CSC 3130: Automata theory and formal languages

# Context-free languages

# Context-free grammar

- This is an a different model for describing languages

- The language is specified by productions (substitution rules) that tell how strings can be obtained, e.g.

$$A \rightarrow 0A1 \qquad \text{A, B are variables}$$
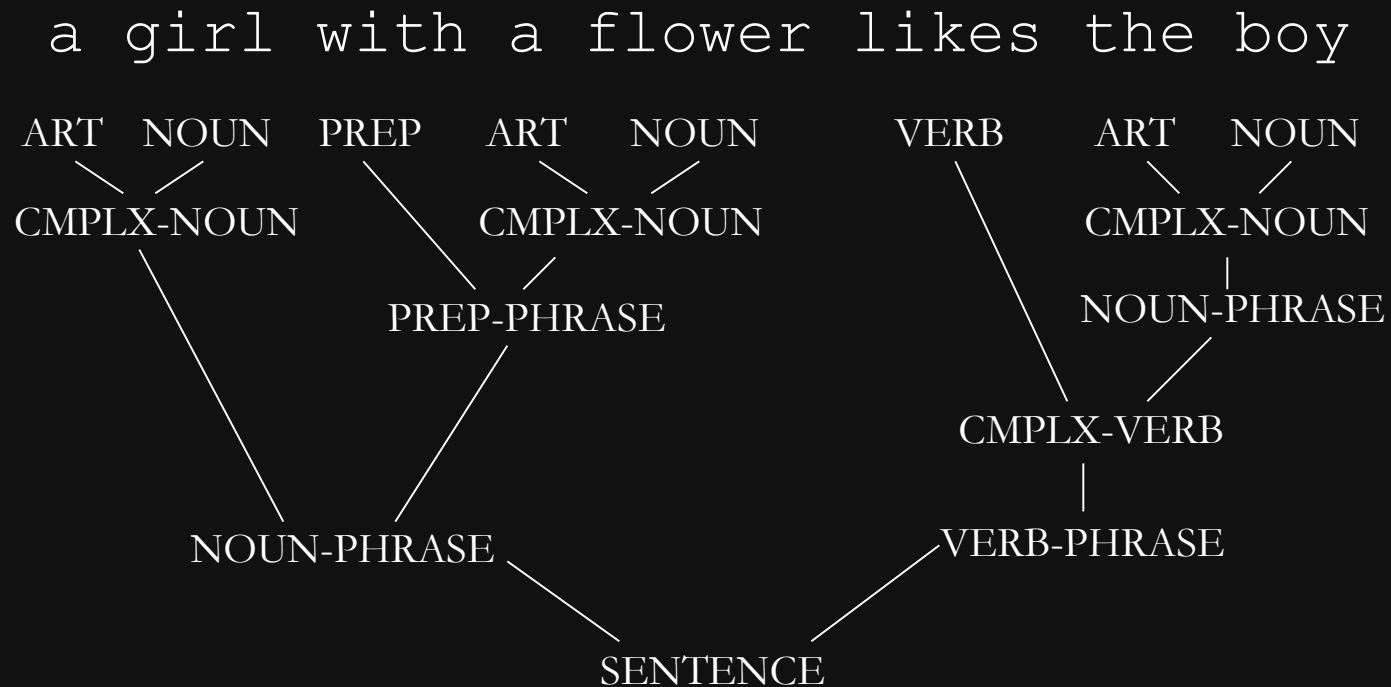$$A \rightarrow B \qquad \text{0, 1, \# are terminals}$$
$$B \rightarrow \# \qquad \text{A is the start variable}$$

- Using these rules, we can derive strings like this:

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

# Some natural examples

- Context-free grammars were first used for natural languages

# Natural languages

- We can describe (some fragments) of the English language by a context-free grammar:

