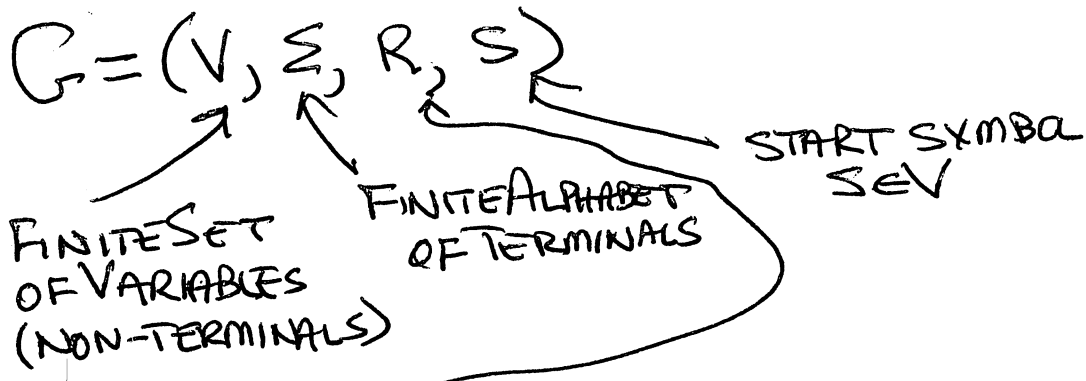


CONTEXT FREE GRAMMARS



RULES, EACH OF FORM

$$A \rightarrow \alpha \quad A \in V, \alpha \in (V \cup \Sigma)^*$$

A- RULES CAN BE LISTED ON SEPARATE LINES OR USING | TO INDICATE ALTERNATIVES, E.G.

$$S \rightarrow a S b \mid \lambda \quad \{a^n b^n \mid n \geq 0\}$$

OR

$$\begin{aligned} S &\rightarrow a S b \\ S &\rightarrow \lambda \end{aligned}$$

OR

$$\begin{aligned} S &\rightarrow a S b \\ &\mid \lambda \end{aligned}$$

SOME EXAMPLE CFLs

$$L_1 = \{a^n b^n \mid n \geq 0\}$$

$$G_1 = (\{S\}, \{a, b\}, R, S)$$

$$R: S \rightarrow aSb \mid \lambda$$

$$L_2 = \{a^n b^m \mid m \geq n\}$$

$$S \rightarrow aSb \mid Sb \mid \lambda$$

$$L_3 = \{a^n b^m \mid m > n\}$$

$$S \rightarrow aSb \mid Sb \mid b$$

$$L_4 = \{w w^R \mid w \in \{a, b\}^*\}$$

$$S \rightarrow aSa \mid bSb \mid \lambda$$

NOTE $\{w w \mid w \in \{a, b\}^*\}$

IS NOT A CFL

$$L_5 = \{w \mid \#a\text{'s in } w = \#b\text{'s in } w\}$$

$$S \rightarrow SaSbS \mid SbSaS \mid \lambda$$

CAN
OMIT

PROOF THAT A GRAMMAR GENERATES DESIRED LANGUAGE

$$L = \{ w \mid w \in \{a,b\}^* \text{ \& } w_a = w_b \}$$

HERE w_a IS # OF a's IN w
 w_b IS # OF b's IN w

$$G = (\{S\}, \{a,b\}, \mathbb{R}, S)$$

$$R: S \rightarrow aSbS \mid bSaS \mid \lambda$$

$$\text{PROVE } L(G) = L$$

1. FIRST $L(G) \subseteq L$
AS ALL RHS OF S HAVE EQUAL #
OF a's \& b's

$$w \in L(G) \Rightarrow w \in L$$

2. NOW SHOW $L \subseteq L(G)$

DO BY INDUCTION ON $|w|$ TO SHOW

$$w \in L \Rightarrow w \in L(G)$$

INDUCTIVE PROOF

BASIS: $|w|=0$. SHOW $\lambda \in \mathcal{L}(G)$
AS $S \rightarrow \lambda \in R$ THEN $S \Rightarrow \lambda$ AND λ
IS ONLY STRING OF LENGTH 0.

