

# Introduction to Parsing

CSCI 3136: Principles of  
Programming Languages

# Nonregular Languages

- Not all languages are regular.

E.g.  $L = \{0^n 1^n \mid n \geq 0\}$  is not regular.

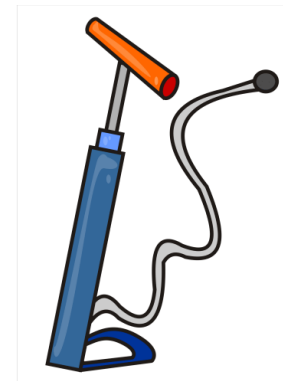
- Why?

- ★ Intuition: We need to keep track of how many 0's we encounter so we can match them with 1's.

- ★ A DFA has a finite number of states – that is its maximum "memory", so beyond that number we cannot keep track.

- How do we prove this formally?

The Pumping Lemma!



# Finiteness

- The language  $L = \{0^i1^i \mid i > 0\}$  is not regular (just believe me for now).
- $L$  is infinite in size as  $i$  is not bounded.
- Strings in  $L$  can be infinite in length.
- A DFA has a finite number of states.
- Hmmmm ... is this useful?

YES!

# Intuition: The Pumping Lemma

- Every Regular Language has a DFA recognizer.
- This recognizer has  $n$  states.
- If a string in the language is longer than  $n$ , the recognizer has to have a cycle.
- Thus, we can break the string into 3 parts:
  1. characters before the cycle
  2. the cyclic part
  3. characters after the cycle
- "Pumping" means "repeating the cycle"

# The Pumping Lemma

For every regular language  $L$ , there exists a constant  $n$  such that every  $\sigma \in L$ , where  $|\sigma| \geq n$ , can be divided into three substrings  $\sigma = \alpha\beta\gamma$  with the following properties:

- $|\alpha\beta| \leq n$
  - $|\beta| > 0$ , and
  - $\alpha\beta^k\gamma \in L, \forall k \geq 0$
- We can use this Lemma to show that a given language is non-regular.

# The Pumping Lemma

1.  $n$  is related to the # of states in the recognizing DFA.

2.  $\sigma = \alpha\beta\gamma$  if  $|\sigma| \geq n$

$\alpha$  is the part before the cycle

$\beta$  is the cyclic part (possibly with more than one iteration)

$\gamma$  is the part after the cycle

3. The properties mean:

- $|\alpha\beta| \leq n$  – getting to and through the cycle once can't use more than  $n$  states and hence can't have more characters

# Applying the Pumping Lemma

Procedure: To show that  $L$  is not regular

- Convince yourself  $L$  is not regular (intuition)
- Assume that  $L$  is regular and that there is a constant  $n$  as stated by the Pumping Lemma
- Select  $\sigma \in L$  such that
  - $|\sigma| > n$
  - $\sigma = \alpha\beta\gamma$ , for all  $\alpha$  and  $\beta$ 
    - $|\alpha\beta| \leq n$
    - $|\beta| > 0$
  - There exists a  $j \geq 0$  such that  $\alpha\beta^j\gamma \in L$

