# Sub: Compiler Construction Syntax Analysis PART 1

Compiled for: 7th Sem, CE, DDU

Compiled by: Niyati J. Buch

Ref. : Compilers: Principles, Techniques, and Tools, 2nd Edition
Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman
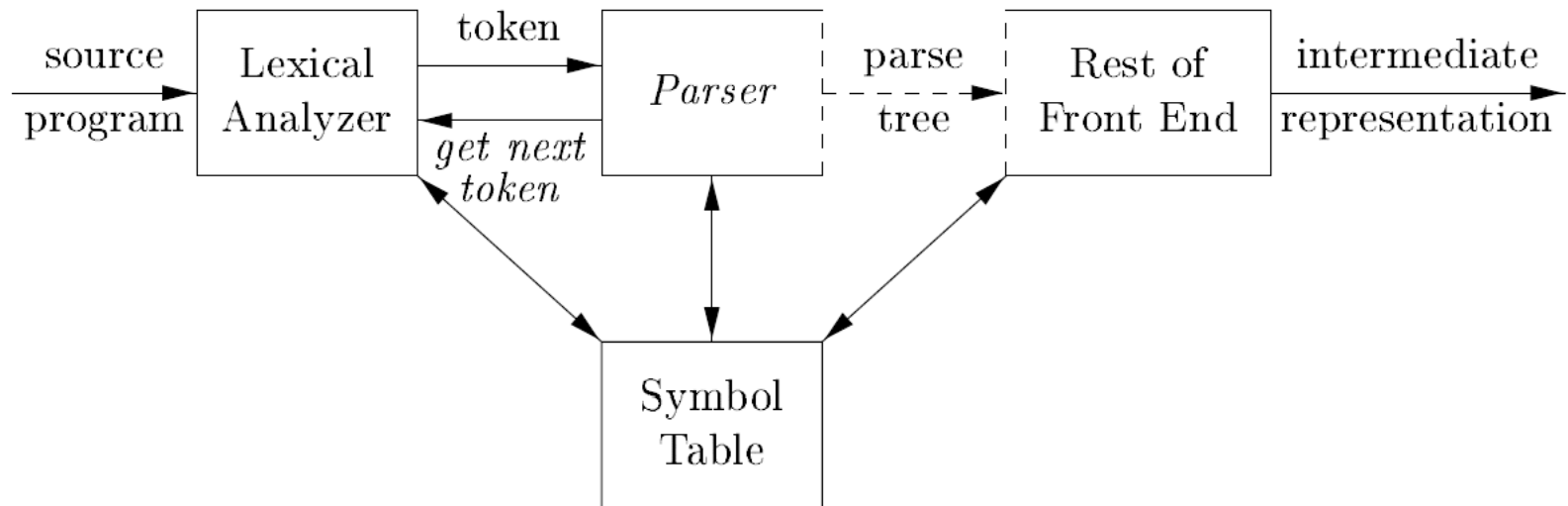
# Topics Covered

- Introduction
  - [Role of a parser](#)
  - [Representative Grammars](#)
  - [Syntax Error Handling](#)
  - [Error-Recovery Strategies](#)
- Context-Free Grammars
  - [Formal Definition](#)
  - [Conventions](#)
  - [Sentinel](#) and [Canonical](#) form
  - [Ambiguity](#)
- Writing a grammar
  - [Eliminating useless variables](#)
  - [Eliminating left recursion](#)
  - [Eliminating left factoring](#)
  - [Elimination of ε productions](#)
  - [Eliminating unit productions](#)

# Introduction to Syntax Analysis

# Role of a parser

- A parser uses a grammar to check structure of tokens.
- It produces a parse tree.
- It checks for syntactic errors and recovery.
- It recognize correct syntax.
- It report errors.

# 3 Types of Parsers

1. **Universal parsing methods**
   – Such as Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar.
   – But they are too inefficient to use in production compilers.
2. **Top down methods**
   – Build the parse tree from top(root) to the bottom(leaves)
3. **Bottom up methods**
   – Builds the parse tree from bottom(leaves) to the top(root)

# Types of parser

- Top-down: **LL** $\rightarrow$ scan **L**eft to right and consider **L**eft most derivative

- Bottom-up: **LR** $\rightarrow$ scan **L**eft to right and consider **R**ight most derivative in reverse

- In both top-down and bottom up, the input to the parser is scanned from **left to right** , one symbol at a time.

- Parsers implemented by hand often use LL grammar( ex. predictive parsing)

- Parsers for the larger class of LR grammars are usually constructed using automated tools.

