

Compilation

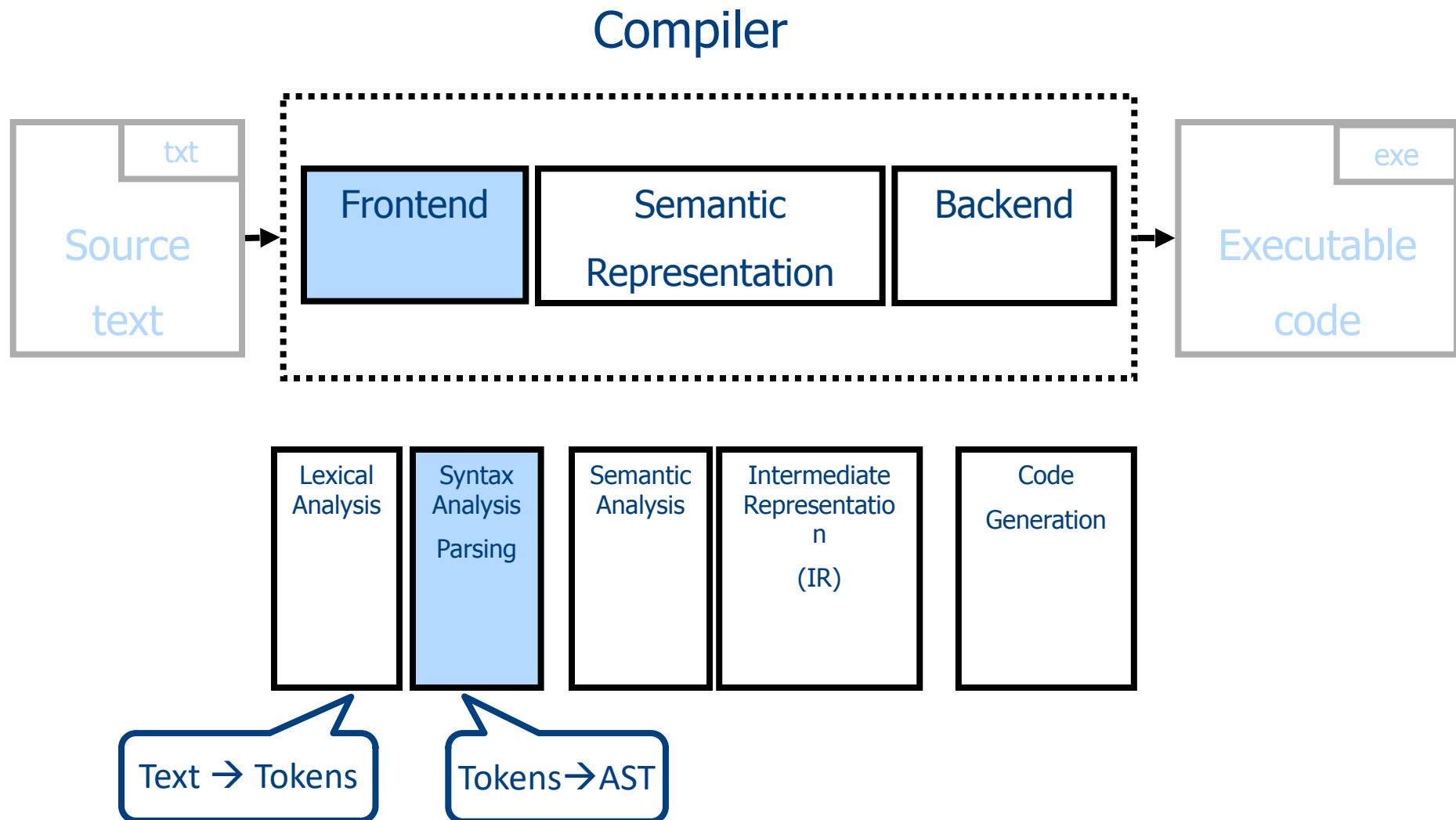
0368-3133

Lecture 2b:

Syntax Analysis:

Top-Down parsing

Conceptual Structure of a Compiler

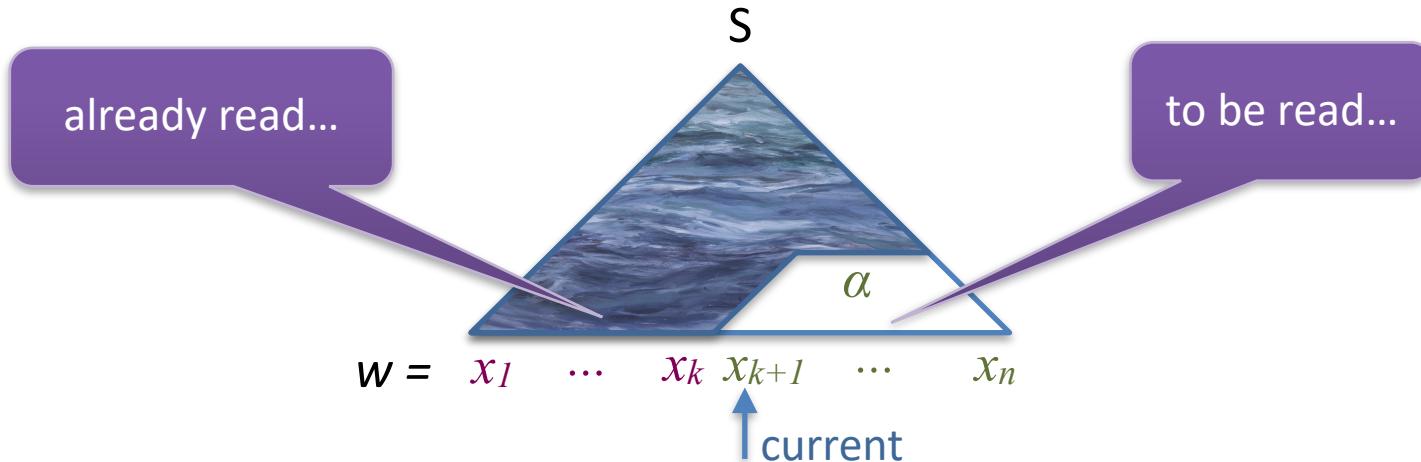


Parsing

- Goals
 - Decide whether the sequence of tokens a valid program in the source language
 - Construct a structured representation of the input
 - Error detection and reporting

Top-Down (Predictive) Parsers

- Construct the leftmost derivation
 - Read input from left to right
 - Apply rules “from left to right”
- Predict what rule to apply based on nonterminal and current token



$$S \Rightarrow A t B \Rightarrow \beta t B \Rightarrow^* x_1 \dots x_k \alpha$$

Top-Down Parsing

- Given a grammar $G = \langle V, T, P, S \rangle$ and a sentence w
- Goal: derive w using G
 - ▶ $S \xrightarrow{*} w$
- Idea
 - ▶ Apply production to leftmost nonterminal
 - ▶ Pick production rule based on next input token