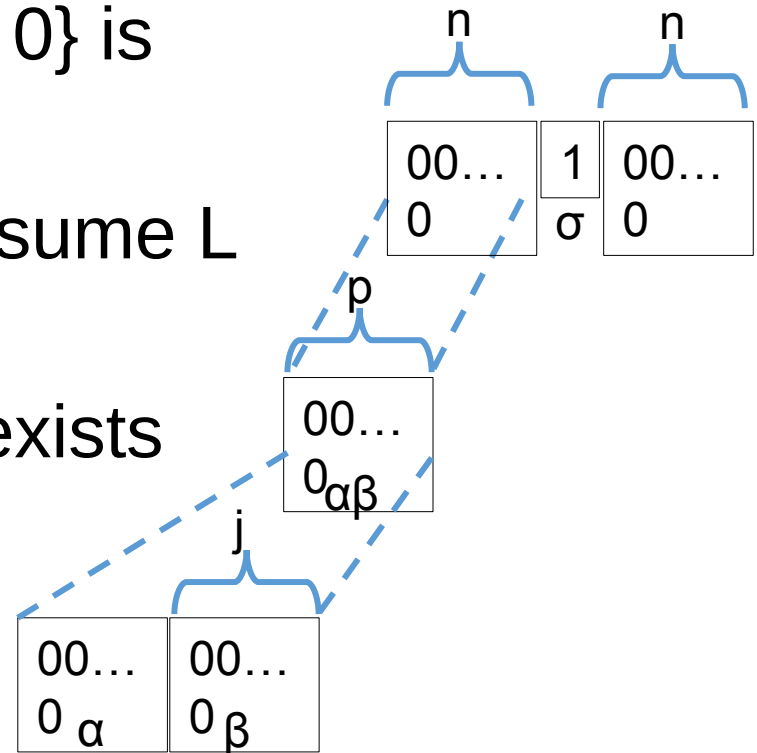
PART

# Example: Use the Pumping Lemma

Intuition: DFA can only keep track of  $n + 1$  things (at most).

Show that  $L = \{0^m 1 0^m \mid m \geq 0\}$  is not regular.

- Proof by contradiction: Assume  $L$  is regular.
- If  $L$  is regular, then there exists an  $n$ , by the PL
- Select  $\sigma = 0^n 1 0^n$
- Therefore, **for all  $\alpha$  and  $\beta$**



- $\alpha\beta = 0^p$ , which satisfies  $|\alpha\beta| \leq n$  as  $p$

- $\beta = 0^i$ ,  $0 \leq i \leq n$



# Or in other words ...

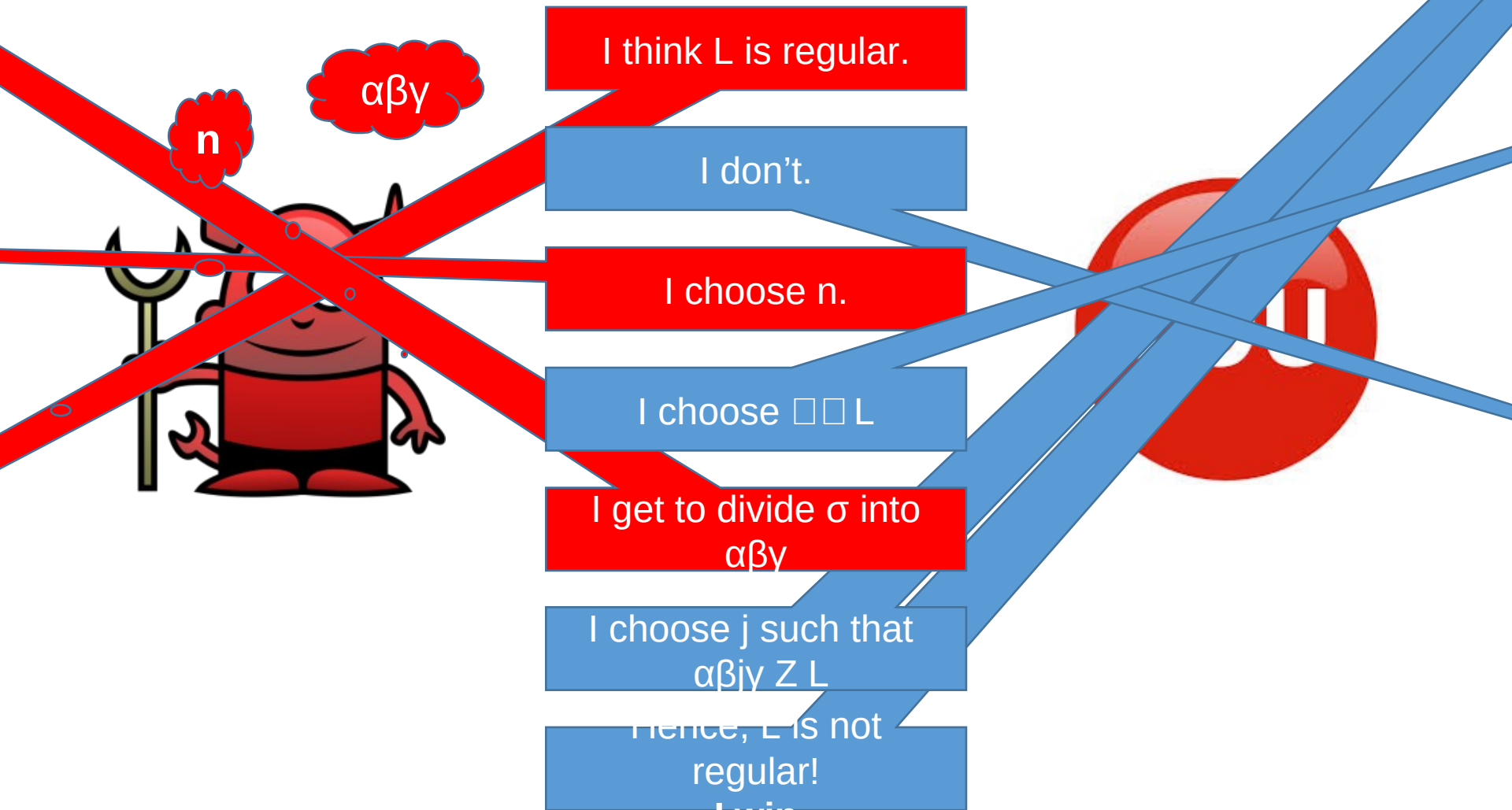
- $L = \{0^m 1 0^m \mid m \geq 0\}$
- So consider a small string for now ...

