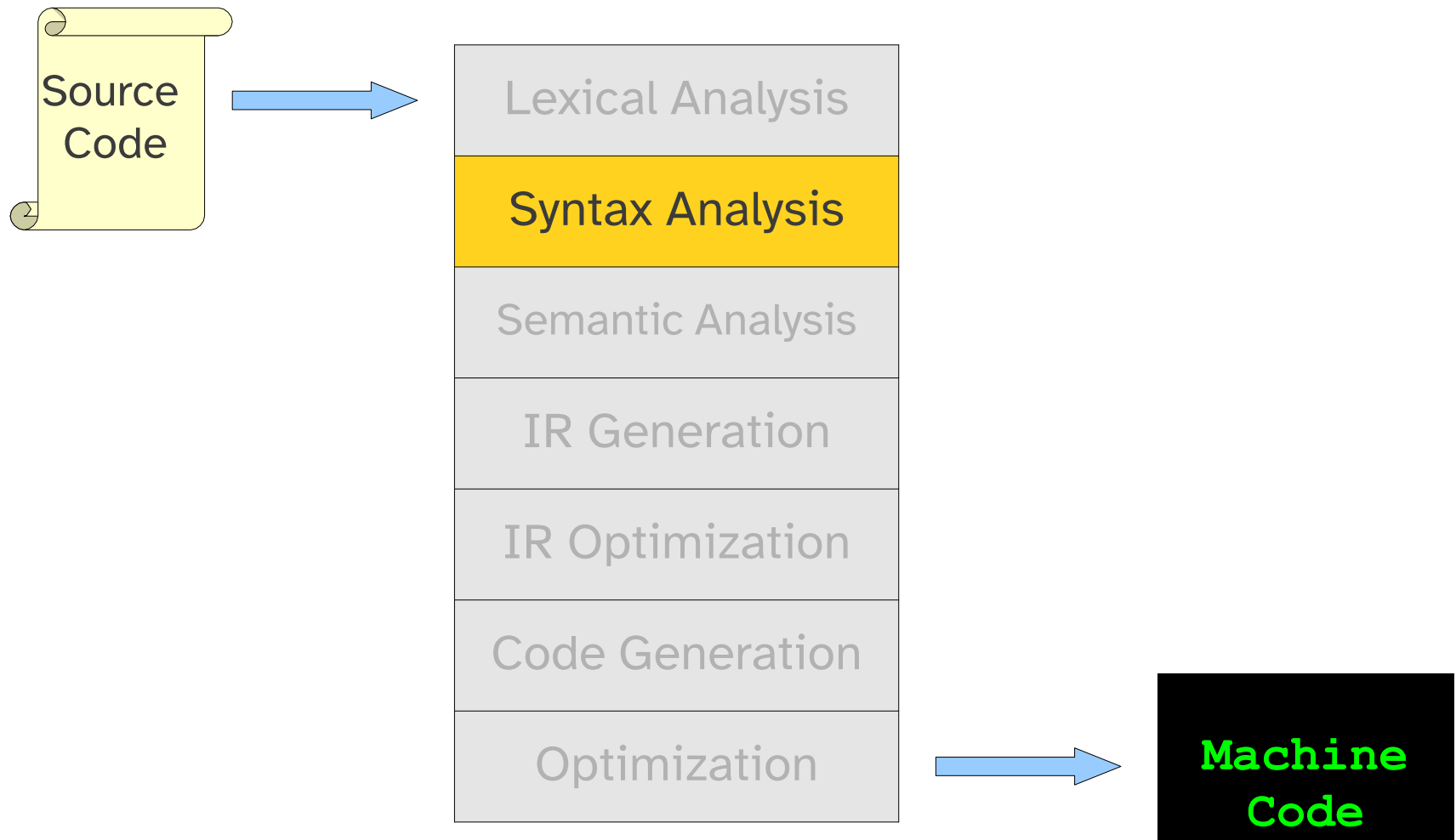


# Top-Down Parsing

# Where We Are



# Review from Last Time

- Goal of syntax analysis: recover the intended structure of the program.
- First step: Use a **context-free grammar** to describe the language.
- Given a sequence of tokens, look for a **parse tree** that generates those tokens.
- Recovering this syntax tree is called parsing

# Different Types of Parsing

- **Top-Down Parsing**

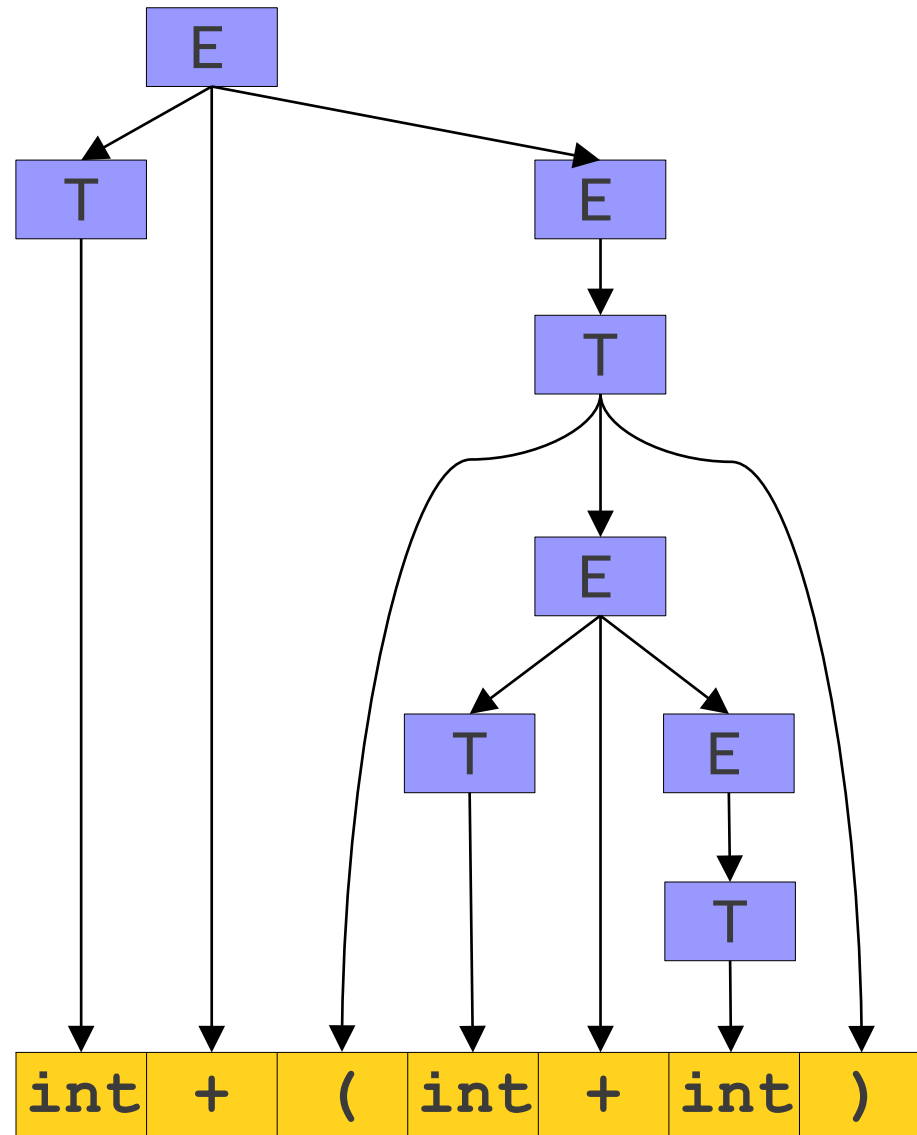
- Beginning with the start symbol, try to guess the productions to apply to end up at the user's program.

