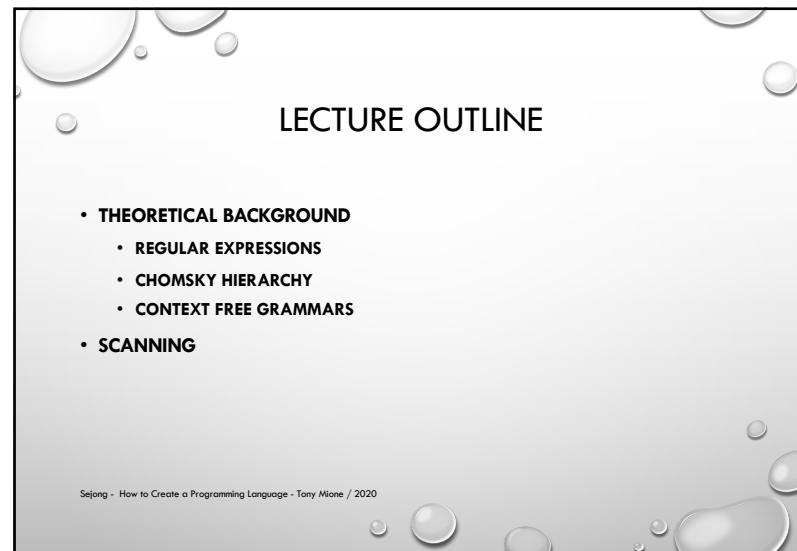
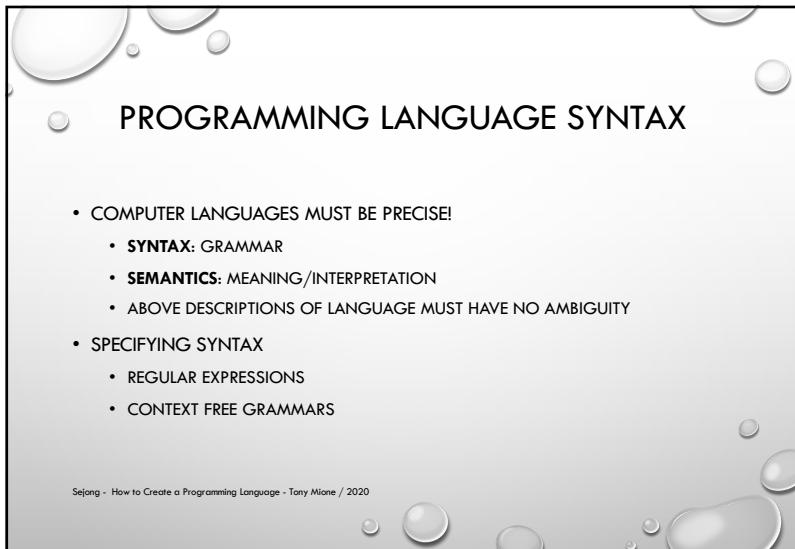


