

## Predictive Parsers

Can we avoid backtracking? Yes, if for a given input symbol and given non-terminal, we can choose the alternative **appropriately**.

This is possible if the first terminal of every alternative in a production is unique:

$A \rightarrow a B D \mid b B B$   
 $B \rightarrow c \mid b c e$   
 $D \rightarrow d$

parsing an input "abcd" has no backtracking.

Left factoring to enable predication:

$A \rightarrow \alpha\beta \mid \alpha\gamma$

change to

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

For predicative parsers, must eliminate left recursion

## LL(k) Parsing

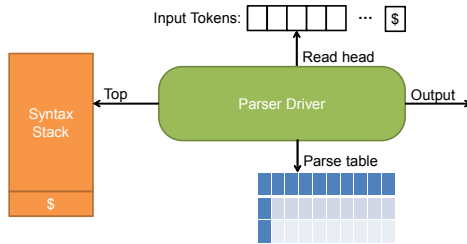
LL(k)

- L — left to right scan
- L — leftmost derivation
- k — k symbols of lookahead

in practice,  $k = 1$

It is table-driven and efficient.

## LL(k) Parser Structure



**Syntax stack** — hold right hand side (RHS) of grammar rules  
**Parse table** —  $M[A,b]$  — an entry containing rule " $A \rightarrow \dots$ " or error  
**Parser driver** — next action based on (*current token, stack top*)

## Sample Parse Table

	int	*	+	(	)	\$
E	$E \rightarrow TX$			$E \rightarrow TX$		
X			$X \rightarrow +E$		$X \rightarrow \epsilon$	$X \rightarrow \epsilon$
T	$T \rightarrow \text{int } Y$			$T \rightarrow ( E )$		
Y		$Y \rightarrow * T$	$Y \rightarrow \epsilon$		$Y \rightarrow \epsilon$	$Y \rightarrow \epsilon$

Implementation with 2-D parse table:

- A row for each non-terminal
- A column for all possible terminals and \$ (the end of input marker)
- Every table entry contains at most one production
  - Required for a grammar to be LL(1)
- No backtracking

Fixed action for each (*non-terminal, input symbol*) combination

## LL(1) Parsing Algorithm

$X$  — symbol at the top of the syntax stack  
 $a$  — current input symbol

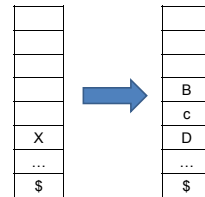
Parsing based on  $(X, a)$ :

- If  $X = a = \$$ , then
  - parser halts with "success"
- If  $X = a \neq \$$ , then
  - pop  $X$  from stack **and** advance input head
- If  $X \neq a$ , then
  - Case (a): if  $X \in T$ , then
    - parser halts with "failed," input rejected
  - Case (b): if  $X \in N$ ,  $M[X,a] = "X \rightarrow RHS"$ 
    - pop  $X$  **and** push RHS to stack in reverse order

## Push RHS in Reverse Order

$X$  — symbol at the top of the syntax stack  
 $a$  — current input symbol

if  $M[X,a] = "X \rightarrow B c D"$ :



## LL(1) Grammars

Remove left recursive and perform left factoring

Given the grammar:

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} * T \mid \text{int} \mid ( E )$$

The grammar has no left recursion but requires left factoring.

After rewriting grammar, we have:

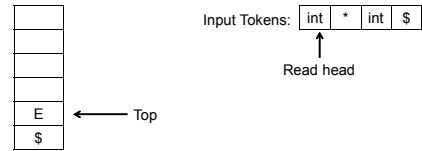
$$E \rightarrow TX$$

$$X \rightarrow + E \mid \epsilon$$

$$T \rightarrow \text{int} Y \mid ( E )$$

$$Y \rightarrow * T \mid \epsilon$$

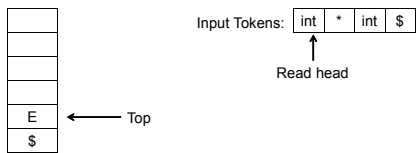
## LL(1) Parsing



Parse table

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

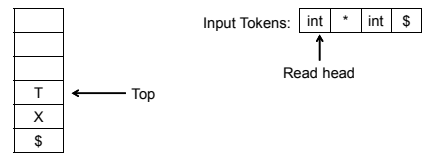
## LL(1) Parsing



Parse table

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

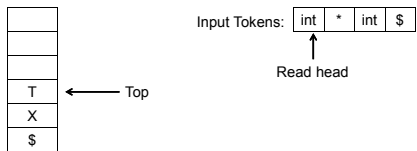
## LL(1) Parsing



Parse table

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

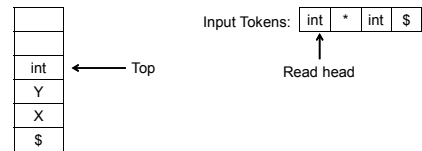
## LL(1) Parsing



Parse table

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

## LL(1) Parsing



Parse table

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε



### LL(1) Parsing

Input Tokens: int \* int \$

Stack: [Y, X, \$] ← Top

Read head: ↑

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

### LL(1) Parsing

Input Tokens: int \* int \$

Stack: [X, \$] ← Top

Read head: ↑

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

### LL(1) Parsing

Input Tokens: int \* int \$

Stack: [\$] ← Top

Read head: ↑

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

### LL(1) Parsing

Input Tokens: int \* int \$

Stack: [\$] ← Top

Read head: ↑

**Accept!**

	int	*	+	(	)	\$
E	E → TX			E → TX		
X			X → +E		X → ε	X → ε
T	T → int Y			T → ( E )		
Y		Y → * T	Y → ε		Y → ε	Y → ε

### Action List

Stack	Input	Action
E \$	int * int \$	E → TX
T X \$	int * int \$	T → int Y
int Y X \$	int * int \$	terminal
Y X \$	* int \$	Y → * T
* T X \$	* int \$	terminal
T X \$	int \$	T → int Y
int Y X \$	int \$	terminal
Y X \$	\$	Y → ε
X \$	\$	X → ε
\$	\$	Halt and accept

### Constructing the Parse Table

We need to know what non-terminals to place our productions in the table?

We know that we have restricted our grammars so that left recursion is eliminated and they have been left factored. That means that each production is uniquely recognizable by the first terminal that production would derive.

Thus, we can construct our table from 2 sets:

- For each symbol A, the set of terminals that can begin a string derived from A. This set is called the **FIRST** set of A
- For each non-terminal A, the set of terminals that can appear after a string derived from A is called the **FOLLOW** set of A

## First( $\alpha$ )

First( $\alpha$ ) = set of terminals that start string of terminals derived from  $\alpha$ .

Apply following rules until no terminal or  $\epsilon$  can be added

1. If  $t \in T$ , then **First**( $t$ ) = { $t$ }.  
For example **First**( $+$ ) = { $+$ }.
2. If  $X \in N$  and  $X \rightarrow \epsilon$  exists (nullable), then add  $\epsilon$  to **First**( $X$ ).  
For example, **First**( $Y$ ) = { $^*$ ,  $\epsilon$ }.
3. If  $X \in N$  and  $X \rightarrow Y_1 Y_2 Y_3 \dots Y_m$ , where  $Y_1, Y_2, Y_3, \dots, Y_m$  are non-terminals, then:  
for each  $i$  from 1 to  $m$   
if  $Y_1, \dots, Y_{i-1}$  are all nullable (or if  $i=1$ )  
**First**( $X$ ) = **First**( $X$ )  $\cup$  **First**( $Y_i$ )

## Follow( $\alpha$ )

**Follow**( $\alpha$ ) = { $t \mid S \Rightarrow ^* \alpha t \beta$ }

- Intuition: if  $X \rightarrow A B$ , then **First**( $B$ )  $\subseteq$  **Follow**( $A$ )
- However,  $B$  may be  $\epsilon$ , i.e.,  $\beta \Rightarrow \epsilon$

Apply following rules until no terminal or  $\epsilon$  can be added

1.  $\$ \in$  **Follow**( $S$ ), where  $S$  is the start symbol.  
e.g., **Follow**( $E$ ) = { $\$$  ... }.
2. Look at the occurrence of a non-terminal on the right hand side of a production which is followed by something  
If  $A \rightarrow \alpha B \beta$ , then **First**( $\beta$ ) - { $\epsilon$ }  $\subseteq$  **Follow**( $B$ )
3. Look at  $N$  on the RHS that is not followed by anything,  
if ( $A \rightarrow \alpha B$ ) or ( $A \rightarrow \alpha B \beta$  and  $\epsilon \in$  **First**( $\beta$ )),  
then **Follow**( $A$ )  $\subseteq$  **Follow**( $B$ )

