

CS 432
Fall 2023

Mike Lam, Professor

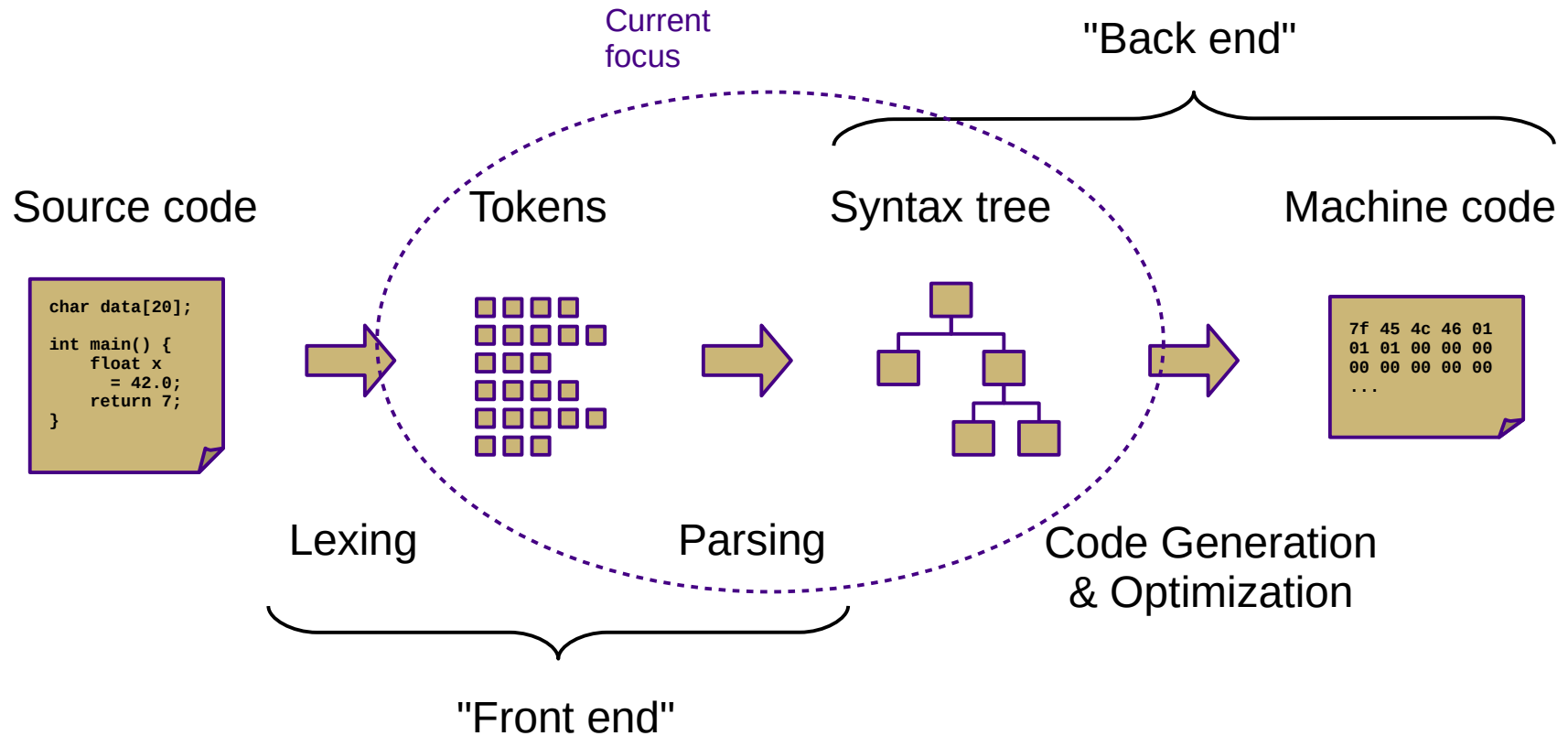
urban
DICTIONARY

recursion

See recursion.

