

Lecture 11

CS 131: COMPILERS

Announcements

- HW3: LLVM lite
 - Available on Blackboard.
 - Due: November 6th at 11:59:59pm

**you should have
*STARTED EARLY!!***

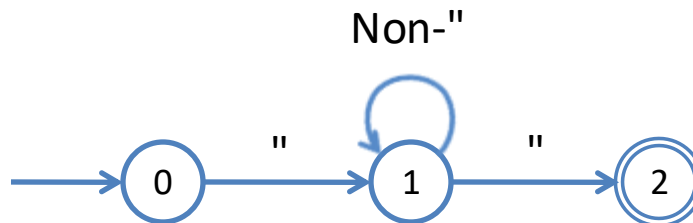
- Midterm: November 19th (tentative)
 - In class
 - One-page, letter-sized, double-sided “cheat sheet” of notes permitted
 - Coverage: interpereters, x86, LLVMlite, lexing, parsing
 - See examples of previous exam on Blackboard

Finite Automata

- Consider the regular expression: `""[^\"]*"''`
- An automaton (DFA) can be represented as:
 - A transition table:

	"	Non-"
0	1	ERROR
1	2	1
2	ERROR	ERROR

- A graph:



RE to Finite Automaton?

- Can we build a finite automaton for every regular expression?
 - Yes!
- Strategy: consider every possible regular expression (by induction on the structure of the regular expressions):

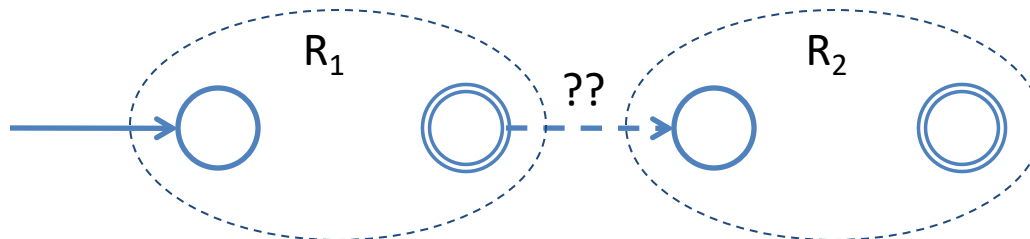
'a'



ϵ



R_1R_2

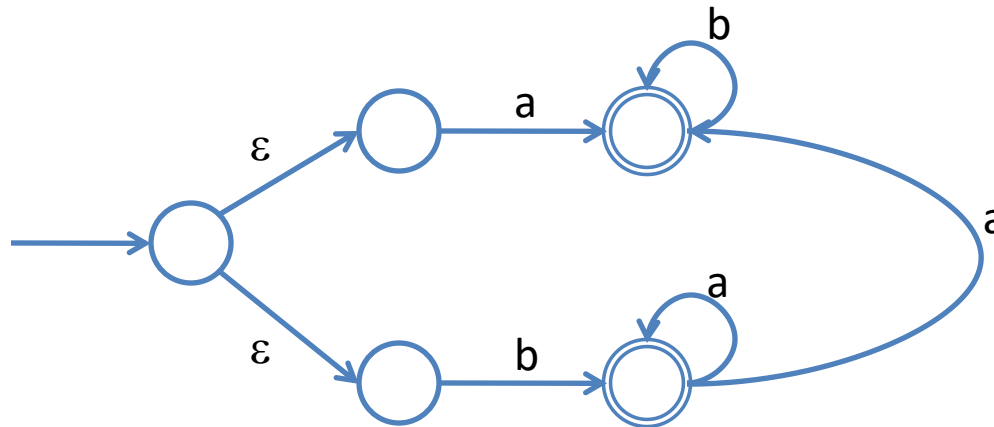


What about?