- **Bottom-Up Parsing**

- Beginning with the user's program, try to apply productions in reverse to convert the program back into the start symbol.

# Top-Down Parsing

$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

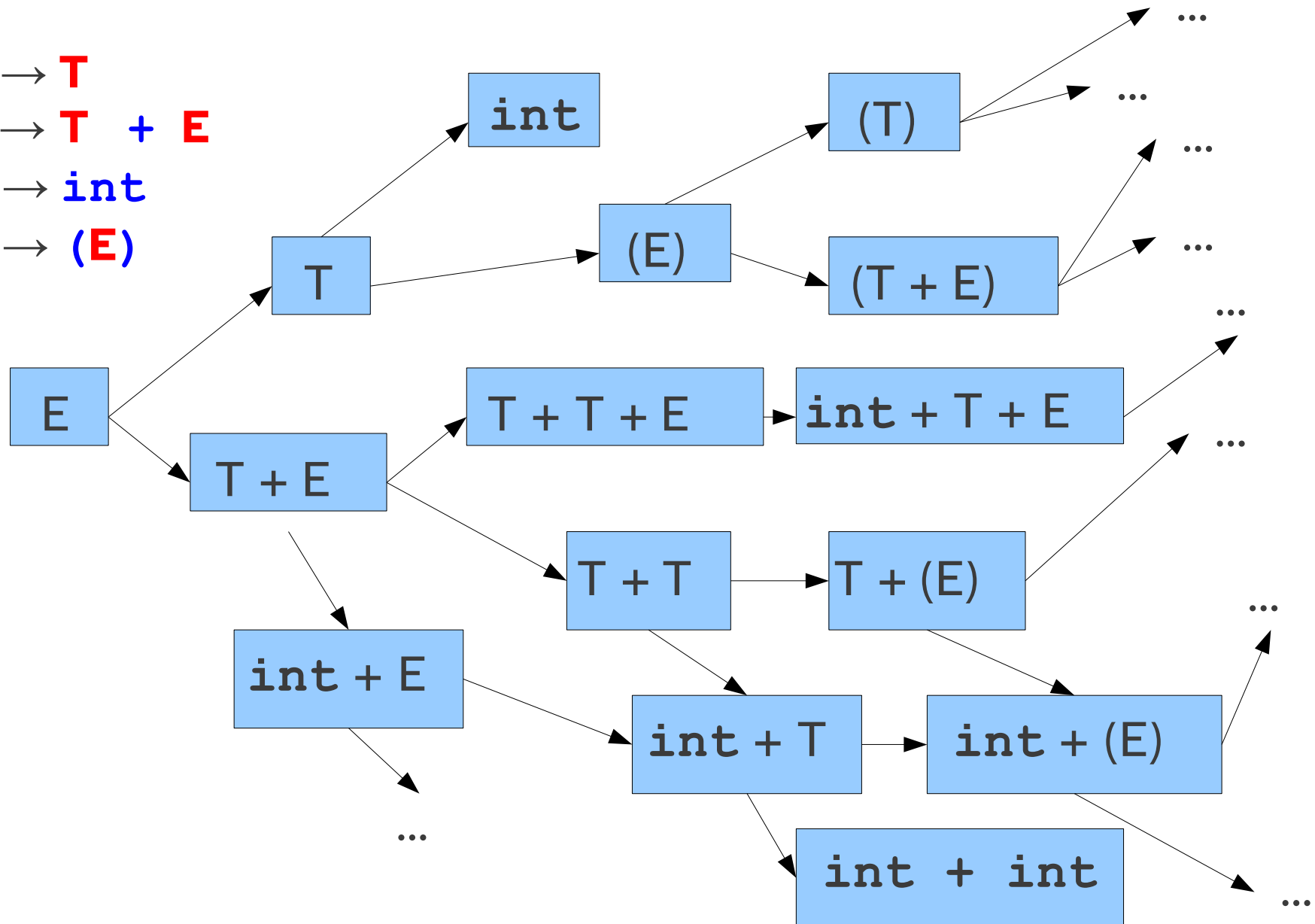


# Parsing as a Search

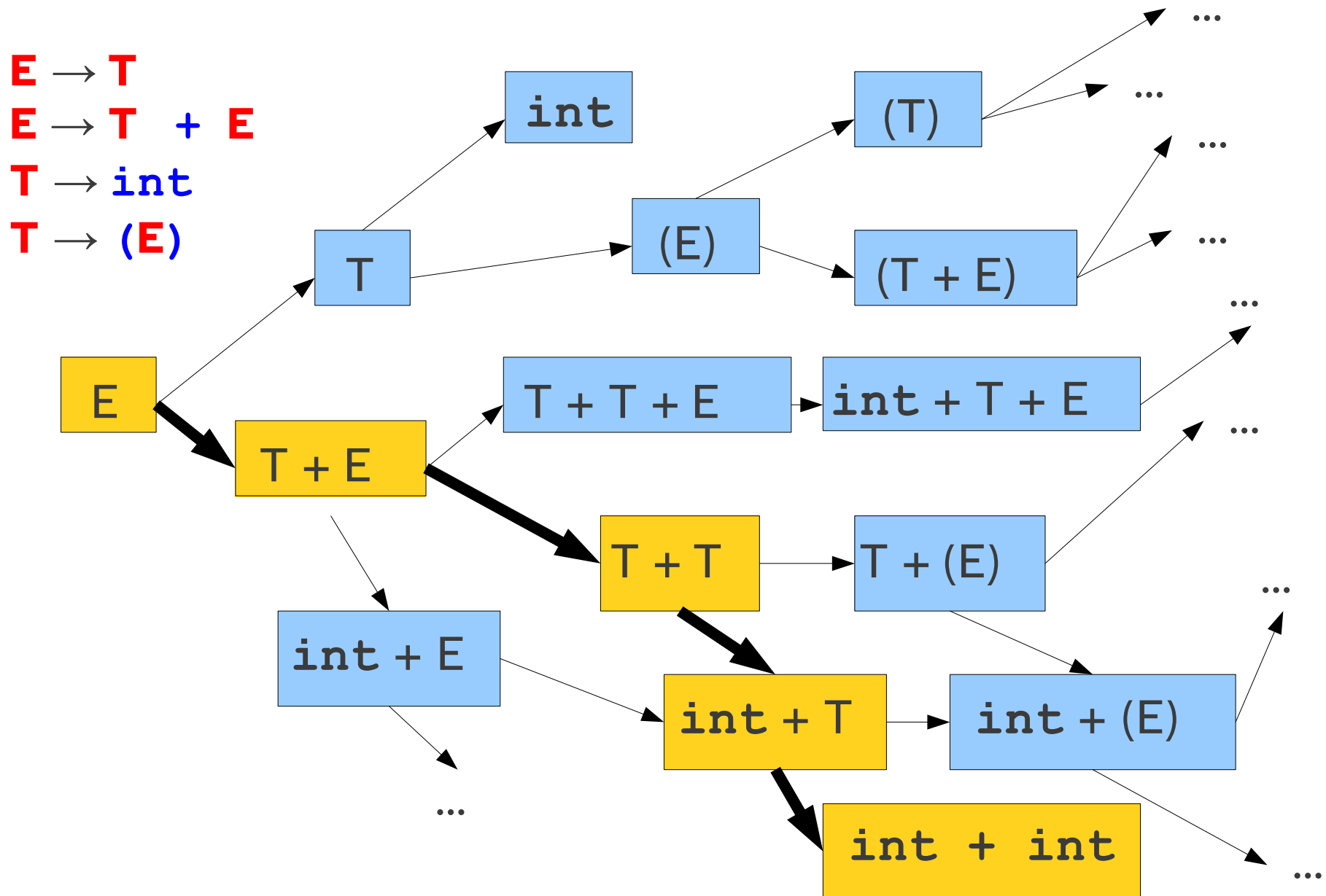
- An idea: **treat parsing as a graph search**
- Each node is a **sentential form**, i.e., a string of terminals and nonterminals derivable from the start symbol.
- There is an edge from node  $\alpha$  to node  $\beta$  iff  $\alpha \Rightarrow \beta$ .

# Parsing as a Search

$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$



# Parsing as a Search





# Our First Top-Down Algorithm

- **Breadth-First Search**
- Maintain a worklist of sentential forms, initially just the start symbol **S**.
- While the worklist isn't empty:
  - Remove an element from the worklist.
  - If it matches the target string, you're done.
  - Otherwise, for each possible string that can be derived in one step, add that string to the worklist.
