# Context free Grammar

- Context-free grammar is a 4-tuple
  $G = (N, T, P, S)$ where
    - $T$ is a finite set of tokens (*terminal* symbols)
    - $N$ is a finite set of *nonterminals*
    - $P$ is a finite set of *productions* of the form
        $\alpha \rightarrow \beta$
      where $\alpha \in N$ and $\beta \in (N \cup T)^*$
    - $S \in N$ is a designated *start symbol*

# Example Grammar

Context-free grammar for simple expressions:

$G$ = <{$expr$,op,digit}, {**+,-,\*,/,0,1,2,3,4,5,6,7,8,9,),(**}, $P$,expr>

with productions $P$ =

$expr\rightarrow expr$  op expr

expr $\rightarrow$(expr)

expr $\rightarrow$ $digit$

$digit \rightarrow$ **0** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9**

$op \rightarrow$ **+** | **-** | **\*** | **/**

# Notational Conventions Used

- <u>Terminals:</u> Lower case letters, operator symbols, punctuation symbols, digits, bolface strings are all terminals

- <u>Non Terminals:</u> Upper case letters, lower case italic names are usually non terminals

- Greek letters such as $\alpha,\beta,\gamma$ represent strings of grammars symbols. Thus a generic production can be written as $A \rightarrow \alpha$

# Example

- Design a CFG for the language

$$L(G) = \{0^n 1^m \mid n <> m\}$$

There are two cases:

- For n>m
- For n<m
- Write two separate set of rules and combine them

# Example

- For n>m

    $S1 \to AB$

    $B \to 0B1 \mid \epsilon$

    $A \to 0A \mid 0$

  For n<m

    $S2 \to XY$

    $X \to 0X1 \mid \epsilon$

    $Y \to 1Y \mid 1$

Combining both:

$S \to S1 \mid S2$

# Examples

- Write a CFG that generates Equal number of a's, b's and c's?

# Derivations

- The *one-step derivation* is defined by
  $$\alpha\ A\ \beta \Rightarrow \alpha\ \gamma\ \beta, \text{ where } A \rightarrow \gamma \text{ is a production in the grammar}$$

- In addition, we define
  - $\Rightarrow$ is *leftmost* $\Rightarrow_{lm}$ if $\alpha$ does not contain a nonterminal
  - $\Rightarrow$ is *rightmost* $\Rightarrow_{rm}$ if $\beta$ does not contain a nonterminal
  - Transitive closure $\Rightarrow^*$ (zero or more steps)
  - Positive closure $\Rightarrow^+$ (one or more steps)

- The *language generated by G* is defined by
  $$L(G) = \{w \in T^* \mid S \Rightarrow^+ w\}$$

# Derivation (Example)

Grammar $G = (\{E\}, \{+, *, (, ), -, id\}, P, E)$ with
productions $P = \begin{aligned} &E \rightarrow E + E \\ &E \rightarrow E * E \\ &E \rightarrow ( E ) \\ &E \rightarrow - E \\ &E \rightarrow id \end{aligned}$

Example derivations:

$$E \Rightarrow - E \Rightarrow - id$$

$$E \Rightarrow_{rm} E + E \Rightarrow_{rm} E + id \Rightarrow_{rm} id + id$$

$$E \Rightarrow^* E$$

$$E \Rightarrow^* id + id$$

$$E \Rightarrow^+ id * id + id$$

# Derivation for the Example Grammar

$list$
$\Rightarrow list$ **+** $digit$
$\Rightarrow list$ **-** $digit$ **+** $digit$
$\Rightarrow digit$ **-** $digit$ **+** $digit$
$\Rightarrow$ **9 -** $digit$ **+** $digit$
$\Rightarrow$ **9 - 5 +** $digit$
$\Rightarrow$ **9 - 5 + 2**

This is an example *leftmost derivation*, because we replaced the leftmost nonterminal (underlined) in each step. Likewise, a *rightmost derivation* replaces the rightmost nonterminal in each step

# Chomsky Hierarchy: Language Classification

- A grammar $G$ is said to be

  - *Regular* if it is *right linear* where each production is of the form
    $$A \rightarrow w\, B \qquad \text{or} \qquad A \rightarrow w$$
    or *left linear* where each production is of the form
    $$A \rightarrow B\, w \qquad \text{or} \qquad A \rightarrow w$$

  - *Context free* if each production is of the form
    $$A \rightarrow \alpha$$
    where $A \in N$ and $\alpha \in (N \cup T)^*$

  - *Context sensitive* if each production is of the form
    $$\alpha\, A\, \beta \rightarrow \alpha\, \gamma\, \beta$$
    where $A \in N$, $\alpha, \gamma, \beta \in (N \cup T)^*$, $|\gamma| > 0$