$R_1 | R_2$

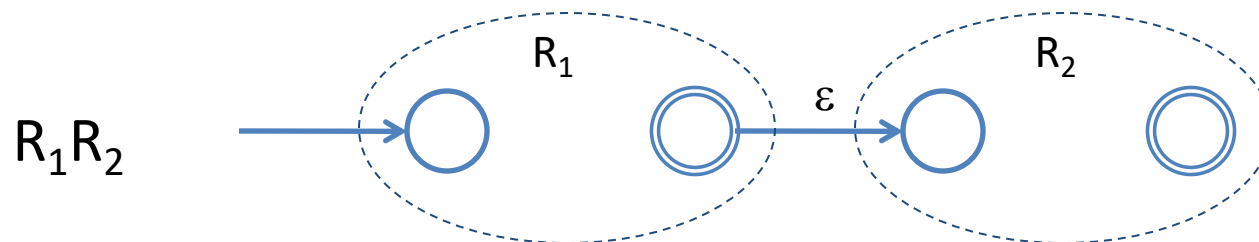
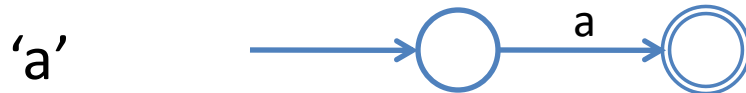
Nondeterministic Finite Automata

- A finite set of states, a start state, and accepting state(s)
- Transition arrows connecting states
 - Labeled by input symbols
 - Or ϵ (which does not consume input)
- *Nondeterministic*: two arrows leaving the same state may have the same label



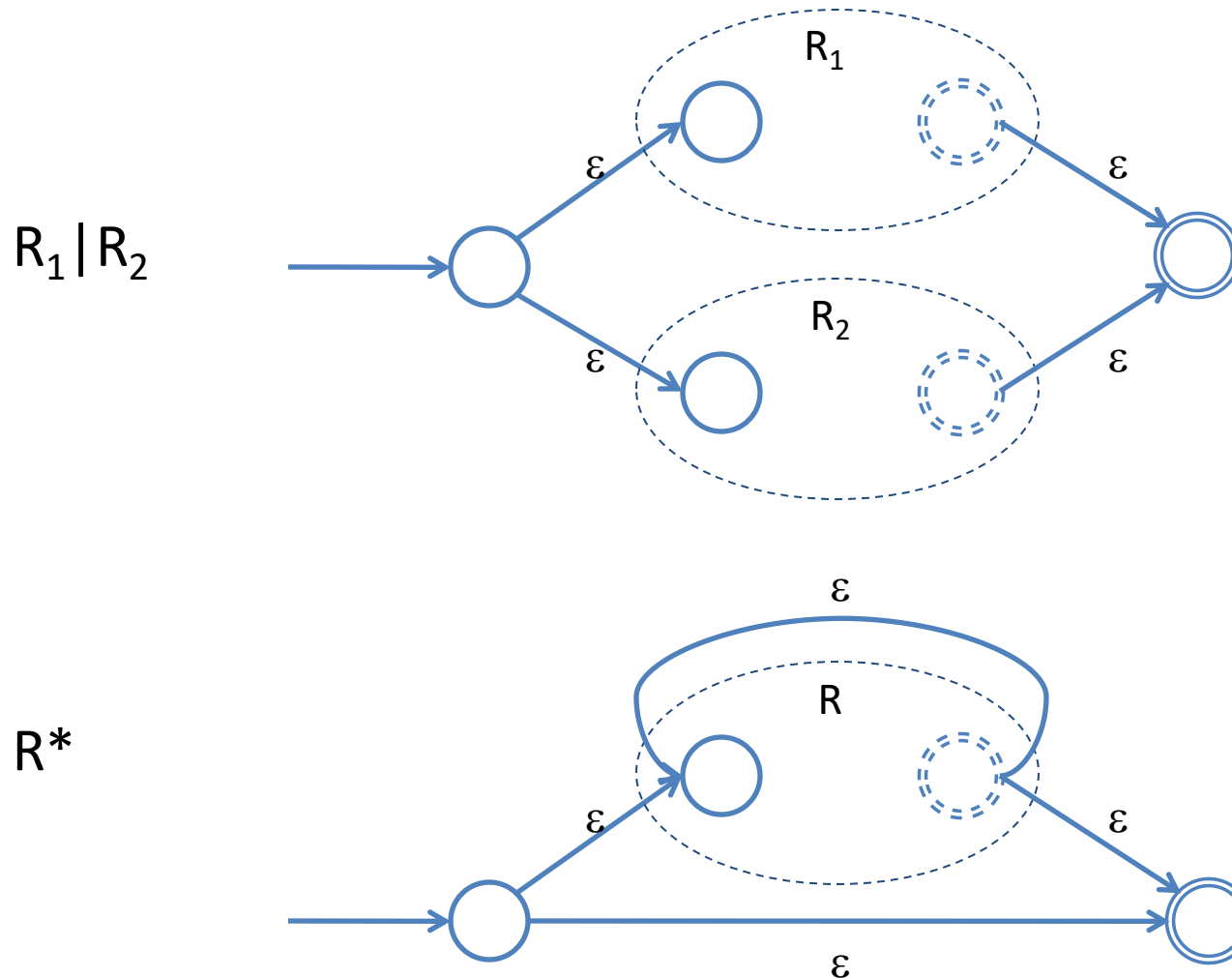
RE to NFA?