- Can recover a parse tree by tracking what productions we applied at each step.

# BFS is Slow

- Enormous time and memory usage:
- Lots of **wasted effort**:
  - Generates a lot of sentential forms that couldn't possibly match.
  - But in general, extremely hard to tell whether a sentential form can match – that's the job of parsing!
- High **branching factor**:
  - Each sentential form can expand in (potentially) many ways for each nonterminal it contains.

# Leftmost Derivations

- Recall: A **leftmost derivation** is one where we always expand the leftmost symbol first.
- Updated algorithm:
- Do a breadth-first search, **only considering leftmost derivations**.
  - Dramatically drops branching factor.
  - Increases likelihood that we get a prefix of nonterminals.
- Prune sentential forms that can't possibly match.
  - Avoids wasted effort.

# Leftmost BFS

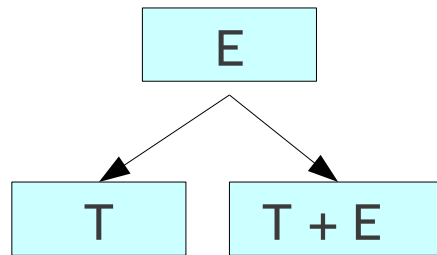


$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

`int + int`

# Leftmost BFS

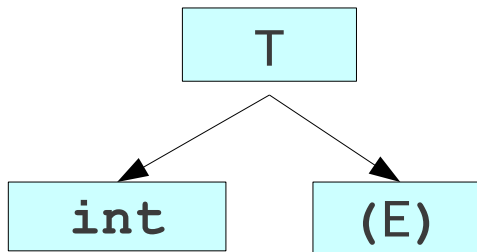
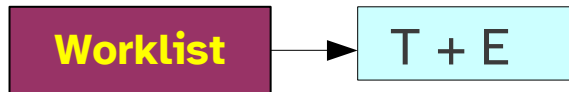
Worklist