SENTENCE → NOUN-PHRASE VERB-PHRASE
NOUN-PHRASE → CMPLX-NOUN
NOUN-PHRASE → CMPLX-NOUN PREP-PHRASE
VERB-PHRASE → CMPLX-VERB
VERB-PHRASE → CMPLX-VERB PREP-PHRASE
PREP-PHRASE → PREP CMPLX-NOUN
CMPLX-NOUN → ARTICLE NOUN
CMPLX-VERB → VERB NOUN-PHRASE
CMPLX-VERB → VERB

ARTICLE → a
ARTICLE → the
NOUN → boy
NOUN → girl
NOUN → flower
VERB → likes
VERB → touches
VERB → sees
PREP → with

**variables:** SENTENCE, NOUN-PHRASE, …

**terminals:** a, the, boy, girl, flower, likes, touches, sees, with

**start variable:** SENTENCE

# Programming languages

- Context-free grammars are also used to describe (parts of) programming languages

- For instance, expressions like $(2 + 3) * 5$ or $3 + 8 + 2 * 7$ can be described by the CFG

$\langle expr \rangle \rightarrow \langle expr \rangle + \langle expr \rangle$

$\langle expr \rangle \rightarrow \langle expr \rangle * \langle expr \rangle$

$\langle expr \rangle \rightarrow (\langle expr \rangle)$

$\langle expr \rangle \rightarrow 0$

$\langle expr \rangle \rightarrow 1$

$\ldots$

$\langle expr \rangle \rightarrow 9$

Variables: $\langle expr \rangle$

Terminals: $+, *, (, ), 0, 1, \ldots, 9$

# Motivation for studying CFGs

- Context-free grammars are essential for understanding the meaning of computer programs

$$\text{code: } (2 + 3) * 5$$

$$\text{meaning: "add } 2 \text{ and } 3, \text{ and then multiply by } 5\text{"}$$

- They are used in compilers

# Definition of context-free grammar

- A context-free grammar (CFG) is a 4-tuple $(V, T, P, S)$ where
  - $V$ is a finite set of variables or non-terminals
  - $T$ is a finite set of terminals ($V \cap T = \varnothing$)
  - $P$ is a set of productions or substitution rules of the form

    $$A \rightarrow \alpha$$

    where $A$ is a symbol in $V$ and $\alpha$ is a string over $V \cup T$
  - $S$ is a variable in $V$ called the start variable

# Shorthand notation for productions

- When we have multiple productions with the same variable on the left like

$$E \rightarrow E + E \qquad N \rightarrow 0N$$
$$E \rightarrow E * E \qquad N \rightarrow 1N$$
$$E \rightarrow (E) \qquad N \rightarrow 0$$
$$E \rightarrow N \qquad N \rightarrow 1$$

Variables: $E$, $N$

Terminals: $+$, $*$, $($, $)$, $0$, $1$

Start variable: $E$

we can write this in shorthand as

$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid 1$$
$$N \rightarrow 0N \mid 1N \mid 0 \mid 1$$

# Derivation

- A derivation is a sequential application of productions:

$$E \Rightarrow E * E$$
$$\Rightarrow (E) * E$$
$$\Rightarrow (E) * N$$
$$\Rightarrow (E + E) * N$$
$$\Rightarrow (E + E) * 1$$
$$\Rightarrow (E + N) * 1$$
$$\Rightarrow (N + N) * 1$$
$$\Rightarrow (N + 1N) * 1$$
$$\Rightarrow (N + 10) * 1$$
$$\Rightarrow (1 + 10) * 1$$

derivation

$$\alpha \Rightarrow \beta$$

means $\beta$ can be obtained from $\alpha$ with one production

$$\alpha \overset{*}{\Rightarrow} \beta$$

means $\beta$ can be obtained from $\alpha$ after zero or more productions

# Language of a CFG

- The language of a CFG $(V, T, P, S)$ is the set of all strings containing only terminals that can be derived from the start variable $S$

$$L = \{\omega \mid \omega \in T^* \text{ and } S \overset{*}{\Rightarrow} \omega \}$$

- This is a language over the alphabet $T$

- A language $L$ is context-free if it is the language of some CFG

# Example 1

$A \to 0A1 \mid B$

$B \to \#$

variables: $A, B$
terminals: $0, 1, \#$
start variable: $A$

- Is the string $00\#11$ in L?