IH: FOR ALL $L \subseteq R$, $k \geq 0$,
IF $w \in L$, $|w|=2k$ THEN $w \in \mathcal{L}(G)$
NOTE: ALL ELEMENTS OF L ARE OF EVEN LENGTH

(S): SHOW FOR $|w|=2k+2$
ASSUME w STARTS WITH AN 'a'
THEN $w = a\alpha b\beta$
WHERE $\alpha \in L$, $\beta \in L$ AND α IS SHORTEST
STRING LEADING TO MATCHING 'b' FOR FIRST 'a'
AS $|\alpha| \leq 2k$ AND $|\beta| \leq 2k$
 $S \xrightarrow{*} \alpha$, $S \xrightarrow{*} \beta$ BY IH
SINCE $S \rightarrow a s b s \in R$
 $S \Rightarrow a s b s \xrightarrow{*} a \alpha b s \xrightarrow{*} a \alpha b \beta$
AND SO $S \Rightarrow w$ AND $w \in \mathcal{L}(G)$
CAN REDO FOR STARTING 'b'
THUS, $w \in L \Rightarrow w \in \mathcal{L}(G)$

COMBINING (1) AND (2)

$$w \in L \Leftrightarrow w \in \mathcal{L}(G)$$

∴ SO

$$L = \mathcal{L}(G)$$

ONE MORE PROOF WRT CFGS, $G = (\{P\}, \{0, 1\}, R, P)$

R: $P \rightarrow \lambda \mid 0 \mid 1 \mid 0P0 \mid 1P1$

$$\mathcal{L}(G) = \{w \mid w \text{ IS A BINARY PALINDROME}\}$$

$$= \{w \mid w \in \{0, 1\}^* \wedge w = w^R\}$$

PROOF: "ALL BINARY PALINDROMES ARE IN $\mathcal{L}(G)$ "

BASIS: $|w| = 0$ OR $|w| = 1$

$$P \rightarrow \lambda \text{ SO } P \xrightarrow{*} \lambda \Rightarrow \lambda \in \mathcal{L}(G)$$

$$P \rightarrow 0 \text{ SO } P \xrightarrow{*} 0 \Rightarrow 0 \in \mathcal{L}(G)$$

$$P \rightarrow 1 \text{ SO } P \xrightarrow{*} 1 \Rightarrow 1 \in \mathcal{L}(G)$$

AND SO BASIS IS SHOWN

IH: ASSUME FOR ALL w , $|w| \leq k, k \geq 1$,
 $P \xrightarrow{*} w$ WHENEVER $w = w^R$

IS: LET $|w| = k+1$. SINCE $k \geq 1, |w| \geq 2$ & SO
 $w = 0x0$ OR $w = 1x1$ FOR $|x| = k-1$ AND $x = x^R$

BUT THEN

$$P \Rightarrow 0P0 \xrightarrow{*} 0x0 \quad \text{OR} \quad P \Rightarrow 1P1 \xrightarrow{*} 1x1 \quad \text{IND. HYP.}$$

& SO $P \xrightarrow{*} w$

PROOF: 2 SHOW $L(G)$ CONTAINS ONLY PALINDROMES
DO INDUCTION OF # OF STEPS IN DERIVATION

BASIS: $k=1$

EITHER $P \Rightarrow \lambda$ OR $P \Rightarrow 0$ OR $P \Rightarrow 1$ AND
ALL ARE PALINDROMES

IH: ASSUME FOR ALL $i \leq k, k \geq 1$,
IF $P \xRightarrow{i} X$ AND $X \in \{0,1\}^*$ THEN
 X IS AN PALINDROME

IS: $P \Rightarrow 010 \xRightarrow{k} 0X0$ } $X \in \{0,1\}^*$
OR $P \Rightarrow 111 \xRightarrow{k} 1X1$

ARE ONLY $k+1$ LENGTH DERIVATIONS
OF TERMINAL STRINGS.