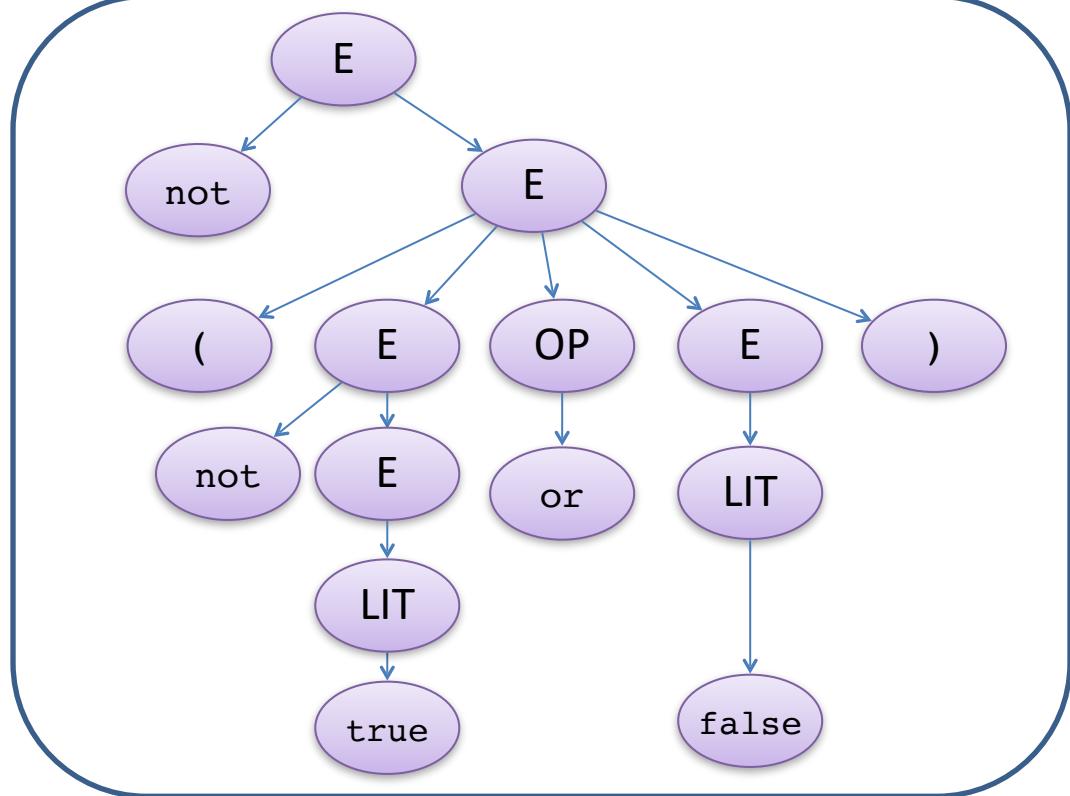


Boolean Expressions Example

not (not true
or false)

$E \rightarrow \text{LIT} \mid (\ E \text{ OP } E) \mid \text{not } E$
 $\text{LIT} \rightarrow \text{true} \mid \text{false}$
 $\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

$E \Rightarrow$
 $\text{not } E \Rightarrow$
 $\text{not } (E \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } E \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{LIT} \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ OP } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } E) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } \text{LIT}) \Rightarrow$
 $\text{not } (\text{not } \text{true} \text{ or } \text{false})$



Production to apply is known from **next input token**

Challenges in top-down parsing

- Top-down parsing begins with virtually no information
 - Begins with just the start symbol, which matches every program
- How can we know which productions to apply?

Which productions to apply?

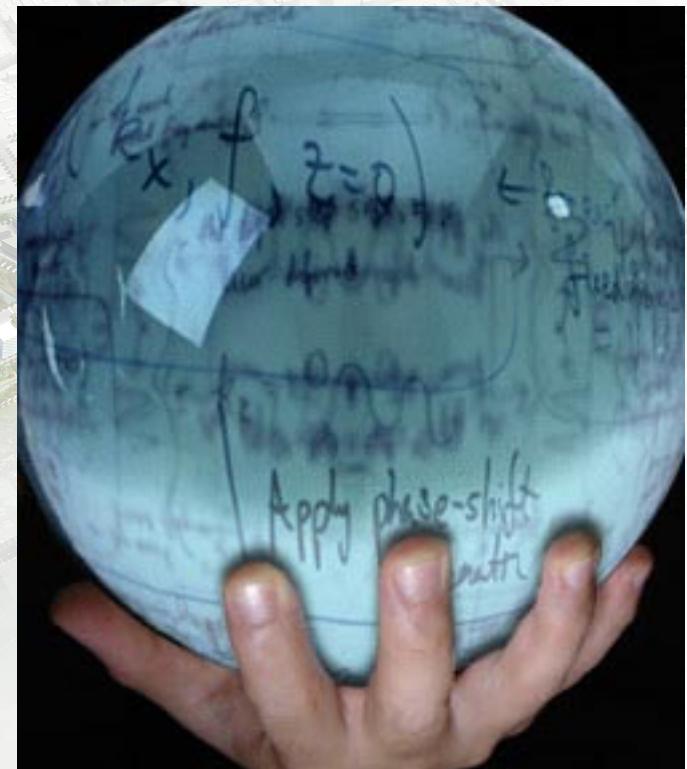
- In general, we can't guess effectively
 - There are some grammars for which the best we can do is guess and backtrack if we're wrong
 - Recall “Brute force” parsing

