# Formal Languages
# Simplifications of CFGs

# A Substitution Rule

$S \rightarrow aB$

$A \rightarrow aaA$

$A \rightarrow abBc$

$B \rightarrow aA$

$B \rightarrow b$

**Substitute** $B \rightarrow b$

Equivalent grammar

$S \rightarrow aB \,|\, ab$

$A \rightarrow aaA$

$A \rightarrow abBc \,|\, abbc$

$B \rightarrow aA$

# A Substitution Rule

$$S \rightarrow aB \,|\, ab$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc \,|\, abbc$$

$$B \rightarrow aA$$

Substitute

$$B \rightarrow aA$$

$$S \rightarrow a\!\!\!\times\!\!\!B \,|\, ab \,|\, aaA$$

$$A \rightarrow aaA$$

$$A \rightarrow ab\!\!\!\times\!\!\!Bc \,|\, abbc \,|\, abaAc$$

Equivalent grammar

In general:

$$A \rightarrow xBz$$

$$B \rightarrow y_1$$

Substitute
$$B \rightarrow y_1$$

$$A \rightarrow xBz \mid xy_1z$$

equivalent
grammar

# Language?

# Nullable Variables

$\lambda-$production: $\qquad A \rightarrow \lambda$

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

# Removing Nullable Variables

Example Grammar:

$$S \rightarrow aMb$$

$$M \rightarrow aMb$$

$$M \rightarrow \lambda$$

Nullable variable

$$S \rightarrow aMb$$

$$M \rightarrow aMb$$

~~$M \rightarrow \lambda$~~

**Substitute**

$M \rightarrow \lambda$

$$S \rightarrow aMb$$

$$S \rightarrow ab$$

$$M \rightarrow aMb$$

$$M \rightarrow ab$$

8

# Unit-Productions

Unit Production:     $A \rightarrow B$

(single variables on both sides)

# Removing Unit Productions

Observation:

$$A \to A$$

Is removed immediately

# Example Grammar:

$$S \rightarrow aA$$

$$A \rightarrow a$$

$$A \rightarrow B$$

$$B \rightarrow A$$

$$B \rightarrow bb$$

$$S \to aA$$

$$A \to a$$

$$\cancel{A \to B}$$

$$B \to A$$

$$B \to bb$$

Substitute
$$A \to B$$

$$S \to aA \mid aB$$

$$A \to a$$

$$B \to A \mid B$$

$$B \to bb$$

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

$$A \rightarrow a$$

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

$$B \rightarrow bb$$

Remove
$$B \rightarrow B$$

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

$$A \rightarrow a$$

$$B \rightarrow A$$

$$B \rightarrow bb$$

13

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

$A \rightarrow a$

~~$B \rightarrow A$~~

$B \rightarrow bb$

**Substitute**
$B \rightarrow A$
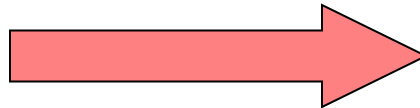
$S \rightarrow aA \,|\, aB \,|\, aA$

$A \rightarrow a$

$B \rightarrow bb$

# Remove repeated productions

$$S \rightarrow aA \mid aB \mid \cancel{aA}$$
$$A \rightarrow a$$
$$B \rightarrow bb$$

$\Longrightarrow$

**Final grammar**

$$S \rightarrow aA \mid aB$$
$$A \rightarrow a$$
$$B \rightarrow bb$$

# Language?

# Useless Productions

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$

$$S \rightarrow A$$

$$A \rightarrow aA$$  Useless Production

Some derivations never terminate…

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

Another grammar:

$$S \rightarrow A$$

$$A \rightarrow aA$$

$$A \rightarrow \lambda$$

$$B \rightarrow bA \quad \text{Useless Production}$$

Not reachable from S

In general:

contains only terminals

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

$$w \in L(G)$$

then variable $A$ is useful

otherwise, variable $A$ is useless

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

$$S \rightarrow aSb$$

$$S \rightarrow \lambda$$   Productions

Variables

$S \rightarrow A$   useless

useless   $A \rightarrow aA$   useless

useless   $B \rightarrow C$   useless

useless   $C \rightarrow D$   useless

20

# Removing Useless Productions

Example Grammar:

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

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

# Remove useless productions

**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:   $\{A, B, S\}$

(other variables are useless)

