



计算机科学与技术学院
College of Computer Science and Technology

Chapter 3 Context-Free Language



Elements Of The Theory Of Computation
Zhejiang University/CS/Course/Xiaogang Jin
E-Mail: xiaogangj@cise.zju.edu.cn



□ So-far: regular languages

- DFA = NFA (language recognizer)
- Regular expressions (language generator)
- Many languages are not regular
 - *Balanced parentheses*
 - *Arithmetic expressions*

□ Next: context-free languages

- PDA = CFG
- Add LIFO (stack) memory

3.1 Context-Free Grammars

Example:

Consider the regular expression $a(a^* \cup b^*)b$.

- Regular expressions can be viewed as language generators.



- The language denoted by $a(a^* \cup b^*)b$ can be defined by the following generator:

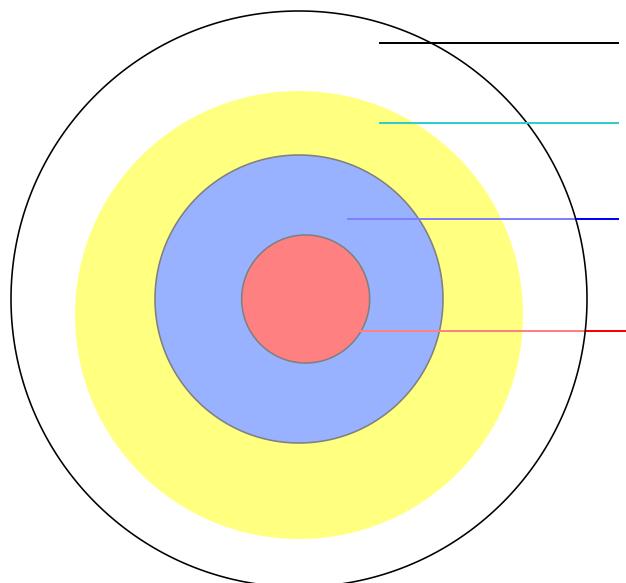
$$\boxed{S \rightarrow aMb, M \rightarrow A, M \rightarrow B, \\ A \rightarrow aA, A \rightarrow e, B \rightarrow bB, B \rightarrow e.}$$



Question: What is the context-free?

(as opposed to **Context-Sensitive Grammar**, CSG.)

Chomsky Hierarchy



Unrestricted Grammar

CSG

CFG

Regular Grammar



Definition: A **context-free grammar (CFG)** is a quadruple $G = (V, \Sigma, R, S)$, where

V is an alphabet;

$\Sigma \subseteq V$ is the set of terminal symbols;

$S \in V - \Sigma$ is the start symbol; and

R is the set of rules, a finite subset of $(V - \Sigma) \times V^*$.

Remark:

- The member of $V - \Sigma$ are called **nonterminals**.

For any $A \in V - \Sigma$ and $u \in V^*$,

$$A \rightarrow_G u \Leftrightarrow (A, u) \in R.$$



- For any strings $u, v \in V^*$,

$u \Rightarrow_G v \Leftrightarrow \exists x, y \in V^*, \text{ and } A \in V - \Sigma, \text{ such that}$

$$u = xAy, v = xv'y, \text{ and } A \rightarrow_G v'.$$

- \Rightarrow_G^* is the reflexive, transitive closure of \Rightarrow_G .
- $w_0 \Rightarrow_G w_1 \Rightarrow_G \cdots \Rightarrow_G w_n$
 - a **derivation** in G of w_n from w_0 . n be the length of the derivation.
- The **language generated by G** ,

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

L is a **context-free language(CFL)** $\Leftrightarrow \exists$ a context-free grammar(CFG) G , such that $L = L(G)$.



Example: Consider the CFG $G = (V, \Sigma, R, S)$ where $V = \{S, a, b\}$, $\Sigma = \{a, b\}$, $R = \{S \rightarrow aSb, S \rightarrow e\}$.

$$L(G) = \{a^n b^n : n \geq 0\}$$

*L(G) is context-free
but not regular*

Example: Let $G = (V, \Sigma, R, S)$ where $V = \{S, (,)\}$, $\Sigma = \{(,)\}$, $R = \{S \rightarrow e, S \rightarrow SS, S \rightarrow (S)\}$.

L(G) is the language containing all strings of balanced parentheses.



Example: Let $G = (\{S, a, b\}, \{a, b\}, R, S)$, where, $R = \{S \rightarrow e, S \rightarrow SS, S \rightarrow aSb, S \rightarrow bSa\}$.

$L(G) = \{w \in \{a, b\}^*: w \text{ has the same number of } a's \text{ and } b's\}$

Proof:

$w \in L(G) \Rightarrow w \text{ has the same number of } a's \text{ and } b's.$

By Induction on length k of derivation.

(a) $k = 1$

The derivation is $S \Rightarrow e$

$\Rightarrow w = e$ has the same number of a 's and b 's. ✓



(b) $k > 1$

Then either:

$$S \Rightarrow SS \Rightarrow^* xy = w$$

$$S \Rightarrow aSb \Rightarrow^* axb = w$$

$$S \Rightarrow bSa \Rightarrow^* bxa = w$$

Since $S \Rightarrow^* x, S \Rightarrow^* y$ by derivations of length $< k$
 x, y have equal number of a 's and b 's(IH)

so do xy, axb , and bxa . ✓



w has the same number of a 's and b 's $\Rightarrow w \in L(G)$.

By induction on $|w|$.

(a) $|w| = 0$

$$w = e \in L(G)$$

$S \rightarrow e$ is a rule.

(b) $|w| = k + 2$: ($|w|$ must be even)

4 subcase depending on first and last symbols of w :

Case 1 $w = axb$ for some $x \in \Sigma^*$

$$|x| = k$$

$$\Rightarrow S \Rightarrow^* x$$

x has the same number of a 's and b 's.

$$S \Rightarrow aSb \Rightarrow^* axb = w$$



Case 2 $w = bxa$ for some $x \in \Sigma^*$ - similar

Case 3 $w = axa$ for some $x \in \Sigma^*$

$|x| = k$, x has 2 more b 's than a 's.

$x = uv$ for some u and v such that

- u has one more b than a .
- v has one more b than a .

$S \Rightarrow^* au$ and $S \Rightarrow^* va$

So $S \Rightarrow^* SS \Rightarrow^* auva = w$.

Case 4 $w = bxb$ for some $x \in \Sigma^*$ - similar



Example: Consider the CFG $G = (V, \Sigma, R, S)$ where

- $V = \{+, *, (,), id, T, F, E\}$
- $\Sigma = \{+, *, (,), id\}$
- $R = \{E \rightarrow E + T, E \rightarrow T, T \rightarrow TF, T \rightarrow F, F \rightarrow (E), F \rightarrow id\}$

E: expression, T: term, F: factor.

Remark:

- Computer programs written in any programming language must be syntactically correct and enable to mechanical interpretation.
- The syntax of most programming language can be captured by CFG.



Example All regular languages are CFL.

Proof: We will encounter several proof of this fact.

- In **section 3.3:**

The CFL are precisely the languages accepted by Push-down Automata, which is a generalization of the FA.

- In **section 3.5:**

The class of CFL is closed under union, concatenation, and Kleene star.

The trivial languages \emptyset ; and $\{a\}$ are definitely context-free.



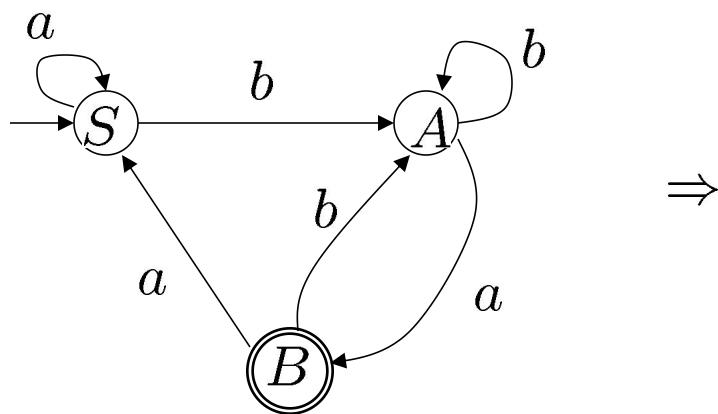
- Now we show that all regular languages are CFL by a **direct construction**.

Consider the regular language accepted by the DFA $M = (K, \Sigma, \delta, s, F)$.

To build a CFG from DFA M



Example:



\Rightarrow

$$\begin{aligned}S &\rightarrow aS, \\S &\rightarrow bA, \\A &\rightarrow aB, \\A &\rightarrow bA, \\B &\rightarrow aS, \\B &\rightarrow bA, \\B &\rightarrow e.\end{aligned}$$

Consider string $aabbaba$:

$$\begin{aligned}S aabbaba &\vdash_M aSabbaba \vdash_M aaSbbaba \vdash_M aabAbaba \vdash_M aabbAaba \\&\vdash_M aabbaBba \vdash_M aabbabAa \vdash_M aabbabaB \\S \Rightarrow_G aS &\Rightarrow_G aaS \Rightarrow_G aabA \Rightarrow_G aabbA \Rightarrow_G aabbaBba \Rightarrow_G \\aabbabAa &\Rightarrow_G aabbabaA \Rightarrow_G aabbaba\end{aligned}$$



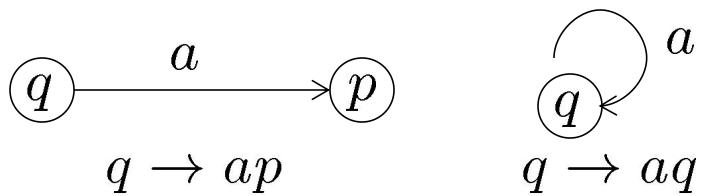
Consider the regular language accepted by the DFA $M = (K, \Sigma, \delta, s, F)$.

