

Context Free Grammar – Quick Review

- Grammar - quadruple
 - A set of tokens (terminals): T
 - A set of non-terminals: N
 - A set of productions $\{ \text{lhs} \rightarrow \text{rhs} , \dots \}$
 - lhs in N
 - rhs is a sequence of $N \cup T$
 - A Start symbol: S (in N)
- Shorthands
 - Provide only the productions
 - All lhs symbols comprise N
 - All other symbols comprise T
 - lhs of first production is S

Using Grammars to derive Strings

- Rewriting rules
 - Pick a non-terminal to replace. Which order?
 - left-to-right
 - right-to-left
- Derives relation: $\alpha A \gamma \Rightarrow \alpha \beta \gamma$
 - When $A \rightarrow \beta$ is a production
- Derivations (a list of productions used to derive a string from a grammar).
- A sentence of G: $L(G)$
 - Start with S
 - $S \Rightarrow^* w$ where w is only terminal symbols
 - all strings of terminals derivable from S in 1 or more steps

CF Grammar Terms

- Parse trees.
 - Graphical representations of derivations.
 - The leaves of a parse tree for a fully filled out tree is a sentence.
- Regular language v.s. Context Free Languages
 - how do CFL compare to regular expressions?
 - Nesting (matched ()) requires CFG,'s RE's are not powerful enough.
- Ambiguity
 - A string has two derivations
 - $E \rightarrow E + E \quad | \quad E * E \quad | \quad id$
 - $x + x * y$
- Left-recursion
 - $E \rightarrow E + E \quad | \quad E * E \quad | \quad id$
 - Makes certain top-down parsers loop

Parsing

- Act of constructing derivations (or parse trees) from an input string that is derivable from a grammar.
- Two general algorithms for parsing
 - Top down - Start with the start symbol and expand Non-terminals by looking at the input
 - Use a production on a left-to-right manner
 - Bottom up - replace sentential forms with a non-terminal
 - Use a production in a right-to-left manner

Top Down Parsing

- Begin with the start symbol and try and derive the parse tree from the root.
- Consider the grammar
 1. $\text{Exp} \rightarrow \text{Id} \mid \text{Exp} + \text{Exp} \mid \text{Exp} * \text{Exp} \mid (\text{Exp})$
 2. $\text{Id} \rightarrow x \mid y$

Some strings derivable from the grammar

x
x+x
x+x+x,
x * y
x + y * z ...

Example Parse (top down)

– stack input

Exp x + y * z

 Exp x + y * z
 / | \
Exp + Exp

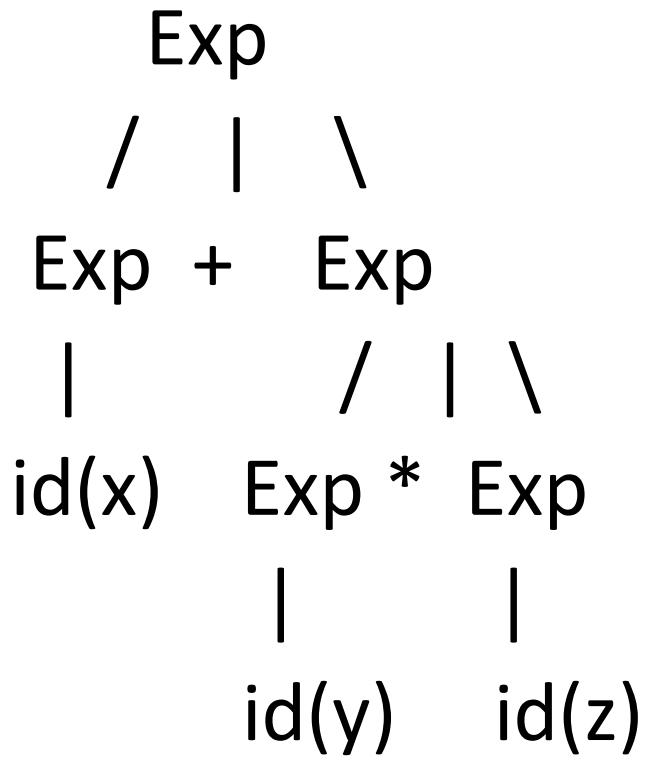
 Exp y * z
 / | \
Exp + Exp
 |
id(x)

Top Down Parse (cont)

Exp y * z
/ | \
Exp + Exp
| / | \
id(x) Exp * Exp

Exp z
/ | \
Exp + Exp
| / | \
id(x) Exp * Exp
 |
 id(y)

Top Down Parse (cont.)



