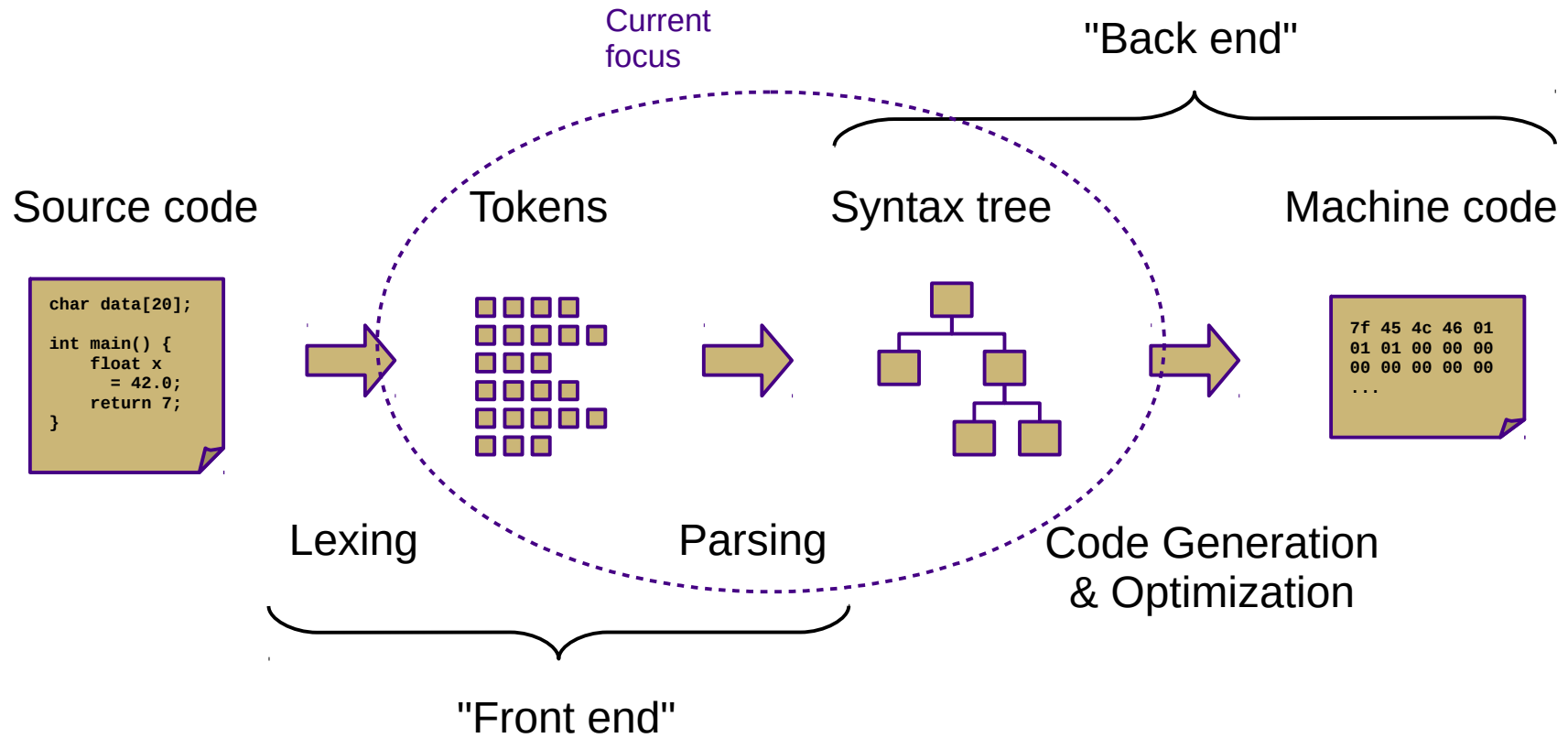


CS 432
Fall 2017

Mike Lam, Professor

Top-Down (LL) Parsing

Compilation



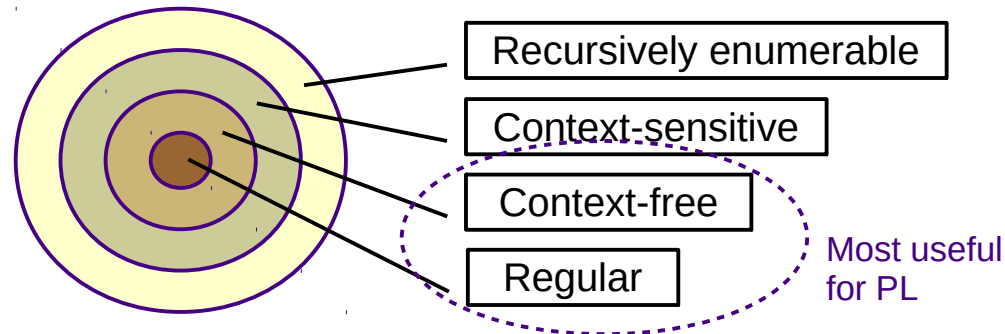
Segue

- Recognize **regular languages** with **finite automata**
 - Described by regular expressions
 - Rule-based transitions, no memory required
- Recognize **context-free languages** with **pushdown automata**
 - Described by context-free grammars
 - Rule-based transitions, MEMORY REQUIRED
 - Add a stack!

Segue

KEY OBSERVATION: Allowing the translator to use memory to track **parse state** information enables a **wider range** of automated machine translation.