$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

`int + int`

# Leftmost BFS

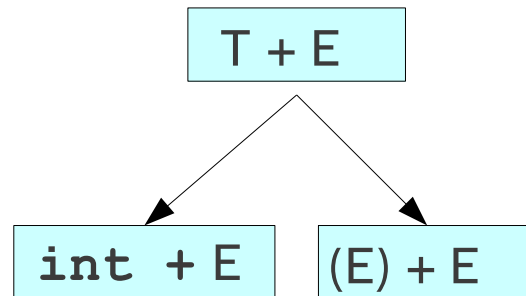


$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

`int + int`

# Leftmost BFS

Worklist



$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

`int + int`

# Leftmost BFS Has Problems

Worklist

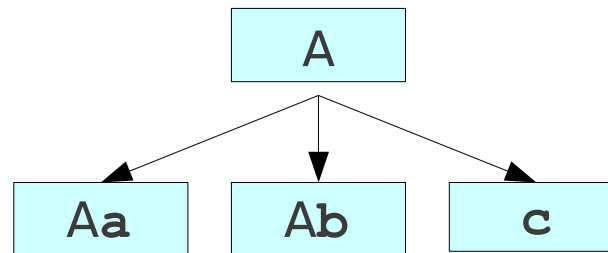
$A \rightarrow Aa \mid Ab \mid c$

caaaaaaaaaa



# Leftmost BFS Has Problems

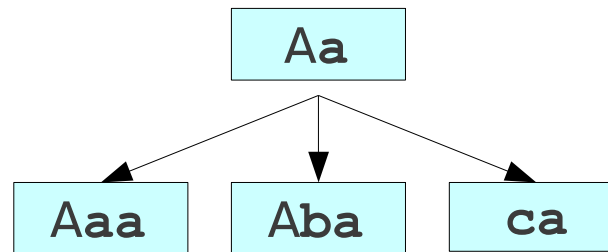
Worklist



$A \rightarrow Aa \mid Ab \mid c$

caaaaaaaaaa

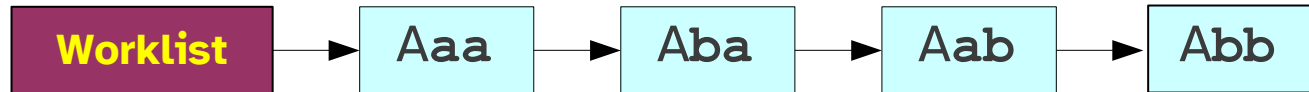
# Leftmost BFS Has Problems



**A** → **A**a | **A**b | c

caaaaaaaaaaa

# Leftmost BFS Has Problems



**A** → **A**a | **A**b | c

caaaaaaaaaa

# Leftmost DFS

- Idea: Use **depth-first** search.
- Advantages:
  - Lower memory usage: Only considers one branch at a time.
  - High performance: On many grammars, runs very quickly.
  - Easy to implement: Can be written as a set of mutually recursive functions.

# Leftmost DFS

$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

E
T + E
int + E
int + T
int + int

int + int

# Problems with Leftmost DFS

**A**  $\rightarrow$  **A**a | c

A
Aa
Aaa
Aaaa
Aaaaa



c

# Summary of Leftmost BFS/DFS

## Leftmost BFS

- Works on **all** grammars.
- Worst-case runtime is **exponential**.
- Worst-case memory usage is **exponential**.
- Rarely used in practice.

## Leftmost DFS

- Works on grammars without left recursion.
- Worst-case runtime is **exponential**.
- Worst-case memory usage is **linear**.
- Often used in a limited form as **recursive descent**.

# Recursive Descent Parsing

- Define a function for every nonterminal
- Every function works as follows
  - Find applicable production rule
  - Terminal function checks match with next input token (if no match reports error)
  - Nonterminal function calls (recursively) other functions
- If there are several applicable productions for a nonterminal, use lookahead



# Matching Tokens

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

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

# Functions for nonterminals

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

```
E() {  
    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;  
}
```

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

# Predictive Parsing

# Predictive Parsing

- The leftmost DFS/BFS algorithms are **backtracking** algorithms.
  - Guess which production to use, then back up if it doesn't work.
  - Try to match a prefix by sheer dumb luck.
- There is another class of parsing algorithms called **predictive** algorithms.
  - Based on remaining input, predict (*without backtracking*) which production to use.

# Exploiting Lookahead

- Given just the start symbol, how do you know which productions to use to get to the input program?
- Idea: Use **lookahead tokens**.
- When trying to decide which production to use, look at some number of tokens of the input to help make the decision.

# Predictive Parsing

**E**  $\rightarrow$  **T**

**E**  $\rightarrow$  **T** + **E**

**T**  $\rightarrow$  int

**T**  $\rightarrow$  (**E**)

int	+	(	int	+	int	)
-----	---	---	-----	---	-----	---

# Predictive Parsing

$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

<b>E</b>
<b>T + E</b>
int + <b>E</b>
int + <b>T</b>
int + ( <b>E</b> )
int + ( <b>T + E</b> )
int + (int + <b>E</b> )
int + (int + <b>T</b> )
int + (int + int)

int	+	(	int	+	int	)
-----	---	---	-----	---	-----	---

# A Simple Predictive Parser: LL(1)

- Top-down, predictive parsing:
  - **L**: Left-to-right scan of the tokens
  - **L**: Leftmost derivation.
  - **(1)**: One token of lookahead
- Construct a leftmost derivation for the sequence of tokens.
  - When expanding a nonterminal, we predict the production to use by looking at the next token of the input. **The decision is forced.**



# LL(1) Parse Tables

**E**  $\rightarrow$  **int**

**E**  $\rightarrow$  (**E Op E**)

**Op**  $\rightarrow$  **+**

**Op**  $\rightarrow$  **\***

	int	(	)	+	*
E	<b>E</b> $\rightarrow$ <b>int</b>	<b>E</b> $\rightarrow$ ( <b>E Op E</b> )			
Op				<b>Op</b> $\rightarrow$ <b>+</b>	<b>Op</b> $\rightarrow$ <b>*</b>

