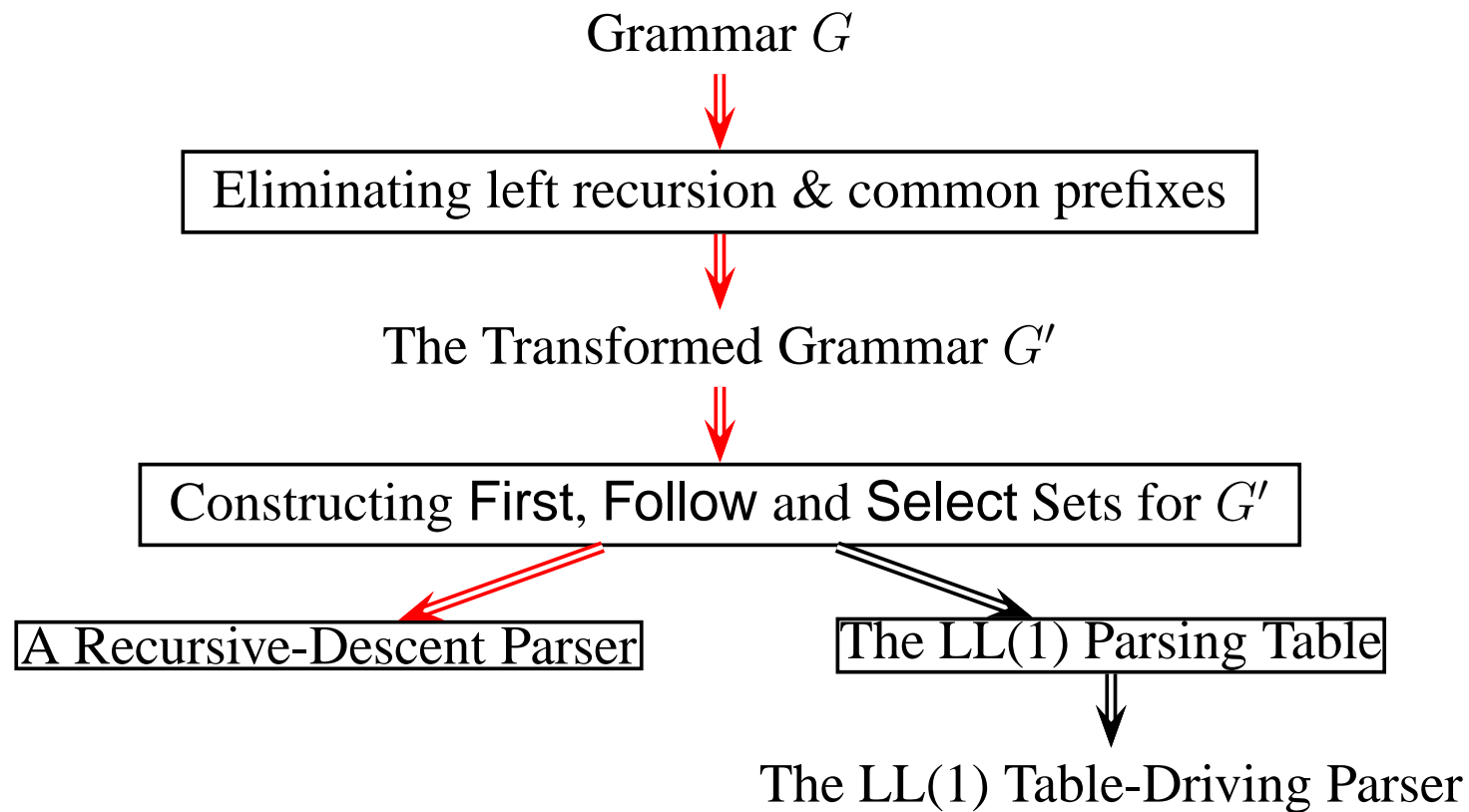# Lecture 4: Top-Down Parsing: Recursive-Descent

1. Compare and contrast top-down and bottom-up parsing
2. Write a predictive (or non-backtracking) top-down parser

Grammar $G$

⇓

Eliminating left recursion & common prefixes

⇓

The Transformed Grammar $G'$

⇓

Constructing **First**, **Follow** and **Select** Sets for $G'$

A Recursive-Descent Parser

The LL(1) Parsing Table

⇓

The LL(1) Table-Driving Parser

# The micro-English Grammar Revisited

| 1 | $\langle$sentence$\rangle$ | $\rightarrow$ | $\langle$subject$\rangle$ $\langle$predicate$\rangle$ |
| 2 | $\langle$subject$\rangle$ | $\rightarrow$ | **NOUN** |
| 3 | | | | **ARTICLE NOUN** |
| 4 | $\langle$predicate$\rangle$ | $\rightarrow$ | **VERB** $\langle$object$\rangle$ |
| 5 | $\langle$object$\rangle$ | $\rightarrow$ | **NOUN** |
| 6 | | | | **ARTICLE NOUN** |

## The English Sentence

**PETER PASSED THE TEST**

# The micro-English Grammar Revisited (Cont'd)

- **The Leftmost Derivation:**

$$\langle sentence \rangle \quad \Longrightarrow_{lm} \quad \langle subject \rangle \, \langle predicate \rangle \qquad \text{by P1}$$
$$\Longrightarrow_{lm} \quad \textbf{NOUN} \, \langle predicate \rangle \qquad \text{by P2}$$
$$\Longrightarrow_{lm} \quad \textbf{NOUN VERB} \, \langle object \rangle \qquad \text{by P4}$$
$$\Longrightarrow_{lm} \quad \textbf{NOUN VERB ARTICLE NOUN} \qquad \text{by P6}$$

- **The Rightmost Derivation:**

$$\langle sentence \rangle \quad \Longrightarrow_{rm} \quad \langle subject \rangle \, \langle predicate \rangle \qquad \text{by P1}$$
$$\Longrightarrow_{rm} \quad \langle subject \rangle \, \textbf{VERB} \, \langle object \rangle \qquad \text{by P4}$$
$$\Longrightarrow_{rm} \quad \langle subject \rangle \, \textbf{VERB ARTICLE NOUN} \qquad \text{by P6}$$
$$\Longrightarrow_{rm} \quad \textbf{NOUN VERB ARTICLE NOUN} \qquad \text{by P2}$$

