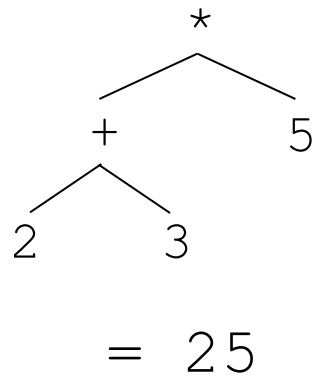


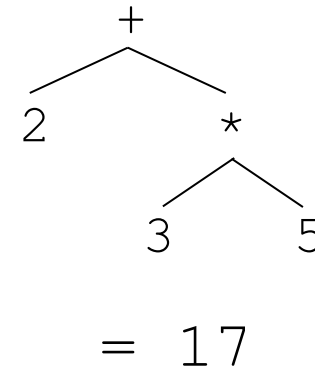
Context-free languages

Precedence in arithmetic expressions

```
bash-3.2$ python  
>>> 2+3*5  
17
```



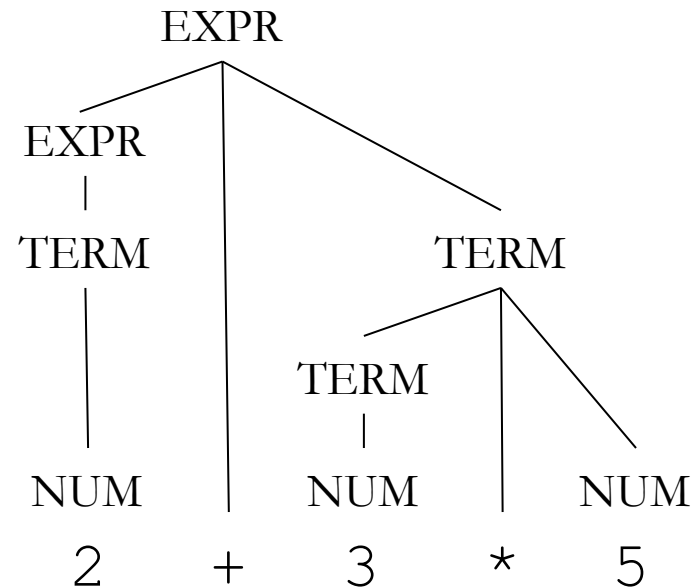
or



Grammars describe meaning

$\text{EXPR} \rightarrow \text{EXPR} + \text{TERM}$
 $\text{EXPR} \rightarrow \text{TERM}$
 $\text{TERM} \rightarrow \text{TERM} * \text{NUM}$
 $\text{TERM} \rightarrow \text{NUM}$
 $\text{NUM} \rightarrow 0-9$

rules for valid (simple)
arithmetic expressions



These rules always yield the correct meaning

The grammar of English

SENTENCE \rightarrow NOUN-PHRASE VERB-PHRASE

$\underbrace{\text{a girl}}_{\text{NOUN-PHRASE}} \underbrace{\text{likes the cat}}_{\text{VERB-PHRASE}}$

NOUN-PHRASE \rightarrow A-NOUN

or \rightarrow A-NOUN PREP-PHRASE

$\underbrace{\text{a girl}}_{\text{A-NOUN}}$

$\underbrace{\text{a girl}}_{\text{A-NOUN}} \underbrace{\text{with a flower}}_{\text{PREP-PHRASE}}$

The grammar of English

NOUN-PHRASE \rightarrow A-NOUN

or \rightarrow A-NOUN PREP-PHRASE

a girl

A-NOUN

a girl

A-NOUN

with a flower

PREP-PHRASE

recursive
structure

PREP-PHRASE \rightarrow PREP NOUN-PHRASE

with a flower

PREP NOUN-PHRASE

The grammar of (parts of) English

SENTENCE → NOUN-PHRASE VERB-PHRASE

NOUN-PHRASE → A-NOUN

NOUN-PHRASE → A-NOUN PREP-PHRASE

VERB-PHRASE → CMPLX-VERB

VERB-PHRASE → CMPLX-VERB PREP-PHRASE

PREP-PHRASE → PREP A-NOUN

A-NOUN → ARTICLE NOUN

CMPLX-VERB → VERB NOUN-PHRASE

CMPLX-VERB → VERB

ARTICLE → a

ARTICLE → the

NOUN → cat

NOUN → girl

NOUN → flower

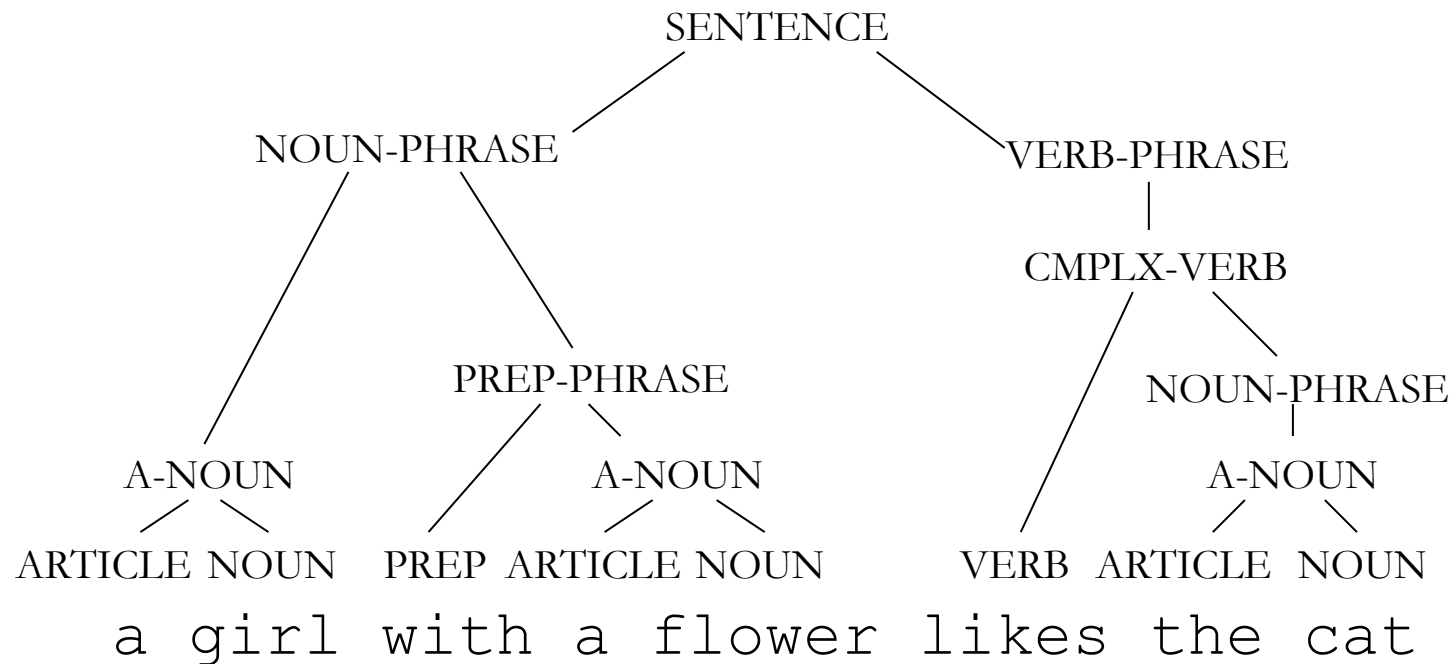
VERB → likes

VERB → touches

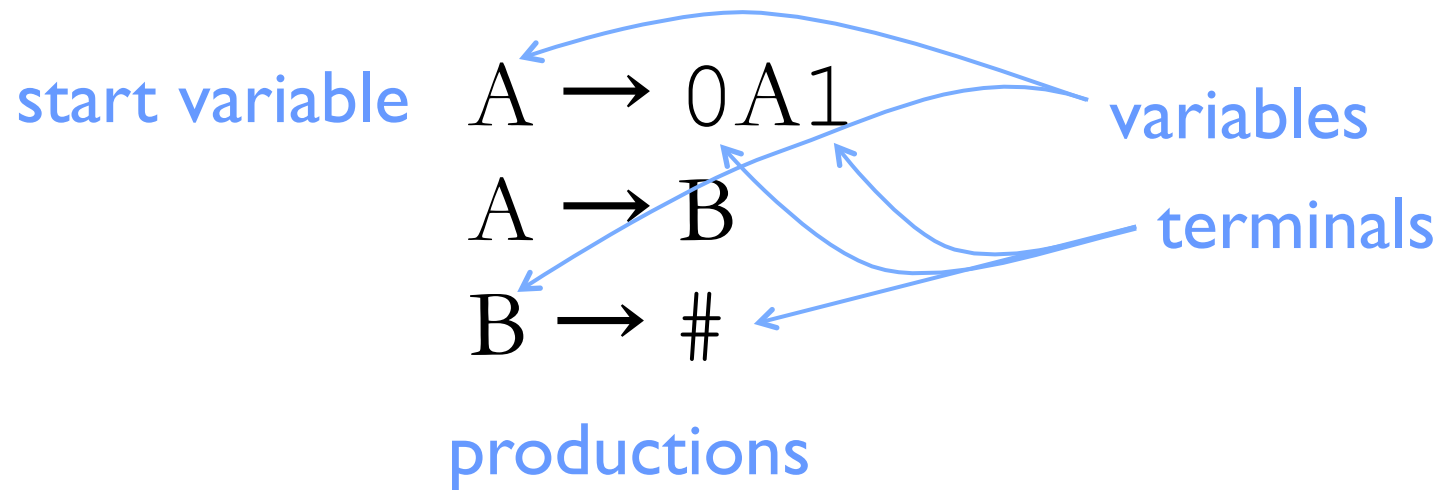
VERB → sees

PREP → with

The meaning of sentences



Context-free grammar



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

derivation

Context-free grammar

- A **context-free grammar** is given by (V, Σ, R, S) where
 - V is a finite set of **variables** or **non-terminals**
 - Σ is a finite set of **terminals**
 - R is a set of **productions** or **substitution rules** of the form

$$A \rightarrow \alpha$$

A is a variable and α is a **string** of variables and terminals

- S is a variable called the **start variable**

Notation and conventions

$E \rightarrow E + E$

$E \rightarrow (E)$

$E \rightarrow N$

$N \rightarrow 0N$

$N \rightarrow 1N$

$N \rightarrow 0$

$N \rightarrow 1$

Variables: E, N

Terminals: +, (,), 0, 1

Start variable: E

shorthand:

$E \rightarrow E + E \mid (E) \mid N$

$N \rightarrow 0N \mid 1N \mid 0 \mid 1$

conventions:

Variables in UPPERCASE

Start variable comes first
(typically denoted by S)

Derivation

- A **derivation** is a sequential application of productions:

$$\begin{aligned}
 E &\Rightarrow E + E \\
 &\Rightarrow (E) + E \\
 &\Rightarrow (E) + N \\
 &\Rightarrow (E + E) + 1 \\
 &\Rightarrow (E + E) + 1 \\
 &\Rightarrow (E + N) + 1 \\
 &\Rightarrow (N + N) + 1 \\
 &\Rightarrow (N + 1N) + 1 \\
 &\Rightarrow (N + 10) + 1 \\
 &\Rightarrow (1 + 10) + 1
 \end{aligned}$$

derivation

$$\begin{aligned}
 E &\rightarrow E + E \mid (E) \mid N \\
 N &\rightarrow 0N \mid 1N \mid 0 \mid 1
 \end{aligned}$$

$\alpha \Rightarrow \beta$ one production

$\alpha \xRightarrow{*} \beta$ **derivation**

$E \xRightarrow{*} (1 + 10) + 1$

(Two derivations are different if they use two different rules at any point.)

Context-free languages

- The *language generated by a CFG G* is the set of all strings at the end of a derivation

$$L(G) = \{ w : w \in \Sigma^* \text{ and } S \xRightarrow{*} w \}$$

- Questions we will ask:

Given a CFG, what is its language?

Given a language, write a CFG for it.

Analysis example I

$$\begin{aligned} A &\rightarrow 0A1 \mid B \\ B &\rightarrow \# \end{aligned}$$

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

- Can you derive:

$$00\#11 \quad A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 00B11 \Rightarrow 00\#11$$

$$\# \quad A \Rightarrow B \Rightarrow \#$$

$$00\#111 \quad \text{No, there is an uneven number of 0s and 1s}$$

$$00\#\#11 \quad \text{No, there are too many \#}$$

Analysis example 2

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

- Can you derive

$$\begin{aligned} S &\Rightarrow (S) \\ &\Rightarrow () \end{aligned}$$

$()$

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

$((())())$

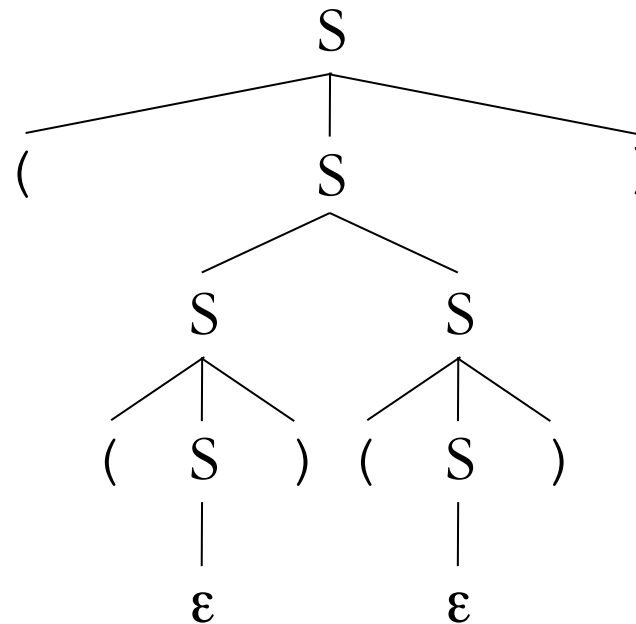
Parse trees

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

- A **parse tree** gives a more compact representation:

$S \Rightarrow (S)$
 $\Rightarrow (SS)$
 $\Rightarrow ((S)S)$
 $\Rightarrow ((S)(S))$
 $\Rightarrow (() (S))$
 $\Rightarrow (() ())$

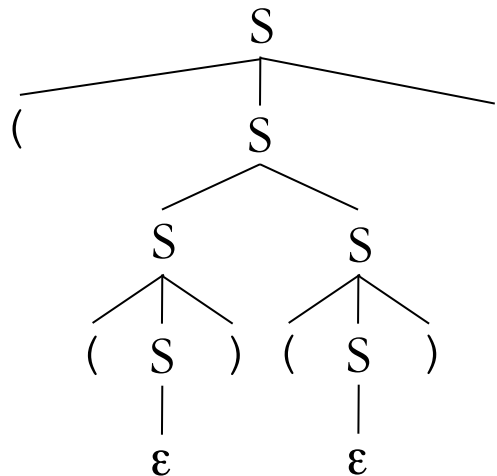
 $(() ())$



Parse trees

$S \Rightarrow (S)$
 $\Rightarrow (SS)$
 $\Rightarrow ((S)S)$
 $\Rightarrow ((S)(S))$
 $\Rightarrow ((())(S))$
 $\Rightarrow ((())())$

$S \Rightarrow (S)$
 $\Rightarrow (SS)$
 $\Rightarrow ((S)S)$
 $\Rightarrow ((())S)$
 $\Rightarrow ((())(S))$
 $\Rightarrow ((())())$



$S \Rightarrow (S)$
 $\Rightarrow (SS)$
 $\Rightarrow (S(S))$
 $\Rightarrow ((S)(S))$
 $\Rightarrow ((())(S))$
 $\Rightarrow ((())())$

$S \Rightarrow (S)$
 $\Rightarrow (SS)$
 $\Rightarrow (S(S))$
 $\Rightarrow (S())$
 $\Rightarrow ((S)())$
 $\Rightarrow ((())())$

One parse tree can represent many derivations

Analysis example 2

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

- Can you derive

(() ()

No, because there is an **uneven** number of (and)

()) (()

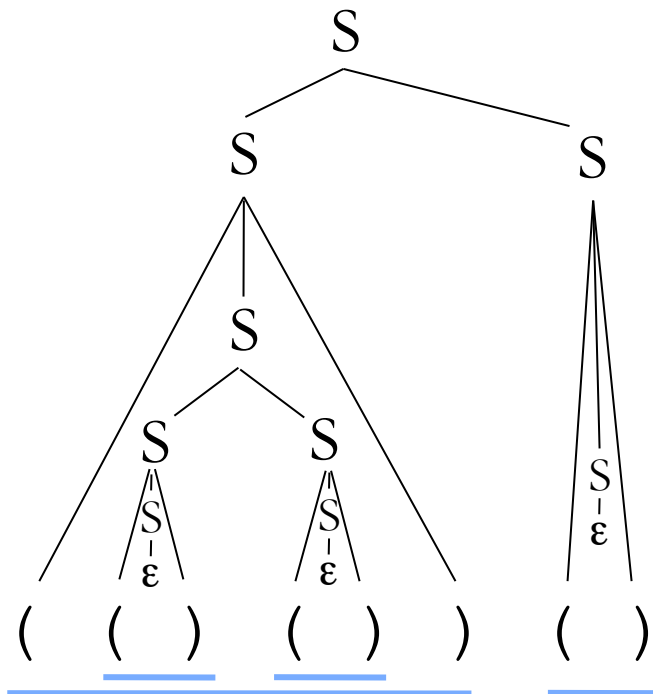
No, because there is a **prefix** with an excess of)

Analysis example 2

$$S \rightarrow SS \mid (S) \mid \varepsilon \quad L(G) = \{w:$$

w has the same number of (and)

no prefix of w has more) than (}



Parsing rules:

Divide w up in **blocks** with same number of (and)

Each block is in $L(G)$

Parse each block **recursively**

Design example I

$$L_1 = \{0^n 1^n \mid n \geq 0\}$$

These strings have recursive structure:

000000111111

0000011111

00001111

000111

0011

01

ϵ

$S \rightarrow 0S1 \mid \epsilon$

Design example 2

L_2 = all natural numbers in decimal notation

0, 109, 2, 23
allowed

ϵ , 01, 003
not allowed

$S \rightarrow 0 | LN$

$N \rightarrow DN | \epsilon$

$D \rightarrow 0 | L$

$L \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

1052870032

↓ $\underbrace{\hspace{10em}}$
any number N
leading digit L

Design examples

$$L_3 = \{0^n 1^n 0^m 1^m \mid n \geq 0, m \geq 0\}$$

010011

00110011

000111

These strings have **two parts**:

$$L_3 = L_p L_s$$

$$L_p = \{0^n 1^n \mid n \geq 0\}$$

$$L_s = \{0^m 1^m \mid m \geq 0\}$$

rules for L_p : $S_1 \rightarrow 0S_11 \mid \epsilon$

L_s is the same as L_p

$$S \rightarrow S_1 S_1$$

$$S_1 \rightarrow 0S_11 \mid \epsilon$$

Design examples

$$L_4 = \{0^n 1^m 0^m 1^n \mid n \geq 0, m \geq 0\}$$

011001

0011

1100

00110011

These strings have **nested structure**:

outer part: $0^n 1^n$

inner part: $1^m 0^m$

$S \rightarrow 0S1 \mid I$

$I \rightarrow 1I0 \mid \epsilon$

Design examples

$L_5 = \{x: x \text{ has two 0-blocks with same number of 0s}\}$

01011, 001011001, 10010101001
allowed

01001000, 01111
not allowed

$\underbrace{1001}_{\text{initial part}} \underbrace{00110100}_{\text{middle part}} \underbrace{10110}_{\text{final part}}$
A B C

A: ε , or ends in 1

C: ε , or begins with 1

Design examples

10010011010010110

A B C

A: ϵ , or ends in 1

C: ϵ , or begins with 1

U: any string

$S \rightarrow ABC$

$A \rightarrow \epsilon \mid U1$

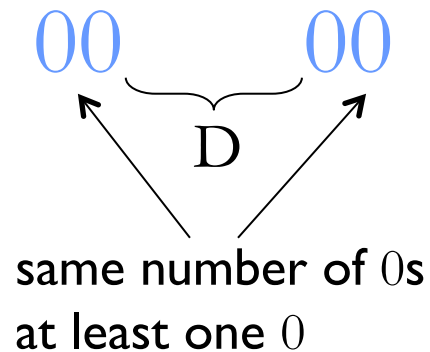
$U \rightarrow 0U \mid 1U \mid \epsilon$

$C \rightarrow \epsilon \mid 1U$

$B \rightarrow 0D0 \mid 0B0$

$D \rightarrow 1U1 \mid 1$

B has recursive structure:



D: begins and ends in 1

Context-free versus regular

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

$$S \rightarrow U111$$

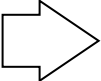
$$U \rightarrow 0U \mid 1U \mid \varepsilon$$

- Can you do so for **every** regular language?

Every regular language is context-free



From regular to context-free

regular expression		CFG
\emptyset		grammar with no rules
ε		$S \rightarrow \varepsilon$
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^*		$S \rightarrow S S_1 \mid \varepsilon$

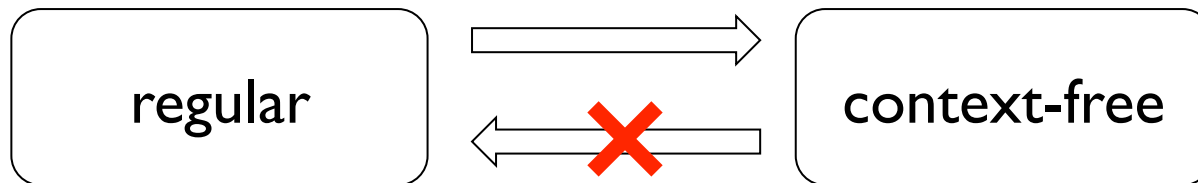
(S becomes the new **start symbol**)

Context-free versus regular

- Is every context-free language regular?

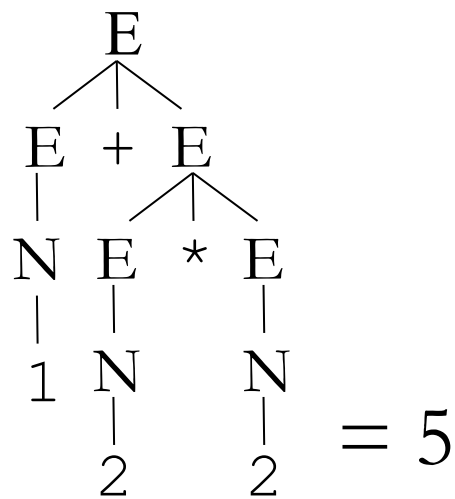
$$S \rightarrow 0S1 \mid \varepsilon \qquad L_1 = \{0^n 1^n : n \geq 0\}$$

Is context-free but not regular

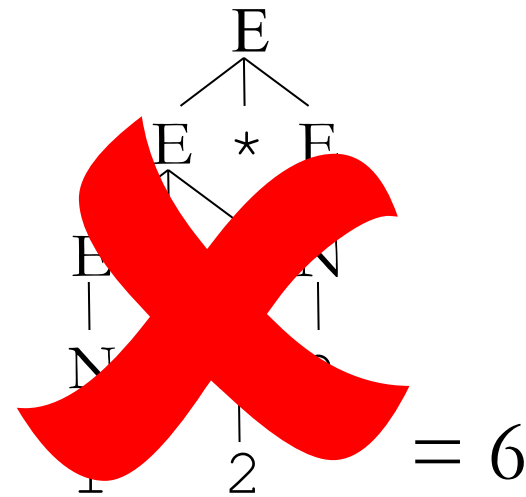


Ambiguity
Parsing algorithms
Probabilistic Context-Free Grammars

Ambiguity

$$E \rightarrow E + E \mid E * E \mid (E) \mid N$$
$$N \rightarrow 1N \mid 2N \mid 1 \mid 2$$


$1+2*2$

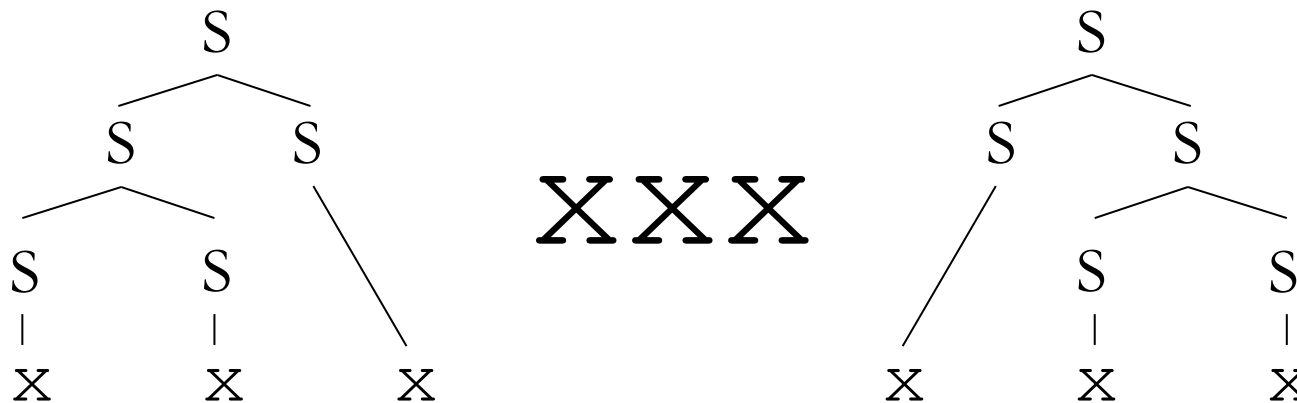


A CFG is **ambiguous** if some string has more than one parse tree

Example

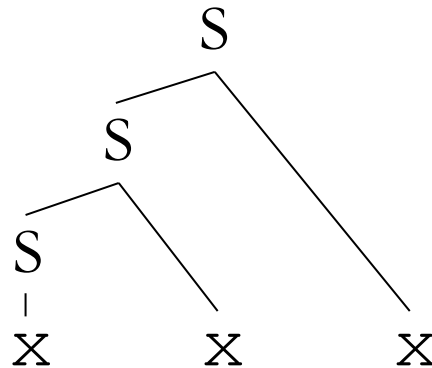
Is $S \rightarrow SS \mid x$ ambiguous?

Yes, because



Disambiguation

$$S \rightarrow \textcircled{SS} \mid x \quad \Rightarrow \quad S \rightarrow Sx \mid x$$



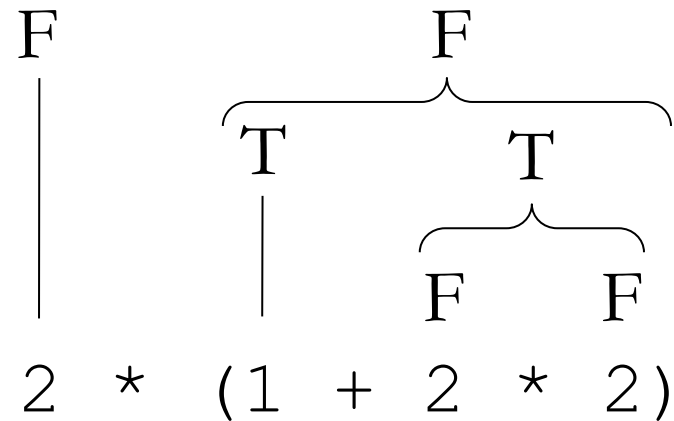
Sometimes we can **rewrite the grammar** to remove ambiguity

Disambiguation

same precedence!

$$\begin{aligned} E &\rightarrow E + E \mid E * E \mid (E) \mid N \\ N &\rightarrow 1N \mid 2N \mid 1 \mid 2 \end{aligned}$$

Divide expression
into **terms** and **factors**



Disambiguation

$$\begin{aligned} E &\rightarrow E + E \mid E * E \mid (E) \mid N \\ N &\rightarrow 1N \mid 2N \mid 1 \mid 2 \end{aligned}$$

An expression is a sum of one or more **terms**

Each term is a product of one or more **factors**

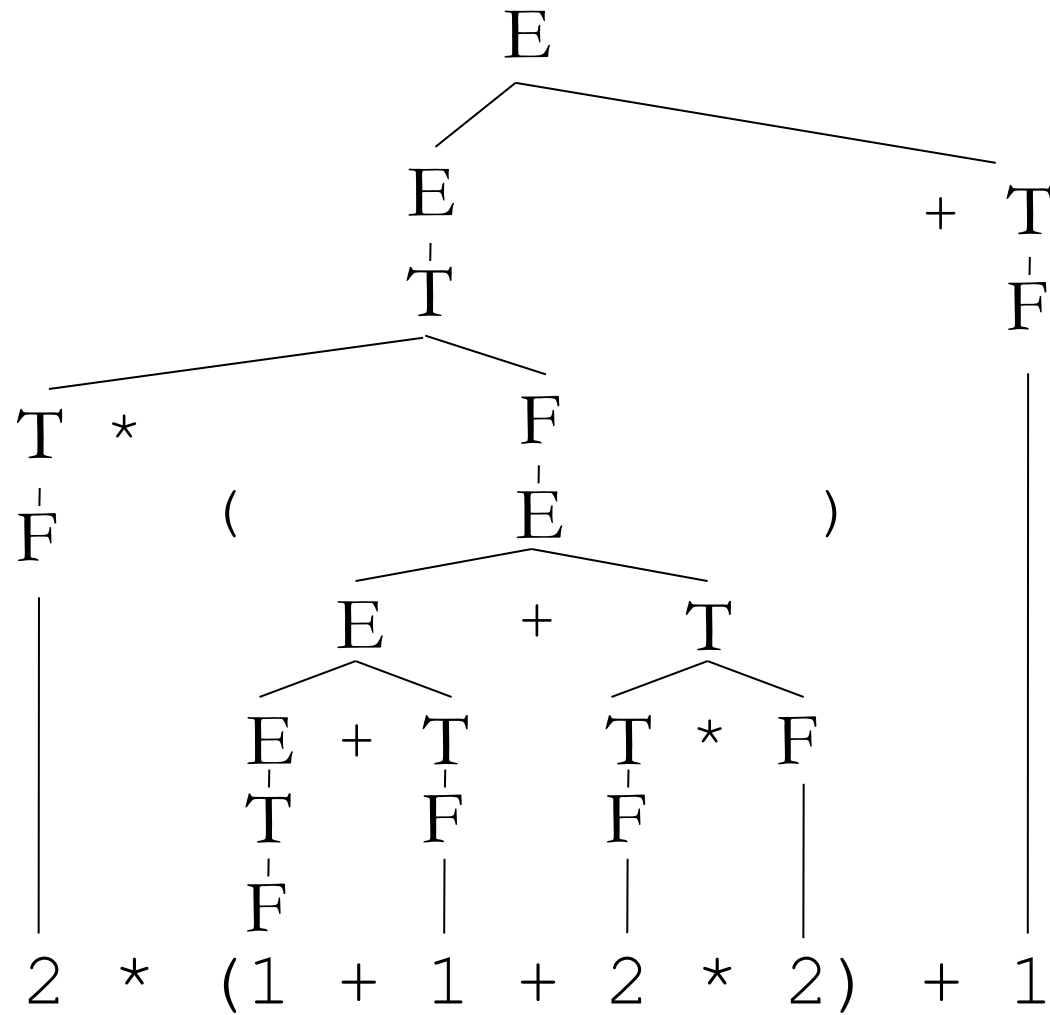
Each factor is a **parenthesized expression** or a **number**

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

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

$$F \rightarrow (E) \mid 1 \mid 2$$

Parsing example


$$\begin{array}{l} E \rightarrow T \mid E + T \\ T \rightarrow F \mid T * F \\ F \rightarrow (E) \mid 1 \mid 2 \end{array}$$

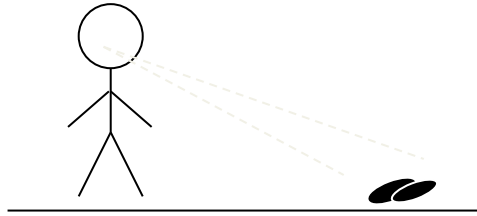
Unique parse tree
(because $+/*$ is
right-binding)

Disambiguation

- Disambiguation is **not always possible** because
 - There exist **inherently ambiguous** languages
 - There is **no general procedure** for disambiguation
- In **programming languages**, ambiguity comes from precedence rules, and we can deal with it like in the previous example
- In English, ambiguity is sometimes a problem:

He ate the cookies on the floor.