# LL(1) Parsing

(int + (int \* int))

- (1) **E** → int
- (2) **E** → (**E Op** E)
- (3) **Op** → +
- (4) **Op** → \*

Example with tables format

# LL(1) Error Detection

- (1) **E** → **int**
- (2) **E** → (**E Op E**)
- (3) **Op** → **+**
- (4) **Op** → **\***

int + int\$

	int	(	)	+	*
E	1	2			
Op				3	4

# LL(1) Error Detection

- (1) **E** → **int**
- (2) **E** → (**E Op E**)
- (3) **Op** → **+**
- (4) **Op** → **\***

E\$	int + int\$
int \$	int + int\$
\$	+ int\$

	int	(	)	+	*
E	1	2			
Op				3	4

???

# LL(1) Error Detection, Part II

- (1) **E** → **int**
- (2) **E** → (**E Op E**)
- (3) **Op** → **+**
- (4) **Op** → **\***

(int (int))\$

	int	(	)	+	*
E	1	2			
Op				3	4

# LL(1) Error Detection, Part II

- (1) **E** → **int**
- (2) **E** → (**E Op E**)
- (3) **Op** → **+**
- (4) **Op** → **\***

	int	(	)	+	*
E	1	2			
Op				3	4

E\$	(int (int))\$
(E Op E) \$	(int (int))\$
E Op E) \$	int (int))\$
int Op E) \$	int (int))\$
Op E) \$	(int))\$

???

# The LL(1) Algorithm

- Suppose a grammar has start symbol **S** and LL(1) parsing table T. We want to parse string  $\omega$
- Initialize a stack containing **S** $\$$ .
- Repeat until the stack is empty:
  - Let the next character of  $\omega$  be **t**.
  - If the top of the stack is a terminal **r**:
    - If **r** and **t** don't match, report an error.
    - Otherwise consume the character **t** and pop **r** from the stack.
  - Otherwise, the top of the stack is a nonterminal **A**:
    - If  $T[\mathbf{A}, \mathbf{t}]$  is undefined, report an error.
    - Replace the top of the stack with  $T[\mathbf{A}, \mathbf{t}]$ .

Can we find an algorithm for constructing LL(1) parse tables?



# Filling in Table Entries

- Intuition: The next character should uniquely identify a production, so we should pick a production that ultimately starts with that character.
- $T[A, t]$  should be a production  $A \rightarrow \omega$  iff  $\omega$  derives something starting with  $t$ .

In what follows, assume that our grammar does not contain any  $\varepsilon$ -productions.

We'll relax this restriction later.

# FIRST Sets

- We want to tell if a particular nonterminal **A** derives a string starting with a particular nonterminal **t**.
- We can formalize this with **FIRST sets**.
  - $\text{FIRST}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{A} \Rightarrow^* \mathbf{t}\omega \text{ for some } \omega \}$
  - Intuitively,  $\text{FIRST}(\mathbf{A})$  is the set of terminals that can be at the start of a string produced by **A**.
- If we can compute FIRST sets for all nonterminals in a grammar, we can efficiently construct the LL(1) parsing table.

# Computing FIRST Sets

- Initially, for all nonterminals **A**, set
$$\text{FIRST}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{A} \rightarrow \mathbf{t}\omega \text{ for some } \omega \}$$
- Then, repeat the following until no changes occur:
  - For each nonterminal **A**, for each production  $\mathbf{A} \rightarrow \mathbf{B}\omega$ , set
$$\text{FIRST}(\mathbf{A}) = \text{FIRST}(\mathbf{A}) \cup \text{FIRST}(\mathbf{B})$$
- This is known a **fixed-point iteration** or a **transitive closure algorithm**.

# Iterative FIRST Computations

**STMT** → if **EXPR** then **STMT**  
| while **EXPR** do **STMT**  
| **EXPR** ;

**EXPR** → **TERM** -> id  
| zero? **TERM**  
| not **EXPR**  
| ++ id  
| -- id

**TERM** → id  
| constant

STMT	EXPR	TERM

# Iterative FIRST Computations

**STMT** → if **EXPR** then **STMT**  
| while **EXPR** do **STMT**  
| **EXPR** ;

**EXPR** → **TERM** -> id  
| zero? **TERM**  
| not **EXPR**  
| ++ id  
| -- id

**TERM** → id  
| constant

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

# Iterative FIRST Computations

**STMT** → if **EXPR** then STMT  
| while **EXPR** do STMT  
| **EXPR** ;

**EXPR** → **TERM** -> id  
| zero? **TERM**  
| not **EXPR**  
| ++ id  
| -- id

**TERM** → id  
| constant

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

# Iterative FIRST Computations

STMT → if **EXPR** then STMT  
| while **EXPR** do STMT  
| **EXPR** ;

**EXPR** → **TERM** → id  
| zero? **TERM**  
| not **EXPR**  
| ++ id  
| -- id

**TERM** → id  
| constant

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



# Iterative FIRST Computations

**STMT** → if **EXPR** then STMT  
| while **EXPR** do STMT  
| **EXPR** ;

**EXPR** → **TERM** → id  
| zero? **TERM**  
| not **EXPR**  
| ++ id  
| -- id

**TERM** → id  
| constant

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

# From FIRST Sets to LL(1) Tables

```

STMT  →  if EXPR then STMT
          |  while EXPR do STMT
          |  EXPR ;

EXPR   →  TERM           -> id
          |  zero? TERM
          |  not EXPR
          |  ++ id
          |  -- id

TERM   →  id
          |  constant

```

	STMT	EXPR	TERM
(1)			
(2)			
(3)	if	zero?	id
(4)	while	not	constant
(5)	zero?	++	
(6)	not	--	
(7)	++	id	
(8)	--	constant	
(9)	id		
(10)	constant		

[illegible]

# $\epsilon$ -Free LL(1) Parse Tables

- The following algorithm constructs an LL(1) parse table for a grammar with no  $\epsilon$ -productions.
- Compute the **FIRST** sets for all nonterminals in the grammar.
  - For each production  $A \rightarrow t\omega$ , set  $T[A, t] = t\omega$ .
  - For each production  $A \rightarrow B\omega$ , set  $T[A, t] = B\omega$  for each  $t \in \text{FIRST}(B)$ .

