



COSC252: Programming Languages Lecture: Parsing

Jeremy Bolton, PhD
Asst Professor

Outline

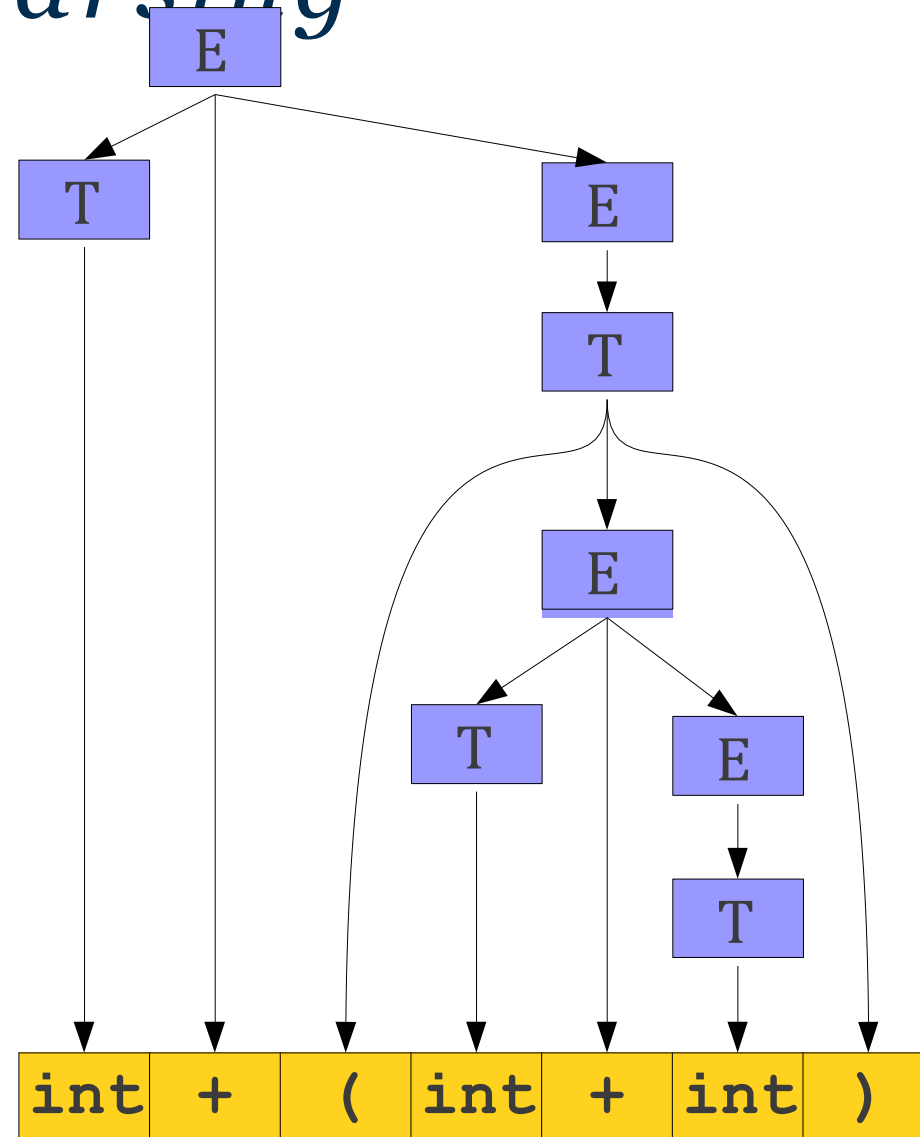
- I. Parsing Intro
- II. Top-Down
 - I. Method
 - I. Issues
 - II. Fixes
 - II. Recursive decent
- III. Table Implementation