# The Role of the Parser

**PETER PASSED THE TEST**

↓

Scanner

↓

$\text{NOUN}_1$ **VERB ARTICLE** $\text{NOUN}_2$

↓

Parser

↓

⟨sentence⟩

⟨subject⟩  ⟨predicate⟩

⟨**NOUN**⟩  **VERB**  ⟨object⟩
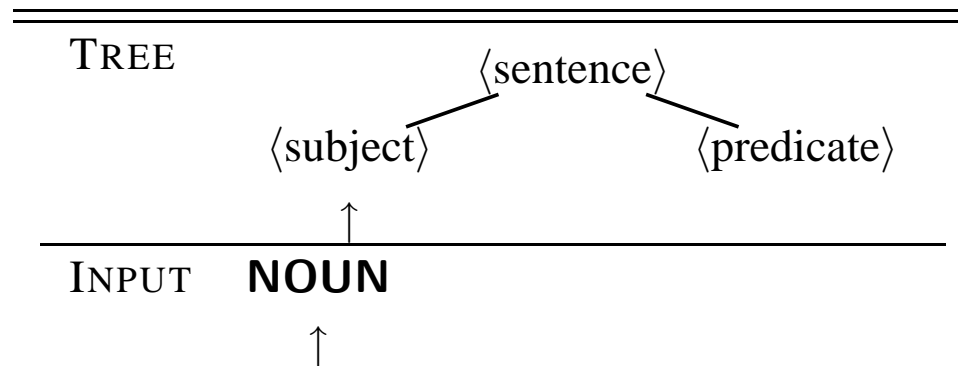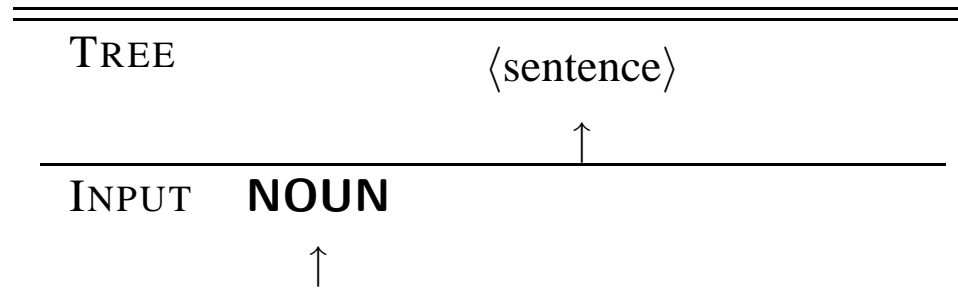
PETER  PASSED  **ARTICLE**  **NOUN**

THE  TEST

# Two General Parsing Methods

1. Top-down parsing – Build the parse tree top-down:

   - Productions used represent the leftmost derivation.

   - The best known and widely used methods:
     - Recursive descent
     - Table-driven
     - LL(k) (Left-to-right scan of input, Leftmost derivation, k tokens of lookahead).
     - Almost all programming languages can be specified by LL(1) grammars, but such grammars may not reflect the structure of a language
     - In practice, LL(k) for small $k$ is used

   - Implemented more easily by hand.

   - Used in parser generators such as JavaCC

2. Bottom-up parsing – Build the parse tree bottom-up:

   - Productions used represent the rightmost derivation in reverse.

   - The best known and widely used method: LR(1) (Left-to- right scan of input, Rightmost derivation in reverse, 1 token of lookahead)

   - More powerful – every LL(1) is LR(1) but the converse is false

   - Used by parser generators (e.g., Yacc and JavaCUP).

# Lookahead Token(s)
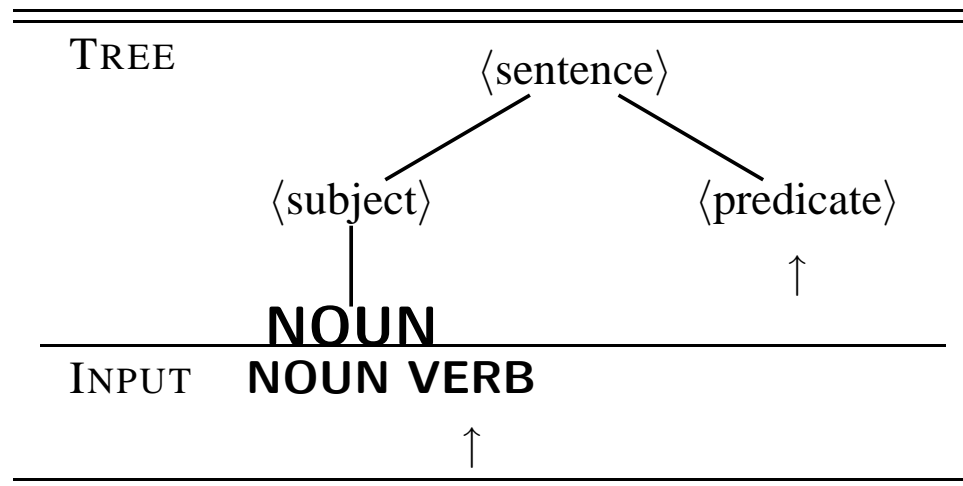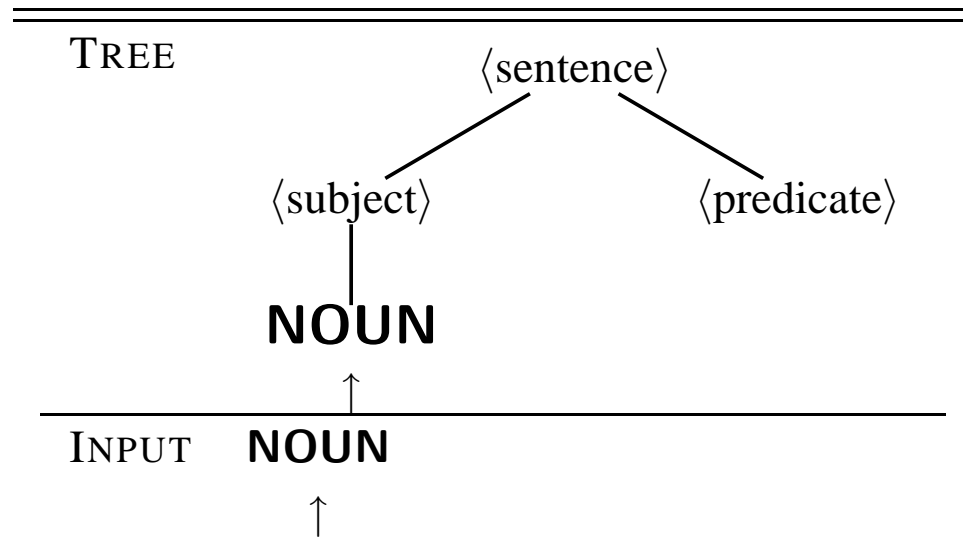
- **Lookahead Token(s)**: The currently scanned token(s) in the input.

- In Recogniser.java, currentToken represents the lookahead token

- For most programming languages, one token lookahead only.

