# COMP 442 / 6421 Compiler Design

## Tutorial 3
## Syntactic Analyser

Instructor:       Dr. Joey Paquet        paquet@cse.concordia.ca
TAs:                 Haotao Lai             h_lai@encs.concordia.ca

# Lab Instructor

Section:  lab hours NNK    M------          20:30-22:20    H819

Name: Haotao Lai (Eric)

Office: EV 8.241

Email: h_lai@encs.concordia

Website: http://laihaotao.me/ta

# The Goal of Assignment 2

- input: a list of tokens from the first assignment
- output: an abstract syntax tree

1. Convert the given CFG to an LL(1) grammar
   a. Use tools to help your transformation procedure
   b. Remove the grammar from EBNF to non-EBNF representation
   c. Remove ambiguities and left recursions
   d. After each transformation step, verify that your grammar was not broken

2. Implement a LL(1) parser
   a. Recursive descent predictive parsing
   b. Table-driven predictive parsing

# Notation

| Names Beginning With | Represent Symbols In | Examples |
|---|---|---|
| Uppercase | $N$ | A, B, C, Prefix |
| Lowercase and punctuation | $\Sigma$ | a, b, c, if, then, (, ; |
| $X, Y$ | $N \cup \Sigma$ | $X_i, Y_3$ |
| Other Greek letters | $(N \cup \Sigma)^{\star}$ | $\alpha, \gamma$ |

If A→γ is a production, then αAβ ⇒ αγβ denotes one step of a derivation using this production.

⇒     denote one step of a derivation of a production

⇒ +   derives in one or more steps

⇒ *   derives in zero or more steps

# Convert CFG an LL(1)G

In context free grammar, all rules are one-to-one, one-to-many or one-to-none. These rule can be applied regardless of context [1]. A CFG can be defined as G(T, N, P, S), it is usually represented by using BNF notation.

Given the following grammar:
S → AA   A → a   A → b

The following input are valid:
aa   ab   ba   bb

[1] https://en.wikipedia.org/wiki/Context-free_grammar

# Convert CFG an LL(1)G

LL(k) grammar is a formal grammar that can be parsed by an LL parser, which parse the input from left to right and constructs a leftmost derivation of the sentence. The k within the parenthesis is the number of token the parser will lookahead when parsing a sentence.

Given a grammar G with three productions: $S \to E$    $E \to (E + E)$    $E \to i$
and input string: $w = ( ( i + i ) + i )$, the leftmost derivation will be:

$S \to E \to (E + E) \to ((E + E) + E) \to ((i + E) + E) \to ((i + i) + E) \to ((i + i) + i)$

# Remove EBNF Representation

We need to remove two notations introduced by the EBNF format which are repetition and optional of the symbol.

Repetition example: A -> B {C} D
Optional example:    A -> B [C] D

They all can be eliminated by introducing an new nonterminal symbol.

```
A -> B {C} D
------------
A  -> B C' D
C' -> C C' | epsilon
```

```
A -> B [C] D
------------
A  -> B C' D
C' -> C | epsilon
```

# LL(1) Parsing

In LL(1) parsing, for each combination of a nonterminal and a input token, there should be only one possible production (if it is syntax valid) or no production (if it is a error state).

So we need to make sure after elimination of the EBNF notation, our grammar should be deterministic for each nonterminal symbol. A non-deterministic example can be: <u>A -> B C | B D</u>

Another situation that can lead to the failure of LL(1) parsing is left recursion of the production, an example can be: <u>A -> A a | b</u>. With such an production, you don't know where will be the ending point during the parsing.

# Three Roadblocks

Quick review

1. Ambiguity
2. Non-deterministic
3. Left recursion

For theoretical detail, see the lecture slide set [syntax analysis: introduction].