- How about $00\#111, 00\#0\#1\#11$?

- What is the language of this CFG?

$$L = \{0^n\#1^n : n \geq 0\}$$

# Example 2

$S \rightarrow SS \mid (S) \mid \varepsilon$

convention: variables in uppercase,
terminals in lowercase, start variable first

- Give derivations of $()$, $(()())$

| | | |
|---|---|---|
| $S$ | $\Rightarrow (S)$ | (rule 2) |
| | $\Rightarrow ()$ | (rule 3) |

| | | |
|---|---|---|
| $S$ | $\Rightarrow (S)$ | (rule 2) |
| | $\Rightarrow (SS)$ | (rule 1) |
| | $\Rightarrow ((S)S)$ | (rule 2) |
| | $\Rightarrow ((S)(S))$ | (rule 2) |
| | $\Rightarrow (()(S))$ | (rule 3) |
| | $\Rightarrow (()())$ | (rule 3) |

- How about $())$?

# Examples: Designing CFGs

- Write a CFG for the following languages
  - Linear equations over $x$, $y$, $z$, like:
    $x + 5y - z = 9$
    $11x - y = 2$

  - Numbers without leading zeros, e.g., $109$, $0$ but not $019$

  - The language $L = \{a^n b^n c^m d^m \mid n \geq 0, m \geq 0\}$

  - The language $L = \{a^n b^m c^m d^n \mid n \geq 0, m \geq 0\}$

# Context-free versus regular

- Write a CFG for the language $(0 + 1)*111$

$$S \rightarrow A111$$
$$A \rightarrow \varepsilon \mid 0A \mid 1A$$

- Can you do so for every regular language?

Every regular language is context-free

- Proof:

# From regular to context-free

| regular expression ⇨ | CFG |
|---|---|
| $\varnothing$ | grammar with no rules |
| $\varepsilon$ | $S \rightarrow \varepsilon$ |
| $a$ (alphabet symbol) | $S \rightarrow a$ |
| $E_1 + E_2$ | $S \rightarrow S_1 \mid S_2$ |
| $E_1 E_2$ | $S \rightarrow S_1 S_2$ |
| $E_1{}^*$ | $S \rightarrow S S_1 \mid \varepsilon$ |

In all cases, $S$ becomes the new start symbol

# Context-free versus regular

- Is every context-free language regular?

- No! We already saw some examples:

$$A \rightarrow 0A1 \mid B$$
$$B \rightarrow \#$$

$$L = \{0^n\#1^n : n \geq 0\}$$

- This language is context-free but not regular

# Parse tree

- Derivations can also be represented using parse trees

$$E \rightarrow E + E \mid E - E \mid (E) \mid V$$
$$V \rightarrow x \mid y \mid z$$

$E \Rightarrow E + E$
$\quad \Rightarrow V + E$
$\quad \Rightarrow x + E$
$\quad \Rightarrow x + (E)$
$\quad \Rightarrow x + (E - E)$
$\quad \Rightarrow x + (V - E)$
$\quad \Rightarrow x + (y - E)$
$\quad \Rightarrow x + (y - V)$
$\quad \Rightarrow x + (y - z)$

# Definition of parse tree

- A parse tree for a CFG $G$ is an ordered tree with labels on the nodes such that
  - Every internal node is labeled by a variable
  - Every leaf is labeled by a terminal or $\varepsilon$
  - Leaves labeled by $\varepsilon$ have no siblings
  - If a node is labeled $A$ and has children $A_1, \ldots, A_k$ from left to right, then the rule

  $$A \rightarrow A_1 \ldots A_k$$

  is a production in $G$.