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

$A \rightarrow a$

$B \rightarrow aa$

$\cancel{C \rightarrow aCb}$

$\Rightarrow$

$S \rightarrow aS \mid A$

$A \rightarrow a$

$B \rightarrow aa$
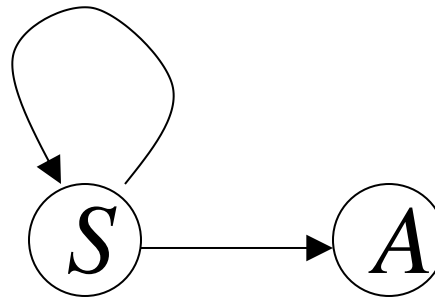
Remove useless productions

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

Use a Dependency Graph

$S \rightarrow aS \mid A$

$A \rightarrow a$

$B \rightarrow aa$



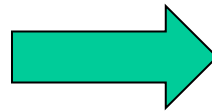not reachable

# Keep only the variables reachable from S

Final Grammar

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

~~$B \rightarrow aa$~~

$$\Longrightarrow$$

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

Remove useless productions
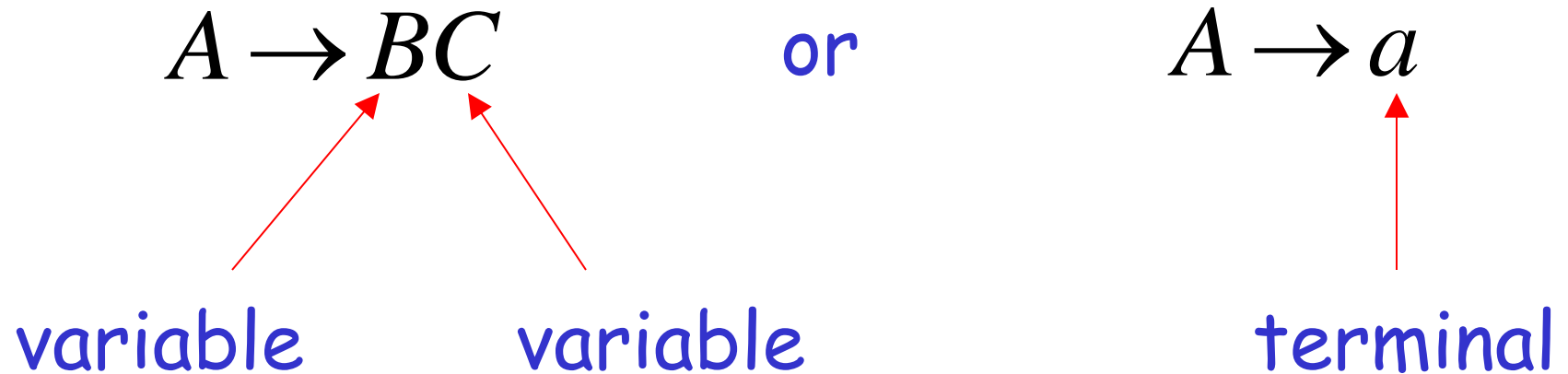
# Removing All

**Step 1:** Remove Nullable Variables

**Step 2:** Remove Unit-Productions

**Step 3:** Remove Useless Variables

# Normal Forms
# for
# Context-free Grammars

# Chomsky Normal Form

Each production has form:

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

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

# Conversion to Chomsky Normal Form

Example:

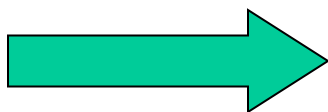$$S \rightarrow ABa$$

$$A \rightarrow aab$$

$$B \rightarrow Ac$$

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: $V_1$

$$S \rightarrow ABT_a$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

$$S \rightarrow AV_1$$

$$V_1 \rightarrow BT_a$$

$$A \rightarrow T_a T_a T_b$$

$$B \rightarrow AT_c$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

$$T_c \rightarrow c$$

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:

$$S \to AV_1$$

$$V_1 \to BT_a$$

$$A \to T_aV_2$$

$$V_2 \to T_aT_b$$

$$B \to AT_c$$

$$T_a \to a$$

$$T_b \to b$$

$$T_c \to c$$

## Initial grammar

$$S \to ABa$$

$$A \to aab$$

$$B \to Ac$$

# 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 $\quad T_a \rightarrow a$

In productions:  replace $\ a\ $ with $\ T_a$

New variable: $\ T_a$

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

with $\quad A \rightarrow C_1 V_1$

$\qquad V_1 \rightarrow C_2 V_2$

$\qquad \cdots$

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

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

# exercise

Find CNF for this grammar:

S -> 0A0 | 1B1 | BB
A -> C
B -> S | A
C -> S | epsilon

(exercise 7.1.3, Hopcroft, Motwani, Ullman)

# Greibach Normal Form

All productions have form:

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

terminal      variables

Examples:

$$S \rightarrow cAB$$

$$A \rightarrow aA \mid bB \mid b$$

$$B \rightarrow b$$

$$S \rightarrow abSb$$

$$S \rightarrow aa$$

Greibach
Normal Form

Not Greibach
Normal Form

# Conversion to Greibach Normal Form:

$$S \rightarrow abSb$$

$$S \rightarrow aa$$

$\Longrightarrow$

$$S \rightarrow aT_bST_b$$

$$S \rightarrow aT_a$$

$$T_a \rightarrow a$$

$$T_b \rightarrow b$$

Greibach
Normal Form

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

# Observations

- Greibach normal forms are very good for parsing

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

# Try to compute Greibach Normal Form for grammar in CNF example

# 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
jmplt 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 | INT

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

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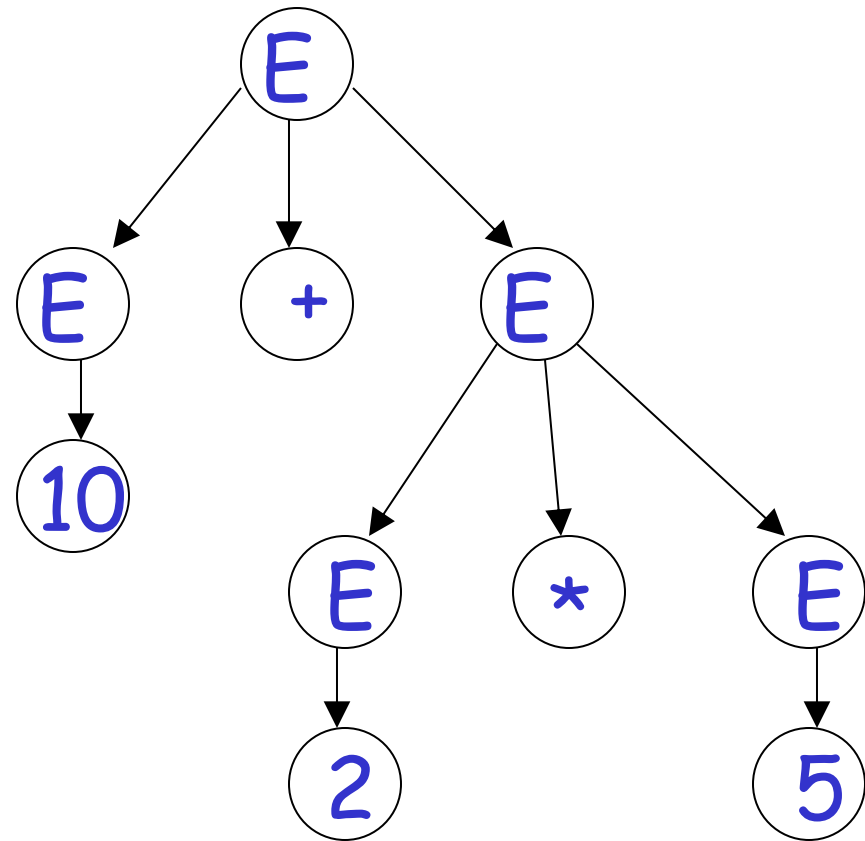
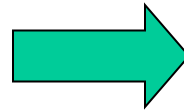# The parser finds the derivation of a particular input

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

⟷

# derivation tree



machine code

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

# Parsing

**Parser**

input string → | grammar | → derivation

58

Example:

Parser

input

$aabb$

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

derivation

?

# Parsing algorithm?

# Exhaustive Search

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

Phase 1:

$$S \Rightarrow SS$$

$$S \Rightarrow aSb$$

$$S \Rightarrow bSa$$

$$S \Rightarrow \lambda$$

Find derivation of

$$aabb$$

All possible derivations of length 1

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

$$S \Rightarrow aSb$$

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

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

Phase 2    $S \rightarrow SS \,|\, aSb \,|\, bSa \,|\, \lambda$

$$S \Rightarrow SS \Rightarrow SSS$$

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

Phase 1

~~$S \Rightarrow SS \Rightarrow bSaS$~~    + 2 more

$S \Rightarrow SS$    $S \Rightarrow SS \Rightarrow S$

$S \Rightarrow aSb$    $S \Rightarrow aSb \Rightarrow aSSb$

$$S \Rightarrow aSb \Rightarrow aaSbb$$

~~$S \Rightarrow aSb \Rightarrow abSab$~~

~~$S \Rightarrow aSb \Rightarrow ab$~~

63

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

**Phase 2**

$$S \Rightarrow SS \Rightarrow SSS$$

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

$$S \Rightarrow SS \Rightarrow S$$

$$S \Rightarrow aSb \Rightarrow aSSb$$

$$S \Rightarrow aSb \Rightarrow aaSbb$$

**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

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

# Is exhaustive search a good parsing algorithm?

# Time complexity of exhaustive search

Suppose there are no productions of the form

$$A \rightarrow \lambda$$

$$A \rightarrow B$$

Number of phases for string $w:$ $2|w|$

For grammar with $k$ rules

Time for phase 1: $k$

$k$    possible derivations

Time for phase 2: $k^2$

$k^2$    possible derivations

Time for phase $2|w|$:  $k^{2|w|}$

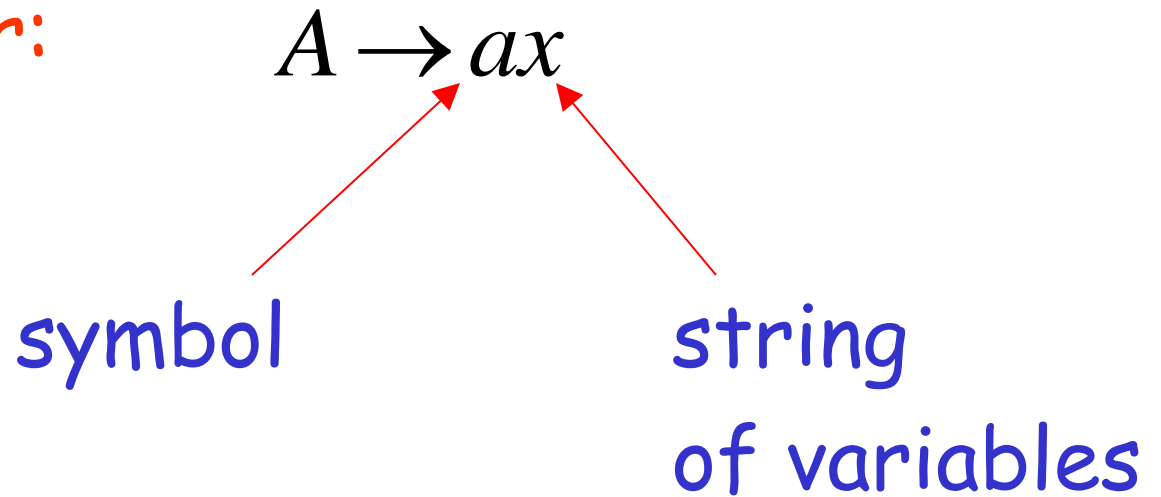$k^{2|w|}$  possible derivations

Total time needed for string $w$:

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

phase 1            phase 2            phase 2|w|

Extremely bad!!!

There exist faster algorithms
for specialized grammars

S-grammar:

$$A \rightarrow ax$$

symbol        string
of variables

Pair $(A, a)$ appears once

S-grammar example:

$$S \rightarrow aS$$

$$S \rightarrow bSS$$

$$S \rightarrow c$$

Each string has a unique derivation

$$S \Rightarrow aS \Rightarrow abSS \Rightarrow abcS \Rightarrow abcc$$

For S-grammars:

In the exhaustive search parsing
there is only one choice in each phase

Time for a phase:  $1$

Total time for parsing string  $w$:  $|w|$

# For general context-free grammars:

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

(we will show this in the next class)

# The CYK Parser

# The CYK Membership Algorithm

Input:

- Grammar $G$ in Chomsky Normal Form

- String $w$

Output:

find if $w \in L(G)$

# The Algorithm

Input example:

- Grammar $G$:

$$S \rightarrow AB$$
$$A \rightarrow BB$$
$$A \rightarrow a$$
$$B \rightarrow AB$$
$$B \rightarrow b$$