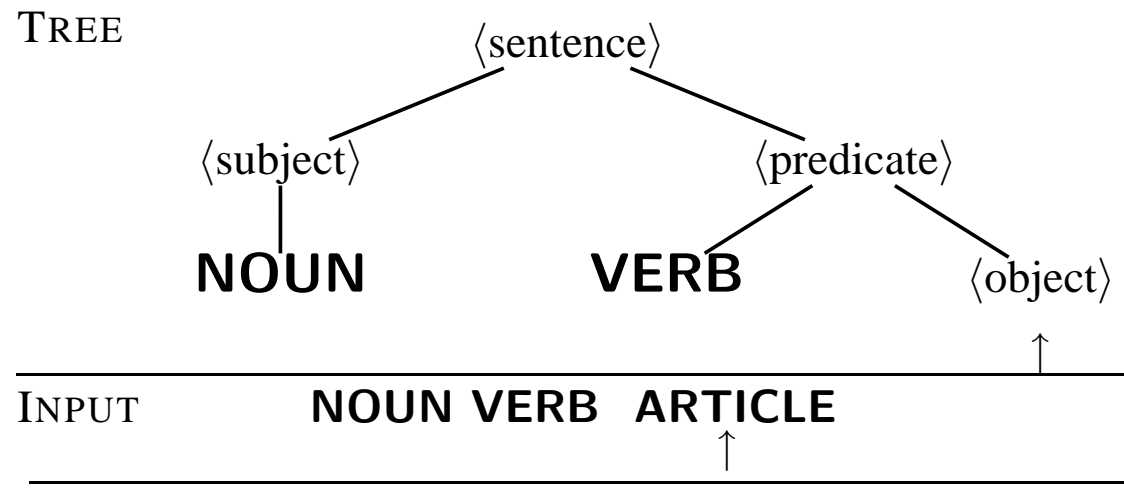- Initially, the lookahead token is the leftmost token in the input.

# Top-Down Parse Of NOUN VERB ARTICLE NOUN

TREE $\langle \text{sentence} \rangle$

$\uparrow$

INPUT **NOUN**

$\uparrow$

---

TREE $\langle \text{sentence} \rangle$

$\langle \text{subject} \rangle$ $\langle \text{predicate} \rangle$

$\uparrow$

INPUT **NOUN**

$\uparrow$

**Notations:**

- $\uparrow$ on the tree indicates the nonterminal being expanded or recognised

- $\uparrow$ on the sentence points to the lookahead token

  - All tokens to the **left** of $\uparrow$ have been read
  - All tokens to the **right** of $\uparrow$ have **NOT** been processed

TREE

⟨sentence⟩

⟨subject⟩                    ⟨predicate⟩

**NOUN**
↑
INPUT    **NOUN**
↑

---

TREE

⟨sentence⟩

⟨subject⟩                    ⟨predicate⟩
↑

**NOUN**
INPUT    **NOUN VERB**
↑

TREE

⟨sentence⟩

⟨subject⟩                    ⟨predicate⟩

**NOUN**          **VERB**          ⟨object⟩
                                ↑
INPUT          **NOUN  VERB**
                        ↑

TREE

⟨sentence⟩

⟨subject⟩                    ⟨predicate⟩

**NOUN**              **VERB**          ⟨object⟩
                                            ↑
INPUT          **NOUN VERB  ARTICLE**
                              ↑

TREE

$\langle$sentence$\rangle$

$\langle$subject$\rangle$          $\langle$predicate$\rangle$

NOUN

VERB          $\langle$object$\rangle$

ARTICLE          NOUN
↑

INPUT          NOUN VERB  ARTICLE
↑

TREE

$\langle$sentence$\rangle$

$\langle$subject$\rangle$          $\langle$predicate$\rangle$

NOUN

VERB          $\langle$object$\rangle$

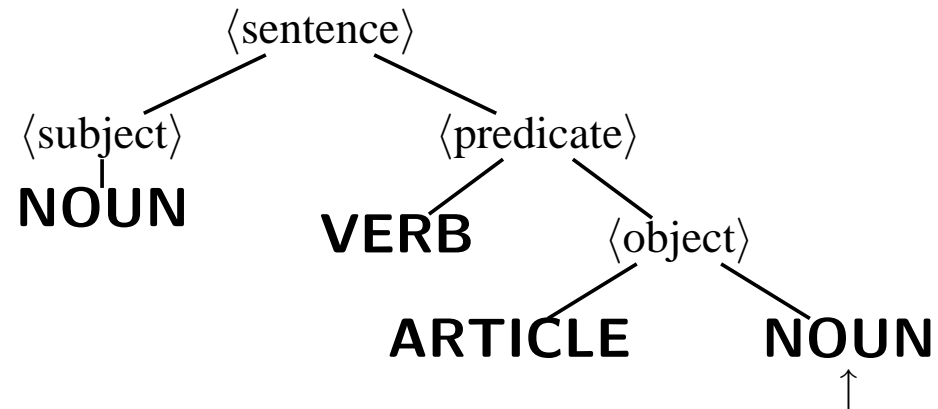ARTICLE          NOUN
↑

INPUT          NOUN VERB ARTICLE  NOUN
↑

# Top-Down Parsing

- Build the parse tree starting with the start symbol (i.e., the root) towards the sentence being analysed (i.e., leaves).

- Use one token of lookahead, in general

- Discover the leftmost derivation

  I.e, the productions used in expanding the parse tree represent a leftmost derivation

# Predictive (Non-Backtracking) Top-Down Parsing

- To expand a nonterminal, the parser always predict (choose) the right alternative for the nonterminal by looking at the lookahead symbol only.

- Flow-of-control constructs, with their distinguishing keywords, are detectable this way, e.g., in the VC grammar:

$$\langle stmt \rangle \ \rightarrow \ \langle compound\text{-}stmt \rangle$$
$$| \quad \textbf{if} \ "(" \ \langle expr \rangle \ ")" \ (\textbf{ELSE} \ \langle stmt \rangle)?$$
$$| \quad \textbf{break} \ ";"$$
$$| \quad \textbf{continue} \ ";"$$
$$\ldots$$

- Prediction happens before the actual match begins.