The same language is generated by the CFG $G = (V, \Sigma, R, S)$ where,

$$V = K \cup \Sigma$$

$$S = s$$

$$\begin{aligned} R = & \{q \rightarrow ap : \delta(q, a) = p\} \\ & \cup \{q \rightarrow e : q \in F\}. \end{aligned}$$



3.2 Parse Tree

Example: If G is the CFG that generates language of balanced parentheses.

Then string $((())()$ can be derived from S by at least two distinct derivations, namely,

$$\begin{array}{l} S \rightarrow e \\ S \rightarrow SS \\ S \rightarrow (S) \end{array}$$

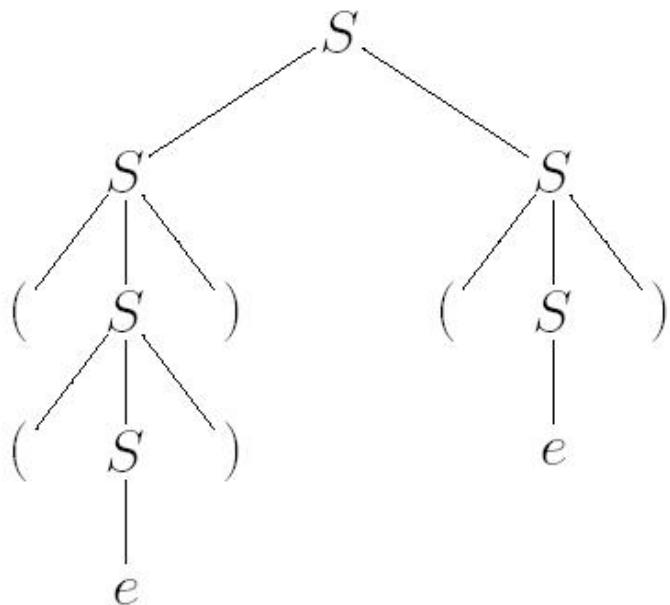
 Same!

$$S \Rightarrow \underline{SS} \Rightarrow (\underline{S})S \Rightarrow ((\underline{S}))S \Rightarrow ((())\underline{S} \Rightarrow ((())(S) \Rightarrow ((())()$$

$$S \Rightarrow S\underline{S} \Rightarrow \underline{S}(S) \Rightarrow (S)(\underline{S}) \Rightarrow (\underline{S})() \Rightarrow ((\underline{S}))() \Rightarrow ((())()$$



Both derivations can be pictured as following:



- Node: has a label in V
- Root: Start symbol
- Leaves: labeled by terminals
- The string
 - *by concatenating the labels of leaves from left to right is called the yield of the parse tree.*

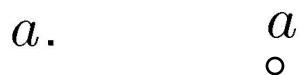
— Parse tree



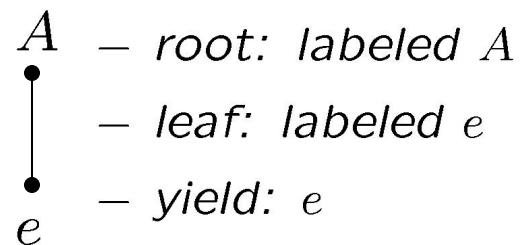
□ Parse tree: formal definition

For an arbitrary CFG $G = (V, \Sigma, R, S)$, we define its parse tree as following:

- This is the parse tree for each $a \in \Sigma$. The single node is both the root and a leaf. The yield of this parse tree is a .

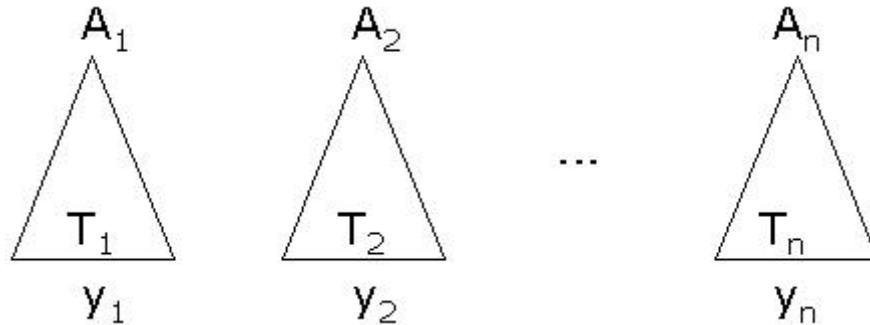


- If $A \rightarrow e$ is a rule in R ,
then is a parse tree.

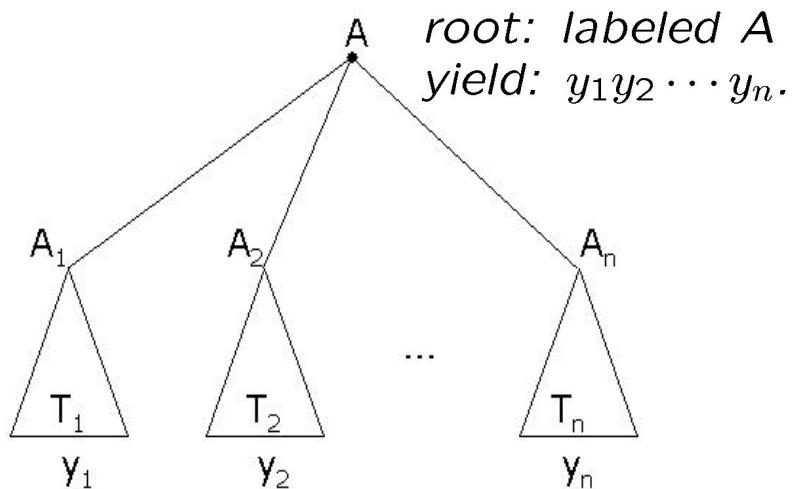




- If



are parse trees, where $n \geq 1$,
with roots labeled A_1, \dots, A_n ,
respectively, and yields
 y_1, \dots, y_n , and $A \rightarrow A_1 \dots A_n$
is a rule in R , then



is a parse tree.



- Nothing else is a parse tree.

Remark:

Parse tree

— Equivalence classes of derivation



□ **Derivations D and D' are similar:** formal definition

Let $G = (V, \Sigma, R, S)$ be a CFG,

$$D = x_1 \Rightarrow x_2 \Rightarrow \cdots \Rightarrow x_n$$

$$D' = x'_1 \Rightarrow x'_2 \Rightarrow \cdots \Rightarrow x'_n$$

be two derivations, where $x_i, x'_i \in V^*$ for $i = 1 \dots n$, $x_1, x'_1 \in V - \Sigma$, and $x_n, x'_n \in \Sigma^*$.

1) D **precedes** D' ($D \prec D'$) $\Leftrightarrow \exists 1 < k < n$, such that

- for all $i \neq k$, we have $x_i = x'_i$.
- $x_{k-1} = x'_{k-1} = uAvBw$, where $u, v, w \in V^*$, and $A, B \in V - \Sigma$.



- $x_k = uyvBw$, where $A \rightarrow y \in R$.
- $x'_k = uAvzw$, where $B \rightarrow z \in R$.
- $x_{k+1} = x'_{k+1} = uyvzw$.

2) D and D' are **similar** $\Leftrightarrow (D, D')$, belong in the reflexive, symmetric, transitive closure of \prec .

Remark:

- Similarity is an equivalence relation.
- Derivations in the same equivalence class under similarity have the same parse tree.