# Left derivation

- Always derive the leftmost variable first:

$$E \Rightarrow E + E$$
$$\Rightarrow V + E$$
$$\Rightarrow x + E$$
$$\Rightarrow x + (E)$$
$$\Rightarrow x + (E - E)$$
$$\Rightarrow x + (V - E)$$
$$\Rightarrow x + (y - E)$$
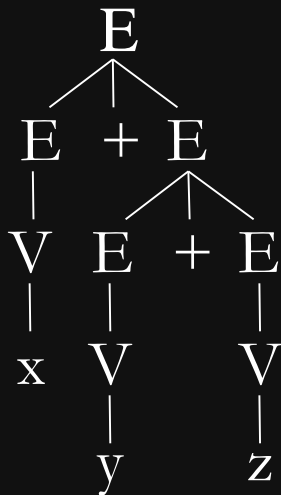$$\Rightarrow x + (y - V)$$
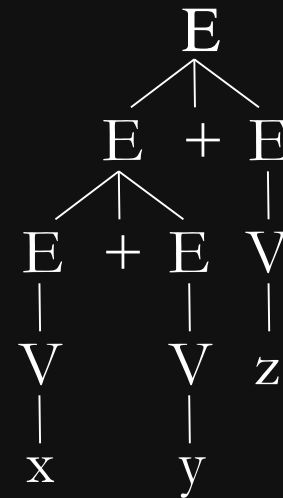$$\Rightarrow x + (y - z)$$



- Corresponds to a left-to-right traversal of parse tree

# Ambiguity

- A grammar is ambiguous if some strings have more than one parse tree

- Example:

$$E \rightarrow E + E \mid E - E \mid (E) \mid V$$
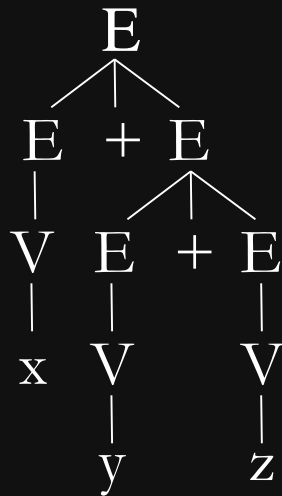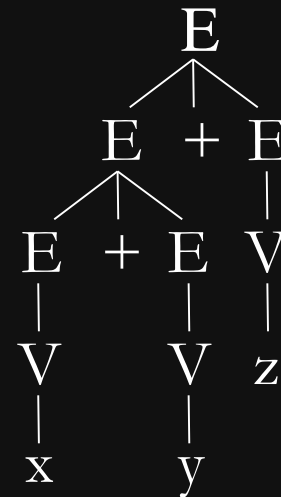$$V \rightarrow x \mid y \mid z$$

$$x + y + z$$

# Why ambiguity matters

- The parse tree represents the intended meaning:

$$x + y + z$$

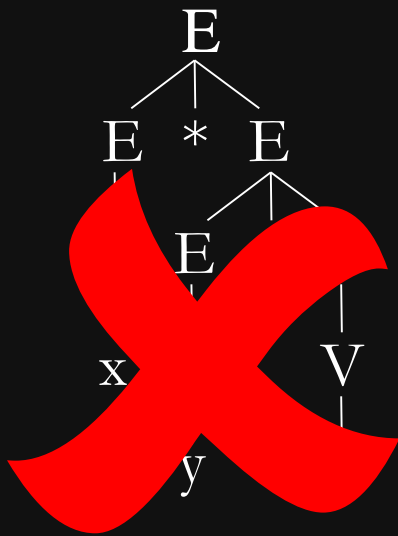"first add $y$ and $z$,
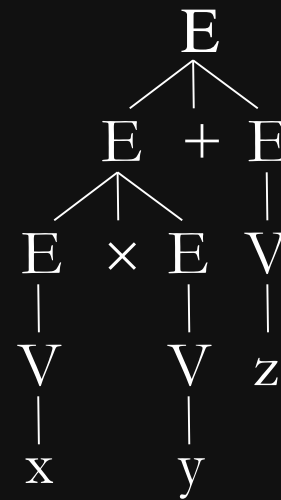and then add this to $x$"

"first add $x$ and $y$,
and then add $z$ to this"