# Bottom-Up Parse Of NOUN VERB ARTICLE NOUN

| TREE | |
|---|---|
| INPUT | **NOUN**<br>↑ |

| TREE | ⟨subject⟩<br>\|<br>**NOUN** |
|---|---|
| INPUT | **NOUN**<br>↑ |

| TREE | ⟨subject⟩<br>\|<br>**NOUN** |
|---|---|
| INPUT | **NOUN VERB**<br>↑ |

TREE   ⟨subject⟩
           |
       **NOUN**

INPUT   **NOUN VERB ARTICLE**
                          ↑

TREE   ⟨subject⟩
           |
       **NOUN**
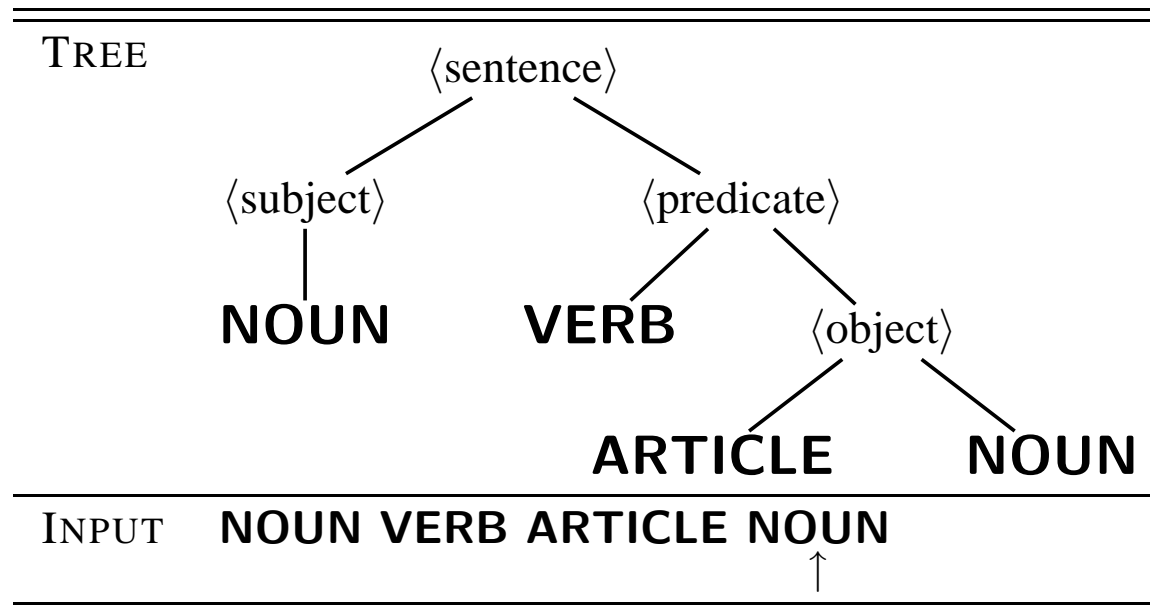
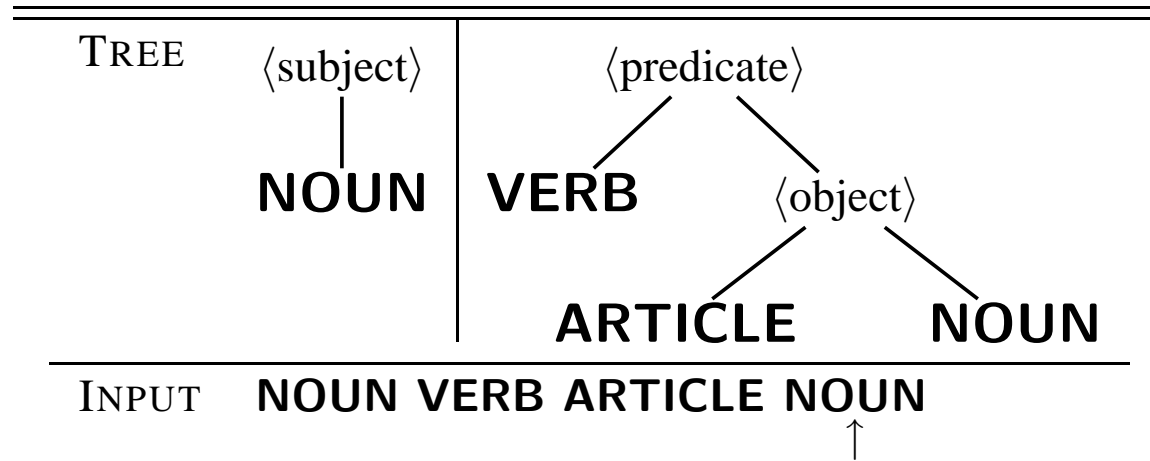INPUT   **NOUN VERB ARTICLE NOUN**
                                ↑

TREE   ⟨subject⟩        ⟨object⟩
           |          ╱        ╲
       **NOUN**  **ARTICLE**        **NOUN**

INPUT   **NOUN VERB ARTICLE NOUN**
                            ↑

**Note:** What if the parser had chosen ⟨subject⟩ →**ARTICLE NOUN** instead of ⟨object⟩ →**ARTICLE NOUN**?

In this case, the parser would not make any further process.

Having read a ⟨subject⟩ and **VERB**, the parser has reached a state in which it should not parse another ⟨subject⟩.