- Each parse tree contains a derivation that is maximal under \prec .
 - leftmost derivation (similarly, rightmost derivation)

Theorem: Let $G = (V, \Sigma, R, S)$ be a CFG, and let $A \in V - \Sigma$, and $w \in \Sigma^*$. Then the following statements are equivalent:

- (a) $A \Rightarrow^* w$.
 - (b) \exists Parse tree with root A and yield w .
 - (c) \exists a leftmost derivation $A \xrightarrow{L^*} w$.
 - (d) \exists a rightmost derivation $A \xrightarrow{R^*} w$.
-



Example: If G is the CFG that generates language of balanced parentheses. Consider the following derivations D_1, D_2 and D_3 :

$$D_1 = S \Rightarrow SS \Rightarrow (S)S \Rightarrow ((S))S \Rightarrow ((())S) \Rightarrow ((())(S) \Rightarrow ((())()$$

$$D_2 = S \Rightarrow SS \Rightarrow (S)S \Rightarrow ((S))S \Rightarrow ((S))(S) \Rightarrow ((())(S) \Rightarrow ((())()$$

$$D_3 = S \Rightarrow SS \Rightarrow (S)S \Rightarrow ((S))S \Rightarrow ((S))(S) \Rightarrow ((S))() \Rightarrow ((())()$$

$D_1 \prec D_2, D_2 \prec D_3$

but not the case $D_1 \prec D_3$

D_1, D_2 , and D_3
are similar.

$$\begin{aligned} D = S &\Rightarrow SS \Rightarrow SSS \Rightarrow S(S)S \Rightarrow S((S))S \Rightarrow S((())S) \\ &\Rightarrow S((())()) \Rightarrow ((())()) \end{aligned}$$

— not similar to the above



Ambiguity

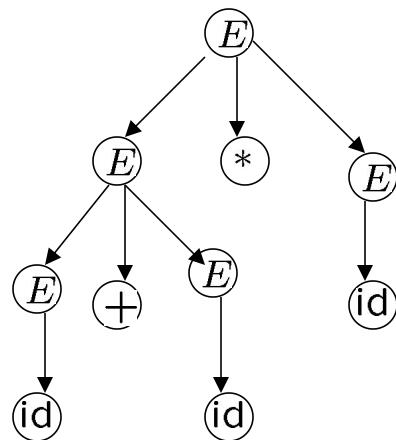
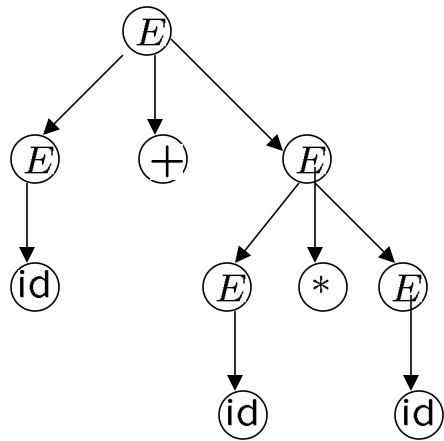
- A leftmost-derivation is a derivation in which a production is always applied to the leftmost symbol.
- In general, a string may have multiple left and rightmost derivations.

Definition: A grammar in which some word has two parse trees is said to be **ambiguous**.



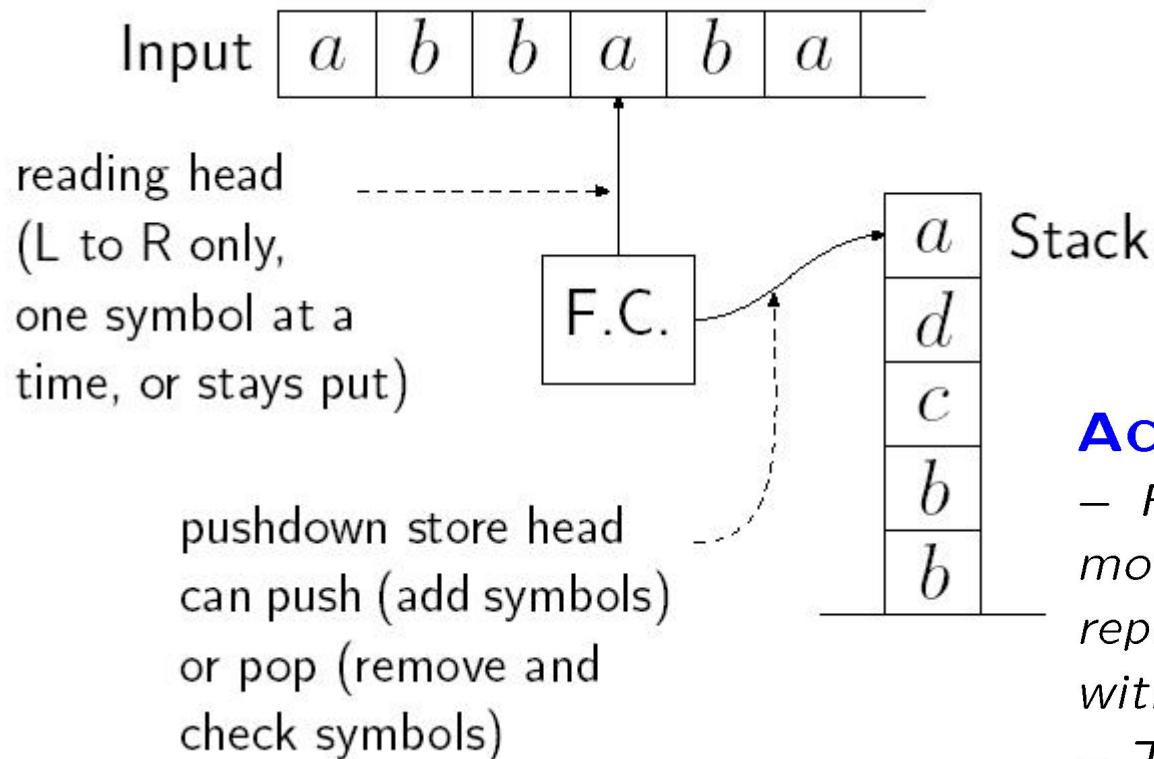
Example: Consider the CFG $G' = (V, \Sigma, R, S)$ where

- $V = \{+, *, (,), id, E\}$
- $\Sigma = \{+, *, (,), id\}$
- $R = \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E), E \rightarrow id\}$



Grammar G'
is ambiguous.

3.3 Pushdown Automata



PDA has

- An *input tape*
- A *finite control*
- A *stack*

Action:

- *Read:* Examine input, move tape head right, and replace the top of the stack with a new symbol
- *Think:* manipulate stack without reading



Definition: A **pushdown automata(PDA)** is a sextuple $M = (K, \Sigma, \Gamma, \Delta, s, F)$, where

- K is a finite set of states
 - Σ is an alphabet (the input symbols)
 - Γ is an alphabet (the stack symbols)
 - $s \in K$ is the initial state
 - $F \subseteq K$ is the set of final states
 - Δ , transition relation, is a subset of
$$(K \times (\Sigma \cup \{e\}) \times \Gamma^*) \times (K \times \Gamma^*).$$
-



- **PDA execution:** reading a symbol

Consider $((p, \alpha, \beta), (q, \gamma)) \in \Delta$, Then the PDA can:

- *enter some state q*
- *replace β by γ on the top of the stack*
- *advance the tape head*

- **PDA execution:** e -transition

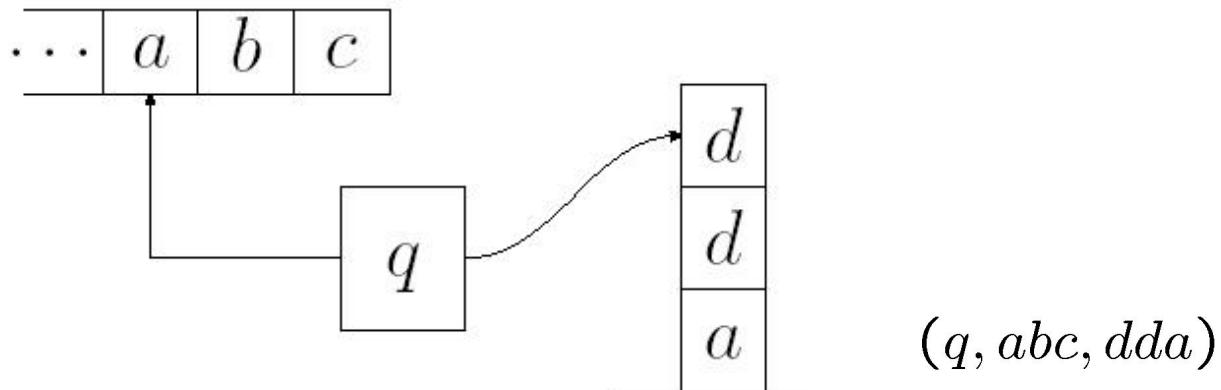
Consider $((p, e, \beta), (q, \gamma)) \in \Delta$, Then the PDA can:

- *enter some state q*
- *replace β by γ on the top of the stack*
- *does not advance the tape head*



Remark:

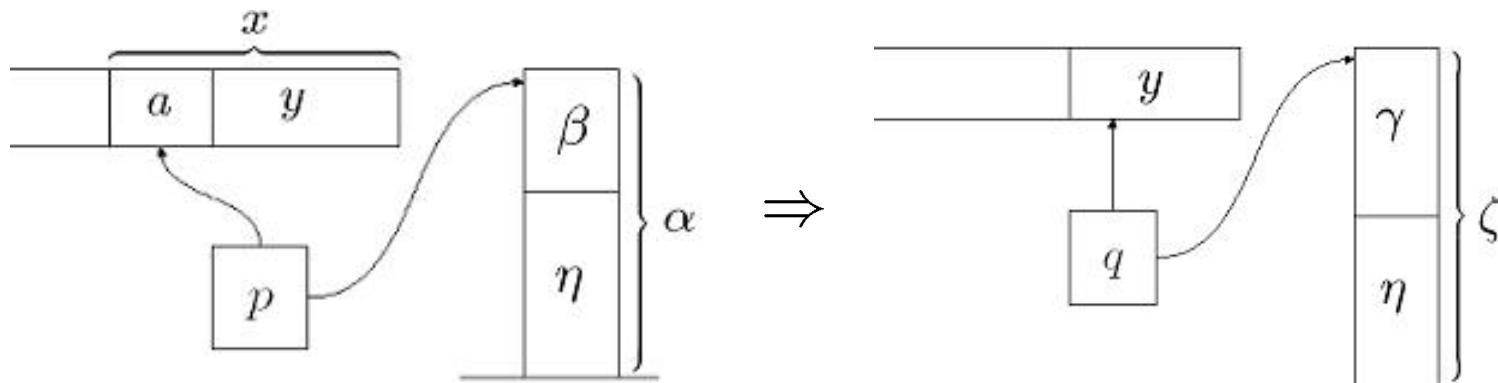
- Since several transition of M may be simultaneously applicable at any point, the machines are **nondeterministic**.
- $((p, u, e), (q, a))$ — *push* a ; $((p, u, a), (q, e))$ — *pop* a .
- **Configuration** of a PDA: a member of $K \times \Sigma^* \times \Gamma^*$.