# Exercise

- Write a CFG for integers without  $\varepsilon$ -productions
- FIRST sets?

# FIRST Sets with $\epsilon$

**Num**  $\rightarrow$  **Sign Digits**  
**Sign**  $\rightarrow$  + | - |  $\epsilon$   
**Digits**  $\rightarrow$  **Digit More**  
**More**  $\rightarrow$  **Digits** |  $\epsilon$   
**Digit**  $\rightarrow$  0 | 1 | 2 | ... | 9

Num	Sign	Digit	Digits	More

# FIRST Sets with $\epsilon$

**Num** → **Sign Digits**  
**Sign** → **+** | **-** |  $\epsilon$   
**Digits** → **Digit More**  
**More** → **Digits** |  $\epsilon$   
**Digit** → **0** | **1** | **2** | ... | **9**

Num	Sign	Digit	Digits	More
<b>+</b> <b>-</b>	<b>+</b> <b>-</b>	<b>0</b> <b>5</b>	<b>0</b> <b>5</b>	<b>0</b> <b>5</b>
		<b>1</b> <b>6</b>	<b>1</b> <b>6</b>	<b>1</b> <b>6</b>
		<b>2</b> <b>7</b>	<b>2</b> <b>7</b>	<b>2</b> <b>7</b>
		<b>3</b> <b>8</b>	<b>3</b> <b>8</b>	<b>3</b> <b>8</b>
		<b>4</b> <b>9</b>	<b>4</b> <b>9</b>	<b>4</b> <b>9</b>

# FIRST Sets with $\epsilon$

**Num**  $\rightarrow$  **Sign Digits**  
**Sign**  $\rightarrow$   $+$  |  $-$  |  $\epsilon$   
**Digits**  $\rightarrow$  **Digit More**  
**More**  $\rightarrow$  **Digits** |  $\epsilon$   
**Digit**  $\rightarrow$   $0$  |  $1$  |  $2$  | ... |  $9$

Num	Sign	Digit	Digits	More
$+$ $-$	$+$ $-$	$0$ $5$	$0$ $5$	$0$ $5$
	$\epsilon$	$1$ $6$	$1$ $6$	$1$ $6$
		$2$ $7$	$2$ $7$	$2$ $7$
		$3$ $8$	$3$ $8$	$3$ $8$
		$4$ $9$	$4$ $9$	$4$ $9$
				$\epsilon$

# FIRST Sets with $\epsilon$

**Num** → **Sign Digits**  
**Sign** → + | - |  $\epsilon$   
**Digits** → **Digit More**  
**More** → **Digits** |  $\epsilon$   
**Digit** → 0 | 1 | 2 | ... | 9

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	$\epsilon$		1	6	1	6	1	6
1	6			2	7	2	7	2	7
2	7			3	8	3	8	3	8
3	8			4	9	4	9	4	9
4	9							$\epsilon$	



# FIRST and $\varepsilon$

- When computing FIRST sets in a grammar with  $\varepsilon$ -productions, we often have to “look through” nonterminals.
- Rationale: Might have a derivation like this:

$$\mathbf{A} \Rightarrow \mathbf{B}t \Rightarrow t$$

So  $t \in \text{FIRST}(\mathbf{A})$ .

# FIRST Computation with $\epsilon$

- Initially, for all nonterminals  $A$ , set
$$\text{FIRST}(A) = \{ t \mid A \rightarrow t\omega \text{ for some } \omega \}$$
- For all nonterminals  $A$  where  $A \rightarrow \epsilon$  is a production, add  $\epsilon$  to  $\text{FIRST}(A)$ .
- Repeat the following until no changes occur:
  - For each production  $A \rightarrow \alpha$ , where  $\alpha$  is a string of nonterminals whose FIRST sets all contain  $\epsilon$ , set  $\text{FIRST}(A) = \text{FIRST}(A) \cup \{ \epsilon \}$ .
  - For each production  $A \rightarrow \alpha t \omega$ , where  $\alpha$  is a string of nonterminals whose FIRST sets all contain  $\epsilon$ , set  $\text{FIRST}(A) = \text{FIRST}(A) \cup \{ t \}$
  - For each production  $A \rightarrow \alpha B \omega$ , where  $\alpha$  is string of nonterminals whose FIRST sets all contain  $\epsilon$ , set  $\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(B) - \{ \epsilon \})$ .

# A Notational Diversion

- Once we have computed the correct FIRST sets for each nonterminal, we can generalize our definition of FIRST sets to strings.
- Define  $\text{FIRST}^*(\omega)$  as follows:
  - $\text{FIRST}^*(\epsilon) = \{ \epsilon \}$
  - $\text{FIRST}^*(t\omega) = \{ t \}$
  - If  $\epsilon \notin \text{FIRST}(\mathbf{A})$ :
    - $\text{FIRST}^*(\mathbf{A}\omega) = \text{FIRST}(\mathbf{A})$
  - If  $\epsilon \in \text{FIRST}(\mathbf{A})$ :
    - $\text{FIRST}^*(\mathbf{A}\omega) = (\text{FIRST}(\mathbf{A}) - \{ \epsilon \}) \cup \text{FIRST}^*(\omega)$

# FIRST Computation with $\varepsilon$

- Initially, for all nonterminals **A**, set
$$\text{FIRST}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{A} \rightarrow \mathbf{t}\omega \text{ for some } \omega \}$$
- For all nonterminals **A** where  $\mathbf{A} \rightarrow \varepsilon$  is a production, add  $\varepsilon$  to  $\text{FIRST}(\mathbf{A})$ .
- Repeat the following until no changes occur:
  - For each production  $\mathbf{A} \rightarrow \alpha$ , set
$$\text{FIRST}(\mathbf{A}) = \text{FIRST}(\mathbf{A}) \cup \text{FIRST}^*(\alpha)$$