- Converting regular expressions to NFAs is easy.
- Assume each NFA has one start state, unique accept state



RE to NFA (cont'd)

- Sums and Kleene star are easy with NFAs



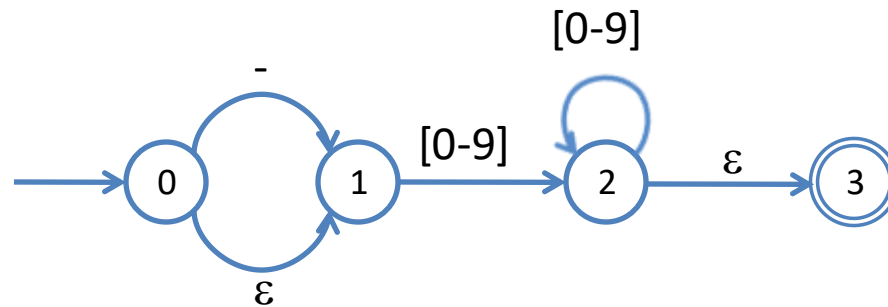
DFA versus NFA

- DFA:
 - Action of the automaton for each input is fully determined
 - Automaton accepts if the input is consumed upon reaching an accepting state
 - Obvious table-based implementation
- NFA:
 - Automaton potentially has a choice at every step
 - Automaton accepts an input string if there *exists* a way to reach an accepting state
 - Less obvious how to implement efficiently

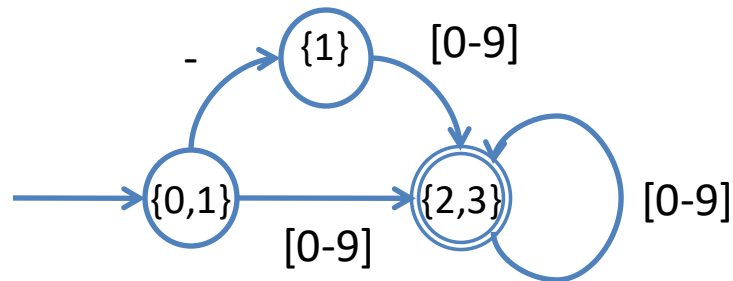
NFA to DFA conversion (Intuition)

- Idea: Run all possible executions of the NFA “in parallel”
- Keep track of a set of possible states: “finite fingers”
- Consider: $-?[0-9]^+$

- NFA representation:



- DFA representation:



Summary of Lexer Generator Behavior

- Take each regular expression R_i and its action A_i
- Compute the NFA formed by $(R_1 \mid R_2 \mid \dots \mid R_n)$
 - Remember the actions associated with the accepting states of the R_i
- Compute the DFA for this big NFA
 - There may be multiple accept states (why?)
 - A single accept state may correspond to one or more actions (why?)
- Compute the minimal equivalent DFA
 - There is a standard algorithm due to Myhill & Nerode
- Produce the transition table
- Implement longest match:
 - Start from initial state
 - Follow transitions, remember last accept state entered (if any)
 - Accept input until no transition is possible (i.e. next state is “ERROR”)
 - Perform the highest-priority action associated with the last accept state; if no accept state there is a lexing error

Lexer Generators in Practice

- Many existing implementations: lex, Flex, Jlex, ocamllex, ...
 - For example ocamllex program
 - see lexlex.mll, olex.mll, piglatin.mll on course website
- Error reporting:
 - Associate line number/character position with tokens
 - Use a rule to recognize ‘\n’ and increment the line number
 - The lexer generator itself usually provides character position info.
- Sometimes useful to treat comments specially
 - Nested comments: keep track of nesting depth
- Lexer generators are usually designed to work closely with parser generators...



Creating an abstract representation of program syntax.

PARSING

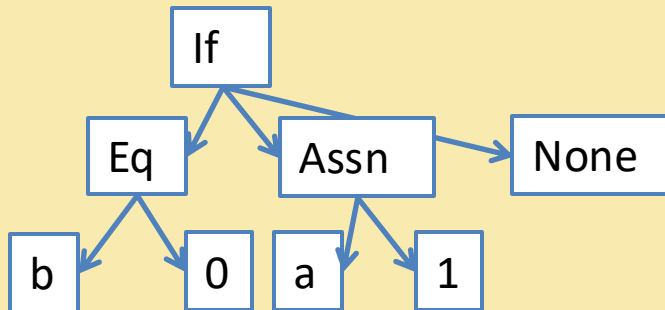
Parsing

Source Code
(Character stream)
if (b == 0) { a = 1; }

Token stream:

if	(b	==	0)	{	a	=	0	;	}
----	---	---	----	---	---	---	---	---	---	---	---

Abstract Syntax Tree:



Intermediate code:

```
l1:  
%cnd = icmp eq i64 %b, 0  
br i1 %cnd, label %l2, label %l3  
l2:  
store i64* %a, 1  
br label %l3  
l3:
```