Ambiguity in English



He ate the cookies on the floor.



Parsing

$$\begin{aligned} S &\rightarrow 0S1 \mid 1S0S \mid T \\ T &\rightarrow S \mid \varepsilon \end{aligned}$$

input: 0011

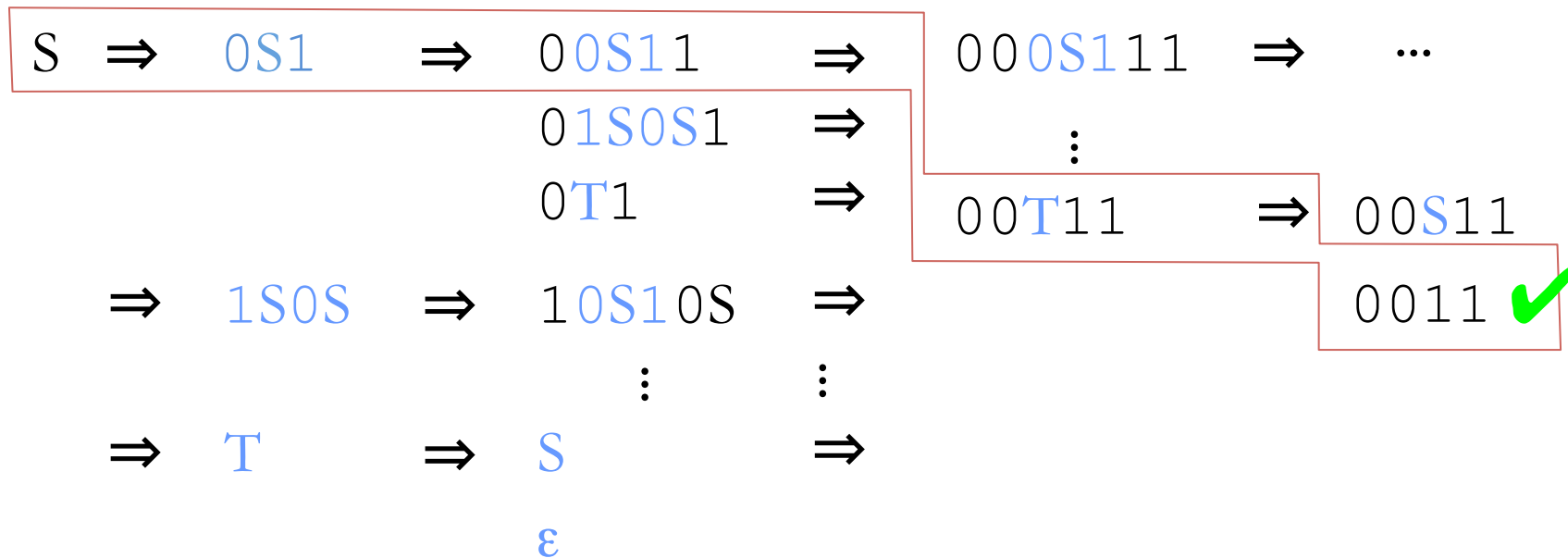
How would we **program** the computer to build a parse tree for us?

Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid \varepsilon$$

input: 0011



First idea: Try **all** derivations

Problems

① Trying all derivations may take a very long time

② If input is **not in the language**, parsing will never stop

When to stop

$$\begin{aligned} S &\rightarrow 0S1 \mid 1S0S \mid T \\ T &\rightarrow S \mid \varepsilon \end{aligned}$$

Idea 2: Stop when

$$|\text{derived string}| > |\text{input}|$$

Problems

$$\begin{array}{ccccccc} S & \Rightarrow & 0S1 & \Rightarrow & 0T1 & \Rightarrow & 01 \\ 1 & & 3 & & 3 & & 2 \end{array}$$

Derived strings may **shrink**
because of “ ε -productions”

$$S \Rightarrow T \Rightarrow S \Rightarrow T \Rightarrow \dots$$

Derivation may **loop**
because of “unit productions”

Task: remove ε - and unit- productions

Removal of ϵ -productions

- A variable N is **nullable** if it can derive the empty string

$$N \xRightarrow{*} \epsilon$$

- ① **Identify** all nullable variables N
- ② **Remove** all ϵ -productions carefully (while adding extra productions to compensate for this)
- ③ If start variable S is nullable:
Add a **new start variable** S'
Add **special productions** $S' \rightarrow S \mid \epsilon$

Example

grammar

$S \rightarrow ACD$

$A \rightarrow a$

$B \rightarrow \varepsilon$

$C \rightarrow ED \mid \varepsilon$

$D \rightarrow BC \mid b$

$E \rightarrow b$

nullable variables

B C D

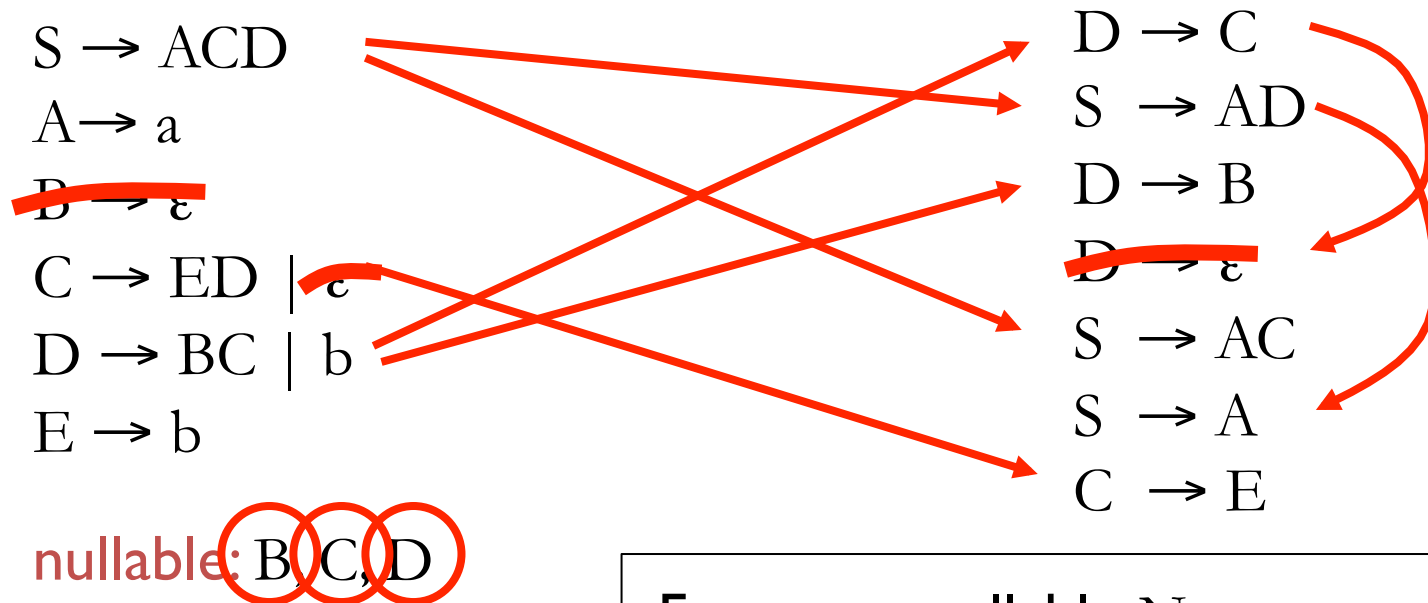
Repeat the following:

If $X \rightarrow \varepsilon$, mark X as nullable

If $X \rightarrow YZ \dots W$, all marked nullable,
mark X as nullable also.

① Identify all nullable variables

Eliminating ϵ -productions



For every nullable N :

If you see $X \rightarrow \alpha N \beta$, add $X \rightarrow \alpha \beta$

If you see $N \rightarrow \epsilon$, remove it.

② Remove all ϵ -productions carefully

The end result

The old grammar

$S \rightarrow ACD$

$A \rightarrow a$

$B \rightarrow \varepsilon$

$C \rightarrow ED \mid \varepsilon$

$D \rightarrow BC \mid b$

$E \rightarrow b$

The new grammar

$S \rightarrow ACD$

$A \rightarrow a$

$C \rightarrow ED$

$D \rightarrow BC \mid b$

$E \rightarrow b$

$D \rightarrow C$

$S \rightarrow AD$

$D \rightarrow B$

$S \rightarrow AC$

$S \rightarrow A$

$C \rightarrow E$

② Remove all ε -productions carefully

Eliminating unit productions

- A **unit production** is a production of the form

$$A \rightarrow B$$

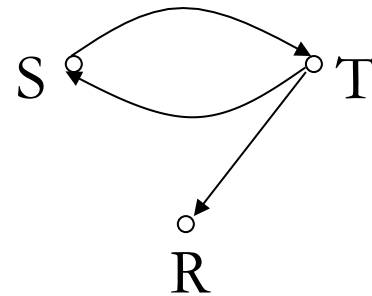
grammar:

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$

unit productions graph:

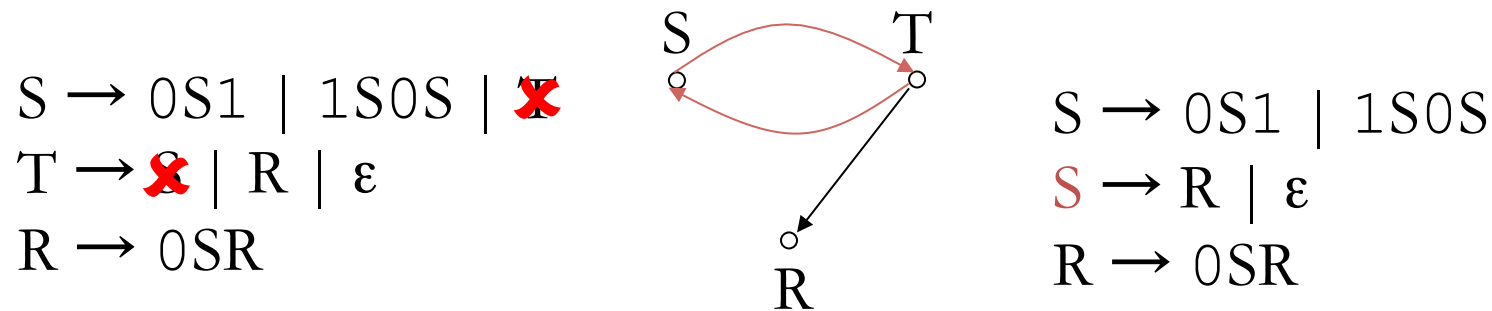


Removal of unit productions

- ① If there is a **cycle** of unit productions

$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow A$$

delete it and **replace** everything with A



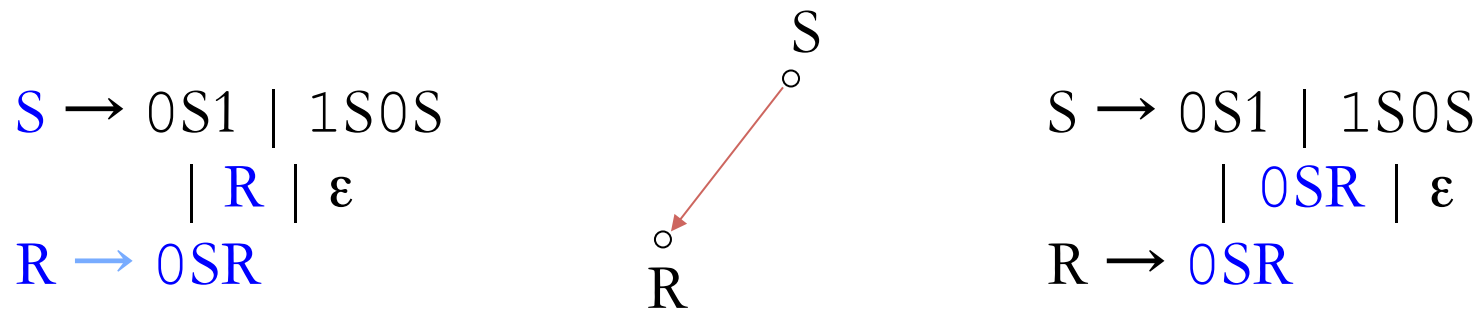
replace T by S

Removal of unit productions

② Replace every chain

$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow \alpha$$

by $A \rightarrow \alpha, B \rightarrow \alpha, \dots, C \rightarrow \alpha$



$S \rightarrow R \rightarrow 0SR$ is replaced by $S \rightarrow 0SR, R \rightarrow 0SR$

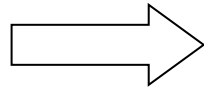
Recap

Problem: If input is **not in the language**, parsing will never stop

Solution: *important to do
in this order* ↓
Eliminate ϵ productions
Eliminate unit productions
Try all possible derivations but
stop parsing when

$$|\text{derived string}| > |\text{input}|$$

Example

$$\begin{aligned} S &\rightarrow 0S1 \mid 0S0S \mid T \\ T &\rightarrow S \mid 0 \end{aligned}$$


$$S \rightarrow 0S1 \mid 0S0S \mid 0$$

input: 0011

conclusion:
0011 $\notin L$

$S \Rightarrow 0 \text{ } \times$

$\Rightarrow 0S1 \Rightarrow 001 \text{ } \times$
 $00S11$ too long
 $00S0S1$ too long

$\Rightarrow 0S0S \Rightarrow 000S \Rightarrow 0000 \text{ } \times$
 $00S10S$ too long $0000S1$ too long
 $00S0S0S$ too long $0000S0S$ too long

Problems

- ① Trying all derivations may take a very long time
- ② If input is not in the language, parsing will never stop

Preparations

A faster way to parse:
Cocke-Younger-Kasami algorithm

- To use it we must prepare the CFG and convert it to Chomsky Normal Form

Step 1. Eliminate ϵ productions

Step 2. Eliminate unit productions

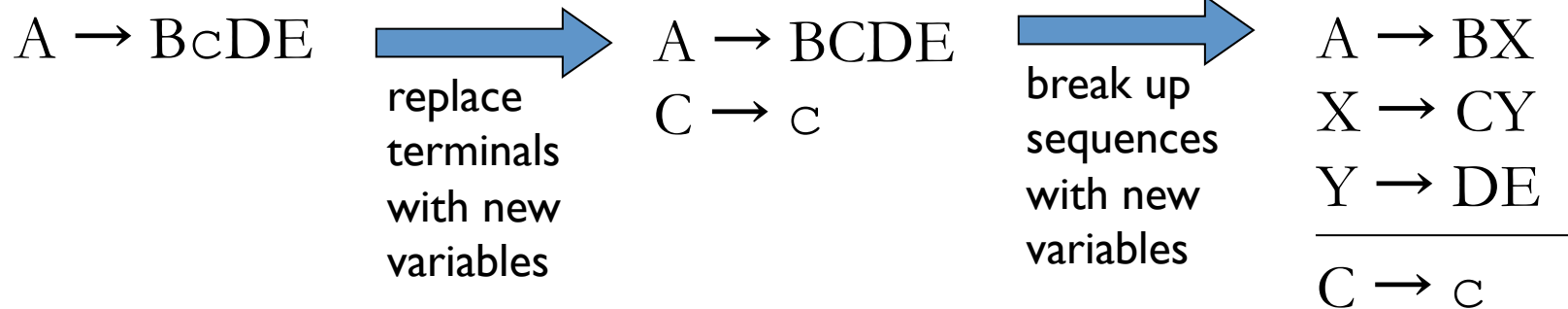
Step 3. Add rules for terminals and split longer sequences of non-terminals (next slide)

Chomsky Normal Form

- A CFG is in **Chomsky Normal Form** if every production* has the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

3rd step of conversion to Chomsky Normal Form:



Noam Chomsky

* **Exception:** We allow $S \rightarrow \varepsilon$ for the start variable only

Cocke-Younger-Kasami (CYK) algorithm

Grammar without ϵ and unit productions
in Chomsky Normal Form

Input string

For all cells in last row

If there is a production $A \rightarrow x_i$

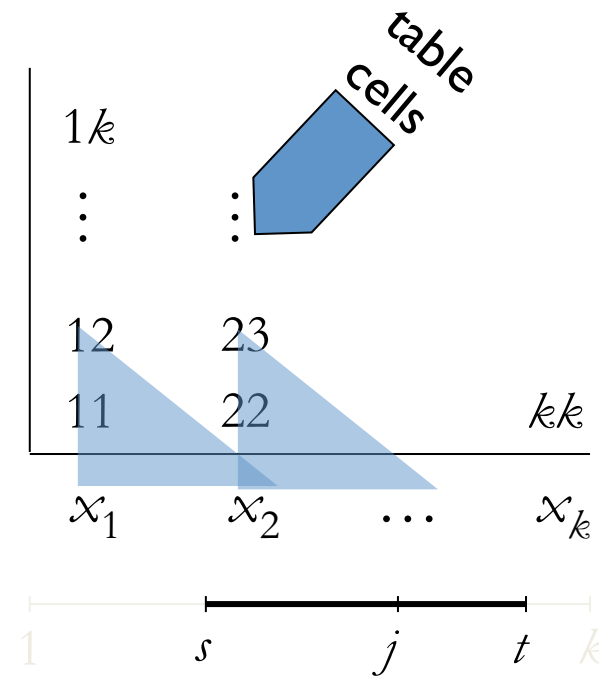
Put A in table cell ii

For cells st in other rows (going bottom-up)

If there is a production $A \rightarrow BC$

where B is in cell sj and C is in cell $(j+1)t$

Put A in cell st



Cell st remembers **all variables able to generate** substring $x_s \dots x_t$

Cocke-Younger-Kasami algorithm

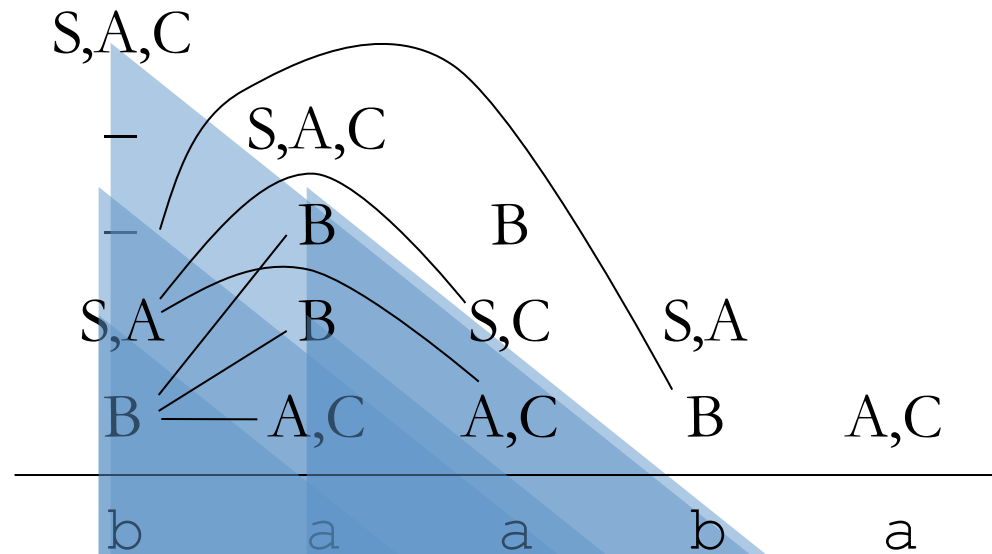
$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

$x = \text{baaba}$



Idea: We generate each substring of x bottom up

Parse tree reconstruction

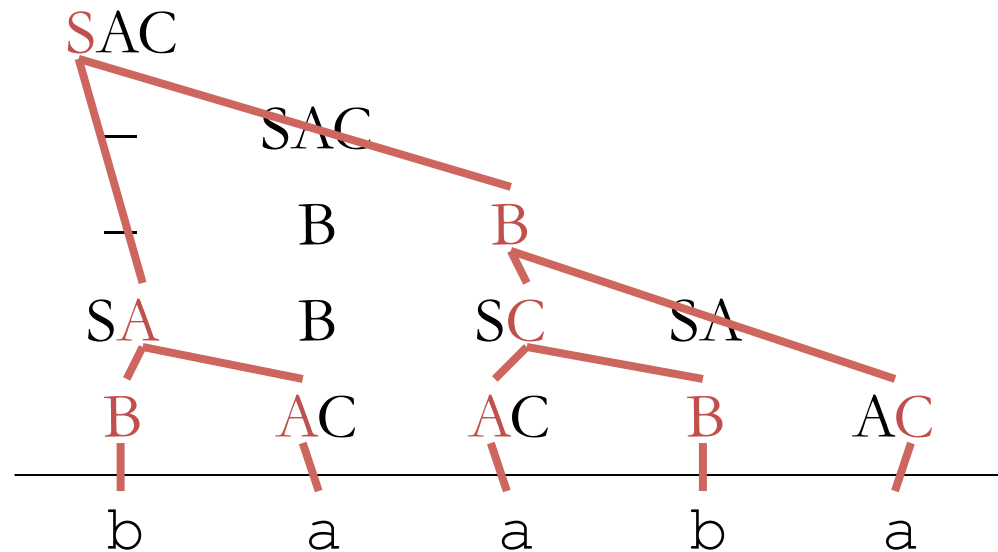
$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

$x = \text{baaba}$



Tracing back the derivations, we obtain the parse tree

Number of different parse trees

Grammar without ε and unit productions
in Chomsky Normal Form

Input string

For all cells in last row

If there is a production $A \rightarrow x_i$

Put $A : 1$ in table cell ii

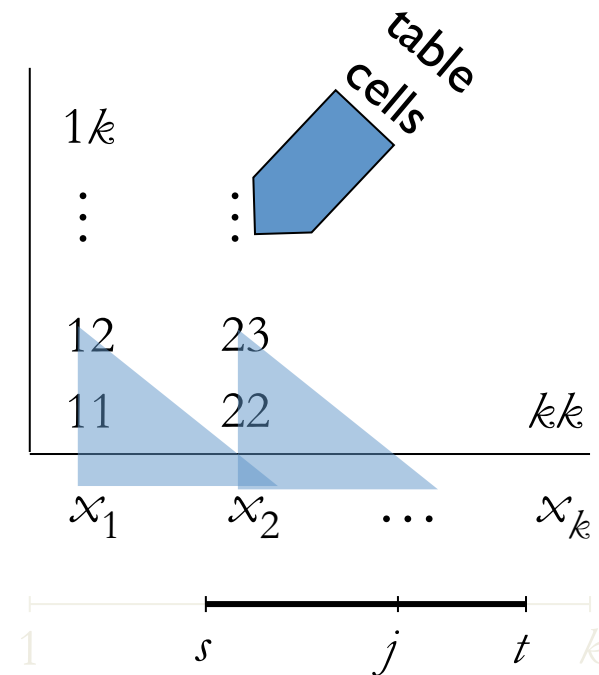
For cells st in other rows (going bottom-up)

If there is a production $A \rightarrow BC$

where $B : n_B$ is in cell sj and $C : n_C$ is in cell $(j+1)t$ and $A : n_A$ in cell st (if not present assume $n_A = 0$)

Update $A : n_A + n_B * n_C$ in cell st

Cell st remembers for each variable the number of parse trees
generating substring $x_s \dots x_t$



Example: Number of parse trees

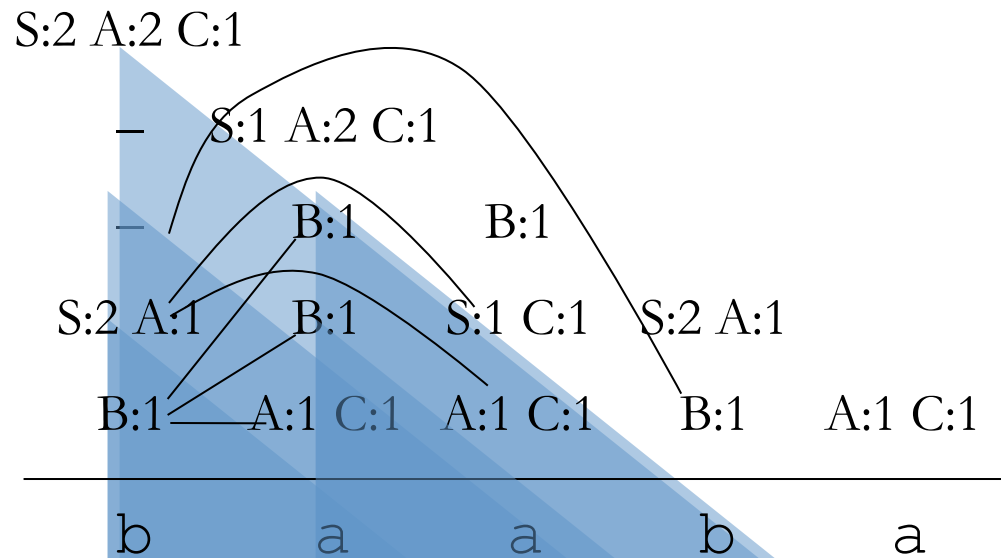
$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

$x = \text{baaba}$



Probabilistic Context Free Grammars (PCFG)

- A PCFG is a probabilistic version of a CFG where each production has a probability.
- Probabilities of all productions rewriting a given non-terminal must add up to 1, defining a distribution for each non-terminal.
- String generation is now probabilistic where production probabilities are used to non-deterministically select a production for rewriting a given non-terminal.

Statistical Parsing

- Statistical parsing uses a probabilistic model of syntax in order to assign probabilities to each parse tree.
- Provides a systematic approach to resolving syntactic ambiguity.

...the notion of “probability of a sentence” is an entirely useless one, under any known interpretation of this term.

— Noam Chomsky (famous linguist)

Every time I fire a linguist, the performance of the recognizer improves.