- $(p, x, \alpha) \vdash_M (q, y, \zeta)$ (yield in one step) iff there is some transitions $((p, a, \beta), (q, \gamma)) \in \Delta$ such that

- $\square x = ay, a \in \Sigma \cup \{e\}$
- $\square \alpha = \beta\eta$
- $\square \zeta = \gamma\eta$ for some $\eta \in \Gamma^*$



- \vdash_M^* be the reflexive, transitive closure of \vdash_M .



□ Acceptance conditions

A PDA M accepts a string $w \in \Sigma^*$ iff

- $(s, w, e) \vdash_M^* (p, e, e)$ for some $p \in F$.
- There is a sequence of configuration $C_0, \dots, C_n (n > 0)$, $(s, w, e) = C_0 \vdash_M C_1 \vdash_M \dots \vdash_M (p, e, e)$ for some $p \in F$.

Note: Two traditional conditions

- Process the input, and accept if the stack is empty
- Accept if the PDA is in a final state

- The language accepted by M :

$$L(M) = \{w | (s, w, e) \vdash_M^* (p, e, e) \text{ for some state } p \in F\}.$$



Example: Design a PDA M to accepted the language

$$L = \{wcw^R : w \in \{a, b\}^*\}.$$

Solution:

Let $M = (K, \Sigma, \Gamma, \Delta, s, F)$, where

- $K = \{s, f\}$
- $\Sigma = \{a, b, c\}$
- $\Gamma = \{a, b\}$
- $F = \{f\}$
- Δ contains five transitions:

$((s, a, e), (s, a))$
$((s, b, e), (s, b))$
$((s, c, e), (f, e))$
$((f, a, a), (f, e))$
$((f, b, b), (f, e))$



Example: Design a PDA M to accepted the language

$$L = \{ww^R : w \in \{a, b\}^*\}.$$

Solution:

Let $M = (K, \Sigma, \Gamma, \Delta, s, F)$, where

- $K = \{s, f\}$
- $\Sigma = \{a, b\}$
- $\Gamma = \{a, b\}$
- $F = \{f\}$
- Δ contains five transitions:

$((s, a, e), (s, a))$
$((s, b, e), (s, b))$
$((s, e, e), (f, e))$
$((f, a, a), (f, e))$
$((f, b, b), (f, e))$



Example: Design a PDA M to accepted the language
 $L = \{w \in \{a, b\}^*: w \text{ has the same number of } a's \text{ and } b's\}$.

Solution:

Let $M = (K, \Sigma, \Gamma, \Delta, s, F)$, where

- $K = \{s, q, f\}$
- $\Sigma = \{a, b\}$
- $\Gamma = \{a, b, c\}$
- $F = \{f\}$
- Δ is listed right.

$((s, e, e), (q, c))$
$((q, a, c), (q, ac))$
$((q, a, a), (q, aa))$
$((q, a, b), (q, e))$
$((q, b, c), (q, bc))$
$((q, b, b), (q, bb))$
$((q, b, a), (q, e))$
$((q, e, c), (f, e))$

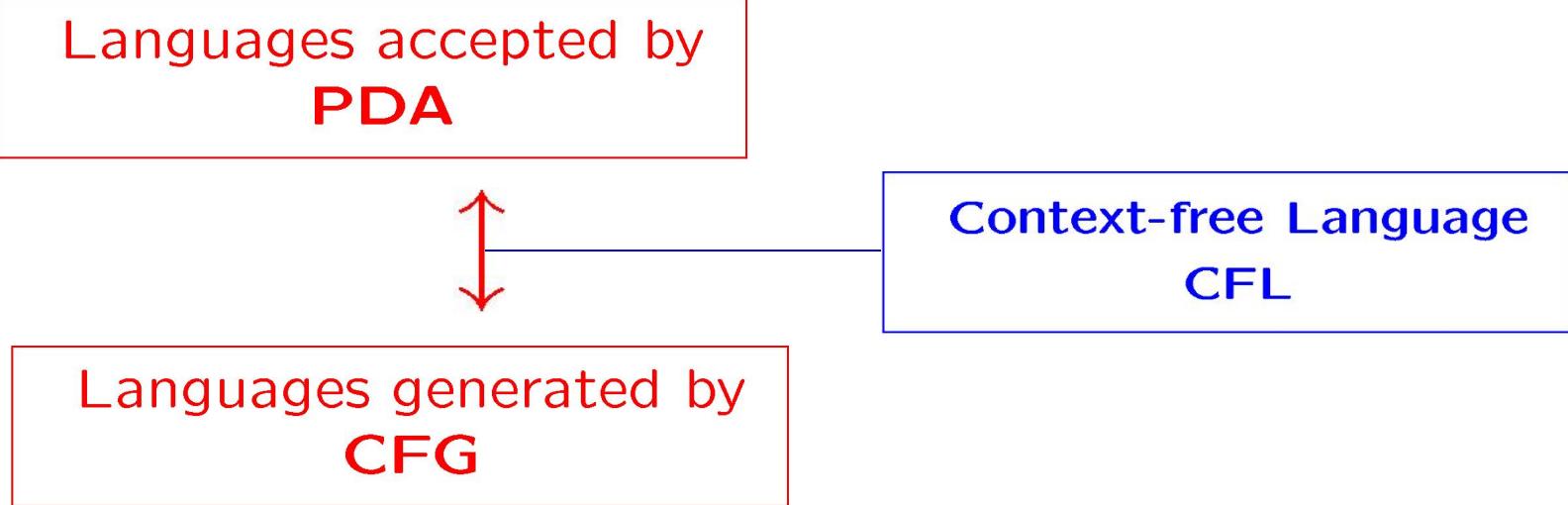


Example: Every FA can be trivially viewed a PDA that never operates on its stack.

- Let $M = (K, \Sigma, \Delta, s, F)$ be a NFA.
- Let $M' = (K, \Sigma, \emptyset, \Delta', s, F)$ be a PDA.
where $\Delta' = \{((p, u, e), (q, e)) : (p, u, q) \in \Delta\}$

Then, $L(M) = L(M')$.

3.4 Pushdown Automata and Context-Free Language



Theorem: The class of languages accepted by PDA is exactly the class of CFL.



I. Building a PDA from a CFG

Lemma: Each Context-Free language is accepted by some PDA.

Proof:

To build the PDA M for CFG $G = (V, \Sigma, R, S)$ such that

$$L(M) = L(G).$$

Main idea:

Define PDA M to mimics a leftmost derivation of the input string.



□ Construction of PDA

Define PDA $M = (K, \Sigma, \Gamma, \Delta, s, F)$.

- PDA M has just 2 states.

$p \sim$ start state

$q \sim$ final state

- Stack alphabet $\Gamma = V$.

- Let Δ contains the following transitions:

1) $((p, e, e), (q, S))$

2) $((q, e, A), (q, x))$ for each rule $A \rightarrow x \in R$

3) $((q, a, a), (q, e)), \forall a \in \Sigma$.



Example: Let CFG $G = (V, \Sigma, R, S)$ with $V = \{S, a, b, c\}$, $\Sigma = \{a, b, c\}$, and $R = \{S \rightarrow aSa, S \rightarrow bSb, S \rightarrow c\}$, then $L(G) = \{wcw^R : w \in \{a, b\}^*\}$.

The corresponding PDA, according to the construction above, is $M = (K, \Sigma, V, \Delta, s, F)$, with

- $K = \{p, q\}$
- $s = p$
- $F = \{q\}$
- Δ contains the following transitions:

I	$((p, e, e), (q, S))$
II	$((q, e, S), (q, aSa)), ((q, e, S), (q, bSb)), ((q, e, S), (q, c))$
III	$((q, a, a), (q, e)), ((q, b, b), (q, e)), ((q, c, c), (q, e))$



- Consider the string $abcbba$.

Derivation:

$$S \Rightarrow aSa \Rightarrow abSba \Rightarrow abbSbba \Rightarrow abcbba$$

Corresponding Computation:

$$\begin{aligned} (p, abcbba, e) &\vdash (q, abcbba, S) \vdash (q, abcbba, aSa) \\ &\vdash (q, bbcbba, Sa) \vdash (q, bbcbba, bSba) \\ &\vdash (q, bcbba, Sba) \\ &\vdash (q, bcbba, bSbba) \\ &\vdash (q, cbba, Sbba) \\ &\vdash (q, cbba, cbba) \\ &\vdash^* (q, e, e) \end{aligned}$$



□ Verify $L(M) = L(G)$

Claim: Let $w \in \Sigma^*$ and $\alpha \in (V - \Sigma)V^* \cup \{e\}$. Then

$$S \xrightarrow{L}^* w\alpha \iff (q, w, S) \vdash_M^* (q, e, \alpha)$$

The claim will suffice to Lemma. Taking $\alpha = e$ that

$$\begin{array}{ccc} S \xrightarrow{L}^* w & \iff & (q, w, S) \vdash_M^* (q, e, e) \\ \downarrow & & \downarrow \\ w \in L(G) & \iff & w \in L(M) \end{array}$$



□ Proof of Claim

⇒

Suppose that $S \xrightarrow{L^*} w\alpha$ where $w \in \Sigma^*$, $\alpha \in (V - \Sigma)V^* \cup \{e\}$.
Induction on the length of **the leftmost derivation of w** .