# Representative Grammars

- Associativity and precedence are captured in the following grammar, for describing expressions, terms, and factors.

- E represents expressions consisting of terms separated by + signs

- T represents terms consisting of factors separated by * signs

- F represents factors that can be either parenthesized expressions or identifiers:

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow ( E ) \mid id$$

# Can you observe any recursion??

E → E + T | T

T → T * F | F

F → ( E ) | id

# Can you observe any recursion??

E → E + T | T

T → T * F | F

F → ( E ) | id

- Recursion on **left side** is observed.
- **LR grammar** is suitable for **bottom up parsing**.
- It cannot be used for top-down parsing because it is left recursive.

# Non-left-recursive variant

E → E + T | T
T → T * F | F
F → ( E ) | id

⟹

E → T E'
E' → + T E' | ε
T → F T'
T' → * F T' | ε
F → ( E ) | id

# Common Programming Errors

- **Lexical Errors**
  - Misspellings of identifiers, keywords, or operators
- **Syntactic Errors**
  - Misplaced semicolons
  - Extra or missing braces
- **Semantic Errors**
  - Type mismatch between operators and operands
- **Logical Errors**
  - Incorrect reasoning like use of assignment operator = instead of the comparison operator ==

# Syntax Error Handling

- The precision of parsing methods allows syntactic errors to be detected very efficiently.

- Several parsing methods, such as the LL and LR methods, detect an error as soon as possible; that is, when the stream of tokens from the lexical analyzer cannot be parsed further according to the grammar for the language.

- More precisely, they have the **viable-prefix property**, meaning that they detect that an error has occurred as soon as they see a prefix of the input that cannot be completed to form a string in the language.

# Syntax Error Handling

- Another reason for emphasizing error recovery during parsing is that many errors appear syntactic, whatever their cause, and are exposed when parsing cannot continue.

- A few semantic errors, such as type mismatches, can also be detected efficiently; however, accurate detection of semantic and logical errors at compile time is in general a difficult task.

# Goals of error handler in a parser

- ✓ Report the presence of errors clearly and accurately.
- ✓ Recover from each error quickly enough to detect subsequent errors.
- ✓ Add minimal overhead to the processing of correct programs.

- It must report the place in the source program where an error is detected, because there is a good chance that the actual error occurred within the previous few tokens.
- A common strategy is **to print the offending line with a pointer to the position** at which an error is detected.

# Error Recovery Strategies

1. Panic Mode Recovery

2. Phrase-Level Recovery

3. Error Productions

4. Global Correction

# Panic-Mode Recovery

- On discovering an error, the parser discards input symbols one at a time until one of a designated set of synchronizing tokens is found.

- The synchronizing tokens are usually delimiters, such as semicolon or }, whose role in the source program is clear and unambiguous.

- The compiler designer must select the synchronizing tokens appropriate for the source language.

- While panic-mode correction often skips a considerable amount of input without checking it for additional errors, it has the advantage of simplicity, and is guaranteed not to go into an infinite loop.

# Phrase-Level Recovery

- On discovering an error, a parser may perform local correction on the remaining input; that is, it may replace a prefix of the remaining input by some string that allows the parser to continue.

- A typical local correction is to replace a comma by a semicolon, delete an extraneous semicolon, or insert a missing semicolon.

- Phrase-level replacement has been used in several error-repairing compilers, as it can correct any input string.

- Its major drawback is the difficulty it has in coping with situations in which the actual error has occurred before the point of detection.

# Error Production

- By anticipating common errors that might be encountered, we can augment the grammar for the language at hand with productions that generate the erroneous constructs.

- A parser constructed from a grammar augmented by these error productions detects the anticipated errors when an error production is used during parsing.

- The parser can then generate appropriate error diagnostics about the erroneous construct that has been recognized in the input.

# Global Correction

- Ideally, we would like a compiler to make as few changes as possible in processing an incorrect input string.

- There are algorithms for choosing a minimal sequence of changes to obtain a globally least-cost correction.

- Given an incorrect input string x and grammar G, these algorithms will find a parse tree for a related string y, such that the number of insertions, deletions, and changes of tokens required to transform x into y is as small as possible.

- Unfortunately, these methods are in general too costly to implement in terms of time and space, so these techniques are currently only of theoretical interest.

- Do note that a closest correct program may not be what the programmer had in mind.