## Algorithm to Compute FIRST, FOLLOW, and nullable

Initialize FIRST and FOLLOW to all empty sets, and nullable to all false.

```

foreach terminal symbol Z
    FIRST[Z] = {}
do
    foreach production X  $\rightarrow$   $Y_1 Y_2 \dots Y_k$ 
        if  $Y_1 \dots Y_k$  are all nullable (or if  $k = 0$ )
            then nullable[X] = true
        foreach i from 1 to k, each j from i + 1 to k
            if  $Y_1 \dots Y_{i-1}$  are all nullable (or if  $i = 1$ )
                then FIRST[X] = FIRST[X]  $\cup$  FIRST[ $Y_i$ ]
            if  $Y_{i+1} \dots Y_k$  are all nullable (or if  $i = k$ )
                then FOLLOW[ $Y_i$ ] = FOLLOW[ $Y_i$ ]  $\cup$  FOLLOW[X]
            if  $Y_{i+1} \dots Y_{j-1}$  are all nullable (or if  $i + 1 = j$ )
                then FOLLOW[ $Y_j$ ] = FOLLOW[ $Y_j$ ]  $\cup$  FIRST[ $Y_j$ ]
    until FIRST, FOLLOW, and nullable did not change in this iteration.
    
```

## Example

**Grammar:**

```

E  $\rightarrow$  T X
X  $\rightarrow$  + E |  $\epsilon$ 
T  $\rightarrow$  int Y | ( E )
Y  $\rightarrow$  * T |  $\epsilon$ 
    
```

**First Set:**      **Follow Set:**  
E  $\rightarrow$  T X      \$  
X  $\rightarrow$  + E      E  $\rightarrow$  T X  
X  $\rightarrow$   $\epsilon$       X  $\rightarrow$  + E  
T  $\rightarrow$  int Y      T  $\rightarrow$  int Y  
T  $\rightarrow$  ( E )      T  $\rightarrow$  ( E )  
Y  $\rightarrow$  \* T      Y  $\rightarrow$  \* T  
Y  $\rightarrow$   $\epsilon$

Symbol	First	Follow
(	(	
)	)	
+	+	
*	*	
int	int	
Y	* , $\epsilon$	\$ , ) , +
X	+ , $\epsilon$	\$ , )
T	( , int	\$ , ) , +
E	( , int	\$ , )

## Constructing LL(1) Parse Table

To construct the parse table, we check each  $A \rightarrow \alpha$

- For each terminal  $a \in$  **First**( $\alpha$ ), add  $A \rightarrow \alpha$  to  $M[A, a]$ .
- If  $\epsilon \in$  **First**( $\alpha$ ), then for each terminal  $b \in$  **Follow**( $A$ ),  
• add  $A \rightarrow \alpha$  to  $M[A, a]$ .
- If  $\epsilon \in$  **First**( $\alpha$ ) and  $\$ \in$  **Follow**( $A$ ), then add  $A \rightarrow \alpha$  to  $M[A, \$]$ .

## Constructing LL(1) Parse Table

For each terminal  $a \in$  **First**( $\alpha$ ), add  $A \rightarrow \alpha$  to  $M[A, a]$ .

**Grammar:**

```

E  $\rightarrow$  T X
X  $\rightarrow$  + E
X  $\rightarrow$   $\epsilon$ 
T  $\rightarrow$  int Y
T  $\rightarrow$  ( E )
Y  $\rightarrow$  * T
Y  $\rightarrow$   $\epsilon$ 
    
```

Symbol	First	Follow
(	(	
)	)	
+	+	
*	*	
int	int	
Y	* , $\epsilon$	\$ , ) , +
X	+ , $\epsilon$	\$ , )
T	( , int	\$ , ) , +
E	( , int	\$ , )

	int	*	+	(	)	\$
E	E $\rightarrow$ T X			E $\rightarrow$ T X		
X			X $\rightarrow$ + E			
T	T $\rightarrow$ int Y			T $\rightarrow$ ( E )		
Y		Y $\rightarrow$ * T				

## Constructing LL(1) Parse Table

If  $\epsilon \in \text{First}(\alpha)$ , then for each terminal  $b \in \text{Follow}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \alpha]$ .