Top-Down (LL) Parsing

Compilation



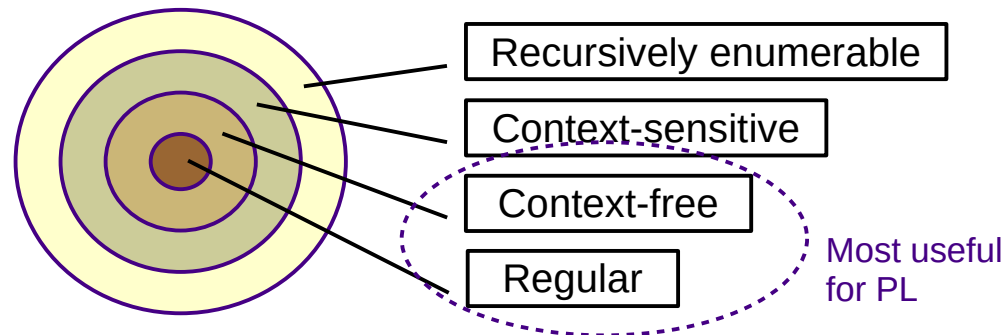
Review

- 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.

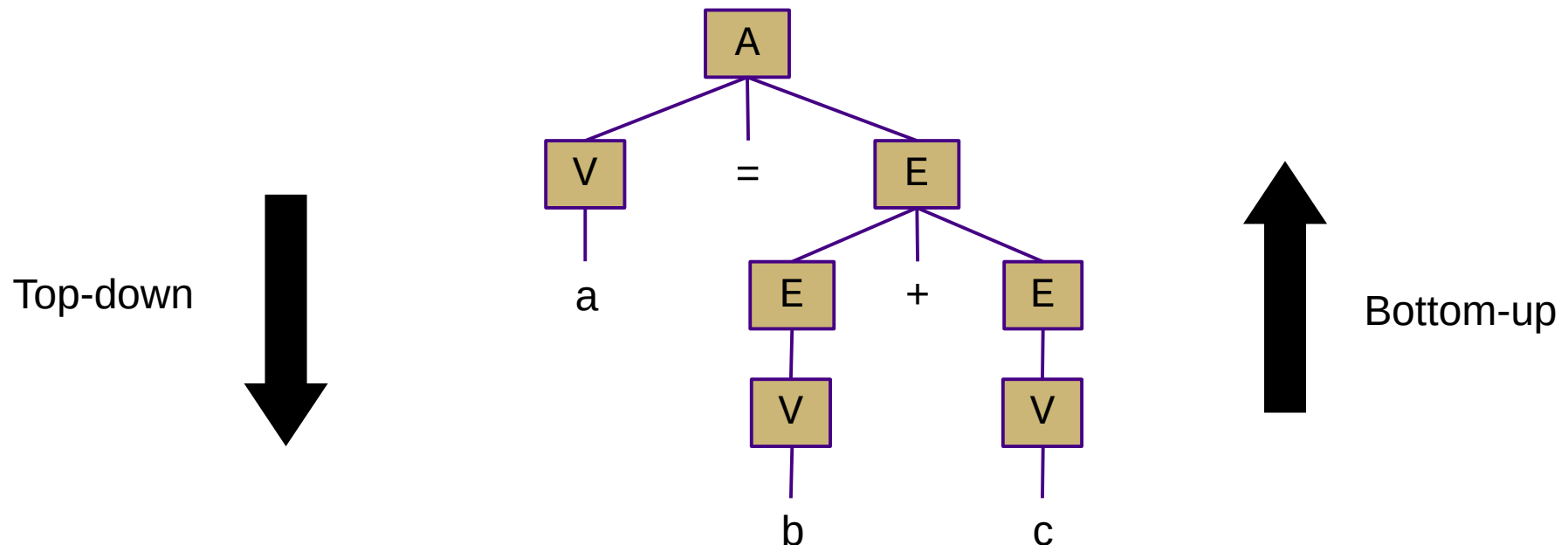
Chomsky Hierarchy of Languages



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 non-terminals that are still being 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