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

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



Challenges in Top-Down Parsing

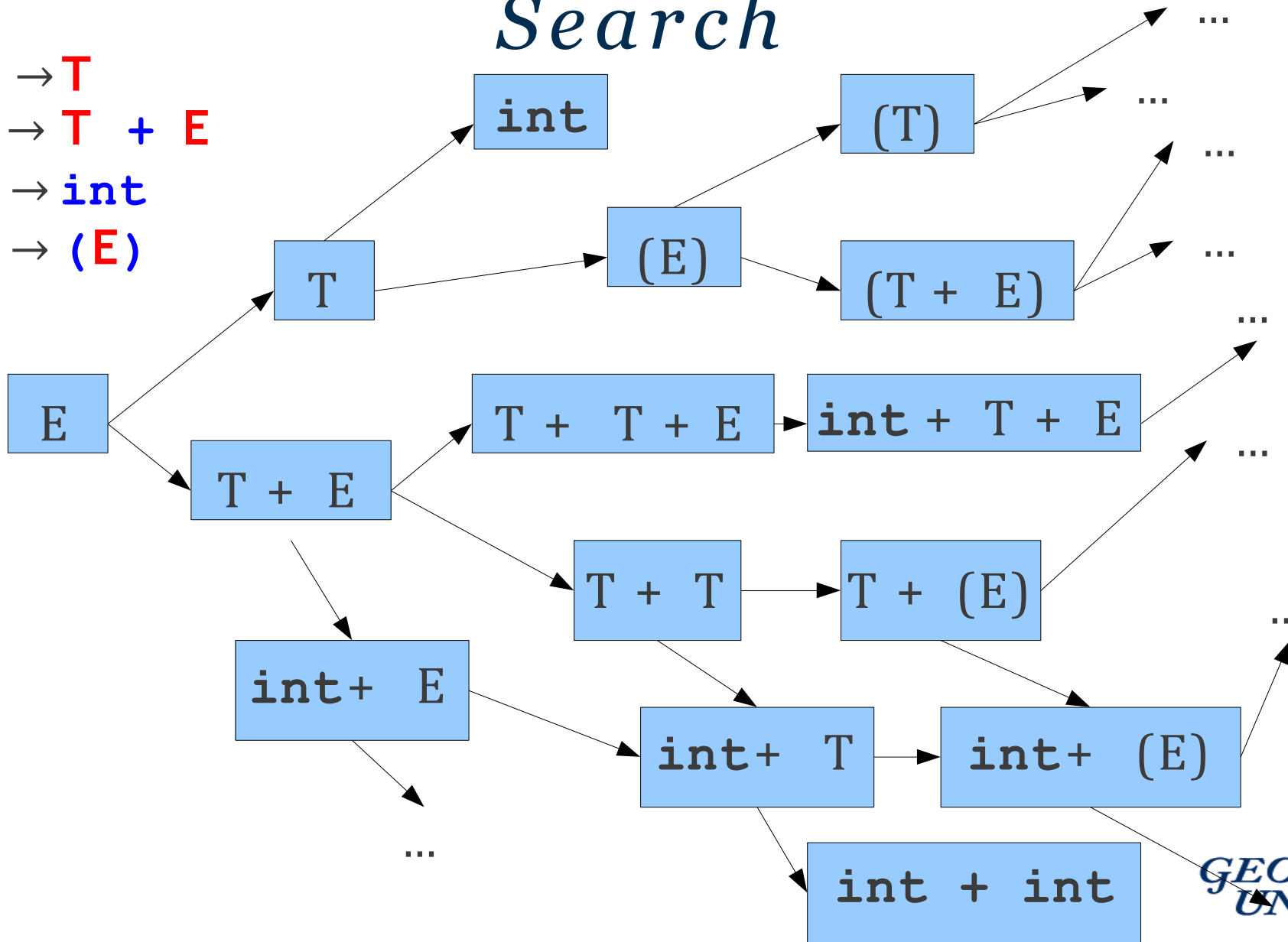
- Top-down parsing begins with virtually no information.
 - Begins with just the start symbol, which matches *every* program.
- How can we know which productions to apply?
- In general, we can't.
 - There are some grammars for which the best we can do is guess and backtrack if we're wrong.
 - If we have to guess, how do we do it?

Parsing as a Search

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

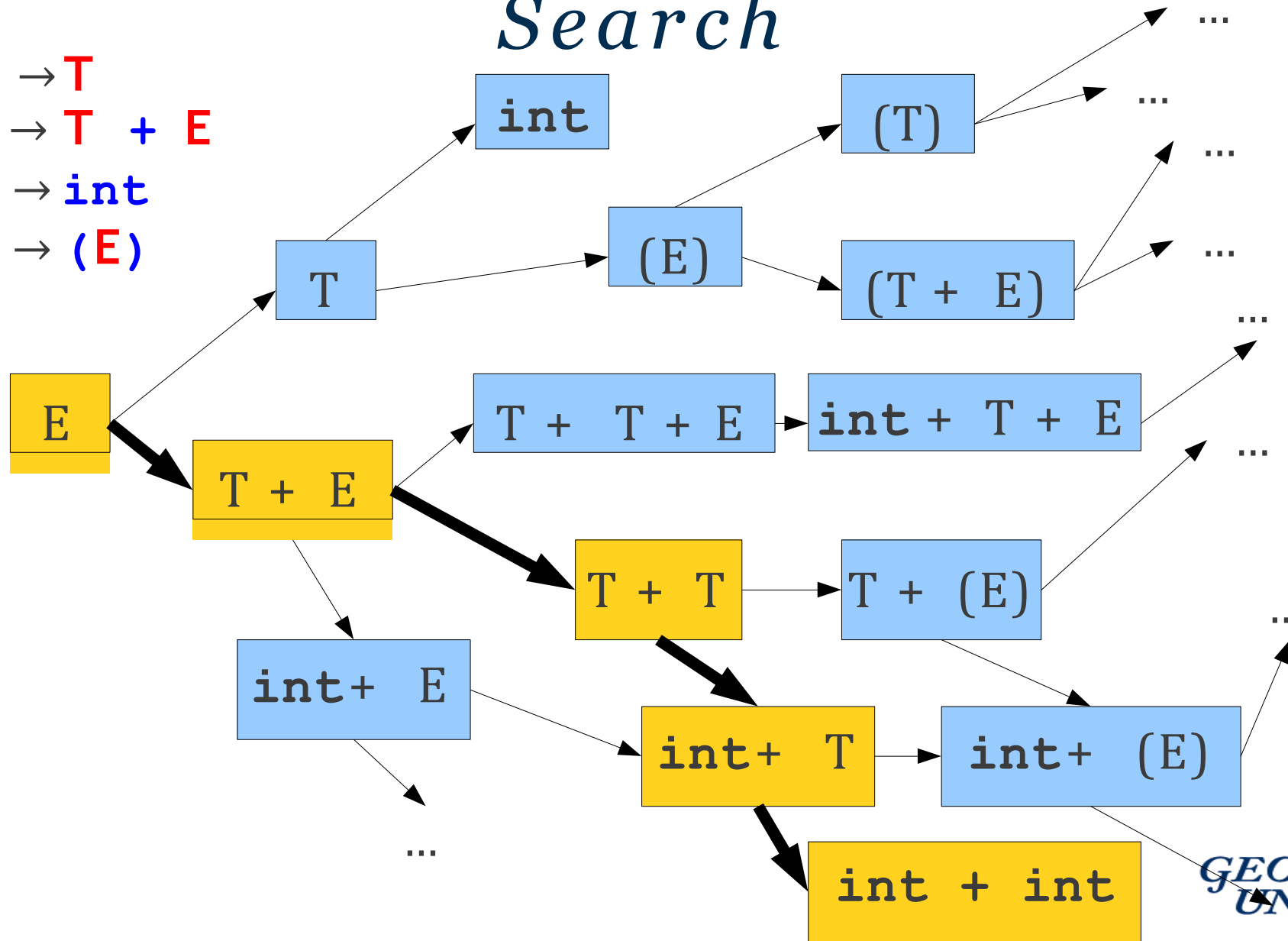
Parsing as a Search

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



Parsing as a Search

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



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.