— Fred Jelinek
(former head of IBM speech recognition group)

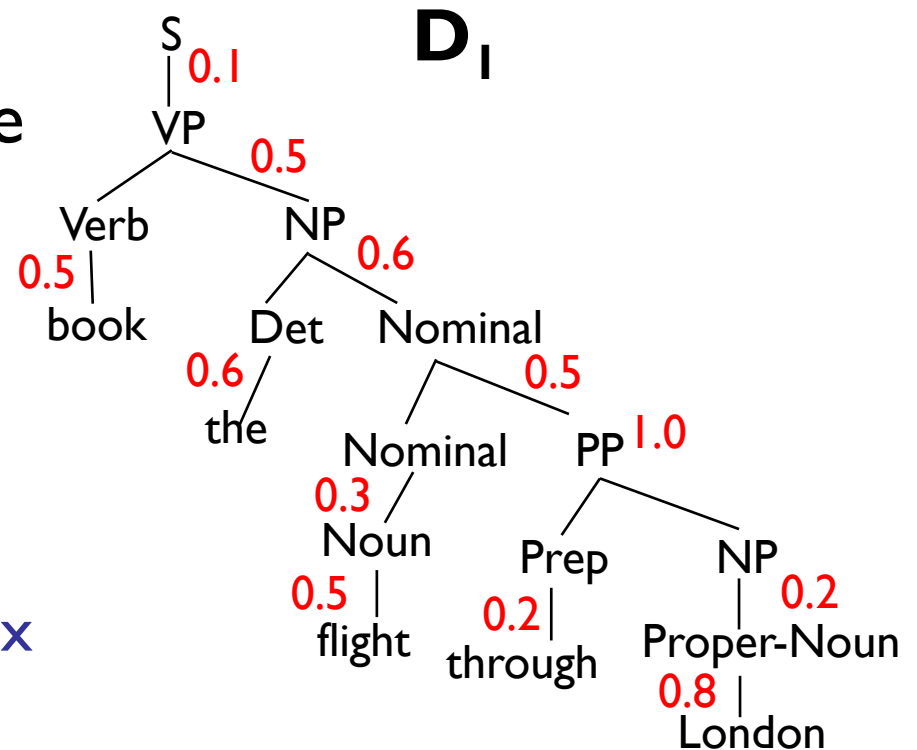
Simple PCFG for English

Grammar	Prob	Lexicon
$S \rightarrow NPVP$	0.8	$Det \rightarrow the \mid a \mid that \mid this$
$S \rightarrow Aux \ NPVP$	0.1	0.6 0.2 0.1 0.1
$S \rightarrow VP$	0.1	$Noun \rightarrow book \mid flight \mid meal \mid money$
$NP \rightarrow Pronoun$	0.2	0.1 0.5 0.2 0.2
$NP \rightarrow Proper-Noun$	0.2	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Det \ Nominal$	0.6	0.5 0.2 0.3
$Nominal \rightarrow Noun$	0.3	$Pronoun \rightarrow I \mid he \mid she \mid me$
$Nominal \rightarrow Nominal \ Noun$	0.2	0.5 0.1 0.1 0.3
$Nominal \rightarrow Nominal \ PP$	0.5	$Proper-Noun \rightarrow London \mid BA$
$VP \rightarrow Verb$	0.2	0.8 0.2
$VP \rightarrow Verb \ NP$	0.5	$Aux \rightarrow does$
$VP \rightarrow VP \ PP$	0.3	1.0
$PP \rightarrow Prep \ NP$	1.0	$Prep \rightarrow from \mid to \mid on \mid near \mid through$
		0.25 0.25 0.1 0.2 0.2

Parse Tree Probability

- Assume productions for each node are chosen independently.
- Probability of a parse tree is the product of the probabilities of its productions.

$$\begin{aligned} P(D_1) &= 0.1 \times 0.5 \times 0.5 \times 0.6 \times 0.6 \times \\ &\quad 0.5 \times 0.3 \times 1.0 \times 0.2 \times 0.2 \times \\ &\quad 0.5 \times 0.8 \\ &= 0.0000216 \end{aligned}$$

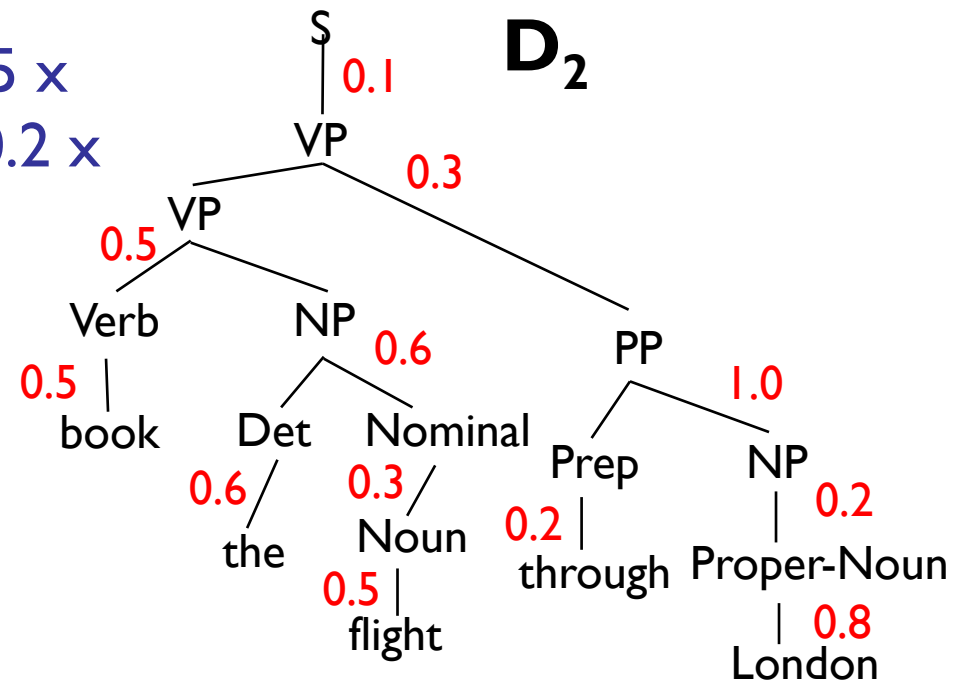


Syntactic Disambiguation

- Resolve ambiguity by picking most probable parse tree.

$$\begin{aligned} P(D_2) &= 0.1 \times 0.3 \times 0.5 \times 0.6 \times 0.5 \times \\ &\quad 0.6 \times 0.3 \times 1.0 \times 0.5 \times 0.2 \times \\ &\quad 0.2 \times 0.8 \\ &= 0.00001296 \end{aligned}$$

In this case D_1 is the best interpretation, because $0.0000216 > 0.00001296$



Sentence Probability

- Probability of a given sentence is the sum of the probabilities of all of its parse trees.

$$\begin{aligned} P(\text{"book the flight through London"}) &= \\ P(D_1) + P(D_2) &= 0.0000216 + 0.00001296 \\ &= 0.00003456 \end{aligned}$$

Probabilistic CYK:

Most Likely Parse Tree

- CYK can be modified for PCFG parsing by including in each cell a probability for each non-terminal.
- Cell ij contains the **most probable** parse tree of each non-terminal that can generate the part of the input word from i through j together with its associated probability.
- When transforming the grammar to CNF, we set the production probabilities to preserve the probability of derivations.

Probabilistic Grammar Conversion

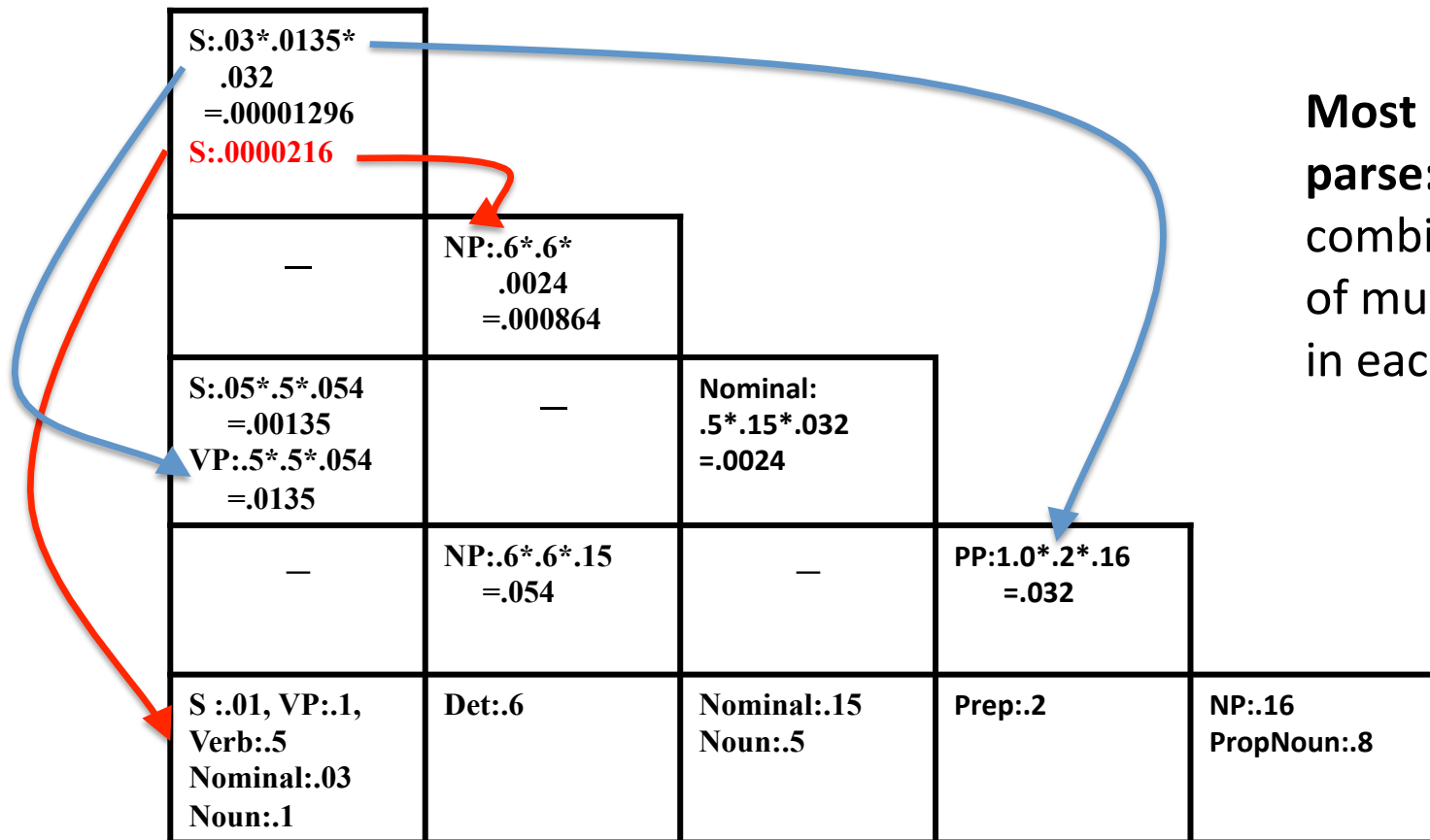
Original Grammar

$S \rightarrow \text{NP VP}$	0.8
$S \rightarrow \text{Aux NP VP}$	0.1
$S \rightarrow \text{VP}$	0.1
$\text{NP} \rightarrow \text{Pronoun}$	0.2
$\text{NP} \rightarrow \text{Proper-Noun}$	0.2
$\text{NP} \rightarrow \text{Det Nominal}$	0.6
$\text{Nominal} \rightarrow \text{Noun}$	0.3
$\text{Nominal} \rightarrow \text{Nominal Noun}$	0.2
$\text{Nominal} \rightarrow \text{Nominal PP}$	0.5
$\text{VP} \rightarrow \text{Verb}$	0.2
$\text{VP} \rightarrow \text{Verb NP}$	0.5
$\text{VP} \rightarrow \text{VP PP}$	0.3
$\text{PP} \rightarrow \text{Prep NP}$	1.0

Chomsky Normal Form

$S \rightarrow \text{NP VP}$	0.8
$S \rightarrow \text{XI VP}$	0.1
$\text{XI} \rightarrow \text{Aux NP}$	1.0
$S \rightarrow \text{book} \mid \text{include} \mid \text{prefer}$ 0.01 0.004 0.006	
$S \rightarrow \text{Verb NP}$	0.05
$S \rightarrow \text{VP PP}$	0.03
$\text{NP} \rightarrow \text{I} \mid \text{he} \mid \text{she} \mid \text{me}$ 0.1 0.02 0.02 0.06	
$\text{NP} \rightarrow \text{London} \mid \text{BA}$ 0.16 .04	
$\text{NP} \rightarrow \text{Det Nominal}$	0.6
$\text{Nominal} \rightarrow \text{book} \mid \text{flight} \mid \text{meal} \mid \text{money}$ 0.03 0.15 0.06 0.06	
$\text{Nominal} \rightarrow \text{Nominal Noun}$	0.2
$\text{Nominal} \rightarrow \text{Nominal PP}$	0.5
$\text{VP} \rightarrow \text{book} \mid \text{include} \mid \text{prefer}$ 0.1 0.04 0.06	
$\text{VP} \rightarrow \text{Verb NP}$	0.5
$\text{VP} \rightarrow \text{VP PP}$	0.3
$\text{PP} \rightarrow \text{Prep NP}$	1.0

Probabilistic CYK Parser: Most Likely Parse Tree



Most probable parse: use max to combine probabilities of multiple parse trees in each cell.

Book

the

flight

through

London

Probabilistic CYK Parser: Sentence Probability

S: .00001296 +.0000216 =.00003456				
—	NP:.6*.6* .0024 =.000864			
S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	—	Nominal: .5*.15*.032 =.0024		
—	NP:.6*.6*.15 =.054	—	PP:1.0*.2*.16 =.032	
S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	Det:.6	Nominal:.15 Noun:.5	Prep:.2	NP:.16 PropNoun:.8

Book the flight through London

The probability of generating a given sentence is the sum all probabilities of multiple parse trees in each cell.

More precisely

Sentence
probability

For all cells in last row

If there is a production $A \rightarrow x_i$ with probability p

Put $A : p$ in table cell ii

For cells st in other rows (going bottom-up)

If there is a production $A \rightarrow BC$ with probability p

where $B : p_B$ is in cell sj and $C : p_C$ is in cell $(j+1)t$ and

$A : p_A$ in cell st (if not present assume $p_A = 0$)

Update $A : p_A + p * p_B * p_C$ in cell st

Most likely
parse tree

For all cells in last row

If there is a production $A \rightarrow x_i$ with probability p

Put $A : p$ in table cell ii

For cells st in other rows (going bottom-up)

If there is a production $A \rightarrow BC$ with probability p

where $B : p_B$ is in cell sj and $C : p_C$ is in cell $(j+1)t$ and

$A : p_A$ in cell st (if not present assume $p_A = 0$)

Update $A : \max(p_A, p * p_B * p_C)$ in cell st

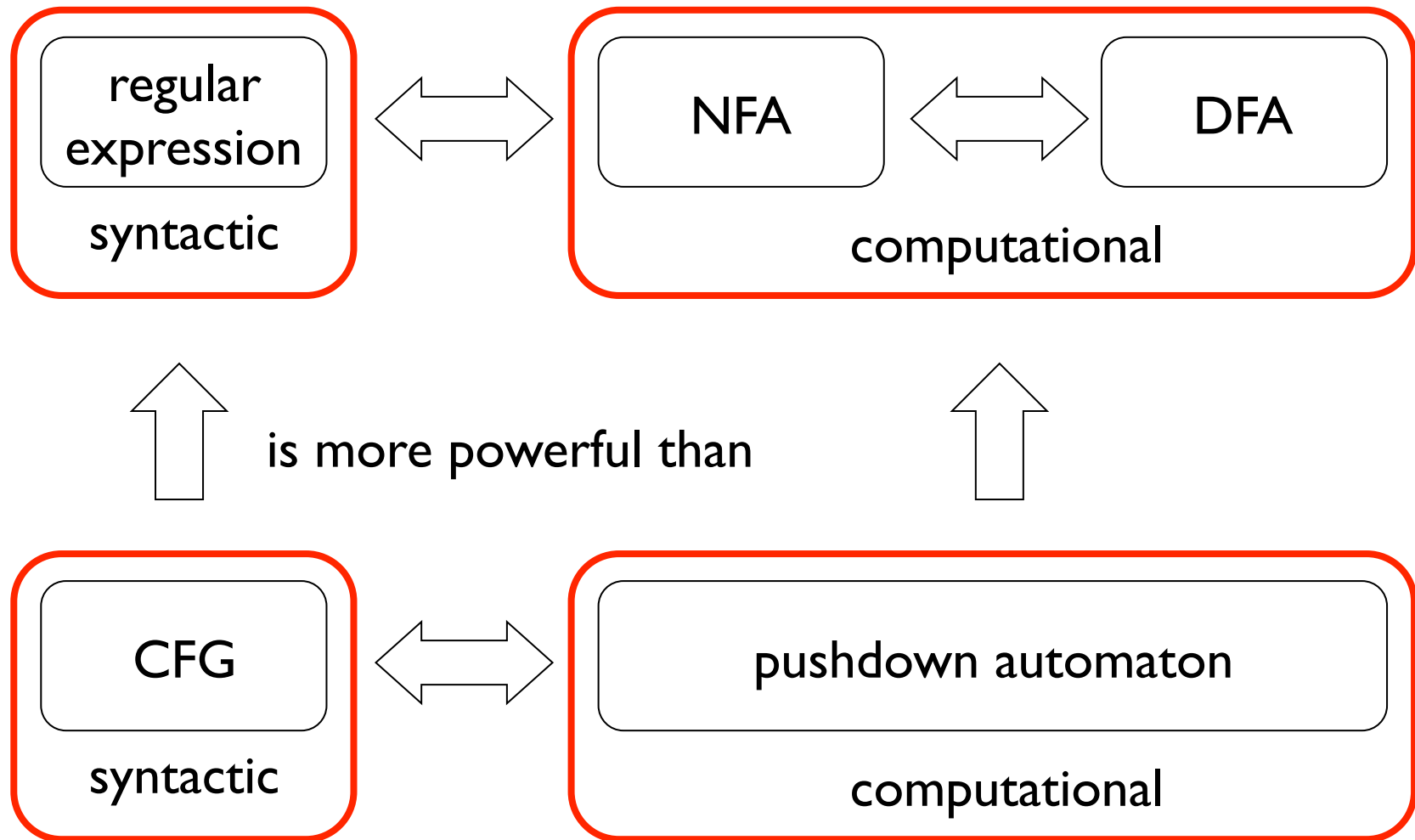
PCFGs and NLP

PCFG is **consistent**
if and only if
it generates a finite word with probability 1.

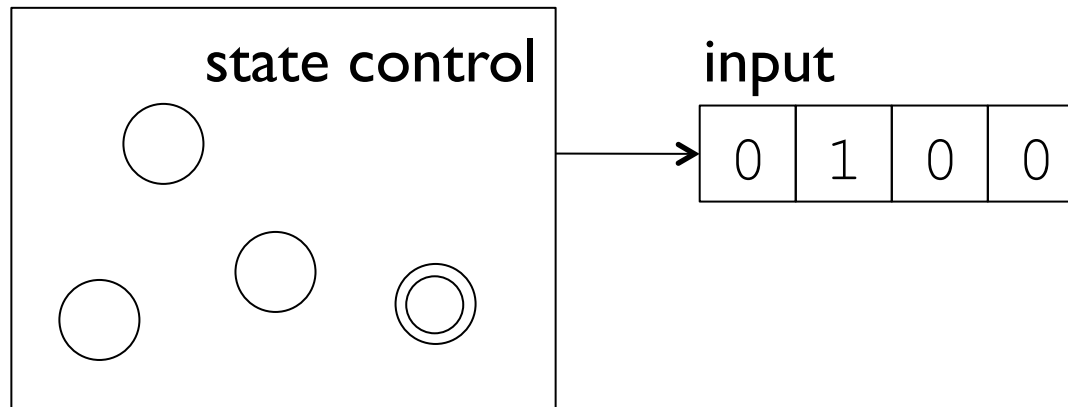
name	#prod	max-scc	Jacobi	Gauss Seidel	SOR $\omega=1.05$	DNewton	SNewton
brown	22866 ✗	448	312.084(9277)	275.624(7866)	diverge	2.106(8)	2.115(9)
lemonde	32885 ✓	527	234.715(5995)	30.420(767)	diverge	1.556(7)	2.037(7)
negra	29297 ✓	518	16.995(610)	4.724(174)	4.201(152)	1.017(6)	0.499(6)
swbd	47578 ✗	1123	445.120(4778)	19.321(202)	25.654(270)	6.435(6)	3.978(6)
tiger	52184 ✓	1173	99.286(1347)	16.073(210)	12.447(166)	5.274(6)	1.871(6)
tuebadz	8932 ✓	293	6.894(465)	1.925(133)	6.878(461)	0.477(7)	0.341(7)
wsj	31170 ✓	765	462.378(9787)	68.650(1439)	diverge	2.363(7)	3.616(8)

Pushdown automata

Motivation

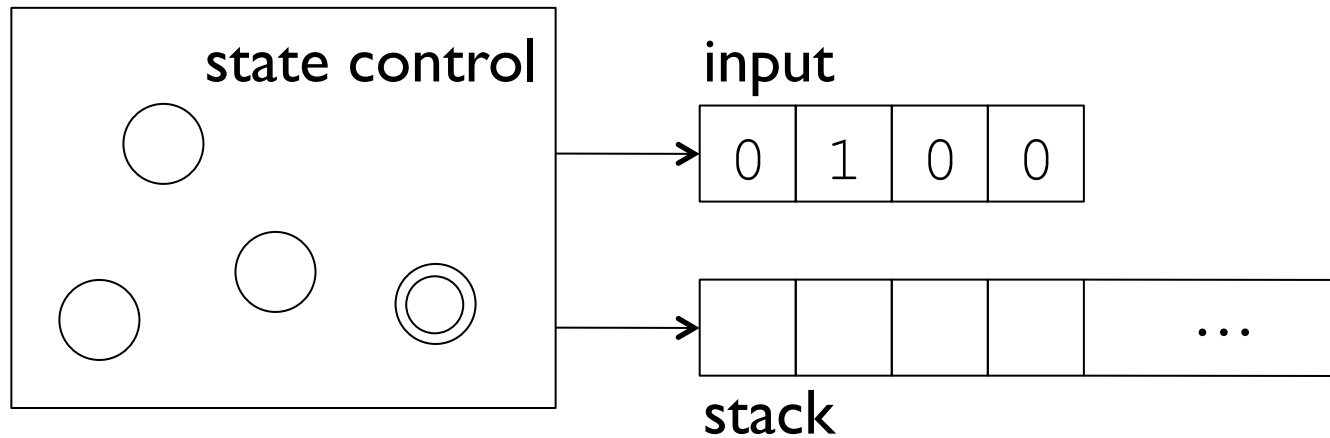


Pushdown automata versus NFA



NFA

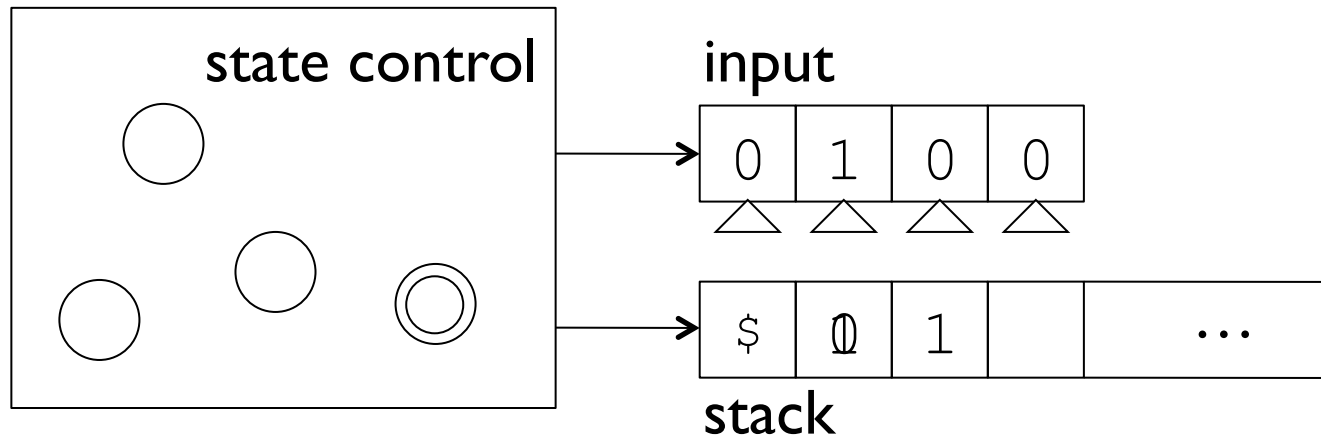
Pushdown automata



pushdown automaton (PDA)

A PDA is like an NFA with but with an infinite **stack**

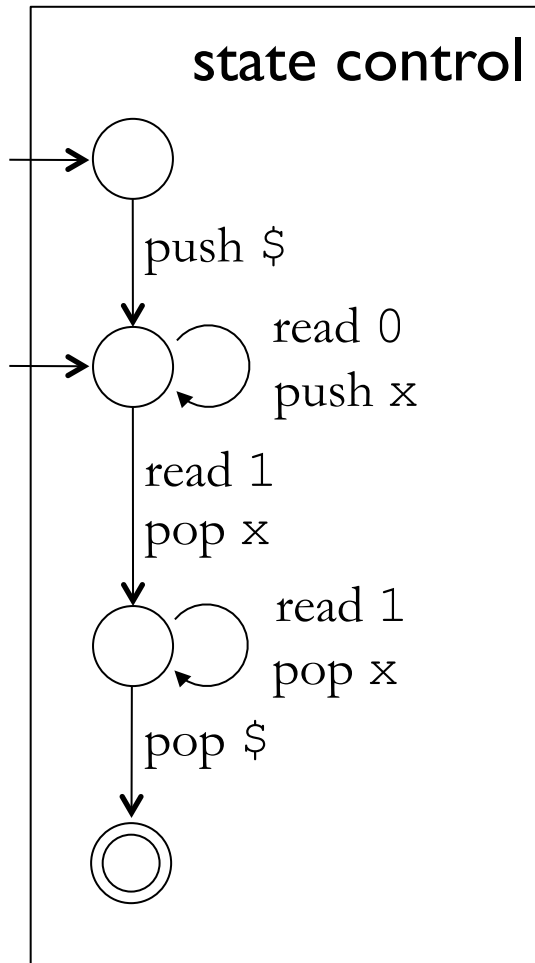
Pushdown automata



pushdown automaton (PDA)

As the PDA is reading the input, it can
push / pop symbols from the top of the stack

Building a PDA



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

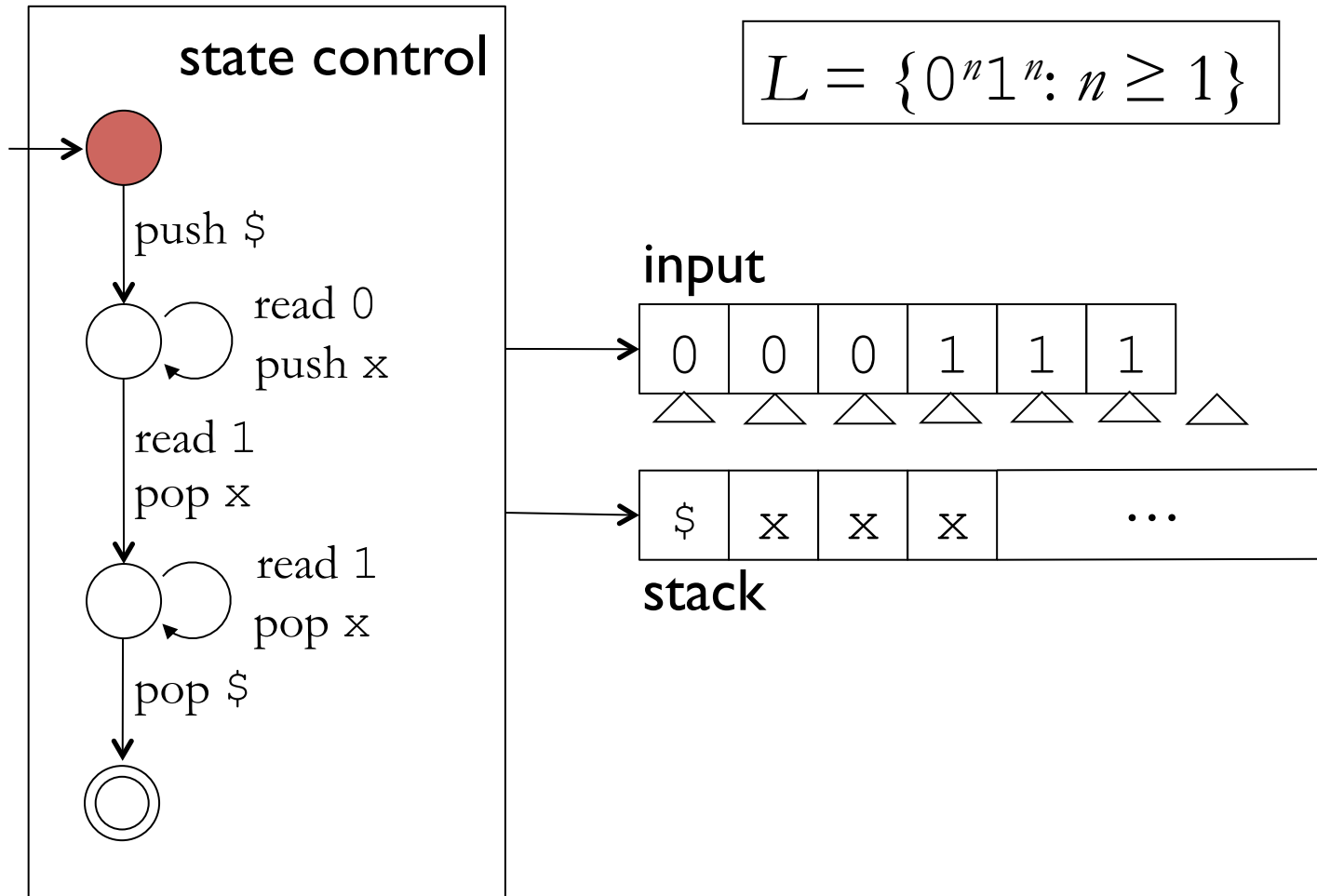
We remember each 0
by **pushing** x onto the stack

When we see a 1, we
pop an x from the stack

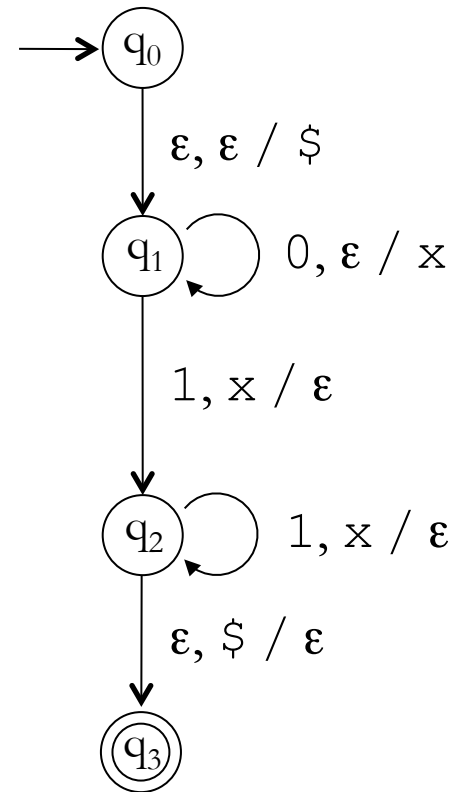
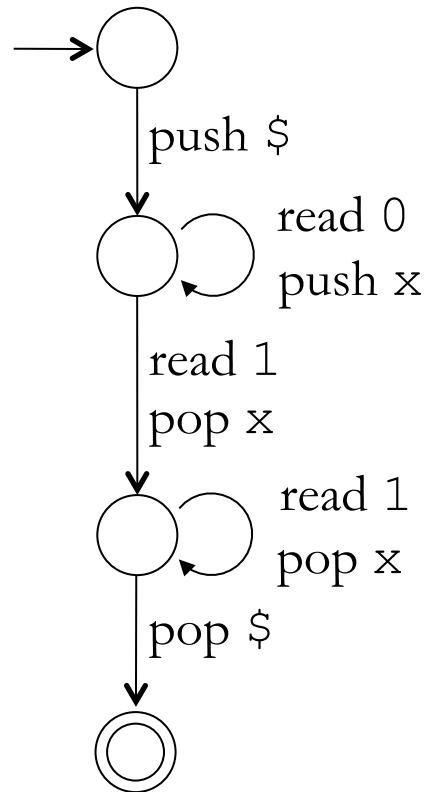
We want to accept when
we hit the **stack bottom**

We will use \$ to **mark** bottom

A PDA in action



Notation for PDAs



read, pop / push

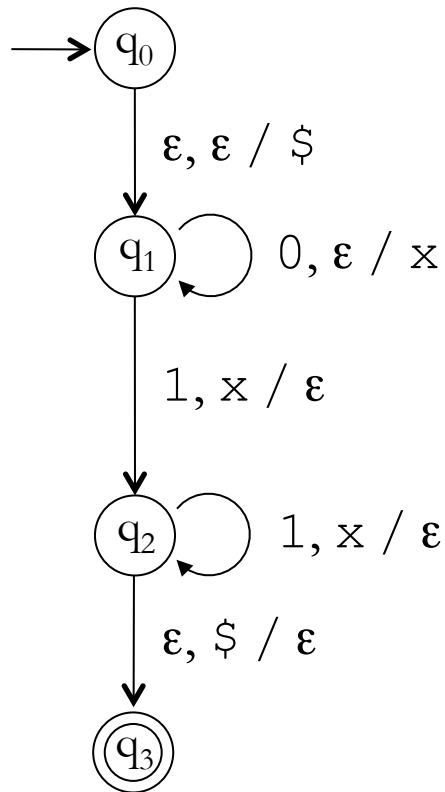
Definition of a PDA

A pushdown automaton is $(Q, \Sigma, \Gamma, \delta, q_0, F)$ where:

- Q is a finite set of **states**;
- Σ is the **input alphabet**;
- Γ is the **stack alphabet**;
- q_0 in Q is the **initial state**;
- $F \subseteq Q$ is a set of **final states**;
- **δ** is the **transition function**.

$$\delta: \underset{\text{state}}{Q} \times \underset{\text{input symbol}}{(\Sigma \cup \{\varepsilon\})} \times \underset{\text{pop symbol}}{(\Gamma \cup \{\varepsilon\})} \rightarrow \text{subsets of } \underset{\text{state}}{Q} \times \underset{\text{push symbol}}{(\Gamma \cup \{\varepsilon\})}$$

Example



$$\Sigma = \{0, 1\}$$

$$\Gamma = \{\$, x\}$$

$$\delta(q_0, \epsilon, \epsilon) = \{(q_1, \$)\}$$

$$\delta(q_0, \epsilon, \$) = \emptyset$$

$$\delta(q_0, \epsilon, x) = \emptyset$$

$$\delta(q_0, 0, \epsilon) = \emptyset$$

...

$$\delta: \underset{\text{state}}{Q} \times \underset{\text{input symbol}}{(\Sigma \cup \{\epsilon\})} \times \underset{\text{pop symbol}}{(\Gamma \cup \{\epsilon\})} \rightarrow \text{subsets of } \underset{\text{state}}{Q} \times \underset{\text{push symbol}}{(\Gamma \cup \{\epsilon\})}$$

The language of a PDA

- A PDA is nondeterministic, i.e. multiple transitions on same pop/input are allowed
- Transitions **may** but **do not have to** push, pop or read input

The **language** of a PDA is the set of all strings in Σ^* that **can lead** the PDA to an accepting state after the whole input is read.