TREE     ⟨subject⟩        ⟨predicate⟩

**NOUN**   **VERB**     ⟨object⟩

**ARTICLE**     **NOUN**

INPUT    **NOUN VERB ARTICLE NOUN**
↑

TREE          ⟨sentence⟩

⟨subject⟩          ⟨predicate⟩

**NOUN**     **VERB**     ⟨object⟩

**ARTICLE**     **NOUN**

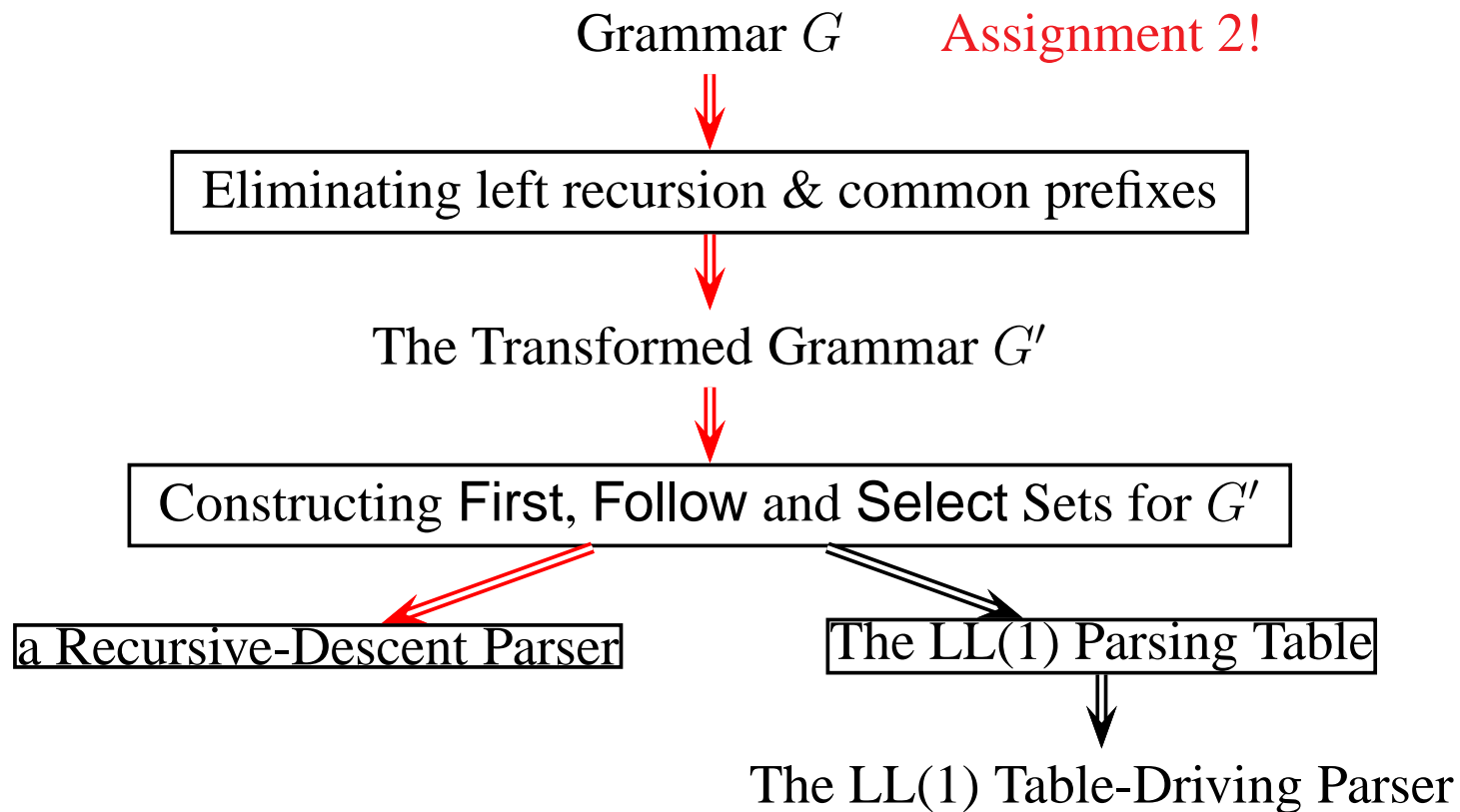INPUT    **NOUN VERB ARTICLE NOUN**
↑

# Bottom-Up Parsing

- Build the parse tree starting with the the sentence being analysed (i.e., leaves) towards the start symbol (i.e., the root).

- Use one token of lookahead, in general.

- The basic (smallest) language constructs recognised first, then they are used to discover more complex constructs.

- Discover the rightmost derivation in reverse — the productions used in expanding the parse tree represent a rightmost derivation in reverse order

## Lecture 4: Top-Down Parsing: Recursive-Descent

1. Compare and contrast top-down and bottom-up parsing $\checkmark$

2. Write a predictive (or non-backtracking) top-down parser

Grammar $G$    Assignment 2!

$\Downarrow$

| Eliminating left recursion & common prefixes |

$\Downarrow$

The Transformed Grammar $G'$

$\Downarrow$

| Constructing **First**, **Follow** and **Select** Sets for $G'$ |

| a Recursive-Descent Parser |        | The LL(1) Parsing Table |

$\Downarrow$

The LL(1) Table-Driving Parser

# Which of the Two Alternatives on $S$ to Choose?

- Grammar:

$$S \to aA \mid bB$$
$$A \to \cdots$$
$$B \to \cdots$$

- Sentence: $a \cdots$

- The leftmost derivation:

$$S \Longrightarrow_{\text{lm}} aA$$
$$\Longrightarrow_{\text{lm}} \cdots$$

Select the first alternative $aA$

# Which of the Two Alternatives on $S$ to Choose?

- Grammar:

$$S \rightarrow Ab \mid Bc$$
$$A \rightarrow Df \mid CA$$
$$B \rightarrow gA \mid e$$
$$C \rightarrow dC \mid c$$
$$D \rightarrow h \mid i$$

- Sentence: $gchfc$

- The leftmost derivation:

$$S \implies_{lm} Bc \implies_{lm} gAc \implies_{lm} gCAc$$
$$\implies_{lm} gcAc \implies_{lm} gcDfc \implies_{lm} gchfc$$
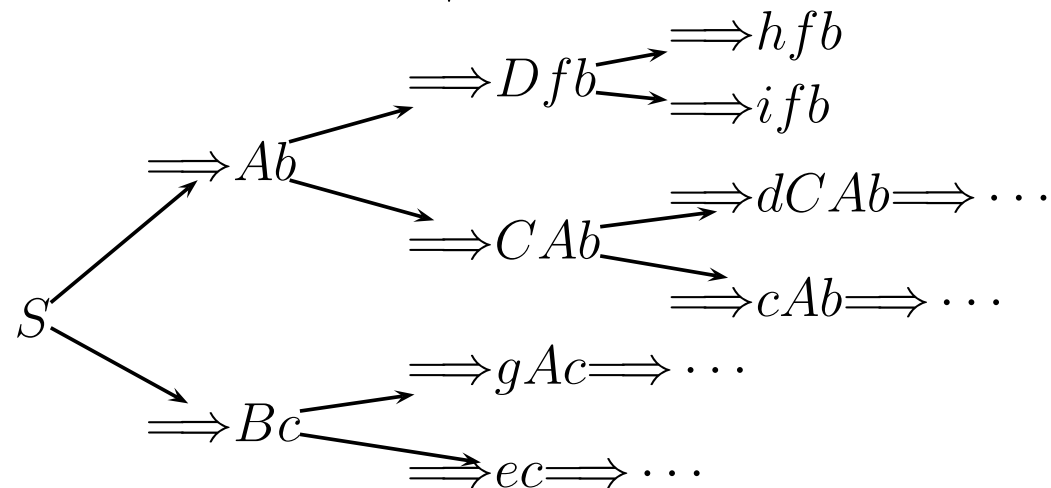
# Intuition behind First Sets

- Grammar:

$$S \rightarrow Ab \mid Bc$$
$$A \rightarrow Df \mid CA$$
$$B \rightarrow gA \mid e$$
$$C \rightarrow dC \mid c$$
$$D \rightarrow h \mid i$$

$$S \Longrightarrow Ab \begin{cases} \Longrightarrow Dfb \begin{cases} \Longrightarrow hfb \\ \Longrightarrow ifb \end{cases} \\ \Longrightarrow CAb \begin{cases} \Longrightarrow dCAb \Longrightarrow \cdots \\ \Longrightarrow cAb \Longrightarrow \cdots \end{cases} \end{cases}$$

$$S \Longrightarrow Bc \begin{cases} \Longrightarrow gAc \Longrightarrow \cdots \\ \Longrightarrow ec \Longrightarrow \cdots \end{cases}$$

- All possible leftmost derivations:

$\mathsf{First}(Ab) = \{c, d, h, i\}$ $\qquad$ $\mathsf{First}(Bc) = \{e, g\}$

# Definition of First Sets

First($\alpha$):

- The set of all terminals that can begin any strings derived from $\alpha$.

- if $\alpha \Longrightarrow^* \epsilon$, then $\epsilon$ is also in First($\alpha$)

## Nullable Nonterminals

A nonterminal $A$ is nullable if $A \Longrightarrow^* \epsilon$.