  - *Unrestricted*

# Chomsky Hierarchy

$L(regular) \subset L(context\ free)$
$\subset L(context\ sensitive) \subset L(unrestricted)$

Examples:

Every *finite language* is regular!
(construct a FSA for strings in $L(G)$)

$L_1 = \{\ a^n b^n\ |\ n \geq 1\ \}$ is context free
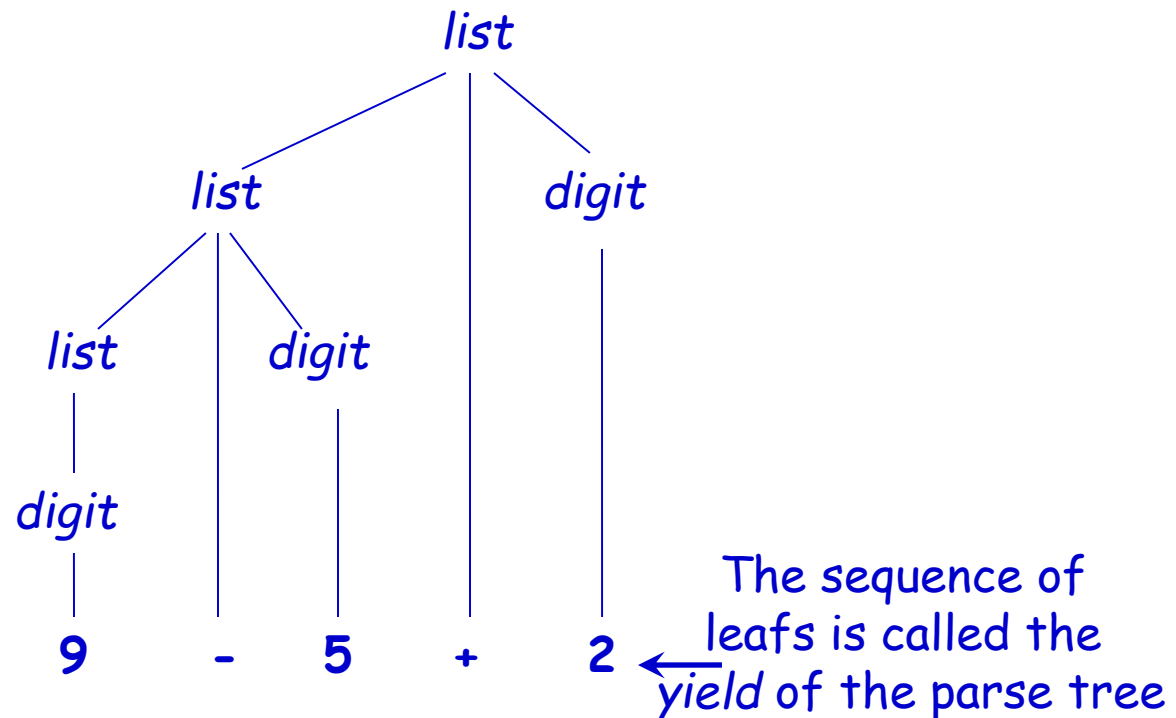
$L_2 = \{\ a^n b^n c^n\ |\ n \geq 1\ \}$ is context sensitive

# Parse Trees

- The *root* of the tree is labeled by the start symbol

- Each *leaf* of the tree is labeled by a terminal (=token) or $\varepsilon$

- Each *interior node* is labeled by a nonterminal

If $A \rightarrow X_1 X_2 \ldots X_n$ is a production, then node $A$ has immediate *children $X_1, X_2, \ldots, X_n$* where $X_i$ is a (non)terminal or $\varepsilon$ ( $\varepsilon$ denotes the *empty string*)

# Parse Tree for the Example Grammar

Parse tree of the string **9-5+2** using grammar *G*



*The sequence of leafs is called the yield of the parse tree*

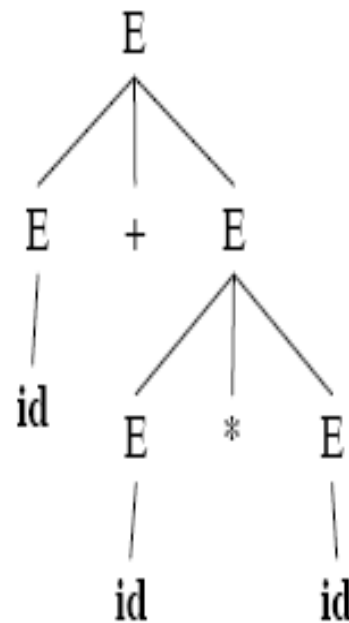# Example of Parse Tree

- Suppose we have the following grammar
  $E \rightarrow E + E$
  $E \rightarrow E * E$
  $E \rightarrow ( E )$
  $E \rightarrow - E$
  $E \rightarrow$ **id**