SINCE $P \xRightarrow{k} X$ IS A PALINDROME
THEN SO ARE $0X0$ AND $1X1$

~~(QED)~~

A CHALLENGING CFL

THE COMPLEMENT OF $\{ww \mid w \in \{a,b\}^*\}$
IS A CFL. IT CAN BE DESCRIBED
AS

$$\{xy \mid |x|=|y|, x \neq y\} \cup \{w \mid |w| \text{ IS ODD}\}$$

WHERE ALL STRINGS ARE OVER $\{a,b\}$

$\{w \mid w \in \{a,b\}^+ \text{ AND } |w| \text{ IS ODD}\}$
IS REGULAR

$$S \rightarrow aT \mid bT$$

$$T \rightarrow aS \mid bS \mid \lambda$$

WE WILL SHOW CFLS ARE CLOSED UNDER \cup

FOCUS ON

$$\{xy \mid |x|=|y|, x \neq y, x,y \in \{a,b\}^*\}$$

CHALLENGE CONTINUED

$$\{xy \mid |x| = |y|, x \neq y\}$$

NOTE THAT WHILE WW REQUIRES THAT THERE BE NO ERRORS TRANSCRIBING (COPYING) FIRST HALF TO SECOND HALF, ITS COMPLEMENT REQUIRES JUST ONE TRANSCRIPTION ERROR. ELEMENTS ARE OF FORM

$$x_1 a x_2 y_1 b y_2 \text{ OR } x_1 b x_2 y_1 a y_2$$

WHERE $|x_1| = |y_1|$ AND $|x_2| = |y_2|$

NOTE THAT WE DO NOT CARE WHAT x_1, x_2, y_1 AND y_2 ARE. ONLY THEIR LENGTHS MATTER. WITH THAT INSIGHT IN MIND, WE CAN SEE THAT

$$x_1 a y_1 x_2 b y_2 \text{ OR } x_1 b y_1 x_2 a y_2$$

WHERE $|x_1| = |y_1|$ AND $|x_2| = |y_2|$ ALSO DESCRIBES ELEMENTS IN LANGUAGE. THIS IS JUST

$$\left. \begin{array}{l} S \rightarrow AB \mid BA \\ A \rightarrow CAC \mid a \\ B \rightarrow CBC \mid b \\ C \rightarrow a \mid b \end{array} \right\} (\{S, A, B, C\}, \{a, b\}, R, S)$$

WHICH IS A CFG.

SIMPLE CLOSURES

LET $G_1 = (V_1, \Sigma, R_1, S_1)$ AND $G_2 = (V_2, \Sigma, R_2, S_2)$

WHERE $V_1 \cap V_2 = \emptyset$

LET $L_1 = \mathcal{L}(G_1)$ AND $L_2 = \mathcal{L}(G_2)$

WHERE G_1 AND G_2 ARE CFGS AND SO

L_1 AND L_2 ARE CFLS

UNION: $G_3 = (V_1 \cup V_2 \cup \{S\}, \Sigma, R_3, S)$

$R_3: S \rightarrow S_1 \mid S_2$

ADD R_1 AND R_2 TO R_3 AND GET

$L_1 \cup L_2$

CONCATENATION:

$R_3: S \rightarrow S_1 S_2$

ADD R_1 AND R_2 TO R_3 AND GET

$L_1 \circ L_2$

STAR:

$R_3: S \rightarrow S S_1 \mid \lambda$

ADD R_1 TO R_3 AND GET

L_1^*

NOTE: $S \rightarrow S_1 S \mid \lambda$ ALSO WORKS

NON-CLOSURES

ASSUME, FOR NOW (PUMPING LEMMA WILL VERIFY) THAT

$$L = \{a^n b^n c^n \mid n \geq 0\}$$

IS NOT A CFL

NOW, CONSIDER TWO CFLS

$$L_1 = \{a^n b^n c^m \mid n, m \geq 0\}$$

$$= a^n b^n c^*$$

$$\text{AND } L_2 = \{a^m b^n c^n \mid n, m \geq 0\}$$

$$= a^* b^n c^n$$

$$L_1 \cap L_2 = L = \{a^n b^n c^n \mid n \geq 0\}$$

BUT $L_1 = \mathcal{L}(G_1)$, $L_2 = \mathcal{L}(G_2)$ WHERE

$$G_1 = (\{S, T\}, \{a, b, c\}, R_1, S) \quad G_2 = (\{S, T\}, \{a, b, c\}, R_2, S)$$

$$R_1: S \rightarrow Sc \mid T \\ T \rightarrow aTb \mid \lambda$$

$$R_2: S \rightarrow aS \mid T \\ T \rightarrow bTc \mid \lambda$$

SO CFLS NOT CLOSED UNDER \cap
PREVIOUSLY. WE SHOWED NOT UNDER
COMPLEMENT, IF WE CAN SHOW

$$L = \{ww \mid w \in \{a, b\}^*\}$$

IS NOT A CFL (LATER)

