

Context-free languages

Stefan D. Bruda

CS 310, Winter 2025

CONTEXT-FREE GRAMMARS



- Recall that languages defined using only union, concatenation, and recursion from symbols from Σ , ε , and \emptyset are called **context-free languages**
 - A language description as above is called a **context-free grammar**
 - Examples: $B = \varepsilon + 0B1$, $L = \varepsilon + 1 + 0 + 0L0 + 1L1$
 - The customary way of writing down a context-free grammar is using a **definition by cases** rather than the recursive equation (the **Backus-Naur form** or **BNF**)

$$\begin{array}{ll} B \rightarrow \varepsilon & \langle \text{balanced} \rangle ::= \\ B \rightarrow 0B1 & \langle \text{balanced} \rangle ::= 0 \langle \text{balanced} \rangle 1 \\ \\ L \rightarrow \varepsilon & \langle \text{palindrome} \rangle ::= \\ L \rightarrow 0 & \langle \text{palindrome} \rangle ::= 0 \\ L \rightarrow 1 & \langle \text{palindrome} \rangle ::= 1 \\ L \rightarrow 0L0 & \langle \text{palindrome} \rangle ::= 0 \langle \text{palindrome} \rangle 0 \\ L \rightarrow 1L1 & \langle \text{palindrome} \rangle ::= 1 \langle \text{palindrome} \rangle 1 \\ \\ \langle \text{palindrome} \rangle ::= \varepsilon \mid 0 \mid 1 \mid 0 \langle \text{palindrome} \rangle 0 \mid 1 \langle \text{palindrome} \rangle 1 \end{array}$$



- Formally a **context-free grammar** is a tuple $G = (N, \Sigma, R, S)$, where
 - Σ is an alphabet of **terminals**
 - N alphabet of symbols called by contrast **nonterminals**
 - Traditionally nonterminals are capitalized or surrounded by \langle and \rangle , everything else being a terminal
 - $S \in N$ is the **axiom** (or the **start symbol**)
 - $R \subseteq N \times (N + \Sigma)^*$ is the set of **(rewriting) rules** or **productions**
 - Common ways of expressing $(\alpha, \beta) \in R$: $\alpha \rightarrow \beta$ or $\alpha ::= \beta$
- Further examples:

$\begin{aligned} \langle \text{exp} \rangle &::= \langle \text{const} \rangle \\ &\quad \langle \text{var} \rangle \\ &\quad \langle \text{exp} \rangle \langle \text{op} \rangle \langle \text{exp} \rangle \\ &\quad (\langle \text{exp} \rangle) \\ \langle \text{op} \rangle &::= + \mid - \mid * \mid / \end{aligned}$	$\begin{aligned} \langle \text{stmt} \rangle &::= ; \\ &\quad \langle \text{var} \rangle = \langle \text{exp} \rangle ; \\ &\quad \text{if } (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ &\quad \text{while } (\langle \text{exp} \rangle) \langle \text{stmt} \rangle \\ &\quad \{ \langle \text{seq} \rangle \} \\ \langle \text{seq} \rangle &::= \varepsilon \mid \langle \text{stmt} \rangle \langle \text{seq} \rangle \end{aligned}$
---	---

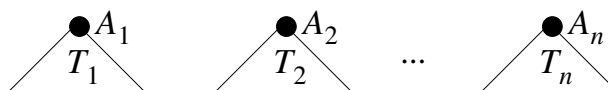
DERIVATIONS



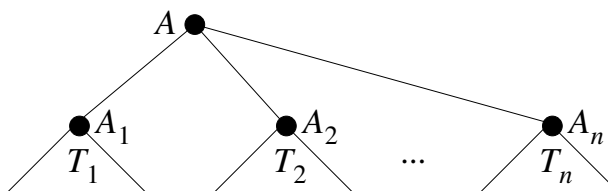
- $G = (N, \Sigma, R, S)$
 - A rewriting rule $A \rightarrow v' \in R$ is used to rewrite its left-hand side (A) into its right-hand side (v'):
 - $u \Rightarrow v$ iff $\exists x, y \in (N + \Sigma)^* : \exists A \in N : u = xAy, v = xv'y, A \rightarrow v' \in R$
 - Rewriting can be chained (\Rightarrow^* , the reflexive and transitive closure of \Rightarrow = **derivation**)
 - $s \Rightarrow^* s'$ iff $s = s'$, $s \Rightarrow s'$, or there exist strings s_1, s_2, \dots, s_n such that $s \Rightarrow s_1 \Rightarrow s_2 \Rightarrow \dots \Rightarrow s_n \Rightarrow s'$
 - $\langle \text{pal} \rangle \Rightarrow 0\langle \text{pal} \rangle 0 \Rightarrow 01\langle \text{pal} \rangle 10 \Rightarrow 010\langle \text{pal} \rangle 010 \Rightarrow 0101010$
- $$\langle \text{pal} \rangle ::= \varepsilon \mid 0 \mid 1 \mid 0 \langle \text{pal} \rangle 0 \mid 1 \langle \text{pal} \rangle 1$$
- The language generated by grammar G : exactly all the **terminal** strings generated from S : $\mathcal{L}(G) = \{w \in \Sigma^* : S \Rightarrow^* w\}$
 - Same language as defined by the respective recursive equation (and the approximation scheme shown earlier)