# Context-Free Grammars

# Chomsky hierarchy

# Chomsky hierarchy



Type-0 → Unristricted Grammar (Recognized by Turing Machine)

Type-1 → Context Sesitive Grammar (Accepted by Linear Bound Automata)

Type-2 → Context Free Grammar (Accepted by Push Down Automata)

Type-3 → Regular Grammar (Accepted By Finite Automata)

https://www.geeksforgeeks.org/chomsky-hierarchy-in-theory-of-computation/

# The formal definition of Context-free grammars

A context-free grammar (grammar for short) consists

1. The **Terminals** are the basic symbols from which strings are formed.
2. The **Nonterminals** are syntactic variables that denote sets of strings.
   - The sets of strings denoted by nonterminals help define the language generated by the grammar.
   - Nonterminals impose a hierarchical structure on the language that is key to syntax analysis and translation.
3. In a grammar, one nonterminal is distinguished as the **start symbol**, and the set of strings it denotes is the language generated by the grammar.
   - Conventionally, the productions for the start symbol are listed first.
4. The **productions of a grammar** specify the manner in which the terminals and nonterminals can be combined to form strings.

# The formal definition of Context-free grammars

Each production consists of:

1. A nonterminal called **the head or left side of the production**; this production defines some of the strings denoted by the head.

2. The **symbol →.**

   - Sometimes ::= has been used in place of the arrow.

3. A **body or right side** consisting of zero or more terminals and nonterminals.

   - The components of the body describe one way in which strings of the nonterminal at the head can be constructed.

# Grammar for simple arithmetic expressions

expression → expression + term

expression → expression - term

expression → term

term → term * factor

term → term / factor

term → factor

factor → ( expression )

factor → id

# Conventions

- These symbols are **terminals:**
  - Lowercase letters early in the alphabet, such as a, b, c.
  - Operator symbols such as +, *, and so on.
  - Punctuation symbols such as parentheses, comma, and so on.
  - The digits 0; 1,…, 9.
  - Boldface strings such as **id** or **if**, each of which represents a single terminal symbol.
- These symbols are **nonterminals:**
  - Uppercase letters early in the alphabet, such as A, B, C.
  - The letter S, which, when it appears, is usually the start symbol.
  - Lowercase, italic names such as *expr* or *stmt*.
  - When discussing programming constructs, uppercase letters may be used to represent nonterminals for the constructs.

# Conventions

- Uppercase letters late in the alphabet, such as X, Y, Z, represent grammar symbols; that is, either nonterminals or terminals.

- Lowercase letters late in the alphabet, chiefly u, v, …,z, represent (possibly empty) strings of terminals.

- Lowercase Greek letters α, β, γ, for example represent (possibly empty) string s of grammar symbols.

- Unless stated otherwise, **the head of the first production is the start symbol.**

# Some context-free grammars

- E → E + E

  E → E * E

  E → (E)

  E → id

- S → aSb

  S → ε

- S → 0S0

  S→ 1S1

  S → 0

  S → 1

  S → ε

- S → aB|bA

  A → a | aS | bAA

  B → b |bS | aBB

# Sentinel Form (Left most derivation)

- E → E + E

  E → E * E

  E → (E)

  E → id

- Here, we replace the left most non-terminal of production by appropriate grammar rule. This is called **left most derivation**.

For input: a + b * c

$$E \xrightarrow{lm} E + E$$

$$\xrightarrow{lm} id + E$$

$$\xrightarrow{lm} id + E * E$$

$$\xrightarrow{lm} id + id * E$$

$$\xrightarrow{lm} id + id * id$$

# Canonical derivation (Right most derivation)

- E → E + E

  E → E * E

  E → (E)

  E → id

- Here, we replace the right most non-terminal of production by appropriate grammar rule. This is called **right most derivation**.

For input: a + b * c

$E \xrightarrow{rm} E + E$

$\xrightarrow{rm} E + E * E$

$\xrightarrow{rm} E + E * id$

$\xrightarrow{rm} E + id * id$

$\xrightarrow{rm} id + id * id$

# Ambiguity

- A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

- Put another way, an ambiguous grammar is one that produces more than one leftmost derivation or more than one rightmost derivation for the same sentence.