TYPES OF DERIVATIONS

- (1) NO CONSTRAINT ON WHICH NON-TERMINAL TO REWRITE WHEN CHOICES EXIST
- (2) RIGHTMOST: ALWAYS REWRITE RIGHTMOST NON-TERMINAL
- (3) LEFTMOST: ALWAYS REWRITE LEFTMOST NON-TERMINAL

WHILE (1), (2) AND (3) MAKE NO DIFFERENCE ON WHAT CAN BE DERIVED AT FRONTIER (LEAVES) OF TREE, THEY DO TYPICALLY AFFECT INTERMEDIATE TREES AND ARE INTEGRALLY TIED TO PARSING STRATEGY (TOP DOWN VS BOTTOM UP)

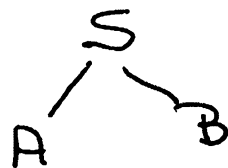
DERIVATIONS & TREES (CONTEXT FREE)

IF A GRAMMAR, G , IS A CFG THEN SYNTACTIC FORMS CAN BE ANY α , $\alpha \in (\Sigma \cup V)^*$

THUS, WHEN $S \xRightarrow{*} \alpha$ AND α CONTAINS

NON-TERMINALS, THERE ARE CHOICES
FIRST OF (1) WHICH NON-TERMINAL TO REWRITE
(2) WHICH RHS OF THAT NONTERM TO USE

A DERIVATION (PARSE) TREE HAS A FORM THAT COULD INCLUDE CASES LIKE

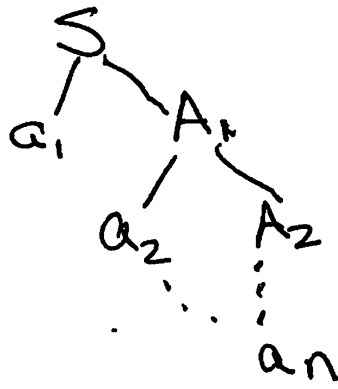


WE CAN NOW CHOOSE TO REWRITE A OR B
IMPORTANTLY, THOUGH, OUR CHOICE HAS NO EFFECT ON THE POSSIBLE OUTCOMES AS EACH SUBTREE EVOLVES INDEPENDENT OF ALL OTHER SUBTREES

DERIVATIONS & TREES (REGULAR)

IF A GRAMMAR, G , IS RIGHT LINEAR THEN EVERY NON-TERMINAL SYNTACTIC FORM IS OF THE CONSTRAINED FORM αB WHERE $\alpha \in \Sigma^*$ AND $B \in V$, THAT IS, IF $S \xRightarrow{*} \beta$ THEN $\beta = \alpha B$ OR $\beta = \alpha$, WHERE $\alpha \in \Sigma^*$ AND $B \in V$.

IF WE DRAW A TREE ASSOCIATED WITH A DERIVATION FOR A REGULAR GRAMMAR, THE TREE LOOKS LIKE



AT EACH STAGE, THE ONLY CHOICE IS AMONG THE RIGHT-HAND SIDES (RHS) OF B -RULES, WHEN WE ARE AT $S \xRightarrow{*} \alpha B$

AMBIGUITY

$$L = \{ a^i b^j c^k \mid k > i \text{ OR } k > j \}$$

IS INHERENTLY AMBIGUOUS (WHY?)

BUT ARITHMETIC EXPRESSIONS ARE NOT

$$E \rightarrow E + T \mid E - T \mid T$$

$$T \rightarrow T * F \mid T / F \mid F$$

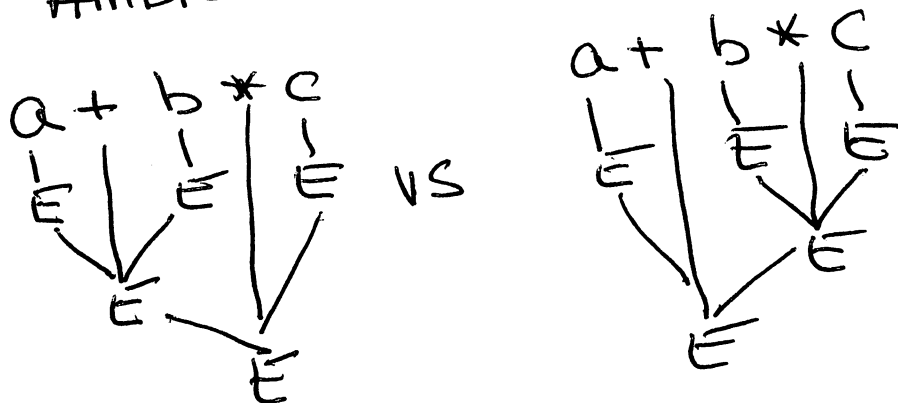
$$F \rightarrow (E) \mid id \mid number$$

A DFA CAN RECOGNIZE ID, NUMBER, OPS

POOR GRAMMAR

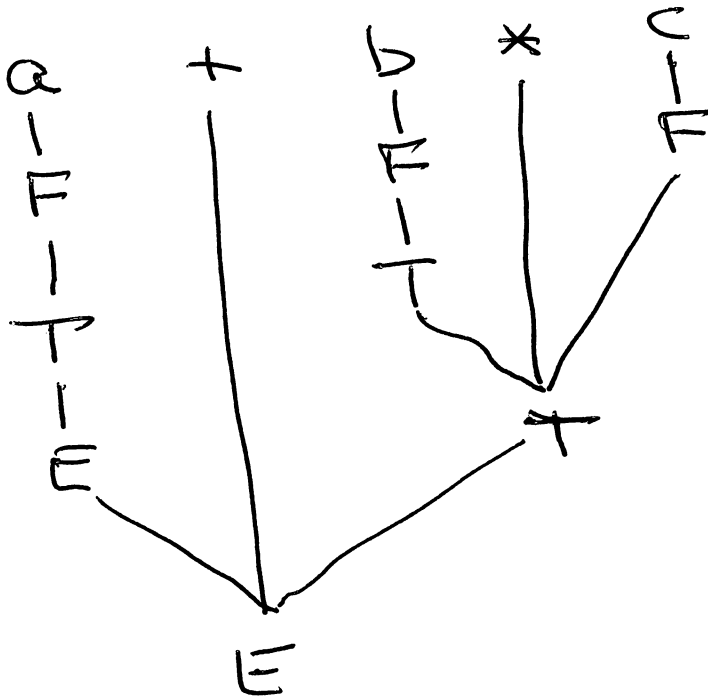
$$E \rightarrow ETE \mid E - T \mid E * E \mid E / E \mid (E) \mid id \mid number$$

IS AMBIGUOUS BUT LANGUAGE IS NOT



LIKE 2ND AND EARLIER GRAMMAR ENSURES THIS PARSE

UNAMBIGUOUS PARSE



IS ONLY POSSIBLE PARSE USING
GRAMMAR WE PREFERRED

LEFTMOST

$$\begin{aligned}
 E &\Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\
 &\Rightarrow a + T \Rightarrow a + T * F \Rightarrow a + F * F \\
 &\Rightarrow a + b * F \Rightarrow a + b * c
 \end{aligned}$$

RIGHTMOST

$$\begin{aligned}
 E &\Rightarrow E + T \Rightarrow E + T * F \\
 &\Rightarrow E + T * c \Rightarrow E + F * c \Rightarrow E + b * c \\
 &\Rightarrow T + b * c \Rightarrow F + b * c \Rightarrow a + b * c
 \end{aligned}$$

AMBIGUITY

A CFG IS AMBIGUOUS IF IT HAS

(i) TWO DISTINCT L.M. DERIVATIONS OF SOME STRING;

OR (ii) TWO DISTINCT R.M. DERIVATIONS OF SOME STRING

OR (iii) TWO DISTINCT PARSE TREES OF SOME STRING

A LANGUAGE IS INHERENTLY AMBIGUOUS IF ALL GRAMMARS FOR THAT LANGUAGE ARE AMBIGUOUS

TO BE SHOWN: THERE IS NO ALGORITHM THAT CAN ANSWER THE QUESTIONS, FOR ANY ARBITRARY CFG, G , "IS G AMBIGUOUS?"