0000001000000

IF:

1. The cycle is before the 1, more cycles make a string with more 0's before the 1 ... it can't be there.
2. If the cycle is after the 1, more cycles make a string with more 0's after the 1 ... it can't be there
3. if the cycle contains the one, more cycles make a string with more 1's ... it can't be there.

# Using the Pumping Lemma is like an Argument with the Devil



# Examples

- $L = \{aibj \mid i < j\}$
- $L = \{ap \mid p \text{ is prime}\}$
- $L = \{aibj \mid i = j \bmod 3\}$

This one is actually regular

- Note: We cannot use the Pumping Lemma to prove a language is regular.
- Question: How do you show a language is regular?
  - Construct a regular expression for the language
  - Construct an NFA that recognizes the language

# Why Do We Need a Parser?

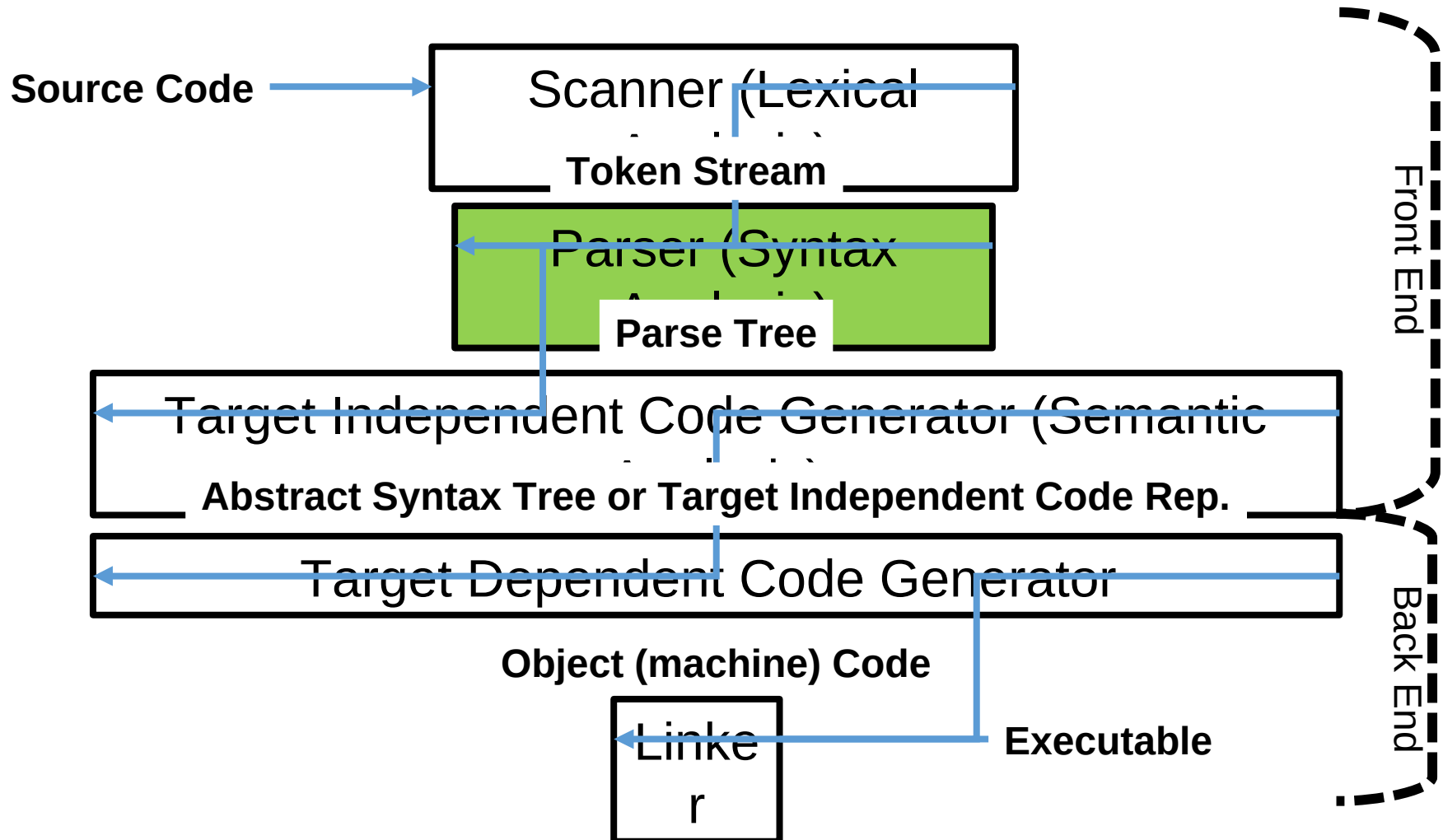
- A scanner yields a stream of tokens
- Q: Is this sufficient to determine if the input is a valid program?
- A: No! Most programming languages are not regular!