# Ambiguity Grammar

Grammar: <u>E -> E + E | E * E | id</u>
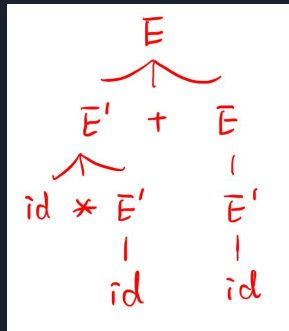
Input string: <u>id * id + id</u>



Requirement of the parse tree:

A tree that its in order traversal should give the string same as the input string

# Ambiguity Grammar

The solution for ambiguity is rewrite the grammar (that's exactly what you need to do in assignment 2) to make it unambiguous.

In this case, we want to enforce precedence of multiplication over addition.



original: E -> E + E | E * E | id

modified:

E -> E' + E | E'

E' -> id * E' | id

**Note**

The modified grammar here is not a LL(1) grammar, the example here just show how to remove ambiguity.

If you look carefully, you will find it is actually a LL(2) grammar

# Non-deterministic Grammar

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \alpha\beta_3$$

1. backtracking can solve this problem, but it is inefficient;
2. introduce a new non-terminal which we refer as left factoring

$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$$

# Left Recursion

Grammar: A -> Aα | β

By analyze these three possibilities, our goal is to construct something like:  A -> βα*

But we don't allow * in the grammar, so we can replace a* with a new non-terminal A', so we have:

A -> βA'
A' -> αA' | ε

# Perform Left Factoring

A -> B C | B D

A -> B'
B' -> C | D

# Eliminate Left Recursion

A -> A a | b

A -> b A'
A' -> a A' | epsilon

# Parsing Example

Assume we have the following grammar:

E → T + E | T
T → int | int * T | ( E )

Left-factored grammar:

E → T X
T → ( E ) | int Y
X → + E | ε
Y → * T | ε

# First Set

Definition

$$\text{First}(X) = \{ t \mid X \to^* t\alpha\} \cup \{\varepsilon \mid X \to^* \varepsilon\}$$

Algorithm sketch:

1. $\text{First}(t) = \{ t \}$
2. $\varepsilon \in \text{First}(X)$
   - if $X \to \varepsilon$
   - if $X \to A_1 \ldots A_n$ and $\varepsilon \in \text{First}(A_i)$ for $1 \le i \le n$
3. $\text{First}(\alpha) \subseteq \text{First}(X)$ if $X \to A_1 \ldots A_n \; \alpha$
   - and $\varepsilon \in \text{First}(A_i)$ for $1 \le i \le n$

# First Set

E → T X
T → ( E ) | int Y
X → + E | ε
Y → * T | ε

First( ( ) = { ( }          First( T ) = {int, ( }
First( ) ) = { ) }          First( E ) = {int, ( }
First( int) = { int }       First( X ) = {+, ε }
First( + ) = { + }          First( Y ) = {*, ε }
First( * ) = { * }

# Follow Set

- Definition:

$$\text{Follow}(X) = \{\ t\ |\ S \rightarrow^* \beta X t \delta\ \}$$

Algorithm sketch:

1. $\$ \in \text{Follow}(S)$
2. $\text{First}(\beta) - \{\epsilon\} \subseteq \text{Follow}(X)$
   - For each production $A \rightarrow \alpha X \beta$
3. $\text{Follow}(A) \subseteq \text{Follow}(X)$
   - For each production $A \rightarrow \alpha X \beta$ where $\epsilon \in \text{First}(\beta)$

# Follow Set

E → T X
T → ( E ) | int Y
X → + E | ε
Y → * T | ε

Follow( + ) = { int, ( }    Follow( * ) = { int, ( }
Follow( ( ) = { int, ( }    Follow( E ) = {), $}
Follow( X ) = {$, ) }    Follow( T ) = {+, ) , $}
Follow( ) ) = {+, ) , $}    Follow( Y ) = {+, ) , $}
Follow( int) = {*, +, ) , $}

# First Set and Follow Set

example 1:

S -> A B C D E
A -> a | ε
B -> b | ε
C -> c
D -> d | ε
E -> e | ε

# First Set and Follow Set

example 2:

S -> B b | C d
B -> a B | ε
C -> c C | ε

# First Set and Follow Set

example 3:

S -> A C B | C b B | B a
A -> d a | B C
B -> g | ε
C -> h | ε

the note shows what I did in the lab can access: http://laihaotao.me/ta/w18_comp442_fst_flw_set.pdf

# How to come up with the proper grammar?

- You receive the initial grammar in EBNF in assignment 2 description already
- You need to remove the EBNF since AtoCC kfgEdit cannot understand this form
- Perform left factoring (if necessary)
- Remove left recursion (if exist, unfortunately, they exist in the given grammar)

It is strongly suggested that every time you make a single transformation step, that you use AtoCC to check whether your transformation broke the grammar or not.

Don't try to correct many errors in one shot, it is easy to get lost. Plus, if you make a mistake in one transformation step and you carry on without checking, your further transformation will be made on a wrong grammar and thus be invalid.

# Installing AtoCC

# Installing AtoCC

- AtoCC can be downloaded at the following web site:

  - http://www.atocc.de
- You can either download an installer, or precompiled applications.

# You don't have Windows machine?

Check the following link out:

http://atocc.de/AtoCCFAQ/index.php?option=com_content&task=category&sectionid=11&id=25&Itemid=34

Works, for example, for macOS High Sierra version 10.13.1

# Install AtoCC without administration rights ?



These are portable executables, but they often crash, so save your work frequently!

# Automated grammar transformation tools



- CyberZHG's Compiler construction toolkit:

    https://cyberzhg.github.io/toolbox/

- Can help you apply specific transformations

- Use in conjunction with kfgEdit

- However, it does not use the same grammar representation conventions

# Example
--- How to use AtoCC for verification

# AtoCC kfgEdit

- Tool that allows you to analyze your grammar and locate possible ambiguities in the grammar.
- After you grammar is entered, it also allows you to enter a string representing a token stream and verify if this token stream is derivable from the grammar.  If it is, it generates a parse tree and a derivation for it.

File  Help

New  Open  Save  Validate Grammar  is regular ?  Export Automaton  Export Compiler

kfG Edit | Language | Grammar | Derivation | LL(1) conditions | Definition

## kfG Edit
### First&Follow

**LL(1) Conditions:**

- Check Condition 1
- Check Condition 2
- is LL(1) Grammar?

$$E \rightarrow \alpha_0 \mid \alpha_1 \mid \alpha_2$$

with:
$\alpha_0 = T$
$\alpha_1 = E - T$
$\alpha_2 = E + T$

First-Sets:
$FIRST(\alpha_0) = \{(, \text{id}\}$
$FIRST(\alpha_1) = \{(, \text{id}\}$
$FIRST(\alpha_2) = \{(, \text{id}\}$

| $\cap$ | $\alpha_0$ | $\alpha_1$ | $\alpha_2$ |
|---|---|---|---|
| $\alpha_0$ | – | $\{(, \text{id}\}$ | $\{(, \text{id}\}$ |
| $\alpha_1$ | $\{(, \text{id}\}$ | – | $\{(, \text{id}\}$ |
| $\alpha_2$ | $\{(, \text{id}\}$ | $\{(, \text{id}\}$ | – |

$$T \rightarrow \alpha_0 \mid \alpha_1 \mid \alpha_2$$

with:
$\alpha_0 = F$
$\alpha_1 = T / F$
$\alpha_2 = T * F$

First-Sets:
$FIRST(\alpha_0) = \{(, \text{id}\}$
$FIRST(\alpha_1) = \{(, \text{id}\}$
$FIRST(\alpha_2) = \{(, \text{id}\}$

| $\cap$ | $\alpha_0$ | $\alpha_1$ | $\alpha_2$ |
|---|---|---|---|
| $\alpha_0$ | – | $\{(, \text{id}\}$ | $\{(, \text{id}\}$ |
| $\alpha_1$ | $\{(, \text{id}\}$ | – | $\{(, \text{id}\}$ |