AMBIGUOUS GRAMMAR

A GRAMMAR IS AMBIGUOUS IFF THERE IS SOME STRING w IN $L(A)$, SUCH THAT

a) w IS THE YIELD OF TWO OR MORE DISTINCT PARSE TREES

OR b) w IS THE YIELD OF TWO OR MORE DISTINCT LEFTMOST DERIVATIONS

OR c) w IS THE YIELD OF TWO OR MORE DISTINCT RIGHTMOST DERIVATIONS

NOTE : TOP DOWN PARSER

ASSOCIATED WITH LEFTMOST DERIV.
DOES NOT LIKE LEFT RECURSION

BOTTOM UP PARSER

ASSOCIATED WITH REVERSING
RIGHTMOST DERIV.

DOES NOT LIKE RIGHT RECURSION

AMBIGUITY EXAMPLES

$$E \rightarrow E + E \mid E * E \mid (E) \mid \epsilon$$

IS AMBIGUOUS BUT EXPRESSION LANGUAGE
IS NOT INHERENTLY AMBIGUOUS

$$L = \{ a^i b^j c^k \mid k > i \text{ OR } k > j \}$$

IS INHERENTLY AMBIGUOUS

$$S \rightarrow S_1 \mid S_2$$

$$S_1 \rightarrow a S_1 c \mid S_1 c \mid T_1 c$$

$$T_1 \rightarrow b T_1 \mid \lambda$$

$$S_2 \rightarrow a S_2 \mid T_2$$

$$T_2 \rightarrow b T_2 c \mid T_2 c \mid c$$

CASES WHERE $k > c$ AND $k > j$
CAN ARISE FROM S_1 AND S_2

NOTE: $\{ a^i b^j c^k \mid k > i \text{ AND } k > j \}$

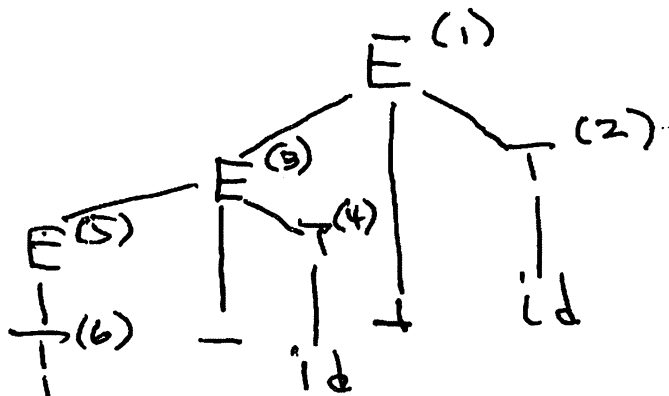
IS NOT A CFL

RIGHTMOST AND BOTTOM UP

RIGHTMOST TREES EVOLVE AS FOLLOWS:

$$E \rightarrow E+T \mid E-T \mid T$$

$$T \rightarrow id \mid (\epsilon)$$



$$id \quad E \Rightarrow E+T \Rightarrow E+id \Rightarrow E-T+id \Rightarrow E-id+id \Rightarrow T-id+id \Rightarrow id-id+id$$

BOTTOM UP is implemented as "SHIFT"/"REDUCE", THIS IS MADE DIFFICULT BY RIGHT RECURSION, SO LOVES OUR CURRENT GRAMMAR.

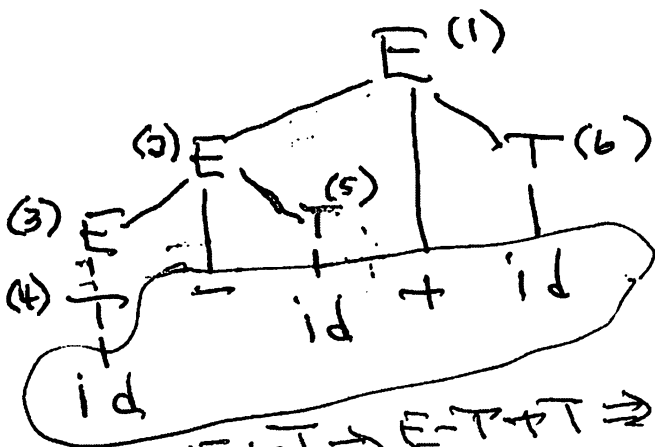
NOTE: THE LAST REWRITE IS OF LEFTMOST NON-TERMINAL AND SO FIRST REDUCE IS ON LEFT IN BOTTOM-UP