Assembly Code

```
l1:  
cmpq %eax, $0  
jeq l2  
jmp l3  
l2:  
...
```

Lexical Analysis

Parsing

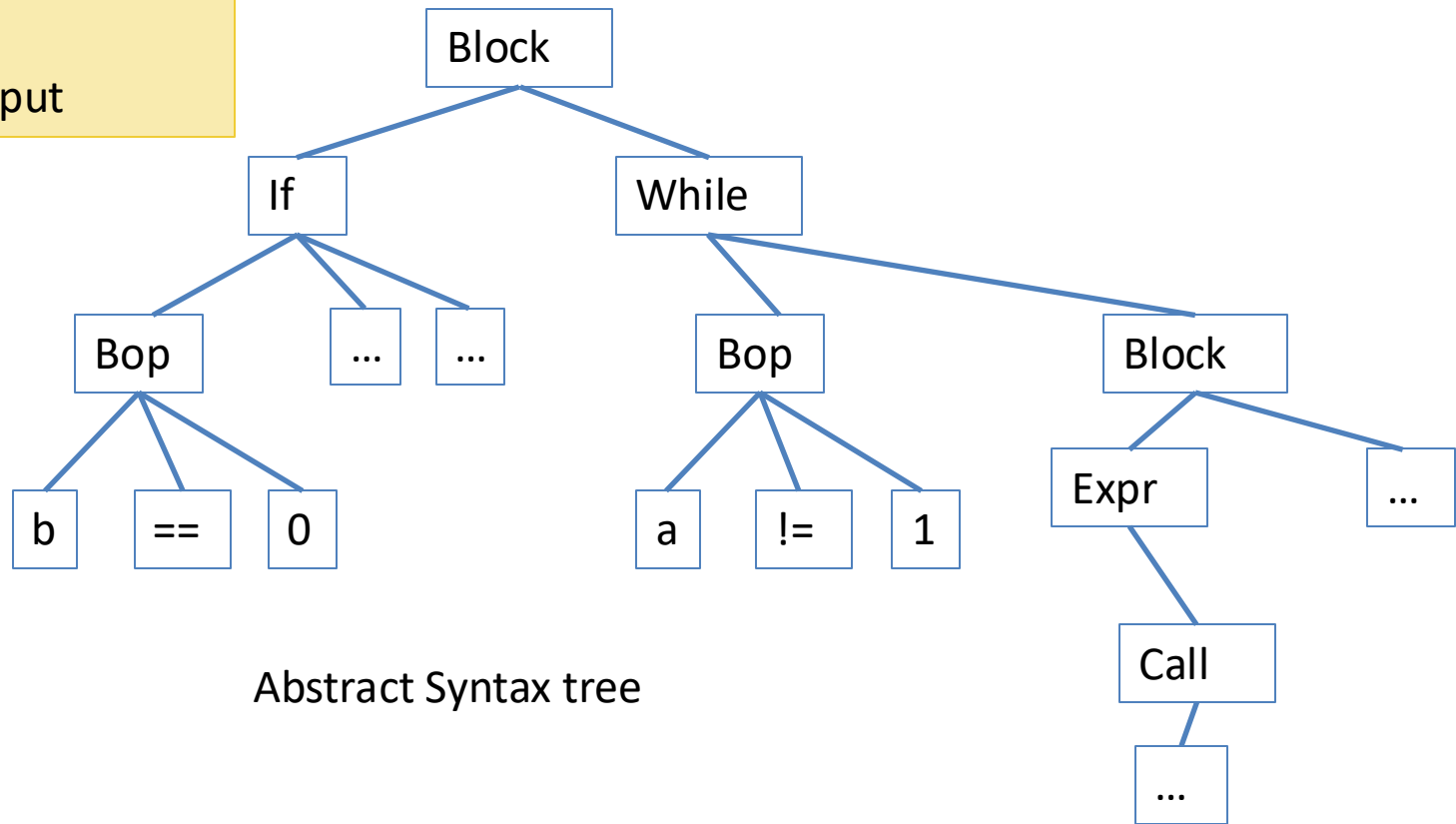
Analysis &
Transformation

Backend

Parsing: Finding Syntactic Structure

```
{  
  if (b == 0) a = b;  
  while (a != 1) {  
    print_int(a);  
    a = a - 1;  
  }  
}
```

Source input



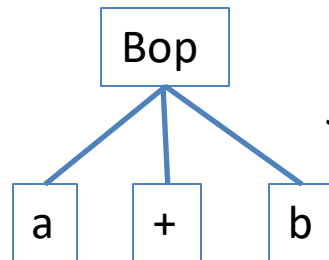
Syntactic Analysis (Parsing): Overview

- Input: stream of tokens (generated by lexer)
- Output: abstract syntax tree
- Strategy:
 - Parse the token stream to traverse the “concrete” syntax
 - During traversal, build a tree representing the “abstract” syntax
- Why abstract? Consider these three *different* concrete inputs:

a + b

(a + ((b)))

((a) + (b))



Same abstract syntax tree

- Note: parsing doesn't check many things:
 - Variable scoping, type agreement, initialization, ...

Specifying Language Syntax

- First question: how to describe language syntax precisely and conveniently?
- Previously we described tokens using regular expressions
 - Easy to implement, efficient DFA representation
 - Why not use regular expressions on tokens to specify programming language syntax?
- Limits of regular expressions:
 - DFA's have only finite # of states
 - So... DFA's can't "count"
 - For example, consider the language of all strings that contain balanced parentheses – easier than most programming languages, but not regular.
- So: we need more expressive power than DFA's



CONTEXT FREE GRAMMARS

Context-free Grammars

- Here is a specification of the language of balanced parens:

$$S \mapsto (S)S$$

$$S \mapsto \varepsilon$$

Note: Once again we have to take care to distinguish meta-language elements (e.g. “S” and “ \mapsto ”) from object-language elements (e.g. “(“).*

- The definition is *recursive* – S mentions itself.
- Idea: “derive” a string in the language by starting with S and rewriting according to the rules:
 - Example:
$$S \mapsto (S)S \mapsto ((S)S)S \mapsto ((\varepsilon)S)S \mapsto ((\varepsilon)S)\varepsilon \mapsto ((\varepsilon)\varepsilon)\varepsilon = (())$$
- You can replace the *nonterminal* S by one of its definitions anywhere
- A context-free grammar accepts a string iff there is a derivation from the start symbol

* And, since we’re writing this description in English, we are careful to distinguish the meta-meta-language (e.g. words) from the meta-language and object-language (e.g. symbols) by using quotes.

CFGs Mathematically

- A Context-free Grammar (CFG) consists of
 - A set of *terminals* (e.g., a lexical token or ϵ)
 - A set of *nonterminals* (e.g., S and other syntactic variables)
 - A designated nonterminal called the *start symbol*
 - A set of productions: $\text{LHS} \mapsto \text{RHS}$
 - LHS is a nonterminal
 - RHS is a *string* of terminals and nonterminals
- Example: The balanced parentheses language:

$$S \mapsto (S)S$$

$$S \mapsto \epsilon$$

- How many terminals? How many nonterminals? Productions?

Another Example: Sum Grammar

- A grammar that accepts parenthesized sums of numbers:

$$\begin{array}{l} S \mapsto E + S \quad | \quad E \\ E \mapsto \text{number} \quad | \quad (S) \end{array}$$

e.g.: $(1 + 2 + (3 + 4)) + 5$

- Note the vertical bar ‘|’ is shorthand for multiple productions:

$S \mapsto E + S$

$S \mapsto E$

$E \mapsto \text{number}$

$E \mapsto (S)$

4 productions

2 nonterminals: S, E

4 terminals: (,), +, number

Start symbol: S

Derivations in CFGs

- Example: derive $(1 + 2 + (3 + 4)) + 5$

- $\underline{S} \mapsto \underline{E} + S$

$$\mapsto (\underline{S}) + S$$

$$\mapsto (\underline{E} + S) + S$$

$$\mapsto (1 + \underline{S}) + S$$

$$\mapsto (1 + \underline{E} + S) + S$$

$$\mapsto (1 + 2 + \underline{S}) + S$$

$$\mapsto (1 + 2 + \underline{E}) + S$$

$$\mapsto (1 + 2 + (\underline{S})) + S$$

$$\mapsto (1 + 2 + (\underline{E} + S)) + S$$

$$\mapsto (1 + 2 + (3 + \underline{S})) + S$$

$$\mapsto (1 + 2 + (3 + \underline{E})) + S$$

$$\mapsto (1 + 2 + (3 + 4)) + \underline{S}$$

$$\mapsto (1 + 2 + (3 + 4)) + \underline{E}$$

$$\mapsto (1 + 2 + (3 + 4)) + 5$$

$$S \mapsto E + S \mid E$$

$$E \mapsto \text{number} \mid (S)$$

For arbitrary strings α, β, γ and production rule $A \mapsto \beta$ a single step of the derivation is:

$$\alpha A \gamma \mapsto \alpha \beta \gamma$$

(*substitute* β for an occurrence of A)

In general, there are many possible derivations for a given string.