⇒ We will sacrifice generality for predictability

- Restricted grammars ($LL(k)$)
 - Know exactly which rule to apply
 - May require k lookahead tokens to decide

Predictive Parsing

- Recursive descent
- LL(k) grammars
- Generated LL(k) parsers
 - Using a pushdown automata



Recursive Descent Parsing

Recursive Descent Parsing

- Define a function for every nonterminal
- Every function work as follows:
 - Select a production rule
 - Go over rule's rhs "from left to right" $(A \rightarrow X_1 \dots X_n)$
 - For terminal symbols — checks **match** with next input token
 - For nonterminal symbols — calls (recursively) their corresponding functions
- When there are several productions for a nonterminal, use **lookahead**

rhs
↓

Matching tokens

```
void match(token t) {  
    if (current == t)  
        current = next_token();  
    else  
        error;  
}
```

- Variable **current** holds the current input token

Functions for Nonterminals

```
E → LIT | ( E OP E ) | not E
```

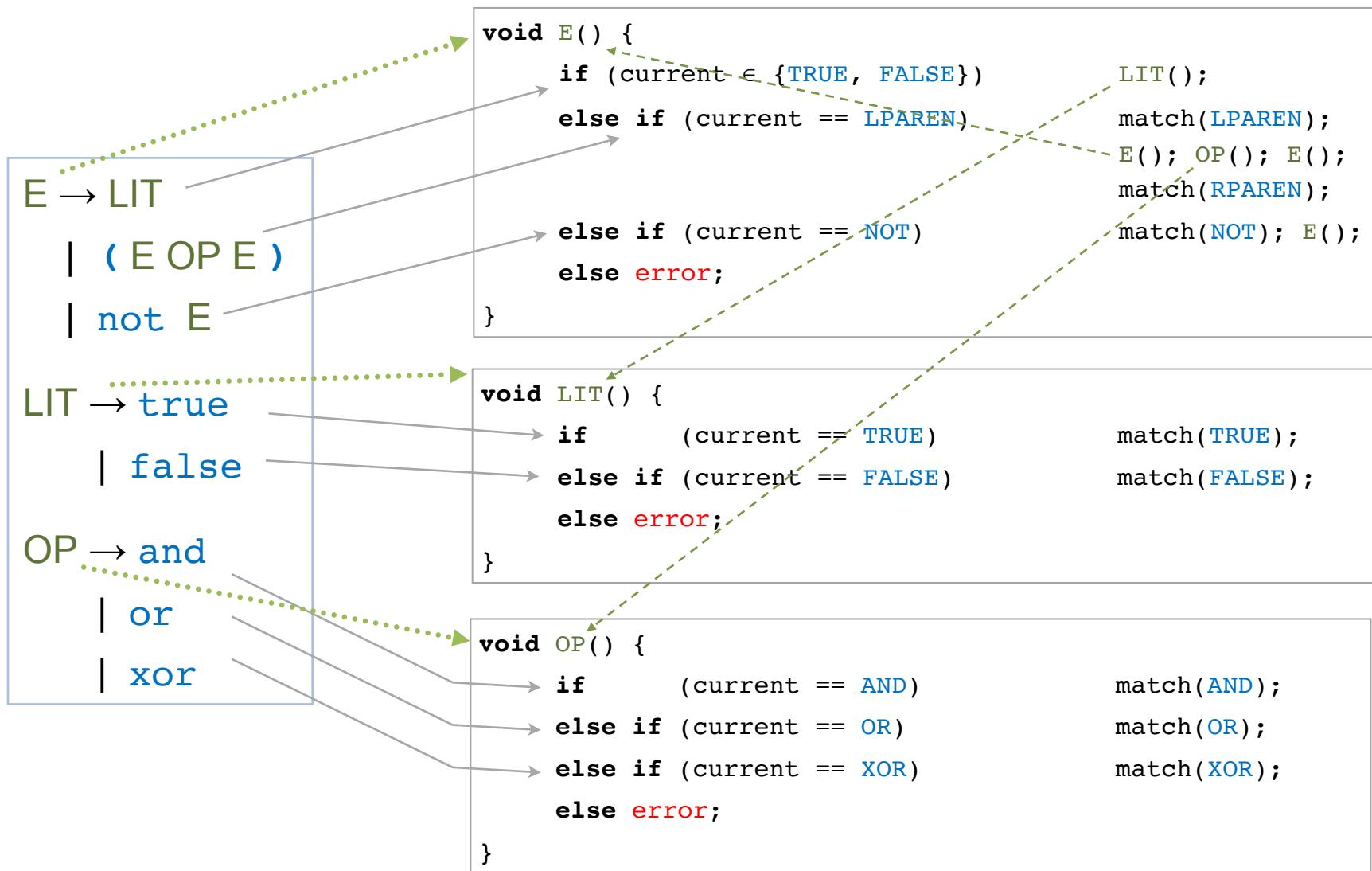
```
LIT → true | false
```

```
OP → and | or | xor
```

```
void E() {
    if (current ∈ {TRUE, FALSE}) // E → LIT
        LIT();
    else if (current == LPAREN) // E → ( E OP E )
        match(LPAREN); E(); OP(); E(); match(RPAREN);
    else if (current == NOT) // E → not E
        match(NOT); E();
    else
        error;
}

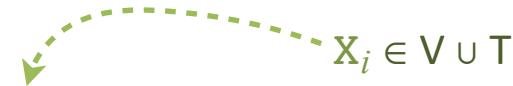
void LIT() {
    if (current == TRUE) match(TRUE);
    else if (current == FALSE) match(FALSE);
    else error;
}
```