# Why ambiguity matters

- Suppose we also had multiplication:

$$E \rightarrow E + E \mid E - E \mid E \times E \mid (E) \mid V$$
$$V \rightarrow x \mid y \mid z$$



$$x \times y + z$$

"first $y + z$, then $x \times$"          "first $x \times y$, then $+ z$"

# Disambiguation

- Sometimes we can rewrite the grammar to remove the ambiguity

$$E \rightarrow E + E \mid E - E \mid E \times E \mid (E) \mid V$$
$$V \rightarrow x \mid y \mid z$$

- Rewrite grammar so $\times$ cannot be broken by $+$:

$$E \rightarrow T \mid E + T \mid E - T$$
$$T \rightarrow F \mid T \times F$$
$$F \rightarrow (E) \mid V$$
$$V \rightarrow x \mid y \mid z$$

T stands for term: $x * (y + z)$
F stands for factor: $x, (y + z)$

A term always splits into factors

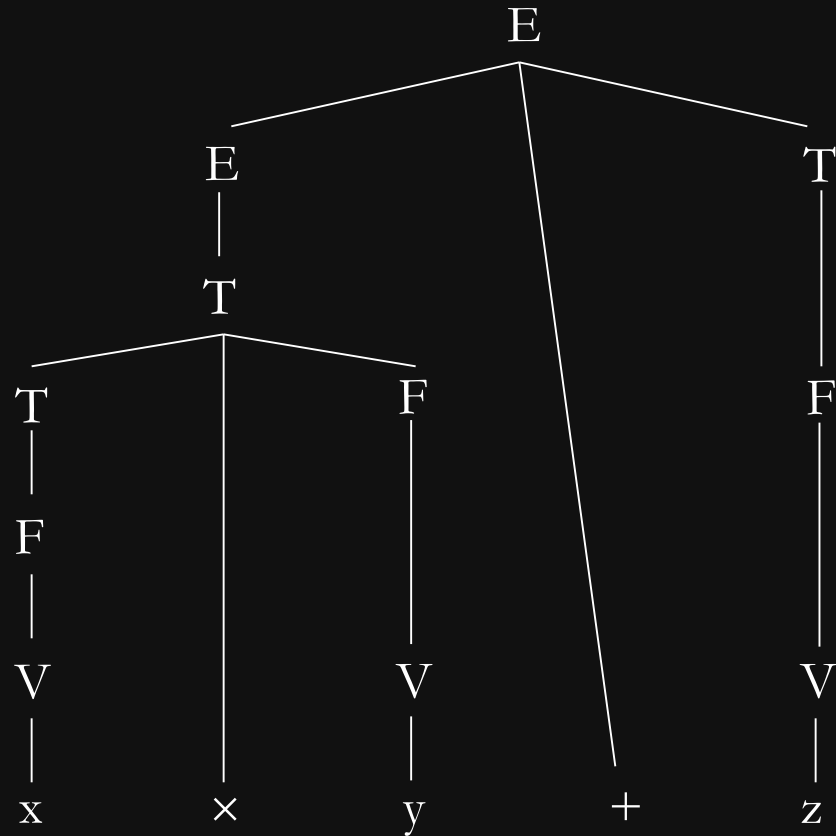A factor is either a variable or a parenthesized expression

# Disambiguation

- Example

$$E \rightarrow T \mid E + T \mid E - T$$
$$T \rightarrow F \mid T \times F$$
$$F \rightarrow (E) \mid V$$
$$V \rightarrow x \mid y \mid z$$

# Disambiguation

- Can we always disambiguate a grammar?

- No, for two reasons
  - There exists an inherently ambiguous context-free $L$: Every CFG for this language is ambiguous
  - There is no general procedure that can tell if a grammar is ambiguous

- However, grammars used in programming languages can typically be disambiguated

# Another Example

$$S \rightarrow aB \mid bA$$
$$A \rightarrow a \mid aS \mid bAA$$
$$B \rightarrow b \mid bS \mid aBB$$

- Is $ab, baba, abbbaa$ in $L$?

- How about $a, bba$?


- What is the language of this CFG?

- Is the CFG ambiguous?