# A Procedure to Compute $\mathsf{First}(\alpha)$

1. Case 1: $\alpha$ is a single symbol or $\epsilon$:

   If $\alpha$ is a terminal $a$, then $\mathsf{First}(\alpha) = \mathsf{First}(a) = \{a\}$

   else if $\alpha$ is $\epsilon$, then $\mathsf{First}(\alpha) = \mathsf{First}(\epsilon) = \{\epsilon\}$

   else if $\alpha$ is a nonterminal and $\alpha \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \mid \cdots$ then
   $$\mathsf{First}(\alpha) = \cup_k \mathsf{First}(\beta_k)$$

2. Case 2: $\alpha = X_1 X_2 \cdots X_n$:

   If $X_1 X_2 \ldots X_i$ is nullable but $X_{i+1}$ is not, then
   $$\mathsf{First}(\alpha) = \mathsf{First}(X_1) \cup \mathsf{First}(X_2) \cup \cdots \cup \mathsf{First}(X_{i+1})$$

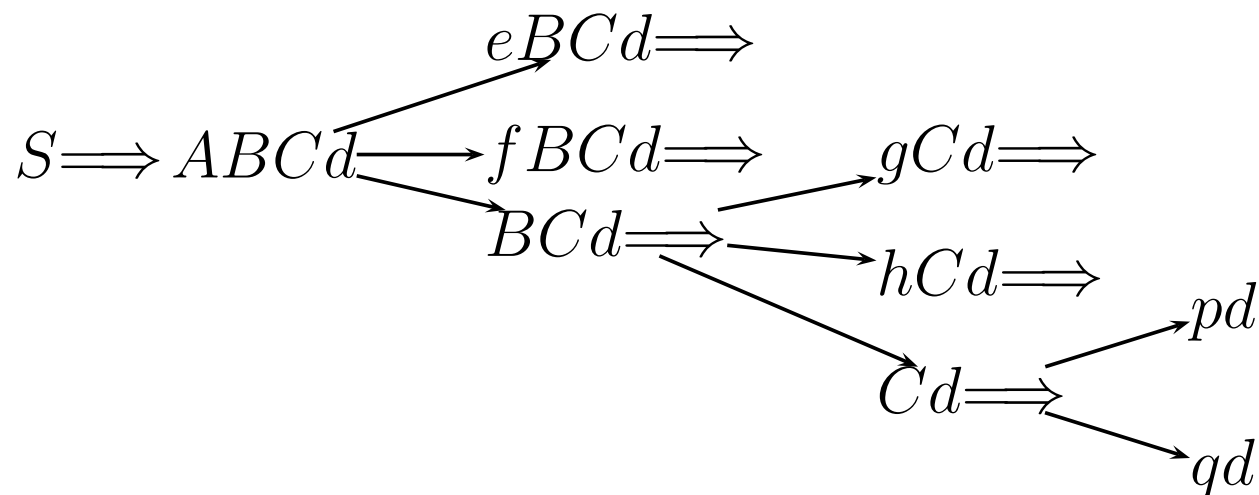   If $i = n$ is nullable, then add $\epsilon$ to $\mathsf{First}(\alpha)$

# Case 2 of the Procedure for Computing First

$$S \rightarrow ABCd$$
$$A \rightarrow e \mid f \mid \epsilon$$
$$B \rightarrow g \mid h \mid \epsilon$$
$$C \rightarrow p \mid q$$



$$\mathsf{First}(ABCd) = \{e, f, g, h, p, q\}$$

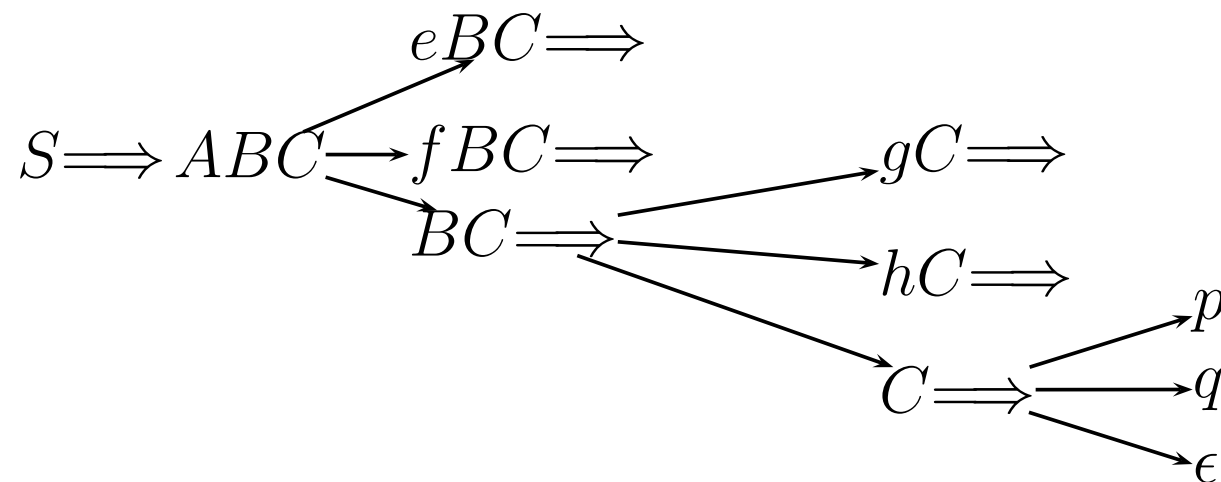# Case 2 of the Procedure for Computing First Again

$S \rightarrow ABC$               $d$ was deleted from the grammar in slide 218

$A \rightarrow e \mid f \mid \epsilon$

$B \rightarrow g \mid h \mid \epsilon$

$C \rightarrow p \mid q \mid \epsilon$

$$S \Longrightarrow ABC \begin{cases} eBC \Longrightarrow \\ fBC \Longrightarrow \\ BC \Longrightarrow \begin{cases} gC \Longrightarrow \\ hC \Longrightarrow \\ C \Longrightarrow \begin{cases} p \\ q \\ \epsilon \end{cases} \end{cases} \end{cases}$$

First$(ABC) = \{e, f, g, h, p, q, \epsilon\}$

# The Expression Grammar

● The grammar with left recursion:

$$\text{Grammar 1:} \quad E \to E + T \mid E - T \mid T$$
$$T \to T * F \mid T/F \mid F$$
$$F \to \textbf{INT} \mid (E)$$

● The transformed grammar without left recursion:

$$\text{Grammar 2:} \quad E \to TQ$$
$$Q \to +TQ \mid -TQ \mid \epsilon$$
$$T \to FR$$
$$R \to *FR \mid /FR \mid \epsilon$$
$$F \to \textbf{INT} \mid (E)$$

# First Sets for Grammar 1 (with Recursion)

$$\text{First}(E) = \text{First}(E + T) = \text{First}(E - T) =$$

$$\text{First}(T) = \text{First}(T * F) = \text{First}(T/F) =$$

$$\text{First}(F) = \{(, \textbf{INT}\}$$

$$\text{First}((E)) = \{(\}$$

$$\text{First}(\textbf{INT}) = \{\textbf{INT}\}$$

The explicit construction is left as an exercise.