Our First Top-Down Algorithm

Breadth-First Search

Bad idea ... trust me!

derived in one step, add that string to the worklist.

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

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

Leftmost DFS

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

`int + int`

Leftmost DFS

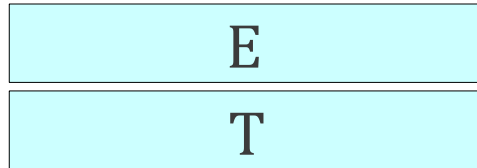
E

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

int + int

Leftmost DFS

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



`int + int`

Leftmost DFS

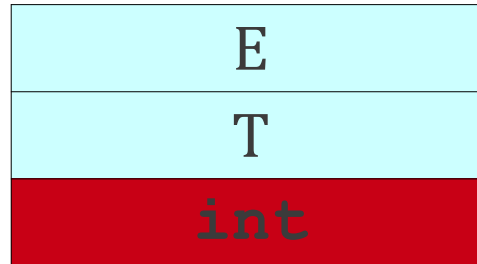
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

E
T
int

int + int

Leftmost DFS

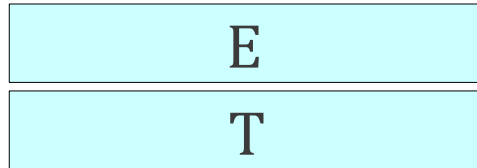
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



int + int

Leftmost DFS

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



`int + int`

Leftmost DFS

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

E
T
(E)

`int + int`

Leftmost DFS

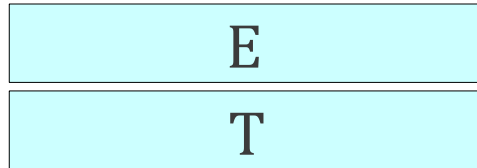
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

E
T
(E)

int + int

Leftmost DFS

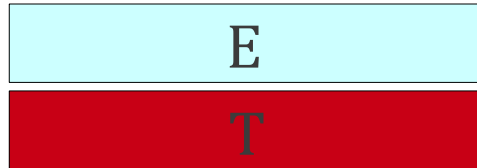
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

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



`int + int`

Leftmost DFS

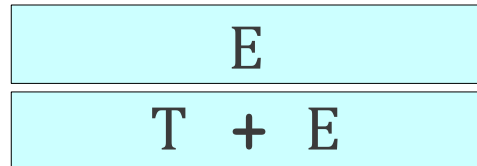
E

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

`int + int`

Leftmost DFS

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



int + int

Leftmost DFS

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

E
T + E
int + E

int + int

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

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

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 → **A**a | c

A
Aa
Aaa
Aaaa
Aaaaa



c

Left Recursion

- A nonterminal **A** is said to be **left-recursive** iff

$$A \Rightarrow^* A\omega$$

for some string ω .

- Leftmost DFS may fail on left-recursive grammars.
- Fortunately, in many cases it is possible to eliminate left recursion (see Handout 08 for details).

Notes concerning 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**.

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.

Tradeoffs in Prediction

- Predictive parsers are *fast*.
 - Many predictive algorithms can be made to run in linear time.
 - Often can be table-driven for extra performance.
- Predictive parsers are *weak*.
 - Not all grammars can be accepted by predictive parsers.
- Trade *expressiveness* for *speed*.

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.

Implementing Predictive Parsing

- Predictive parsing is only possible if we can predict which production to use given some number of lookahead tokens.
- Increasing the number of lookahead tokens increases the number of grammars we can parse, but complicates the parser.
- Decreasing the number of lookahead tokens decreases the number of grammars we can parse, but simplifies the parser.

Predictive Parsing

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

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

Predictive Parsing

E

$E \rightarrow T$

$E \rightarrow T + E$

$T \rightarrow \text{int}$

$T \rightarrow (E)$

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

Predictive Parsing

E

$E \rightarrow T$

$E \rightarrow T + E$

$T \rightarrow \text{int}$

$T \rightarrow (E)$

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

Predictive Parsing

E
T + E

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

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

Predictive Parsing

E
T + E

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

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

Predictive Parsing

E
T + E

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

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

Predictive Parsing

E
T + E
int + E

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

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

Predictive Parsing

E
T + E
int + E

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

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

Predictive Parsing

E
T + E
int + E

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

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

Predictive Parsing

E
T + E
int + E

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

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

Predictive Parsing

E
T + E
int + E
int + T

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

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

Predictive Parsing

E
T + E
int + E
int + T

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

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

Predictive Parsing

E
T + E
int + E
int + T

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

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)

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

Predictive Parsing

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

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

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)

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

Predictive Parsing

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

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

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)
int + (int + T)

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

Predictive Parsing

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

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

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

Predictive Parsing

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

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

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

Predictive Parsing

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

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

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)
int + (int + T)
int + (int + int)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)
int + (int + T)
int + (int + int)

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

Predictive Parsing

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

E
T + E
int + E
int + T
int + (E)
int + (T + E)
int + (int + E)
int + (int + T)
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.**

Common LL Parser Implementations

- Use stack to maintain non-terminals encountered which still need to be resolved.
 - Table Driven Parser
 - Recursive descent parser