Remark: Sometimes acceptance is defined by emptying the stack or even both: empty stack and accepting state. All of these variants correspond to the same class of languages.

Example I

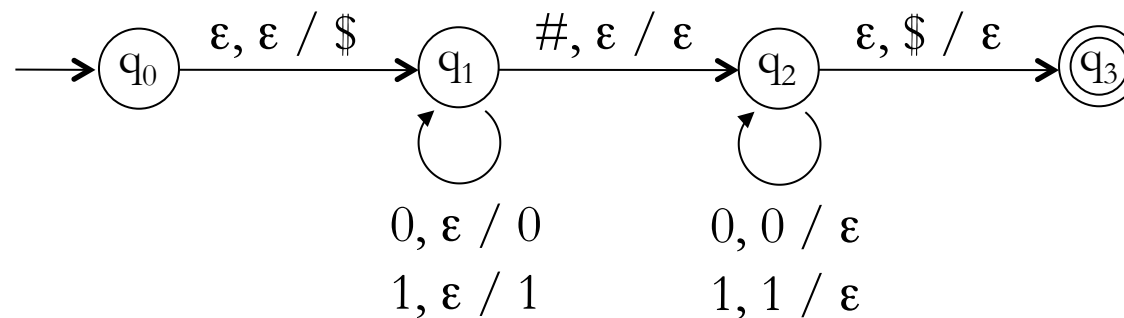
$$L_1 = \{w\#w^R: w \in \{0, 1\}^*\}$$

$$\Sigma = \{0, 1, \#\}$$

$$\Gamma = \{\$, 0, 1\}$$

$$\#, 0\#0, 01\#10 \in L_1$$

$$\epsilon, 01\#1, 0\#\#0 \notin L_1$$



write w on stack
 read w^R from stack

Example 2

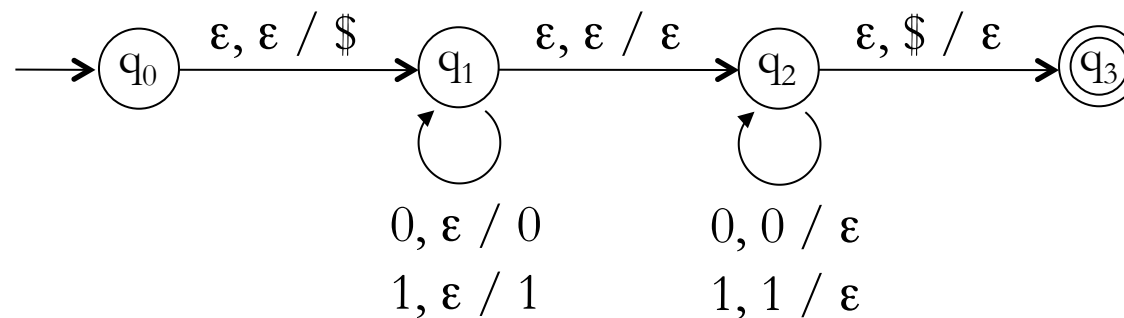
$$L_2 = \{ww^R: w \in \Sigma^*\}$$

$$\Sigma = \{0, 1\}$$

$$\Gamma = \{\$, 0, 1\}$$

$$\epsilon, 00, 0110 \in L_2$$

$$1, 011, 010 \notin L_2$$



guess the middle of a string

Example 3

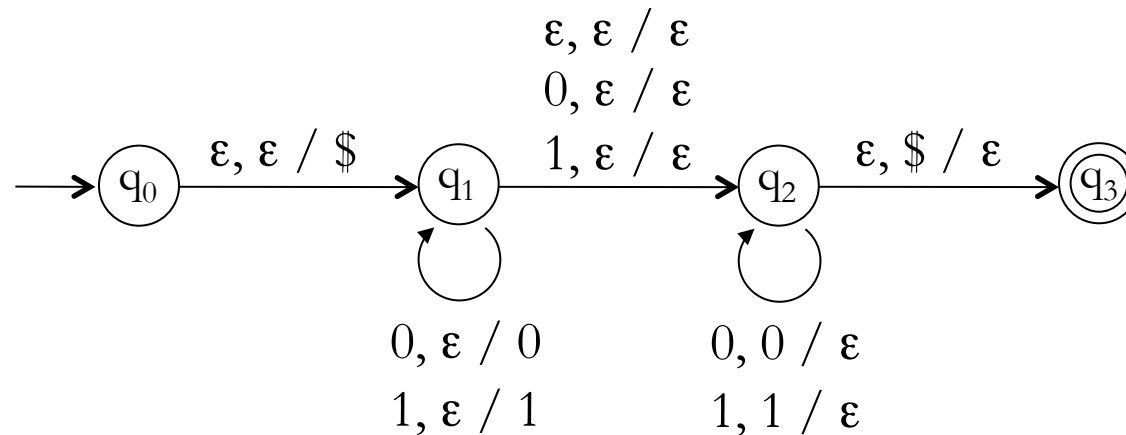
$$L_3 = \{w: w = w^R, w \in \Sigma^*\}$$

$$\Sigma = \{0, 1\}$$

$$\varepsilon, 1, 00, 010, 0110 \in L_3$$

$$011 \notin L_3$$

$$\underbrace{011011}_{x} \underbrace{10110}_{x^R} \text{ or } \underbrace{01101}_{x} \underbrace{10110}_{x^R}$$

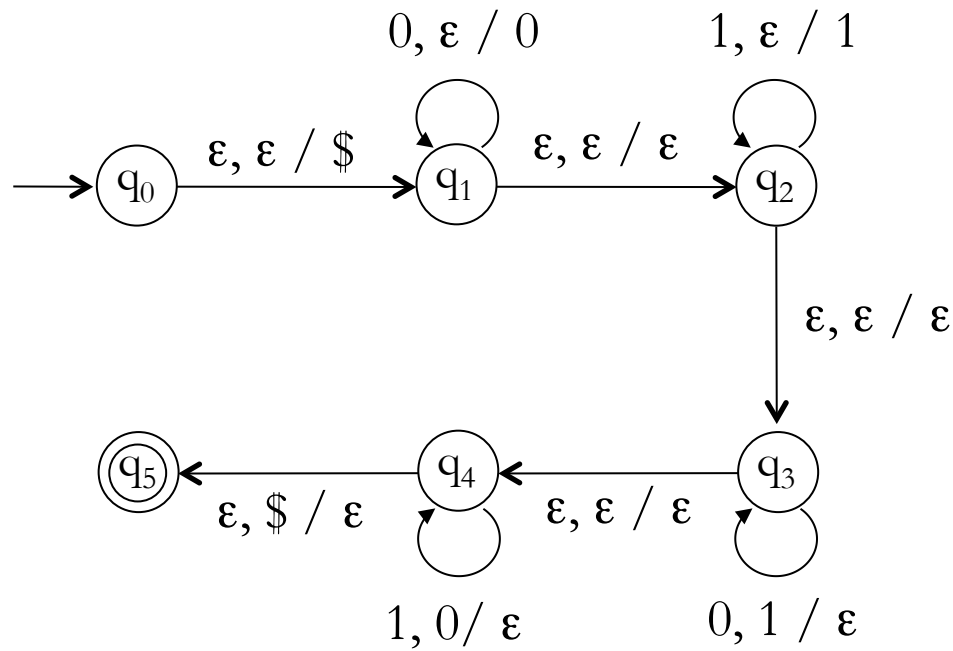


the middle symbols can be ε , 0, or 1

Example 4

$$L_4 = \{0^n 1^m 0^m 1^n \mid n \geq 0, m \geq 0\}$$

$$\Sigma = \{0, 1\}$$



input: $0^n 1^m 0^m 1^n$

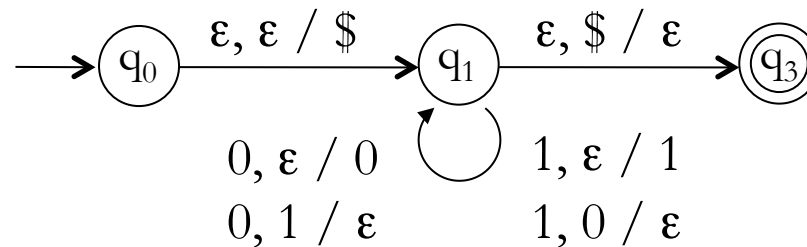
stack: $0^n 1^m$

Example 5

$$L_5 = \{w: w \text{ has the same number of 0s and 1s}\}$$

$$\Sigma = \{0, 1\}$$

Strategy: Stack keeps track of **excess** of 0s or 1s
If at the end, stack is empty, the numbers are equal

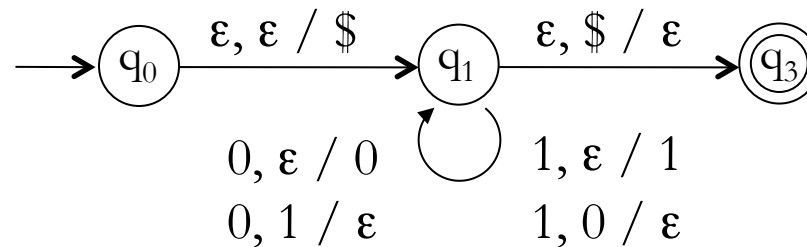


Example 5

$$L_5 = \{w: w \text{ has the same number 0s and 1s}\} \quad \Sigma = \{0, 1\}$$

Invariant: In every execution of the PDA:

$$\#1 - \#0 \text{ on stack} = \#1 - \#0 \text{ in input so far}$$



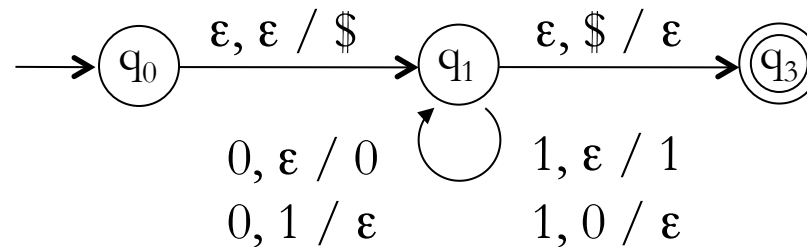
If w is not in L_5 , it must be rejected

Example 5

$$L_5 = \{w: w \text{ has the same number 0s and 1s}\} \quad \Sigma = \{0, 1\}$$

Property: In some execution of the PDA:

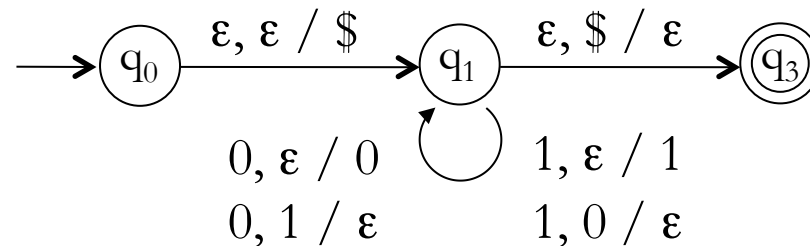
stack consists only of 0s or only of 1s (or ϵ)



If w is in L_5 , some execution will accept

Example 5

$$L_5 = \{w: w \text{ has the same number 0s and 1s}\} \quad \Sigma = \{0, 1\}$$



$w = 001110$

read	stack
0	\$0
0	\$00
1	\$0
1	\$
1	\$1
0	\$

Example 6

$$L_6 = \{w: w \text{ has two 0-blocks with the same number of 0s}\}$$

01011, 001011001, 10010101001

allowed

01001000, 01111

not allowed

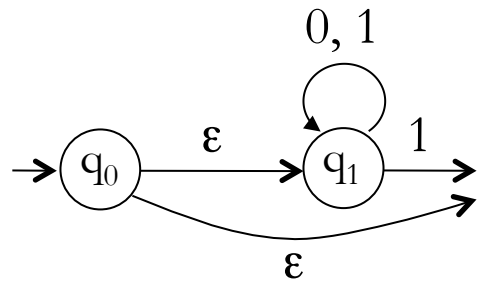
Strategy: Detect start of the first 0-block

Push 0s on the stack

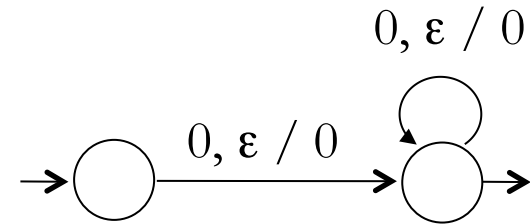
Detect start of the second 0-block

Pop 0s from the stack

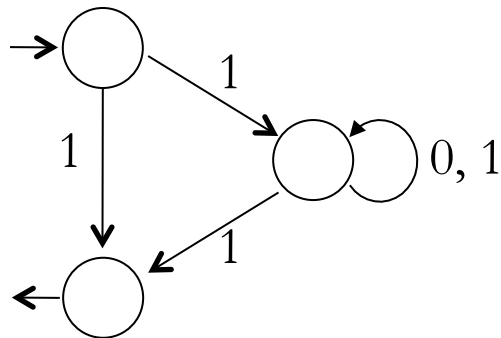
Example 6



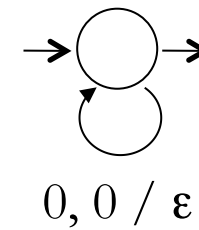
1 Detect start of first 0-block



2 Push 0s on stack



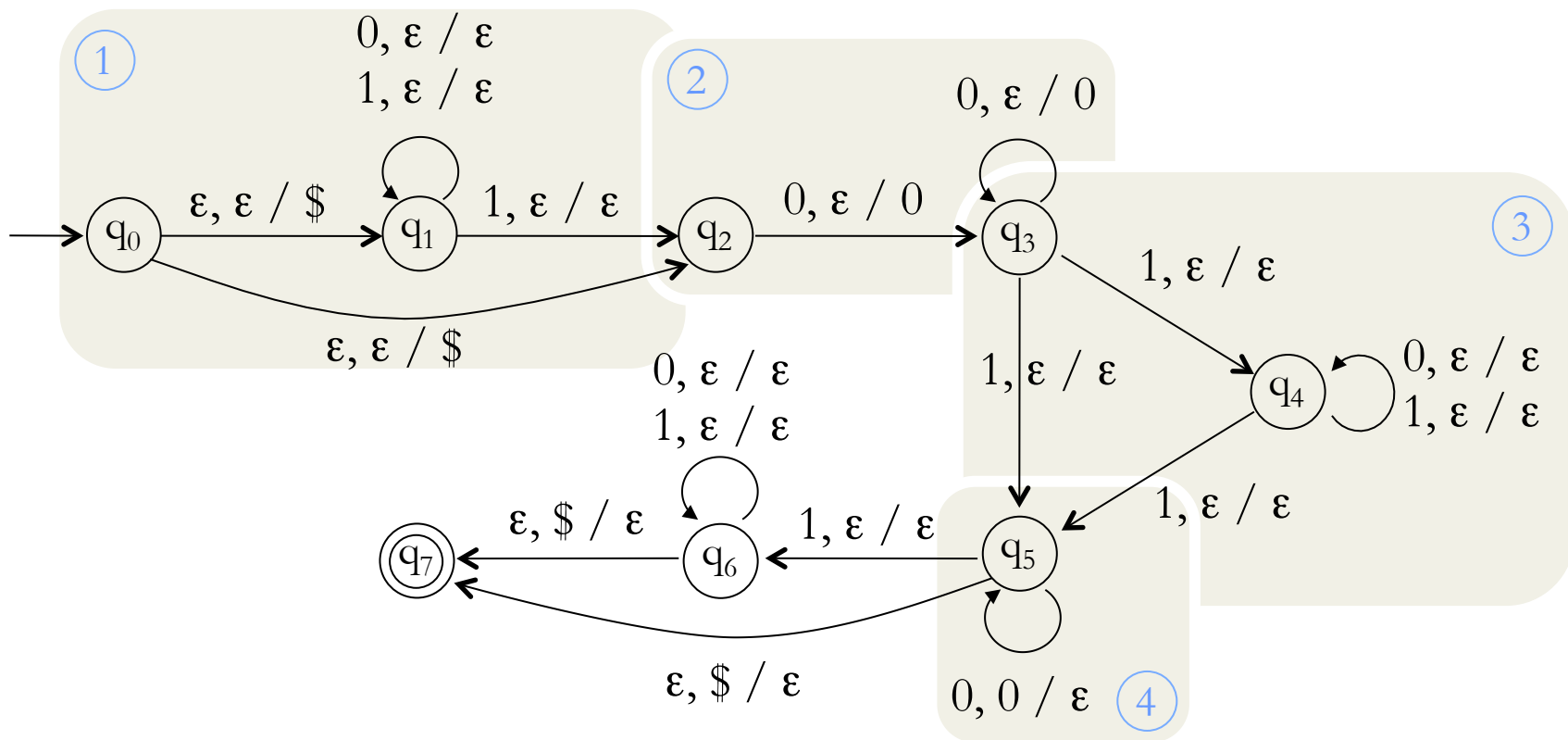
3 Detect start of second 0-block



4 Pop 0s from stack

Example 6

$$L_6 = \{w: w \text{ has two 0-blocks with the same number of 0s}\}$$

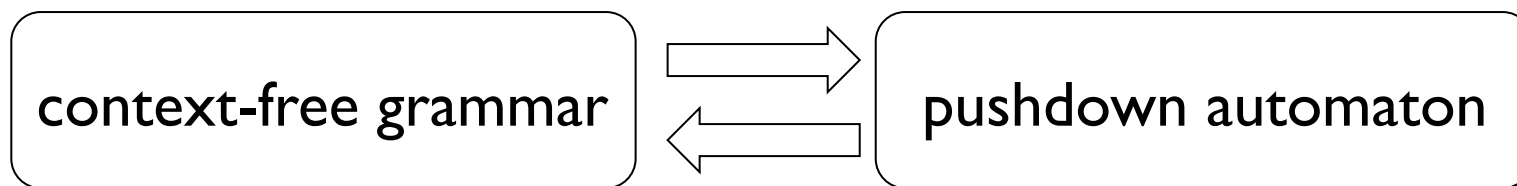


CFG \Leftrightarrow PDA conversions



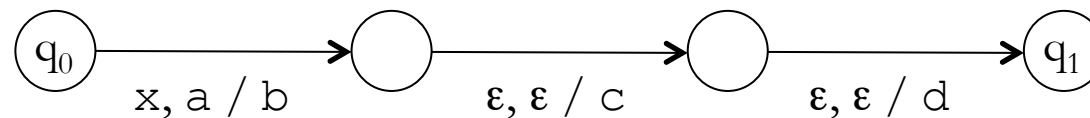
CFGs and PDAs

L has a context-free grammar **if and only if** it is accepted by some pushdown automaton.



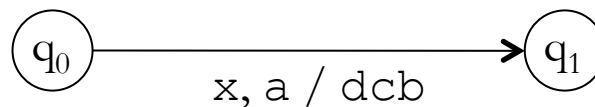
A convention

- When we have a **sequence of transitions** like:



pop a , then push b , c , and d

- We will abbreviate it like this:



replace a by dcb on the top of the stack
(notice the reverse order: the first symbol
of the word is at the top of the stack)

Converting a CFG to a PDA

- **Idea:** Use PDA to simulate derivations

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 00B11 \Rightarrow 00\#11$

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

PDA control:	stack:	input:
write start variable	\$A	00#11
replace production in reverse	\$1A0	00#11
pop terminal and match	\$1A	0#11
replace production in reverse	\$11A0	0#11
pop terminal and match	\$11A	#11
replace production in reverse	\$11B	#11

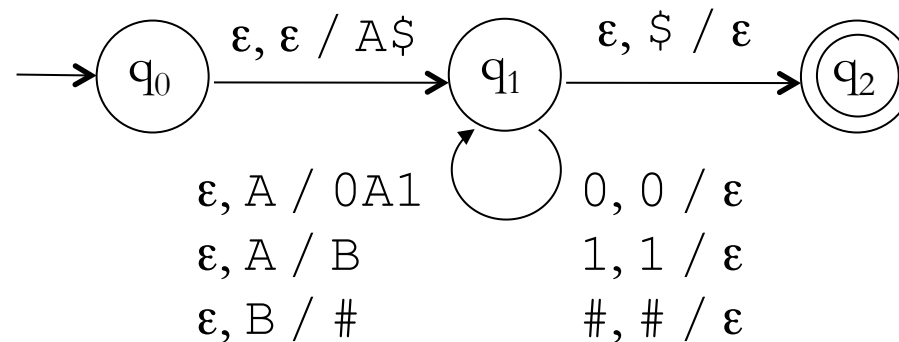
Converting a CFG to a PDA

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

CFG



input
stack

00#11
\$A

00#11
\$1A0

00#11
\$1A

00#11
\$11A0

00#11
\$11A

00#11
\$11B

00#11
\$11#

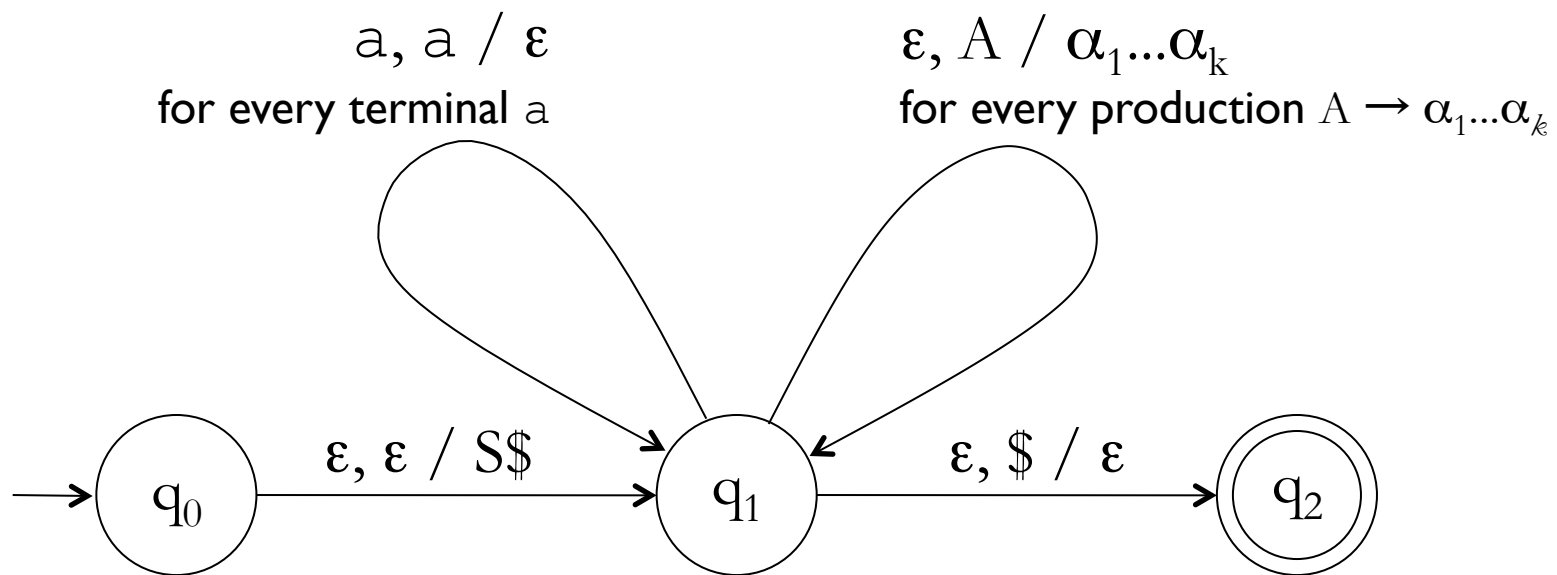
00#11
\$11

00#11
\$1

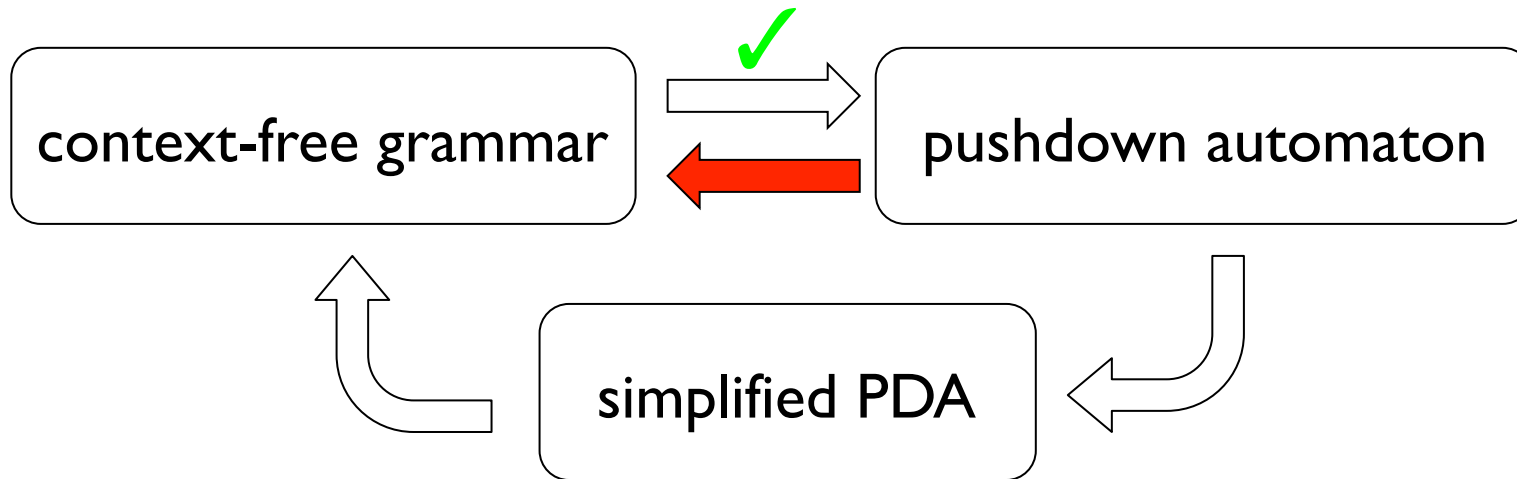
00#11
\$

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 00B11 \Rightarrow 00\#11$

General CFG to PDA conversion

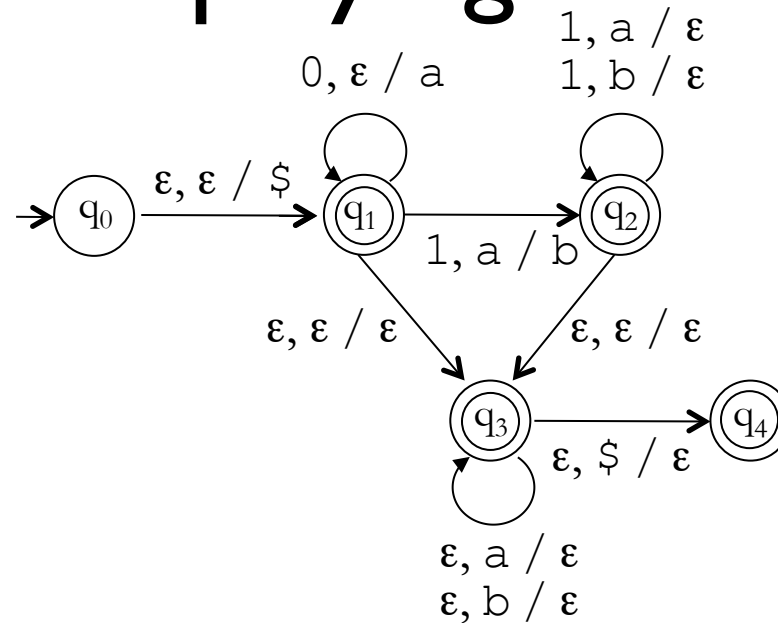


From PDAs to CFGs



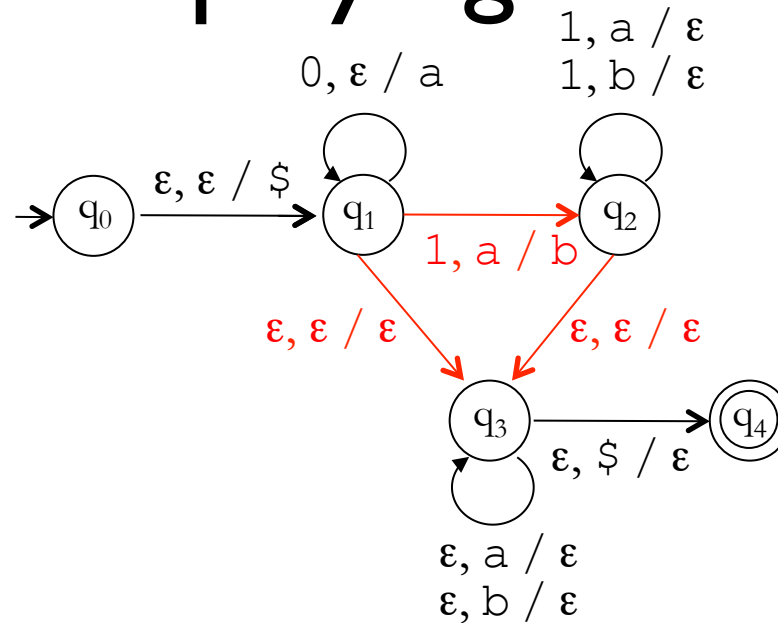
- A simplified PDA:
 - Has a **single accept state**
 - **Empties its stack** before accepting
 - Each transition is either a push, or a pop, but not both

