

Transforming Grammars

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

Grammar Transformations

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

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

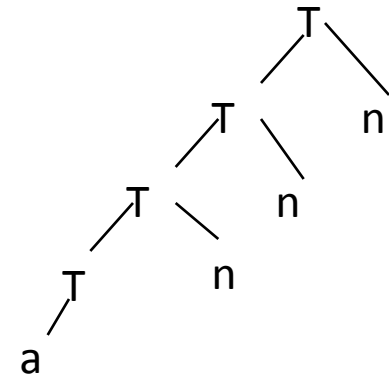
$$T \rightarrow T n \mid T m \quad (\text{left recursive productions})$$
$$\mid a \mid b \quad (\text{non-left recursive productions})$$

- 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.

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



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
```

Removing Λ -productions

- It is possible to write every CFL that does not contain Λ (the empty string) without any Λ in any RHS
- $S \rightarrow aDaE$
- $D \rightarrow bD \mid E$
- $E \rightarrow cE \mid \Lambda$

Rules

1. Find all non-terminal, N , such that N derives Λ
2. For each production, $A \rightarrow w$, create new productions, $A \rightarrow w'$, where w' is derived from w by removing non-terminals that derive Λ (found in rule 1 above)
3. Create a new grammar from the original productions, together with the productions formed in step 2, removing any productions of the form, $A \rightarrow \Lambda$.

$S \rightarrow aDaE$

$D \rightarrow bD \mid E$

$E \rightarrow cE \mid \Lambda$

1. Find all non-terminal, N , such that N derives Λ

$E \rightarrow \Lambda$

2. For each production, $A \rightarrow w$, create new productions, $A \rightarrow w'$, where w' is derived from w by removing non-terminals that derive Λ (found in rule 1 above)

$S \rightarrow aDa$

$D \rightarrow \Lambda$

$E \rightarrow c$

3. Create a new grammar from the original productions, together with the productions formed in step 2, removing any productions of the form, $A \rightarrow \Lambda$.

$S \rightarrow aDaE$

$D \rightarrow bD \mid E$

$E \rightarrow cE$

$S \rightarrow aDa$

$E \rightarrow c$

Top to bottom example

- Start with an easy (to understand) grammar
- Transform it to one that is easier to parse
- Apply some of the transformation rules

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

A datatype suitable for representing 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 )
              |  ^
```