1

# Simplifications
# of
# Context-Free Grammars

# A Substitution Rule

$S \rightarrow aB$

$A \rightarrow aaA$

$A \rightarrow abBc$

$B \rightarrow aA$

$B \rightarrow b$

Substitute
$B \rightarrow b$

Equivalent grammar

$S \rightarrow aB \mid ab$

$A \rightarrow aaA$

$A \rightarrow abBc \mid abbc$

$B \rightarrow aA$

# A Substitution Rule

$$S \to aB \mid ab$$

$$A \to aaA$$

$$A \to abBc \mid abbc$$

$$B \to aA$$

Substitute

$$B \to aA$$

$$S \to \cancel{aB} \mid ab \mid aaA$$

$$A \to aaA$$

$$A \to \cancel{abBc} \mid abbc \mid abaAc$$

Equivalent grammar

In general:

$$A \rightarrow xBz$$

$$B \rightarrow y_1$$

Substitute
$$B \rightarrow y_1$$

$$A \rightarrow xBz \mid xy_1z$$

equivalent grammar

$$\lambda - \text{production}: \qquad A \to \lambda$$

Nullable Variable: $\qquad A \Rightarrow \dots \Rightarrow \lambda$

Example Grammar:

$$S \rightarrow aMb$$

$$M \rightarrow aMb$$

$$M \rightarrow \lambda$$

Nullable variable

Final Grammar

$S \rightarrow aMb$

$M \rightarrow aMb$

~~$M \rightarrow \lambda$~~

Substitute
$M \rightarrow \lambda$

$S \rightarrow aMb$

$S \rightarrow ab$

$M \rightarrow aMb$

$M \rightarrow ab$

Unit Production: $A \to B$

(a single variable in both sides)

Observation:

$$A \rightarrow A$$

Is removed immediately

# Example Grammar:

$$S \rightarrow aA$$

$$A \rightarrow a$$

$$A \rightarrow B$$

$$B \rightarrow A$$

$$B \rightarrow bb$$

$$S \rightarrow aA$$

$$A \rightarrow a$$

~~$A \rightarrow B$~~

$$B \rightarrow A$$

$$B \rightarrow bb$$

Substitute
$$A \rightarrow B$$

$$S \rightarrow aA \,|\, aB$$

$$A \rightarrow a$$

$$B \rightarrow A \,|\, B$$

$$B \rightarrow bb$$

$$S \rightarrow aA \,|\, aB$$

$$A \rightarrow a$$

$$B \rightarrow A \,|\, \cancel{B}$$

$$B \rightarrow bb$$

Remove
$$B \rightarrow B$$

$$S \rightarrow aA \,|\, aB$$

$$A \rightarrow a$$

$$B \rightarrow A$$

$$B \rightarrow bb$$

$$S \rightarrow aA \mid aB$$

$$A \rightarrow a$$

~~$$B \rightarrow A$$~~

$$B \rightarrow bb$$

Substitute
$$B \rightarrow A$$
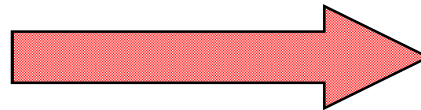
$$S \rightarrow aA \mid aB \mid aA$$

$$A \rightarrow a$$

$$B \rightarrow bb$$

# Remove repeated productions

$$S \rightarrow aA \mid aB \mid \cancel{aA}$$

$$A \rightarrow a$$

$$B \rightarrow bb$$

$$S \rightarrow aA \mid aB$$

$$A \rightarrow a$$

$$B \rightarrow bb$$

# Useless Productions

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

$$S \rightarrow A$$

$$A \rightarrow aA$$ Useless Production

Some derivations never terminate...

$$S \Rightarrow A \Rightarrow aA \Rightarrow aaA \Rightarrow \dots \Rightarrow aa\dots aA \Rightarrow \dots$$

Another grammar:

$$S \rightarrow A$$

$$A \rightarrow aA$$

$$A \rightarrow \lambda$$

$$B \rightarrow bA$$  Useless Production

Not reachable from S

In general:

contains only terminals

if $\quad S \Rightarrow \ldots \Rightarrow xAy \Rightarrow \ldots \Rightarrow w$

$$w \in L(G)$$

then variable $A$ is useful

otherwise, variable $A$ is useless

17

A production $A \rightarrow x$ is useless
if any of its variables is useless

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

Productions

Variables $\quad S \rightarrow A$    useless

useless $\quad A \rightarrow aA$   useless

useless $\quad B \rightarrow C$   useless

useless $\quad C \rightarrow D$   useless

Example Grammar:

$$S \rightarrow aS \mid A \mid C$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