Problems with Top Down Parsing

- Backtracking may be necessary:
 - $S ::= ee \mid bAc \mid bAe$
 - $A ::= d \mid cA$try on string “bcde”
- Infinite loops possible from (indirect) left recursive grammars.
 - $E ::= E + id \mid id$
- Ambiguity is a problem when a unique parse is not possible.
- These often require extensive grammar restructuring (grammar debugging).

Grammar Transformations

- Backtracking and Factoring
- Removing ambiguity.
 - Simple grammars are often easy to write, but might be ambiguous.
- Removing Left Recursion

Backtracking and Factoring

- Backtracking may be necessary:

$$\begin{array}{l} S \rightarrow ee \quad | \quad bAc \quad | \quad bAe \\ A \rightarrow d \quad | \quad cA \end{array}$$

- try on string “bcde”

$$\begin{array}{l} S \rightarrow bAc \quad \text{(by } S \rightarrow bAc) \\ \rightarrow bcAc \quad \text{(by } A \rightarrow cA) \\ \rightarrow bc dc \quad \text{(by } A \rightarrow d) \end{array}$$

- But now we are stuck, we need to backtrack to
 - $S \rightarrow bAc$
 - And then apply the production ($S \rightarrow bAe$)

- Factoring a grammar

- Factor common prefixes and make the different postfixes into a new non-terminal

$$\begin{array}{l} S \rightarrow ee \quad | \quad bAQ \\ Q \rightarrow c \quad | \quad e \\ A \rightarrow d \quad | \quad cA \end{array}$$

Removing ambiguity.

- Adding levels to a grammar

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

Transform to an equivalent grammar

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

Levels make formal the notion of precedence. Operators that bind “tightly” are on the lowest levels

The dangling else grammar.

- $st \rightarrow$ if exp then st else st
| if exp then st
| id := exp
- Note that the following has two possible parses
if x=2 then if x=3 then y:=2 else y := 4

if x=2 then (if x=3 then y:=2) else y := 4
if x=2 then (if x=3 then y:=2 else y := 4)

Adding levels (cont)

- Original grammar

```
st ::=  if exp  then st  else st
      |  if exp then st
      |  id := exp
```

- Assume that every `st` between then and else must be matched, i.e. it must have both a then and an else.
- New Grammar with additional levels

```
st      -> match | unmatched
match   -> if exp then match else match
        |  id := exp
unmatch -> if exp then st
        |  if exp then match else unmatched
```

Top Down Recursive Descent Parsers

- One function (procedure) per non-terminal
- Functions call each other in a mutually recursive way.
- Each function “consumes” the appropriate input.
- If the input has been completely consumed when the function corresponding to the start symbol is finished, the input is parsed.

Example Recursive Descent Parser

$E \rightarrow T + E \mid T$

$T \rightarrow F * T \mid F$

$F \rightarrow x \mid (E)$

```
expr =  
  do { term  
      ; iff (match '+') expr }
```

```
term =  
  do { factor  
      ; iff (match '*') term }
```

```
factor =  
  pCase  
  [ 'x' :=> return ()  
  , '(' :=> do { expr; match ')'; return () }  
  ]
```


Removing Left Recursion

- Top down recursive descent parsers require non-left recursive grammars
- **Technique: Left Factoring**

$$E \rightarrow E + E \quad | \quad E * E \quad | \quad id$$
$$E \rightarrow id E'$$
$$E' \rightarrow + E E'$$
$$| * E E'$$
$$| \Lambda$$

General Technique to remove direct left recursion

- Every Non terminal with productions

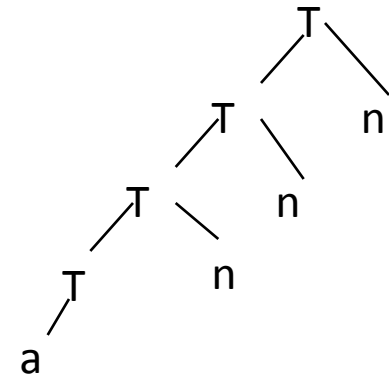
$$T \rightarrow T n \mid T m \quad (\text{left recursive productions})$$
$$\mid a \mid b \quad (\text{non-left recursive productions})$$

- Make a new non-terminal T'
- Remove the old productions
- Add the following productions

$$T \rightarrow a T' \mid b T'$$
$$T' \rightarrow n T' \mid m T' \mid \Lambda$$

“a” and “b” because they are the rhs of the non-left recursive productions.