Functions for Nonterminals



Recursive Descent – General Structure

```
void A() {  
    choose an A-production, A → X1X2⋯Xk ;  
  
    for (i = 1; i ≤ k; i++) {  
        if (Xi is a nonterminal) call function Xi();  
        else call match(Xi);  
    }  
}
```



$X_i \in V \cup T$

Recursive Descent – General Structure

```
void A() {  
    choose an A-production,  $A \rightarrow X_1 X_2 \cdots X_k$ ;  
    for (i = 1; i ≤ k; i++) {  
        if ( $X_i$  is a nonterminal) call function  $X_i()$ ;  
        else if ( $X_i$  == current) current = next_token();  
        else error;  
    }  
}
```

A diagram illustrating the general structure of Recursive Descent parsing. A red oval encloses the 'choose' step. A green curved arrow points from this oval to the production rule $A \rightarrow X_1 X_2 \cdots X_k$. A red arrow points from the 'choose' step down to the 'else' block.

- How do you pick the right A-production?
 - Can we generalize the rationale behind all those **if** statements we saw in the example?

LL(k) Grammars

LL(k) Grammars

- A grammar is in the class LL(k) when it can be parsed via:
 - Top-down analysis
 - Scanning the input from **left to right** (L)
 - Producing the **leftmost derivation** (L)
 - With **lookahead** of k tokens (k)
- A language is said to be LL(k) if it has an LL(k) grammar
- The simplest case is LL(1), which we discuss next

FIRST Sets

- To formalize the property (of a grammar) that we can determine a rule using a single lookahead token (LL(1)), we define the FIRST sets.
- For a **sequence** of symbols $\alpha \in (V \cup T)^*$
 - $\text{FIRST}(\alpha)$ = all terminals that α **can start with**
 - i.e., $t \in \text{FIRST}(\alpha) \Leftrightarrow$ there exists a derivation sequence $\alpha \xrightarrow{*} tw$
 - Special case: $\varepsilon \in \text{FIRST}(\alpha) \Leftrightarrow \alpha \xrightarrow{*} \varepsilon$ (α is *nullable*)

⋮
(not really
a terminal)

