

G53CMP: Recap of Basic Formal Language Notions

Henrik Nilsson

University of Nottingham, UK

G53CMP: Recap of Basic Formal Language Notions – p.1/52

Languages (1)

- A **symbol** is a basic indivisible entity. Concrete examples of symbols are letters and digits.
- A **string** or **word** is a **finite** sequence of juxtapositioned symbols. For example: a , b , and c are symbols and $abcb$ is a string.
- An **alphabet** is a **finite** set of symbols. For example: $\{a, b, c\}$, \emptyset .

G53CMP: Recap of Basic Formal Language Notions – p.4/52

Concatenation of Words

- Concatenation of words is denoted by juxtaposition. For example: Concatenation of ab and ba yields $abba$.
- Concatenation is associative and has unit ϵ :

$$u(vw) = (uv)w \\ \epsilon u = u = u\epsilon$$

where u , v , w are words.

G53CMP: Recap of Basic Formal Language Notions – p.7/52

About These Slides

The following slides give a brief recap on some central notions from the theory of formal languages, along with illustrative examples of specific relevance to G53CMP (including the coursework). This is material that has been covered in G52LAC and should be familiar to students taking G53CMP. This material will thus not be covered in detail in the G53CMP lectures, but is offered here for your convenience if you need to refresh these concepts. You may want to go back the G52LAC lecture notes if you need even more details.

G53CMP: Recap of Basic Formal Language Notions – p.2/52

Languages (2)

- ϵ denotes the word of length 0, the **empty word**.
- A **language** (over alphabet Σ) is a set of words (over alphabet Σ). For example: $\Sigma = \{a\}$; one possible language is $L = \{\epsilon, a, aa, aaa\}$.
- Σ^* denotes the set of **all** words over an alphabet Σ , including ϵ .

G53CMP: Recap of Basic Formal Language Notions – p.5/52

Concatenation of Languages (1)

Concatenation of words is extended to languages by:

$$MN = \{uv \mid u \in M \wedge v \in N\}$$

Example:

$$M = \{\epsilon, a, aa\} \\ N = \{b, c\} \\ MN = \{uv \mid u \in \{\epsilon, a, aa\} \wedge v \in \{b, c\}\} \\ = \{eb, ec, ab, ac, aab, aac\} \\ = \{b, c, ab, ac, aab, aac\}$$

G53CMP: Recap of Basic Formal Language Notions – p.8/52

Content

- Formal Languages
- Context-Free Grammars
- Ambiguous Grammars
- Eliminating Ambiguity
 - Dangling else
 - Operator associativity
 - Operator precedence

G53CMP: Recap of Basic Formal Language Notions – p.3/52

Languages: Examples

alphabet	$\Sigma = \{a, b\}$
words	$\epsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, abb, baa, bab, \dots$
languages	$\emptyset, \{\epsilon\}, \{a\}, \{b\}, \{a, aa\}, \{\epsilon, a, aa, aaa\}, \{a^n \mid n \geq 0\}, \{a^n b^n \mid n \geq 0, n \text{ even}\}$

G53CMP: Recap of Basic Formal Language Notions – p.6/52

Concatenation of Languages (2)

- Concatenation of languages is associative:

$$L(MN) = (LM)N$$

- Concatenation of languages has unit $\{\epsilon\}$:

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

- Concatenation distributes through set union:

$$L(M \cup N) = LM \cup LN \\ (L \cup M)N = LN \cup MN$$

G53CMP: Recap of Basic Formal Language Notions – p.9/52

Context-Free Grammars (1)

A **Context-Free Grammar** (CFG) is a way of formally describing **Context-Free Languages** (CFL):

- The CFLs captures ideas common in programming languages such as
 - nested structure
 - balanced parentheses
 - matching keywords like `begin` and `end`.
- Most “reasonable” CFLs can be recognised by a fairly simple machine: a **deterministic pushdown automaton**.