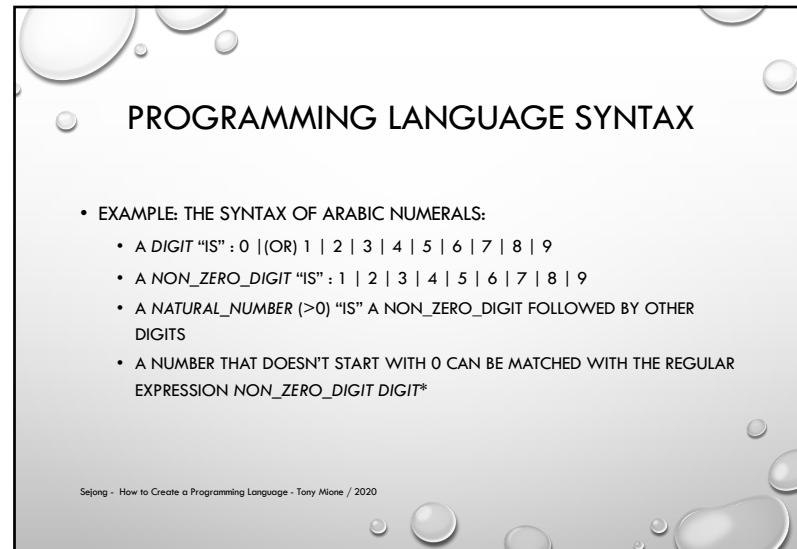
1



2



3



4

## REGULAR EXPRESSIONS

- A REGULAR EXPRESSION IS:  $\epsilon$
- THE EMPTY STRING, DENOTED BY
- A SINGLE CHARACTER
- TWO REGULAR EXPRESSIONS CONCATENATED
  - A IS A REGULAR EXPRESSION, B IS A REGULAR EXPRESSION, THEN AB IS A REGULAR EXPRESSION
- ALTERNATION OF TWO REGULAR EXPRESSIONS [ SEPARATED BY 'OR' | ]
  - GIVEN A AND B ARE REGULAR EXPRESSIONS,  $(A|B)$  IS A REGULAR EXPRESSION
- A REGULAR EXPRESSION FOLLOWED BY KLEENE \* (CLOSURE)
  - GIVEN A IS A REGULAR EXPRESSION,  $A^*$  IS A REGULAR EXPRESSION

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

5

## REGULAR EXPRESSIONS (EXAMPLE)

- DECIMAL NUMERIC LITERAL
  - digit is 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  - integer is digit digit\*
  - real is integer exponent | decimal (exponent |  $\epsilon$ )
  - decimal is digit\*(.digit|digit.)digit\*
  - exponent is  $(e|E)(+|-|\epsilon)$  integer
  - number is (integer | real)

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

6

## REGULAR EXPRESSIONS

- WORK WELL FOR DEFINING TOKENS
- UNABLE TO SPECIFY NESTED CONSTRUCTIONS
  - EXAMPLE: CONTEXT FREE GRAMMAR BNF
    - $EXPR \rightarrow ID \mid NUMBER \mid -EXPR \mid (\text{EXPR}) \mid EXPR \text{ OP } EXPR$
    - $OP \rightarrow + \mid - \mid * \mid /$
    - (SAME NUMBER OF OPEN AND CLOSE PARENSES CANNOT BE DESCRIBED)

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

7

## CHOMSKY HIERARCHY

- Context Free Languages more powerful than RE
- REs are faster to recognize
- Chomsky Hierarchy
  - **Type 3** – Regular Language – Finite Automata
  - **Type 2** – Context Free Languages – Pushdown Automata
  - **Type 1** - Context Sensitive Languages – Turing Machine w/Tape \* Input
  - **Type 0** – Unrestricted Language – Turing Machine
    - Type 0/1 too slow for practical use
    - Type 2 – It depends
    - Type 3 – Fast!

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

8

## CONTEXT FREE GRAMMARS (CFG)

- BACKUS-NAUR FORM (BNF) NOTATION:
 
$$\text{EXPR} \rightarrow \text{ID} \mid \text{NUMBER} \mid \text{-EXPR} \mid (\text{EXPR}) \mid \text{EXPR OP EXPR}$$

$$\text{OP} \rightarrow + \mid - \mid * \mid /$$
- EACH 'RULE' IS CALLED A *PRODUCTION*
- SYMBOLS ON LEFT ARE NON-TERMINALS
- A CFG INCLUDES:
  - A SET OF TERMINALS T
  - A SET OF NON-TERMINALS N
  - A START SYMBOL S (NON-TERMINAL)
  - A SET OF PRODUCTIONS

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

9

## CONTEXT FREE GRAMMARS

- GIVEN EXPRESSION GRAMMAR
  - $\text{EXPR} := \text{EXPR OP EXPR}$
  - $\text{EXPR} := \epsilon$
  - $\text{OP} := + \mid * \mid - \mid /$
- PREVIOUS GRAMMAR AMBIGUOUS
  - CAN DERIVE 2 DIFFERENT PARSE TREES!
    - TREE 1:  $3 + (4 * 5)$
    - TREE 2:  $(3 + 4) * 5$
- BETTER VERSION
  - ...INCLUDE PRECEDENCE AND ASSOCIATIVITY
$$\text{EXPR} \rightarrow \text{TERM} \mid \text{EXPR ADD\_OP TERM}$$

$$\text{TERM} \rightarrow \text{FACTOR} \mid \text{TERM MULT\_OP FACTOR}$$

$$\text{FACTOR} \rightarrow \text{ID} \mid \text{NUMBER} \mid \text{-FACTOR} \mid (\text{EXPR})$$

$$\text{ADD\_OP} \rightarrow + \mid -$$

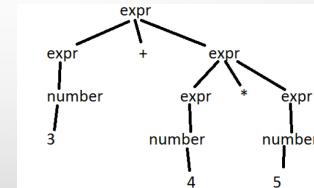
$$\text{MULT\_OP} \rightarrow * \mid /$$