**First:** find all variables that can produce strings with only terminals

$$S \rightarrow aS \mid A \mid C$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

Round 1: $\{A, B\}$

$$S \rightarrow A$$

Round 2: $\{A, B, S\}$

Keep only the variables
that produce terminal symbols: $\{A, B, S\}$

(the rest variables are useless)

$$S \rightarrow aS \mid A \mid \cancel{C}$$

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$\cancel{C \rightarrow aCb}$$

$$\Longrightarrow$$

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$
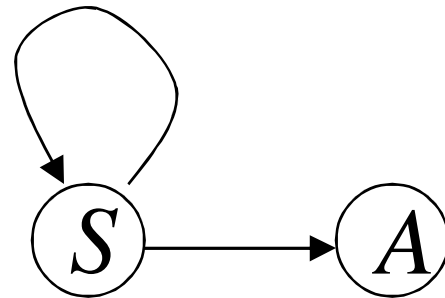
$$B \rightarrow aa$$

Remove useless productions

**Second:** Find all variables reachable from $S$

Use a Dependency Graph

$S \rightarrow aS \mid A$

$A \rightarrow a$

$B \rightarrow aa$



not reachable

# Keep only the variables reachable from S

(the rest variables are useless)

**Final Grammar**

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

~~$B \rightarrow aa$~~

$\Rightarrow$

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$
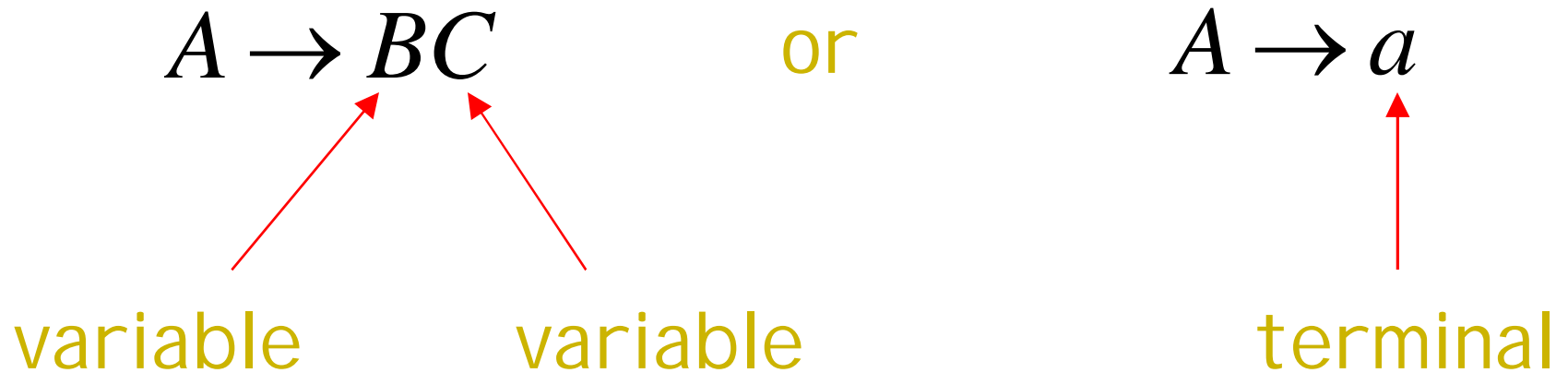
Remove useless productions

# Removing All

- **Step 1:**  Remove Nullable Variables

- **Step 2:**  Remove Unit-Productions

- **Step 3:**  Remove Useless Variables

# Normal Forms
# for
# Context-free Grammars

Each productions has form:

$$A \to BC \qquad \text{or} \qquad A \to a$$

variable     variable     terminal

Examples:

$S \rightarrow AS$

$S \rightarrow a$

$A \rightarrow SA$

$A \rightarrow b$

Chomsky
Normal Form

$S \rightarrow AS$

$S \rightarrow \boxed{AAS}$

$A \rightarrow SA$

$A \rightarrow \boxed{aa}$

Not Chomsky
Normal Form

- Example:

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

Not Chomsky
Normal Form

# Introduce variables for terminals: $T_a, T_b, T_c$

$S \rightarrow ABa$

$A \rightarrow aab$

$B \rightarrow Ac$

$\Longrightarrow$

$S \rightarrow ABT_a$

$A \rightarrow T_a T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

Introduce intermediate variable: $V_1$

$S \to ABT_a$

$A \to T_a T_a T_b$

$B \to AT_c$

$T_a \to a$

$T_b \to b$

$T_c \to c$

$\Longrightarrow$

$S \to AV_1$

$V_1 \to BT_a$

$A \to T_a T_a T_b$

$B \to AT_c$

$T_a \to a$

$T_b \to b$

$T_c \to c$

Introduce intermediate variable: $V_2$

$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$A \rightarrow T_a T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