E.g. braces and brackets must match:  $((1 + 3) \times (3 + 2))$

Think of this as  $(i)i$ ,  $i > 0$ , which isn't regular as we know.

- Scanners are useful for
  - Checking if program's tokens are correct.

# Recall: Phases of Compilation



# Meet the Parser

- Parsing takes a stream of tokens
  - Checks whether the tokens represent a syntactically correct program.
  - Creates a parse tree (a high level representation of the program).
- Question: How do we know what the correct syntax is?
- Answer: Based on the language specification.
- Question: How do we specify the syntax?
- Answer: By a grammar.

# Grammars

- Idea: Grammars specify the syntax of a language:
- Example: English Sentences
  - *Sentence* → *Phrase Verb Phrase*
  - *Phrase* → *Noun* | *Adjective Phrase*
  - *Adjective* → big | small | green
  - *Noun* → boss | cheese
  - *Verb* → is | jumps | eats

Valid Sentences:

★Boss is big cheese.

# Example: Arithmetic Expressions

Grammar

Valid Sentences

$E \rightarrow E \text{ Op } E$

$(1 + 2 - 3) * 4$

$E \rightarrow - E$

$--3$

$E \rightarrow ( E )$

$a + b$

$E \rightarrow \text{Number}$

$E \rightarrow \text{Identifier}$

$\text{Op} \rightarrow +$

$\text{Op} \rightarrow -$

$\text{Op} \rightarrow /$

$\text{Op} \rightarrow *$

Typically programming languages are specified by Context Free Grammars (CFG)



# Context Free Grammars (CFG)

A CFG  $G$  is a 4-tuple  $G = (V, \Sigma, P, S)$  where

- $V$  is the set of non-terminals
  - Also known as “Variables”
  - Denoted by Capitalized letters/words
- $\Sigma$  is the set of terminals
  - The tokens returned by the scanner
- $P$  is the set of productions
  - Of the form  $N \rightarrow (\Sigma \cup V)^*$ ,  $N \in V$
  - Also known as “Rewriting Rules”
- $S$  is the start symbol,  $S \in V$

# A CFG Example: Expressions

- $V = \{E, Op\}$
- $\Sigma = \{\text{identifier, number, } (, ), +, -, X, /\}$
- $P = \{$