- Consider the grammar: $E \rightarrow E + E \mid E * E \mid ( E ) \mid$ **id**

- Input **id + id \*id** generates two distinct leftmost derivatives:

$$
\begin{aligned}
E &\Rightarrow E + E & E &\Rightarrow E * E \\
  &\Rightarrow \mathbf{id} + E & &\Rightarrow E + E * E \\
  &\Rightarrow \mathbf{id} + E * E & &\Rightarrow \mathbf{id} + E * E \\
  &\Rightarrow \mathbf{id} + \mathbf{id} * E & &\Rightarrow \mathbf{id} + \mathbf{id} * E \\
  &\Rightarrow \mathbf{id} + \mathbf{id} * \mathbf{id} & &\Rightarrow \mathbf{id} + \mathbf{id} * \mathbf{id}
\end{aligned}
$$

$$E \to E + E \mid E * E \mid ( E ) \mid \mathbf{id}$$
input: id + id *id

$$
\begin{aligned}
E &\Rightarrow E + E \\
&\Rightarrow \mathbf{id} + E \\
&\Rightarrow \mathbf{id} + E * E \\
&\Rightarrow \mathbf{id} + \mathbf{id} * E \\
&\Rightarrow \mathbf{id} + \mathbf{id} * \mathbf{id}
\end{aligned}
\qquad
\begin{aligned}
E &\Rightarrow E * E \\
&\Rightarrow E + E * E \\
&\Rightarrow \mathbf{id} + E * E \\
&\Rightarrow \mathbf{id} + \mathbf{id} * E \\
&\Rightarrow \mathbf{id} + \mathbf{id} * \mathbf{id}
\end{aligned}
$$

E → E + E | E * E | ( E ) | **id**
input: id + id *id

$$E \Rightarrow E + E \qquad E \Rightarrow E * E$$
$$\Rightarrow \textbf{id} + E \qquad\quad \Rightarrow E + E * E$$
$$\Rightarrow \textbf{id} + E * E \qquad \Rightarrow \textbf{id} + E * E$$
$$\Rightarrow \textbf{id} + \textbf{id} * E \qquad \Rightarrow \textbf{id} + \textbf{id} * E$$
$$\Rightarrow \textbf{id} + \textbf{id} * \textbf{id} \qquad \Rightarrow \textbf{id} + \textbf{id} * \textbf{id}$$

# Why use regular expressions to define the lexical syntax of a language?

1. Separating the syntactic structure of a language into lexical and nonlexical parts provides a convenient way of modularizing the front end of a compiler into two manageable-sized components.

2. The lexical rules of a language are frequently quite simple, and to describe them we do not need a notation as powerful as grammars.

3. Regular expressions generally provide a more concise and easier-to-understand notation for tokens than grammars.

4. More efficient lexical analyzers can be constructed automatically from regular expressions than from arbitrary grammars.

# Observe the grammar

*stmt*       →      **if** *expr* **then** *stmt*

           |      **if** *expr* **then** *stmt* **else** *stmt*

           |      other

Here "other" stands for any other statement.

$$stmt \rightarrow \quad \textbf{if } expr \textbf{ then } stmt$$
$$| \quad \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt$$
$$| \quad \text{other}$$

if $E_1$ then $S_1$ else if $E_2$ then $S_2$ else $S_3$

$$stmt \rightarrow \quad \textbf{if } expr \textbf{ then } stmt$$
$$| \quad \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt$$
$$| \quad \text{other}$$

**if** $E_1$ **then if** $E_2$ **then** $S_1$ **else** $S_2$

# Re-writing dangling else grammar

| *stmt* | → | **if** *expr* **then** *stmt* |
|---|---|---|
| | \| | **if** *expr* **then** *stmt* **else** *stmt* |
| | \| | other |

$$stmt \rightarrow matched\_stmt$$
$$\mid open\_stmt$$
$$matched\_stmt \rightarrow \textbf{if } expr \textbf{ then } matched\_stmt \textbf{ else } matched\_stmt$$
$$\mid \textbf{other}$$
$$open\_stmt \rightarrow \textbf{if } expr \textbf{ then } stmt$$
$$\mid \textbf{if } expr \textbf{ then } matched\_stmt \textbf{ else } open\_stmt$$

- Note:
  - No general techniques for handling ambiguity
  - Impossible to convert automatically an ambiguous grammar to an unambiguous one

# Removing Useless Variables

S → abS | abA | abB

A → cd

B → aB

C → dc

# Removing Useless Variables

S → abS | abA | abB

A → cd

B → aB

C → dc

⟹

S → abS | abA | abB

A → cd

B → aB

C is Non-reachable so remove it.

# Removing Useless Variables

S → abS | abA | abB
A → cd
B → aB
C → dc

➡️

S → abS | abA | abB
A → cd
B → aB

⬇️ B will never lead to a terminal.