# LL(1) Tables with $\epsilon$

**Num** → **Sign Digits**

**Sign** → **+** | **-** |  **$\epsilon$**

**Digits** → **Digit More**

**More** → **Digits** |  **$\epsilon$**

**Digit** → **0** | **1** | **2** | ... | **9**

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	$\epsilon$		1	6	1	6	1	6
1	6			2	7	2	7	2	7
2	7			3	8	3	8	3	8
3	8			4	9	4	9	4	9
4	9							$\epsilon$	

	+	-	#
Num			
Sign			
Digits			
More			
Digit			

# LL(1) Tables with $\epsilon$

**Num** → **Sign Digits**  
**Sign** → **+** | **-** |  **$\epsilon$**   
**Digits** → **Digit More**  
**More** → **Digits** |  **$\epsilon$**   
**Digit** → **0** | **1** | **2** | ... | **9**

Num \$	- 5 \$
Sign Digits \$	- 5 \$
- Digits \$	- 5 \$
Digits \$	5 \$
Digit More \$	5 \$
5 More \$	5 \$
More \$	\$

	+	-	#
Num			
Sign			
Digits			
More			
Digit			

?????

# $\epsilon$ is Complicated

- When constructing LL(1) tables with  $\epsilon$ -productions, we need to have an extra column for  $\$$ .

	+	-	#	\$
Num				
Sign				
Digits				
More				
Digit				

# It Gets Trickier

**Num** → **Sign Digits**  
**Sign** → **+** | **-** | **ε**  
**Digits** → **Digit More**  
**More** → **Digits** | **ε**  
**Digit** → **0** | **1** | **2** | ... | **9**

Num \$	5 \$
Sign Digits \$	5 \$

?????
-------

	+	-	#	\$
Num				
Sign				
Digits				
More				
Digit				



# FOLLOW Sets

- With  $\epsilon$ -productions in the grammar, we may have to “look past” the current nonterminal to what can come after it.
- The **FOLLOW set** represents the set of terminals that might come after a given **nonterminal**. Formally:
  - $\text{FOLLOW}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{S} \Rightarrow^* \alpha \mathbf{A} \mathbf{t} \omega \text{ for some } \alpha, \omega \}$   
where  $\mathbf{S}$  is the start symbol of the grammar.
- Informally, every terminal that can ever come after  $\mathbf{A}$  in a derivation.

# Computation of FOLLOW Sets

- Initially, for each nonterminal **A**, set  
$$\text{FOLLOW}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{B} \rightarrow \alpha \mathbf{A} \mathbf{t} \omega \text{ is a production} \}$$
- Add **\$** to  $\text{FOLLOW}(\mathbf{S})$ , where **S** is the start symbol.
- Repeat the following until no changes occur:
  - If  $\mathbf{B} \rightarrow \alpha \mathbf{A} \omega$  is a production, set  
$$\text{FOLLOW}(\mathbf{A}) = \text{FOLLOW}(\mathbf{A}) \cup \text{FIRST}^*(\omega) - \{ \epsilon \}.$$
  - If  $\mathbf{B} \rightarrow \alpha \mathbf{A} \omega$  is a production and  $\epsilon \in \text{FIRST}^*(\omega)$ , set  
$$\text{FOLLOW}(\mathbf{A}) = \text{FOLLOW}(\mathbf{A}) \cup \text{FOLLOW}(\mathbf{B}).$$

# The Final LL(1) Table Algorithm

- Compute  $\text{FIRST}(\mathbf{A})$  and  $\text{FOLLOW}(\mathbf{A})$  for all nonterminals  $\mathbf{A}$ .
- For each rule  $\mathbf{A} \rightarrow \omega$ , for each terminal  $\mathbf{t} \in \text{FIRST}^*(\omega)$ , set  $T[\mathbf{A}, \mathbf{t}] = \omega$ .
  - Note that  $\epsilon$  is not a terminal.
- For each rule  $\mathbf{A} \rightarrow \omega$ , if  $\epsilon \in \text{FIRST}^*(\omega)$ , set  $T[\mathbf{A}, \mathbf{t}] = \omega$  for each  $\mathbf{t} \in \text{FOLLOW}(\mathbf{A})$ .

# Final LL(1) Table

**Num** → **Sign Digits**  
**Sign** → **+** | **-** |  $\epsilon$   
**Digits** → **Digit More**  
**More** → **Digits** |  $\epsilon$   
**Digit** → **0** | **1** | ... | **9**

FIRST				
Num	Sign	Digit	Digits	More
<b>+</b> <b>-</b> <b>#</b>	<b>+</b> <b>-</b> $\epsilon$	<b>#</b>	<b>#</b>	<b>#</b> $\epsilon$

FOLLOW				
Num	Sign	Digit	Digits	More
<b>\$</b>	<b>#</b>	<b>#</b> <b>\$</b>	<b>\$</b>	<b>\$</b>

	<b>+</b>	<b>-</b>	<b>#</b>	<b>\$</b>
Num	<b>Sign Digits</b>	<b>Sign Digits</b>	<b>Sign Digits</b>	
Sign	<b>+</b>	<b>-</b>	$\epsilon$	
Digits			<b>Digits More</b>	
More			<b>Digits</b>	$\epsilon$
Digit			<b>#</b>	

# Exercise

**Num** → **Digits** **Frac**  
**Digits** → **Digit** **More**  
**Frac** → **.** **Digits** |  $\epsilon$   
**More** → **Digits** |  $\epsilon$   
**Digit** → **0** | **1** | **2** | ... | **9**

FIRST				
Num	Digits	Frac	More	Digit

FOLLOW				
Num	Digits	Frac	More	Digit

	.	#	\$
Num			
Digits			
Frac			
More			
Digit			

# The Limits of LL(1)

# A Grammar that is Not LL(1)

- Consider the following (left-recursive) grammar:
  - $A \rightarrow Ab \mid c$
  - $\text{FIRST}(A) = \{c\}$
- However, we cannot build a valid LL(1) parse table.
- **Why?**