$\Longrightarrow$

$S \rightarrow AV_1$

$V_1 \rightarrow BT_a$

$A \rightarrow T_a V_2$

$V_2 \rightarrow T_a T_b$

$B \rightarrow AT_c$

$T_a \rightarrow a$

$T_b \rightarrow b$

$T_c \rightarrow c$

## Final grammar in Chomsky Normal Form:

Initial grammar

$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a V_2$$

$$V_2 \rightarrow T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

**In general:**

From any context-free grammar
(which doesn't produce $\lambda$)
not in Chomsky Normal Form

we can obtain:

An equivalent grammar
in Chomsky Normal Form

# The Procedure

First remove:

Nullable variables

Unit productions

Then, for every symbol $a$ :

Add production $T_a \rightarrow a$

In productions: replace $a$ with $T_a$

New variable: $T_a$

Replace any production $A \rightarrow C_1 C_2 \cdots C_n$

with

$$A \rightarrow C_1 V_1$$

$$V_1 \rightarrow C_2 V_2$$

$$\cdots$$

$$V_{n-2} \rightarrow C_{n-1} C_n$$

New intermediate variables: $V_1, V_2, \ldots, V_{n-2}$

**Theorem:** For any context-free grammar (which doesn't produce $\lambda$ ) there is an equivalent grammar in Chomsky Normal Form

# Observations

- Chomsky normal forms are good
  for parsing and proving theorems

- It is very easy to find the Chomsky normal
  form for any context-free grammar

All productions have form:

$$A \rightarrow a\, V_1 V_2 \cdots V_k \qquad k \geq 0$$

symbol          variables

# Observations

- Greinbach normal forms are very good for parsing

- It is hard to find the Greinbach normal form of any context-free grammar

# Compilers

## Program

```
v = 5;
if (v>5)
    x = 12 + v;
while (x !=3) {
  x = x - 3;
  v = 10;
}

......
```

Compiler

## Machine Code

```
Add v,v,0
cmp v,5
jmpIt ELSE
THEN:
  add x, 12,v
ELSE:
WHILE:
cmp x,3
...
```

# Compiler



Lexical analyzer → parser

input

program

output

machine code

A parser knows the grammar
of the programming language

# Parser

PROGRAM → STMT_LIST

STMT_LIST → STMT; STMT_LIST | STMT;

STMT → EXPR | IF_STMT | WHILE_STMT
| { STMT_LIST }

EXPR → EXPR + EXPR | EXPR - EXPR | ID

IF_STMT → if (EXPR) then STMT
| if (EXPR) then STMT else STMT

WHILE_STMT → while (EXPR) do STMT
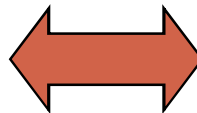
# The parser finds the derivation of a particular input

input
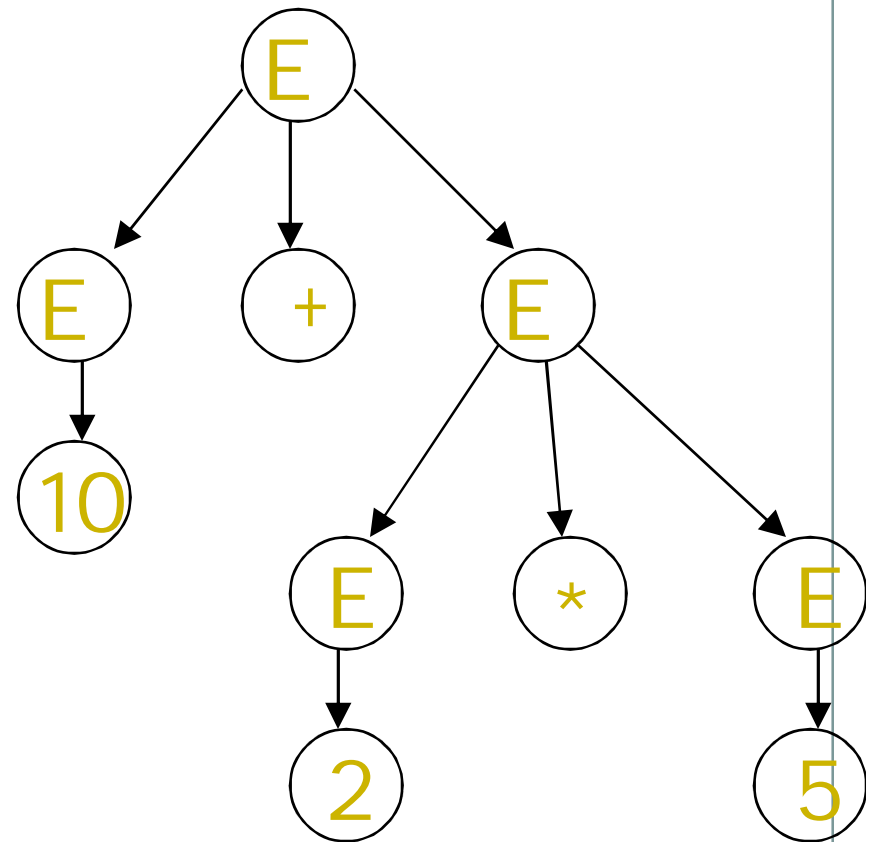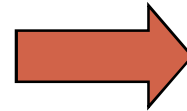
$10 + 2 * 5$