- Definition:

- 1 For every $a \in N + \Sigma$ the following is a parse tree (with yield a): $\bullet a$
- 2 For every $A \rightarrow \varepsilon \in R$ the following is a parse tree (with yield ε): $\bullet A$
 $\bullet \varepsilon$
- 3 If the following are parse trees (with yields y_1, y_2, \dots, y_n , respectively):



and $A \rightarrow A_1 A_2 \dots A_n \in R$, then the following is a parse tree (with yield $y_1 y_2 \dots y_n$):



- Yield: concatenation of leaves in inorder

DERIVATIONS AND PARSE TREES

- Every derivation starting from some nonterminal has an associated parse tree (rooted at the starting nonterminal)
- Two derivations are **similar** iff only the order of rule application varies = can obtain one derivation from the other by repeatedly flipping **consecutive** rule applications
 - Two similar derivations have identical parse trees
 - Can use a “standard” derivation: leftmost ($A \xRightarrow{L}^* w$) or rightmost ($A \xRightarrow{R}^* w$)

Theorem

The following statements are equivalent:

- there exists a parse tree with root A and yield w
- $A \Rightarrow^* w$
- $A \xRightarrow{L}^* w$
- $A \xRightarrow{R}^* w$

- Ambiguity** of a grammar: there exists a string that has two derivations that are not similar (i.e., two derivations with different parse trees)
 - Can be **inherent** or not — impossible to determine algorithmically



- Regular grammar: $G = (N, \Sigma, R, S)$ with N, Σ, S as before, and $R \subseteq N \times (\epsilon + \Sigma + \Sigma N)$
 - Special form of context-free grammar (only rules of form $A \rightarrow \epsilon$, $A \rightarrow a$, and $A \rightarrow aB$ allowed)

Theorem

Exactly all the regular languages are generated by regular grammars

- Let $M = (K, \Sigma, \Delta, s, F)$ be some finite automaton
- We construct the grammar $G = (K, \Sigma, s, R)$ with

$$R = \{q \rightarrow ap : (q, a, p) \in \Delta\} + \{q \rightarrow \epsilon : q \in F\}$$

Corollary

All regular languages are context-free

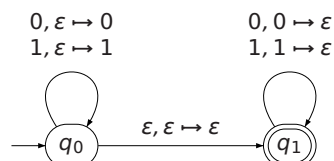
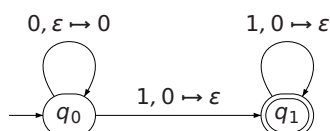
- However, there are more context-free than regular languages

$$S \rightarrow aSb \quad S \rightarrow \epsilon$$

PUSH-DOWN AUTOMATA



- **Push-down automaton**: finite automaton + “push-down store” (or stack)
- $M = (K, \Sigma, \Gamma, \Delta, s, F)$
 - K, Σ, s, F as before (for finite automata)
 - Γ is the stack alphabet
 - $\Delta \subseteq \{(K \times (\Sigma + \{\epsilon\}) \times \Gamma^*) \times (K \times \Gamma^*)\}$
 - Transition: $((q, a, \gamma), (q', \gamma'))$ with a the current input symbol (or ϵ), γ the old stack top, and γ' the replacement top
 - Graphical representation: $\xrightarrow{a, \gamma \mapsto \gamma'} q \rightarrow q'$
- Acceptance: there exists a path (or run) that spells the input and ends in a final state **and** the stack is empty at the beginning as well as at the end of the path