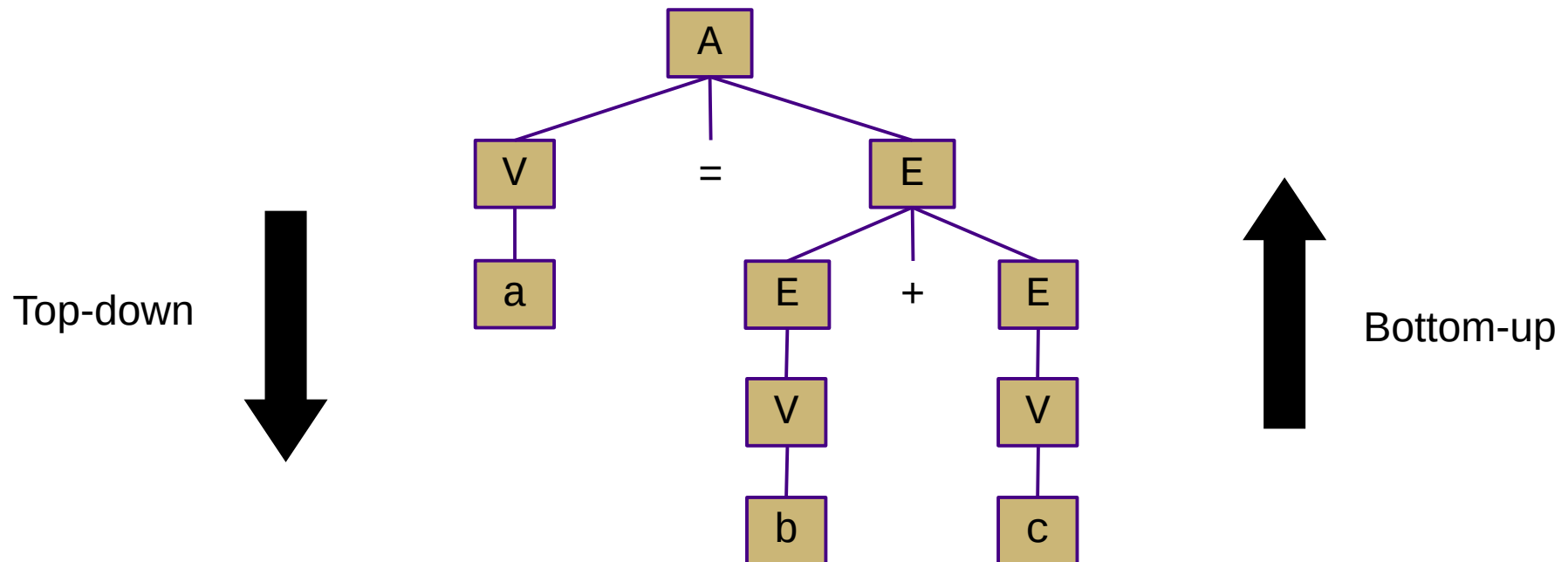
Chomsky Hierarchy of Languages



Grammar	Languages	Automaton	Production rules (constraints)
Type-0	Recursively enumerable	Turing machine	$\alpha \rightarrow \beta$ (no restrictions)
Type-1	Context-sensitive	Linear-bounded non-deterministic Turing machine	$\alpha A \beta \rightarrow \alpha \gamma \beta$
Type-2	Context-free	Non-deterministic pushdown automaton	$A \rightarrow \gamma$
Type-3	Regular	Finite state automaton	$A \rightarrow a$ and $A \rightarrow aB$

Parsing Approaches

- **Top-down**: begin with start symbol (root of parse tree), and gradually expand non-terminals
 - Stack contains leaves that still need to be expanded
- **Bottom-up**: begin with terminals (leaves of parse tree), and gradually connect using non-terminals
 - Stack contains roots of subtrees that still need to be connected



Top-Down Parsing