S → abS | abA
A → cd

# Removing Useless Variables

S → abS | abA | abB
A → cd
B → aB
C → dc

S → abS | abA | abB
A → cd
B → aB

B will never
lead to a terminal.

S → abS | abA
A → cd

# Remove useless variables

1. S → AB/a

   A → BC/b

   B → aB/C

   C → aC/B


2. S → AB/AC

   A → aAb/bAa/a

   B → bbA/aaB/AB

   C → abCA/aDb

   D → bD/aC


3. S → aB / bX

   A → Bad / bSX / a

   B → aSB / bBX

   X → SBD / aBx / ad


4. S → AB|a

   A → b


5. S → AB|CA

   B → BC | AB

   A → a

   C → AB|b

# Remove useless variables

1. S → AB/a
   A → BC/b
   B → aB/C
   C → aC/B

→

S → AB/a
A → BC/b
B → aB/C
C → aC/B

→

S → AB/a
A → BC/b
B → aB/C
C → aC/B

S → a
A → b

S → a

# Remove useless variables

2. S → AB/AC
   A → aAb/bAa/a
   B → bbA/aaB/AB
   C → abCA/aDb
   D → bD/aC

S → AB/AC
A → aAb/bAa/a
B → bbA/aaB/AB
C → abCA/aDb
D → bD/aC

S → AB/AC
A → aAb/bAa/a
B → bbA/aaB/AB
C → abCA/aDb
D → bD/aC

S → AB
A → aAb/bAa/a
B → bbA/aaB/AB

# Remove useless variables

3. S → aB / bX
   A → Bad / bSX / a
   B → aSB / bBX
   X → SBD / aBx / ad

S → aB / bX
A → Bad / bSX / a
B → aSB / bBX
X → SBD / aBx / ad

S → aB / bX
A → Bad / bSX / a
B → aSB / bBX
X → SBD / aBx / ad

S → bX
A → bS X / a
X → ad

S → bX
X → ad

# Remove useless variables

4. S → AB|a
   A → b

S → AB|a
A → b

S → AB|a

A → b

S → a

# Remove useless variables

5. S → AB|CA
   B → BC | AB
   A → a
   C → AB|b

S → AB|CA
B → BC | AB
A → a
C → AB|b

S → AB|CA
B → BC | AB
A → a
C → AB|b

S → CA
A → a
C → b

# Elimination of Left Recursion

- A Grammar G (V, T, P, S) is left recursive if it has a production in the form.

  A → A α |β

- The above Grammar is left recursive because the left of production is occurring at a first position on the right side of production.

- It can eliminate left recursion by replacing a pair of production with

  A → βA'

  A' → αA'|ϵ

# Elimination of Left Recursion

- E → E + T|T

  T → T * F|F

  F → (E)|id

↓

E → TE'

E' → +TE'|ε

T → FT'

T'→ *FT'|ε

F → (E)|id

- Comparing E → E +T |T with A → A α |β.

- Here, A is E, α is +T and β is T

- On eliminating left recursion, using

  A → βA'

  A' → αA'|ε

E → TE'

E' → +TE'|ε

- Similarly for T → T * F|F

T → FT'

T' → *FT'|ε

# Removal of left recursion

- S → a|^|(T)

  T → T, S|S

# Removal of left recursion

- S → a|^|(T)

  T → T, S|S



A → A α |β
After removal,
A → βA'
A' → αA'|ϵ

- Comparing T → T, S|S with

  A → A α |β.

- Here, A is T, α is ,S and β is S

- On eliminating left recursion,

  T → ST'

  T'→ ,ST' |ϵ

- So, finally

  **S → a|^|(T)**

  **T → ST'**

  **T'→ ,ST' |ϵ**

# Remove left recursion

- S → Aa | b
  A→ Ac | Sd | ϵ

# Remove left recursion

- S → Aa | b

  A→ Ac | Sd | ϵ

  A → A α |β

  After removal,

  A → βA'

  A' → αA'|ϵ

- Eliminating the indirect left recursion.

  S → Aa|b

  A→ Ac | Aad | bd|ϵ

- In A→ Ac | Aad | bd|ϵ,

  A is A

  α is c, ad

  β is bd|ϵ

- So,

  A → (bd | ϵ )A' → A' | bdA'

  A' → cA' |adA'|ϵ

# Remove left recursion

- S → Aa | b

  A→ Ac | Sd | ϵ

- Finally,

  S → Aa | b

  A → A' | bdA'

  A' → cA' |adA'|ϵ

- Eliminating the indirect left recursion.