$E \rightarrow E Op E$

$E \rightarrow -E$

$E \rightarrow ( E )$

$E \rightarrow \text{number}$

$E \rightarrow \text{identifier}$

$Op \rightarrow +$

$Op \rightarrow -$

$Op \rightarrow X$

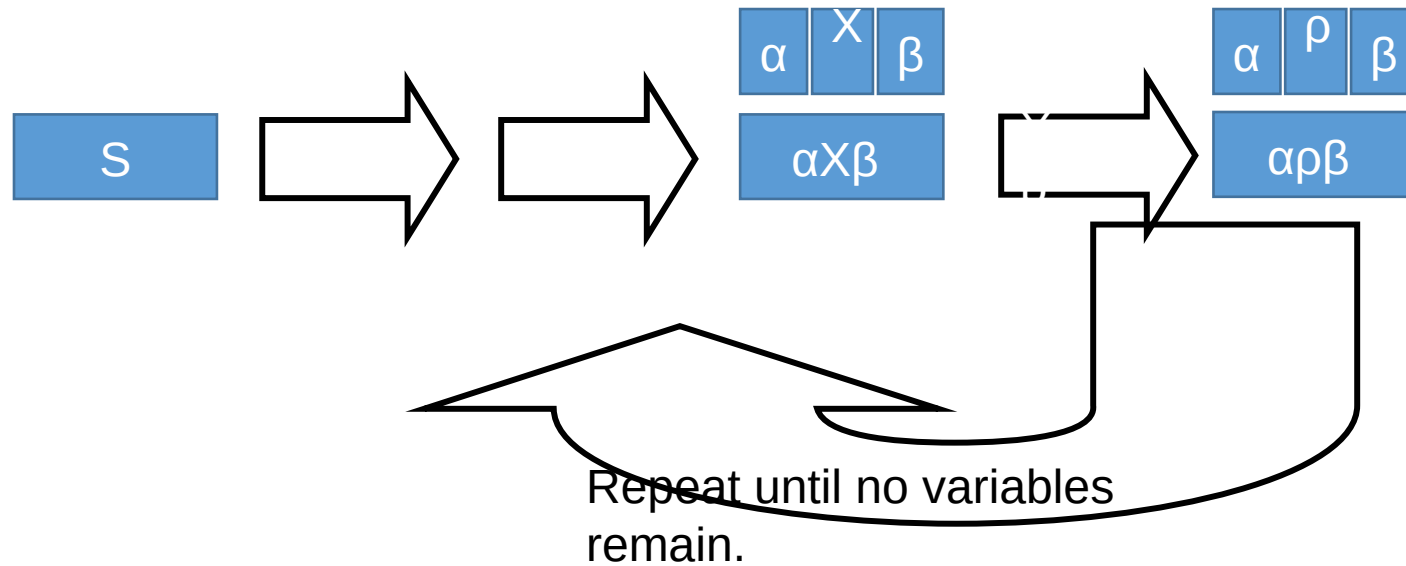
# Notes on CFG Notation

- Note: Alternative productions can be merged using |
  - E.g.,  $Op \rightarrow + \mid - \mid * \mid /$
- Several different notations are in use:
  - **Backus-Naur Form (BNF)** uses  $::=$  instead of  $\rightarrow$
  - **Optional Components notation** Nopt means that N is optional in the production
  - **Regular Expressions in RHS notation** allows regular expressions of terminals and nonterminals
- Question: How do we use a grammar?
- We determine whether a program is *derivable*

# Derivations

- A derivation is a sequence of rewriting operations that starts with the string  $\sigma = S$  and then repeats the following until  $\sigma$  contains only terminals:
  - Select a non-terminal in  $XEV$ , such that  $\sigma = \alpha X \beta$   
where  $\alpha, \beta \in (V \cup \Sigma)^*$
  - Select a production in  $(X \rightarrow \rho) \in P$ ,
  - Replace  $X$  with  $\rho$  in the partial derivation  $\sigma$   
i.e.,  $\sigma = \alpha \rho \beta$
- Eventually,  $\sigma$  will consist of only terminals, meaning the derivation is complete.