OSICMP: Recap of Basic Formal Language Notions – p.11/52

Context-Free Grammar: Example

$$G = (\{S, A\}, \{a, b\}, P, S)$$

where P consists of the productions

$$\begin{aligned} S &\rightarrow \epsilon \\ S &\rightarrow aA \\ A &\rightarrow bS \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions – p.13/52

The Directly Derives Relation (2)

When it is clear which grammar G is involved, we use \Rightarrow instead of \xRightarrow{G} .

Example: Given the grammar

$$\begin{aligned} S &\rightarrow \epsilon \mid aA \\ A &\rightarrow bS \end{aligned}$$

we have

$$\begin{aligned} S &\Rightarrow \epsilon & aA &\Rightarrow abS \\ S &\Rightarrow aA & SaAaa &\Rightarrow SabSaa \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions – p.16/52

Context-Free Grammars (2)

Thus, describing a programming language by a “reasonable” CFG

- allows context-free constraints to be expressed
- imparts a hierarchical structure to the words in the language
- allows simple and efficient **parsing**:
 - determining if a word belongs to the language
 - determining its **phrase structure** if so.

OSICMP: Recap of Basic Formal Language Notions – p.11/52

Context-Free Grammars: Notation

- Productions with the same LHS are usually grouped together. For example, the productions for S from the previous example:

$$S \rightarrow \epsilon \mid aA$$

This is (roughly) what is known as **Backus-Naur Form**.

- Another common way of writing productions is

$$A ::= \alpha$$

OSICMP: Recap of Basic Formal Language Notions – p.14/52

The Derives Relation (1)

The relation $\xRightarrow{*}_G$, read “derives in grammar G ”, is the reflexive, transitive closure of \xRightarrow{G} .

That is, $\xRightarrow{*}_G$ is the least relation on strings over $N \cup T$ such that:

- $\alpha \xRightarrow{*}_G \beta$ if $\alpha \xRightarrow{G} \beta$
- $\alpha \xRightarrow{*}_G \alpha$ (reflexive)
- $\alpha \xRightarrow{*}_G \beta$ if $\alpha \xRightarrow{*}_G \gamma \wedge \gamma \xRightarrow{*}_G \beta$ (transitive)

OSICMP: Recap of Basic Formal Language Notions – p.17/52

Context-Free Grammars (3)

A **Context-Free Grammar** is a 4-tuple (N, T, P, S) where

- N is a finite set of **nonterminals**
- T is a finite set of **terminals** (the **alphabet** of the language being described)
- $N \cap T = \emptyset$ (N and T are disjoint)
- S , the **start symbol**, is a distinguished element of N
- P is a finite set of productions, written $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in (N \cup T)^*$

OSICMP: Recap of Basic Formal Language Notions – p.15/52

The Directly Derives Relation (1)

To formally define the language generated by

$$G = (N, T, P, S)$$

we first define a binary relation \xRightarrow{G} on strings over $N \cup T$, read “directly derives in grammar G ”, being the least relation such that

$$\alpha A \gamma \xRightarrow{G} \alpha \beta \gamma$$

whenever $A \rightarrow \beta$ is a production in G .

Note: a production can be applied regardless of context, hence **context-free**.

OSICMP: Recap of Basic Formal Language Notions – p.15/52

The Derives Relation (2)

Again, we use $\xRightarrow{*}$ instead of $\xRightarrow{*}_G$ when G is obvious.

Example: Given the grammar

$$\begin{aligned} S &\rightarrow \epsilon \mid aA \\ A &\rightarrow bS \end{aligned}$$

we have

$$\begin{aligned} S &\xRightarrow{*} \epsilon & S &\xRightarrow{*} abS \\ S &\xRightarrow{*} aA & S &\xRightarrow{*} ababS \\ aA &\xRightarrow{*} abS & S &\xRightarrow{*} abab \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions – p.18/52

Language Generated by a Grammar

The language generated by a context-free grammar

$$G = (N, T, P, S)$$

denoted $L(G)$, is defined as follows:

$$L(G) = \{w \mid w \in T^* \wedge S \xrightarrow{*}_G w\}$$

A language L is a **Context-Free Language** (CFL) iff $L = L(G)$ for some CFG G .