$(a \mid b) (n \mid m)^*$



Recursive Descent Parsing

- One procedure (function) for each non-terminal.
- Procedures are often (mutually) recursive.
- They can return a bool (the input matches that non-terminal) or more often they return a data-structure (the input builds this parse tree)
- Need to control the lexical analyzer (requiring it to “back-up” on occasion)

A grammar suitable for Recursive Descent parsers for Regular Expressions

- Build an instance of the datatype:

```
data RegExp a
  = Lambda                -- the empty string ""
  | Empty                 -- the empty set
  | One a                 -- a singleton set {a}
  | Union (RegExp a) (RegExp a) -- union of two RegExp
  | Cat (RegExp a) (RegExp a)  -- Concatenation
  | Star (RegExp a)          -- Kleene closure
```

Ambiguous grammar

```
RE -> RE bar RE
RE -> RE RE
RE -> RE *
RE -> id
RE -> ^
RE -> ( RE )
```

- Transform grammar by layering
- Tightest binding operators (*) at the lowest layer
- Layers are Alt, then Concat, then Closure, then Simple.

```
Alt -> Alt bar Concat
Alt -> Concat
Concat -> Concat Closure
Concat -> Closure
Closure -> simple star
Closure -> simple
simple -> id | ( Alt ) | ^
```

Left Recursive Grammar

```
Alt -> Alt bar Concat
Alt -> Concat
Concat -> Concat Closure
Concat -> Closure
Closure -> simple star
Closure -> simple
simple -> id | (Alt ) | ^
```

For every Non terminal with productions

$$T ::= T_n \mid T_m \quad (\text{left recursive prods})$$
$$\mid a \mid b \quad (\text{non-left recursive prods})$$

Make a new non-terminal T'

Remove the old productions

Add the following productions

$$T ::= a T' \mid b T'$$
$$T' ::= n T' \mid m T' \mid \Lambda$$

```
Alt          -> Concat moreAlt
moreAlt      -> Bar Concat moreAlt
              |  Λ
Concat       -> Closure moreConcat
moreConcat   -> Closure moreConcat
              |  Λ
Closure      -> Simple Star
              |  Simple
Simple       -> Id
              |  ( Alt )
              |  ^
```

Predictive Parsers

- Using a stack to avoid recursion. Encoding the diagrams in a table
- The Nullable, First, and Follow functions
 - Nullable: Can a symbol derive the empty string. False for every terminal symbol.
 - First: all the terminals that a non-terminal could possibly derive as its first symbol.
 - term or nonterm \rightarrow set(term)
 - sequence(term + nonterm) \rightarrow set(term)
 - Follow: all the terminals that could immediately follow the string derived from a non-terminal.
 - non-term \rightarrow set(term)

Example First and Follow Sets

$$\begin{array}{lcl}
 E & \rightarrow & T E' \$ \\
 E' & \rightarrow & + T E' \\
 E' & \rightarrow & \Lambda \\
 T & \rightarrow & F T' \\
 T' & \rightarrow & * F T' \\
 T' & \rightarrow & \Lambda \\
 F & \rightarrow & (E) \\
 F & \rightarrow & id
 \end{array}$$

First E = { "(", "id" }	Follow E = { ")", "\$" }
First F = { "(", "id" }	Follow F = { "+", "*", ")", "\$" }
First T = { "(", "id" }	Follow T = { "+", ")", "\$" }
First E' = { "+", ϵ }	Follow E' = { ")", "\$" }
First T' = { "*", ϵ }	Follow T' = { "+", ")", "\$" }

- First of a terminal is itself.
- First can be extended to sequence of symbols.

Nullable

- if Λ is in $\text{First}(\text{symbol})$ then that symbol is nullable.
- Sometime rather than let Λ be a symbol we derive an additional function nullable.

- $\text{Nullable}(E') = \text{true}$

- $\text{Nullable}(T') = \text{true}$

- Nullable for all other symbols is false

E	->	T	E'	\$
E'	->	+	T	E'
E'	->	Λ		
T	->	F	T'	
T'	->	*	F	T'
T'	->	Λ		
F	->	(E)
F	->	id		

Computing First

- Use the following rules until no more terminals can be added to any FIRST set.
 - 1) if X is a term. $FIRST(X) = \{X\}$
 - 2) if $X \rightarrow \Lambda$ is a production then add Λ to $FIRST(X)$, (Or set nullable of X to true).
 - 3) if X is a non-term and
 - $X \rightarrow Y_1 Y_2 \dots Y_k$
 - add a to $FIRST(X)$
 - if a in $FIRST(Y_i)$ and
 - for all $j < i$ Λ in $FIRST(Y_j)$
- E.g.. if Y_1 can derive Λ then if a is in $FIRST(Y_2)$ it is surely in $FIRST(X)$ as well.