# Eliminating Left Recursion

- In general, left recursion can be converted into **right recursion** by a mechanical transformation.
- Consider the grammar

$$\mathbf{A} \rightarrow \mathbf{A}\omega \mid \alpha$$

- This will produce  $\alpha$  followed by some number of  $\omega$ 's.
- Can rewrite the grammar as

$$\mathbf{A} \rightarrow \alpha \mathbf{B}$$

$$\mathbf{B} \rightarrow \epsilon \mid \omega \mathbf{B}$$



# Algorithm to Remove Left Recursion

- Arrange  $n$  nonterminals in some order
- **for**  $i: 1 \dots n$  {  
    **for**  $j: 1 \dots i-1$  {  
        ‘inline’ all  $A_j \rightarrow \delta$  into  $A_i \rightarrow A_j \gamma$   
        i.e., will have by  $A_i \rightarrow \delta \gamma$  productions  
    }  
    eliminate immediate left recursion among  $A_i$  productions  
}
- Example:  
     $S \rightarrow A a \mid b$   
     $A \rightarrow A c \mid S d \mid \epsilon$

Note: this CFG is not LL(1).  
Removing left recursion doesn't make it LL(1)

# Another Non-LL(1) Grammar

- Consider the following grammar:

$E \rightarrow T$

$E \rightarrow T + E$

$T \rightarrow \text{int}$

$T \rightarrow (E)$

- $\text{FIRST}(E) = \{ \text{int}, ( \}$
- $\text{FIRST}(T) = \{ \text{int}, ( \}$
- Why is this grammar not LL(1)?

# Left-Factoring

- A grammar of the form

$$\mathbf{A} \rightarrow \alpha \beta \mid \alpha \gamma \mid \omega$$

Can be transformed as

$$\mathbf{A} \rightarrow \alpha \mathbf{B} \mid \omega$$

$$\mathbf{B} \rightarrow \beta \mid \gamma$$

# Left-Factoring

$E \rightarrow T$   
 $E \rightarrow T + E$   
 $T \rightarrow \text{int}$   
 $T \rightarrow (E)$

# A Formal Characterization of LL(1)

- A grammar  $G$  is LL(1) iff for any productions  $A \rightarrow \omega_1$  and  $A \rightarrow \omega_2$ , the sets
- $\text{FIRST}(\omega_1)$  and  $\text{FIRST}(\omega_2)$  are disjoint
- and**
- if  $\epsilon \in \text{FIRST}(\omega_1)$ :  $\text{FIRST}(\omega_2)$  and  $\text{FOLLOW}(A)$  are disjoint.
- Likewise, if  $\epsilon \in \text{FIRST}(\omega_2)$

This condition is equivalent to saying that there are no conflicts in the table.

# Exercise: Are they LL(1)? If not, fix

- $\text{term} = \text{id} \mid \text{indexed\_elem} .$   
 $\text{indexed\_elem} = \text{id} \text{ “[” expr “]” } .$
- $S \rightarrow A a b$   
 $A \rightarrow a \mid \varepsilon$
- $E \rightarrow E - T \mid T$

# The Strengths of LL(1)

# LL(1) is Straightforward

- Can be implemented quickly with a table- driven design.
- Can be implemented by **recursive descent**:
  - Define a function for each nonterminal.
  - Have these functions call each other based on the lookahead token.



# LL(1) & Ambiguity

- Ground truth: **if a grammar is LL(1), it is not ambiguous**
- Consider:  
**if a grammar is not LL(1), it is ambiguous**  
and  
**if a grammar is ambiguous, it is not LL(1)**
- Example:  
$$\mathbf{S} \rightarrow \mathbf{A} x y \mid \mathbf{B} x z$$
$$\mathbf{A} \rightarrow a$$
$$\mathbf{B} \rightarrow a$$

# Summary

- **Top-down parsing** tries to derive the user's program from the start symbol.
- **Leftmost BFS** is one approach to top-down parsing; it is mostly of theoretical interest.
- **Leftmost DFS** is another approach to top-down parsing that is uncommon in practice.
- **LL(1)** parsing scans from left-to-right, using one token of lookahead to find a leftmost derivation.
- **FIRST sets** contain terminals that may be the first symbol of a production.
- **FOLLOW sets** contain terminals that may follow a nonterminal in a production.
- **Left recursion** and **left factorability** cause LL(1) to fail and can be mechanically eliminated in some cases.

# Recap: FIRST calculation

- $\text{FIRST}(\omega)$  is the set of **terminals** that begin all strings given by  $\omega$ , including  $\epsilon$  iff  $\omega \Rightarrow \epsilon$
- $\omega$  is either:
  - terminal  $a$   $\text{FIRST}(a) = \{a\}$
  - nonterminal  $X$ 
    - repeat
      - for each rule  $X \rightarrow Y_1 Y_2 \dots Y_k$ 
        - Add  $a$  to  $\text{FIRST}(X)$  if  $a \in \text{FIRST}(Y_1)$  **or**  $a \in \text{FIRST}(Y_n)$  and  $Y_1 \dots Y_{n-1} \Rightarrow^* \epsilon$
        - If  $Y_1 \dots Y_k \Rightarrow^* \epsilon$  then add  $\epsilon$  to  $\text{FIRST}(X)$
    - until no changes
- sentential form  $\alpha$ 
  - for each symbol  $Y_1 Y_2 \dots Y_k$  in  $\alpha$ 
    - Add  $a$  to  $\text{FIRST}(\alpha)$  if  $a \in \text{FIRST}(Y_1)$  **or**  $a \in \text{FIRST}(Y_n)$  and  $Y_1 \dots Y_{n-1} \Rightarrow^* \epsilon$
    - If  $Y_1 \dots Y_k \Rightarrow^* \epsilon$  then add  $\epsilon$  to  $\text{FIRST}(\alpha)$