A string $\alpha \in (N \cup T)^*$ is a **sentential form** iff $S \xrightarrow{*} \alpha$.

OSICMP: Recap of Basic Formal Language Notions - p.19/52

Derivation Tree

A tree is a **derivation** or **parse tree** for CFG $G = (N, T, P, S)$ if:

- every vertex has a **label** from $N \cup T \cup \{\epsilon\}$
- the label of the root is S
- labels of interior vertices belong to N
- if vertex n has label A and vertices n_1, n_2, \dots, n_k are the children of n , from left to right, with labels X_1, X_2, \dots, X_k , then $A \rightarrow X_1X_2 \dots X_k$ is a production in P
- if a vertex n has label ϵ , then n is a leaf and the only child of its parent.

OSICMP: Recap of Basic Formal Language Notions - p.20/52

Regular Grammars

- Lexical syntax** is usually defined through **Regular Languages**.
- The regular languages are a proper subset of the context-free languages.
- Context-free grammars can thus be used to describe regular languages.
- If a grammar G is **left-linear** or **right-linear**, then G is a **regular grammar** and $L(G)$ is a regular language.
- Regular languages are easy to recognize (DFA).

OSICMP: Recap of Basic Formal Language Notions - p.25/52

Language Generation: Example

Given the grammar

$G = (N = \{S, A\}, T = \{a, b\}, P, S)$ where P are the productions

$$\begin{aligned} S &\rightarrow \epsilon \mid aA \\ A &\rightarrow bS \end{aligned}$$

we have

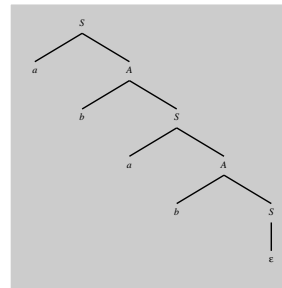
$$\begin{aligned} L(G) &= \{(ab)^i \mid i \geq 0\} \\ &= \{\epsilon, ab, abab, ababab, \dots\} \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions - p.20/52

Derivation Tree: Example

Derivation tree for the string $abab \in L(G)$:

$$\begin{aligned} G: \quad S &\rightarrow \epsilon \mid aA \\ A &\rightarrow bS \end{aligned}$$



OSICMP: Recap of Basic Formal Language Notions - p.20/52

Right-linear Grammar

A CFG $G = (N, T, P, S)$ is **right-linear** if all its productions are of the forms

$$\begin{aligned} A &\rightarrow wB \\ A &\rightarrow w \end{aligned}$$

where $A, B \in N$ and $w \in T^*$.

Example: The regular language $0(10)^*$ is generated by the right-linear grammar

$$\begin{aligned} S &\rightarrow 0A \\ A &\rightarrow 10A \mid \epsilon \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions - p.26/52

Equivalence of Grammars

Two grammars G_1 and G_2 are **equivalent** iff $L(G_1) = L(G_2)$.

Example:

$$\begin{aligned} G_1: \quad S &\rightarrow \epsilon \mid A \\ A &\rightarrow a \mid aA \end{aligned} \quad \begin{aligned} G_2: \quad S &\rightarrow A \\ A &\rightarrow \epsilon \mid Aa \end{aligned}$$

$$L(G_1) = \{a\}^* = L(G_2)$$

Note: the equivalence of CFGs is in general **undecidable**.

OSICMP: Recap of Basic Formal Language Notions - p.21/52

Derivations and Derivation Trees

Given a derivation tree for a grammar G :

- The string of leaf labels read from left to right is the **yield** of the tree.
- The yield is a sentential form of G .

The derives relation and derivation trees are related as follows:

A string α is the yield of some derivation tree for a grammar G iff $S \xrightarrow{*}_G \alpha$.

OSICMP: Recap of Basic Formal Language Notions - p.24/52

Left-linear Grammar

A CFG $G = (N, T, P, S)$ is **left-linear** if all its productions are of the forms

$$\begin{aligned} A &\rightarrow BA \\ A &\rightarrow w \end{aligned}$$

where $A, B \in N$ and $w \in T^*$.