Theorem

Theorem: Push-down automata accept exactly all the context-free languages

- \supseteq : Given the grammar $G = (N, \Sigma, S, R)$ construct the push-down automaton $M = (K, \Sigma, \Gamma, \Delta, s, F)$ such that:

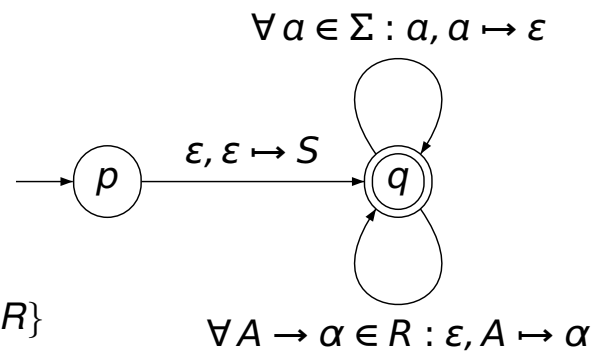
$$\Gamma = N + \Sigma$$

$$K = \{p, q\}$$

$$s = p$$

$$F = \{q\}$$

$$\begin{aligned} \Delta = & \{((p, \varepsilon, \varepsilon), (q, S))\} \\ & + \{((q, \varepsilon, A), (q, \alpha)) : A \rightarrow \alpha \in R\} \\ & + \{((q, a, a), (q, \varepsilon)) : a \in \Sigma\} \end{aligned}$$



- \subseteq : Proof left for a full course in formal languages

CLOSURE PROPERTIES

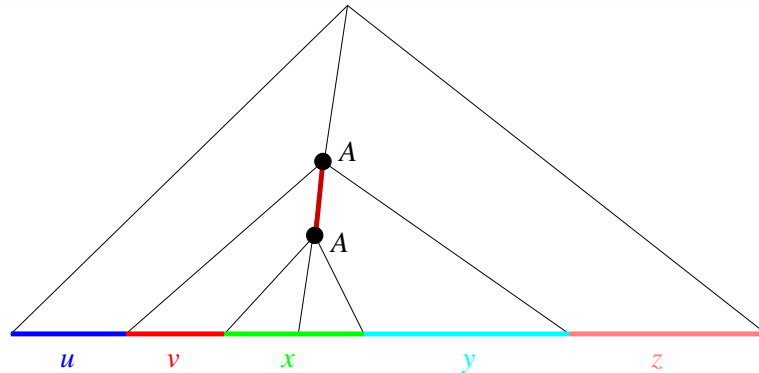


- Consider two grammars with axioms S_1 and S_2 ; construct a grammar with axiom S
- Context-free languages are closed under
 - Union: Add rules $S \rightarrow S_1$ and $S \rightarrow S_2$
 - Concatenation: Add rule $S \rightarrow S_1 S_2$
 - Kleene star: Add rules $S \rightarrow \varepsilon$ and $S \rightarrow SS_1$
 - Intersection and complement:
 - If context-free languages are not closed under one of these then they are not closed under the other either

- Let $\Phi(G)$ be the maximum fanout (branching factor) of any node in any parse tree constructed based on grammar G
- A parse tree of height h has a yield of size no more than $\Phi(G)^h$

Theorem (Pumping context-free languages)

Let $n = \Phi(G)^{|N|}$. For any $w \in \mathcal{L}(G)$ such that $|w| \geq n$ we can write w as $uvxyz$ such that $vy \neq \epsilon$, $|vxy| \leq n$ and $uv^i xy^i z \in \mathcal{L}(G)$ for any $i \geq 0$



PUMPING CONTEXT-FREE LANGUAGES (CONT'D)

- Some interesting non-context-free languages:
 - $\{a^n b^n c^n : n \geq 0\}$
 - $\{a^n : n \text{ is prime}\}$
 - $\{w \in \{a, b, c\}^* : |w|_a = |w|_b = |w|_c\}$

Corollary

Context-free languages are **not** closed under intersection and complement

- Indeed, $\{a^n b^n c^n : n \geq 0\} = \{a^n b^n c^m : n, m \geq 0\} \cap \{a^m b^n c^n : n, m \geq 0\}$
- That $\{a^n b^n c^m : n, m \geq 0\}$ is context free can be shown by constructing a grammar/automaton or by using closure properties
- Tricky language: $\{w \in \{a, b, c\}^* : |w|_a = |w|_b = |w|_c\}$ satisfies the pumping lemma yet is not context-free