Simplifying the PDA



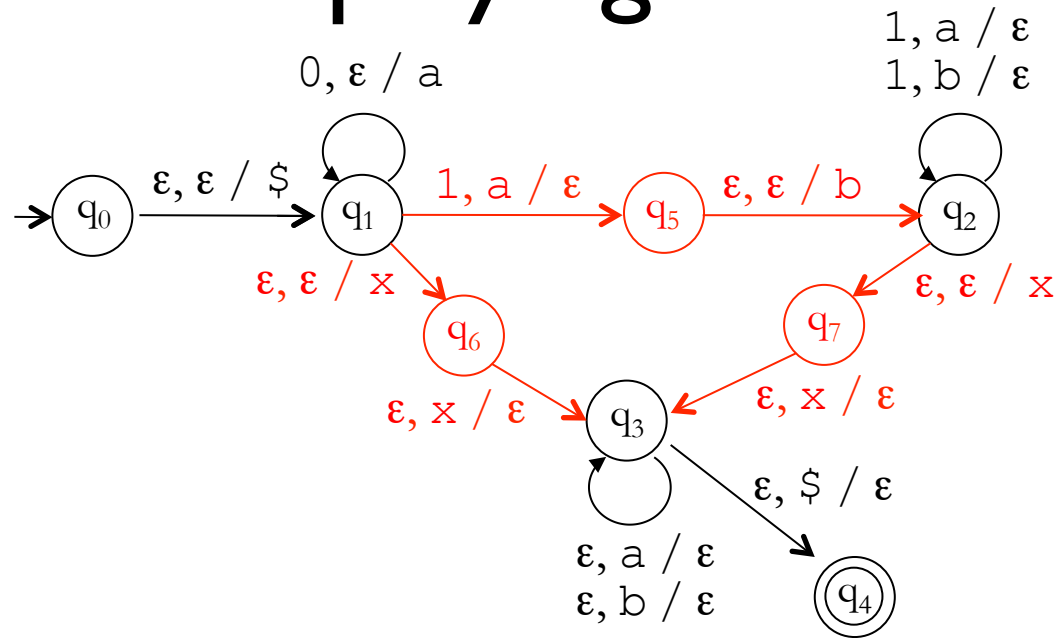
- A simplified PDA:
 - Has a **single accept state** ✓
 - **Empties its stack** before accepting ✓
 - Each transition is either a push, or a pop, but not both

Simplifying the PDA



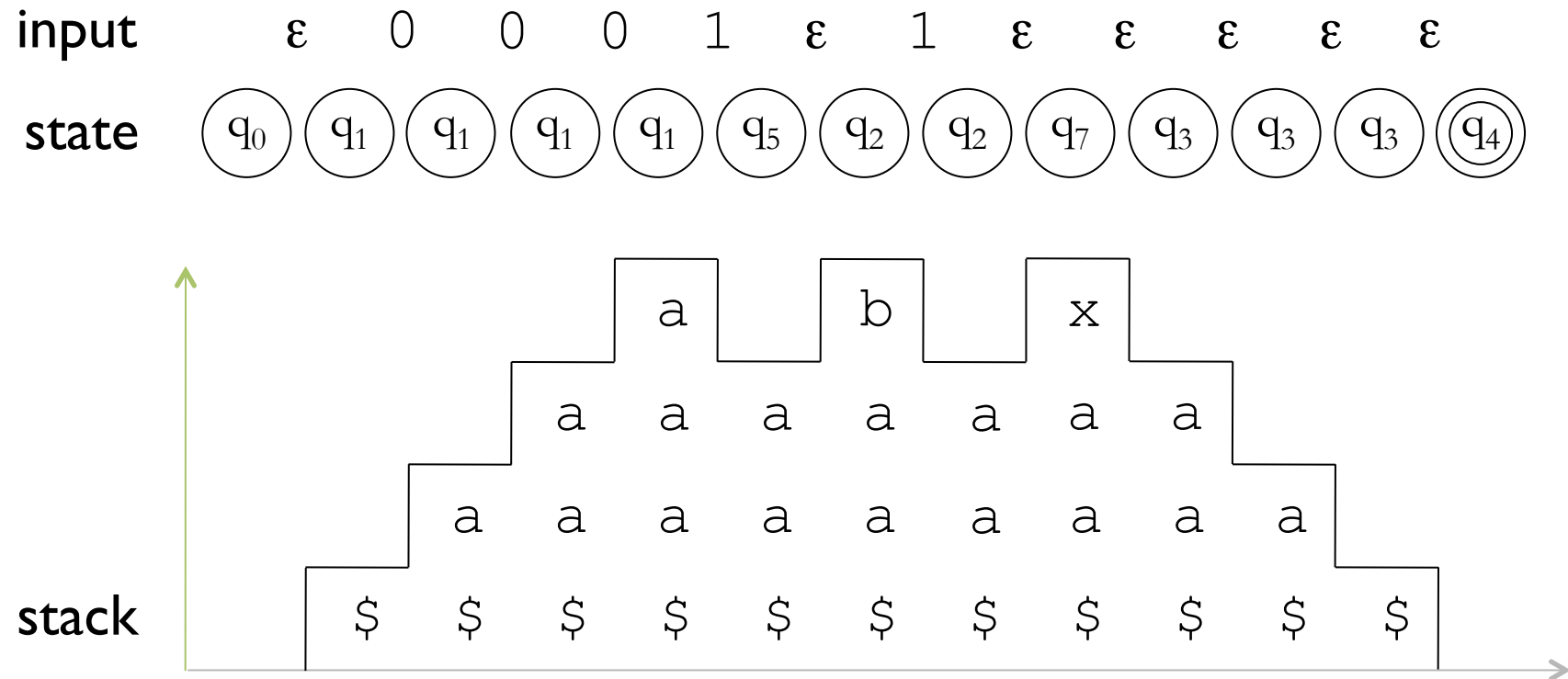
- A simplified PDA:
 - Has a **single accept state** ✓
 - **Empties its stack** before accepting ✓
 - Each transition is either a push, or a pop, but not both

Simplifying the PDA



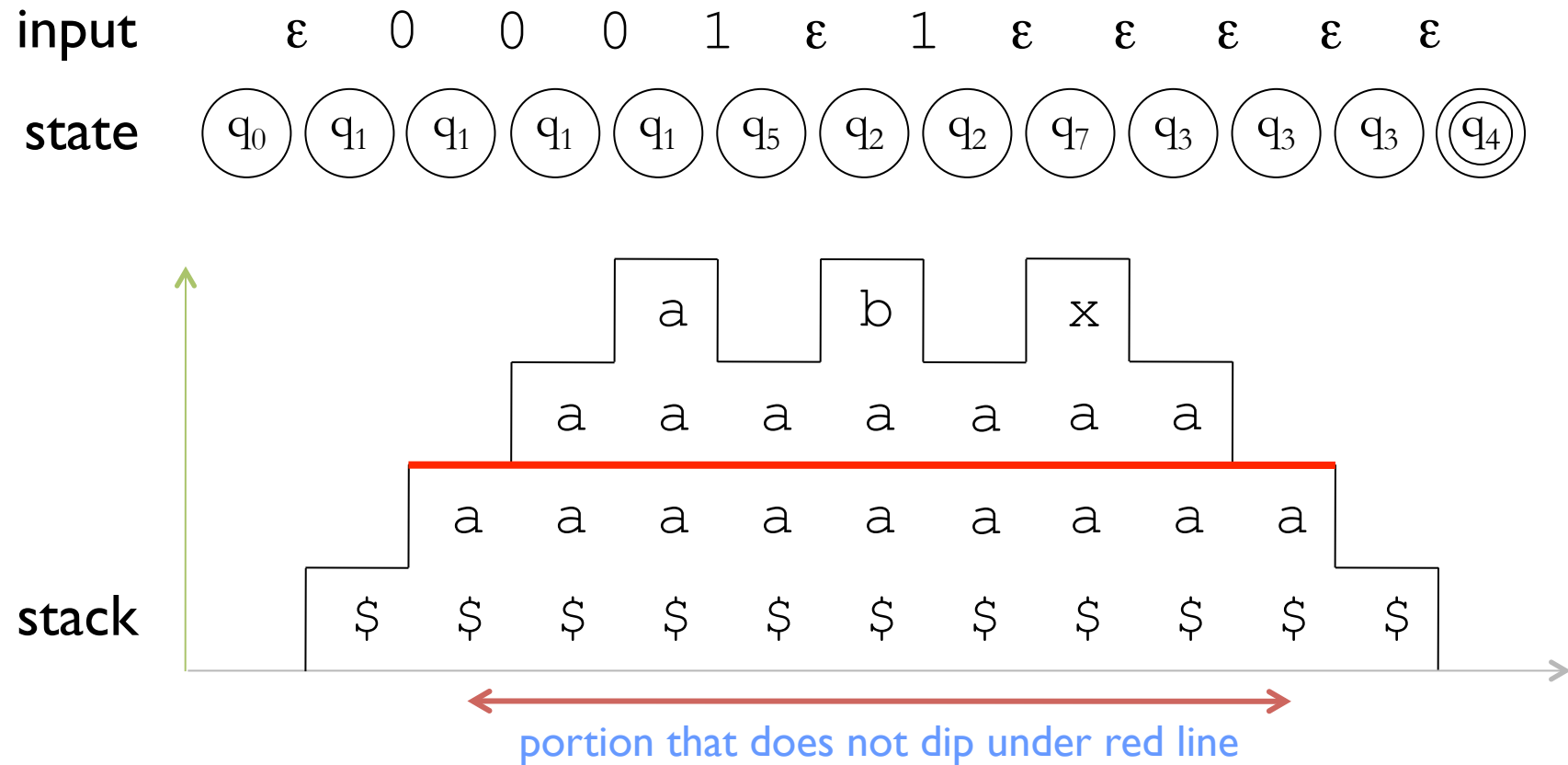
- A simplified PDA:
 - Has a **single accept state** ✓
 - **Empties its stack** before accepting ✓
 - Each transition is either a push, or a pop, but not both ✓

Simplified PDA to CFG



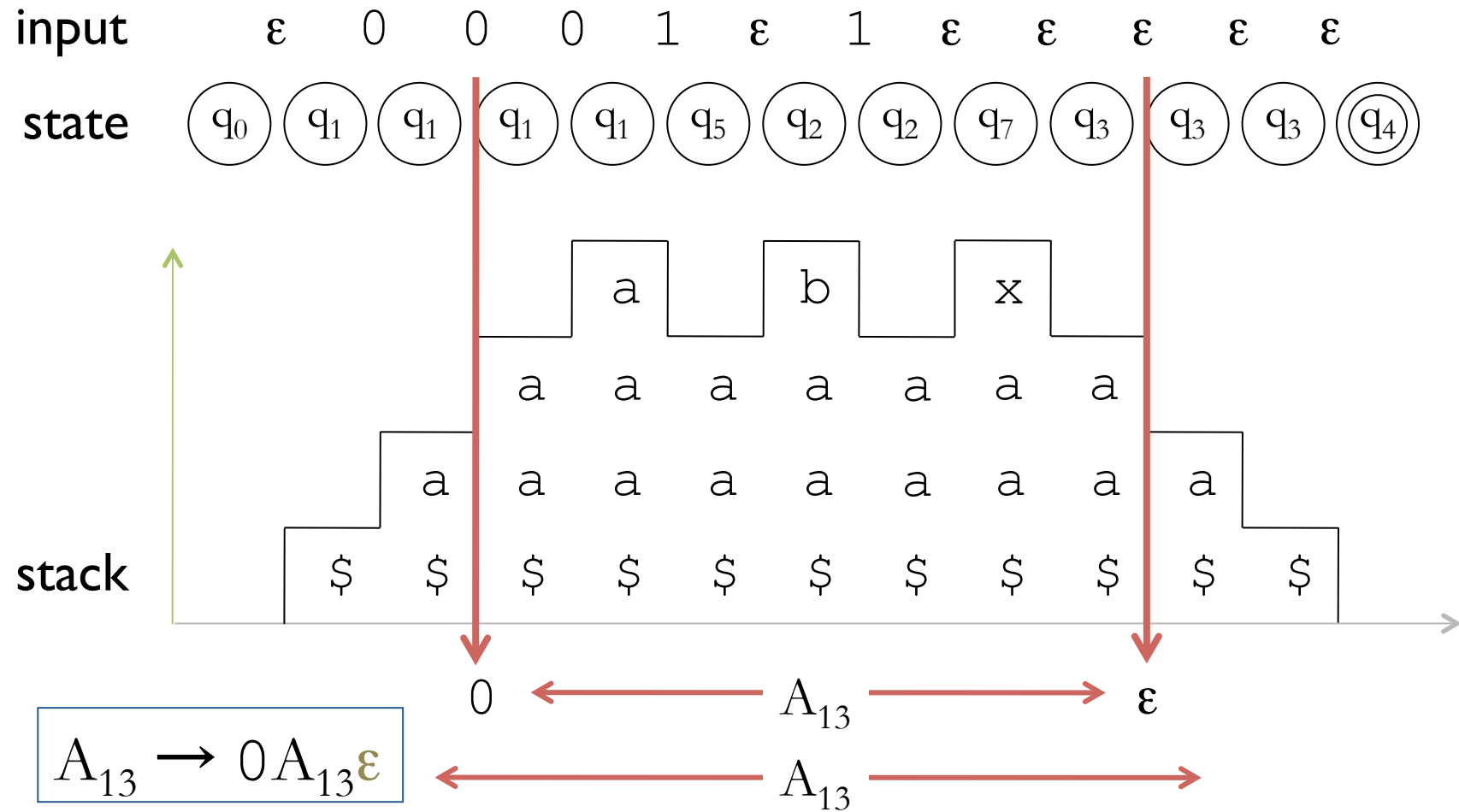
A sample run of the PDA on input 00011

Simplified PDA to CFG

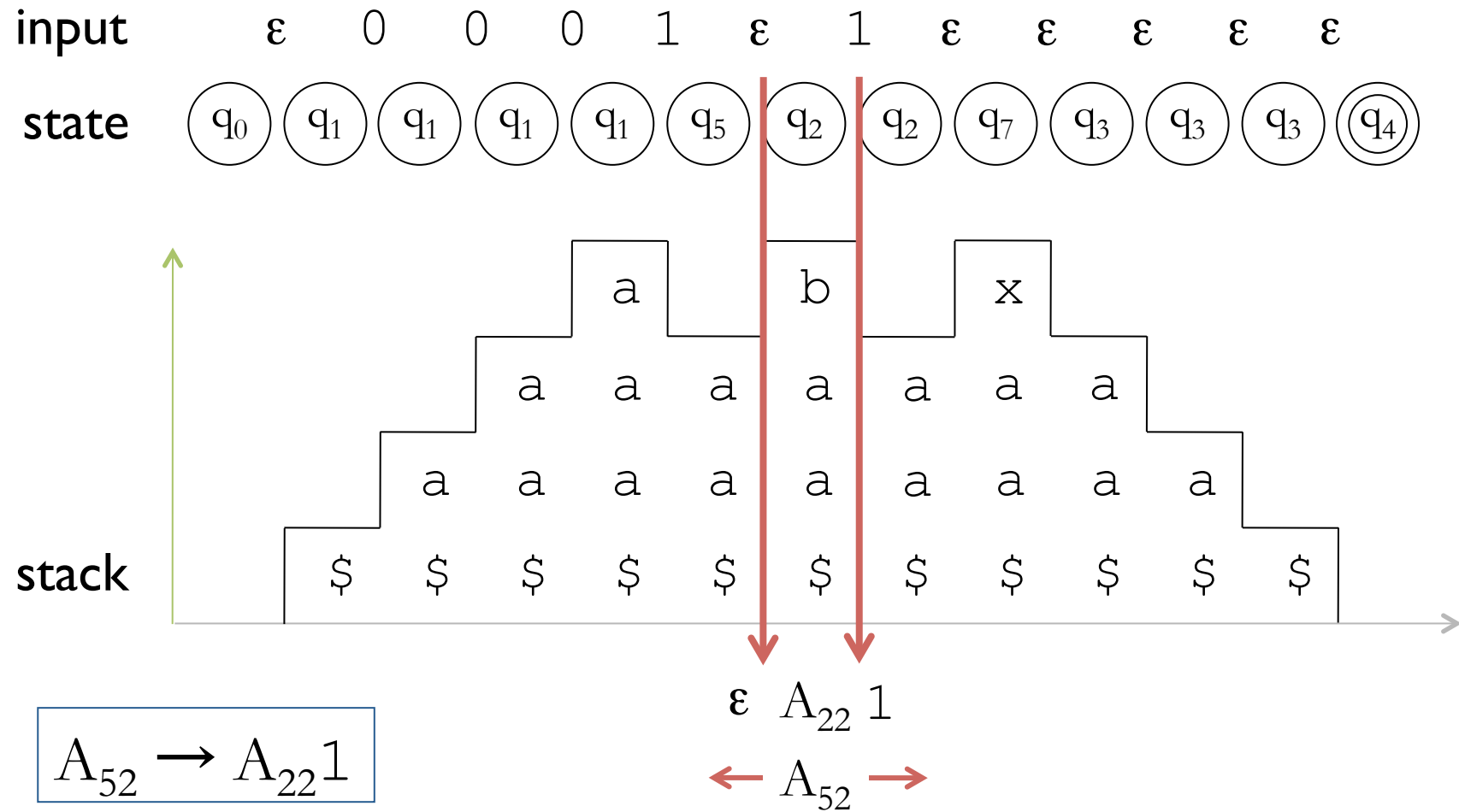


$$A_{13} = \{x: x \text{ leads from } q_1 \text{ to } q_3 \text{ and does not dip under red line}\}$$

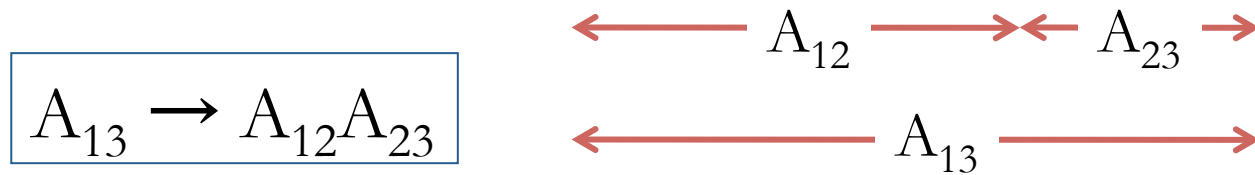
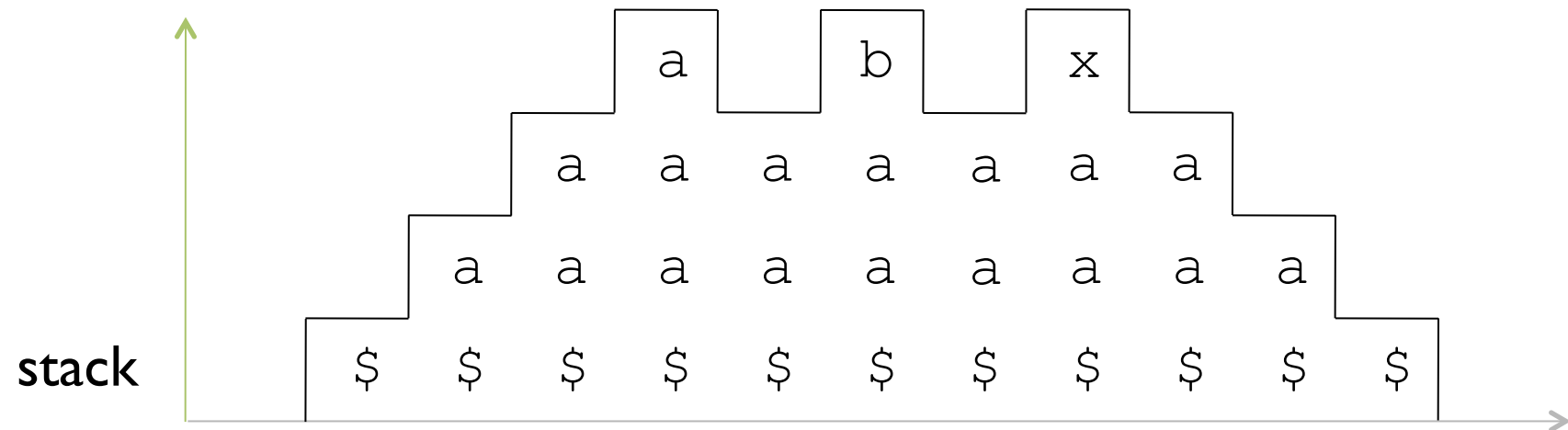
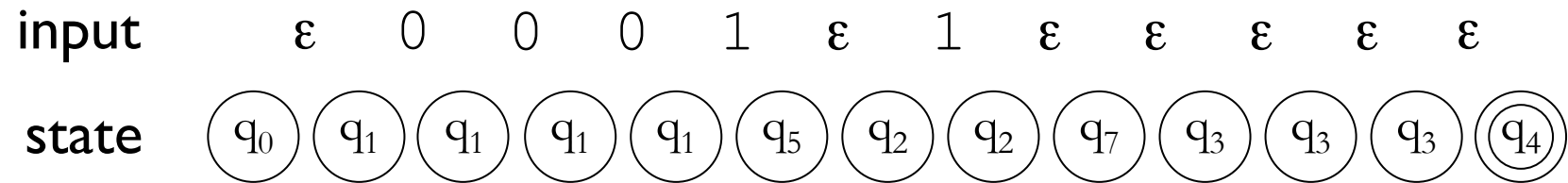
Simplified PDA to CFG



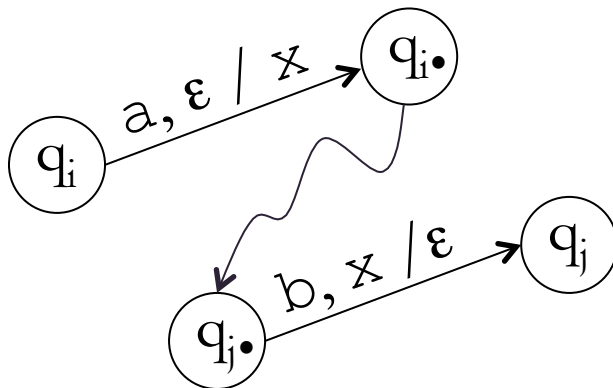
Simplified PDA to CFG



Simplified PDA to CFG



Simplified PDA to CFG



variables: A_{ij}
 start variable: A_{0f}

$$A_{ij} \rightarrow aA_{i\bullet j\bullet}b$$

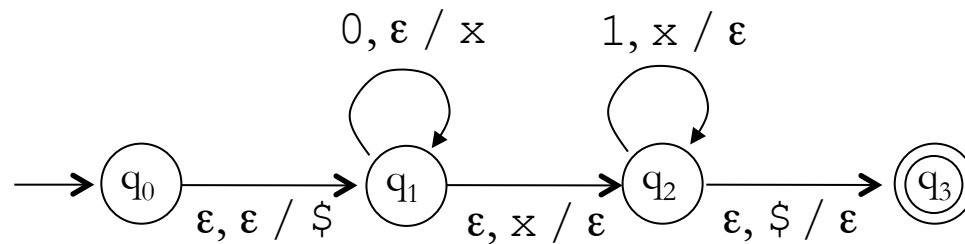


$$A_{ik} \rightarrow A_{ij}A_{jk}$$



$$A_{ii} \rightarrow \varepsilon$$

Example: Simplified PDA to CFG



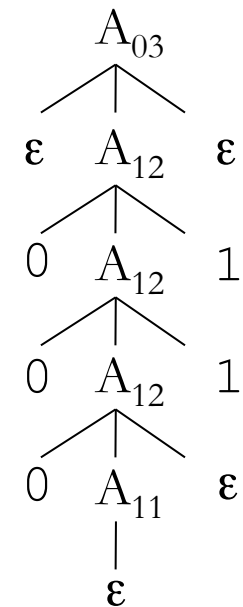
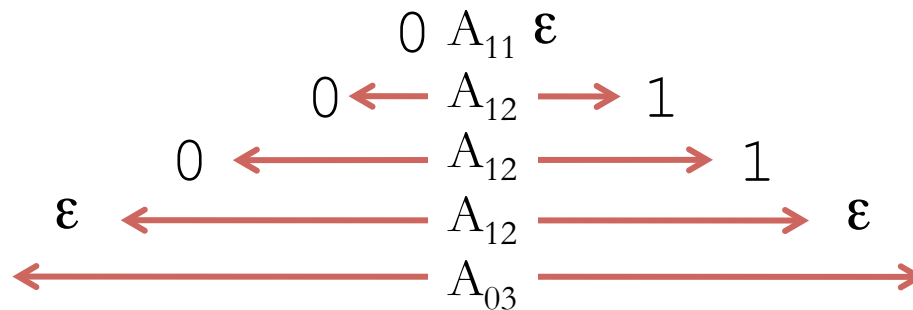
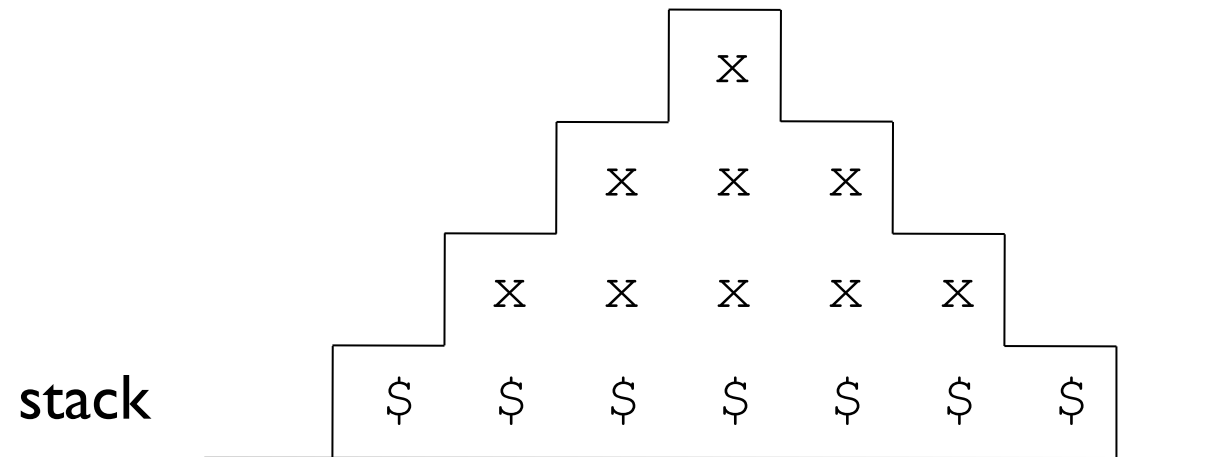
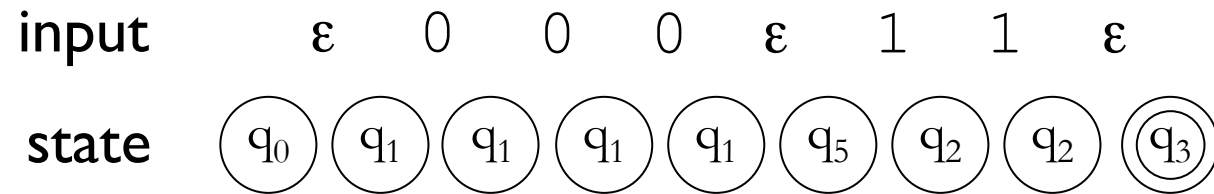
variables: $A_{01}, A_{02}, A_{03},$
 $A_{11}, A_{12}, A_{13}, A_{22}, A_{23}$

start variable: A_{03}

productions:

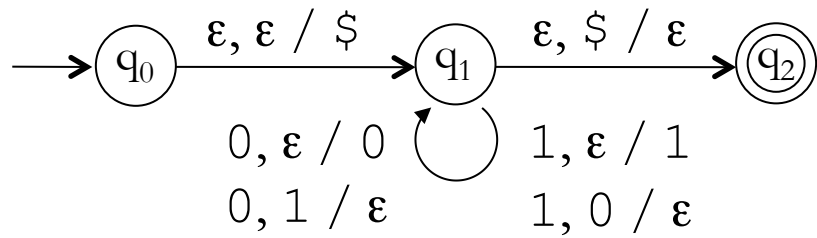
$A_{01} \rightarrow A_{01}A_{11}$	$A_{12} \rightarrow A_{11}A_{12}$	$A_{00} \rightarrow \epsilon$	$A_{12} \rightarrow 0A_{12}1$
$A_{02} \rightarrow A_{01}A_{12}$	$A_{12} \rightarrow A_{12}A_{22}$	$A_{11} \rightarrow \epsilon$	$A_{12} \rightarrow 0A_{11}$
$A_{02} \rightarrow A_{02}A_{22}$	$A_{13} \rightarrow A_{11}A_{13}$	$A_{22} \rightarrow \epsilon$	$A_{03} \rightarrow A_{12}$
$A_{03} \rightarrow A_{01}A_{13}$	$A_{13} \rightarrow A_{12}A_{23}$	$A_{33} \rightarrow \epsilon$	
$A_{03} \rightarrow A_{02}A_{23}$	$A_{22} \rightarrow A_{22}A_{22}$		
$A_{11} \rightarrow A_{11}A_{11}$	$A_{23} \rightarrow A_{22}A_{23}$		