# Derivations in a Nutshell



# Derivation Example of an Expression

$\sigma = \mathbf{E}$

- ☐  $\mathbf{E} \text{ Op } \mathbf{E}$
- ☐  $( \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( \mathbf{E} \text{ Op } \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 \mathbf{Op} \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 + \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 + 13 ) \mathbf{Op} \mathbf{E}$
- ☐  $( 42 + 13 ) \mathbf{X} \mathbf{E}$
- ☐  $( 42 + 13 ) \mathbf{X} 11$

*Grammar*

1.  $E \rightarrow E \text{ Op } E$
2.  $E \rightarrow - E$
3.  $E \rightarrow ( E )$
4.  $E \rightarrow \text{Number}$
5.  $E \rightarrow \text{Identifier}$
6.  $\text{Op} \rightarrow +$

# Definitions

- Definition: We write  $S \sqsubseteq^* \sigma$  if there exists a derivation

$$S \sqsubseteq \sigma_1 \sqsubseteq \sigma_2 \sqsubseteq \dots \sqsubseteq \sigma$$

- Definition: Every grammar  $G$  defines a language:

$$L(G) = \{\sigma \in \Sigma^* \mid S \sqsubseteq \sigma\}$$

- Definition: If  $G$  is a context-free grammar then  $L(G)$  is a context-free language.
- Example: What is the language defined by  $G = (V, \Sigma, P, S)$

# Example 2

- What is the language defined by  $G = (V, \Sigma, P, S)$

- $V = \{S\}$

- $\Sigma = \{0, 1, \varepsilon\}$

- $P = \{$

- $S \rightarrow \varepsilon$

- $S \rightarrow 0S0$

- $S \rightarrow 1S1$

- $\}$

- $S = S$

The language  $L(G) = \{\sigma\sigma^r \mid \sigma \in \Sigma^+\}$



# Parse Trees

- A program is syntactically correct if it can be derived from the grammar of the language it is written in.
- To analyze the program we need a better representation of it.
  - I.e., tokens are the input to the parser
- So, each derivation can be represented by a parse tree.

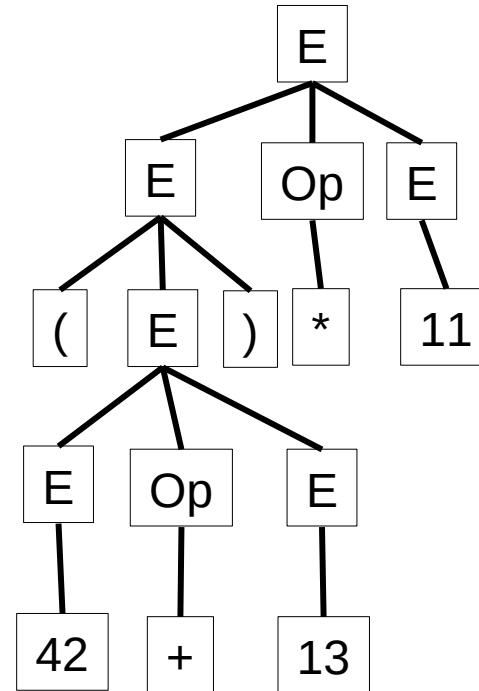
# Structure of Parse Trees

- Root:  $S$ , the start nonterminal
- Internal nodes: nonterminals
- Leaf nodes: terminals (called the *yield* of the tree)
- Edge( $X, w$ ) :  $XEV, wE\alpha$ , where  $(X \rightarrow \alpha) \in P$ .

# Parse Tree Example of an Expression

$\sigma = \mathbf{E}$

- ☐  $\mathbf{E} \text{ Op } \mathbf{E}$
- ☐  $( \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( \mathbf{E} \text{ Op } \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 \text{ **Op** } \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 + \mathbf{E} ) \text{ Op } \mathbf{E}$
- ☐  $( 42 + 13 ) \text{ **Op** } \mathbf{E}$
- ☐  $( 42 + 13 ) \times \mathbf{E}$
- ☐  $( 42 + 13 ) \times 11$