Basis step:

If the derivation is of length 0, then $w = e$, and $\alpha = S$,
hence indeed $(q, w, S) \vdash_M^* (q, e, \alpha)$.

Induction hypothesis:

Assume that if $S \xrightarrow{L^*} w\alpha$ by a derivation of length n or less,
 $n \geq 0$, then $(q, w, S) \vdash_M^* (q, e, \alpha)$.



Induction step:

$$S = u_0 \xrightarrow{L} u_1 \xrightarrow{L} \dots \xrightarrow{L} u_n \xrightarrow{L} u_{n+1} = w\alpha.$$

— be a leftmost derivation of $w\alpha$ from S .

Let $u_n = xA\beta$, A — the leftmost nonterminal of u_n .

$\Rightarrow u_{n+1} = x\gamma\beta$, where $x \in \Sigma^*$, $\beta, \gamma \in V^*$, and $A \rightarrow \gamma$ is in R .

- $S \xrightarrow{L} u_1 \dots \xrightarrow{L} u_n = xA\beta$ —a leftmost derivation of length n .

By the induction hypothesis,

$$(q, x, S) \vdash_M^* (q, e, A\beta) \tag{1}$$

- Since $A \rightarrow \gamma$ is a rule in R ,

$$(q, e, A\beta) \vdash_M (q, e, \gamma\beta) \tag{2}$$



- Note that $u_{n+1} = w\alpha = x\gamma\beta$.
 $\Rightarrow \exists$ string $y \in \Sigma^*$ such that $w = xy$ and $y\alpha = \gamma\beta$.

- Rewrite (1) and (2)

$$(q, w, S) \vdash_M^* (q, y, \gamma\beta) \quad (3)$$

- Since $y\alpha = \gamma\beta$

$$(q, y, \gamma\beta) \vdash_M^* (q, e, \alpha) \quad (4)$$

Combining (3) and (4) completes the induction step.

\Leftarrow (Omitted)



II. Building a CFG from a PDA

Lemma: If a language is accepted by a PDA, it is CFL.

Outline of the Proof

- Define the simple PDA
- Convert a PDA to an equivalent simple PDA
- Building a CFG from the simple PDA



□ Definition of Simple PDA

A PDA is **simple** if the following is true:

Whenever $((q, a, \beta), (p, \gamma))$ is a transition of the PDA and q is not the start state, then $\beta \in \Gamma$ and $|\gamma| \leq 2$

In other words, A simple PDA always

- consults its topmost stack symbol, and
- replace it either e , or with single stack symbol, or with two stack symbols.



□ Convert a PDA to an equivalent simple PDA

Let $M = (K, \Sigma, \Gamma, \Delta, s, F)$ be any PDA,

\Rightarrow construct a simple PDA $M' = (K', \Sigma, \Gamma \cup \{Z\}, \Delta', s', \{f'\})$
such that $L(M) = L(M')$.

$-s', f' \notin K$ be two new states, $Z \notin \Gamma$ be the stack bottom symbol.

$-K' = K \cup \dots$

$-\Delta'$ contains:

- the transition $((s', e, e), (s, Z))$ (start transition)
- for each $f \in F$, $((f, e, Z), (f', e))$ (final transition)
- all transition of Δ
 - replace with equivalent transitions that satisfy the simplicity condition.



Replace transitions with equivalent transitions that satisfy the simplicity condition:

- Get rid of transitions with $\beta \geq 2$.

Consider any transition $((q, a, \beta), (p, \gamma)) \in \Delta'$, where $\beta = B_1 \cdots B_n$, with $n > 1$.

Replace with the following transitions:

$$((q, e, B_1), (q_{B_1}, e))$$

$$((q_{B_1}, e, B_2), (q_{B_1 B_2}, e))$$

⋮

$$((q_{B_1 \cdots B_{n-2}}, e, B_{n-1}), (q_{B_1 \cdots B_{n-1}}, e))$$

$$((q_{B_1 \cdots B_{n-1}}, a, B_n), (p, \gamma))$$



- Get rid of transitions with $\gamma > 2$, without introducing any transitions with $\beta \geq 2$.

Consider any transition $((q, u, \beta), (p, \gamma)) \in \Delta'$, where $\gamma = C_1 \cdots C_m$, with $m \geq 2$.

Replace with the following transitions:

$$((q, u, \beta), (r_1, C_m))$$

$$((r_1, e, e), (r_2, C_{m-1}))$$

⋮

$$((r_{m-2}, e, e), (r_{m-1}, C_2))$$

$$((r_{m-1}, e, e), (p, C_1))$$



- Get rid of transitions with $\beta = e$, without introducing any transitions with $\beta \geq 2$ or $\gamma > 2$.

Consider any transition $((q, a, e), (p, \gamma))$ with $q \neq s'$.

Replace any such transitions by all transitions of the form

$$((q, a, A), (p, \gamma A)) \text{ for all } A \in \Gamma \cup \{Z\}$$

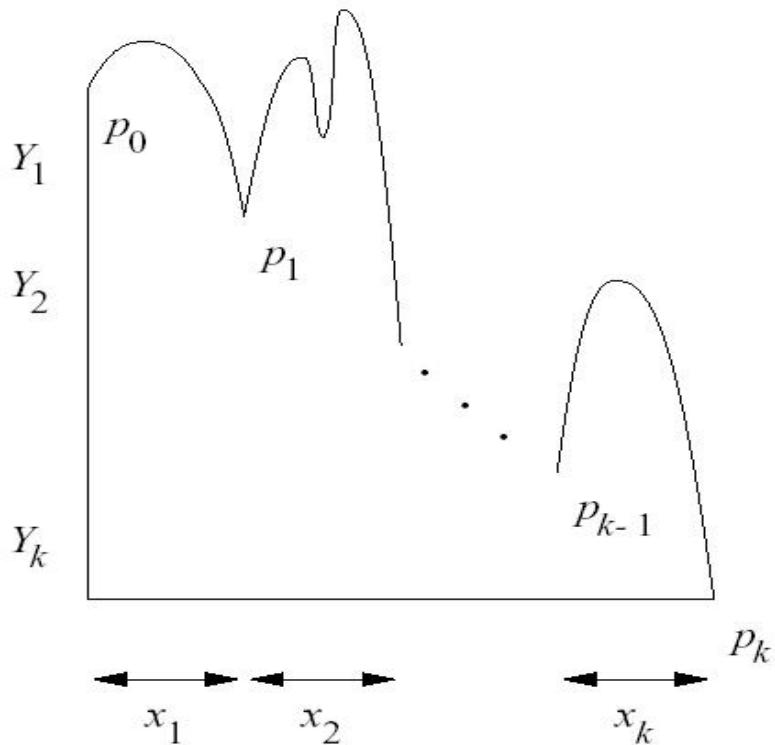
It is easy to see that $L(M) = L(M')$.



□ Construct A CFG from a simple PDA

⇒ Construct a CFG $G = (V, \Sigma, R, S)$ such that $L(G) = L(M')$.

- Let's look at how a PDA can consume $x = x_1x_2 \cdots x_k$ and empty the stack.



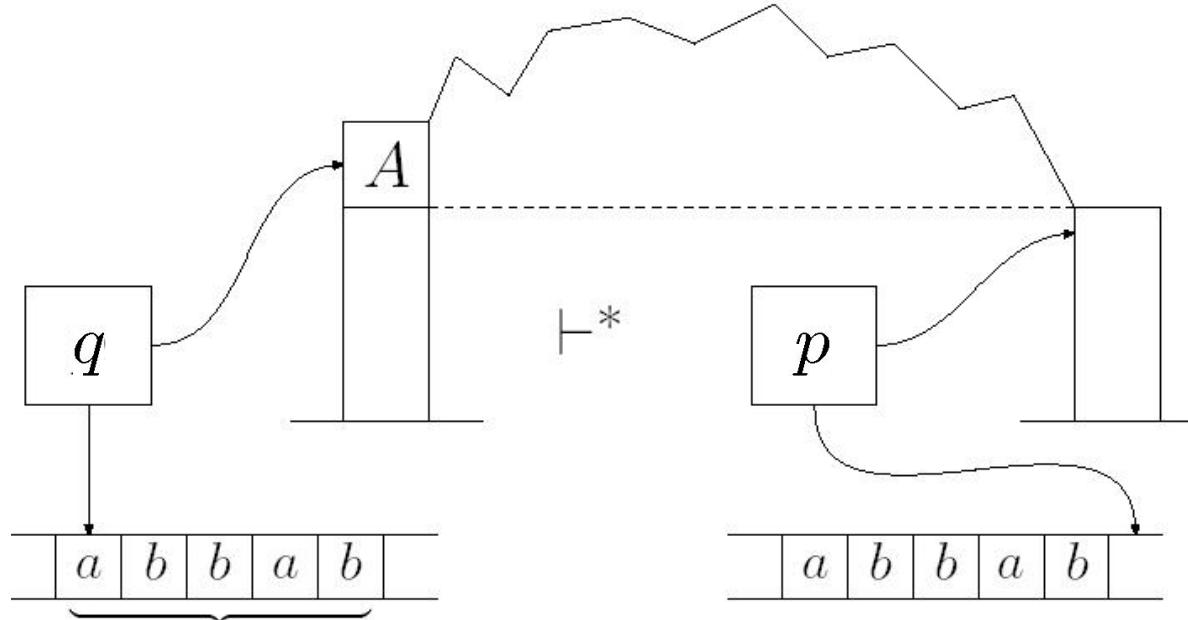


Definition:

The **nonterminals** $\langle q, A, p \rangle$:

represents any portion of the input string that might be read between a point when M is in state q with A on top of stack, and a point in time when M removes the occurrence of A from the stack and enters state p .