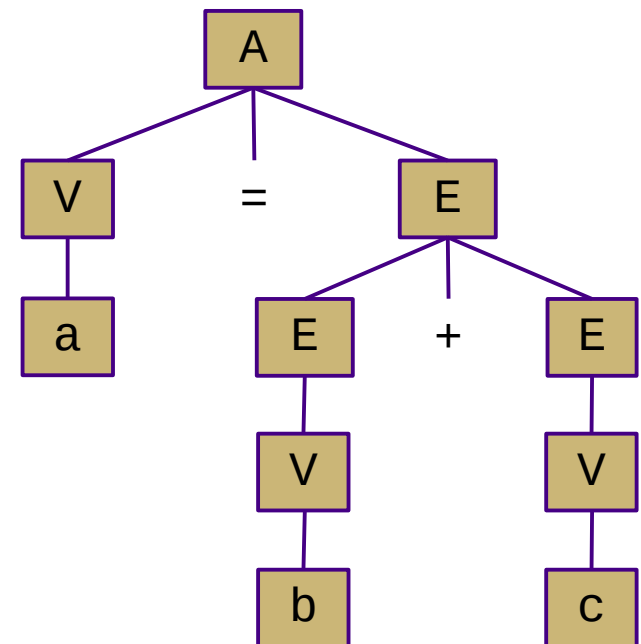
```
root = createNode(S)
focus = root
push(null)
token = nextToken()

loop:
  if (focus is non-terminal):
    B = chooseRuleAndExpand(focus)
    for each b in B.reverse():
      focus.addChild(createNode(b))
      push(b)
      focus = pop()

  else if (token == focus):
    token = nextToken()
    focus = pop()

  else if (token == EOF and focus == null):
    return root

  else:
    exit(ERROR)
```

$$\begin{array}{l} A \rightarrow V = E \\ V \rightarrow a \mid b \mid c \\ E \rightarrow E + E \\ \quad \mid V \end{array}$$


Recursive Descent Parsing

- Idea: use the system stack rather than an explicit stack
 - One function for each non-terminal
 - Encode productions with function calls and token checks
 - Use recursion to track current “state” of the parse
 - Easiest kind of parser to write manually

