

1

LECTURE OUTLINE

- PARSING
 - TOP-DOWN
 - BOTTOM-UP
 - LL AND LR PARSING

Sejong - How to Create a Programming Language - Tony Mione / 2020

2

PARSING

- PARSER – TYPICAL FLOW
 - CALL SCANNER (LEXICAL ANALYZER) TO GET TOKENS
 - ASSEMBLE TOKENS INTO A SYNTAX TREE
 - PASS THE TREE TO LATER PHASES OF THE COMPILER (SYNTAX-DIRECTED TRANSLATION)
- MOST PARSERS HANDLE A CONTEXT-FREE GRAMMAR (CFG)

Sejong - How to Create a Programming Language - Tony Mione / 2020

3

PARSING

- TERMINOLOGY:
 - SYMBOLS
 - TERMINALS – SYMBOLS THAT HAVE NO ADDITIONAL BREAKDOWN
 - NON-TERMINALS – SYMBOLS THAT REPRESENT A SEQUENCE OF OTHER SYMBOLS
 - PRODUCTIONS – A POSSIBLE EXPANSION OF A NON-TERMINAL INTO A SEQUENCE OF TERMINALS AND NON-TERMINALS
 - DERIVATION – A SEQUENCE OF PRODUCTION APPLICATIONS DERIVING A PARSE OR SYNTAX TREE
 - (LEFT-MOST AND RIGHT-MOST – CANONICAL)
 - PARSE TREES – A TREE DATA STRUCTURE REPRESENTING THE SYNTAX OF A PROGRAM
 - SENTENTIAL FORMS – INTERMEDIATE EXPANSIONS (APPLICATIONS OF PRODUCTIONS) OF A PROGRAM'S CODE

Sejong - How to Create a Programming Language - Tony Mione / 2020

4

PARSING

- THERE ARE LARGE CLASSES OF GRAMMARS FOR WHICH PARSERS CAN RUN IN LINEAR TIME
- THE TWO MOST IMPORTANT: LL AND LR
 - LL: 'LEFT-TO-RIGHT, LEFTMOST DERIVATION'
 - LR: 'LEFT-TO-RIGHT, RIGHTMOST DERIVATION'
- LEFTMOST DERIVATION
 - WORK ON THE LEFT SIDE OF THE PARSE TREE
 - EXPAND LEFT-MOST NON-TERMINAL IN A SENTENTIAL FORM
- RIGHTMOST DERIVATION
 - WORK ON THE RIGHT SIDE OF THE TREE
 - EXPAND RIGHT-MOST NON-TERMINAL IN A SENTENTIAL FORM

Sejong - How to Create a Programming Language - Tony Mione / 2020

6

PARSING

- LL PARSERS
 - TOP-DOWN – APPLY PRODUCTIONS STARTING AT START SYMBOL
 - PREDICTIVE – MUST PREDICT WHICH PRODUCTION TO USE
- LR PARSERS
 - BOTTOM-UP – GATHER TOKENS AND ‘COMBINE’ THEM BY APPLYING PRODUCTION IN REVERSE
 - SHIFT-REDUCE
 - SHIFT TOKENS ONTO A STACK
 - REDUCE GROUPS OF TERMINALS AND NON-TERMINALS TO A SINGLE NON-TERMINAL ACCORDING TO PRODUCTIONS
- SEVERAL IMPORTANT SUB-CLASSES OF LR PARSERS
 - SLR
 - LALR
 - “FULL LR”

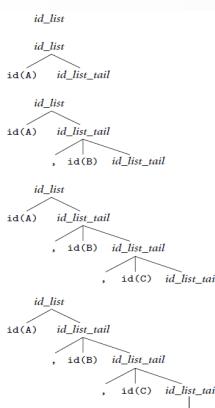
Sejong - How to Create a Programming Language - Tony Mione / 2020

7

TOP-DOWN PARSING (LL)