LL(1) Parse Tables

LL(1) Parse Tables

E → **int**

E → (**E Op E**)

Op → **+**

Op → *****

LL(1) Parse Tables

E → **int**

E → (**E Op E**)

Op → **+**

Op → *****

	int	()	+	*
E	int	(E Op E)			
Op				+	*

LL(1) Parsing

(int + (int * int))

(1) **E** → **int**

(2) **E** → (**E Op E**)

(3) **Op** → **+**

(4) **Op** → *****

LL(1) Parsing

E

(int + (int * int))

(1) **E** → **int**

(2) **E** → (**E Op E**)

(3) **Op** → **+**

(4) **Op** → *****

LL(1) Parsing

E	(int + (int * int))
---	---------------------

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

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

E\$

(int + (int * int))\$

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

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$ (int + (int * int))\$

The **\$** symbol is the end-of-input marker and is used by the parser to detect when we have reached the end of the input. It is not a part of the grammar.

LL(1) Parsing

E\$

(int + (int * int))\$

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

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

E\$

(int + (int * int))\$

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

	int	()	+	*
E	1	2			
Op				3	4

The first symbol of our guess is a nonterminal. We then look at our parsing table to see what production to use.

This is called a **predict** step.

LL(1) Parsing

E\$

(int + (int * int))\$

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

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$

The first symbol of our guess is now a terminal symbol. We thus match it against the first symbol of the string to parse.

This is called a **match** step.

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$

	int	()	+	*
E	1	2			
Op				3	4

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$
int))\$	int))\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$
int))\$	int))\$
)\$)\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$
int))\$	int))\$
)\$)\$
)\$)\$

LL(1) Parsing

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

	int	()	+	*
E	1	2			
Op				3	4

E\$	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
+ E)\$	+ (int * int))\$
E)\$	(int * int))\$
(E Op E))\$	(int * int))\$
E Op E))\$	int * int))\$
int Op E))\$	int * int))\$
Op E))\$	* int))\$
* E))\$	* int))\$
E))\$	int))\$
int))\$	int))\$
)\$)\$
)\$)\$
\$	\$

LL(1) Parsing

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

	int	()
E	1	2	
Op			



E	(int + (int * int))\$
(E Op E)\$	(int + (int * int))\$
E Op E)\$	int + (int * int))\$
int Op E)\$	int + (int * int))\$
Op E)\$	+ (int * int))\$
	+ (int * int))\$
	(int * int))\$
	(int * int))\$
	int * int))\$
	int * int))\$
	* int))\$
	* int))\$
	int))\$
int))\$	int))\$
)\$)\$
)\$)\$
\$	\$

Parse Errors

- Arriving at an empty location in the parse table indicates a parse error.

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** → *****

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

	int	()	+	*
E	1	2			
Op				3	4

Parsing Algorithms

- Algorithms: LL(k)
- There are different LL algorithms for different k
 - Table construction
 - Table usage

The LL(1) Algorithm

- Suppose a grammar has start symbol **S** and LL(1) parsing table T. We want to parse string ω
- Initialize a stack containing **S**\$.
- Repeat until the stack is empty:
 - Let the next character of ω 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[**A**, **t**] is undefined, report an error.
 - Replace the top of the stack with T[**A**, **t**].

But how do we build the table??

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

ϵ -Free LL(1) Parse Tables

- The following algorithm constructs an LL(1) parse table for a grammar with no ϵ -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 $t \in \text{FIRST}(B)$.

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 ω derives something starting with t .
- More rigorously:

$$T[A, t] = B\omega \text{ iff } A \rightarrow \omega \text{ and } \omega \Rightarrow^* t\omega'$$

LL(1) Parse Tables

E → **int**

E → (**E Op E**)

Op → **+**

Op → *****

	int	()	+	*
E	int	(E Op E)			
Op				+	*

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**.

$$\text{FIRST}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{A} \Rightarrow^* \mathbf{t}\omega \text{ for some } \omega \}$$

- We also include **ε** in $\text{FIRST}(\mathbf{A})$ if A can produce the empty string.
- Intuitively, $\text{FIRST}(\mathbf{A})$ is the set of terminals that can be at the start of a string produced by **A**.
- We can generalize FIRST to strings with $\text{FIRST}^*(\omega)$ being the set of all terminals (or **ε**) that can appear at the start of a string derived from **ω**.

FIRST Computation with ϵ

- 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 ϵ to $\text{FIRST}(A)$.
- Repeat the following until no changes occur:
 - For each production $A \rightarrow a$, set
$$\text{FIRST}(A) = \text{FIRST}(A) \cup \text{FIRST}^*(a)$$

LL(1) Tables with ϵ

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

LL(1) Tables with ϵ

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

	+	-	#	\$
Num				
Sign				
Digits				
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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								ϵ

	+	-	#	\$
Num				
Sign				
Digits				
More				
Digit				

LL(1) Tables with ϵ

Num → **Sign Digits**

Sign → + | - | ϵ

Digits → **Digit More**

More → **Digits** | ϵ

Digit → 0 | 1 | ... | 9

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num				
Sign				
Digits				
More				
Digit				

LL(1) Tables with ϵ

Num	→ Sign Digits	Num	Sign	Digit	Digits	More
Sign	→ + - ϵ	+ -	+ -	0 5	0 5	0 5
Digits	→ Digit More	0 5	ϵ	1 6	1 6	1 6
More	→ Digits ϵ	1 6		2 7	2 7	2 7
Digit	→ 0 1 ... 9	2 7		3 8	3 8	3 8
		3 8		4 9	4 9	4 9
		4 9				ϵ

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign				
Digits				
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign				
Digits				
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits				
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits				
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More				
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit				

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

Num	→ Sign Digits	Num	Sign	Digit	Digits	More
Sign	→ + - ϵ	+ -	+ -	0 5	0 5	0 5
Digits	→ Digit More	0 5	ϵ	1 6	1 6	1 6
More	→ Digits ϵ	1 6		2 7	2 7	2 7
Digit	→ 0 1 ... 9	2 7		3 8	3 8	3 8
		3 8		4 9	4 9	4 9
		4 9				ϵ

	+	-	#	\$
Num	Sign Digits	Sign Digits		
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

Num → **Sign Digits**

Sign → **+** | **-** | ϵ

Digits → **Digit More**

More → **Digits** | ϵ

Digit → **0** | **1** | ... | **9**

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-		
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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								ϵ

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	ϵ
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	ϵ
Digit			#	

LL(1) Tables with ϵ

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

Num		Sign		Digit		Digits		More	
+	-	+	-	0	5	0	5	0	5
0	5	ϵ		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							ϵ	

	+	-	#	\$
Num	Sign Digits	Sign Digits	Sign Digits	
Sign	+	-	ϵ	
Digits			Digits More	
More			Digits	ϵ
Digit			#	

FIRST and ϵ

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

$$A \Rightarrow Bt \Rightarrow t$$

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

FIRST Computation with ϵ

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

FIRST Computation with ϵ

- 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 ϵ to $\text{FIRST}(A)$.
- Repeat the following until no changes occur:
 - For each production $A \rightarrow \alpha$, set
$$\text{FIRST}(A) = \text{FIRST}(A) \cup \text{FIRST}^*(\alpha)$$

FOLLOW Sets

- With ϵ -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:

$$\text{FOLLOW}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{S} \Rightarrow^* \mathbf{aAt}\omega \text{ for some } \mathbf{a}, \omega \}$$

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

- Informally, every nonterminal that can ever come after **A** in a derivation.

Computation of FOLLOW Sets

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

An Egregious Abuse of Notation

- 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 $\mathbf{t} \in \text{FIRST}^*(\omega \text{FOLLOW}(\mathbf{A}))$, set $\text{T}[\mathbf{A}, \mathbf{t}] = \omega$.

Example
LL(1)
Construction

*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 an LL(1) parse table.
- **Why?**

A Grammar that is Not $LL(1)$

- Consider the following (left-recursive) grammar:

$$A \rightarrow Ab \mid c$$

- $FIRST(A) = \{c\}$
- However, we cannot build an $LL(1)$ parse table.
- Why?

	b	c
A		$A \rightarrow Ab$ $A \rightarrow c$

A Grammar that is Not $LL(1)$

- Consider the following (left-recursive) grammar:

$$A \rightarrow Ab \mid c$$

- $FIRST(A) = \{c\}$
- However, we cannot build an $LL(1)$ parse table.
- **Why?**

	b	c
A		$A \rightarrow Ab$ $A \rightarrow c$

- **Cannot uniquely predict production!**
- This is called a **FIRST/FIRST conflict**.

Eliminating Left Recursion

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

$$A \rightarrow A\omega \mid a$$

- This will produce a followed by some number of ω 's.
- Can rewrite the grammar as

$$A \rightarrow a B$$

$$B \rightarrow \epsilon \mid \omega B$$

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)?

Another Non-LL(1) Grammar

- Consider the following grammar:

$E \rightarrow T$

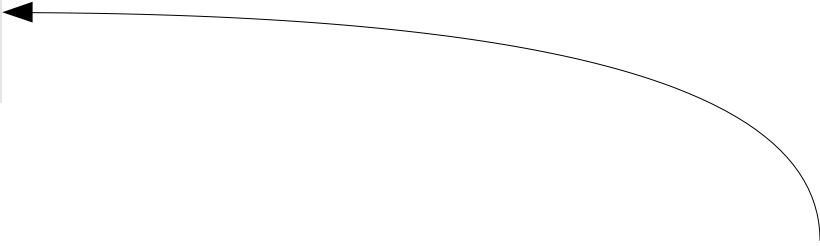
$E \rightarrow T + E$

$T \rightarrow \text{int}$

$T \rightarrow (E)$

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

How do you
predict which of
these to use?



Left- Factoring

E	\rightarrow	T
E	\rightarrow	$T + E$
T	\rightarrow	int
T	\rightarrow	(E)

Left- Factoring

E	\rightarrow	$T\epsilon$
E	\rightarrow	$T + E$
T	\rightarrow	int
T	\rightarrow	(E)

Left- Factoring

E → **TY**

T → **int**

T → **(E)**

Left- Factoring

E → **TY**

T → **int**

T → **(E)**

Y → **+E**

Y → **ε**

Left- Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+E	4
Y	\rightarrow	ϵ	5

Left- Factoring

E	→	TY	1
T	→	int	2
T	→	(E)	3
Y	→	+E	4
Y	→	ε	5

Left- Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
	FOLLOW	
E	T	Y

Left- Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
	int (
	FOLLOW	
E	T	Y

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

FIRST

E	T	Y
	int	+
	(ε

FOLLOW

E	T	Y

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y
\$		

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y
\$		
)		

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y
\$	+	
)		

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y
\$	+	\$
))

Left- Factoring

E → TY	1
T → int	2
T → (E)	3
Y → + E	4
Y → ε	5

	FIRST	
E	T	Y
int	int	+
((ε
	FOLLOW	
E	T	Y
\$	+	\$
)	\$)
)	