A → 'if' C 'then' S
| 'goto' L



```
class A {  
    public enum Type  
        { IFTHEN, GOTO }  
    public Type type  
    public C cond  
    public S stmt  
    public L lbl  
}
```

```
parseA(tokens):  
    node = new A()  
    next = tokens.next()  
    if next == "if":  
        node.type = IFTHEN  
        node.cond = parseC()  
        matchToken("then")  
        node.stmt = parseS()  
    else if next == "goto":  
        node.type = GOTO  
        node.lbl = parseL()  
    else  
        error ("expected 'if' or 'goto'")  
    return node
```

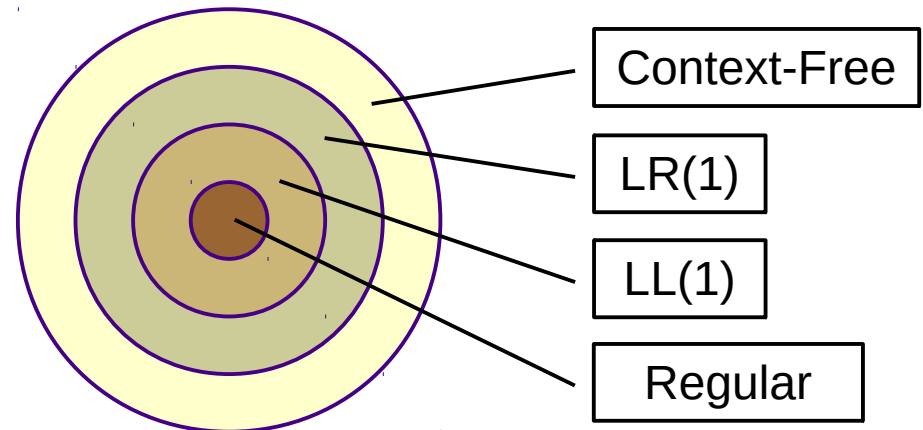
Top-Down Parsing

- Main issue: choosing which rule to use
 - With full lookahead, it would be relatively easy
 - This would be very inefficient
 - Can we do it with a single lookahead?
 - That would be much faster

LL(1) Parsing

- **LL(1)** grammars and parsers
 - **Left-to-right** scan of the input string
 - **Leftmost** derivation
 - **1 symbol** of lookahead
 - Highly restricted form of context-free grammar
 - No left recursion
 - No backtracking

**Context-Free
Hierarchy**



LL(1) Grammars

- We can convert many practical grammars to be LL(1)
 - Must remove left recursion
 - Must remove common prefixes (i.e., left factoring)

$$\begin{array}{l} A \rightarrow A \alpha \\ | \beta \end{array}$$

Grammar with left recursion

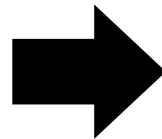
$$\begin{array}{l} A \rightarrow \alpha \beta_1 \\ | \alpha \beta_2 \end{array}$$

Grammar with common prefixes

Eliminating Left Recursion

- **Left recursion:** $A \rightarrow A \alpha \mid \beta$
 - Often a result of left associativity (e.g., expression grammar)
 - Leads to infinite looping/recursion in an LL(1) parser
 - To fix, unroll the recursion into a new non-terminal
 - Practical note: A and A' can be a single method in your code
 - Parse one β , then continue parsing α 's until there are no more
 - Keep adding the previous parse tree as a left subnode of the new parse tree

$$\begin{array}{l} A \rightarrow A \alpha \\ | \quad \beta \end{array}$$

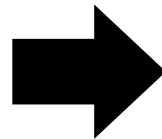


$$\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \\ | \quad \varepsilon \end{array}$$

Left Factoring

- **Common prefix:** $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$
 - Leads to ambiguous rule choice in LL(1) parser
 - One lookahead (α) is not enough to pick a rule; backtracking is required
 - To fix, **left factor** the choices into a new non-terminal
 - Practical note: A and A' can be a single method in your code
 - Parse and save data about α in temporary variables until you have enough information to choose

$$\begin{array}{l} A \rightarrow \alpha \beta_1 \\ | \alpha \beta_2 \end{array}$$