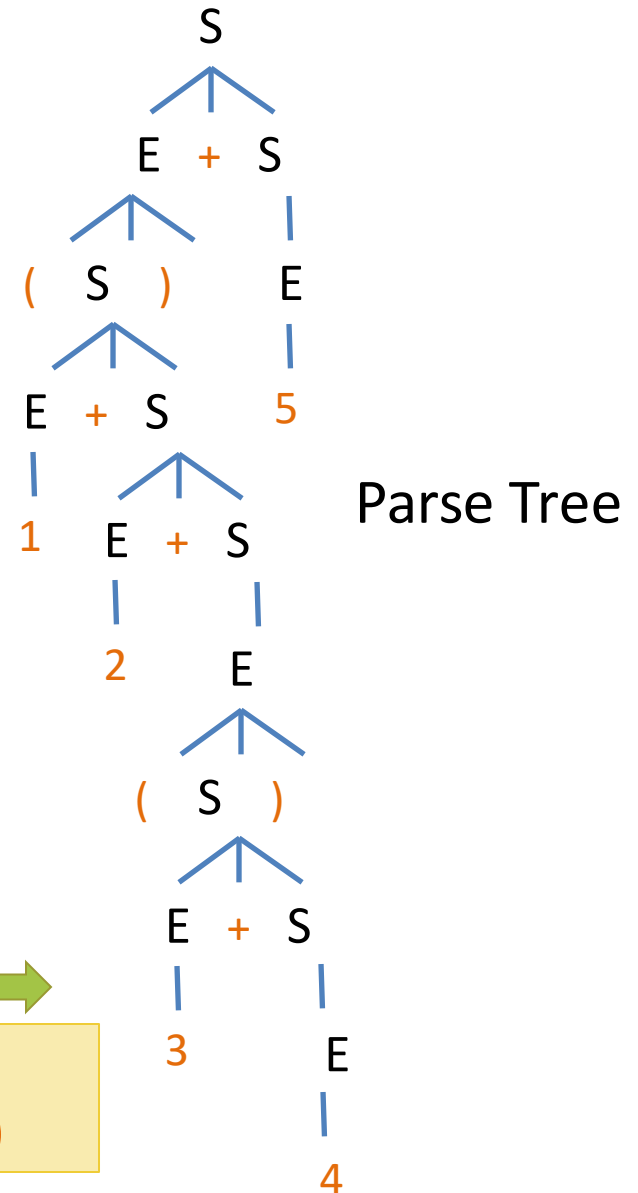
Note: Underline indicates symbol being expanded.

From Derivations to Parse Trees

- Tree representation of the derivation
- Leaves of the tree are terminals
 - In-order traversal yields the input sequence of tokens
- Internal nodes: nonterminals
- No information about the *order* of the derivation steps