LEFTMOST AND TOP DOWN

LEFTMOST TREES EVOLVE AS FOLLOWS :

$$E \rightarrow E+T \mid E-T \mid T$$

$$T \rightarrow id \mid (E)$$



id-id+id

$E \Rightarrow E+T \Rightarrow E-T+T \Rightarrow T-T+T \Rightarrow id-T+T \Rightarrow id-id+T \Rightarrow id-id+id$

PROBLEM HERE IS THE PROCESS OF BUILDING TREE, BASED ON INPUT, IS "PREDICTIVE" (TOP DOWN) AND IS MADE VERY DIFFICULT WHEN GRAMMAR IS LEFT-RECURSIVE.

TOP DOWN IS OFTEN IMPLEMENTED BY "RECURSIVE DESCENT"

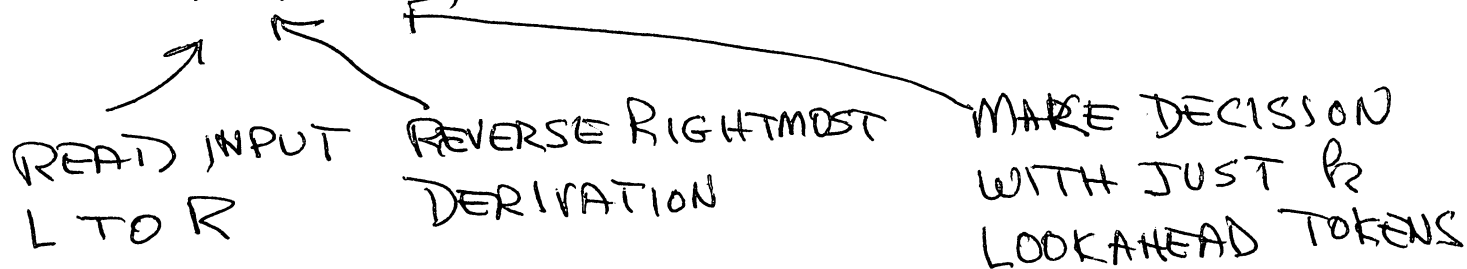
THEORY DRIVES PRACTICE

SINCE AMBIGUITY OF A GRAMMAR IS UNSOLVABLE NOT JUST UNSOLVED

CS COMPILER WRITERS COULD NOT USE ARBITRARY CFGs.

ENTER DONALD KNUTH

LR(k) PARSERS AND LR(k) GRAMMARS



LR(0): NOT GOOD ENOUGH

LR(1): GOOD ENOUGH

LR(1) LANGUAGES ARE ALL UNAMBIGUOUS CFLS

LR(1) = LR(2) = ... WHEN CONSIDERING LANGUAGES

LR(1) \neq LR(2) \neq ... FOR GRAMMARS

LR VS LL

LL(k) PARSERS AND LL(k) GRAMMARS

↑
USE LEFTMOST
DERIVATION

LL(1) \neq LL(2) \neq ...