$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \\ | \beta_2 \end{array}$$

LL(1) Parsing

- LL(1) parsers can also be auto-generated
 - Similar to auto-generated lexers
 - Tables created by a *parser generator* using **FIRST** and **FOLLOW** helper sets
 - These sets are also useful when building hand-written recursive descent parsers
 - And (as we'll see next week), when building SLR parsers

LL(1) Parsing

- **FIRST(α)**
 - Set of terminals (and ϵ) that can appear at the start of a sentence derived from α (can be a terminal or non-terminal)
- **FOLLOW(A)** set
 - Set of terminals (and $\$$) that can occur immediately after non-terminal A in a sentential form
- **FIRST⁺(A \rightarrow β)**
 - If ϵ is not in FIRST(β)
 - FIRST⁺(A) = FIRST(β)
 - Otherwise
 - FIRST⁺(A) = FIRST(β) \cup FOLLOW(A)

Useful for choosing which rule to apply when expanding a non-terminal

Calculating FIRST(α)

- Rule 1: α is a terminal \mathbf{a}
 - $\text{FIRST}(\mathbf{a}) = \{ \mathbf{a} \}$
- Rule 2: α is a non-terminal X
 - Examine all productions $X \rightarrow Y_1 Y_2 \dots Y_k$
 - add $\text{FIRST}(Y_1)$ if not $Y_1 \rightarrow^* \epsilon$
 - add $\text{FIRST}(Y_i)$ if $Y_1 \dots Y_{i-1} \rightarrow^* \epsilon$, where $j = i-1$ (i.e., skip disappearing symbols)
 - $\text{FIRST}(X)$ is union of all of the above
- Rule 3: α is a non-terminal X and $X \rightarrow \epsilon$
 - $\text{FIRST}(X)$ includes ϵ

Calculating FOLLOW(B)

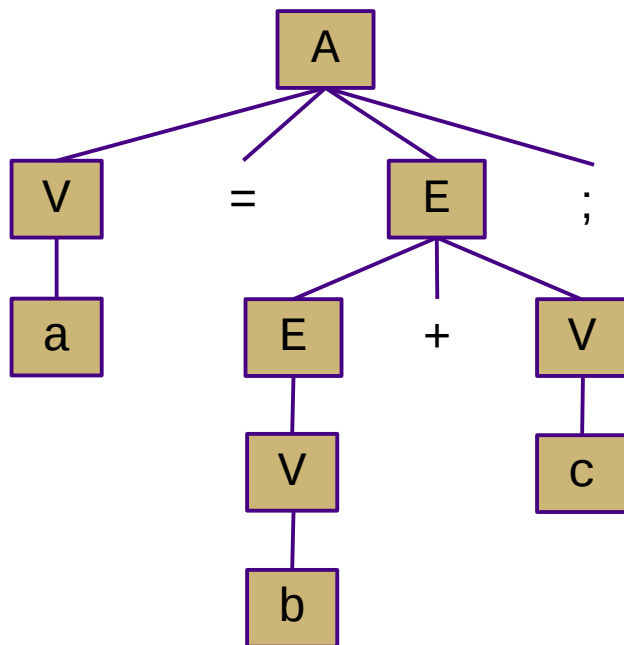
- Rule 1: **FOLLOW(S)** includes **EOF / \$**
 - Where S is the start symbol
- Rule 2: for every production $A \rightarrow \alpha B \beta$
 - **FOLLOW(B)** includes everything in **FIRST(β)** except ϵ
- Rule 3: if $A \rightarrow \alpha B$ or ($A \rightarrow \alpha B \beta$ and **FIRST(β)** contains ϵ)
 - **FOLLOW(B)** includes everything in **FOLLOW(A)**

Aside: abstract syntax trees

Grammar:

$$\begin{aligned} A &\rightarrow V = E ; \\ E &\rightarrow E + V \\ &\quad | V \\ V &\rightarrow a \quad | \quad b \quad | \quad c \end{aligned}$$

Parse tree:



Abstract syntax tree:

