

The Cocke-Kasami-Younger algorithm

Given  $w \in \Sigma^*$  & G a CFG

Want to know  $w \in L(G)$ ?

We assume G is in Chomsky Normal Form. The parse trees are binary so for a string of length n we will have a tree with  $(2n-1)$  variables. There are exponentially many such trees so we can generate them all, check if they are valid trees that generate w. Exponential time!

We can get this down to  $O(n^3)$  using dynamic programming.

Given  $G = (V, \Sigma, P, S)$

Input  $w = a_1 \dots a_n \in \Sigma^*$   $a_i \in \Sigma$

We work bottom-up to construct a possible derivation of w. We first ask how we got each individual symbol.

Since G is in CNF we have to have used rules of the form  $A \rightarrow a$

We define inductively a 2-indexed family of subsets of V:

$$X_{ij} := \{ A \in V \mid A \xrightarrow{*} a_i \dots a_j \}$$

BASE CASE  $X_{ii} = \{ A \in V \mid A \Rightarrow a_i \}$   $x_{11}, x_{22}, \dots x_{nn}$

Next row will have  $x_{12}, x_{23}, \dots x_{i(i+1)}, \dots x_{(n-1)n}$

Next row will have  $x_{13}, x_{24}, \dots x_{i(i+2)}, \dots x_{(n-2)n}$

so  $X_{ij}$  will be in row  $(j-i)+1$ .

We compute and fill in the table ~~top~~ bottom to top.

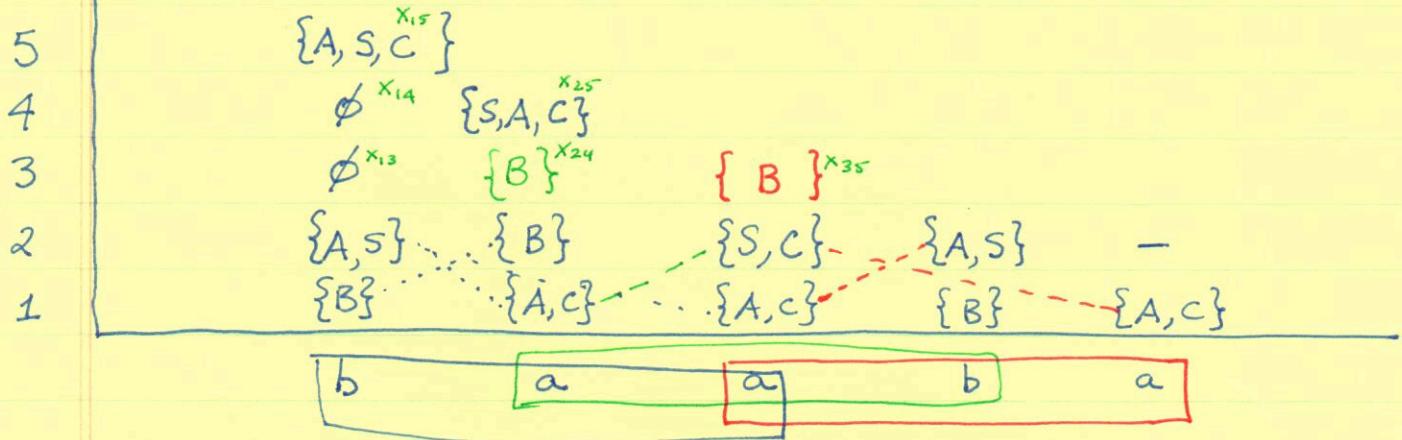
When we compute  $X_{ij}$  we know  $X_{ik}$  &  $X_{kj}$  for all  $i \leq k \leq j$

Now if  $B \in X_{ik}$ ,  $C \in X_{(k+1)j}$  &  $A \rightarrow BC$  is a rule  
we know  $A \in X_{ij}$ . Why?  $B \xrightarrow{*} a_i \dots a_k$   
 $C \xrightarrow{*} a_{k+1} \dots a_j$  so since  $A \rightarrow BC \xrightarrow{*} a_i \dots a_j$ .

(2)

$$S \rightarrow AB \mid BC \quad A \rightarrow BA \mid a \quad B \rightarrow CC \mid b \quad C \rightarrow AB \mid a$$

We want to know  $baaba \in L(G)$ ?