# Another Example: $1 + 2 * 3$

## This is ambiguous!

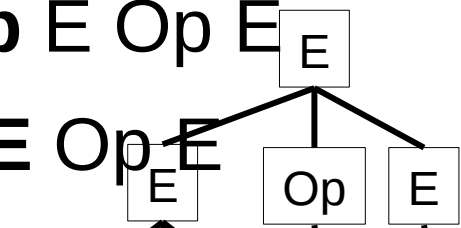
$\sigma = E$

☐  $E \text{ Op } E$

☐  $E \text{ Op } E \text{ Op } E$

☐  $1 \text{ Op } E \text{ Op } E$   $(1 + 2) * 3$

☐  $1 + E \text{ Op } E$



A partial parse tree for the expression  $(1 + 2) * 3$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The left  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The leftmost  $E$  node has a single child  $1$ . The rightmost  $E$  node has a single child  $2$ . The rightmost  $E$  child of the root has a single child  $3$ .

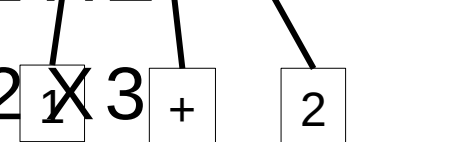
☐  $1 + 2 \text{ Op } E$



A partial parse tree for the expression  $1 + (2 * 3)$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The left  $E$  child has a single child  $1$ . The right  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The leftmost  $E$  node has a single child  $2$ . The rightmost  $E$  node has a single child  $3$ .

☐  $1 + 2 \text{ X } E$

☐  $1 + 2 \text{ X } 3$



A partial parse tree for the expression  $1 + 2 * 3$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The left  $E$  child has a single child  $1$ . The right  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The leftmost  $E$  node has a single child  $2$ . The rightmost  $E$  node has a single child  $3$ .

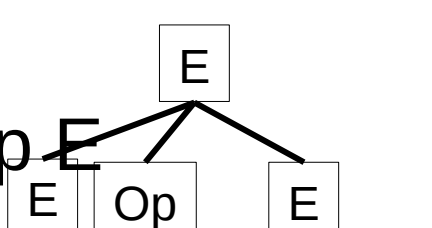
$\sigma = E$

☐  $E \text{ Op } E$

☐  $1 \text{ Op } E$

☐  $1 + E$   $1 + (2 * 3)$

☐  $1 + E \text{ Op } E$



A partial parse tree for the expression  $1 + (2 * 3)$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The left  $E$  child has a single child  $1$ . The right  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The leftmost  $E$  node has a single child  $2$ . The rightmost  $E$  node has a single child  $3$ .

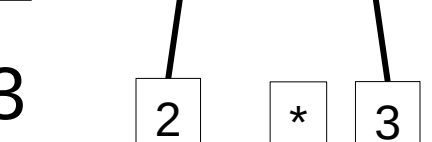
☐  $1 + 2 \text{ Op } E$



A partial parse tree for the expression  $1 + 2 * 3$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The left  $E$  child has a single child  $1$ . The right  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The leftmost  $E$  node has a single child  $2$ . The rightmost  $E$  node has a single child  $3$ .

☐  $1 + 2 \text{ X } E$

☐  $1 + 2 \text{ X } 3$



A partial parse tree for the expression  $1 + 2 * 3$ . The root node is  $E$ , which has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $+$ . The left  $E$  child has a single child  $1$ . The right  $E$  child has three children:  $E$ ,  $\text{Op}$ , and  $E$ . The middle  $\text{Op}$  node has a single child  $*$ . The leftmost  $E$  node has a single child  $2$ . The rightmost  $E$  node has a single child  $3$ .

# Ambiguity

- Observations:
  - There are infinitely many grammars to specify the same language
  - There may be multiple parse trees for the same sentence!
- Definition: If multiple parse trees can be generated by  $G$  for the same sentence, then  $G$  is *ambiguous*.
- Definition: If  $L$  does not have an unambiguous grammar, then  $L$  is *inherently ambiguous*
  - Usually not the case for programming languages!

# An Unambiguous Expression Grammar

Grammar

· Try deriving  $1 + 2 + 3$

- $E \rightarrow T$
- $E \rightarrow E + T$
- $E \rightarrow E - T$
- $T \rightarrow F$
- $T \rightarrow T \times F$
- $T \rightarrow T / F$
- $F \rightarrow \text{number}$
- $F \rightarrow \text{identifier}$