- String $w$: $aabbb$

# *aabbb*

a        a        b        b        b

aa       ab       bb       bb

aab      abb      bbb

aabb     abbb

aabbb

$$S \rightarrow AB$$

$$A \rightarrow BB$$

$$A \rightarrow a$$

$$B \rightarrow AB$$

$$B \rightarrow b$$

| a | a | b | b | b |
| A | A | B | B | B |
|---|---|---|---|---|
| aa | ab | bb | bb | |
| aab | abb | bbb | | |
| aabb | abbb | | | |
| aabbb | | | | |

$$S \rightarrow AB$$

$$A \rightarrow BB$$

$$A \rightarrow a$$

$$B \rightarrow AB$$

$$B \rightarrow b$$

| a | a | b | b | b |
|---|---|---|---|---|
| A | A | B | B | B |

| aa | ab | bb | bb |
|---|---|---|---|
| | S,B | A | A |

| aab | abb | bbb |
|---|---|---|

| aabb | abbb |
|---|---|

aabbb

$S \rightarrow AB$

$A \rightarrow BB$

$A \rightarrow a$

$B \rightarrow AB$

$B \rightarrow b$

| a | a | b | b | b |
|---|---|---|---|---|
| A | A | B | B | B |

| aa | ab | bb | bb |
|---|---|---|---|
| | S,B | A | A |

| aab | abb | bbb |
|---|---|---|
| S,B | A | S,B |

| aabb | abbb |
|---|---|
| A | S,B |

| aabbb |
|---|
| S,B |

Therefore: $aabbb \in L(G)$

Time Complexity: $|w|^3$

Observation: The CYK algorithm can be easily converted to a parser (bottom up parser)

The following slides are courtesy of Professor Papp, University of Debrecen.

1,7

1,1

7,7

a a b b a b a

## Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



| A | A | B,C | B,C | A | B,C | A |

a    a    b    b    a    b    a

## Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA

|  | | S | A,S | D | S | D |
|---|---|---|---|---|---|---|
| A | A | B,C | B,C | A | B,C | A |

a    a    b    b    a    b    a

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



The CYK parsing chart:

|  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|
|  |  | B | S |  |  |  |
|  | S | A,S | D | S | D |  |
| A | A | B,C | B,C | A | B,C | A |

a   a   b   b   a   b   a

## Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |
|   |   |   |   |   |   |   |
|   | S |   |   |   |   |   |
|   |   | B | S |   |   |   |
|   |   | S | A,S | D | S | D |
| A | A | B,C | B,C | A | B,C | A |

**a    a    b    b    a    b    a**

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



|  | S |  |  |  |  |  |
|---|---|---|---|---|---|---|
|  | B | S |  |  |  |  |
| S | A,S | D | S | D |  |  |
| A | A | B,C | B,C | A | B,C | A |

a a | b b a b a

## Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



| S | | | | | | |
| | B | S | | | | |
| S | A,S | D | S | D | | |
| A | A | B,C | B,C | A | B,C | A |

a  a  b | b  a  b  a

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



| | | | | D | | | | |
| | | | S | | D | | | |
| | | | | B | | S | | |
| | | S | | A,S | | D | | S | | D |
| A | | A | | B,C | | B,C | | A | | B,C | | A |

a  a  b  b  |  a  b  a

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



Row 4: S | D | S,B | C

Row 3: B | S

Row 2: S | A,S | D | S | D

Row 1: A | A | B,C | B,C | A | B,C | A

a a b b a b a

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



CYK parsing chart:

Row 7 (top): S
Row 6: A,D
Row 5: S, S,B
Row 4: S, D, C
Row 3: B, S
Row 2: S, A,S, D, S, D
Row 1 (bottom): A, A, B,C, B,C, A, B,C, A

Input string: a a b b a b a

# Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



CYK parsing triangle:

Row 1: S | B,C
Row 2: S | A,D
Row 3: S | D | S,B | C
Row 4: B | S
Row 5: S | A,S | D | S | D
Row 6: A | A | B,C | B,C | A | B,C | A

a  a  b  b  a  b  a

## Grammar

S→AB
S→CD
S→CB
S→SS
A→BC

A→a
B→SC
B→b
C→DD
C→b
D→BA



CYK parsing chart:

Row 7: S

Row 6: S | B,C

Row 5: S | A,D

Row 4: S | D | S,B | C

Row 3: B | S

Row 2: S | A,S | D | S | D

Row 1: A | A | B,C | B,C | A | B,C | A

Input: a  a  b  b  a  b  a

# exercise

Parse "baaba" for this grammar

S -> AB | BC
A -> BA | a
B -> CC | b
C -> AB | a

(Hopcroft, Motwani, Ullman, p301)