Left-Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
int	int	+
((ϵ
E	T	Y
\$	+	\$
)	\$)
))	

	int	()	+	\$
E					
T					
Y					

Left- Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
int	int	+
((ϵ
E	T	Y
\$	+	\$
)	\$)
))	

	int	()	+	\$
E	1	1			
T					
Y					

Left-Factoring

E	→	TY	1
T	→	int	2
T	→	(E)	3
Y	→	+ E	4
Y	→	ε	5

	FIRST	
E	T	Y
int	int	+
((ε
E	T	Y
\$	+	\$
)	\$)
))	

	int	()	+	\$
E	1	1			
T	2	3			
Y					

Left-Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
int	int	+
((ϵ
E	T	Y
\$	+	\$
)	\$)
))	

	int	()	+	\$
E	1	1			
T	2	3			
Y				4	

Left-Factoring

E	\rightarrow	TY	1
T	\rightarrow	int	2
T	\rightarrow	(E)	3
Y	\rightarrow	+ E	4
Y	\rightarrow	ϵ	5

	FIRST	
E	T	Y
int	int	+
((ϵ
E	T	Y
\$	+	\$
)	\$)
))	

	int	()	+	\$
E	1	1			
T	2	3			
Y			5	4	5

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

$$\text{FIRST}(\omega_1 \text{ FOLLOW}(A))$$

and

$$\text{FIRST}(\omega_2 \text{ FOLLOW}(A))$$

are disjoint.

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

*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) is Fast

- Both table-driven LL(1) and recursive-descent-powered LL(1) are fast.
- Can parse in $O(n |G|)$ time, where n is the length of the string and $|G|$ is the size of the grammar.

Recursive Descent Parsing



GEORGETOWN UNIVERSITY

Parsing

- Goals
 - Find syntax errors
 - Produce or trace parse tree
- Types of parsers
 - Top-down: build tree from top down
 - Bottom-up: build tree from bottom up

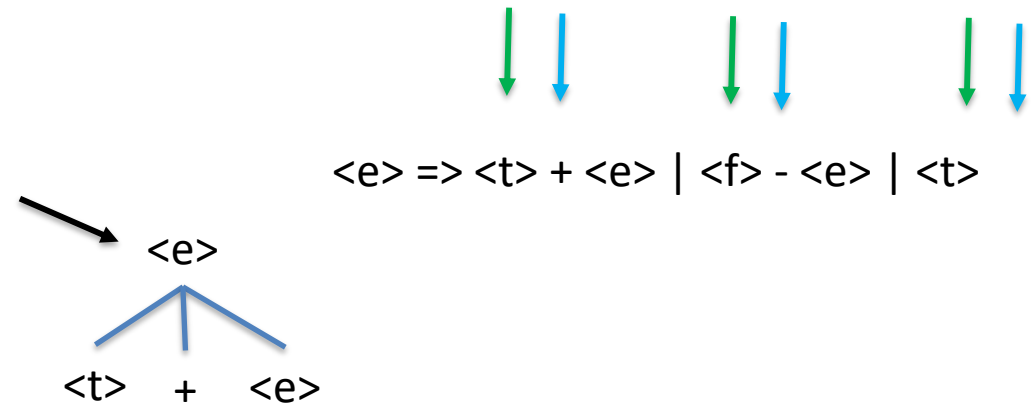
Top-Down Parsers

- Produces a left-most derivation
- Traces or builds tree in pre-order
- Crux: Must choose / **predict** correct RHS when expanding a node in the parse tree.

– Eg: 3+4



- For efficient design best to use next token or **lookahead** only one token



Top Down Parsing

- Formally, a top down parser must do the following:
 - Given a sentence in the form $xA\alpha$, the parser must choose the correct $A \Rightarrow$ rule .
 - Notation
 - Lowercase: terminal sequence
 - Uppercase: non-terminal
 - Greek: sequence of terminals and non-terminals
 - Its best if this decision can be made based on the first tokens of A's RHSs.
 - This is true if all first terminals of A productions are different.
- Implementation of Top Down Parser
 - Recursive decent
 - Table driven

Building a Recursive Descent Parser:

Must remove left recursion. Two options shown here.

```
/* term
Parses strings in the language generated by the rule:
<term> -> <term> * <factor> | ...
*/
void term() {

/* Parse the first factor */
term();

/* As long as the next token is * or /,
next token and parse the next factor */
if(nextToken == MULT_OP || nextToken == DIV_OP) {
    MultDiv();
    factor();
}
} /* End of function term */
```

Remove Left Recursion: Recursive Descent Parser Option 1

```
/* term
Parses strings in the language generated by the rule:
<term> -> <factor> <term_prime>
<term_prime> -> * <factor> <term_prime> | epsilon

*/
void term() {

/* Parse the first factor */
factor();

/* Lookahead: As long as the next token is * or /,
next token and parse the next factor */
if(nextToken == MULT_OP || nextToken == DIV_OP) {
    MultDiv(); // match star
    factor();
    // NOTE: CAN EVAL BEFORE RECURSIVE CALL TO ENFORCE LEFT ASSOCIATIVE OPERATORS
    term_prime(); // recursion
}
} /* End of function term */
```


Remove Left Recursion: Using EBNF as template

Option 2.