FIRST Sets - Example

```
E → LIT | ( E OP E ) | not E | LIT maybe  
LIT → true | false | ε  
OP → and | or | xor
```

- In our Boolean expressions example:

- ▶ $\text{FIRST}(\text{ LIT }) = \boxed{\quad ? \quad}$
 - ▶ $\text{FIRST}(\text{ (E OP E) }) = \boxed{\quad ? \quad}$
 - ▶ $\text{FIRST}(\text{ not E }) = \boxed{\quad ? \quad}$
- $\text{FIRST}(\text{ E }) = \boxed{\quad ? \quad}$
- These are all the alternatives for $E \rightarrow \alpha$

Recursive Descent (revised)

Special case: No ϵ -productions

```
void A() {
```

find an A-production, $A \rightarrow \overbrace{X_1 X_2 \dots X_k}^{\alpha}$,

such that `current` \in `FIRST`(α)

if no such production exists `error`; // If there are multiple applicable
// productions, the grammar is not LL(1)

for ($i = 1; i \leq k; i++$) {

if (X_i is a nonterminal) **call** function `Xi()`;

else if ($X_i == \text{current}$) `current = next_token();`

else `error`;

}

}

FIRST Sets

- No intersection between FIRST sets of rules of the same nonterminal \Rightarrow can *always** pick a single rule
- If the FIRST sets intersect, may need longer lookahead
 - ▶ $\text{LL}(1) \subset \text{LL}(2) \subset \text{LL}(3) \subset \dots$
 - ▶ But the size of FIRST sets grow exponentially with k

* except when null (ε) productions exist — see next slides

Recursive Descent (Naive)

Handling ϵ -productions

```
void A() {
```

find an A-production, $A \rightarrow \overbrace{X_1 X_2 \dots X_k}^{\alpha}$,

such that `current` $\in \text{FIRST}(\alpha)$ or $\alpha == \epsilon$

(can be chosen
regardless of `current`)

if no such production exists **error**;

```
if ( $\alpha == \epsilon$ ) return;
```

```
for (i = 1; i <= k; i++) {
```

if (X_i is a nonterminal)

call function $X_i()$;

else if ($X_i == \text{current}$)

`current = next_token();`

else

error;

```
}
```

```
}
```

Handling null-productions

- What's the problem with $A \rightarrow \epsilon$?
- Example:

$S \rightarrow Ab \mid c \quad A \rightarrow a \mid \epsilon$

$L(S) = \boxed{?}$

$\text{FIRST}(A) = \boxed{?}$

$\text{FIRST}(Ab) = \boxed{?}$

- Consider inputs: ab and b
- Need to know what comes after A in the sentence to select the right production

FOLLOW sets

FOLLOW Sets

- For *every* nonterminal A
 - $\text{FOLLOW}(A)$ = set of "terminals" that can immediately follow A (in some derivation)

FOLLOW Sets

- For every nonterminal A
 - $\text{FOLLOW}(A) = \text{set of "terminals" that can immediately follow } A \text{ (in } \underline{\text{some}} \text{ derivation)}$
 - i.e., $t \in \text{FOLLOW}(A) \Leftrightarrow \text{there exists some derivation sequence } S \xrightarrow{*} \alpha A t \beta$
 - Special case: $\$ \in \text{FOLLOW}(A) \Leftrightarrow S \xrightarrow{*} \alpha A$ (i.e., A at the end)
(not really a terminal)

$$\begin{aligned}\text{FOLLOW}(A) = & \{ \text{t} \in T \mid S \xrightarrow{*} \alpha A t \beta \} \cup \\ & \{ \$ \mid S \xrightarrow{*} \alpha A \}\end{aligned}$$

FOLLOW Sets - Example

$$\begin{array}{l} S \rightarrow Ab \mid c \\ A \rightarrow \epsilon \mid a \end{array}$$
$$\begin{array}{l} S \rightarrow Aab \mid c \\ A \rightarrow \epsilon \mid a \end{array}$$

$\text{FOLLOW}(A) =$?

$\text{FOLLOW}(A) =$?