  S → Aa|b

  A→ Ac | Aad | bd|ϵ

- In A→ Ac | Aad | bd|ϵ,

  A is A

  α is c, ad

  β is bd|ϵ

- So,

  A → (bd | ϵ )A' → A' | bdA'

  A' → cA' |adA'|ϵ

# Remove left recursion (try yourself)

- S → Sa |Sb | c | d

- A → Br

  B → Cd

  C → At

# Elimination of Left factoring

- Left factoring is removing the common left factor that appears in two productions of the same non-terminal.

- It is done to avoid back-tracing by the parser.

  - A → aα1 | aα2 | aα3

    Here, a is a common prefix or factor.

  - After removal,

    A → aA'

    A' → α1 | α2 |α3

# Remove left factoring

- A → aAB / aBc / aAc

- Here, common prefix a is observed.

  A → aA'

  A' → AB | Bc | Ac

- Again common prefix A is observed in A' → AB | Bc | Ac

  A' → AD|Bc

  D → B | c

A → aα1 | aα2 | aα3

After removal,

A → aA'

A' → α1 | α2 |α3

# Remove left factoring

- A → aAB / aBc / aAc

- Finally we have,

  A → aA'

  A'→ AD|Bc

  D → B | c

- Here, common prefix a is observed.

  A → aA'

  A' → AB | Bc | Ac

- Again common prefix A is observed in A' → AB | Bc | Ac

  A'→ AD|Bc

  D → B | c

# Remove left factoring

- S → bSSaaS/bSSaSb/bSb/a

# Remove left factoring

- S → bSSaaS/bSSaSb/bSb/a

- Here common prefix bS is observed.

  S → bSS' | a

  S' → SaaS | SaSb |b

- Again, Sa common prefix is observed

  S' → SaA | b

  A → aS | Sb

A → aα1 | aα2 | aα3

After removal,

A → aA'

A' → α1 | α2 |α3

# Remove left factoring

- S → bSSaaS/bSSaSb/bSb/a

- Finally, we have

  S → bSS' | a

  S' → SaA | b

  A → aS | Sb

- Here common prefix bS is observed.

  S → bSS' | a

  S' → SaaS | SaSb |b

- Again, Sa common prefix is observed

  S' → SaA | b

  A → aS | Sb

# Remove left factoring

- S → iEtS | iEtSeS | a

  E → b

# Remove left factoring

- S → iEtS | iEtSeS | a

  E → b

- Here, common prefix iEtS is observed

  S → iEtSS' | a

  S' → eS | ε

A → aα1 | aα2 | aα3

After removal,

A → aA'

A' → α1 | α2 |α3

- So, finally we have,

  S → iEtSS' | a

  S' → eS | ε

  E → b

# Remove left factoring

- S → aSSbS | aSaSb | abb |b

# Remove left factoring

- S → aSSbS | aSaSb | abb |b

- Common prefix a is observed.

  S → aS' | b

  S' → SSbS | SaSb | bb

- Again common prefix S is observed.

  S' → SA | bb

  A → SbS | aSb

A → aα1 | aα2 | aα3

After removal,

A → aA'

A' → α1 | α2 |α3

# Remove left factoring

- S → aSSbS | aSaSb | abb | b

- Finally we have,

  S → aS' | b

  S' → SA | bb

  A → SbS | aSb

- Common prefix a is observed.

  S → aS' | b

  S' → SSbS | SaSb | bb

- Again common prefix S is observed.

  S' → SA | bb

  A → SbS | aSb

# Remove left factoring

- S → a | ab | abc | abcd

# Remove left factoring

- S → a | ab | abc | abcd

- Common prefix a is observed

  S → aS'

  S' → ε | b | bc |bcd

A → aα1 | aα2 | aα3

After removal,

A → aA'

A' → α1 | α2 |α3

- Common prefix b is observed

  S' → ε | bA

  A → ε | c |cd

- Common prefix c is observed

  A → ε | cB

  B → ε |d

# Remove left factoring

- S → a | ab | abc | abcd

- Finally we have
  S → aS'
  S' → bA |ε
  A → cB|ε
  B → d| ε

- Common prefix a is observed
  S → aS'
  S' → ε | b | bc |bcd

- Common prefix b is observed
  S' → ε | bA
  A → ε | c |cd

- Common prefix c is observed
  A → ε | cB
  B → ε |d

# Remove left factoring

- S → aAd | aB

  A → a | ab

  B → ccd | ddc

# Remove left factoring

- S → aAd | aB
  A → a | ab
  B → ccd | ddc

- S → aS'
  S' → Ad | B
  A → aA'
  A' → b | ∈
  B → ccd | ddc

# Elimination of ε productions

- Steps:
  - To remove A → ε, look for all productions whose right side contains A
  - Replace each occurrence of 'A' in each of these productions with ε
  - Add the resultant productions to the grammar