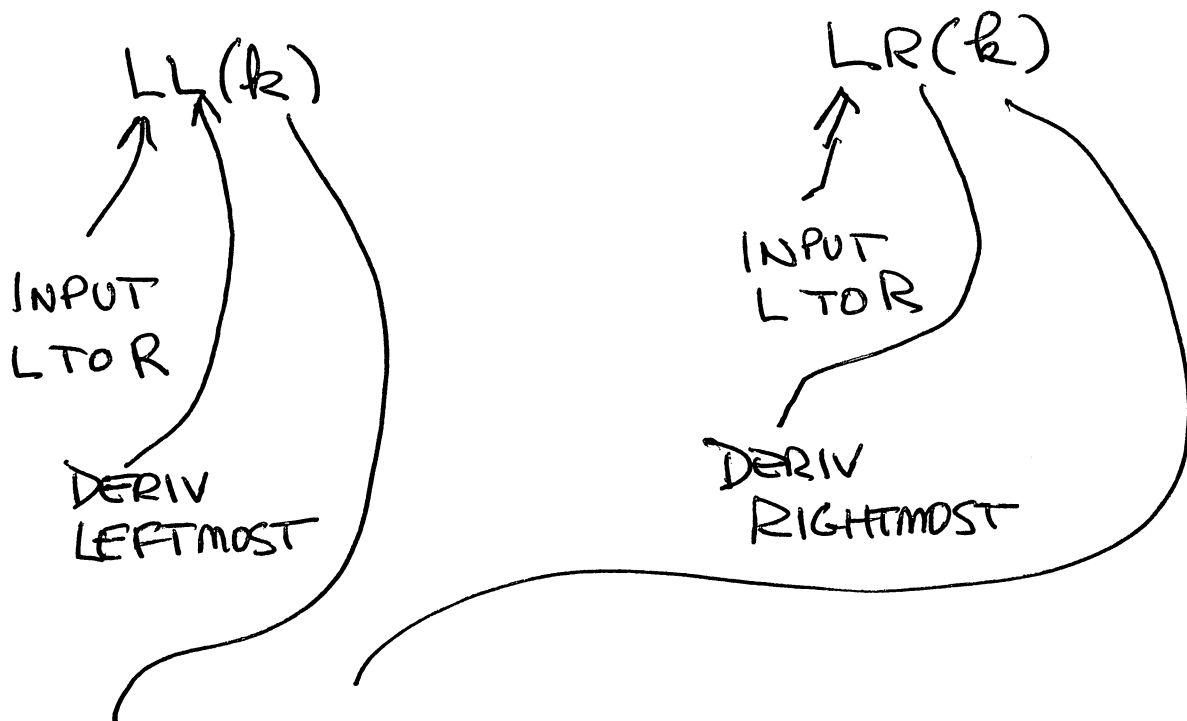
FOR BOTH GRAMMARS AND LANGUAGES

$\lim_{k \rightarrow \infty}$ LL(k) IS ALL UNAMBIGUOUS
LANGUAGES

BUT

LL(1) IS USUALLY GOOD ENOUGH
FOR COMPUTER LANGUAGES

PARSERS AND GRAMMAR CLASSES (CFG SUBCLASSES)



TOKENS NEEDED
TO LOOKAHEAD

LR(1) LANGUAGES = DET. CFLS

LR(1) GRAMMARS \subset DET CFGS

LR(k+1) LANGUAGES = LR(k) LANGS, $k > 0$

LR(k+1) GRAMMARS $\not\equiv$ LR(k) GRAMMARS, $k > 0$

LL(k+1) LANGS $\not\equiv$ LL(k) LANGS

LL(k+1) GRAMS $\not\equiv$ LL(k) GRAMS

$\lim_{k \rightarrow \infty}$ LL(k) LANGS. = DET CFLS

HOWEVER, LL(1) IS, IN PRACTISE,
ALL WE NEED FOR PROGR. LANGS.

TOP DOWN VS BOTTOM UP

TOP DOWN PARSERS

RECURSIVE DESCENT

AKA PREDICTIVE

BASED ON L.M. DERIVATION

CHALLENGE IS WHICH RHS

BOTTOM UP PARSERS

SHIFT/REDUCE

CHALLENGE IS

SHIFT OR REDUCE

IF REDUCE, WHICH ONE

CONFLICTS IN BOTH STRATEGIES

BUT LL(k) CAN MAKE RIGHT DECISION

WITH k LOOKAHEAD

AND LR(k) CAN MAKE RIGHT DECISION

WITH k LOOKAHEAD

BOTH USE STACK (PDA LATER)

CHOMSKY NORMAL FORM

ALL RULES OF FORM

$$A \rightarrow BC \quad B, C \in V$$

$$A \rightarrow a \quad a \in \Sigma$$

OR, IF λ IN LANGUAGE CAN HAVE

$$S \rightarrow \lambda$$

PROVIDE S IS NEVER USED IN
RHS OF ANY RULE

EVERY CFG CAN BE CONVERTED
TO AN EQUIVALENT CNF VERSION

WHY?

- PUMPING LEMMA

- POLYNOMIAL PARSER (CKY)

GREAT EXAMPLE

OF DYNAMIC PROGRAMMING

NOTES: 173-186