{aabBb} \Rightarrow^4 \underline{aabb\bar{b}}.
 \end{aligned}$$

Handling/Removing Grammar Ambiguity

$$E \rightarrow E + E \mid E - E$$

$$E \rightarrow E * E \mid E / E$$

$$E \rightarrow a \mid (E)$$

- *Precedence and Associativity* guide decision:
- ambiguity resolved by **parsing directives**,
- or by rewriting the grammar.

What are the problems:

- Ambiguous derivation of $a + a * a$

Handling/Removing Grammar Ambiguity

$$E \rightarrow E + E \mid E - E$$

$$E \rightarrow E * E \mid E / E$$

$$E \rightarrow a \mid (E)$$

- *Precedence and Associativity* guide decision:
- ambiguity resolved by **parsing directives**,
- or by rewriting the grammar.

What are the problems:

- Ambiguous derivation of $a + a * a$ can be resolved by setting *the precedence* of $*$ higher than $+$: $a + (a * a)$.
- Ambiguous derivation of $a - a - a$

Handling/Removing Grammar Ambiguity

$$E \rightarrow E + E \mid E - E$$

$$E \rightarrow E * E \mid E / E$$

$$E \rightarrow a \mid (E)$$

- *Precedence and Associativity* guide decision:
- ambiguity resolved by *parsing directives*,
- or by rewriting the grammar.

What are the problems:

- Ambiguous derivation of $a + a * a$ can be resolved by setting *the precedence* of $*$ higher than $+$: $a + (a * a)$.
- Ambiguous derivation of $a - a - a$ can be resolved by fixing *a left-associative* derivation: $(a - a) - a$.

Defining/Resolving Operator Precedence

- Introduce precedence levels to set operator priorities
- for example precedence of $*$ and $/$ over (higher than) $+$ and $-$,
- and more precedence levels can be added, e.g., exponentiation.

Defining/Resolving Operator Precedence

- Introduce precedence levels to set operator priorities
- for example precedence of $*$ and $/$ **over (higher than)** $+$ and $-$,
- and more precedence levels can be added, e.g., exponentiation.

At grammar level: this can be accomplished by introducing one nonterminal for each level of precedence:

$$E \rightarrow E + E \mid E - E$$

$$E \rightarrow E * E \mid E / E$$

$$E \rightarrow a \mid (E)$$

$$E \rightarrow E + E \mid E - E \mid T$$

$$T \rightarrow T * T \mid T / T$$

$$T \rightarrow a \mid (E)$$

Defining/Resolving Operator Associativity

A binary operator \oplus is called:

- *left associative* if expression $x \oplus y \oplus z$ should be grouped from left to right: $(x \oplus y) \oplus z$
- *right associative* if expression $x \oplus y \oplus z$ should be grouped from right to left: $x \oplus (y \oplus z)$
- *non-associative* if expressions such as $x \oplus y \oplus z$ are disallowed,
- *associative* if both left-to-right and right-to-left groupings lead to the same result (a *semantic*, not syntactic, property).

Examples:

- *left associative* operators: $-$ and $/$,
- *right associative* operators

Defining/Resolving Operator Associativity

A binary operator \oplus is called:

- *left associative* if expression $x \oplus y \oplus z$ should be grouped from left to right: $(x \oplus y) \oplus z$
- *right associative* if expression $x \oplus y \oplus z$ should be grouped from right to left: $x \oplus (y \oplus z)$
- *non-associative* if expressions such as $x \oplus y \oplus z$ are disallowed,
- *associative* if both left-to-right and right-to-left groupings lead to the same result (a *semantic*, not syntactic, property).

Examples:

- *left associative* operators: $-$ and $/$,
- *right associative* operators: exponentiation, assignment (in C, Java: $a = b = c$), arrow (in F# *types*: $\text{int} \rightarrow \text{int} \rightarrow \text{int}$).

Establishing Intended Associativity

- Can be declared in the parser file via directives
- when operators are semantically associative (e.g., $+$) , use same associativity as comparable operators (e.g., $-$)
- cannot mix left- and right-associative operators at the same precedence level.

Establishing Intended Associativity

- Can be declared in the parser file via directives
- when operators are semantically associative (e.g., +) , use same associativity as comparable operators (e.g., -)
- cannot mix left- and right-associative operators at the same precedence level.

At grammar level: this can be accomplished by introducing new nonterminals that establish explicitly operator's associativity :

$$\begin{aligned}
 E &\rightarrow E + E \mid E - E \mid T \\
 T &\rightarrow T * T \mid T / T \\
 T &\rightarrow a \mid (E)
 \end{aligned}$$

$$\begin{aligned}
 E &\rightarrow E + T \mid E - T \mid T \\
 T &\rightarrow T * F \mid T / F \mid F \\
 F &\rightarrow a \mid (E)
 \end{aligned}$$

Establishing Intended Associativity

- Can be declared in the parser file via directives
- when operators are semantically associative (e.g., +) , use same associativity as comparable operators (e.g., -)
- cannot mix left- and right-associative operators at the same precedence level.

At grammar level: this can be accomplished by introducing new nonterminals that establish explicitly operator's associativity :

$$E \rightarrow E + E \mid E - E \mid T$$

$$T \rightarrow T * T \mid T / T$$

$$T \rightarrow a \mid (E)$$

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T * F \mid T / F \mid F$$

$$F \rightarrow a \mid (E)$$

- Left associative \Rightarrow Left-recursive grammar production.
- Right associative \Rightarrow Right-recursive grammar production.

Working with context-free grammars

Recognition problem:

- Given grammar G and word w , does $w \in L(G)$?
- Not obvious that even decidable, let alone efficiently
- Clever dynamic-programming algorithms work in time $O(|w|^3)$.
- But for “well behaved” grammars, can get $O(|w|)$.
 - Will see later in the course.
- (Some related problems are actually **undecidable**, e.g., given G , does $\forall w \in \Sigma^*. w \in L(G)$?)

Parsing problem:

- Given G and w , construct **parse tree** if $w \in L(G)$ (and produce error message otherwise).
 - Can **annotate** productions with instructions for constructing AST.
- When multiple parses possible, need G to be **disambiguated** (by fixity declarations and/or transformations).
- Most recognition algorithms generalize readily to efficient parsing. Tool: (fs)yacc.

Summary

- Have seen basic concepts of lexing and parsing.
 - Corresponding to *ItCD* secs. 1.1 and 2.1–2.3.
- Should be enough to get you started on the group project.
 - More details in today's LAB
- Will go into considerably more depth with lexing and parsing later in the course.
 - Theory of lexing/parsing and formal languages also useful outside of PL implementation.