- $V = \{S\} \cup \Sigma \cup \{\langle q, A, p \rangle | \forall q, p \in K, A \in \Gamma \cup \{e, Z\}\}$



portion of input read from when M is in state q with A on stack to when M enters state p and pops A .



- The rules in R are of four types.
 - (1) The rules $S \rightarrow \langle s, Z, f' \rangle$:
 - s the start state of original PDA and f' the new final state.
 - (2) For each transition $((q, a, B), (r, e))$, where $q, r \in K$, $a \in \Sigma \cup \{e\}$, $B, C \in \Gamma \cup \{e\}$, and for each $p \in K$, we add the rule $\langle q, B, p \rangle \rightarrow a \langle r, C, p \rangle$.
 - (3) For each transition $((q, a, B), (r, C_1 C_2))$, where $q, r \in K$, $a \in \Sigma \cup \{e\}$, $B \in \Gamma \cup \{e\}$, and $C_1, C_2 \in \Gamma$ and for each $p, p' \in K$ we add the rule $\langle q, B, p \rangle \rightarrow a \langle r, C_1, p' \rangle \langle p', C_2, p \rangle$.
 - (4) For each $q \in K$, the rule $\langle q, e, q \rangle \rightarrow e$.



Claim: For any $q, p \in K$, $A \in \Gamma \cup \{e\}$, and $x \in \Sigma^*$,

$$< q, A, p > \Rightarrow_G^* x \Leftrightarrow (q, x, A) \vdash_M^* (p, e, e)$$

The claim will suffice to Lemma.

$$\begin{array}{ccc} < s, e, f > \Rightarrow_G^* x, \text{ for } f \in F & \Leftrightarrow & (s, x, e) \vdash_M^* (f, e, e) \\ \downarrow & & \downarrow \\ x \in L(G) & \Leftrightarrow & x \in L(M) \end{array}$$

3.5 Languages that are and are not Context-Free

Closure Properties

Theorem: The CFL are closed under union, concatenation, and Kleene star.

Proof:

Let $G_1 = (V_1, \Sigma_1, R_1, S_1)$ and $G_2 = (V_2, \Sigma_2, R_2, S_2)$ be two CFG.

Without loss generality, assume that $(V_1 - \Sigma_1)$ and $(V_2 - \Sigma_2)$ are disjoint.



a) Union

Let $G = (V_1 \cup V_2 \cup \{S\}, \Sigma_1 \cup \Sigma_2, R, S)$,

where

$R = R_1 \cup R_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}$

Then $L(G) = L(G_1) \cup L(G_2)$.

b) Concatenation

Let $G = (V_1 \cup V_2 \cup \{S\}, \Sigma_1 \cup \Sigma_2, R, S)$,

where

$R = R_1 \cup R_2 \cup \{S \rightarrow S_1 S_2\}$

Then $L(G) = L(G_1)L(G_2)$.



c) Kleene star

Let $G = (V_1 \cup \{S\}, \Sigma_1, R, S)$,
where
 $R = R_1 \cup \{S \rightarrow e, S \rightarrow SS_1\}$
Then $L(G) = L(G_1)^*$.

Remark:

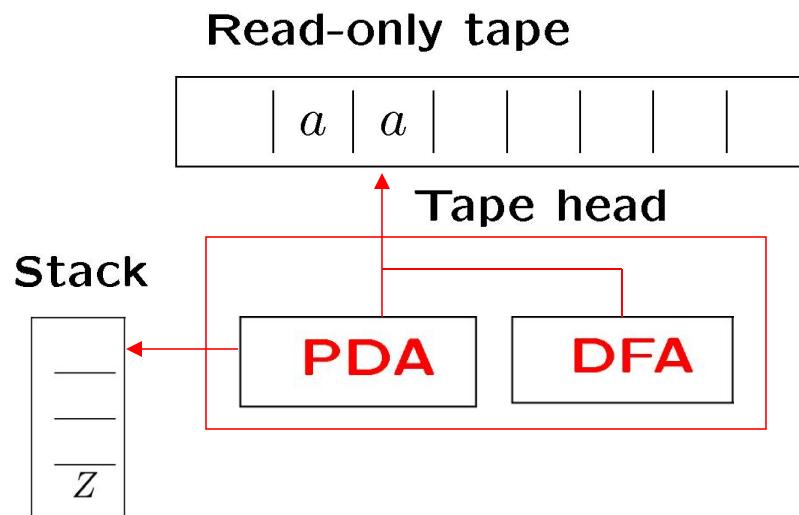
CFLs are not closed under intersection and complement.



Theorem: The Intersection of a CFL with a regular language is a CFL.

Proof:

Build a new machine
that simulates both
automata.





Let $M_1 = (K_1, \Sigma, \Gamma_1, \Delta_1, s_1, F_1)$ be a PDA and $M_2 = (K_2, \Sigma, \delta, s_2, F_2)$ be a DFA.

Bulid PDA $M = (K, \Sigma, \Gamma, \Delta, s, F)$, where

- $K = K_1 \times K_2$
- $\Gamma = \Gamma_1$
- $s = (s_1, s_2)$
- $F = F_1 \times F_2$
- Δ : For each $((q_1, a, \beta), (p_1, \gamma)) \in \Delta_1$, and $q_2 \in K_2$,
 $((q_1, q_2), a, \beta), ((p_1, \delta(q_2, a)), \gamma) \in \Delta$



Example:

$L = \{w : w \in \{a, b\}^*, w \text{ has equal numbers of } a's \text{ and } b's \text{ but containing no substring } abaa \text{ or } babb\}$.

Then L is context-free.

Solution:

$L_1 = \{w : w \in \{a, b\}^*, w \text{ has equal numbers of } a's \text{ and } b's\}$

L_1 be a CFL
accepted by a PDA

$L_2 = \{w \in \{a, b\}^* : w \text{ containing no substring } abaa \text{ or } babb\}.$
 $= \{a, b\}^* - \{a, b\}^*(abaa \cup babb)\{a, b\}^*$.

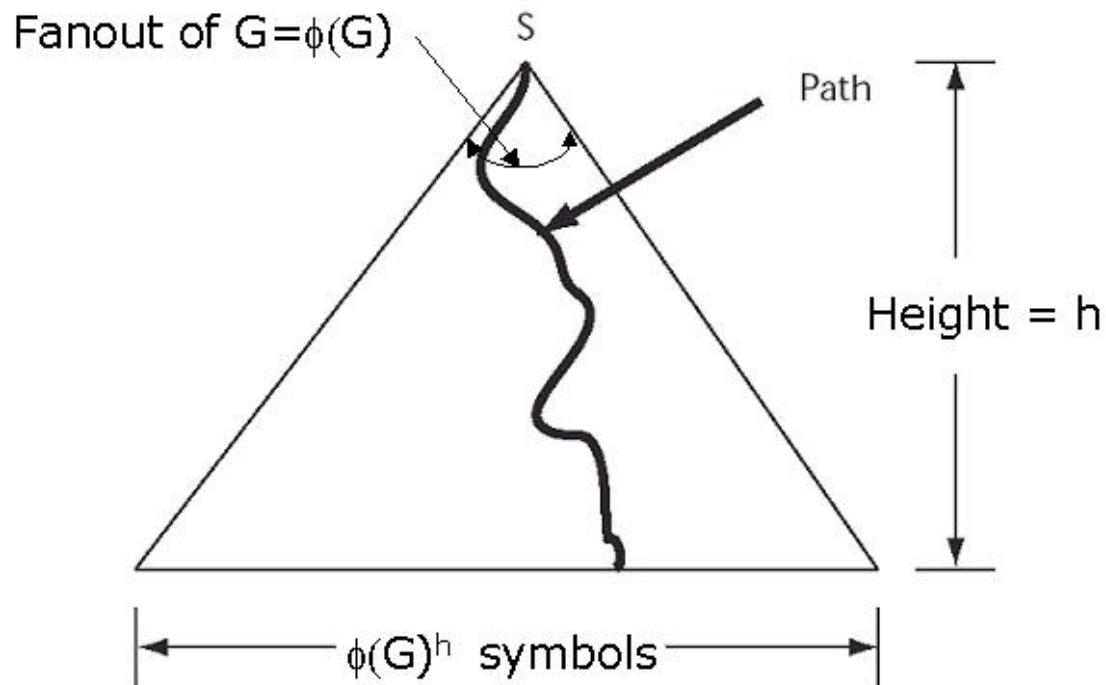
regular

$L = L_1 \cap L_2$ be a context-free language.