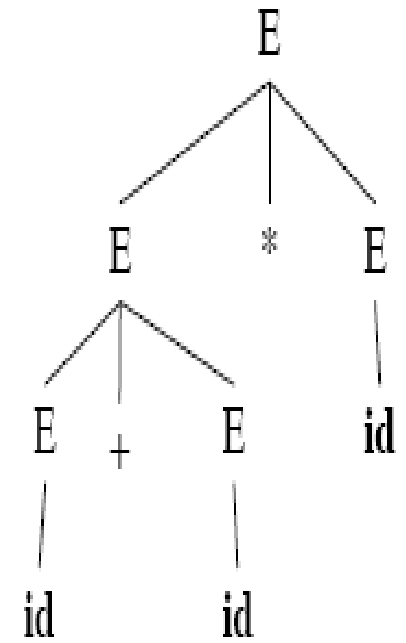Perform Left most derivation, right most derivation and construct a parse tree for the string

**id+id*id**

# Two possible Parse Trees using Leftmost derivation

- $E \Rightarrow E + E$

$\Rightarrow id + E$

$\Rightarrow id + E * E$

$\Rightarrow id + id * E$

$\Rightarrow id + id * id$



- $E \Rightarrow E * E$

$\Rightarrow E + E * E$

$\Rightarrow id + E * E$

$\Rightarrow id + id * E$

$\Rightarrow id + id * id$

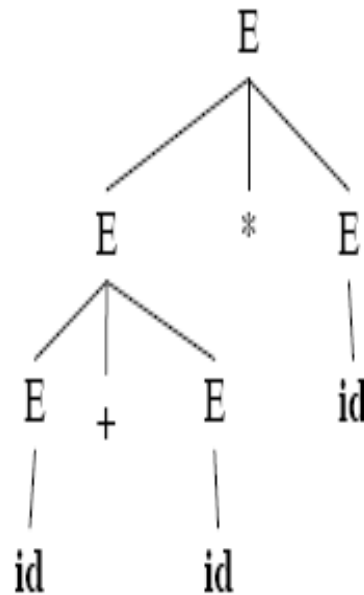# Parse Tree via Right most derivation

$$E \Rightarrow E * E$$
$$\Rightarrow E * id$$
$$\Rightarrow E + E * id$$
$$\Rightarrow E + id * id$$
$$\Rightarrow id + id * id$$

# Ambiguity

- Grammar is ambiguous if more than one parse tree is possible for some string as shown in the previous example. If there are more than one left most derivations or more than one right most derivations.

- Ambiguity is not acceptable
  - Unfortunately, it's undecidable to check whether a given CFG is ambiguous
  - Some CFLs are inherently ambiguous (do not have an unambiguous CFG)

# Ambiguity (cont'd)
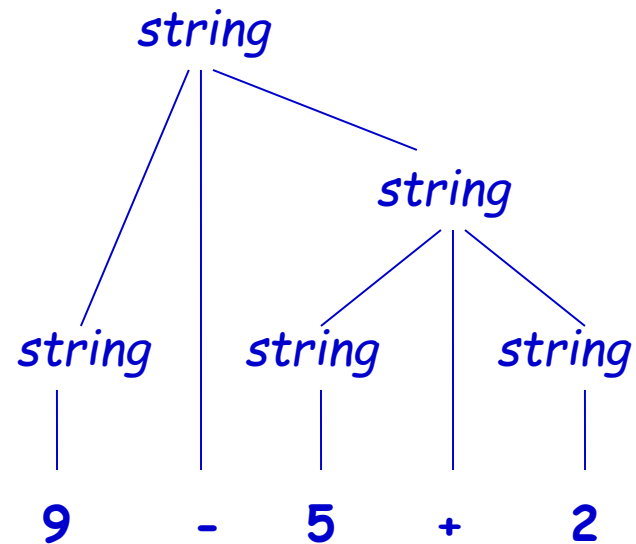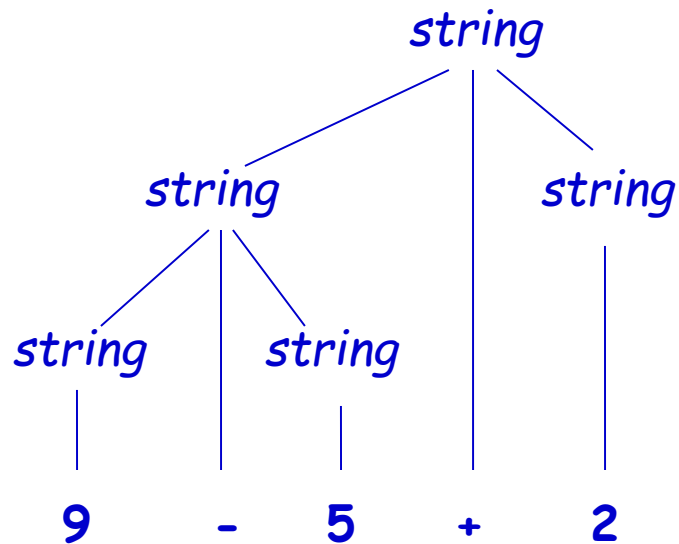
Consider the following context-free grammar:

$G$ = <{*string*}, {**+**,**−**,**0**,**1**,**2**,**3**,**4**,**5**,**6**,**7**,**8**,**9**}, $P$, *string*>

with production $P$ =

*string* → *string* **+** *string* | *string* **−** *string* | **0** | **1** | ... | **9**

This grammar is *ambiguous*, because more than one parse tree represents the string **9-5+2**

# Two Parse Trees for the same string

# Practice

- Show that the following grammar is ambiguous: (Find out strings and two parse trees)

1) S→AB | aaB

   A→ a | Aa

   B→b

2) S→ a | abSb |aAb

   A→ bS | aAAb

3) S→aSb | SS | ε

# 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

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

$A \rightarrow aaA$

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

$B \rightarrow aA$

# A Substitution Rule

$S \rightarrow aB \mid ab$

$A \rightarrow aaA$

$A \rightarrow abBc \mid abbc$

$B \rightarrow aA$

Substitute

$$B \rightarrow aA$$

$S \rightarrow a\cancel{B} \mid ab \mid aaA$

$A \rightarrow aaA$

$A \rightarrow a\cancel{bBc} \mid abbc \mid abaAc$

Equivalent grammar

In general:

$$A \rightarrow xBz$$

$$B \rightarrow y_1$$

Substitute
$$B \rightarrow y_1$$

$$A \rightarrow xBz \mid xy_1z$$

equivalent grammar

# Nullable Variables

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

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
$M \rightarrow \lambda$

Final Grammar

$S \rightarrow aMb$

$S \rightarrow ab$

$M \rightarrow aMb$

$M \rightarrow ab$

# Unit-Productions

Unit Production:     $A \rightarrow B$

(a single variable in both sides)

# Removing Unit Productions

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
$A \rightarrow B$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow A \mid B$

$B \rightarrow bb$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow A \mid \cancel{B}$

$B \rightarrow bb$

**Remove**
$B \rightarrow B$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow A$

$B \rightarrow bb$

$$S \rightarrow aA \mid aB$$

$$A \rightarrow a$$

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

$$B \rightarrow bb$$

Substitute
$$B \rightarrow A$$

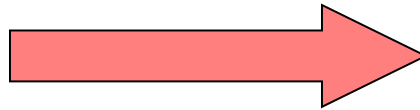$$S \rightarrow aA \mid aB \mid aA$$

$$A \rightarrow a$$

$$B \rightarrow bb$$

# Remove repeated productions

### Final grammar

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

$A \rightarrow a$

$B \rightarrow bb$

$\Longrightarrow$

$S \rightarrow aA \mid aB$

$A \rightarrow a$

$B \rightarrow bb$

# 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:

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

contains only terminals

$w \in L(G)$

then variable $A$ is useful

otherwise, variable $A$ is useless

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

$$S \rightarrow aSb$$

$$S \rightarrow \lambda \quad \text{Productions}$$

Variables $\quad \boxed{S \rightarrow A} \quad$ useless

useless $\quad \boxed{A \rightarrow aA} \quad$ useless

useless $\quad \boxed{B \rightarrow C} \quad$ useless

useless $\quad \boxed{C \rightarrow D} \quad$ useless

# Removing Useless Productions

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)

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

$A \rightarrow a$

$B \rightarrow aa$

$\cancel{C \rightarrow aCb}$

$\Rightarrow$

$S \rightarrow aS \mid A$

$A \rightarrow a$

$B \rightarrow aa$
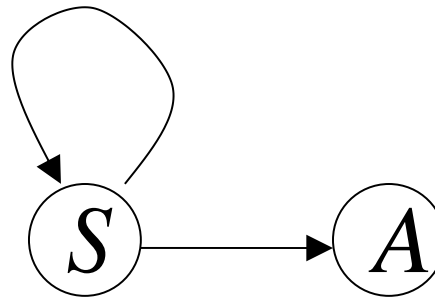
Remove useless productions

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

Use a Dependency Graph

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$B \rightarrow aa$$



$S$ → $A$     $B$

not reachable

# Keep only the variables reachable from S

(the rest variables are useless)

Final Grammar

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

$$\cancel{B \rightarrow aa}$$

➡️

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

Remove useless productions

# Removing All

**Step 1:** Remove Nullable Variables

**Step 2:** Remove Unit-Productions

**Step 3:** Remove Useless Variables

Do it yourself: Why in this order only??