Example First Computation

- Terminals
 - $\text{First}(\$) = \{\$\}$, $\text{First}(\ast) = \{\ast\}$, $\text{First}(+) = \{+\}$, ...
- Empty Productions
 - add Λ to $\text{First}(E')$, add Λ to $\text{First}(T')$
- Other NonTerminals
 - Computing from the lowest layer (F) up
 - $\text{First}(F) = \{\text{id}, (\}$
 - $\text{First}(T') = \{\Lambda, \ast\}$
 - $\text{First}(T) = \text{First}(F) = \{\text{id}, (\}$
 - $\text{First}(E') = \{\Lambda, +\}$
 - $\text{First}(E) = \text{First}(T) = \{\text{id}, (\}$

E	->	T E' \$
E'	->	+ T E'
E'	->	Λ
T	->	F T'
T'	->	\ast F T'
T'	->	Λ
F	->	(E)
F	->	id

Computing Follow

- Use the following rules until nothing can be added to any follow set.
- 1) Place \$ (the end of input marker) in FOLLOW(S) where S is the start symbol.
 - 2) If $A \rightarrow a B b$
then everything in $FIRST(b)$ except Λ is in FOLLOW(B)
 - 3) If there is a production $A \rightarrow a B$
or $A \rightarrow a B b$ where $FIRST(b)$
contains Λ (i.e. b can derive the empty string) then
everything in FOLLOW(A) is in FOLLOW(B)

Ex. Follow Computation

- Rule 1, Start symbol
 - Add \$ to Follow(E)
- Rule 2, Productions with embedded nonterms
 - Add First() = {) } to follow(E)
 - Add First(\$) = { \$ } to Follow(E')
 - Add First(E') = { +, Λ } to Follow(T)
 - Add First(T') = { *, Λ } to Follow(F)
- Rule 3, Nonterm in last position
 - Add follow(E') to follow(E') (doesn't do much)
 - Add follow(T) to follow(T')
 - Add follow(T) to follow(F) since $T' \rightarrow \Lambda$
 - Add follow(T') to follow(F) since $T' \rightarrow \Lambda$

E	->	T	E'	\$
E'	->	+	T	E'
E'	->	Λ		
T	->	F	T'	
T'	->	*	F	T'
T'	->	Λ		
F	->	(E)
F	->	id		

Table from First and Follow

1. For each production $A \rightarrow \alpha$ do 2 & 3
2. For each a in First α do add $A \rightarrow \alpha$ to $M[A,a]$
3. if ϵ is in First α , add $A \rightarrow \alpha$ to $M[A,b]$ for each terminal b in Follow A . If ϵ is in First α and $\$$ is in Follow A add $A \rightarrow \alpha$ to $M[A,\$]$.

First $E = \{(","id")\}$ Follow $E = \{")","$\}$
 First $F = \{(","id")\}$ Follow $F = \{"+","*","(",",$")\}$
 First $T = \{(","id")\}$ Follow $T = \{"+","(",",$")\}$
 First $E' = \{"+",\epsilon\}$ Follow $E' = \{")","$\}$
 First $T' = \{"+",\epsilon\}$ Follow $T' = \{"+","(",",$")\}$

1.	E	\rightarrow	T	E'	$\$$
2.	E'	\rightarrow	$+$	T	E'
3.	E'	\rightarrow	Λ		
4.	T	\rightarrow	F	T'	
5.	T'	\rightarrow	$*$	F	T'
6.	T'	\rightarrow	Λ		
7.	F	\rightarrow	$($	E	$)$
8.	F	\rightarrow	id		

M[A,t] terminals

		+	*)	(id	\$
n	E				1	1	
o	E'	2		3			3
n	T				4	4	
t	T'	6	5	6			6
e	F				7	8	
r							
m							
s							

Predictive Parsing Table

	id	+	*	()	\$
E	TE'			TE'		
E'		+TE'			ϵ	ϵ
T	FT'			FT'		
T'		ϵ	*FT'		ϵ	ϵ
F	id			(E)		

Table Driven Algorithm

push start symbol

Repeat

begin

let X top of stack, A next input

if terminal(X)

then if X=A

then pop X; remove A

else error()

else (* nonterminal(X) *)

begin

if $M[X,A] = Y_1 Y_2 \dots Y_k$

then pop X;

push $Y_k Y_{k-1} \dots Y_1$

else error()

end

until stack is empty, input = \$

Example Parse

Stack

E
 E' T
 E' T' F
 E' T' id
 E' T'
 E'
 E' T +
 E' T
 E' T' F
 E' T' id
 E' T'
 E'

Input

x + y \$
 x + y \$
 x + y \$
 x + y \$
 + y \$
 + y \$
 + y \$
 y \$
 y \$
 y \$
 \$
 \$
 \$

	id	+	*	()	\$
E	TE'			TE'		
E'		+TE'			ε	ε
T	FT'			FT'		
T'		ε	* FT'		ε	ε
F	id			(E)		