Choosing Rules using FIRST & FOLLOW Sets - Example

$$S \rightarrow Ab \mid c$$
$$A \rightarrow \epsilon \mid a$$

$$S \rightarrow Aab \mid c$$
$$A \rightarrow \epsilon \mid a$$
$$\text{FIRST}(\epsilon) = \{ \epsilon \}$$
$$\text{FIRST}(a) = \{ a \}$$
$$\text{FOLLOW}(A) = \{ b \}$$
$$\text{FIRST}(\epsilon) = \{ \epsilon \}$$
$$\text{FIRST}(a) = \{ a \}$$
$$\text{FOLLOW}(A) = \{ a \}$$

FIRST-FOLLOW Conflict - Example

$$\begin{array}{l} S \rightarrow Ab \mid c \\ A \rightarrow \epsilon \mid a \end{array}$$
$$\begin{array}{l} S \rightarrow Aab \mid c \\ A \rightarrow \epsilon \mid a \end{array}$$

$$\text{FIRST}(\epsilon) = \{ \epsilon \}$$

$$\text{FIRST}(a) = \{ a \}$$

$$\text{FOLLOW}(A) = \{ b \}$$

$$\text{FIRST}(\epsilon) = \{ \epsilon \}$$

$$\text{FIRST}(a) = \{ a \}$$

$$\text{FOLLOW}(A) = \{ a \}$$

Consider ab and aab

LL(1) grammars

- A grammar is in the class LL(1) iff
 - For every two productions $A \rightarrow \alpha$ and $A \rightarrow \beta$ we have
 - $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \{\}$ // including ϵ
 - If $\epsilon \in \text{FIRST}(\alpha)$ then $\text{FIRST}(\beta) \cap \text{FOLLOW}(A) = \{\}$
 - If $\epsilon \in \text{FIRST}(\beta)$ then $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \{\}$

Recursive Descent (revised again)

```
void A() {  
    find an A-production,  $A \rightarrow \overbrace{X_1 X_2 \dots X_k}^{\alpha}$ ,  
    such that current  $\in \text{FIRST}(\alpha)$   
    or  $\epsilon \in \text{FIRST}(\alpha)$  and current  $\in \text{FOLLOW}(A)$   
  
    if no such production exists error;  
  
    if ( $\alpha == \epsilon$ ) return;  
  
    for (i = 1; i  $\leq k$ ; i++) {  
        if ( $X_i$  is a nonterminal) call function  $X_i()$ ;  
        else if ( $X_i == \text{current}$ ) current = next_token();  
        else  
    }  
}  
}  
This is the basis for an LL(1) parser.  
(but it is still recursive — stay tuned)
```

Computing FIRST sets

FIRST Sets - Recursive Definition

$$\text{FIRST}(\alpha) = \{ t \in T \mid \alpha \Rightarrow^* t\beta \} \cup \{ \varepsilon \mid \alpha \Rightarrow^* \varepsilon \}$$

- For every terminal t and nonterminal N :
 - ▶ $\text{FIRST}(\varepsilon) = \{\varepsilon\}$
 - ▶ $\text{FIRST}(t) = \{t\}$ t terminal
 - ▶ $\text{FIRST}(N) = \bigcup \{\text{FIRST}(\alpha) \mid N \rightarrow \alpha \in P\}$ N nonterminal

FIRST Sets - Recursive Definition

$$\text{FIRST}(\alpha) = \{ t \in T \mid \alpha \Rightarrow^* t\beta \} \cup \{\epsilon \mid \alpha \Rightarrow^* \epsilon\}$$

- For every terminal t and nonterminal N :
 - $\text{FIRST}(\epsilon) = \{\epsilon\}$
 - $\text{FIRST}(t) = \{t\}$ t terminal
 - $\text{FIRST}(N) = \bigcup \{\text{FIRST}(\alpha) \mid N \rightarrow \alpha \in P\}$ N nonterminal
- β any sequence of terminals and nonterminals ↓
- X is terminal or nonterminal →
- $\text{FIRST}(X\beta) = \begin{cases} \text{FIRST}(X) & \epsilon \notin \text{FIRST}(X) \\ \text{FIRST}(X) \setminus \{\epsilon\} \cup \text{FIRST}(\beta) & \epsilon \in \text{FIRST}(X) \end{cases}$ ↑
- Possible that $\beta=\epsilon$