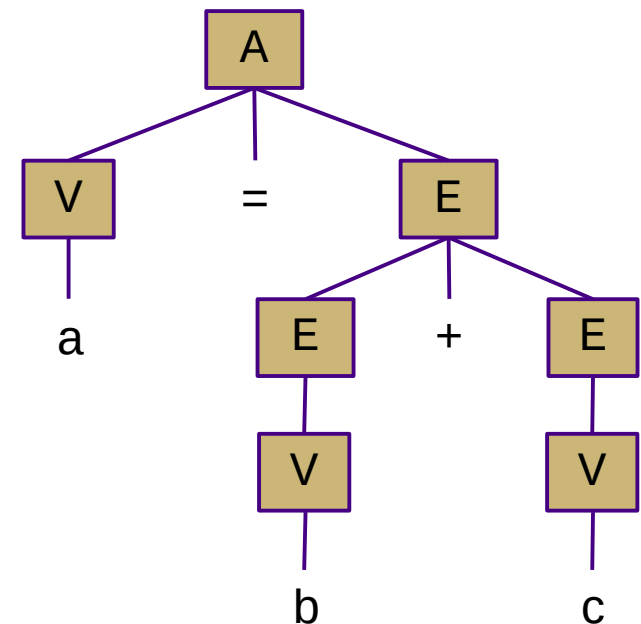
```
root = createNode(S)
focus = root
push(null)
token = nextToken()

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

$A \rightarrow \text{'if' } C \text{'then' } S$
 $\quad \quad \quad | \text{'goto' } L$



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

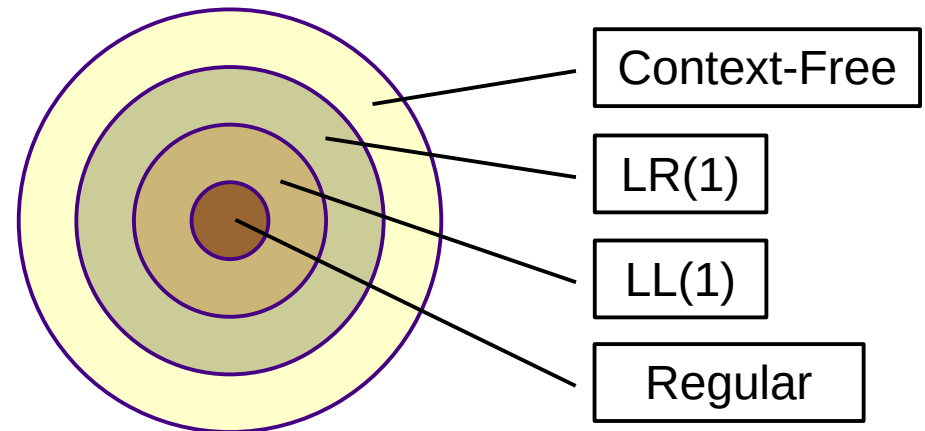
Top-Down Parsing

- Main issue: choosing which rule to use
 - In previous example, we just looked for ‘if’ or ‘goto’
 - 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
 - Must be careful to avoid backtracking

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 grammars to be LL(1)
 - Must remove left recursion
 - Must remove common prefixes (i.e., left factoring)
 - Easy (relatively) to hand-write a parser
 - **Practical** solution to real-world translation problems

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

Grammar with left recursion

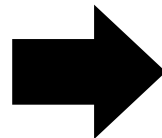
$$\begin{array}{l} A \rightarrow \alpha \beta_1 \\ | \alpha \beta_2 \end{array}$$

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 a top-down parser
 - To fix, unroll the recursion into a new non-terminal
 - **Practical note (P2):** A and A' can be a single function 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:** $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \dots$
 - Leads to ambiguous rule choice in a top-down 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 (P2):** A and A' can be a single function in your code
 - Parse and save data about α in temporary variables until you have enough information to choose

$$\begin{array}{l} A \rightarrow \alpha \beta_1 \\ \quad \mid \alpha \beta_2 \\ \quad \dots \end{array} \quad \longrightarrow \quad \begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \\ \quad \mid \beta_2 \\ \quad \dots \end{array}$$

Examples

- Eliminating left recursion:

$$\begin{array}{l} E \rightarrow E + T \\ \quad | E - T \\ \quad | T \end{array} \quad \longrightarrow \quad \begin{array}{l} E \rightarrow T E' \\ E' \rightarrow + T E' \\ \quad | - T E' \\ \quad | \varepsilon \end{array}$$

- Left factoring:

$$\begin{array}{l} C \rightarrow \text{if } E \text{ B else B fi} \\ \quad | \text{if } E \text{ B fi} \end{array} \quad \longrightarrow \quad \begin{array}{l} C \rightarrow \text{if } E \text{ B } C' \\ C' \rightarrow \text{else B fi} \\ \quad | \text{fi} \end{array}$$

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 (or ϵ) that can appear at the start of a sentence derived from α (a terminal or non-terminal)
- **FOLLOW(A)** set
 - Set of terminals (or $\$$) 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)**

Example

- FIRST and FOLLOW sets:

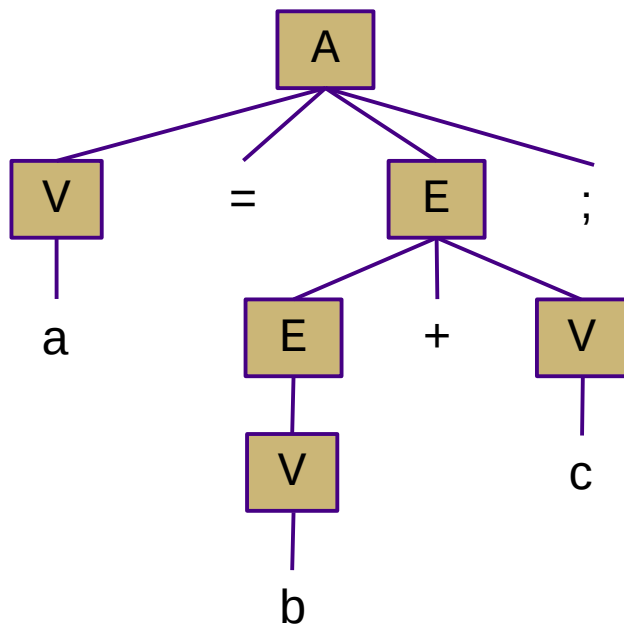
$$\begin{array}{l} A \rightarrow x A x \\ \quad | y B y \\ B \rightarrow C m \\ \quad | C \\ C \rightarrow t \end{array}$$
$$\begin{aligned} \text{FIRST}(x) &= \{ x \} \\ \text{FIRST}(y) &= \{ y \} \end{aligned}$$
$$\begin{aligned} \text{FIRST}(A) &= \{ x, y \} \\ \text{FIRST}(B) &= \{ t \} \\ \text{FIRST}(C) &= \{ t \} \end{aligned}$$
$$\begin{aligned} \text{FIRST}^+(A \rightarrow x A x) &= \{ x \} \\ \text{FIRST}^+(A \rightarrow y B y) &= \{ y \} \\ &\text{(disjoint: this is ok)} \end{aligned}$$
$$\begin{aligned} \text{FIRST}^+(B \rightarrow C m) &= \{ t \} \\ \text{FIRST}^+(B \rightarrow C) &= \{ t \} \\ &\text{(not disjoint: requires backtracking!)} \end{aligned}$$
$$\begin{aligned} \text{FOLLOW}(A) &= \{ x, \$ \} \\ \text{FOLLOW}(B) &= \{ y \} \\ \text{FOLLOW}(C) &= \{ y, m \} \end{aligned}$$

Aside: abstract syntax trees

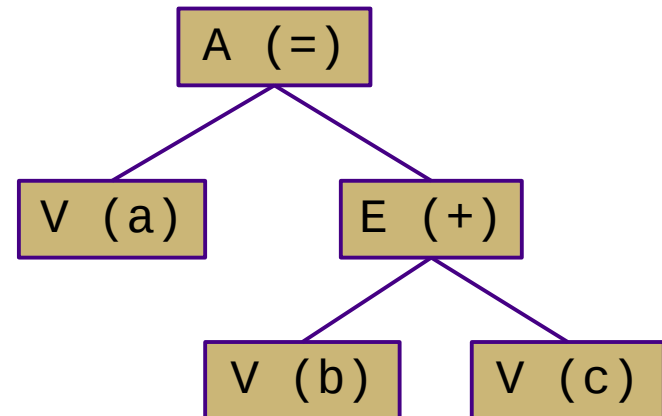
Grammar:

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

Parse tree:



Abstract syntax tree:



In P2, you will build an AST, not a parse tree!

Aside: Parser combinators

- A **parser combinator** is a higher-order function for parsing
 - Takes several parsers as inputs, returns new parser as output
 - Allows parser code to be very close to grammar
 - (Relatively) recent development: '90s and '00s
 - Example: [Parsec](#) in Haskell

```
whileStmt :: Parser Stmt
whileStmt =
  do keyword "while"
     cond <- expression
     keyword "do"
     stmt <- statement
     return (While cond stmt)
```

```
assignStmt :: Parser Stmt
assignStmt =
  do var <- identifier
     operator " := "
     expr <- expression
     return (Assign var expr)
```