# Elimination of ε productions

- S → ABA
  A → aA | ε
  B → bB | ε

# Elimination of ε productions

- S → ABA

  A → aA | ε

  B → bB | ε

- As A and B are directly nullable variables:

  A → aA | a

  B → bB | b

  S → ABA |AB | BA |AA |A |B

# Elimination of ε productions

- S → aS | AB | a

  A → ε

  B → ε

  D → b

# Elimination of ε productions

- S → aS | AB |a
  A → ε
  B → ε
  D → b

- As A and B are directly nullable variables:

  S → aS | a
  D → b

  S → aS | a

# Remove null productions

- S → ABAC
  A → aA / ϵ
  B → bB / ϵ
  C → c

# Remove null productions

- S → ABAC
  A → aA / ϵ
  B → bB / ϵ
  C → c

- S → ABAC / ABC / BAC / BC / AAC / AC / C
  A → aA / a
  B → bB / b
  C → c

# Remove null productions

- S → ABCd
  A → BC
  B → bB | ϵ
  C → cC | ϵ

# Remove null productions

- S → ABCd

  A → BC

  B → bB | ϵ

  C → cC | ϵ

- S → ABCd | ABd | ACd | BCd | Ad | Bd |Cd | d

  A → BC | B | C

  B → bB | b

  C → cC | c

# Remove null productions

- S→a|Ab|aBa

  A→b|ϵ

  B→b|A

# Remove null productions

- S→a|Ab|aBa

  A→b|ϵ

  B→b|A



- S→a|Ab|b|aBa|aa

  A → b
  B → b

# Remove null productions & unit productions

- S → a | Ab | aBa

  A → b | ϵ

  B → b | A

⬇

- S → a | Ab | b | aBa | aa   ➡   S → a | bb | b | aba | aa

  A → b

  B → b

# Remove unit productions

- Steps:

  1. To remove X → Y, add production X→a to the grammar rule whenever Y→a occurs in the grammar

  2. Now delete X →Y from the grammar

  3. Repeat Step 1 and 2 until all unit productions are removed

# Remove unit productions

- S → 0A | 1B | C
  A → 0S |00
  B → 1 | A
  C → 01

# Remove unit productions

- S → 0A | 1B | C

  A → 0S | 00

  B → 1 | A

  C → 01

- S → C is a unit production

- S → 0A | 1B | 01

  A → 0S | 00

  B → 1 | A

  C → 01

# Remove unit productions

- S → 0A | 1B | C

  A → 0S |00

  B → 1 | A

  C → 01

- S → C is a unit production
- S → 0A | 1B | 01

  A → 0S |00

  B → 1 | A

  C → 01

- B → A is also a unit production
- S → 0A |1B |01
- A → 0S | 00
- B → 1 |0S | 00
- C → 01

