

# Lec 22

## Context Free Grammar

ID ↪ Human  
 Lexing ↪ list of char  
 ↪ list of tokens  
 parsing ↪ expression

fact(3)  
 (ID "fact") :: LPAREN :: INT 3 :: RPAREN  
 App(Var "fact", Num 3)

← abstract syntax tree

### # Chomsky hierarchy

level	representation	application
General grammar	turing machine	computing
Context-sensitive grammar	whatever	whatever
Context-free grammar	pushdown automata	parsing
Regular grammar	finite automata	lexing / search

### # Context free grammar

← programme

$$P \rightarrow \epsilon \mid E ; P$$

← empty

$$E \rightarrow E + E \mid E * E \mid E \text{ and also } E \mid \text{case } E \text{ of } M \mid \dots$$

← expression

$$M \rightarrow Q \Rightarrow E \mid Q \Rightarrow E \mid M$$

$$Q \rightarrow () \mid \dots$$

Def A context free grammar (CFG) is a tuple  $(\Sigma, V, S, R)$  s.t.

- $\Sigma$  is set of terminals viz. chars we see in program,
- $V$  is set of non-terminals ( $\Sigma \cap V = \emptyset$ )
- $S \in V$  is the start symbol
- $R$  is set of rules of form  $N \rightarrow w$  where  $N \in V, w \in (\Sigma \cup V)^*$  ← star from regex

Suppose  $\alpha, \beta \in (\Sigma \cup V)^*$

- $\beta$  is derivable from  $\alpha$  in one step ( $\alpha \Rightarrow \beta$ ) if:  
 $\exists \delta, \epsilon$  s.t.  $\alpha = \delta N \delta$  and  $\beta = \delta \omega \delta$  where  $N \rightarrow \omega \in R$
- $\beta$  is derivable from  $\alpha$  if  $\alpha = \beta$  or  $\alpha \Rightarrow \sigma_1 \Rightarrow \sigma_2 \Rightarrow \dots \Rightarrow \beta$
- $L(G) = \{ \omega \in \Sigma^* \mid S \Rightarrow \omega \}$

Ex. Let  $G = (\Sigma, V, S, R)$

$\Sigma = \{a, b\}$

$S = \{S, A\}$

$R = \left\{ \begin{array}{l} S \rightarrow AbA \\ A \rightarrow \epsilon \\ A \rightarrow a \\ A \rightarrow aA \end{array} \right\}$  ← actually redundant

Derive

$S \Rightarrow AbA \Rightarrow abA \Rightarrow abaA \Rightarrow aba$  ← left most derivation: always expand left most  
 $S \Rightarrow AbA \Rightarrow AbaA \Rightarrow Aba \Rightarrow aba$

using rule 3. Multiple ways. Then we get multiple parse trees: (may cause inconsistency) ← Ambiguous

Def A grammar is ambiguous if it has multiple left-most derivation

- Figuring out if a CFG is ambiguous is undecidable
- There are languages s.t. all CFG to describe it is ambiguous

Thm All regex can be expressed as CFG

0 -  $(\Sigma, \{S\}, S, \emptyset)$

1 -  $(\Sigma, \{S\}, S, \{S \rightarrow \epsilon\})$

a - ( $\Sigma, \{S\}, S, \{S \rightarrow a\}$ )

⋮

Ex.  $\{a^n \mid n \equiv 0 \pmod{3}, n \geq 0\}$  ← regular. finite storage to accept.  
 $\{a^n b^n \mid n \geq 0\}$  ← not regular, but context free. finite storage + stack to accept  
 $\{a^n b^n c^n \mid n \geq 0\}$  ← not context free. Need finite storage + at least two stacks  
 $\{a^n b^m c^m d^n \mid m, n \geq 0\} \cup \{a^n b^n c^m d^m \mid m, n \geq 0\}$  ← context free, but always ambiguous

Pumping Lemma — for proving some lang is not in a Chomsky class

If  $L$  is infinite regular language,  $\exists a, w, \beta$  s.t.  
 $w \neq \epsilon$  and  $aw^k\beta \in L$  for every  $k$  viz. there's some loop

Consider  $\{a^n b^n \mid n \geq 0\}$ . Suppose it's regular.  
Then  $\exists w \in \epsilon, aw^k\beta \in L$  for all  $k$ .  
idea is that we can't repeat any  $ab, aabb, \dots$

Consider  $S \rightarrow \epsilon \mid aSa \mid bSb \mid a \mid b$

# Simple Parser ← (usually we can use parser generator these days)

← probably better  
→ Bottom up: try rebuild to  $S$  from string  
→ Top down: start from  $S$  and try make target string  
← for now

$G_1 = \{S \rightarrow X \mid \lambda X.S \mid (SS)\}$

datatype token = LAMBDA | LPAREN | RPAREN | ID of string | DOT  
datatype exp = Var of string  
| Lam of string \* exp  
| App of exp \* exp

```

(* parseExp : token list → (exp * token list → 'a) → 'a
fun parseExp (ID x::ts) k = Var (x, ts)
  | parseExp (LPAREN::ts) k =
    parseExp ts (fn (e1, ts') ⇒
      parseExp ts' (fn (e2, RPAREN::ts'') ⇒
        k (App (e1, e2), ts''))
    )
  | parseExp (LAMBDA::ID x::DOT::ts) k =
    parseExp ts (fn (e, ts') ⇒ k (Lam (x, e), ts'))
  | ... - - - = raise SyntaxError

```