**Parser**

E -> E + E
| E * E
| INT

derivation

E => E + E
=> E + E * E
=> 10 + E*E
=> 10 + 2 * E
=> 10 + 2 * 5

derivation tree

derivation

E => E + E
=> E + E * E
=> 10 + E*E
=> 10 + 2 * E
=> 10 + 2 * 5

E
E  +  E
10
E  *  E
2      5

derivation tree

machine code

mult a, 2, 5
add b, 10, a

49

# Parsing

Parser

input string → grammar → derivation

Example:

input

$aabb$

Parser

$$S \rightarrow SS$$
$$S \rightarrow aSb$$
$$S \rightarrow bSa$$
$$S \rightarrow \lambda$$

derivation

?

# Exhaustive Search

$$S \rightarrow SS \mid aSb \mid bSa \mid \lambda$$

Phase 1:

$$S \Rightarrow SS$$

$$S \Rightarrow aSb$$

$$S \Rightarrow bSa$$

$$S \Rightarrow \lambda$$

Find derivation of

$$aabb$$

All possible derivations of length 1

$$S \Rightarrow SS \qquad\qquad\qquad aabb$$

$$S \Rightarrow aSb$$

$$\cancel{S \Rightarrow bSa}$$

$$\cancel{S \Rightarrow \lambda}$$

$S \rightarrow SS \mid aSb \mid bSa \mid \lambda$

$$S \Rightarrow SS \Rightarrow SSS$$

$$S \Rightarrow SS \Rightarrow aSbS$$

$aabb$

Phase 1

$$S \Rightarrow \cancel{SS \Rightarrow bSaS}$$

$$S \Rightarrow SS$$

$$S \Rightarrow SS \Rightarrow S$$

$$S \Rightarrow aSb$$

$$S \Rightarrow aSb \Rightarrow aSSb$$

$$S \Rightarrow aSb \Rightarrow aaSbb$$

$$S \Rightarrow \cancel{aSb \Rightarrow abSab}$$

$$S \Rightarrow \cancel{aSb \Rightarrow ab}$$

$$S \rightarrow SS \mid aSb \mid bSa \mid \lambda$$

$$S \Rightarrow SS \Rightarrow SSS$$

$$S \Rightarrow SS \Rightarrow aSbS$$

$$S \Rightarrow SS \Rightarrow S$$

$$aabb$$

$$S \Rightarrow aSb \Rightarrow aSSb$$

$$S \Rightarrow aSb \Rightarrow aaSbb$$

Phase 3

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$$

# Final result of exhaustive search
## (top-down parsing)

Parser

$$S \rightarrow SS$$
$$S \rightarrow aSb$$
$$S \rightarrow bSa$$
$$S \rightarrow \lambda$$

input

$aabb$

derivation

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$$

# Time complexity of exhaustive search

Suppose there are no productions of the form

$$A \rightarrow \lambda$$

$$A \rightarrow B$$

Number of phases for string $w$ : approx. |w|

For grammar with $k$ rules

Time for phase 1: $k$

$k$    possible derivations

Time for phase 2:  $k^2$

$$k^2 \quad \text{possible derivations}$$

Time for phase |w| is k$^{|w|}$:

A total of k$^{|w|}$ possible derivations

Total time needed for string $w$:

$$k + k^2 + \cdots + k^{|w|}$$

phase 1     phase 2     phase |w|

Extremely bad!!!

# For general context-free grammars:

There exists a parsing algorithm
that parses a string $|w|$
in time $|w|^3$

The CYK parser