# First Sets for Grammar 2 (without Recursion)

$$\mathsf{First}(E) \;=\; \mathsf{First}(TQ) \;\;=\; \{(, \mathbf{INT}\}$$
$$\mathsf{First}(T) \;=\; \mathsf{First}(FR) \;\;=\; \{(, \mathbf{INT}\}$$
$$\mathsf{First}(Q) \;\;\;\;\;=\; \{+, -, \epsilon\}$$
$$\mathsf{First}(R) \;\;\;\;\;=\; \{*, /, \epsilon\}$$
$$\mathsf{First}(F) \;\;\;\;\;=\; \{\mathbf{INT}, (\}$$
$$\mathsf{First}(+TQ) \;=\; \{+\}$$
$$\mathsf{First}(-TQ) \;=\; \{-\}$$
$$\mathsf{First}(*FR) \;=\; \{*\}$$
$$\mathsf{First}(/FR) \;=\; \{/\}$$
$$\mathsf{First}((E)) \;=\; \{(\}$$
$$\mathsf{First}(\mathbf{INT}) \;=\; \{\mathbf{INT}\}$$

# Why Follow Sets?

- **First** sets do not tell us when to apply $A{\rightarrow}\alpha$ such that $\alpha {\Longrightarrow}^* \epsilon$ (the important special case is $A{\rightarrow}\epsilon$)

- **Follow** sets do

- **Follow** sets constructed only for nonterminals

- By convention, assume every input is terminated by a special end marker (i.e., the EOF marker), denoted $\$$

- **Follow** sets do not contain $\epsilon$

## Definition of Follow Sets

Let $A$ be a nonterminal. Define $\mathsf{Follow}(A)$ to be the set of terminals that can appear immediately to the right of $A$ in some sentential form. That is,

$$\mathsf{Follow}(A) \;\; = \;\; \{a \mid S \Longrightarrow^* \cdots Aa \cdots \}$$

where $S$ is the start symbol of the grammar.

# A Procedure to Compute Follow Sets

1. If $A$ is the start symbol, add \$ to $\mathsf{Follow}(A)$.

2. Look through the grammar for all occurrences of $A$ on the right of productions. Let a typical production be:

$$B \;\rightarrow\; \alpha A \beta$$

   There are two cases – both may be applicable:
   
   (a) $\mathsf{Follow}(A)$ includes $\mathsf{First}(\beta) - \{\epsilon\}$.
   
   (b) If $\beta \Longrightarrow^* \epsilon$, then include $\mathsf{Follow}(B)$ in $\mathsf{Follow}(A)$.

# Follow Sets for Grammar 1 (with Recursion)

$$\text{Follow}(E) = \{+, -, ), \$\}$$
$$\text{Follow}(T) = \text{Follow}(F) = \{+, -, *, /, ), \$\}$$

The explicit construction is left as an exercise.

# Follow Sets for Grammar 2 (without Recursion)

$$
\begin{aligned}
\text{Follow}(E) &= \{ ), \$ \} \\
\text{Follow}(Q) &= \{ ), \$ \} \\
\text{Follow}(T) &= \{ +, -, ), \$ \} \\
\text{Follow}(R) &= \{ +, -, ), \$ \} \\
\text{Follow}(F) &= \{ +, -, *, /, ), \$ \}
\end{aligned}
$$

J. Xue

# Select Sets for Productions

- One **Select** set for every production in the grammar:

- The **Select** set for a production of the form $A \rightarrow \alpha$ is:

  - If $\epsilon \in \mathsf{First}(\alpha)$, then

    $$\mathsf{Select}(A \rightarrow \alpha) = (\mathsf{First}(\alpha) - \{\epsilon\}) \cup \mathsf{Follow}(A)$$

  - Otherwise:

    $$\mathsf{Select}(A \rightarrow \alpha) = \mathsf{First}(\alpha)$$

- The **Select** set predicts $A \rightarrow \alpha$ to be used in a derivation.

- Thus, the **Select** not needed if $A$ has has one alternative

COMP3131/9102 Page 228 March 22, 2010

# Select Sets for Grammar 1

Follow sets not used since the grammar has no $\epsilon$-productions

$$\mathsf{Select}(E{\rightarrow}E + T) = \mathsf{First}(E + T) = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(E{\rightarrow}E - T) = \mathsf{First}(E - T) = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(E{\rightarrow}T) \quad = \mathsf{First}(T) \quad = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(T{\rightarrow}T * F) = \mathsf{First}(T * F) = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(T{\rightarrow}T/F) \quad = \mathsf{First}(T/F) \quad = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(T{\rightarrow}F) \quad = \mathsf{First}(F) \quad = \{(, \mathbf{INT}\}$$
$$\mathsf{Select}(F{\rightarrow}\mathbf{INT}) \quad = \mathsf{First}(\mathbf{INT}) \quad = \{\mathbf{INT}\}$$
$$\mathsf{Select}(F{\rightarrow}(E)) \quad = \mathsf{First}((E)) \quad = \{(\}$$

# Select Sets for Grammar 2

$$
\begin{aligned}
\mathsf{Select}(E \rightarrow TQ) &= \mathsf{First}(TQ) &&= \{(, \mathbf{INT}\} \\
\mathsf{Select}(Q \rightarrow +\,TQ) &= \mathsf{First}(+TQ) &&= \{+\} \\
\mathsf{Select}(Q \rightarrow -\,TQ) &= \mathsf{First}(-TQ) &&= \{-\} \\
\mathsf{Select}(Q \rightarrow \epsilon) &= (\mathsf{First}(\epsilon) - \{\epsilon\}) \cup \mathsf{Follow}(Q) &&= \{), \$\} \\
\mathsf{Select}(T \rightarrow FR) &= \mathsf{First}(FR) &&= \{(, \mathbf{INT}\} \\
\mathsf{Select}(R \rightarrow *\,FR) &= \mathsf{First}(+FR) &&= \{*\} \\
\mathsf{Select}(R \rightarrow /FR) &= \mathsf{First}(/FR) &&= \{/\} \\
\mathsf{Select}(R \rightarrow \epsilon) &= (\mathsf{First}(\epsilon) - \{\epsilon\}) \cup \mathsf{Follow}(T) &&= \{+, -, ), \$\} \\
\mathsf{Select}(F \rightarrow \mathbf{INT}) &= \mathsf{First}(\mathbf{INT}) &&= \{\mathbf{INT}\} \\
\mathsf{Select}(F \rightarrow (E)) &= \mathsf{First}((E)) &&= \{(\}
\end{aligned}
$$

## Outline for the Rest of the Lecture

1. Definition of LL(1) grammar

2. One simplification in the presence of a nullable alternative

3. Eliminate left recursion and common prefixes

4. Write parsing methods in the presence of regular operators

5. $LL(k)$ for small $k$ often necessary for some constructs

# Definition of LL(1) Grammar

- A grammar is LL(1) if for every nonterminal of the form:

$$A \; \rightarrow \; \alpha_1 \mid \cdots \mid \alpha_n$$

the select sets are pairwise disjoint, i.e.:

$$\mathsf{Select}(A \rightarrow \alpha_i) \cap \mathsf{Select}(A \rightarrow \alpha_j) \; = \; \emptyset$$

for all $i$ and $j$ such that $i \neq j$.