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

11

## DERIVATIONS AND PARSE TREES

- A CFG SHOWS HOW A STRING OF TERMINALS IS GENERATED BY A LANGUAGE
  - USING PRODUCTIONS OF A LANGUAGE
  - $\rightarrow$  SYNTACTICALLY VALID IN THE 'LANGUAGE'
- EX: PARSE TREE FOR  $3 + 4 * 5$

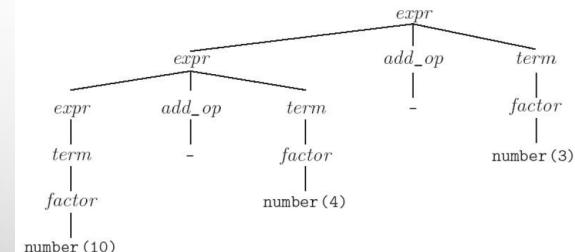


Sejong - How to Create a Programming Language - Tony Mione / 2020

10

## CONTEXT FREE GRAMMARS

- PARSE TREE FOR EXPRESSION GRAMMAR (WITH LEFT ASSOCIATIVITY) FOR  $10 - 4 - 3$



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

12

## SCANNING

- THE (LEXICAL) SCANNER AND PARSER (TOGETHER) RESPONSIBLE FOR DISCOVERING SYNTACTIC STRUCTURE OF PROGRAM
  - THE SCANNER IS RESPONSIBLE FOR
    - TOKENIZING SOURCE
    - REMOVING COMMENTS, (OFTEN) DEALING WITH PRAGMAS (I.E., SIGNIFICANT COMMENTS))
    - SAVING TEXT OF IDENTIFIERS, NUMBERS, STRINGS
    - SAVING SOURCE LOCATIONS (FILE, LINE, COLUMN) FOR ERROR MESSAGES

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

13

## SCANNING

Consider an ad-hoc (hand-written) scanner for a Calculator:

```

assign → :=  

plus → +  

minus → -  

times → *  

div → /  

lparen → (  

rparen → )  

id → letter (letter | digit)*  

number → digit digit*  

| digit * ( . digit | digit . ) digit *  

comment → /* ( non-* | * non-/ )* */  

| // ( non-newline )* newline
  
```

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

14

## SCANNING

- Read characters one at a time with look-ahead
    - skip initial white space (spaces, tabs, and newlines)
- ```

if cur_char ∈ {‘(’, ‘)’, ‘+’, ‘-’, ‘*’}  

  return the corresponding single-character token  

if cur_char = ‘:’  

  read the next character  

  if it is '=' then return assign else announce an error  

if cur_char = ‘/’  

  peek at the next character  

  if it is '*' or '/'  

    read additional characters until “*/” or  

    newline is seen, respectively  

  jump back to top of code  

else return div
  