Example: Simplified PDA to CFG



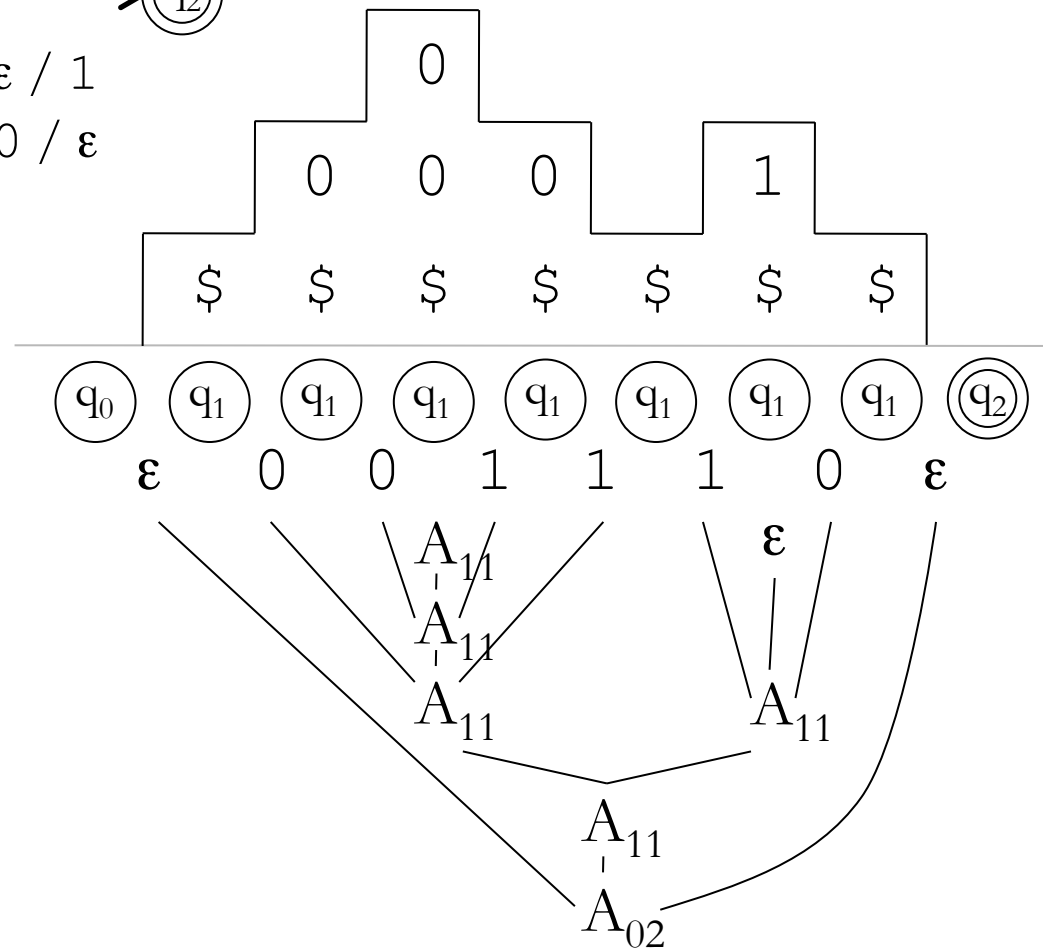
parse tree

Another example

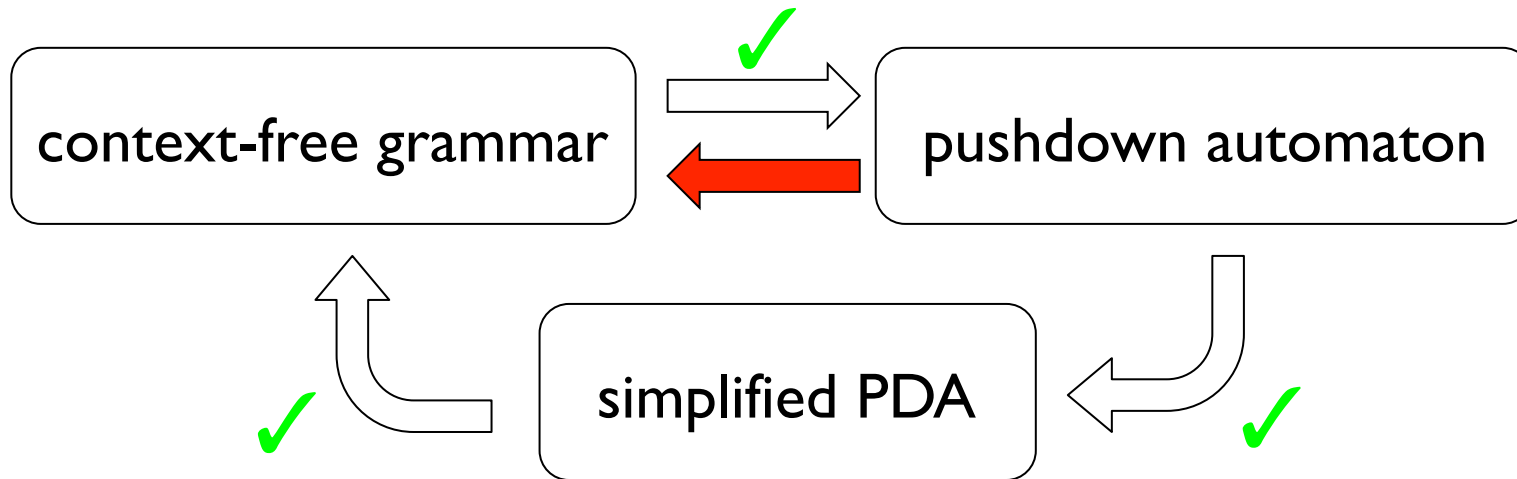


productions:

$A_{01} \rightarrow A_{01}A_{11}$
 $A_{02} \rightarrow A_{01}A_{12}$
 $A_{12} \rightarrow A_{11}A_{12}$
 $A_{11} \rightarrow A_{11}A_{11}$
 $A_{11} \rightarrow 0A_{11}1$
 $A_{11} \rightarrow 1A_{11}0$
 $A_{02} \rightarrow A_{11}$
 $A_{00}, A_{11}, A_{22} \rightarrow \epsilon$



From PDAs to CFGs



- A simplified PDA:
 - Has a **single accept state**
 - **Empties its stack** before accepting
 - Each transition is either a push, or a pop, but not both

Limitations of pushdown automata

Non context-free languages

$$L_1 = \{a^n b^n : n \geq 0\} \quad \checkmark$$

$$L_2 = \{s : s \text{ has the same number of } a\text{'s and } b\text{'s}\} \quad \checkmark$$

$$L_3 = \{a^n b^n c^n : n \geq 0\} \quad ?$$

$$L_4 = \{ss^R : s \in \{a, b\}^*\}$$

$$L_5 = \{ss : s \in \{a, b\}^*\}$$

These are not regular

Are they context-free?

An attempt

$$L_3 = \{a^n b^n c^n : n \geq 0\}$$

- Let's try to design a CFG or PDA

$S \rightarrow aBc \mid \epsilon$

$B \rightarrow ??$

read a / push x

read b / pop x

???

What would happen if...

- Suppose we could construct some CFG for L_3 , e.g.

$S \rightarrow BC$

$B \rightarrow CS \mid b$

$C \rightarrow SB \mid a$

...

$S \Rightarrow BC$

$\Rightarrow CSC$

$\Rightarrow aSC$

$\Rightarrow aBCC$

$\Rightarrow abCC$

$\Rightarrow abaC$

$\Rightarrow abaSB$

$\Rightarrow abaBCB$

$\Rightarrow ababCB$

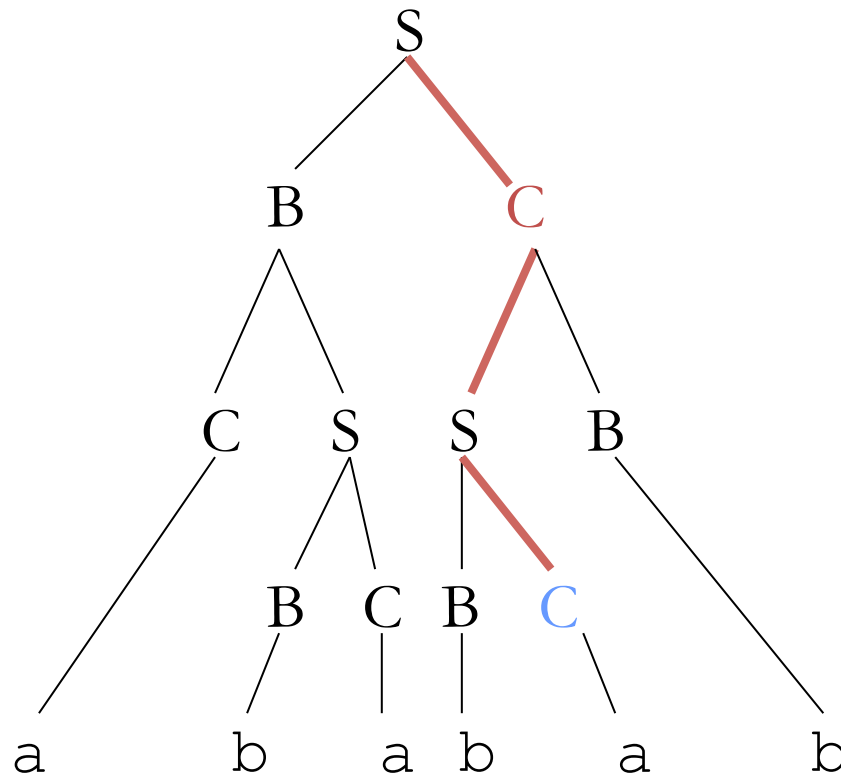
$\Rightarrow ababaB$

$\Rightarrow ababab$

- Let's do some long derivations

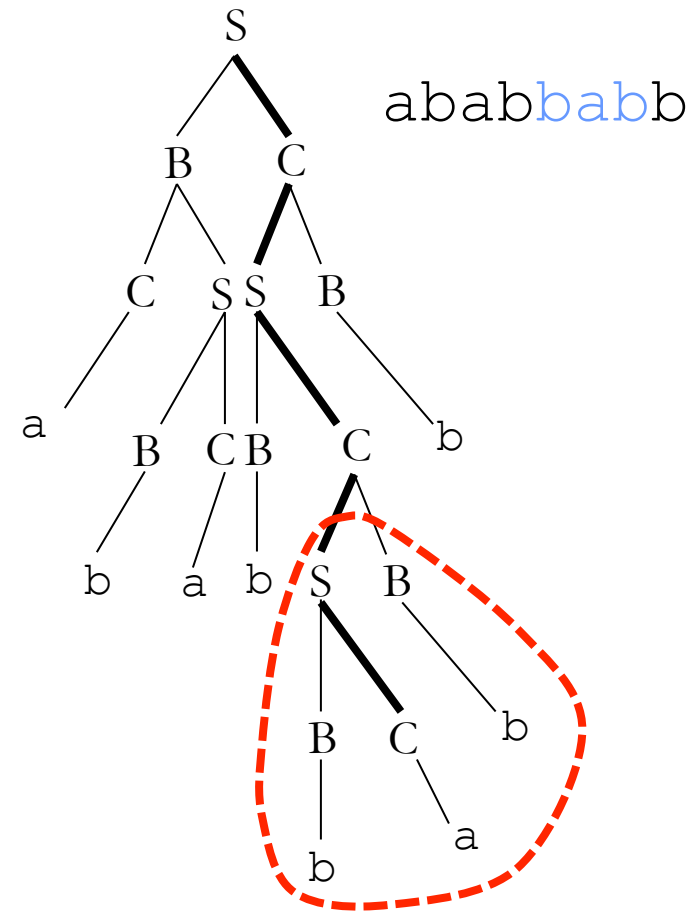
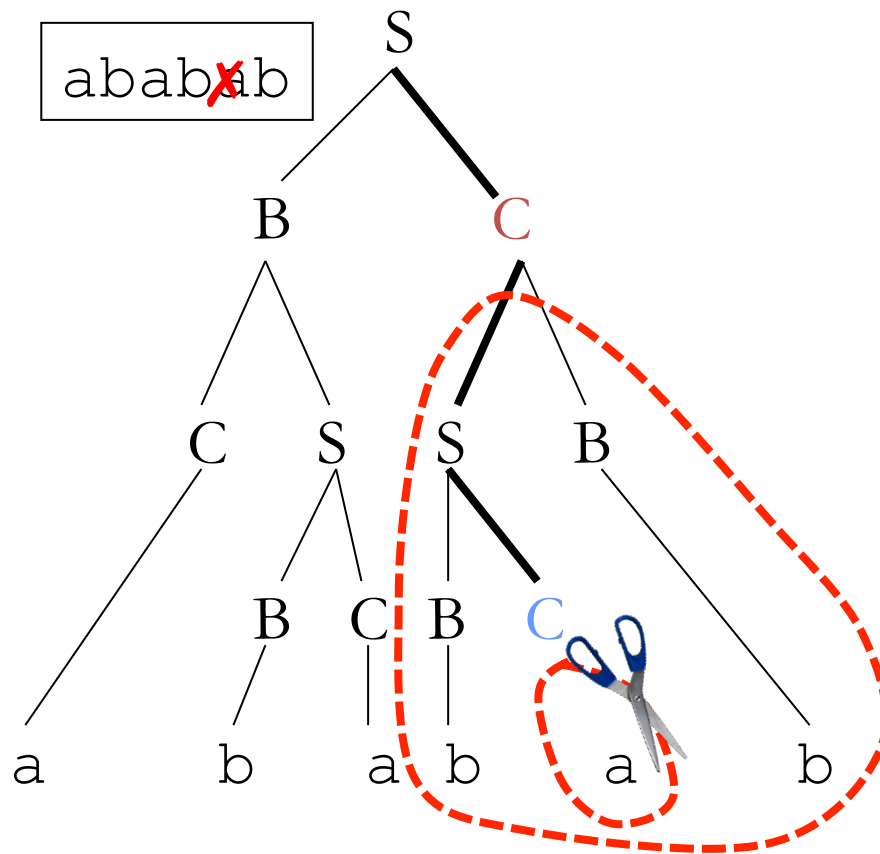
Repetition in long derivations

- If a derivation is long enough, some variable must appear **twice on the same path** in a parse tree

$$\begin{aligned} S &\Rightarrow BC \\ &\Rightarrow CSC \\ &\Rightarrow aSC \\ &\Rightarrow aBCC \\ &\Rightarrow abCC \\ &\Rightarrow abaC \\ &\Rightarrow abaSB \\ &\Rightarrow abaBCB \\ &\Rightarrow ababCB \\ &\Rightarrow ababaB \\ &\Rightarrow ababab \end{aligned}$$


Pumping example

- Then we can “cut and paste” part of the parse tree



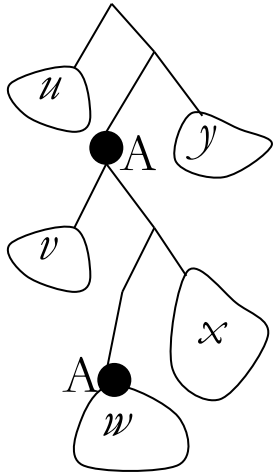
Pumping example

- We can repeat this many times

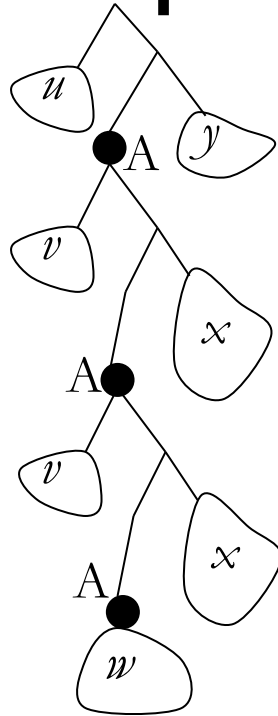
$abab\cancel{a}b \Rightarrow abab\cancel{b}bb \Rightarrow ababbbabbb$
 $\Rightarrow aba\cancel{b}^n\cancel{a}^nb\cancel{b}^nbb$

- Every sufficiently large derivation will have a middle part that can be repeated indefinitely

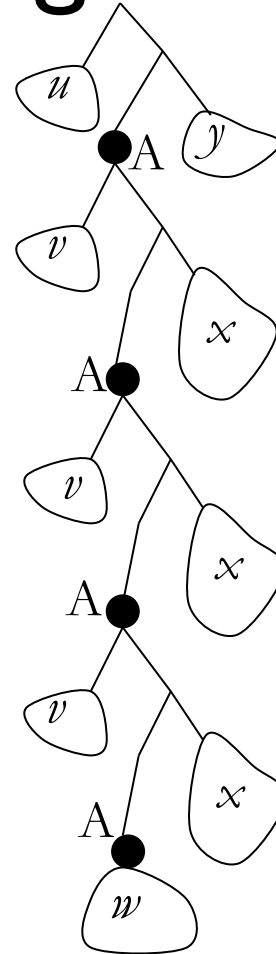
Pumping in general



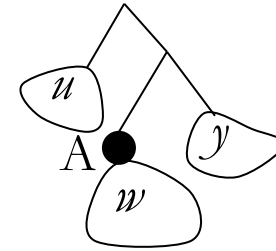
$uvwx y$



uv^2wx^2y



uv^3wx^3y



uwy

Example

$$L_3 = \{a^n b^n c^n : n \geq 0\}$$

- If L_3 has a context-free grammar G , then

If $uvwx$ is in G , so are uv^2wx^2y , uv^3wx^3y , uvw , ...

- What happens for $a^n b^n c^n$?

$a \ a \ a \ \dots \ a \ a \ b \ b \ b \ \dots \ b \ b \ c \ c \ c \ \dots \ c \ c$
 $\leftarrow u \ \rightarrow \leftarrow v \ \rightarrow \leftarrow w \ \rightarrow \leftarrow x \ \rightarrow \leftarrow y \ \rightarrow$

- No matter how it is split, $uv^2wx^2y \notin L_3$!

Pumping lemma for context-free languages

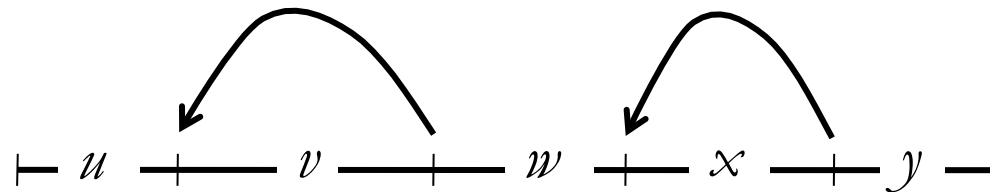
- **Pumping lemma:** For every context-free language L

There exists a number n such that for every string z in L longer than n , we can write $z = uvwx y$ where

① $|vwx| \leq n$

② $|vx| \geq 1$

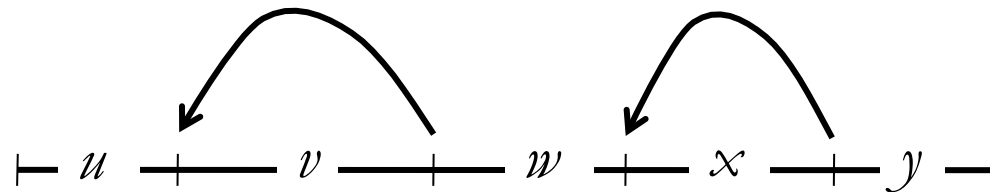
③ For every $i \geq 0$, the string uv^iwx^iy is in L .



Pumping lemma for context-free languages

- So to prove L is **not context-free**, it is enough that

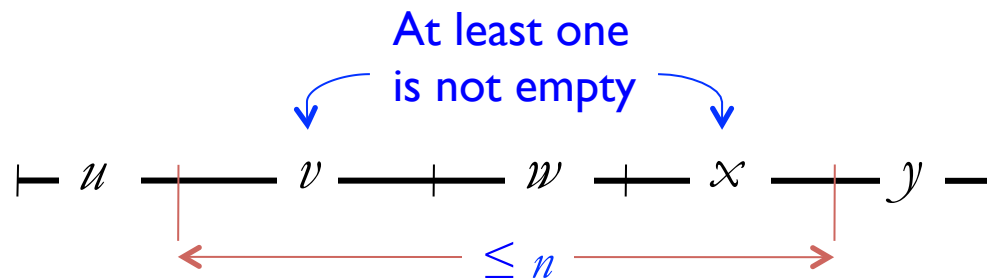
For all n there exists z in L longer than n , such that for all ways of writing $z = uvwx$ where
① $|vwx| \leq n$ and ② $|vx| \geq 1$, the string uv^iwx^iy is not in L for some $i \geq 0$.



Proving a language is not context-free

- Just like for regular languages, Eve needs a **strategy** that wins her this game **no matter** what Adam chooses

Adam	Eve
1 chooses n	chooses $z \in L$ ($ z > n$)
2 writes $z = uvwx$ ($ vwx \leq n, vx \geq 1$)	chooses i Eve wins if $uv^iwx^iy \notin L$



Example

Adam	Eve
1 chooses n	chooses $z \in L$ ($ z > n$)
2 writes $z = uvwx y$ ($ vwx \leq n, vx \geq 1$)	chooses i Eve wins if $uv^iwx^iy \notin L$

$$L_3 = \{a^n b^n c^n : n \geq 0\}$$

1 chooses n	$z = a^n b^n c^n$
2 writes $z = uvwx y$	$i = ?$

$a \ a \ a \ \dots \ a \ a \ b \ b \ b \ \dots \ b \ b \ c \ c \ c \ \dots \ c \ c$
 $\leftarrow u \ \rightarrow \leftarrow v \ \rightarrow \leftarrow w \ \rightarrow \leftarrow x \ \rightarrow \leftarrow y \ \rightarrow$

Example

- **Case 1:** v or x contains two kinds of symbols

$a\ a\ a\ \dots\ a\ a\ b\ b\ b\ \dots\ b\ b\ c\ c\ c\ \dots\ c\ c$
 $\longleftarrow v \longrightarrow \qquad \qquad \qquad \longleftarrow x \longrightarrow$

Then uv^2wx^2y not in L_3 because pattern is wrong

- **Case 2:** v and x both contain one kind of symbol

$a\ a\ a\ \dots\ a\ a\ b\ b\ b\ \dots\ b\ b\ c\ c\ c\ \dots\ c\ c$
 $\longleftarrow v \longrightarrow \qquad \qquad \qquad \longleftarrow x \longrightarrow$

Then uv^2wx^2y does not have same number of a s,
 b s, c s

More examples

$$L_1 = \{a^n b^n : n \geq 0\} \quad \checkmark$$

$$L_2 = \{s : s \text{ has same number of } a\text{'s and } b\text{'s}\} \quad \checkmark$$

$$L_3 = \{a^n b^n c^n : n \geq 0\} \quad \times$$

$$L_4 = \{ss^R : s \in \{a, b\}^*\} \quad \checkmark$$

$$L_5 = \{ss : s \in \{a, b\}^*\}$$

Which is context-free?

Example

$$L_5 = \{s : s \in \{a, b\}^*\}$$

1 chooses n

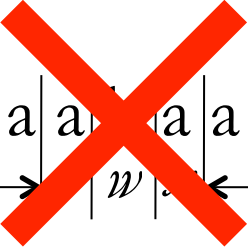
$$z = a^n b a^n b$$

2 writes $z = uvwxy$

$$i = ?$$

a a a a a a a a a b a a a a a a a a b
 $\leftarrow u \rightarrow \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \leftarrow y \rightarrow$

What if:


 $\left| \begin{array}{c} a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ a \ b \\ \leftarrow u \rightarrow \leftarrow w \rightarrow \leftarrow y \rightarrow \end{array} \right|$

Example

$$L_5 = \{s : s \in \{a, b\}^*\}$$

1 chooses n

$$z = a^n b^n a^n b^n$$

2 writes $z = uvwx$

$$i = ?$$

a a a a a b b b b b b a a a a a b b b b b b
 $\longleftarrow u \quad \longrightarrow \times v \quad \rightarrow \leftarrow w \quad \rightarrow \leftarrow x \quad \rightarrow \leftarrow y \quad \longrightarrow$

Recall that $|vwx| \leq n$

Example

Three cases

Case 1: $\left| \begin{array}{cccccccccccc} a & a & a & a & a & a & b & b & b & b & b & b \end{array} \right| a & a & a & a & a & a & b & b & b & b & b & b$
 $\leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow$
 $vw x$ is in the **first half** of $a^n b^n a^n b^n$

Case 2: $a & a & a & a & a & a \left| b & b & b & b & b & b & a & a & a & a & a & a \right| b & b & b & b & b & b$
 $\leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow$
 $vw x$ is in the **middle part** of $a^n b^n a^n b^n$

Case 3: $a & a & a & a & a & a & b & b & b & b & b & b \left| a & a & a & a & a & a & b & b & b & b & b & b \right|$
 $\leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow$
 $vw x$ is in the **second half** of $a^n b^n a^n b^n$

Example

Apply pumping with $i = 0$

Case 1: $\left| \begin{array}{c} a \ a \ a \ a \ a \ a \ b \ b \ b \ b \ b \ b \\ \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \end{array} \right| a \ a \ a \ a \ a \ a \ b \ b \ b \ b \ b \ b$

$uv^jw^kx^n$ looks like $a^ja^kb^na^nb^n$, where $j < n$ or $k < n$

Case 2: $\begin{array}{c} a \ a \ a \ a \ a \ a \\ \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \end{array} \left| \begin{array}{c} b \ b \ b \ b \ b \ b \ a \ a \ a \ a \ a \ a \\ \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \end{array} \right| b \ b \ b \ b \ b \ b$

$uv^jw^kx^n$ looks like $a^nb^ja^kb^n$, where $j < n$ or $k < n$

Case 3: $\begin{array}{c} a \ a \ a \ a \ a \ a \ b \ b \ b \ b \ b \ b \\ \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \end{array} \left| \begin{array}{c} a \ a \ a \ a \ a \ a \ b \ b \ b \ b \ b \ b \\ \leftarrow v \rightarrow \leftarrow w \rightarrow \leftarrow x \rightarrow \end{array} \right|$

$uv^jw^kx^n$ looks like $a^nb^na^jb^k$, where $j < n$ or $k < n$

Example

Apply pumping with $i = 0$

$$L_5 = \{ss : s \in \{a, b\}^*\}$$

Case 1: w^0wx^0y looks like $a^jb^ka^nb^n$, where $j < n$ or $k < n$

Not of the form ss

Case 2: w^0wx^0y looks like $a^nb^ja^kb^n$, where $j < n$ or $k < n$

Not of the form ss

Case 3: w^0wx^0y looks like $a^nb^na^jb^k$, where $j < n$ or $k < n$

Not of the form ss

This covers all the cases, so L_5 is not context-free.