- This implies there can be at most one nullable alternative

# Left Recursion

- **Direct** left-recursion:

$$A \quad \rightarrow \quad A\alpha$$

- **Non-direct** left-recursion:

$$A \quad \rightarrow \quad B\alpha$$
$$B \quad \rightarrow \quad A\beta$$

  - Algorithm 4.1 of text eliminates both kinds of left recursion
  - In real programming languages, non-direct left-recursion is rare
  - Not required

- A grammar with left recursion is not LL(1)

# Eliminating Direct Left Recursion

- The grammar $G_1$:

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_n \text{ // } \alpha_i \text{ does not beging with } A$$
$$A \rightarrow A\beta_1 \mid A\beta_2 \mid \cdots \mid A\beta_m$$

- The transformed grammar $G_2$:

$$A \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \cdots \mid \alpha_n A'$$
$$A' \rightarrow \beta_1 A' \mid \beta_2 A' \mid \cdots \mid \beta_m A' \mid \epsilon$$

- $G_1$ and $G_2$ define the same language: $L(G_1) = L(G_2)$

- Example: in Slide 220, Grammar 2 is the transformed version of Grammar 1

# Eliminating Direct Left Recursion: Special Case

- The grammar $G_1$:

$$A \rightarrow \alpha \qquad // \ \alpha \text{ does not beging with } A$$
$$A \rightarrow A\beta$$

- The transformed grammar $G_2$:

$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta A' \mid \epsilon$$

# Eliminating Common Prefixes: Left-Factoring

- The grammar $G_1$:

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_m$$
$$A \rightarrow \gamma$$

- The transformed grammar $G_2$:

$$A \rightarrow \alpha A'$$
$$A \rightarrow \gamma$$
$$A' \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m$$

- $L(G_1) = L(G_2)$

- A grammar with common prefixes is not LL(1)

# Example: Eliminating Common Prefixes

$$\langle\text{stmt}\rangle \rightarrow \textbf{IF} \text{ ''('' } \langle\text{expr}\rangle \text{ '')'' } \langle\text{stmt}\rangle \textbf{ ELSE } \langle\text{stmt}\rangle$$
$$\langle\text{stmt}\rangle \rightarrow \textbf{IF} \text{ ''('' } \langle\text{expr}\rangle \text{ '')'' } \langle\text{stmt}\rangle$$
$$\langle\text{stmt}\rangle \rightarrow \textbf{other}$$

$$\Downarrow$$

$$\langle\text{stmt}\rangle \rightarrow \textbf{IF} \text{ ''('' } \langle\text{expr}\rangle \text{ '')'' } \langle\text{stmt}\rangle \langle\text{else-cluase}\rangle$$
$$\langle\text{else-clause}\rangle \rightarrow \textbf{ELSE } \langle\text{stmt}\rangle \mid \epsilon$$
$$\langle\text{stmt}\rangle \rightarrow \textbf{other}$$

# Eliminating Direct Left Recursion Using Regular Operators

- The grammar $G_1$:

$$A \;\rightarrow\; \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_n$$
$$A \;\rightarrow\; A\beta_1 \mid A\beta_2 \mid \cdots \mid A\beta_m$$

- The transformed grammar $G_2$:

$$A \;\rightarrow\; (\alpha_1 \mid \cdots \mid \alpha_n)(\beta_1 \mid \cdots \mid \beta_m)^*$$

- $G_1$ and $G_2$ define the same language: $L(G_1) = L(G_2)$

- Recommended to use in Assignment 2, where $n = m = 1$ for most left-recursive cases:

$$A \rightarrow \alpha\beta^*$$

# The Expression Grammar

- The grammar with left recursion:

$$\text{Grammar 1:} \quad E \to E + T \mid E - T \mid T$$
$$T \to T * F \mid T/F \mid F$$
$$F \to \textsf{INT} \mid (E)$$

- Eliminating left recursion using the Kleene Closure

$$\text{Grammar 3:} \quad E \to T \ (\text{"+"} \ T \mid \text{"-"} \ T)^*$$
$$T \to F \ (\text{"*"} \ F \mid \text{"/"} \ F)^*$$
$$F \to \textsf{INT} \mid \text{"("} \ E \ \text{")"}$$

All tokens are enclosed in double quotes to distinguish them for the regular operators: (, ) and $^*$

- Compare with Slide 220

# Eliminating Common Prefixes using Choice Operator

- The grammar $G_1$:

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_m$$
$$A \rightarrow \gamma$$

- The transformed grammar $G_2$:

$$A \rightarrow \alpha(\beta_1 \mid \beta_2 \mid \cdots \mid \beta_m)$$
$$A \rightarrow \gamma$$

- Recommended to use in Assignment 2

# Example: Eliminating Common Prefixes

$$\begin{aligned}
\langle \text{stmt} \rangle &\rightarrow \textbf{IF } \text{"("} \langle \text{expr} \rangle \text{")"} \langle \text{stmt} \rangle \textbf{ ELSE } \langle \text{stmt} \rangle \\
\langle \text{stmt} \rangle &\rightarrow \textbf{IF } \text{"("} \langle \text{expr} \rangle \text{")"} \langle \text{stmt} \rangle \\
\langle \text{stmt} \rangle &\rightarrow \textbf{other}
\end{aligned}$$

$\Downarrow$

$$\begin{aligned}
\langle \text{stmt} \rangle &\rightarrow \textbf{IF } \text{"("} \langle \text{expr} \rangle \text{")"} \langle \text{stmt} \rangle \text{ ( } \textbf{ELSE } \langle \text{stmt} \rangle \text{)?} \\
\langle \text{stmt} \rangle &\rightarrow \textbf{other}
\end{aligned}$$

Compare with Slide 251

# LL(k) Grammar and Parsing

- A grammar is LL($k$) if it can be parsed deterministically using $k$ tokens of lookahead

- A formal definition for LL(k) grammars can be found in Grune and Jacobs' book                    es

- Grammar 1 in Slide 220 is not LL($k$) for any $k$!

- However, Grammar 2 in Slide 220 is LL(1)

- Only a understanding of LL(1) is required this year ear