$E \rightarrow \alpha_0 \mid \alpha_1 \mid \alpha_2$

with:
$\alpha_0$ = T
$\alpha_1$ = E − T
$\alpha_2$ = E + T

First-Sets:
FIRST($\alpha_0$) = {(, id}
FIRST($\alpha_1$) = {(, id}
FIRST($\alpha_2$) = {(, id}

| ∩ | $\alpha_0$ | $\alpha_1$ | $\alpha_2$ |
|---|---|---|---|
| $\alpha_0$ | − | {(, id} | {(, id} |
| $\alpha_1$ | {(, id} | − | {(, id} |
| $\alpha_2$ | {(, id} | {(, id} | − |

first set intersection

$T \rightarrow \alpha_0 \mid \alpha_1 \mid \alpha_2$

with:
$\alpha_0$ = F
$\alpha_1$ = T / F
$\alpha_2$ = T * F

First-Sets:
FIRST($\alpha_0$) = {(, id}
FIRST($\alpha_1$) = {(, id}
FIRST($\alpha_2$) = {(, id}

| ∩ | $\alpha_0$ | $\alpha_1$ | $\alpha_2$ |
|---|---|---|---|
| $\alpha_0$ | − | {(, id} | {(, id} |
| $\alpha_1$ | {(, id} | − | {(, id} |
| $\alpha_2$ | {(, id} | {(, id} | − |

$$F \rightarrow \alpha_0 \mid \alpha_1$$

with:
  $\alpha_0$ = id
  $\alpha_1$ = ( E )

First-Sets:
  FIRST($\alpha_0$) = {id}
  FIRST($\alpha_1$) = {(}

| $\cap$ | $\alpha_0$ | $\alpha_1$ |
|--------|------------|------------|
| $\alpha_0$ | – | $\varnothing$ |
| $\alpha_1$ | $\varnothing$ | – |

go to the very end of the page

LL(1) first condition not fulfilled!

35

# What you should do?



there is something wrong with this prodution

1. Locate a specific error and identify the faulty productions (shown in red)
2. Copy the related productions into the grammar transformation tool mentioned above (https://cyberzhg.github.io/toolbox/cfg2ll).
3. Copy the correction from the tool and paste it into AtoCC
4. Do some modification to adapt to AtoCC format
5. Check the grammar again

Note: Don't try to solve more than one production at a time. When you solve one production's error, use the tool to check to make sure you are not bringing new errors.

```
14 E -> T E''
15 T -> F T''
16 F -> ( E )
17    | id
18 E' -> + T
19    | - T
20 T' -> * F
21    | / F
22 E'' -> E' E''
23     | ?
24 T'' -> T' T''
25     | ?
```

```
1 E -> T ETailTail
2 T -> F TTailTail
3 F -> ( E )
4    | id
5 ETail -> + T
6     | - T
7 TTail -> * F
8     | / F
9 ETailTail -> ETail ETailTail
10           | EPSILON
11 TTailTail -> TTail TTailTail
12           | EPSILON
13
```

result from the tool                    after modification, adapted to AtoCC

## LL(1) first condition fulfilled!

```
FIRST (ETailTail) = {+, -, EPSILON}
FOLLOW(ETailTail) = {$, )}
FIRST (ETailTail) ∩ FOLLOW(ETailTail) = ∅
```

```
FIRST (TTailTail) = {*, /, EPSILON}
FOLLOW(TTailTail) = {$, ), +, -}
FIRST (TTailTail) ∩ FOLLOW(TTailTail) = ∅
```

## LL(1) second condition fulfilled!

# Tool

If you plan to use the table-driven approach, you will need a parse table. Of course you can generate your own parse table, or put a proper grammar into a tool and it will give you the table.

We propose an online tool to do that:
http://hackingoff.com/compilers/ll-1-parser-generator

note: you need to use EPSILON to repersent ε

# Thanks