Which language is context-free?

$$L_1 = \{a^n b^n : n \geq 0\} \quad \checkmark$$

$$L_2 = \{s : s \text{ has same number of } a\text{'s and } b\text{'s}\} \quad \checkmark$$

$$L_3 = \{a^n b^n c^n : n \geq 0\} \quad \times$$

$$L_4 = \{ss^R : s \in \{a, b\}^*\} \quad \checkmark$$

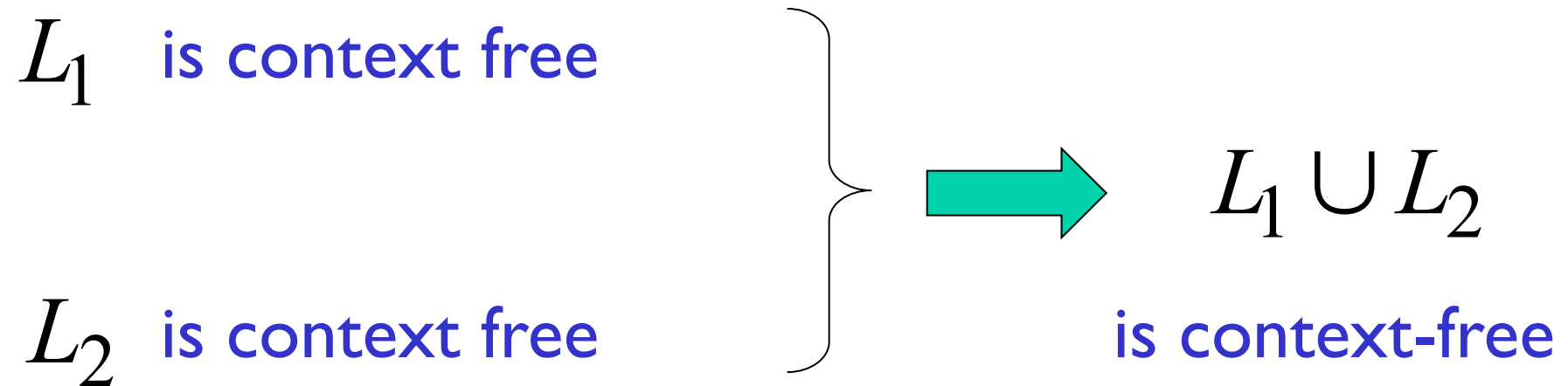
$$L_5 = \{ss : s \in \{a, b\}^*\} \quad \times$$

Properties of Context-Free languages

Union

Context-free languages
are closed under:

Union



Example

Language

Grammar

$$L_1 = \{a^n b^n\}$$

$$S_1 \rightarrow aS_1b \mid \varepsilon$$

$$L_2 = \{ww^R\}$$

$$S_2 \rightarrow aS_2a \mid bS_2b \mid \varepsilon$$

Union

$$L = \{a^n b^n\} \cup \{ww^R\}$$

$$S \rightarrow S_1 \mid S_2$$

In general:

For context-free languages
with context-free grammars
and start variables

L_1, L_2

G_1, G_2

S_1, S_2

The grammar of the **union**
has new start variable
and additional production

$L_1 \cup L_2$

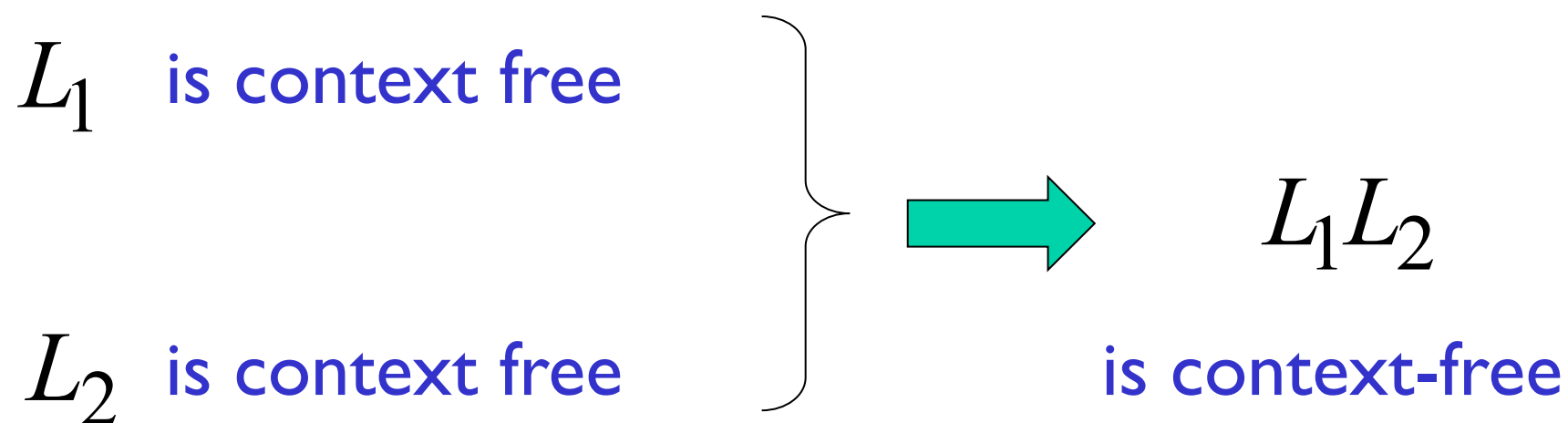
S

$S \rightarrow S_1 \mid S_2$

Concatenation

Context-free languages
are closed under:

Concatenation



Example

Language

Grammar

$$L_1 = \{a^n b^n\}$$

$$S_1 \rightarrow aS_1b \mid \varepsilon$$

$$L_2 = \{ww^R\}$$

$$S_2 \rightarrow aS_2a \mid bS_2b \mid \varepsilon$$

Concatenation

$$L = \{a^n b^n\} \{ww^R\}$$

$$S \rightarrow S_1 S_2$$

In general:

For context-free languages
with context-free grammars
and start variables

L_1, L_2

G_1, G_2

S_1, S_2

The grammar of the **concatenation**
has new start variable
and additional production

L_1L_2

S

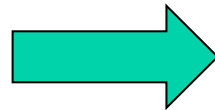
$S \rightarrow S_1S_2$

Star Operation

Context-free languages
are closed under:

Star-operation

L is context free



L^* is context-free

Example

Language

$$L = \{a^n b^n\}$$

Grammar

$$S \rightarrow aSb \mid \varepsilon$$

Star Operation

$$L = \{a^n b^n\}^*$$

$$S_1 \rightarrow SS_1 \mid \varepsilon$$

In general:

For context-free language
with context-free grammar
and start variable

 L G S

The grammar of the **star operation**
has new start variable
and additional production

 L^* S_1 $S_1 \rightarrow SS_1 \mid \varepsilon$

Negative Properties of Context-Free Languages

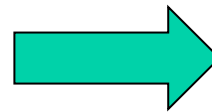
Intersection

Context-free languages
are not closed under:

intersection

L_1 is context free

L_2 is context free



$L_1 \cap L_2$

not necessarily
context-free

Example

$$L_1 = \{a^n b^n c^*\}$$

Context-free:

$$S \rightarrow AC$$

$$A \rightarrow aAb \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

$$L_2 = \{a^* b^m c^m\}$$

Context-free:

$$S \rightarrow AB$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$


Intersection

$$L_1 \cap L_2 = \{a^n b^n c^n\} \quad \textbf{NOT} \text{ context-free}$$

Complement

Context-free languages
are not closed under:

complement

L is context free  \bar{L} not necessarily
context-free

Example

$$L_1 = \{a^n b^n c^m\}$$

$$L_2 = \{a^n b^m c^m\}$$

Context-free:

$$S \rightarrow AC$$

$$A \rightarrow aAb \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

Context-free:

$$S \rightarrow AB$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$

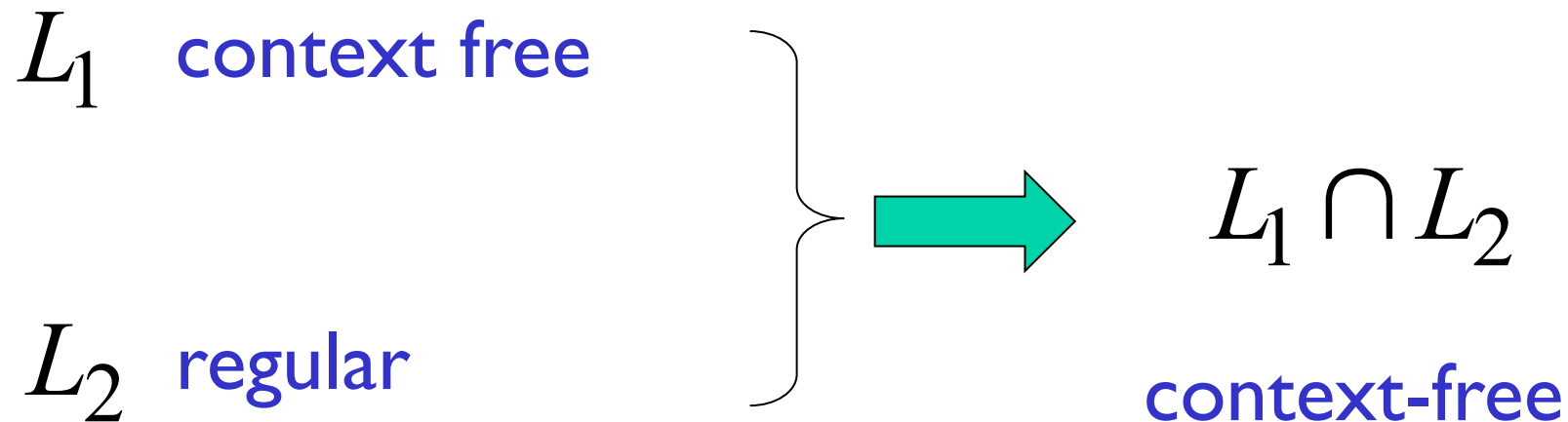
Complement

$$\overline{\overline{L_1} \cup \overline{L_2}} = L_1 \cap L_2 = \{a^n b^n c^n\}$$

NOT context-free

Intersection
of
Context-free languages
and
Regular Languages

The intersection of
a context-free language and
a regular language
is a context-free language



Machine M_1

PDA for L_1
context-free

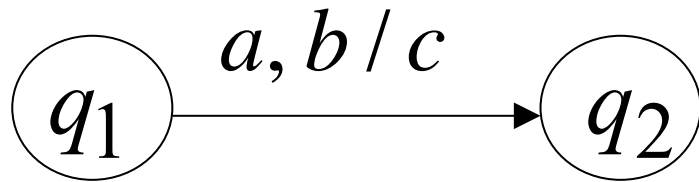
Machine M_2

DFA for L_2
regular

Construct a new PDA M
that accepts $L_1 \cap L_2$

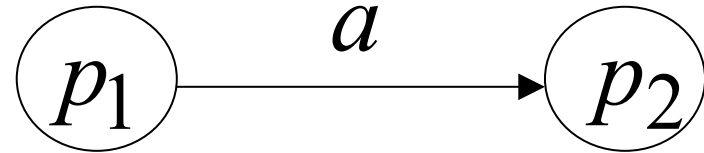
M will simulate in parallel M_1 and M_2

PDA M_1

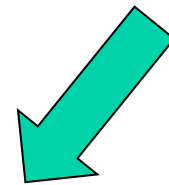


transition

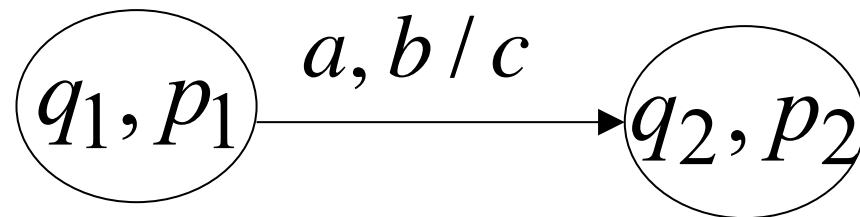
DFA M_2



transition

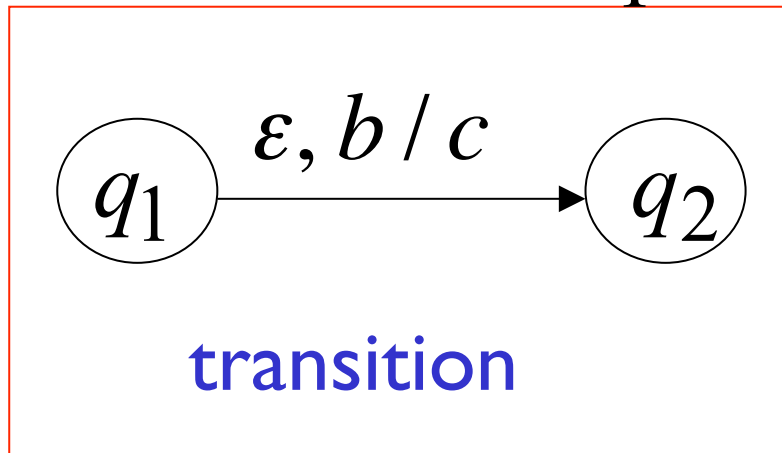


PDA M

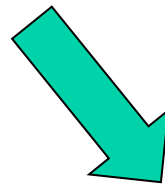


transition

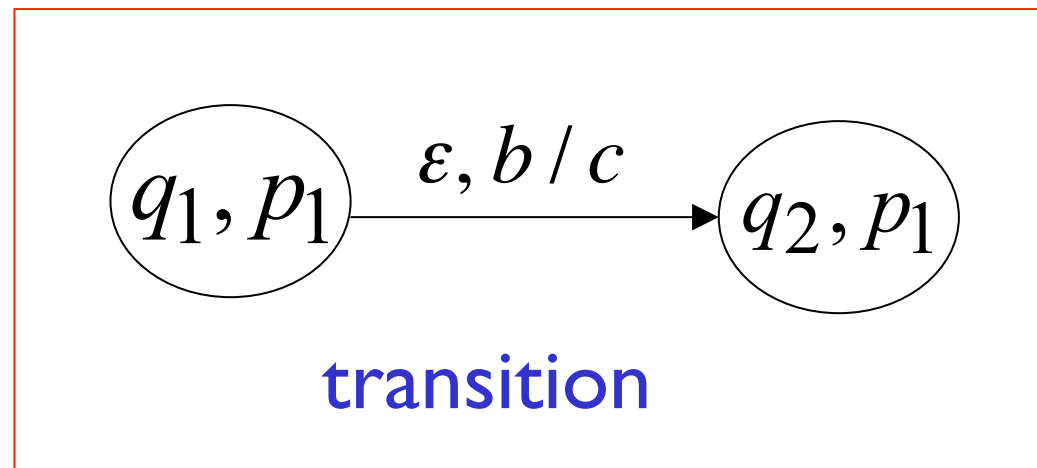
M_1



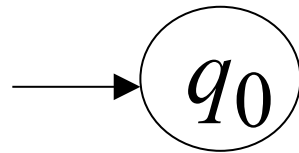
DFA M_2



M

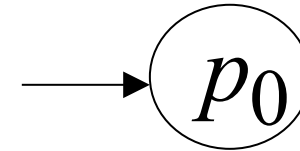


M_1

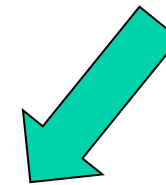
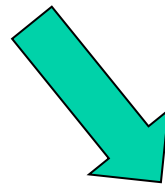


initial state

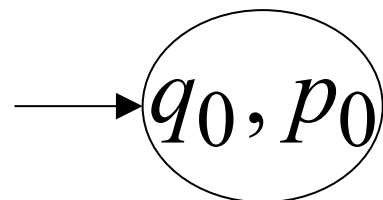
DFA M_2



initial state



M

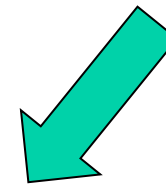
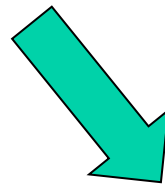
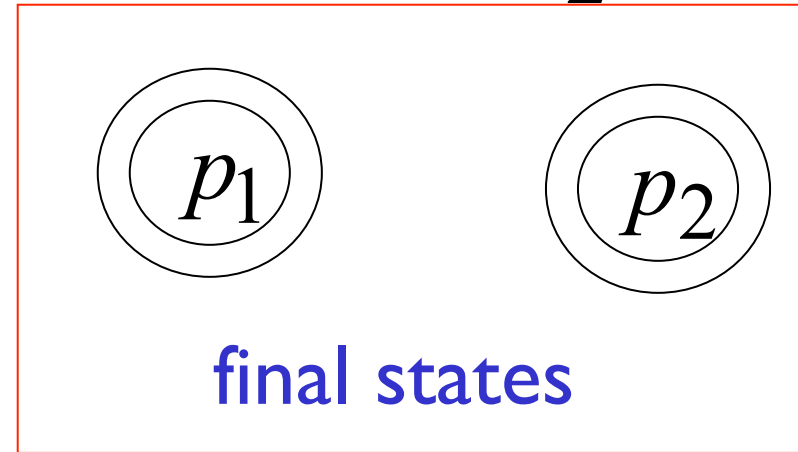


Initial state

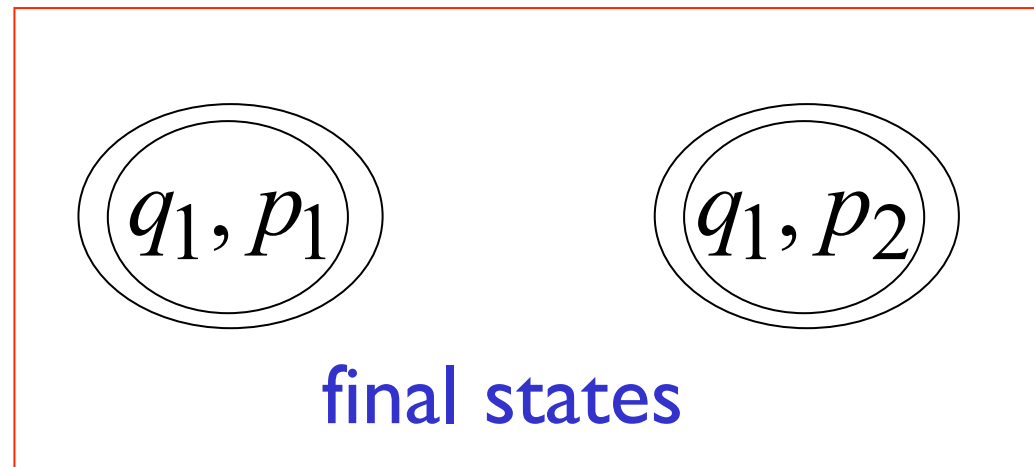
M_1



DFA M_2



M

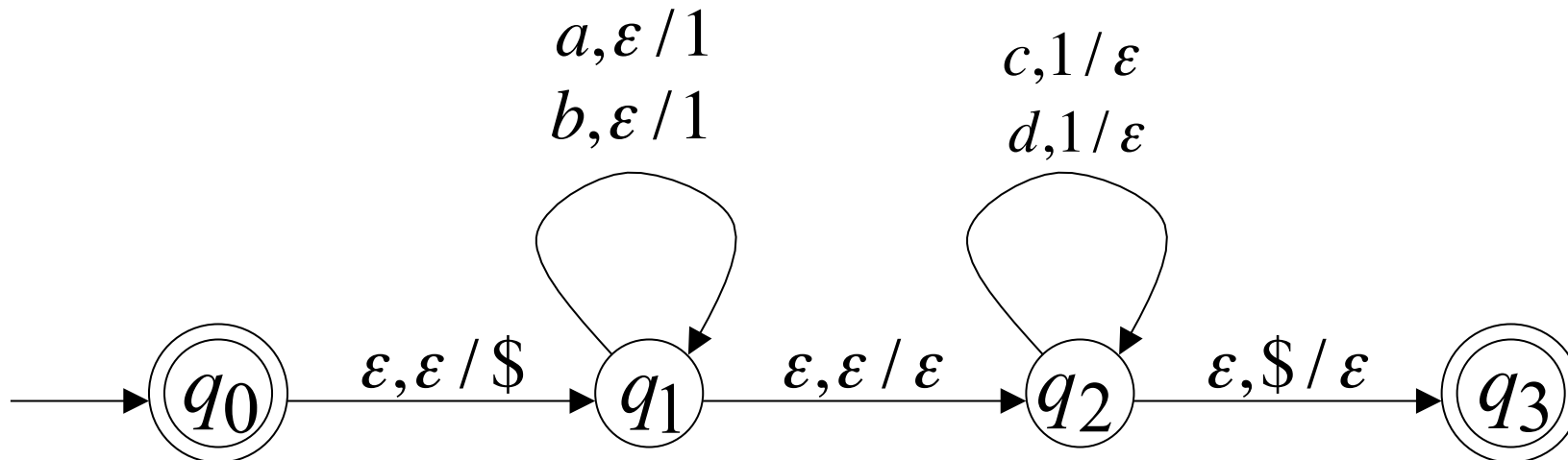


Example:

context-free

$$L_1 = \{w_1 w_2 : |w_1| = |w_2|, w_1 \in \{a, b\}^*, w_2 \in \{c, d\}^*\}$$

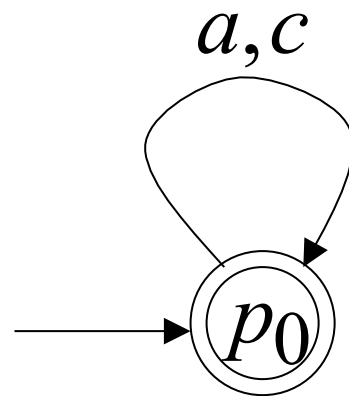
NPDA M_1



regular

$$L_2 = \{a, c\}^*$$

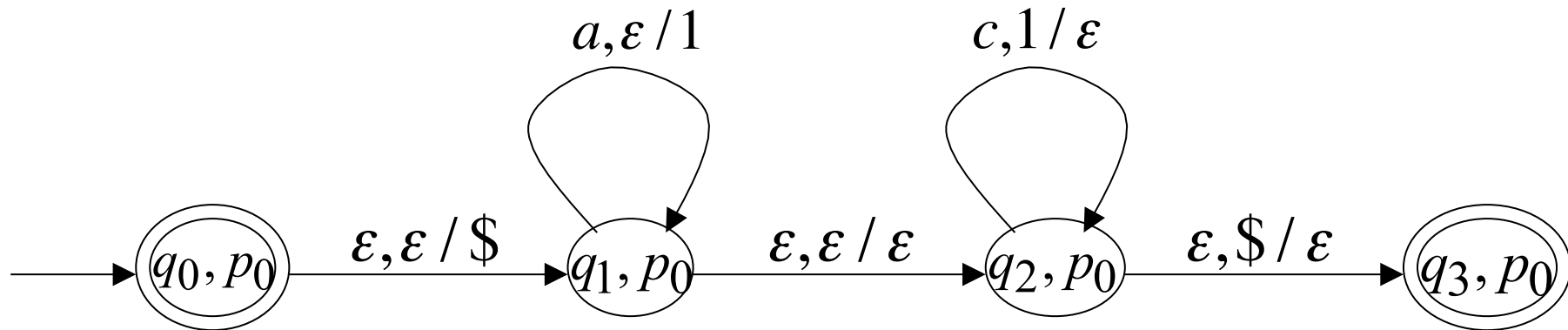
DFA M_2



context-free

Automaton for: $L_1 \cap L_2 = \{a^n c^n : n \geq 0\}$

PDA M



In General:

M simulates in parallel M_1 and M_2

M accepts string w if and only if

M_1 accepts string w and

M_2 accepts string w , i.e.

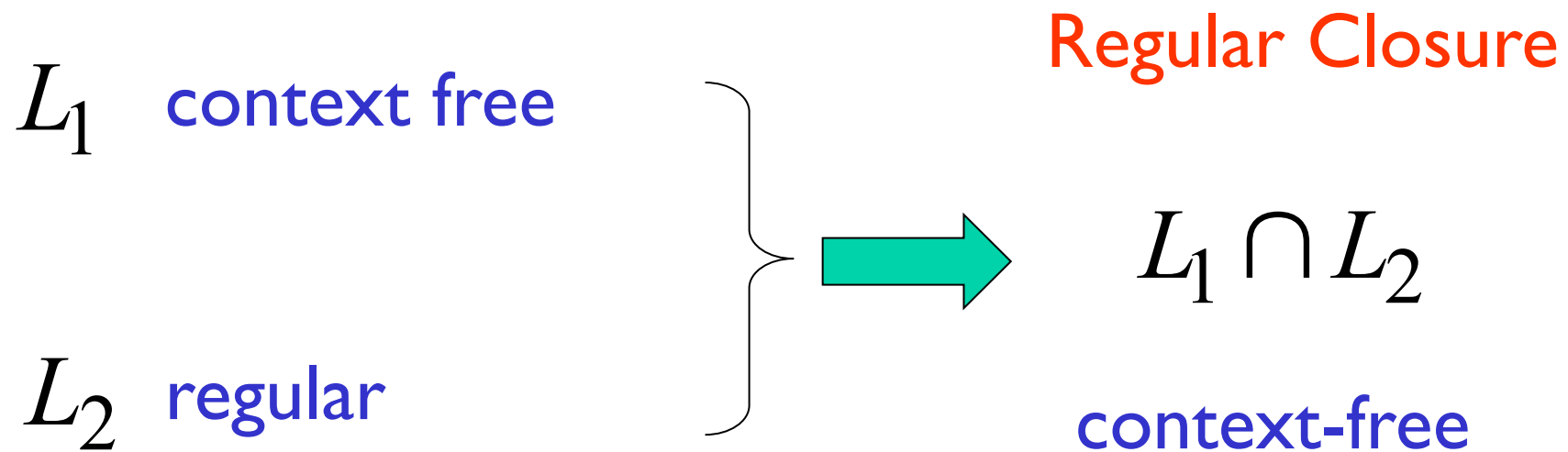
$$L(M) = L_1 \cap L_2$$

As M is a PDA, then

$L(M) = L_1 \cap L_2$ is context-free.

Applications of Regular Closure

The intersection of
a context-free language and
a regular language
is a context-free language



An Application of Regular Closure

Prove that: $L = \{a^n b^n : n \neq 100, n \geq 0\}$

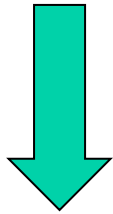
is context-free

We know:

$\{a^n b^n : n \geq 0\}$ is context-free

We also know:

$$L_1 = \{a^{100}b^{100}\} \quad \text{is regular}$$



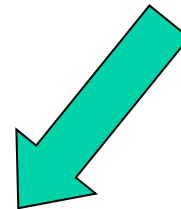
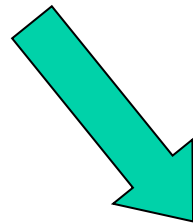
$$\overline{L_1} = \{(a+b)^*\} - \{a^{100}b^{100}\} \quad \text{is regular}$$

$$\{a^n b^n\}$$

context-free

$$\overline{L_1} = \{(a+b)^*\} - \{a^{100}b^{100}\}$$

regular



(regular closure)

$$\{a^n b^n\} \cap \overline{L_1}$$

context-free



$$\{a^n b^n\} \cap \overline{L_1} = \{a^n b^n : n \neq 100, n \geq 0\} = L$$

is context-free

Another Application of Regular Closure

Prove that: $L = \{w : n_a = n_b = n_c\}$

is **not** context-free

If $L = \{w : n_a = n_b = n_c\}$ is context-free

(regular closure)

Then $L \cap \{a^* b^* c^*\} = \{a^n b^n c^n\}$

context-free

regular

context-free

Impossible!!!

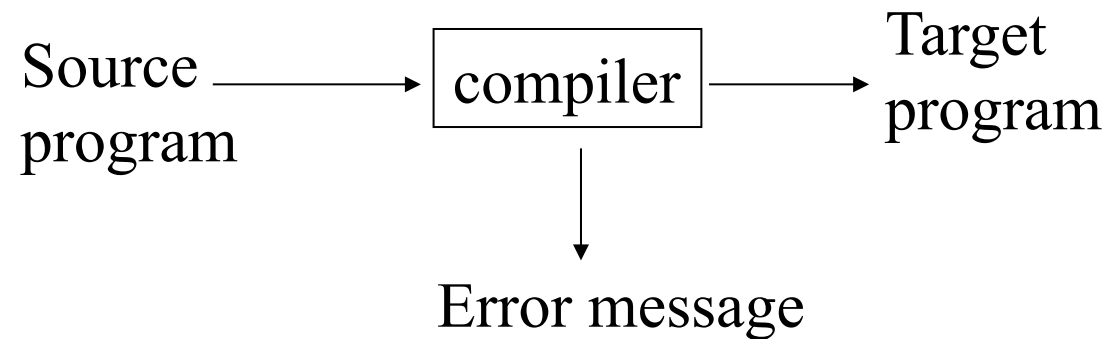
Therefore, L is **not** context free

LR(0) grammars



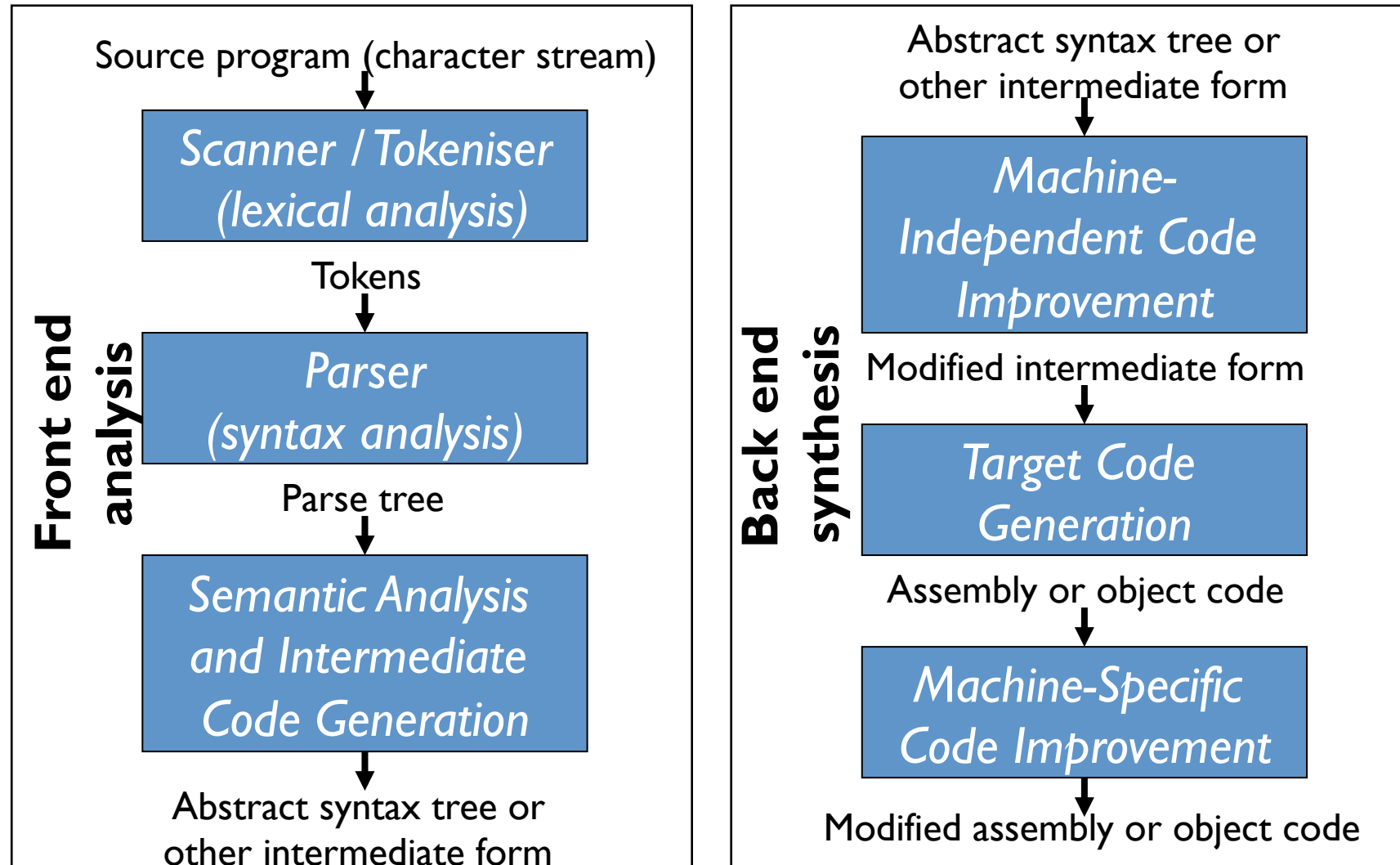
Compiler

A program that reads a program written in one language (**source language**) and translates it into an equivalent program in another language (**target language**).



Typically the source language is a high-level language and the target language is a low-level language (machine code).

Compiler Front-end and Back-end



Tokenising computer programs

```
if (n == 0) { return x; }
```

- First the `javac` compiler does a **lexical analysis**:

```
if (ID == INT_LIT) { return ID; }
```

ID = identifier (name of variable, procedure, class, ...)