FIRST Sets - Constraint Definition

$$\text{FIRST}(X) = \{ t \in T \mid X \Rightarrow^* t\beta \} \cup \{ \epsilon \mid X \Rightarrow^* \epsilon \}$$

- $\text{FIRST}(\epsilon) = \{\epsilon\}$
- $\text{FIRST}(t) = \{t\}$
- For every rule $X \rightarrow A_1 \dots A_k \alpha$
such that $\epsilon \in \text{FIRST}(A_i)$ for every $i=1\dots k$
 $\text{FIRST}(\alpha) \subseteq \text{FIRST}(X)$
 - α can be ϵ

if $S \Rightarrow^* \gamma X \mu \Rightarrow \gamma A_1 \dots A_k \alpha \mu \wedge A_1 \Rightarrow^* \epsilon \wedge \dots \wedge A_k \Rightarrow^* \epsilon \wedge \alpha \Rightarrow^* t\beta$
then $S \Rightarrow^* \gamma X \mu \Rightarrow^* \gamma t\beta \mu$

Computing FIRST sets

$$G = (V, T, P, S)$$

1. For all $A \in V$, $\text{FIRST}(A) = \{t \mid A \rightarrow t\beta \in P\} \cup \{\epsilon \mid A \rightarrow \epsilon \in P\}$
2. Repeat until no $\text{FIRST}(X)$ changes:
foreach $A \in V$
foreach $(A \rightarrow B_1 B_2 \dots B_k B_{k+1} \dots B_n) \in P$
if $(\epsilon \in \text{FIRST}(B_1) \wedge \dots \wedge \epsilon \in \text{FIRST}(B_n))$
$$\text{FIRST}(A) = \text{FIRST}(A) \cup \{\epsilon\}$$

for $(k = 1 \dots n-1)$
if $(\epsilon \in \text{FIRST}(B_1) \wedge \dots \wedge \text{FIRST}(B_k)) \ // k < n$
$$\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(B_{k+1}) \setminus \{\epsilon\})$$
- This is known as fixed-point computation

FIRST sets computation example

STMT	EXPR	TERM

$\text{STMT} \rightarrow \text{if EXPR then STMT}$
| $\text{while EXPR do STMT}$
| EXPR ;
 $\text{EXPR} \rightarrow \text{TERM}$
| zero? TERM
| not EXPR
| ++ id
| -- id
 $\text{TERM} \rightarrow \text{id}$
| constant

1. Initialization

STMT	EXPR	TERM
if	zero?	id
while	Not ++ --	constant

$\text{STMT} \rightarrow \text{if EXPR then STMT}$
| $\text{while EXPR do STMT}$
| EXPR ;
 $\text{EXPR} \rightarrow \text{TERM}$
| zero? TERM
| not EXPR
| ++ id
| -- id
 $\text{TERM} \rightarrow \text{id}$
| constant

2. Iterate 1

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --		

$\text{STMT} \rightarrow \text{if EXPR then STMT}$
 | $\text{while EXPR do STMT}$
 | EXPR ;
 $\text{EXPR} \rightarrow \text{TERM}$
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 $\text{TERM} \rightarrow \text{id}$
 | constant

2. Iterate 2

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	

$\text{STMT} \rightarrow \text{if EXPR then STMT}$
 | $\text{while EXPR do STMT}$
 | EXPR ;
 $\text{EXPR} \rightarrow \text{TERM}$
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 $\text{TERM} \rightarrow \text{id}$
 | constant

2. Iterate 3 – fixed-point

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	
id constant		

$\text{STMT} \rightarrow \text{if EXPR then STMT}$
 | $\text{while EXPR do STMT}$
 | EXPR ;
 $\text{EXPR} \rightarrow \text{TERM}$
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
 $\text{TERM} \rightarrow \text{id}$
 | constant