

## What we'll talk about today

## CS164 Sections 1 & 4

- LL1 parsing
- LL1 on a really nice grammar
- LL1 on a nice grammar
  - o First and Follow Sets
- LL1 on an almost nice grammar
  - o Left Factoring
- LL1 failing miserably with a not quite nice grammar

### LL1 parsing, the idea:

You will start with a single non-terminal, the Start production.

You want to expand your non-terminals by using the production rules. i.e. build the tree top down.

You would like to be able to be able to decide which production to use based on the next token.

Given a non-terminal with several possible productions, how do we tell which one to use next? That's the question at the heart of LL1 parsing.

### Really nice grammar:

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

	+	(	*	int	)
E					
T					

Here are some sample strings you can produce from this grammar.

+ int int  
+ \* int \* int (+ int int) int  
(\* ((int)) \* int int)

In this case, for a given non-terminal, it's easy to tell which rule to use if we know what the next token of the input is. Simply pick the production whose first token matches the next token in the input. We'll generalize this idea for cases where the grammar is not so nice.

## Nice grammar:

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

	+	(	*	int	)
E					
X					
T					
Y					

In this case, we can use the same rule as before. But how do we know what the first token will be for a production like  $E \rightarrow TX$ .

## First and Follow Sets.

In order to solve this problem, you want to find the set of all terminals that could possibly start a string corresponding to this production.

For example, you would like to find all possible terminals that could start a string that can be produced by  $T X$  in our above example. Such a set is called the First set of  $T X$ , and is denoted as  $\text{First}(T X)$ .

You can define it recursively:

Suppose you are given a sentence of terminals and non-terminals  $A_1 A_2 A_3 \dots A_n$

Then there are four cases we need to consider:

- If  $n=1$  and  $A_1$  is a Terminal,  $\text{First}(A_1) = \{ A_1 \}$

- If  $n \neq 1$  and  $\text{First}(A_1)$  doesn't contain  $\epsilon$ ,  
 $\text{First}(A_1 A_2 A_3 \dots A_n) = \text{First}(A_1)$

- If  $n=1$  and  $A_1$  is a non-Terminal, Find all productions of the form

$$A_1 \rightarrow B_{1,1} B_{2,1}, \dots B_{n_1,1}$$

$$A_1 \rightarrow B_{1,2} B_{2,2}, \dots B_{n_2,2}$$

...

$$A_1 \rightarrow B_{1,k} B_{2,k}, \dots B_{n_k,k}$$

Then,  $\text{First}(A_1) = \text{Union}(\text{First}(B_{1,i}, \dots, B_{n_i,i}))$

- If  $n \neq 1$  and  $\text{First}(A_1)$  contains  $\epsilon$ ,  
 $\text{First}(A_1 A_2 A_3 \dots A_n) = \text{Union}(\text{First}(A_1) - \epsilon, \text{First}(A_2 A_3 \dots A_n))$

First sets allow you to know where to use each production. From our above example,

$\text{First}(TX) = \{ (, \text{int} \}$

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

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

$\text{First}(\text{int } Y) = \{ \text{int} \}$

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

So from this we can fill the following entries in the table:

	+	(	*	int	)
E		TX		TX	
X	+E				
T		(E)		int Y	
Y			*T		

But how do we know when to use the  $\epsilon$  rules?

Then you need Follow Sets, which are the set of tokens that could potentially come after a non Terminal.

To compute  $\text{Follow}(X)$ ,

Start by setting  $\text{Follow}X$  to empty

Find all rules of the form  $Y \rightarrow \dots X A_1 \dots A_n$

Then, for each of this rule, add  $\text{First}(A_1 \dots A_n)$  to  $\text{Follow}(X)$

And if  $\epsilon$  is in  $\text{First}(A_1 \dots A_n)$ , you must also add  $\text{Follow}(Y)$  to  $\text{Follow}(X)$ .

Note that  $\text{Follow}(Y)$  may require that you know  $\text{Follow}(X)$ , which leads to a chicken-egg problem. How do you fix it? You have to do it iteratively.

Initialize the Follow sets to  $\{\}$ . Then, update them according to the following equations (which you can derive from the above rules), and repeat until they stop changing.

$\text{Follow}(E) = \{ ), \$ \} \cup \text{Follow}(X)$

$\text{Follow}(X) = \text{Follow}(E)$

$\text{Follow}(T) = \text{First}(X) \cup \text{Follow}(E) \cup \text{Follow}(Y)$

$\text{Follow}(Y) = \text{Follow}(T)$

From this, you will get

Follow( E ) = { }, \$ }

Follow( X ) = { }, \$ }

Follow( T ) = { +, ε, ), \$ }

Follow( Y ) = { +, ε, ), \$ }

Now, we can handle the cases we were missing.

	+	(	*	int	)
E	ERROR	TX	ERROR	TX	ERROR
X	+E	ERROR	ERROR	ERROR	ε
T	ERROR	( E )	ERROR	int Y	ERROR
Y	ε	ERROR	* T	ERROR	ε

## Almost Nice grammar

In some cases, we may be given grammars for which it seems you can not use this technique.

$E \rightarrow T + E \mid T$

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

Many times, you can rescue them by using **Left Factoring**

Factor out common prefixes: T, int

$E \rightarrow T X$

$X \rightarrow + E \mid \epsilon$

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

$Y \rightarrow * T \mid \epsilon$

## L1 failing miserably with a not quite nice grammar

Mostly for your own entertainment, here is an example where left factoring is not enough, and LL1 will fail miserably:

$$\begin{aligned} E &\rightarrow T + E \mid T \\ T &\rightarrow \text{int} \mid \text{int} * T \mid (E) \mid \text{int} + \end{aligned}$$

After left factoring, we get:

$$\begin{aligned} E &\rightarrow T X \\ X &\rightarrow + E \mid \varepsilon \\ T &\rightarrow ( E ) \mid \text{int} Y \mid T \\ Y &\rightarrow * T \mid \varepsilon \mid + \end{aligned}$$

This is a perfectly legal context free grammar, and it is in fact unambiguous. The problem is that the Follow set of Y contains +, which means that on seeing a + you should use the epsilon rule for Y, but + is also in the First set for +, so this creates a conflict; you can't tell which rule to use by just looking ahead by one.