- CONSIDER A GRAMMAR FOR
 - COMMA SEPARATED LIST OF IDENTIFIERS
 - TERMINATED BY A SEMICOLON
- $\text{id_list} \rightarrow \text{id id_list_tail}$
- $\text{id_list_tail} \rightarrow , \text{id id_list_tail}$
- $\text{id_list_tail} \rightarrow ;$
- TOP-DOWN CONSTRUCTION OF A PARSE TREE FOR THE STRING: “A, B, C;” STARTS FROM THE ROOT AND APPLIES RULES

Sejong - How to Create a Programming Language - Tony Mione / 2020

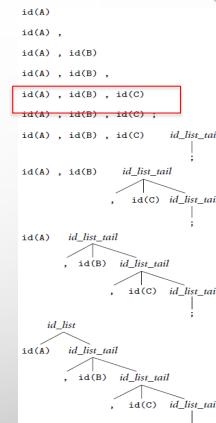


8

BOTTOM-UP PARSING (LR)

- $\text{id_list} \rightarrow \text{id id_list_tail}$
- $\text{id_list_tail} \rightarrow , \text{id id_list_tail}$
- $\text{id_list_tail} \rightarrow ;$
- BOTTOM-UP CONSTRUCTION OF A PARSE TREE FOR: “A, B, C;”
 - THE PARSER FINDS:
 - LEFT-MOST LEAF OF THE TREE IS AN ID
 - NEXT LEAF IS A COMMA
 - PARSER CONTINUES SHIFTING NEW LEAVES INTO A FOREST OF PARTIALLY COMPLETED PARSE TREE FRAGMENTS

Sejong - How to Create a Programming Language - Tony Mione / 2020



9

BOTTOM-UP PARSING (LR)

THE BOTTOM-UP CONSTRUCTION REALIZES SOME FRAGMENTS CONSTITUTE A COMPLETE RIGHT-HAND SIDE.

- I.E. OCCURS WHEN THE PARSER HAS SEEN THE SEMICOLON—THE RIGHT-HAND SIDE OF *ID_LIST_TAIL*
- THE PARSER REDUCES THE SEMICOLON TO AN *ID_LIST_TAIL*
- THEN REDUCES ", ID *IDLIST TAIL*" INTO ANOTHER *ID LIST TAIL*.
- AFTER ONE MORE TIME IT CAN REDUCE "ID *IDLIST TAIL*" INTO PARSE TREE ROOT: *ID_LIST*.

Sejong - How to Create a Programming Language - Tony Mione / 2020

10

PARSING

- NUMBER IN LL(1), LL(2), LALR(1), ...
- HOW MANY TOKENS OF LOOK-AHEAD ARE REQUIRED
- ALMOST ALL REAL COMPILERS USE **ONE TOKEN** OF LOOK-AHEAD
- EARLIER EXPRESSION GRAMMAR (WITH PRECEDENCE AND ASSOCIATIVITY) IS LR(1), BUT NOT LL(1)
- EVERY LL(1) GRAMMAR IS ALSO LR(1)
- CAVEAT: RIGHT RECURSION IN PRODUCTION REQUIRES VERY DEEP STACKS AND COMPLICATES SEMANTIC ANALYSIS

Sejong - How to Create a Programming Language - Tony Mione / 2020

11

EXAMPLE: AN LL(1) GRAMMAR

```

program   → stmt_list $$ (end of file)
stmt_list → stmt stmt_list
           | ε
stmt      → id := expr
           | read id
           | write expr
expr      → term term_tail
term      → factor fact_tail
fact      → mult_op factor fact_tail
           | ε
term_tail → add_op term term_tail
           | ε
add_op    → +
mult_op   → *
           | /

```

Sejong - How to Create a Programming Language - Tony Mione / 2020

12

LL PARSING

- ADVANTAGES OF THIS GRAMMAR
 - CAPTURES ASSOCIATIVITY AND PRECEDENCE
 - SIMPLICITY OF PARSING ALGORITHM
- DISADVANTAGES OF THE GRAMMAR
 - NOT AS 'PRETTY' AS OTHER EQUIVALENT GRAMMARS
 - OPERANDS OF AN OPERATOR NOT ON SAME RIGHT-HAND SIDE (RHS)

Sejong - How to Create a Programming Language - Tony Mione / 2020

13

LL PARSING

- EXAMPLE (THE AVERAGE PROGRAM):


```

READ A
READ B
SUM := A + B
WRITE SUM
WRITE SUM / 2    $$
```
- PARSING PROCESS
 - START AT TOP
 - PREDICT NEXT NEEDED PRODUCTION BASED ON:
 - CURRENT LEFT-MOST NON-TERMINAL
 - CURRENT INPUT TOKEN

Sejong - How to Create a Programming Language - Tony Mione / 2020

14

LL PARSING

- TABLE-DRIVEN PARSING:
 - LOOP: LOOK UP NEXT ACTION IN 2D TABLE BASED ON:
 - CURRENT LEFT-MOST NON-TERMINAL
 - CURRENT INPUT TOKEN
 - POSSIBLE ACTIONS:
 - MATCH A TERMINAL => PREDICT A PRODUCTION
 - SYNTAX ERROR

Sejong - How to Create a Programming Language - Tony Mione / 2020

15

LL PARSING

PREDICT

```

1. program → stmt_list $$ {id.read, write,$$}
2. stmt_list → stmt stmt_list {id.read, write}
3. stmt_list → ε {$$}
4. stmt → id := expr {id}
5. stmt → read id {read}
6. stmt → write expr {write}
7. expr → term term_tail {(), id, number}
8. term_tail → add_op term term_tail {+, -}
9. term_tail → ε {(), id, read, write, $$}
10. term → factor factor_tail {(), id, number}
11. factor_tail → mult_op factor factor_tail {*, /}
12. factor_tail → ε {+, -, }, id, read, write, $$}
13. factor → ( expr ) {}
14. factor → id {id}
15. factor → number {number}
16. add_op → + {+}
17. add_op → - {-}
18. mult_op → * {*}
19. mult_op → / {/}
```

Sejong - How to Create a Programming Language - Tony Mione / 2020

16

LL PARSING

LL(1) PARSE TABLE FOR CALCULATOR LANGUAGE

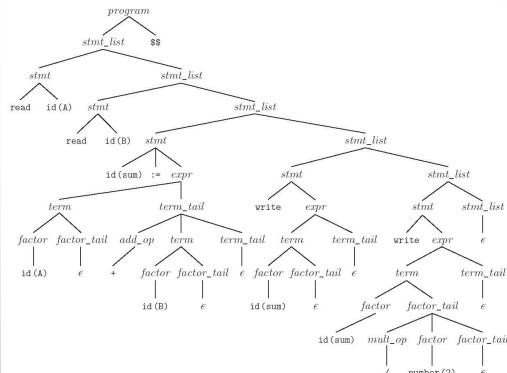
READ A
READ B
SUM := A + B
WRITE SUM
WRITE SUM / 2 \$\$

Top-of-stack nonterminal	id	number	read	Current input token	write	:=	()	+	*	/	\$\$
program	1	-	1	1	-	-	-	-	-	-	-	1
stmt_list	2	-	2	2	-	-	-	-	-	-	-	3
stmt	4	-	5	6	-	-	-	-	-	-	-	-
expr	7	7	-	-	-	7	-	-	-	-	-	-
term_tail	9	-	9	9	-	9	8	8	-	-	-	9
term	10	10	-	-	-	10	-	-	-	-	-	-
factor_tail	12	-	12	12	-	-	12	12	12	11	11	12
factor	14	15	-	-	-	13	-	-	-	-	-	-
add_op	-	-	-	-	-	-	-	-	16	17	-	-
mult_op	-	-	-	-	-	-	-	-	-	18	19	-

Sejong - How to Create a Programming Language - Tony Mione / 2020

17

PARSE TREE FOR THE AVERAGE PROGRAM



Sejong - How to Create a Programming Language - Tony Mione / 2020

18

LL PARSING

- PROBLEMS TRYING TO MAKE A GRAMMAR LL(1)