Example: The regular language $0(10)^*$ is generated by the left-linear grammar

$$S \rightarrow S10 \mid 0$$

OSICMP: Recap of Basic Formal Language Notions - p.27/52

Leftmost and Rightmost Derivations

- A derivation is **leftmost** if productions are always applied to the leftmost nonterminal at each step in a derivation.
- A derivation is **rightmost** if productions are always applied to the rightmost nonterminal at each step in a derivation.

Leftmost derivation:

$$G: \begin{array}{l} S \rightarrow AB \mid BA \\ A \rightarrow a \\ B \rightarrow Ab \end{array} \quad \begin{array}{l} S \xRightarrow{lm} BA \xRightarrow{lm} AbA \\ \xRightarrow{lm} abA \xRightarrow{lm} aba \end{array}$$

OSICMP: Recap of Basic Formal Language Notions – p.28/52

Ambiguous Grammars (3)

- Most CFLs are not inherently ambiguous; i.e., an ambiguous CFG G for a language L can often be **transformed** into an **equivalent** but unambiguous grammar G' .
- The ambiguity of a CFG is in general **undecidable**.

OSICMP: Recap of Basic Formal Language Notions – p.31/52

Elim. Ambiguity: Dangling-Else (3)

Note that the distinction is important, as the two trees suggest **different semantics**.

For example, suppose $expr_1$ evaluates to true, and $expr_2$ evaluates to false. Which, if any, of $stmt_1$ and $stmt_2$ gets executed?

OSICMP: Recap of Basic Formal Language Notions – p.34/52

Ambiguous Grammars (1)

A CFG G is **ambiguous** if some word in $L(G)$ has **more than one derivation tree**.

A derivation tree determines a unique leftmost and a unique rightmost derivation.

Thus, equivalently: A CFG G is **ambiguous** if some word in $L(G)$ has

- more than one leftmost derivation**, or
- more than one rightmost derivation**.

OSICMP: Recap of Basic Formal Language Notions – p.29/52

Eliminating Ambiguity: Dangling-Else

Consider the following “dangling-else” grammar:

$$\begin{array}{l} Stmt \rightarrow \text{if } Expr \text{ then } Stmt \\ \quad \quad | \text{if } Expr \text{ then } Stmt \text{ else } Stmt \\ \quad \quad | \text{other} \end{array}$$

and the following program fragment:

$$\text{if } expr_1 \text{ then if } expr_2 \text{ then } stmt_1 \text{ else } stmt_2$$

Two possible parse trees!
Hence the grammar is ambiguous!

OSICMP: Recap of Basic Formal Language Notions – p.32/52

Elim. Ambiguity: Dangling-Else (4)

Preferred interpretation:

“Match each `else` with the closest previous unmatched `then`”

That is, **Tree 1** is preferred.

Q: How can that be achieved?

A: Transform the grammar into an **equivalent** but **unambiguous** grammar.

Exercise: convince yourself that the following grammar indeed is equivalent!

OSICMP: Recap of Basic Formal Language Notions – p.35/52

Ambiguous Grammars (2)

- A CFL for which every CFG is ambiguous is **inherently ambiguous**.

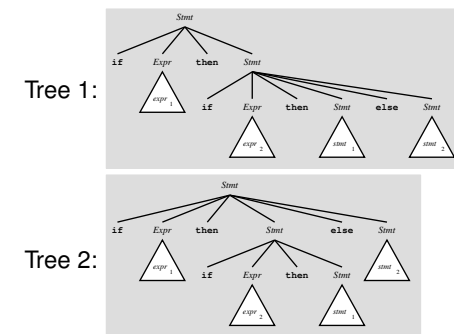
- The following language L is inherently ambiguous:

$$L = \{a^n b^n c^m d^m \mid n \geq 1, m \geq 1\} \cup \{a^n b^m c^n d^n \mid n \geq 1, m \geq 1\}$$

- Reason: All but a finite number of strings of the form $a^n b^n c^n d^n$ must be generated in two different ways. (The proof is not easy!)

OSICMP: Recap of Basic Formal Language Notions – p.30/52

Elim. Ambiguity: Dangling-Else (2)



OSICMP: Recap of Basic Formal Language Notions – p.33/52

Elim. Ambiguity: Dangling-Else (5)

Idea: a statement appearing between a `then` and an `else` must be a “matched” statement.

$$\begin{array}{l} Stmt \rightarrow MatchedStmt \\ \quad \quad | UnmatchedStmt \\ MatchedStmt \rightarrow \text{if } Expr \text{ then } MatchedStmt \\ \quad \quad \quad \text{else } MatchedStmt \\ \quad \quad \quad | \text{other} \\ UnmatchedStmt \rightarrow \text{if } Expr \text{ then } Stmt \\ \quad \quad \quad | \text{if } Expr \text{ then } MatchedStmt \\ \quad \quad \quad | \text{else } UnmatchedStmt \end{array}$$

OSICMP: Recap of Basic Formal Language Notions – p.36/52

Elim. Ambiguity: Dangling-Else (6)

Compare with the grammar for `if`-statements given in section 14.9 of the Java Language Specification, Third Edition:

<http://java.sun.com/docs/books/jls>

It uses the grammar structure of the previous slide to solve the dangling-else problem, even if the names of the non-terminals are somewhat different.

OSICMP: Recap of Basic Formal Language Notions – p.37/52

Elim. Ambiguity: Associativity (3)

If we make the **choice** of letting the parse tree structure impart the bracketing structure, we see that the two parse trees correspond to

- $(1 + 2) + 3$
- $1 + (2 + 3)$

Similarly, $47 - 3 - 2$ can be parsed in two ways:

- $(47 - 3) - 2$
- $47 - (3 - 2)$

Clearly the choice affects the of the code!

OSICMP: Recap of Basic Formal Language Notions – p.40/52

Elim. Ambiguity: Associativity (6)

To disambiguate, we want to make both $+$ and $-$ **left-associative**.

That can be achieved by making the relevant grammar productions **left-recursive**:

$$\begin{aligned} \text{Expr} &\rightarrow \text{PrimExpr} \\ &| \text{Expr} + \text{PrimExpr} \\ &| \text{Expr} - \text{PrimExpr} \\ \text{PrimExpr} &\rightarrow \text{integer} \\ &| (\text{Expr}) \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions – p.43/52

Eliminating Ambiguity: Associativity

It is standard practice to leave out unnecessary parentheses when writing down mathematical expressions:

$1 + 2 + 3$ instead of $(1 + 2) + 3$
 $47 - 3 - 2$ instead of $(47 - 3) - 2$

We would like to do the same when writing programs!

OSICMP: Recap of Basic Formal Language Notions – p.38/52

Elim. Ambiguity: Associativity (4)

- The choice might not seem important for $+$ since, mathematically, $+$ is **associative**:

$$(1 + 2) + 3 = 1 + (2 + 3) = 6$$

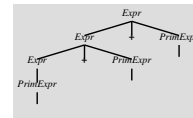
But the **computer implementation** of $+$ might not be so well-behaved!

- Floating-point addition is **not** associative!
- Integer addition is not associative if e.g. overflow is trapped.

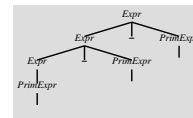
OSICMP: Recap of Basic Formal Language Notions – p.41/52

Elim. Ambiguity: Associativity (7)

Thus, $1 + 2 + 3$ is parsed as $(1 + 2) + 3$:



And $47 - 3 - 2$ is parsed as $(47 - 3) - 2$:



OSICMP: Recap of Basic Formal Language Notions – p.44/52

Elim. Ambiguity: Associativity (2)

The following grammar achieves that:

$$\begin{aligned} \text{Expr} &\rightarrow \text{integer} \\ &| \text{Expr} + \text{Expr} \\ &| \text{Expr} - \text{Expr} \\ &| (\text{Expr}) \end{aligned}$$

But ambiguous! Parse trees for $1 + 2 + 3$:



(Slightly simplified: 1, 2, etc. considered terminals.)