We see that  $S \in X_{15}$  so

$$S \xrightarrow{*} a_1 \dots a_5$$

In general  $w \in L(G)$  iff  $S \in X_{1n}$

Size of table  $O(n^2)$

The time taken to find  $X_{ij}$  is  $O(j-i)$

The time taken to compare  $X_{ik} \& X_{(k+1)j}$  & find a variable that generates them is  $O(1)$ : it depends on the size of the grammar but not on n. So time to compute each  $X_{ij}$  is  $O(n)$  & so overall  $O(n^3)$ .

There are better algorithms possible. Best is  $O(n^{2.8})$ .

$\longrightarrow x \longrightarrow$

Is a CFL empty? Given  $G$  we want to know  $L(G) = \emptyset$ ?

~~Call a variable  $A$  productive if there is a rule  $A \Rightarrow x$  where  $x \in \Sigma^*$~~

Call a variable  $A$  productive if  $A \xrightarrow{*} x$  where  $x \in \Sigma^*$ .

Check every variable  $A$  & see if  $\exists A \rightarrow x$  with  $x \in \Sigma^*$

If so label  $A$  productive. Then repeat.

(3)

go through all variables  $B$  in  $V - \text{Prod}$  & check if there is a rule  $B \rightarrow \beta$  where  $\beta \in (\Sigma \cup \text{Prod})^*$ . If yes add  $B$  to Prod. Keep repeating until there are no more changes. If  $S \in \text{Prod}$  then  $L(G) \neq \emptyset$ , & if  $S \notin \text{Prod}$  then  $L(G) = \emptyset$ .

Naively this  $O(n^2)$  where  $n$  is the # of variables. A more careful implementation does it in  $O(n)$ .

### Parsing Hugely practical

Extract the parse tree & build it up.

We design DCFL in order to do this :  $LR(k)$ ,  $LL(k)$ ,  $LALR(k)$ , Rec-descent etc.

$S \rightarrow (S)S \mid \epsilon$  Unambiguous & deterministic

This is really an inductive definition so a parser for this is naturally recursive:

parse( $\epsilon$ ):

if  $x = \epsilon$  return nil

else if getsym = '(' then parse {

$t_0 := \text{parse}();$

if getsym = ')'

$t_1 := \text{parse}();$  ~~return mktree('(', t<sub>0</sub>, ')')~~

else error

else error;

$t_2 := \text{parse}();$

if  $t_2 = \text{nil}$  return  $t_1$

else return mktree('(', t<sub>0</sub>, ')', t<sub>2</sub>)

$$\{0^i 1^j \mid j = i^2\}$$

$$A \rightarrow p$$

$$I \rightarrow 0^p 1^{p^2}$$

$$A \rightarrow uvwxy = 0^p 1^{p^2}; |vwx| \leq p \quad |vx| > 0 \quad \text{**}$$

I → case analysis

0...01-----1,

- (a)  $vwx \rightarrow$  all zeros ; pick any  $i \neq 1$
- (b)  $vwx \rightarrow$  all ones, pick any  $i \neq 1$
- (c)  $V$  overlaps 0's & 1's,  $x$  all ones; pick any  $i \neq 1$ , 0's 1's out of order same story with  $\Rightarrow$  all zero  $x$  overlaps
- (d)  $V$  all zeros,  $x$  all ones pick  $i = 2$

$$|v| = m > 0, |x| = q > 0$$

$$0^{p+m} 1^{p^2+q},$$

Question is  $(p+m)^2 = p^2 + q^2$  ?

$$p+m^2 = p^2 + 2pm + m^2$$

$$2pm > p \quad \text{so } 2pm + \cancel{m^2} > p \cancel{+} q$$

Thus we have left the language.

So  $0^p 1^{p^2}$  is not context free.

Suppose  $L \subseteq a^*$  is not regular.  $\Rightarrow L$  not context free

$A \rightarrow p$ I choose $s$ . $A \rightarrow s = xyz \quad  y  \leq p$	$A \rightarrow p$ I choose as before $s \rightarrow uvwxy \quad  vx  \leq p$
---	--

whatever I did for the regular case can be mimiced!!

Thus  $L$  not regular  $\Rightarrow L$  not context-free

or  $L$  context-free  $\Rightarrow L$  regular