- LEFT-RECURSION

- EXAMPLE

$\text{id_list} \rightarrow \text{id}$
| $\text{id_list}, \text{id}$

- CAN REMOVE LEFT RECURRENCE MECHANICALLY

$\text{id_list} \rightarrow \text{id id_list_tail}$
 $\text{id_list_tail} \rightarrow , \text{id id_list_tail}$
| ϵ

Sejong - How to Create a Programming Language - Tony Mione / 2020

19

LL PARSING

- MORE PROBLEMS
 - COMMON PREFIXES

$\text{stmt} \rightarrow \text{id} := \text{expr}$
| $\text{id} (\text{arg_list})$

- CAN ELIMINATE COMMON FACTOR MECHANICALLY
 - LEFT-FACTORING

$\text{stmt} \rightarrow \text{id id_stmt_tail}$
 $\text{id_stmt_tail} \rightarrow := \text{expr}$
| (arg_list)

Sejong - How to Create a Programming Language - Tony Mione / 2020

20

LL PARSING

- ELIMINATING LEFT RECURRENCE AND COMMON PREFIXES DOES NOT MAKE A GRAMMAR LL

- THERE ARE INFINITELY MANY NON-LL LANGUAGE

- MECHANICAL TRANSFORMATIONS SHOWN WORK ON THEM JUST FINE
- GRAMMAR IS STILL NOT LL(1)

Sejong - How to Create a Programming Language - Tony Mione / 2020

21

LL PARSING

- MORE PROBLEMS MAKING GRAMMAR LL(1)
 - THE DANGLING-ELSE PROBLEM:
 - EXAMPLE (PASCAL): IF C1 THEN IF C2 THEN S1 ELSE S2
 - GRAMMAR:


```
stmt → if cond then_clause else_clause
                | other_stuff
        then_clause → then stmt
        else_clause → else stmt | ε
```
 - THE ELSE CAN BE PAIRED WITH EITHER IF-THEN IN THE EXAMPLE!
 - FIX (DISAMBIGUATION RULE): PAIR AN ELSE WITH THE CLOSEST "ELSE-LESS" NESTED IF-THEN

Sejong - How to Create a Programming Language - Tony Mione / 2020

22

LL PARSING

- A LESS NATURAL GRAMMAR FRAGMENT:


```
stmt → balanced_stmt | unbalanced_stmt
            balanced_stmt → if cond then
                            balanced_stmt
                            else balanced_stmt
                            | other_stuff
            unbalanced_stmt → if cond then stmt
                            | if cond then
                                balanced_stmt
                                else unbalanced_stmt
```
- BALANCED_STMT HAS EQUAL NUMBER OF THENS AND ELSESES
- UNBALANCED_STMT HAS MORE THENS

Sejong - How to Create a Programming Language - Tony Mione / 2020

23

LL PARSING

- USUAL APPROACH [REGARDLESS OF TOP-DOWN OR BOTTOM-UP]
 - USE AMBIGUOUS GRAMMAR
 - ADD DISAMBIGUATING RULES
 - I.E., ELSE GOES WITH CLOSEST UNMATCHED THEN
 - GENERALLY: THE FIRST OF TWO POSSIBLE PRODUCTION IS PREDICTED [OR SELECTED FOR REDUCTION]

Sejong - How to Create a Programming Language - Tony Mione / 2020

24

LL PARSING

- NEWER LANGUAGES (SINCE PASCAL) INCLUDE EXPLICIT END MARKERS
- MODULA-2:


```
if A = B then
        if C = D then E := F end
      else
        G := H
      end
```

ADA USES END IF; SOME LANGUAGES USE FI

Sejong - How to Create a Programming Language - Tony Mione / 2020

25

LL PARSING

- PROBLEM WITH END MARKERS: THEY ACCUMULATE IN PASCAL:

```
if A = B then ...
else if A = C then ...
else if A = D then ...
else if A = E then ...
else ...;
```