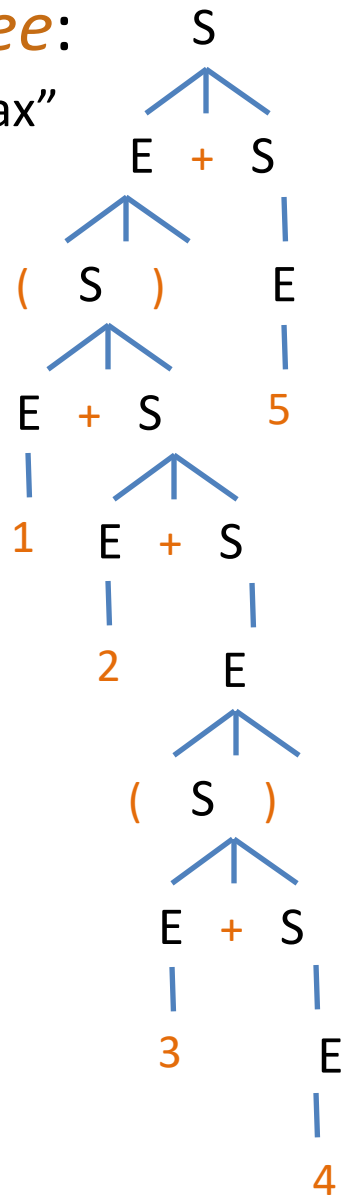
(1 + 2 + (3 + 4)) + 5

$S \mapsto E + S \mid E$
 $E \mapsto \text{number} \mid (S)$

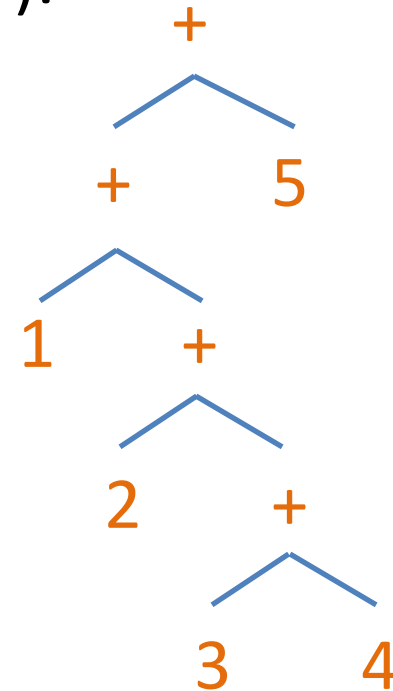


From Parse Trees to Abstract Syntax

- *Parse tree*:
“concrete syntax”



- *Abstract syntax tree*
(AST):



- Hides, or *abstracts*,
unneded information.

Derivation Orders

- Productions of the grammar can be applied in any order.
- There are two standard orders:
 - *Leftmost derivation*: Find the left-most nonterminal and apply a production to it.
 - *Rightmost derivation*: Find the right-most nonterminal and apply a production there.
- Note that both strategies (and any other) yield the same parse tree!
 - Parse tree doesn't contain the information about what order the productions were applied.

Example: Left- and rightmost derivations

- Leftmost Derivation

- $\underline{S} \mapsto \underline{E} + S$
 $\mapsto (\underline{S}) + S$
 $\mapsto (\underline{E} + S) + S$
 $\mapsto (1 + \underline{S}) + S$
 $\mapsto (1 + \underline{E} + S) + S$
 $\mapsto (1 + 2 + \underline{S}) + S$
 $\mapsto (1 + 2 + \underline{E}) + S$
 $\mapsto (1 + 2 + (\underline{S})) + S$
 $\mapsto (1 + 2 + (\underline{E} + S)) + S$
 $\mapsto (1 + 2 + (3 + \underline{S})) + S$
 $\mapsto (1 + 2 + (3 + \underline{E})) + S$
 $\mapsto (1 + 2 + (3 + 4)) + \underline{S}$
 $\mapsto (1 + 2 + (3 + 4)) + \underline{E}$
 $\mapsto (1 + 2 + (3 + 4)) + 5$

- Rightmost derivation:

- $\underline{S} \mapsto E + \underline{S}$
 $\mapsto E + \underline{E}$
 $\mapsto \underline{E} + 5$
 $\mapsto (\underline{S}) + 5$
 $\mapsto (E + \underline{S}) + 5$
 $\mapsto (E + E + \underline{S}) + 5$
 $\mapsto (E + E + \underline{E}) + 5$
 $\mapsto (E + E + (\underline{S})) + 5$
 $\mapsto (E + E + (E + \underline{S})) + 5$
 $\mapsto (E + E + (E + \underline{E})) + 5$
 $\mapsto (E + E + (\underline{E} + 4)) + 5$
 $\mapsto (E + \underline{E} + (3 + 4)) + 5$
 $\mapsto (\underline{E} + 2 + (3 + 4)) + 5$
 $\mapsto (1 + 2 + (3 + 4)) + 5$