INT_LIT = integer literal (value)

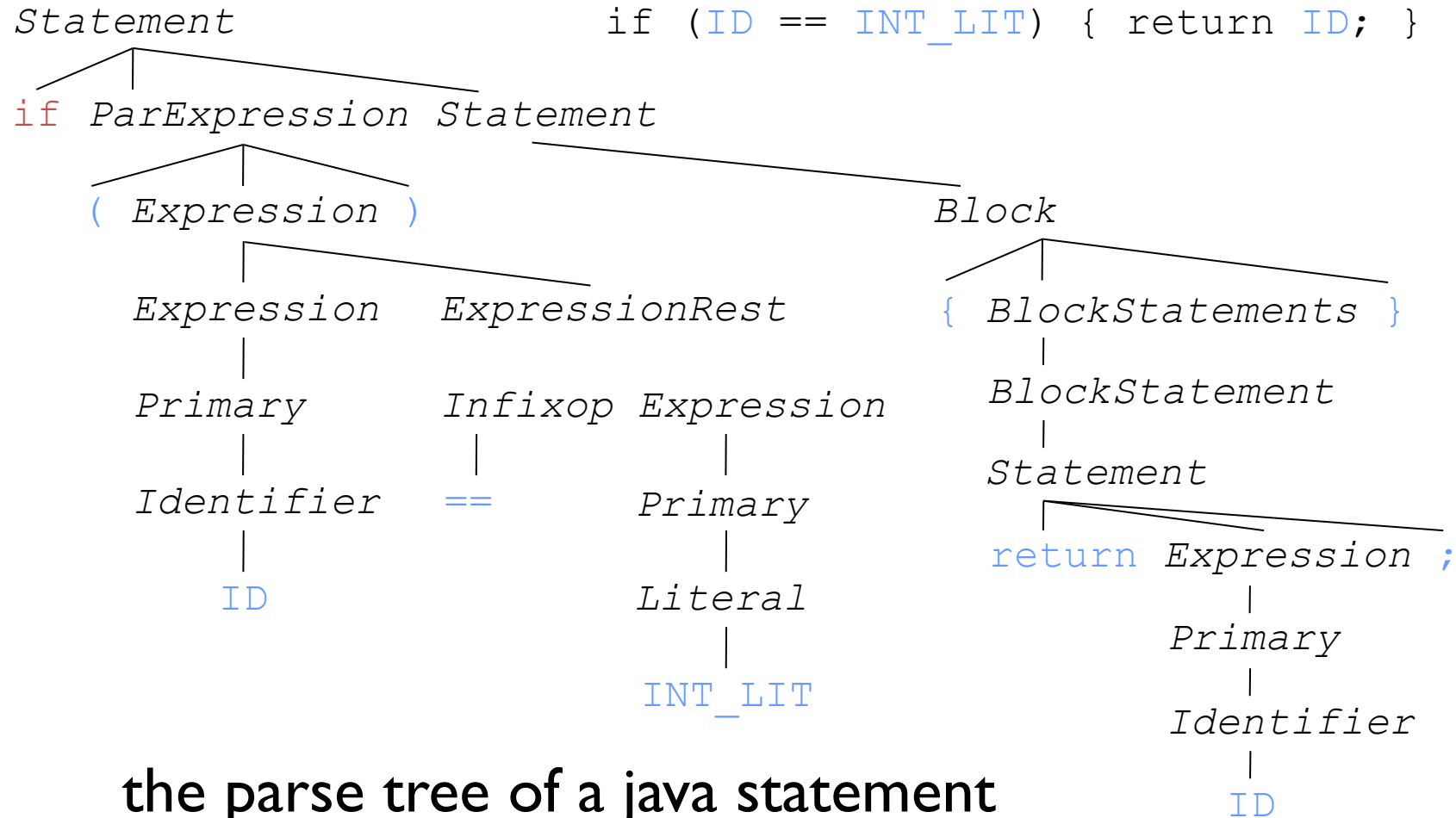
- The **alphabet** of java CFG consists of symbols like:

$$\Sigma = \{ \text{if, return, (,) \{, \}, ;, ==, ID, INT_LIT, ...} \}$$

Parsing computer programs

```
if (n == 0) { return x; }
```

```
if (ID == INT_LIT) { return ID; }
```



the parse tree of a java statement

CFG of the java programming language

Identifier:

ID

QualifiedIdentifier:

Identifier { . Identifier }

Literal:

IntegerLiteral

FloatingPointLiteral

CharacterLiteral

StringLiteral

BooleanLiteral

NullLiteral

Expression:

Expression1 [AssignmentOperator Expression1]]

AssignmentOperator:

=

+=

-=

**=*

/=

&=

|=

...

from http://java.sun.com/docs/books/jls/second_edition/html/syntax.doc.html#52996

Parsing java programs

```
class Point2d {
    /* The X and Y coordinates of the point--instance variables */
    private double x;
    private double y;
    private boolean debug;                                // A trick to help with debugging

    public Point2d (double px, double py) {               // Constructor
        x = px;
        y = py;

        debug = false;                                    // turn off debugging
    }

    public Point2d () {                                    // Default constructor
        this (0.0, 0.0);                                  // Invokes 2 parameter Point2D constructor
    }
    // Note that a this() invocation must be the BEGINNING of
    // statement body of constructor

    public Point2d (Point2d pt) {                          // Another constructor
        x = pt.getX();
        y = pt.getY();
    }

    ...
}
```

Simple java program: about 500 symbols

Parsing algorithms

- How long would it take to parse this program?

try all parse trees

about 10^{80} years

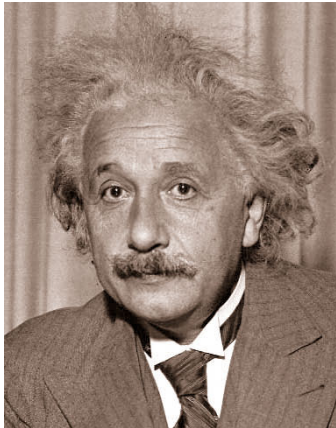
CYK algorithm

about 1 week!

- Can we parse faster?
- **No!** CYK $O(|G|n^3)$ is essentially the fastest known **general-purpose** parsing algorithm for CFGs

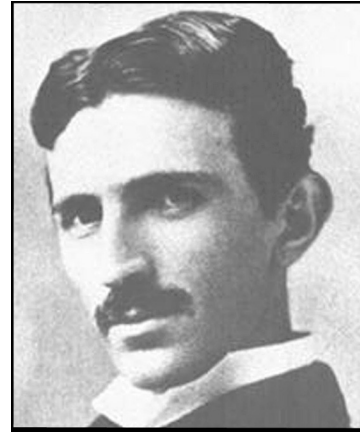
Remark: Coppersmith–Winograd algorithm for multiplying matrices gives an asymptotic worst-case running time $O(n^{2.38} \cdot |G|)$. However, the hidden constant is so high that it is non-practical for present-day computers (Knuth 1997).

Another way of thinking



Scientist:

Find an algorithm that
can parse any CFG



Engineer:

Design your CFG
so it can be parsed
very quickly

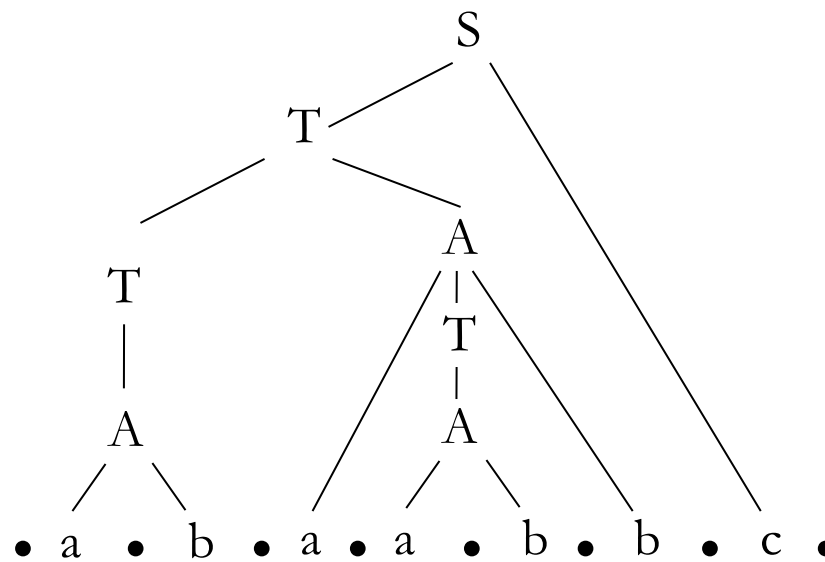
Left-to-right parsing

$$S \rightarrow Tc^{(1)}$$

$$T \rightarrow TA^{(2)} \mid A^{(3)}$$

$$A \rightarrow aTb^{(4)} \mid ab^{(5)}$$

input: abaabbcc



Try to match **to the left** of •

Items

$S \rightarrow Tc^{(1)}$	$T \rightarrow TA^{(2)}$	$T \rightarrow A^{(3)}$	$A \rightarrow aTb^{(4)}$	$A \rightarrow ab^{(5)}$
$S \rightarrow \bullet Tc$	$T \rightarrow \bullet TA$	$T \rightarrow \bullet A$	$A \rightarrow \bullet aTb$	$A \rightarrow \bullet ab$
$S \rightarrow T\bullet c$	$T \rightarrow T\bullet A$	$T \rightarrow A\bullet$	$A \rightarrow a\bullet Tb$	$A \rightarrow a\bullet b$
$S \rightarrow Tc\bullet$	$T \rightarrow TA\bullet$		$A \rightarrow aT\bullet b$	$A \rightarrow ab\bullet$
			$A \rightarrow aTb\bullet$	

- An **item** is a production augmented with a •
- The item is **complete** if the • is the last symbol

Meaning of items

$$A \rightarrow a \bullet T b$$

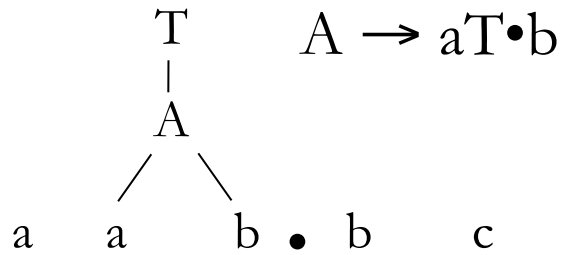
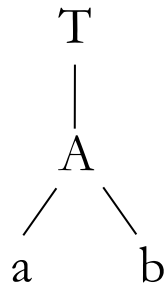
$$A \rightarrow a \bullet b$$

$$S \rightarrow T c^{(1)}$$

$$T \rightarrow T A^{(2)} \mid A^{(3)}$$

$$A \rightarrow a T b^{(4)} \mid a b^{(5)}$$

a • b a a b b c



$$A \rightarrow a T \bullet b$$

Items represent **possibilities**
at various stages of the parsing process

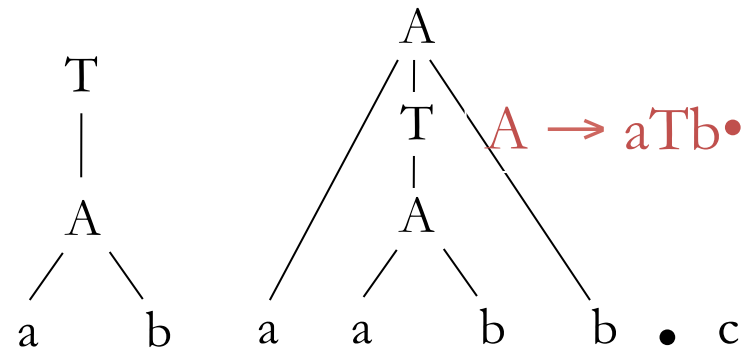
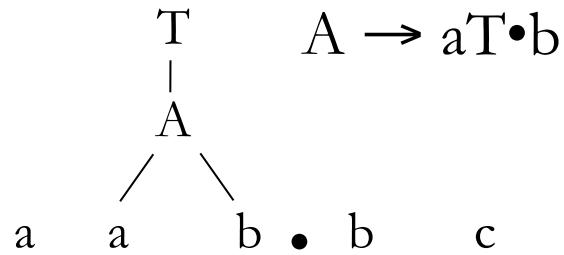
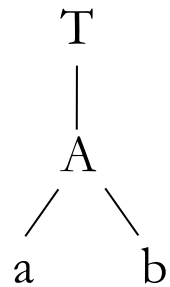
Meaning of items

$S \rightarrow Tc^{(1)}$
$T \rightarrow TA^{(2)} \mid A^{(3)}$
$A \rightarrow aTb^{(4)} \mid ab^{(5)}$

$A \rightarrow a \bullet T b$

$A \rightarrow a \bullet b$

a • b a a b b c



When a **complete item** occurs,
a part of the parse tree is discovered

LR(0) parsing

Move from left to right

$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build
part of parse tree

valid items registry	
$A \rightarrow \bullet aAb$	$A \rightarrow \bullet ab$

$\bullet a \quad a \quad b \quad b$

LR(0) parsing

Move from left to right

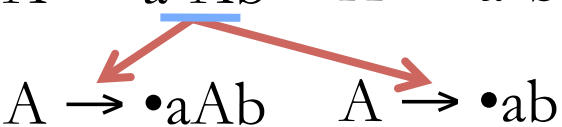
$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build part of parse tree

valid items registry	
$A \rightarrow a \cdot Ab$	$A \rightarrow a \cdot b$
$A \rightarrow \cdot aAb$	$A \rightarrow \cdot ab$



a • a b b

LR(0) parsing

Move from left to right

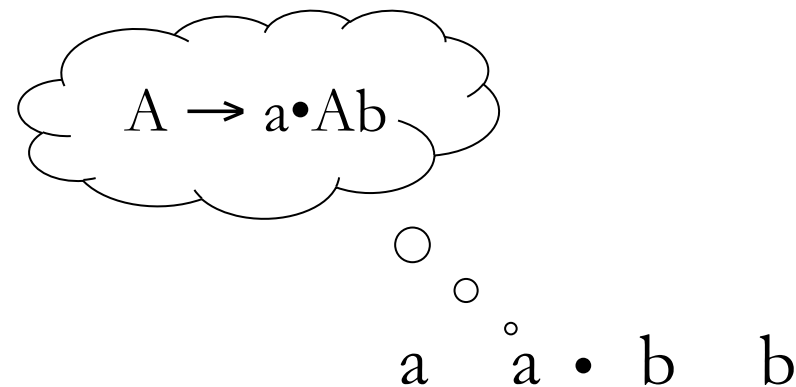
$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build
part of parse tree

valid items registry	
$A \rightarrow a \bullet Ab$	$A \rightarrow a \bullet b$
$A \rightarrow a \bullet Ab$	$A \rightarrow a \bullet b$



LR(0) parsing

Move from left to right

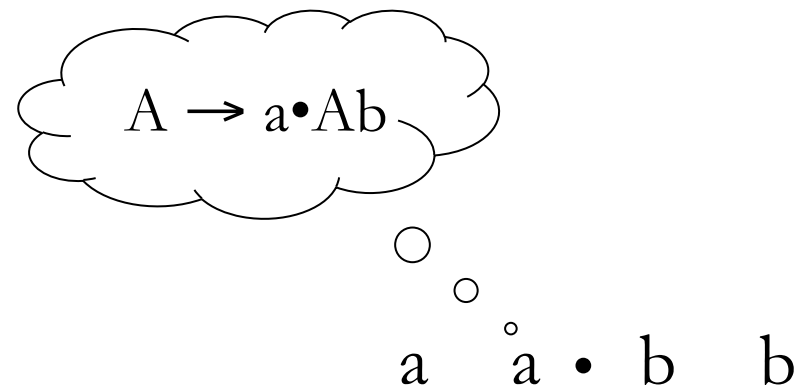
$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build
part of parse tree

valid items registry	
$A \rightarrow a \bullet Ab$	$A \rightarrow a \bullet b$
$A \rightarrow \bullet aAb$	$A \rightarrow \bullet ab$



LR(0) parsing

Move from left to right

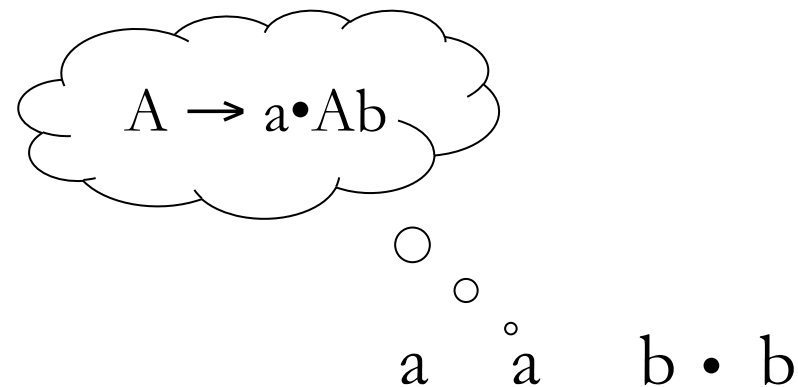
$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build
part of parse tree

valid items registry	
$A \rightarrow aAb$	$A \rightarrow ab\bullet$
$A \rightarrow \bullet aAb$	$A \rightarrow \bullet ab$



LR(0) parsing

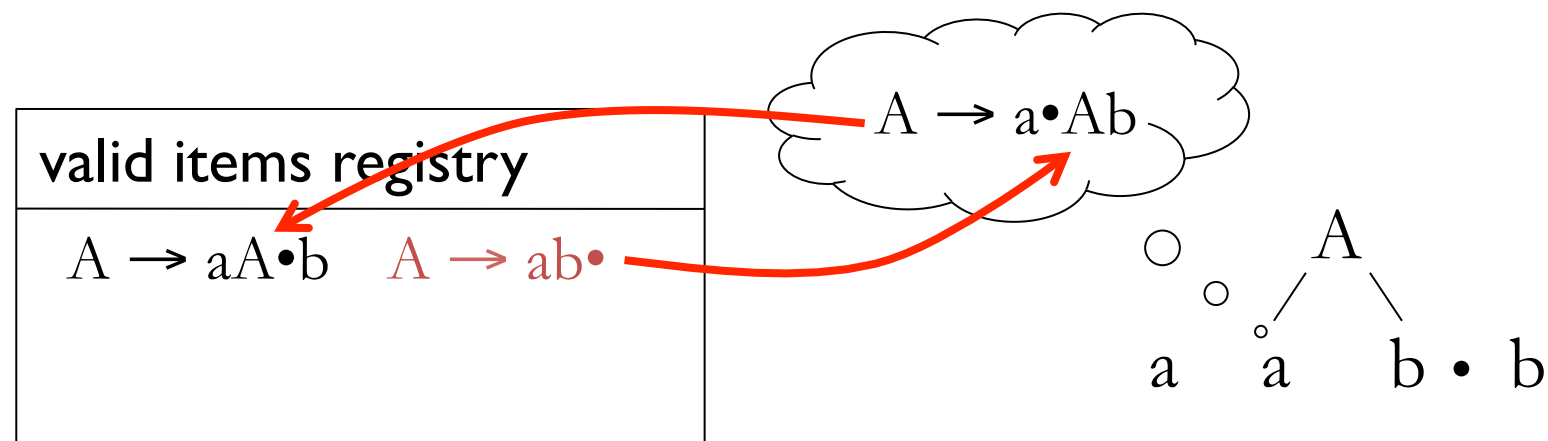
Move from left to right

$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build
part of parse tree



LR(0) parsing

Move from left to right

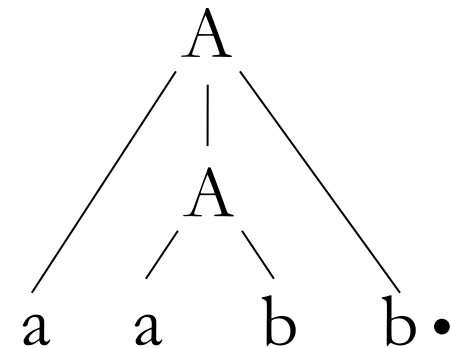
$$A \rightarrow aAb \mid ab$$

Keep track of all possible valid items

Prune the invalid items

When a complete item occurs, build part of parse tree

valid items registry
$A \rightarrow aAb \bullet$

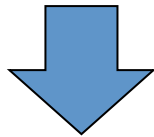


two kinds of actions

valid items registry			
$A \rightarrow a \bullet Ab$	$A \rightarrow a \bullet b$		
$A \rightarrow \bullet aAb$	$A \rightarrow \bullet ab$		

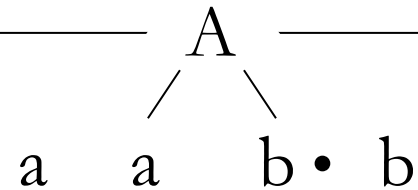
$a \bullet a \bullet b \quad b$

no **complete item**
in registry

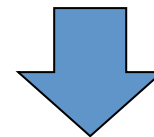


shift

valid items registry	
$A \rightarrow ab \bullet$	



exactly one
complete item



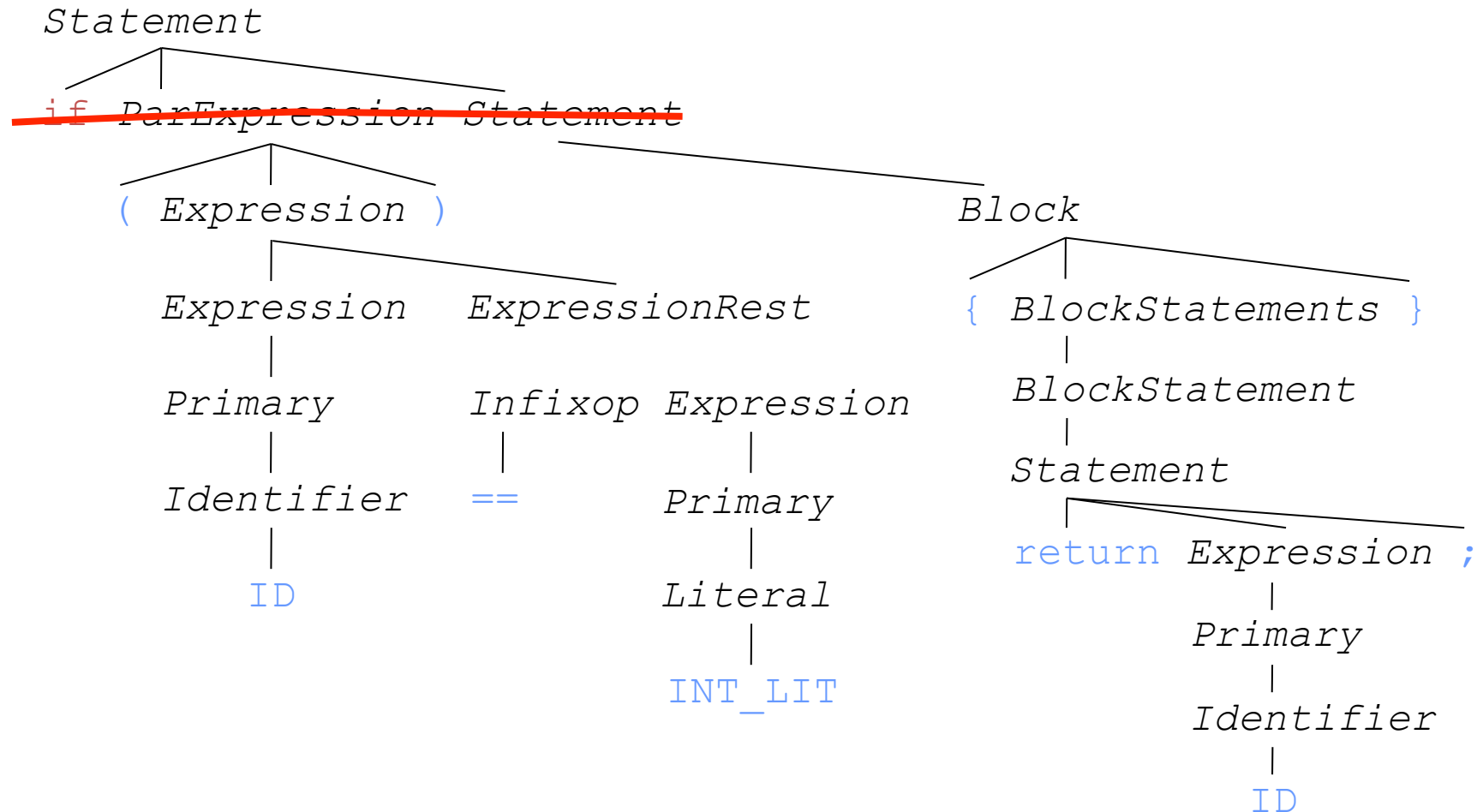
reduce

LR(0) grammars and deterministic PDAs

- The PDA for LR(0) parsing is **deterministic**
- Some context-free languages require non-deterministic PDAs, e.g. $L = \{ww^R : w \in \{a, b\}^*\}$
- The class of deterministic CFGs/PDAs is **less expressive** than CFGs/PDAs
- What can go wrong with LR(0) parsing?

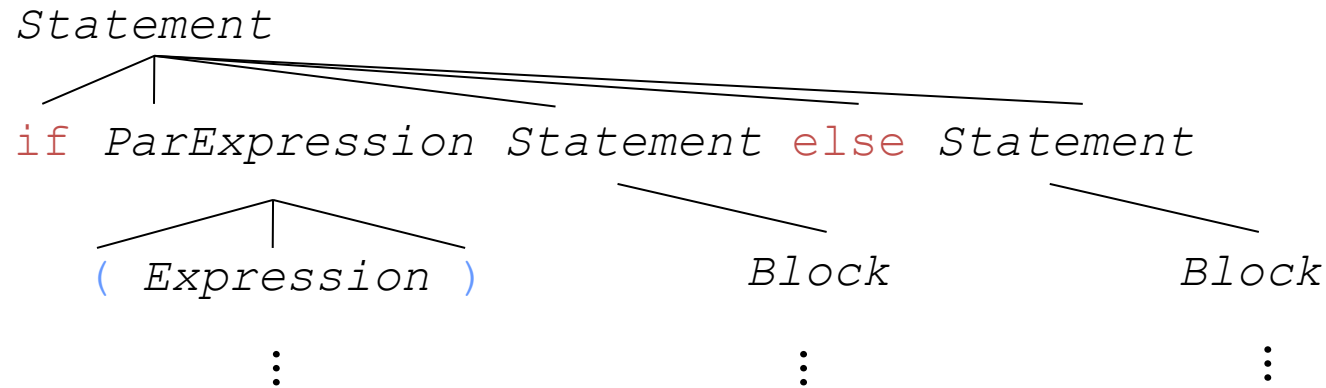
Parsing computer programs

```
if (n == 0) { return x; }  
else { return x + 1; }
```



Parsing computer programs

```
if (n == 0) { return x; }  
else { return x + 1; }
```



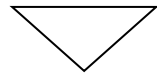
LR(0) parsers cannot tell apart

if ... then **from** if ... then ... else

When you can't LR(0) parse

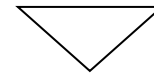
- LR(0) parser can perform two actions:

no complete item
is valid



shift (S)

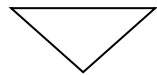
there is one valid item,
and it is complete



reduce (R)

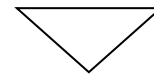
- What if:

some valid items
complete, some not



S / R conflict

more than one valid
complete item



R / R conflict

Dangling else

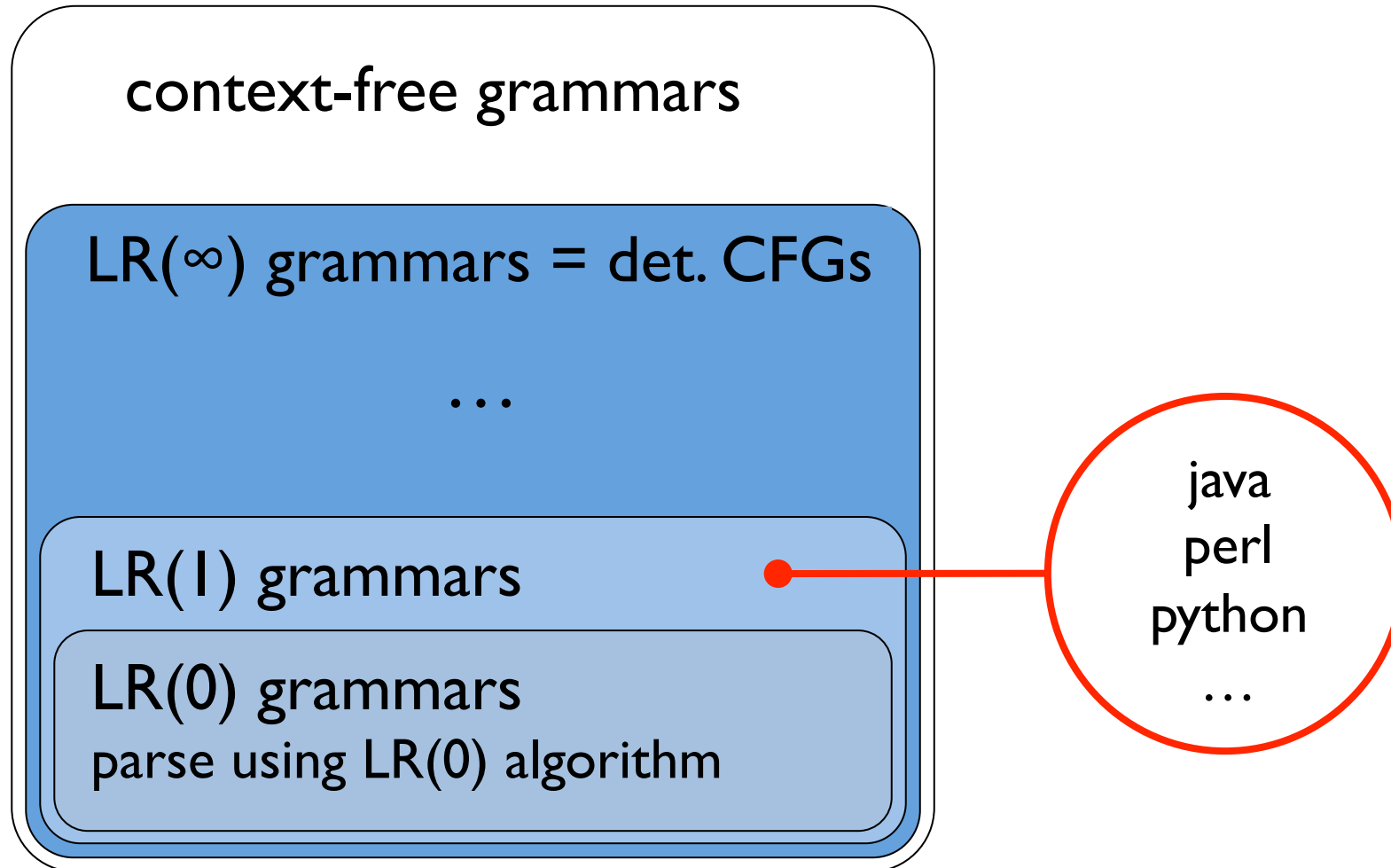
if a then if b then s else s2

Can be interpreted as

- 1. if a then (if b then s) else s2**
- 2. if a then (if b then s else s2)**

shift-reduce conflict

Hierarchy of context-free grammars



LR(k) = grammars that can be parsed like LR(0) but the decision whether to shift or reduce can be based on the next k symbols