- WITH END MARKERS, THIS BECOMES:

```
if A = B then ...
else if A = C then ...
else if A = D then ...
else if A = E then ...
else ...;
end; end; end; end; end; ...
```

Sejong - How to Create a Programming Language - Tony Mione / 2020

26

LR PARSING

- LR PARSERS ALMOST ALWAYS TABLE-DRIVEN
 - LIKE A TABLE-DRIVEN LL PARSER, LR PARSER:
 - USES A LOOP
 - REPEATEDLY INSPECTS A TWO-DIMENSIONAL TABLE TO FIND NEXT ACTION
 - UNLIKE THE LL PARSER, LR DRIVER:
 - HAS NON-TRIVIAL STATE (LIKE A DFA)
 - HAS TABLE INDEXED BY CURRENT INPUT TOKEN & CURRENT STATE
 - STACK CONTAINS A RECORD OF WHAT HAS BEEN SEEN SO FAR (NOT WHAT IS EXPECTED)

Sejong - How to Create a Programming Language - Tony Mione / 2020

27

LR PARSING

- A SCANNER IS A DFA
- AN LL OR LR PARSER IS A PDA [PUSH DOWN AUTOMATA]
 - EARLY'S & COCKE-YOUNGER-KASAMI ALGORITHMS DO NOT USE PDAs
- A PDA CAN BE SPECIFIED WITH A STATE DIAGRAM AND A STACK
 - STATE DIAGRAM LOOKS JUST LIKE A DFA STATE DIAGRAM
 - EXCEPT ARCS ARE LABELED WITH <INPUT SYMBOL, TOP-OF-STACK SYMBOL> PAIRS
 - IN ADDITION TO MOVING TO A NEW STATE, PDA HAS THE OPTION OF PUSHING OR POPPING A FINITE NUMBER OF SYMBOLS ONTO/OFF THE STACK

Sejong - How to Create a Programming Language - Tony Mione / 2020

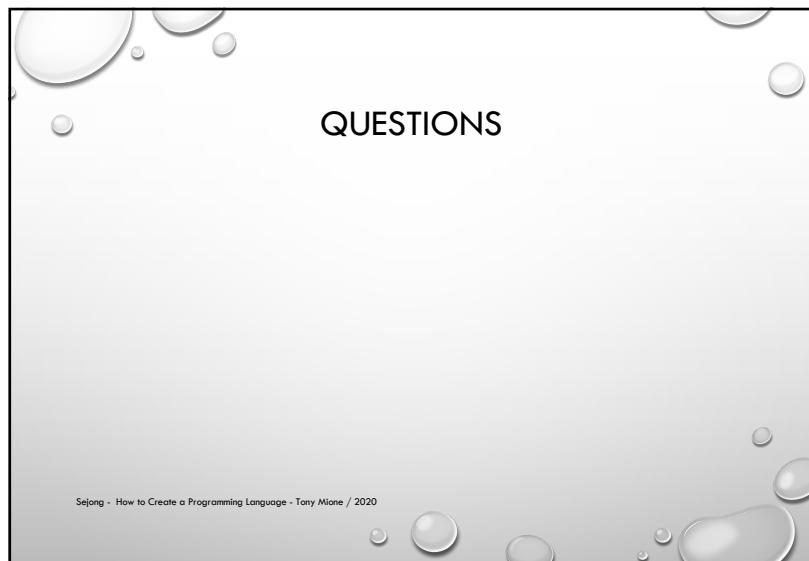
28

ACTIONS

- WE CAN RUN ACTIONS WHEN A RULE TRIGGERS:
 - OFTEN USED TO CONSTRUCT AN AST (ABSTRACT SYNTAX TREE) FOR A COMPILER
 - FOR SIMPLE LANGUAGES, CAN INTERPRET CODE DIRECTLY [AND POSSIBLY GENERATE A TRANSLATION]
 - CAN USE ACTIONS TO FIX THE TOP-DOWN PARSING PROBLEMS

Sejong - How to Create a Programming Language - Tony Mione / 2020

29



30