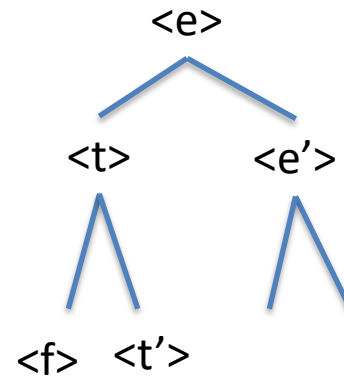
```
/* term
Parses strings in the language generated by the EBNF rule:
<term> -> <factor> { (* | /) <factor>}
*/
void term() {

    /* Parse the first factor */
    factor();

    /* Lookahead:  As long as the next token is * or /,
       next token and parse the next factor */
    // Uses while loop to allow for arbitrary repetition. Not as elegant, but practical.
    while (nextToken == MULT_OP || nextToken == DIV_OP) {
        MultDiv();
        factor();
    }
} /* End of function term */
```

Recursive Decent Parser a “limited” implementation of a PDA

- A Recursive Decent Parser performs a Predictive DFS of the parse tree
 - DFS inherently uses a stack to traverse tree and “store” unresolved symbols to be traversed later on



Recursive Decent Parser and the parse tree

- Recursive Decent Parser Traces out Parse Tree with recursive call chain

Example Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ (+ | -) \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ (* | /) \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow 0 | 1 | \dots | 9$

```
void expr() {  
    term(); // nextToken is Lookahead  
    while (nextToken == PLUS_OP || nextToken == SUB_OP) {  
        PlusSub();  
        term();  
    }  
}
```

```
void term() {  
    factor(); // nextToken is Lookahead  
    while (nextToken == MULT_OP || nextToken == DIV_OP) {  
        MultDiv();  
        factor();  
    }  
}
```

```
void PlusSub() { // expects to match '+' or '-'  
  
    if(thisToken == PLUS_OP || thisToken == SUB_OP) {  
        incrementTokenCounter() // Terminal Matched!  
        //Proceed to next token  
    }  
    else  
        throw parseError("Parse Error: + or - expected,  
        but instead found ", nextToken )  
}
```

```
void factor() { // expects a numeric terminal  
  
    if(thisToken == NUM) {  
        incrementTokenCounter() // Terminal Matched!  
        //Proceed to next token  
    }  
    else  
        throw parseError("Parse Error: NUM expected, but  
        instead found ", nextToken )  
}
```

Recursive Decent Parser Example

Assume input 3 + 5 * 6

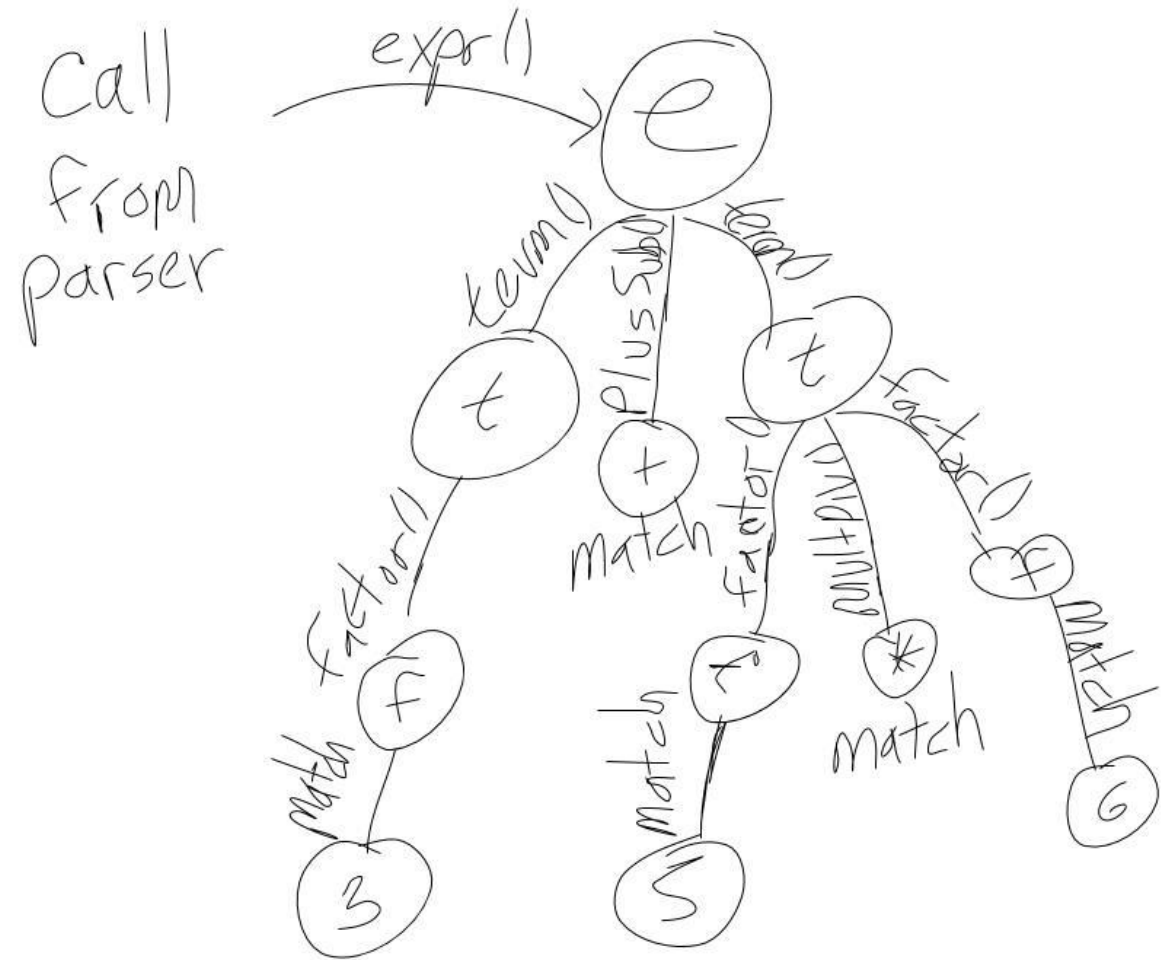
Example Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ (+ | -) \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ (* | /) \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle \rightarrow 0 | 1 | \dots | 9$

```
void expr() {  
    term();  
    while (nextToken == PLUS_OP || nextToken == SUB_OP) {  
        PlusSub();  
        term();  
    }  
}
```

```
void term() {  
    factor();  
    while (nextToken == MULT_OP || nextToken == DIV_OP) {  
        MultDiv();  
        factor();  
    }  
}
```

```
void factor() {  
    if(nextToken == NUM) {  
        incrementTokenCounter() // Terminal Matched! //Proceed to next token  
    }  
    else  
        throw parseError("Parse Error: NUM expected, but instead found ", nextToken )  
}
```