**Grammar:**

$E \rightarrow T X$   
 $X \rightarrow + E$   
 $X \rightarrow \epsilon$   
 $T \rightarrow \text{int } Y$   
 $T \rightarrow ( E )$   
 $Y \rightarrow * T$   
 $Y \rightarrow \epsilon$

Symbol	First	Follow
(	(	
)	)	
+	+	
*	*	
int	int	
Y	*, $\epsilon$	), +
X	+, $\epsilon$	), \$
T	(, int	), +
E	(, int	), \$

	int	*	+	(	)	\$
E	$E \rightarrow T X$			$E \rightarrow T X$		
X			$X \rightarrow + E$		$X \rightarrow \epsilon$	
T	$T \rightarrow \text{int } Y$			$T \rightarrow ( E )$		
Y		$Y \rightarrow * T$	$Y \rightarrow \epsilon$		$Y \rightarrow \epsilon$	

## Constructing LL(1) Parse Table

If  $\epsilon \in \text{First}(\alpha)$  and  $\$ \in \text{Follow}(A)$ , then add  $A \rightarrow \alpha$  to  $M[A, \$]$ .

**Grammar:**

$E \rightarrow T X$   
 $X \rightarrow + E$   
 $X \rightarrow \epsilon$   
 $T \rightarrow \text{int } Y$   
 $T \rightarrow ( E )$   
 $Y \rightarrow * T$   
 $Y \rightarrow \epsilon$

Symbol	First	Follow
(	(	
)	)	
+	+	
*	*	
int	int	
Y	*, $\epsilon$	), +
X	+, $\epsilon$	), \$
T	(, int	), +
E	(, int	), \$

	int	*	+	(	)	\$
E	$E \rightarrow T X$			$E \rightarrow T X$		
X			$X \rightarrow + E$		$X \rightarrow \epsilon$	$X \rightarrow \epsilon$
T	$T \rightarrow \text{int } Y$			$T \rightarrow ( E )$		
Y		$Y \rightarrow * T$	$Y \rightarrow \epsilon$		$Y \rightarrow \epsilon$	$Y \rightarrow \epsilon$

## Is a Grammar LL(1)?

**Observation**

If a grammar is LL(1), then each of its LL(1) table entries contain at most one rule. Otherwise, it is not LL(1)

Two methods to determine if a grammar is LL(1) or not:

1. Construct LL(1) table, and check if there is a multi-rule entry or
2. Check each rule as if the table were being constructed:

G is LL(1) iff for a rule  $A \rightarrow \alpha | \beta$

- a)  $\text{First}(\alpha) \cap \text{First}(\beta) = \emptyset$
- b) at most one of  $\alpha$  and  $\beta$  can derive  $\epsilon$
- c) If  $\beta$  derives  $\epsilon$ , then  $\text{First}(\alpha) \cap \text{Follow}(\beta) = \emptyset$

## Ambiguous Grammars

Some grammars may need more than one token of lookahead (k). However, some grammars are not LL regardless of how the grammar is changed.

$S \rightarrow \text{if } C \text{ then } S \mid \text{if } C \text{ then } S \text{ else } S \mid \dots$   
 $C \rightarrow b$

change to

$S \rightarrow \text{if } C \text{ then } S X \mid \dots$   
 $X \rightarrow \text{else } S \mid \epsilon$   
 $C \rightarrow b$

**problem sentence:** "if b then if b then a else a"

- "else"  $\in \text{First}(X)$
- $\text{First}(X) - \epsilon \subseteq \text{Follow}(S)$
- $X \rightarrow \text{else } \dots \mid \epsilon$
- "else"  $\in \text{Follow}(X)$

## Removing Ambiguity

To remove ambiguity, it is possible to rewrite the grammar.

For the "if-then-else" example, how to rewrite?

May not even need to rewrite in this case, we can just use the  $X \rightarrow \text{else } S$  production over the  $X \rightarrow \epsilon$

However, by changing the grammar,

- It might make the other phases of the compiler more difficult
- It becomes harder to determine semantics and generate code
- It is less appealing to programmers

## LL(1) Summary

LL(1) parsers operate in linear time and at most linear space relative to the length of input because:

Time — each input symbol is processed constant number of times

Space — stack is smaller than the input (in case we remove  $X \rightarrow \epsilon$ )

## Summary

---

**First** and **Follow** sets are used to construct predictive parsing tables

Intuitively, **First** and **Follow** sets guide the choice of rules:

- For non-terminal **A** and input **t**, use a production rule  $A \rightarrow \alpha$  where  $t \in \mathbf{First}(\alpha)$
- For non-terminal **A** and input **t**, if  $A \rightarrow \alpha$  and  $t \in \mathbf{Follow}(A)$ , use the production  $A \rightarrow \alpha$  where  $\varepsilon \in \mathbf{First}(\alpha)$

What is LL(0)?

Why are LL(2) ... LL(k) are not widely used ?