□ Pumping Theorem

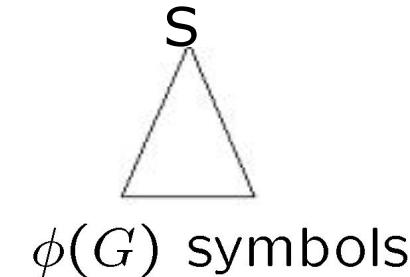
Lemma The yield of any parse tree of G of height h has length at most $\phi(G)^h$.





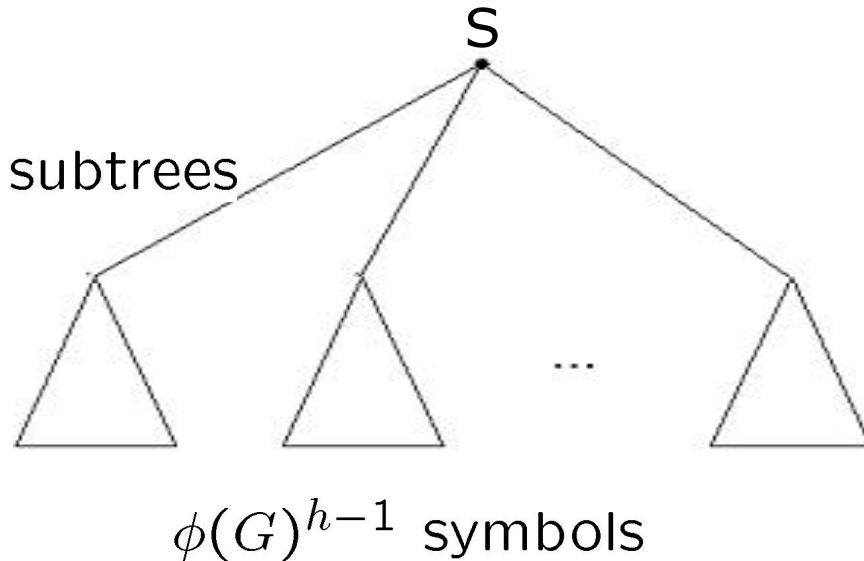
Proof: The proof is by induction on h .

Basis step: $h = 1$



Induction step:

at most $\phi(G)$ subtrees





Theorem (Pumping Theorem) Let $G = (V, \Sigma, R, S)$ be a CFG. Then any string $w \in L(G)$ of length greater than $\phi(G)^{|V-\Sigma|}$ can be rewritten as $w = uvxyz$ in such way that

- $|vy| \geq 1$
 - $uv^nxy^nz \in L(G)$ for every $n \geq 0$.
-

Proof

- Let w be such a string.
- Let T be the parse tree with root labeled S and with yield w that has the smallest number of leaves.



$$|w| > \phi(G)^{|V - \Sigma|}$$

\Rightarrow The height of $T > |V - \Sigma|$ (by Lemma)

$\Rightarrow \exists$ a path of length at least $|V - \Sigma| + 1$, with at least $|V - \Sigma| + 2$ nodes.

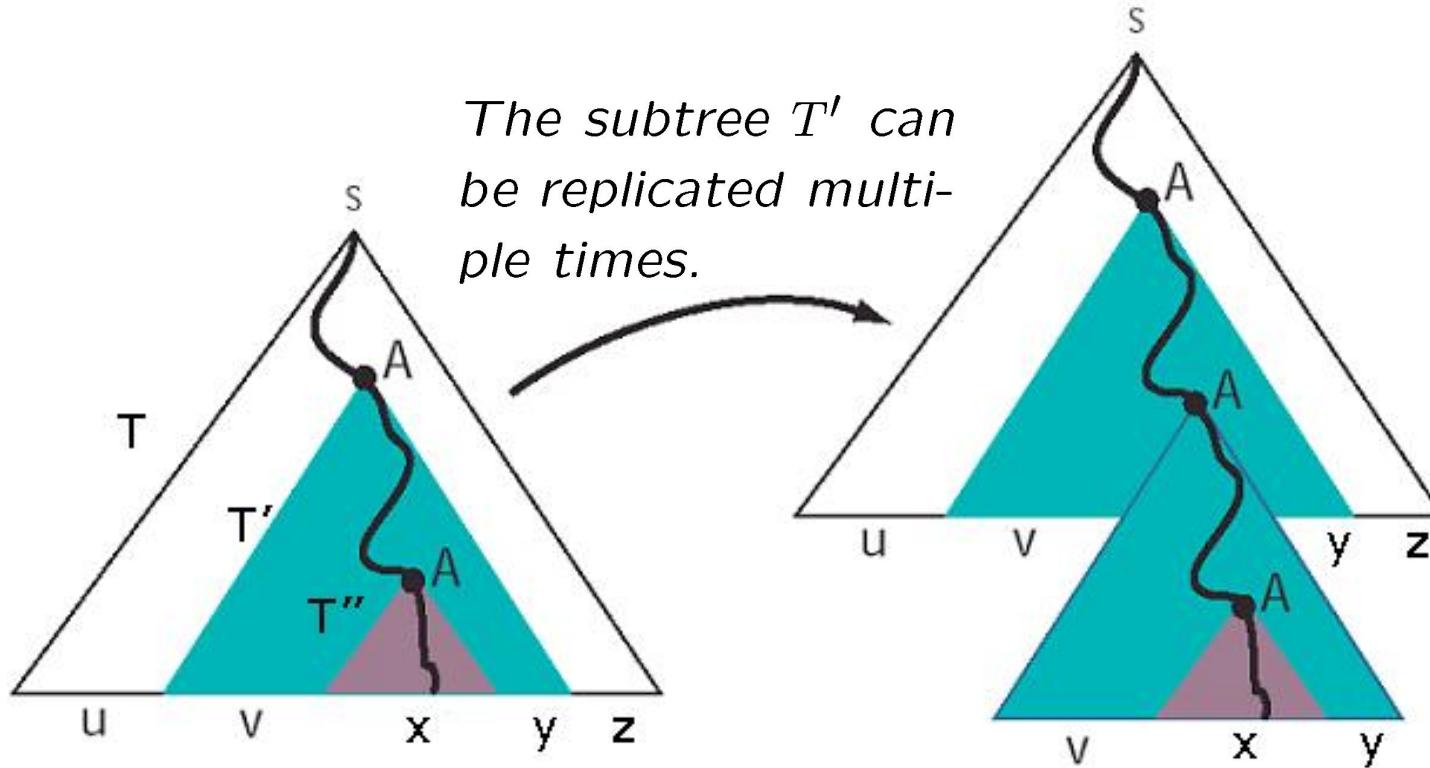
(only one labeled by terminal, the remaining are labeled by nonterminal).

$\Rightarrow \exists$ two nodes on the path labeled with the same symbol.

If $vy = e$, then there is a parse tree with root S and yield w with fewer leaves.

But T is the smallest tree of this kind.

—contradiction!

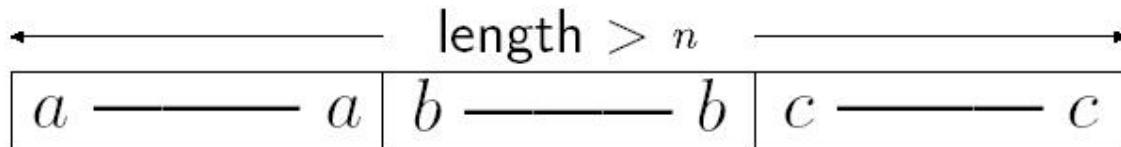




Example:

Show that $L = \{a^n b^n c^n : n \geq 0\}$ is not CFL.

Proof: Suppose that $L = L(G)$ for some CFG $G = (V, \Sigma, R, S)$.



Let $n > \frac{\phi(G)^{|V|-\sum|}}{3}$.

Then $w = a^n b^n c^n \in L(G)$ and can has a representation $w = uvxyz$ such that

- $vy \neq e$
- $uv^n xy^n z \in L(G)$ for each $n = 0, 1, 2, \dots$



Question: what about v, y ?

Case 1. Contain all three kinds of symbols

Case 2. Contain 2 kinds of symbols

- v, y contains occurrences of all three symbols a, b, c .
⇒ at least one of v, y must contain at least two of them.
⇒ order error in uv^2xy^2z .
- v, y contains occurrences of some but not all of them
⇒ uv^2xy^2z has unequal number of a 's, b 's and c 's.

Contradiction!!



Remark:

Proving that a language is not CFL

- Let L be the proposed CFL
- There is some n , by the pumping lemma
- Choose a string s , longer than n symbols, in the language L
- Using the pumping lemma, construct a new string s' that is not in the language



Example: Show $L = \{a^n | n \text{ is prime}\}$ is not Context-free.

Proof: Take a prime $p > \phi(G)^{|V-\Sigma|}$, where $G = (V, \Sigma, R, S)$ is CFG and $L = L(G)$.

Then $w = a^p$ can be written as prescribed by Pumping theorem, $w = uvxyz$ and $vy \neq e$.