```

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

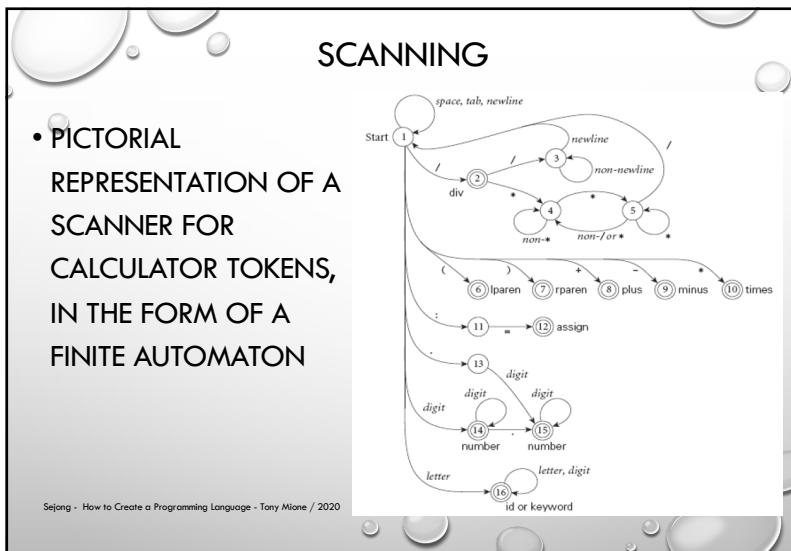
15

## SCANNING

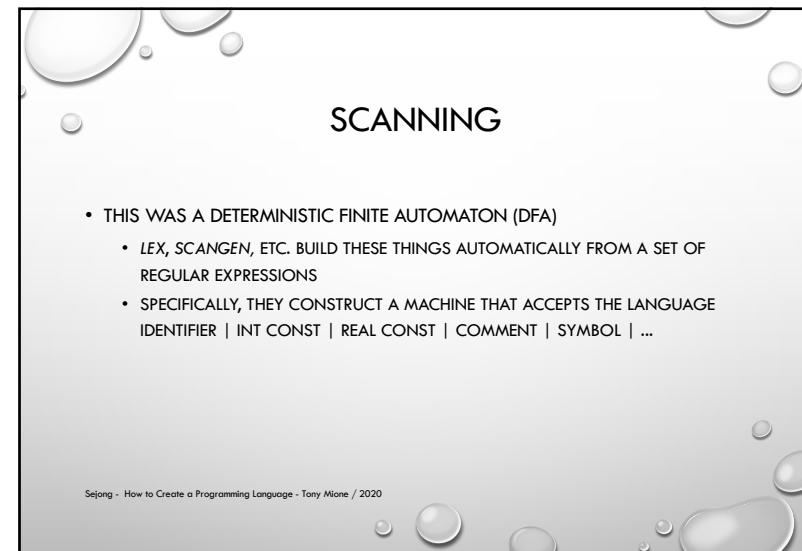
- if **cur\_char** = .
  - read the next character
  - if it is a digit
    - read any additional digits
    - return number
  - else announce an error
- if **cur\_char** is a digit
  - read any additional digits and at most one decimal point
  - return number
- if **cur\_char** is a letter
  - read any additional letters and digits
  - check to see whether the resulting string is **read** or **write**
  - if so then return the corresponding token
  - else return id
- else announce an error

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

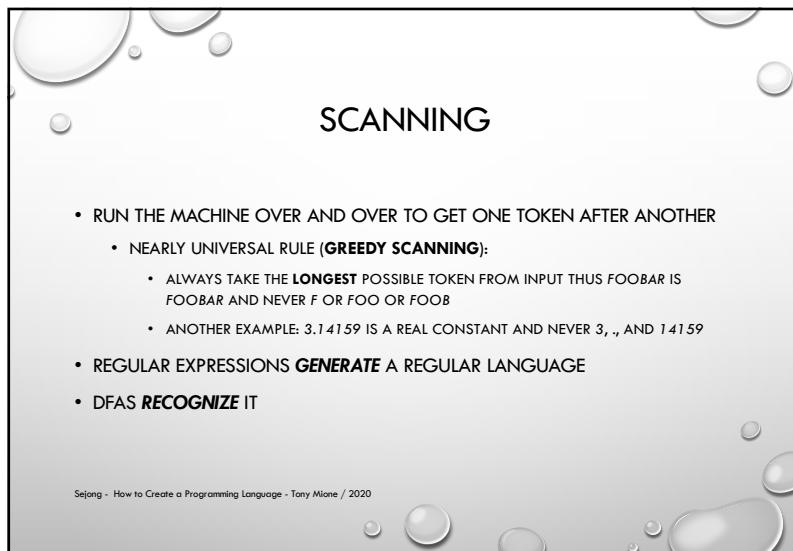
16



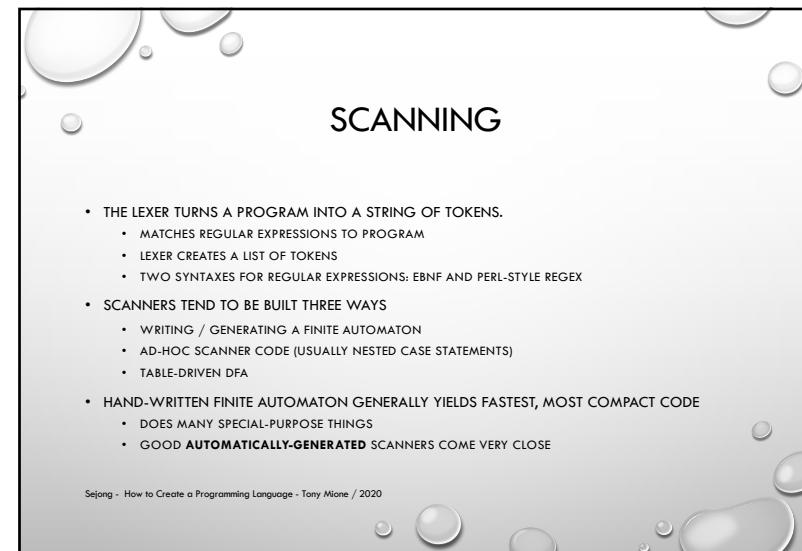
17



18



19



20

## SCANNING

- WRITING A PURE DFA AS A SET OF NESTED CASE STATEMENTS IS A SURPRISINGLY USEFUL TECHNIQUE
  - USE PERL, AWK, SED
  - TABLE-DRIVEN DFA IS WHAT LEX AND SCANGEN PRODUCE
    - LEX(FLEX) IN THE FORM OF C CODE
    - SCANGEN IN THE FORM OF NUMERIC TABLES AND A SEPARATE DRIVER

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

21

## SCANNING

- RULE ABOUT LONGEST-POSSIBLE TOKENS
  - YOU RETURN **ONLY** WHEN THE NEXT CHARACTER CAN'T BE USED TO CONTINUE THE CURRENT TOKEN
  - NEXT CHARACTER GENERALLY NEEDS TO BE SAVED FOR NEXT TOKEN
- MAY NEED TO PEEK AT MORE THAN ONE CHARACTER TO KNOW WHETHER TO PROCEED
  - IN PASCAL, FOR EXAMPLE, WHEN YOU HAVE A 3 AND YOU SEE A DOT
    - DO YOU PROCEED (IN HOPES OF GETTING 3.14)? OR
    - DO YOU STOP (IN FEAR OF GETTING 3..5)? (DECLARATION OF ARRAYS IN PASCAL, E.G., "ARRAY [1..6] OF INTEGER")

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

22

## PERL-STYLE REGEXP

- Learning by examples:
  - abcd - concatenation
  - a(b|c)d - grouping
  - a(b|c)\*d - can apply # of repeats to char or group.
    - ? = 0-1
    - \* = 0-infinity
    - + = 1-infinity
  - [bc] -character class – b or c
  - [a-zA-Z0-9\_] –ranges- any alpha numeric character
  - . -matches any character – one character only
  - \alpha -alpha
  - \d -numeric
  - \w -word (alpha, num, \_)
  - \s -whitespace

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

23

## PERL-STYLE REGEXP

Learning by examples:

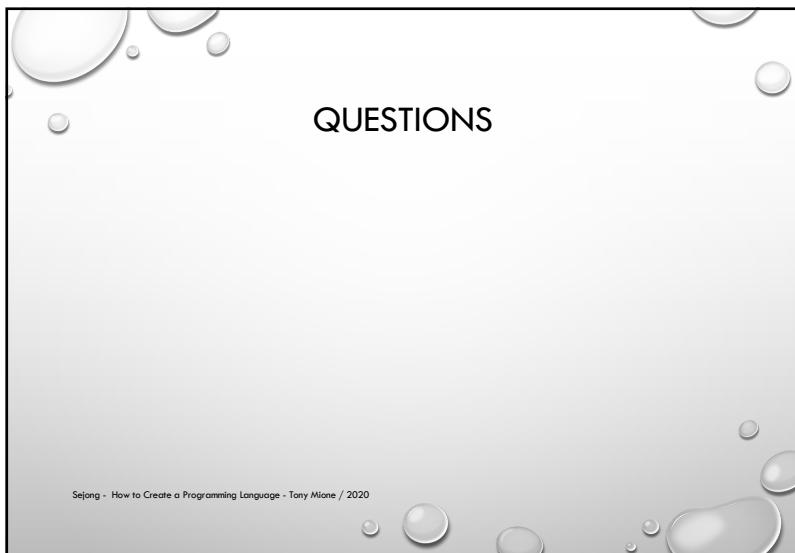
**How do we write a regexp that matches floats?**

$\text{digit}^*(.\text{digit}|\text{digit.})\text{digit}^*$

$\backslash d^* \backslash . \backslash d | \backslash d \backslash . \backslash d \backslash d^*$

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

24



25