# Remove unit productions

- S → 0A | 1B | C

  A → 0S | 00

  B → 1 | A

  C → 01

  ⬇

- S → 0A | 1B | 01

  A → 0S | 00

  B → 1 | 0S | 00

  C → 01

- S → C is a unit production

- S → 0A | 1B | 01

  A → 0S | 00

  B → 1 | A

  C → 01

- B → A is also a unit production

- S → 0A | 1B | 01

  A → 0S | 00

  B → 1 | 0S | 00

  C → 01

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

- First writing the productions with unit productions
- S→B

  B→A

  A→B

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

- First writing the productions with unit productions
- For the production S→B

  S→B →bb

  S →B→A→a

  S →B→A→bc

  B→A

  A→B

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

- First writing the productions with unit productions
- For the production S→B

  S→B →bb

  S →B→A→a

  S →B→A→bc

- For the production B→A

  B→A→a

  B→A→bc

  A→B

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

- First writing the productions with unit productions
- For the production S→B

  S→B →bb

  S →B→A→a

  S →B→A→bc
- For the production B→A

  B→A→a

  B→A→bc
- For the production A→B

  A→B→bb

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B

- First writing the productions with unit productions
- For the production S→B

  S→B →bb

  S →B→A→a

  S →B→A→bc

- For the production B→A

  B→A→a

  B→A→bc

- For the production A→B

  A→B→bb

# Remove unit productions

- S → Aa/B/c

  B → A/bb

  A → a/bc/B



- S →Aa | c |bb | a | bc

  B →bb | a | bc

  A →a | bc | bb

- First writing the productions with unit productions
- For the production S→B

  S→B →bb

  S →B→A→a

  S →B→A→bc
- For the production B→A

  B→A→a

  B→A→bc
- For the production A→B

  A→B→bb

# Sequence of steps for elimination

1. Eliminate/Remove null production

2. Eliminate/Remove unit production

3. Eliminate/Remove useless symbol

# Simplify the grammar

- S → Aa | B
  B → a | bC
  C → a | ε

# Simplify the grammar
# Eliminate/Remove null production

- S → Aa | B

  B → a | bC

  C → a | ε

- C → ε is a null production
- To remove it, add the production
- B → bC → b ε → b


- So we have,
- S → Aa | B

  B → a | bC | b

  C → a
- Now, no null productions.

# Simplify the grammar
## Eliminate/Remove null production

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇩

- S → Aa | B

  B → a | bC |b

  C → a

- C → ε is a null production
- To remove it, add the production
- B →bC →b ε → b

- So we have,
- S → Aa | B

  B → a | bC |b

  C → a

- Now, no null productions.

# Simplify the grammar
# Eliminate/Remove unit production

- S → Aa | B
  B → a | bC
  C → a | ε

  ⇩

- S → Aa | B
  B → a | bC |b
  C → a

- Identifying unit productions
  S→B

- Removing it, gives:
  S→B → a | bC |b

- So grammar is:
  S → Aa |a |bC |b
  B → a |bC |b
  C →a

# Simplify the grammar
# Eliminate/Remove unit production

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇩

- S → Aa | B

  B → a | bC |b

  C → a

  ⇩

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- Identifying unit productions

  S→B

- Removing it, gives:

  S→B → a | bC |b

- So grammar is:

  S → Aa |a |bC |b

  B → a |bC |b

  C →a

- Now, no unit productions.

# Simplify the grammar
# Eliminate/Remove unit production

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇓

- S → Aa | B

  B → a | bC |b

  C → a

  ⇓

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- Identifying unit productions

  S→B

- Removing it, gives:

  S→B → a | bC |b

- So grammar is:

  S → Aa |a |bC |b

  B → a |bC |b

  C →a

- Now, no unit productions.

# Simplify the grammar
## Eliminate/Remove useless variables

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇩

- S → Aa | B

  B → a | bC |b

  C → a

  ⇩

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- B is a useless variable as cannot reach B from S.
- So remove it.

  S → Aa |a |bC |b

  C →a

# Simplify the grammar
# Eliminate/Remove useless variables

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇩

- S → Aa | B

  B → a | bC |b

  C → a

  ⇩

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- B is a useless variable as cannot reach B from S.
- So remove it.

  S → Aa |a |bC |b

  C →a


- Now, Aa is useless as no production head A.
- So remove it.

  S → a |bC |b

  C →a

# Simplify the grammar
# Eliminate/Remove useless variables

- S → Aa | B

  B → a | bC

  C → a | ε

  ⇓

- S → Aa | B

  B → a | bC |b

  C → a

  ⇓

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- B is a useless variable as cannot reach B from S.
- So remove it.

  S → Aa |a |bC |b

  C →a

- Now, Aa is useless as no production head A.
- So remove it.

  S → a |bC |b

  C →a

- Now, C → a can be substituted

  S → a | ba |b

# Simplify the grammar
## Eliminate/Remove useless variables

- S → Aa | B

  B → a | bC

  C → a | ε

  ⬇

- S → Aa | B

  B → a | bC |b

  C → a

  ⬇

- S → Aa |a |bC |b

  B → a |bC |b

  C →a

- B is a useless variable as cannot reach B from S.

- So remove it.

  S → Aa |a |bC |b

  C →a

- Now, Aa is useless as no production head A.

- So remove it.

  S → a |bC |b

  C →a

- Now, C → a can be substituted

  S → a | ba |b

# Simplify the grammar
# Eliminate/Remove useless variables

- S → Aa | B
  B → a | bC
  C → a | ε

  ⬇

- S → Aa | B
  B → a | bC |b
  C → a

  ⬇

- S → Aa |a |bC |b
  B → a |bC |b
  C →a

  ➡  S → a | ba |b