Recursive Decent Parser Example (In Class)

- Try

3 + 5 + -

3 * 5

3 -- 3

Example Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ (+ | -) \langle \text{term} \rangle \}$

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ (* | /) \langle \text{factor} \rangle \}$

$\langle \text{factor} \rangle \rightarrow 0 | 1 | \dots | 9$

Top-Down Parsing Table Implementation

- Non – recursive implementation of Top-Down Parser
 - More “explicitly” builds a limited PDA using a Table and Stack
- Table Driven Implementation
 1. Eliminate Left Recursion ** (discussed)
 2. Build Table
 3. Perform parse algorithm using table and stack

Build Top-Down Parse Table M

Table M: numAbstractions x numCharacters. Each entry is a production

\$: end of input symbol

1. For each production $A \Rightarrow \alpha$ do the following
2. For each terminal a in $\text{FIRST}(\alpha)$, add $A \Rightarrow \alpha$ to $M[A, a]$
3. If ϵ is in $\text{FIRST}(\alpha)$, add $A \Rightarrow \alpha$ to $M[A, b]$ for each b in $\text{FOLLOW}(A)$. If ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$, add $A \Rightarrow \alpha$ to $M[A, \$]$

3?: If α can be the empty string, then the next input char may be from $\text{FOLLOW}(A)$

$\text{FOLLOW}(A)$: All terminals that can directly follow A (not necessarily finish A) in a derivation. Be aware of ϵ productions.

Creating A Top-Down Parse Table. EG

- Start with Grammar:

$\langle e \rangle \Rightarrow \langle e \rangle + \langle t \rangle$

$\langle t \rangle \Rightarrow x \mid y$

- Step 1: Remove Left Recursion + Left Factor

$\langle e \rangle \Rightarrow \langle t \rangle \langle e' \rangle$

$\langle e' \rangle \Rightarrow + \langle t \rangle \langle e' \rangle \mid \varepsilon$

$\langle t \rangle \Rightarrow x \mid y$

Creating A Top-Down Parse Table. EG

FIRST($\langle t \rangle \langle e' \rangle$) = {x,y}

$\langle e \rangle \Rightarrow \langle t \rangle \langle e' \rangle$

$\langle e' \rangle \Rightarrow + \langle t \rangle \langle e' \rangle \mid \epsilon$

$\langle t \rangle \Rightarrow x \mid y$

M	x	y	+	\$
$\langle e \rangle$	$\langle e \rangle \Rightarrow \langle t \rangle \langle e' \rangle$	$\langle e \rangle \Rightarrow \langle t \rangle \langle e' \rangle$		
$\langle e' \rangle$			$\langle e' \rangle \Rightarrow + \langle t \rangle \langle e' \rangle$	$\langle e' \rangle \Rightarrow \epsilon$
$\langle t \rangle$	$\langle t \rangle \Rightarrow x$	$\langle t \rangle \Rightarrow y$		

Table M: numAbstractions x numCharacters. Each entry is a production

For each production $A \Rightarrow \alpha$ do the following

For each terminal a in FIRST(α), add $A \Rightarrow \alpha$ to $M[A,a]$

If ϵ is in FIRST(α), add $A \Rightarrow \alpha$ to $M[A,b]$ for each b in FOLLOW(A). If ϵ is in FIRST(α) and $\$$ is in FOLLOW(A), add $A \Rightarrow \alpha$ to $M[A,\$]$

Example: Perform a Top-Down Parse

M	x	y	+	\$
<e>	<e> => <t><e'>	<e> => <t><e'>		
<e'>			<e'>=> + <t><e'>	<e'>=> ε
<t>	<t> => x	<t> => y		

- Step 3: Using a Top-Down Parse Table

set pos to first symbol

set stack to S\$ // S = startSymbol

Let: X = peek(); a = input [pos] ;

do {

 if X is terminal or \$ {

 if X == a

 {pop X; pos++; }

 else

 syntax error

 }

 else if M[X, a] == X=>α₁ , ... α_k {

 pop X;

 push α_k , ... α₁ ;}

 else error

 } while(X != \$)

// accept

Example:
x+y

Stack	Input
\$<e>	x+y\$
\$<e'><t>	x+y\$
\$<e'>x	x+y\$
\$<e'>	+y\$
\$<e'> <t>+	+y\$
\$<e'> <t>	y\$
\$<e'> y	y\$
\$<e'>	\$
\$	\$

Try Example (20)

$A \Rightarrow Ab \mid AD \mid a \mid b$

$D \Rightarrow c \mid d$

Step 1: Remove Left Recursion + Left Factor

LL Grammars

- All Grammars that can be parsed using a top-down parser with k lookahead tokens are called LL(k) grammars.
 - Our top-down, table parsing algorithm can parse any LL(1) grammar. 1 lookahead token
 - Note: multiple entries in the table implies an ambiguous grammar (not LL(1))
 - LL: Left to right scan, leftmost derivation
 - LL(1) have the following characteristics,
 - Grammars cannot be left recursive
 - Whenever $A \Rightarrow \alpha \mid \beta$
 - » $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \phi$
 - » At most only one of α and β can produce ϵ
 - » If $\beta \Rightarrow^* \epsilon$, then α does not produce $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \phi$
 - LL(k) is a strict subset of CFG