OSICMP: Recap of Basic Formal Language Notions – p.38/52

Elim. Ambiguity: Associativity (5)

- The choice clearly matters for $-$:

$$(47 - 3) - 2 \neq 47 - (3 - 2)$$

OSICMP: Recap of Basic Formal Language Notions – p.42/52

Elim. Ambiguity: Associativity (8)

Some operators are usually considered **right-associative**.

Consider an arithmetic exponentiation operator $^$. We would like

$$3 \wedge 2 \wedge 3$$

to be parsed as

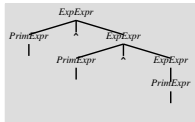
$$3 \wedge (2 \wedge 3)$$

so that the meaning is $3^{2^3} = 3^{(2^3)} = 6561$ rather than $(3^2)^3 = 729$.

OSICMP: Recap of Basic Formal Language Notions – p.45/52

Elim. Ambiguity: Associativity (9)

An operator can be made **right-associative** through **right-recursive** grammar productions:

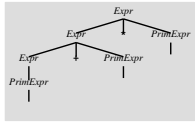
$$\begin{aligned} \text{ExpExpr} &\rightarrow \text{PrimExpr} \\ &\quad | \text{PrimExpr} \wedge \text{ExpExpr} \\ \text{PrimExpr} &\rightarrow \text{integer} \\ &\quad | (\text{Expr}) \end{aligned}$$


OSICMP: Recap of Basic Formal Language Notions - p.46/52

Eliminating Ambiguity: Precedence (3)

However, the meaning is **not** what we want!

$1 + 2 * 3$ gets parsed as $(1 + 2) * 3$:



OSICMP: Recap of Basic Formal Language Notions - p.46/52

Other ways of dealing with ambiguity

Transforming a grammar to eliminate ambiguity is not always desirable:

- Can be quite hard to do correctly.
- The transformed grammar might be less easy to understand than the original.

Parser generator tools often provide alternative disambiguation mechanisms:

- Meta-rules that favours the longest RHS among a group of conflicting productions.
- Explicit declaration of operator precedence.

OSICMP: Recap of Basic Formal Language Notions - p.52/52

Eliminating Ambiguity: Precedence (1)

We would also like to be able to rely on standard rules for **operator precedence** to make it clear what is meant.

For example, it should be possible to write

$$1 + 2 * 3$$

instead of having to write out the fully parenthesized version

$$1 + (2 * 3)$$

OSICMP: Recap of Basic Formal Language Notions - p.47/52

Eliminating Ambiguity: Precedence (4)

We rewrite the grammar so that expressions involving **high-precedence** operators only can occur as **subexpressions** of expressions involving low-precedence operators.

$$\begin{aligned} \text{Expr} &\rightarrow \text{MulExpr} \\ &\quad | \text{Expr} + \text{MulExpr} \\ \text{MulExpr} &\rightarrow \text{PrimExpr} \\ &\quad | \text{MulExpr} * \text{PrimExpr} \\ \text{PrimExpr} &\rightarrow \text{integer} \\ &\quad | (\text{Expr}) \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions - p.50/52

Eliminating Ambiguity: Precedence (2)

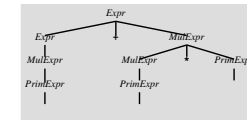
We chose to make $*$ left-associative (standard). The following grammar accepts expressions like $1 + 2 * 3$:

$$\begin{aligned} \text{Expr} &\rightarrow \text{PrimExpr} \\ &\quad | \text{Expr} + \text{PrimExpr} \\ &\quad | \text{Expr} * \text{PrimExpr} \\ \text{PrimExpr} &\rightarrow \text{integer} \\ &\quad | (\text{Expr}) \end{aligned}$$

OSICMP: Recap of Basic Formal Language Notions - p.48/52

Eliminating Ambiguity: Precedence (5)

Now $1 + 2 * 3$ gets parsed as $1 + (2 * 3)$:



OSICMP: Recap of Basic Formal Language Notions - p.51/52