Suppose that $vy = a^q$ and $uxz = a^r$ where p and q are natural numbers and $q > 0$.

Then the theorem states that $r + nq$ is prime, for all $n \geq 0$.

Contradiction!!

Remark:

Any CFL over a single-letter alphabet is regular.



Example: Show

$L = \{w \in \{a, b, c\}^* | w \text{ has an equal number of } a's, b's \text{ and } c's\}$
is not Context-free.

Note:

$$\{a^n b^n c^n | n \geq 0\} = L \cap a^* b^* c^*$$



Theorem: The CFL are not closed under intersection or complementation.

Proof: 1) Intersection

Counterexample

$$L_1 = \{a^n b^n c^m : m, n \geq 0\}$$

$$= \{a^n b^n : n \geq 0\} \circ c^*$$

$$L_2 = \{a^m b^n c^n : m, n \geq 0\}$$

L_1 and L_2 are CFL.

$\{a^n b^n c^n : n \geq 0\} = L_1 \cap L_2$
not CFL.

2) Complementation.

Note:

$$L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$$

3.6 Algorithms for Context-Free Grammar

Theorem (a) There is a polynomial algorithm which, given a CFG, constructs an equivalent PDA.
(b) There is a polynomial algorithm which, given a PDA, constructs an equivalent CFG.
(c) There is a polynomial algorithm which, given a CFG and a string x , decides whether $x \in L(G)$.

Remark:

(a) There are no algorithms for two CFGs(or PDA) are equivalent(see Chapter 5).



(b) For **regular language**, there was no need for an analog of part (c) above, since regular languages are represented in terms of an efficient algorithm for deciding precisely the membership question in (c): a DFA.

For **CFL**, we have so far introduced only non-algorithmic, nondeterministic acceptor—PDA.

We will construct a deterministic acceptor—algorithm in the next subsection that for any CFL.



Question:

Given CFG G and string w to determine if $w \in L(G)$

- Note that constructing a PDA from G does not in general solve this problem since the PDA will generally be *non-deterministic*.

- Brute-Force Method:

Check all parse trees of height up to some upper limit depending on G and $|w|$ (Exponentially costly)

- The dynamic programming algorithm:

1. Transform G into Chomsky normal form (CNF)
2. Apply a special algorithm for CNF grammars



Definition: A grammar is in **Chomsky normal form**(CNF) iff

$$R \subseteq (V - \Sigma) \times V^2.$$

Remark:

- The right-hand side of a rule in a CFG in Chomsky normal form must have length two.
- No grammar in CNF would be able to produce strings of length less than 2 and CFL containing such strings cannot be generated by grammar in CNF.



Theorem: For any CFG G there is a CFG G' in CNF such that $L(G') = L(G) - (\Sigma \cup \{e\})$. Furthermore, the construction of G' can be carried out in time polynomial in size of G .

Proof:

We shall show how to transform any given CFG $G = (V, \Sigma, R, S)$ into CNF.

We eliminate the following in order:

1. *Long rules:* of the form $A \rightarrow \alpha$ where $|\alpha| > 2$
2. *e-rules:* of the form $A \rightarrow e$
3. *Short rules:* of the form $A \rightarrow a$ or $A \rightarrow B$



□ Transforming a CFG into CNF

Step 1: Chop up all long rules

Let $A \rightarrow B_1B_2 \cdots B_n \in R$, where $B_1, \dots, B_n \in V$ and $n \geq 3$.
Replace this rule with $n - 1$ new rules, namely:

$$A \rightarrow B_1 A_1,$$

$$A_1 \rightarrow B_2 A_2,$$

...

$$A_{n-2} \rightarrow B_{n-1} B_{n-2}.$$

where A_1, \dots, A_n are newly introduced nonterminals.



Step 2: To eliminate e -rules

- a. Determine *the set of erasable nonterminals*

$$\mathcal{E} = \{A \in V - \Sigma : A \Rightarrow^* e\}$$

— is the set of all nonterminals that may drive the empty string.

Algorithm: Computing the set of erasable nonterminals

$$\mathcal{E} := \emptyset$$

While there is a rule $A \Rightarrow \alpha$ with $\alpha \in \mathcal{E}$ and $A \notin \mathcal{E}$ **do**
 add A to \mathcal{E}



-
- b. Given any rule $A \rightarrow \alpha$ where contains n occurrences of symbols from \mathcal{E} , replace it with 2^n rules in which $0, \dots, n$ occurrences are replaced with e .

e.g.

$$X \rightarrow AYZBY$$

$$Y \in \mathcal{E}.$$

$$\begin{aligned} & X \rightarrow AYZBY \\ \Rightarrow & X \rightarrow AYZB \\ & X \rightarrow AZBY \\ & X \rightarrow AZB \end{aligned}$$

- c. Eliminate all e -rules.

Remark:

Any derivation in the original grammar can be simulated in the new, and vice versa — with one exception: e can not be derived in the language any longer.



Step 3: To get rid of the short rules

For each $A \in V$ we compute by a simple closure algorithm, *the set $\mathcal{D}(A)$ of symbols that can be derived from A in grammar*, $\mathcal{D}(A) = \{B \in V : A \Rightarrow^* B\}$

Algorithm: Computing the set $\mathcal{D}(A)$

$$\mathcal{D}(A) := \{A\}$$

While there is a rule $B \Rightarrow C$ with $B \in \mathcal{D}(A)$ and $C \notin \mathcal{D}(A)$ **do**
 add C to $\mathcal{D}(A)$



Remark:

- For all symbols A , $A \in \mathcal{D}(A)$;
- If a be a terminal, then $\mathcal{D}(a) = \{a\}$.

If there is a rule $X \rightarrow Y$ where X, Y are nonterminal:

For any nonterminal $Z \in \mathcal{D}(Y)$ and any rule $Z \rightarrow \alpha$ where α is not a single nonterminal, add the rule $X \rightarrow \alpha$.

Then drop all rules $X \rightarrow Y$.



Example:

Transform the following grammar to CNF:

$$S \rightarrow SS$$

$$S \rightarrow (S)$$

$$S \rightarrow e$$

(1) *Eliminate long rules:*

$$S \rightarrow SS$$

$$\underline{S \rightarrow (S)}$$

$$S \rightarrow e$$

$$S \rightarrow (S_1, S_1 \rightarrow S)$$



(2) *Eliminate e rules:* $\mathcal{E} = \{S\}$

$$S \rightarrow SS, \quad S \rightarrow S$$

$$S \rightarrow (S_1, S_1 \rightarrow S), \quad S_1 \rightarrow)$$

$$S \rightarrow e$$

(3) *Get rid of the short rules:*

$$S \rightarrow SS,$$

$$S \rightarrow (S_1, S_1 \rightarrow S), \quad S_1 \rightarrow)$$

- $\mathcal{D}(S_1) = \{S_1, \) \} , $\mathcal{D}(A) = \{A\}, \forall A \in V - \{S_1\}.$$

- *Final Chomsky normal form grammar:*

$$S \rightarrow SS,$$

$$S \rightarrow (S_1, S_1 \rightarrow S), \quad S \rightarrow () .$$



□ Determining $x \in L(G)$ for G in CNF

- Let $G = (V, \Sigma, R, S)$ be a CFG in CNF.
- Let $x = x_1x_2 \cdots x_n$, $x_i \in \Sigma$ and $n \geq 2$.
- For each i and s such that $1 \leq i \leq i + s \leq n$, define:

$$N[i, i + s] = \{X : X \in V, X \Rightarrow^* x_i \cdots x_{i+s}\}.$$

- $x \in L(G)$ iff $S \in N[1, n]$.
- Calculate $N[i, i + s]$ by induction on s :

$A \in N[i, i + s]$ iff \exists rule $A \rightarrow BC$

such that $B \in N[i, k]$

$C \in N[k + 1, i + s]$.



Algorithm: Determine $x \in L(G)$

For $i := 1$ **to** n **do** $N[i, i] := \{x_i\}$; all other $N[i, j]$ are initially empty

For $s := 1$ **to** $n - 1$ **do**

For $i := 1$ **to** $n - s$ **do**

For $k := i$ **to** $i + s - 1$ **do**

If there is a rule $A \rightarrow BC \in R$ with $B \in N[i, k]$ and $C \in N[k + 1, i + s]$ then add A to $N[i, i + s]$

Accept x if $S \in N[1, n]$.

— Dynamic programming



Example:

G be a CFG in CNF:

$$S \rightarrow SS$$

$$S \rightarrow (S_1$$

$$S_1 \rightarrow S)$$

$$S \rightarrow () .$$

Determine whether

$$((())()) \in L(G) ?$$

(1								
\emptyset	(2							
\emptyset	S)	3						
\emptyset	\emptyset	\emptyset	(4					
\emptyset	\emptyset	\emptyset	\emptyset	(5				
\emptyset	\emptyset	\emptyset	\emptyset	S)	6			
\emptyset	S	\emptyset	S	S_1	\emptyset)	7		
S	S_1	\emptyset	S_1	\emptyset	\emptyset	\emptyset)	8



Homework 3:	
P120	3.1.3(c) 3.1.9 (a)(b)
P135	3.3.2 (c)(d)
P142	3.4.1
P148	3.5.1 (b)(c)(d) 3.5.2 (c) 3.5.14 (a)(b)(c)(d) 3.5.15