

# **Context-Free Languages**

# The Limits of Regular Languages

- The **pumping lemma for regular languages** can be used to establish limits on what languages are regular.
- If we want to describe more complex languages, we need a more powerful formalism.

# Context-Free Grammars

- A **context-free grammar** (or **CFG**) is an entirely different formalism for defining certain languages.
- CFGs are best explained by example...

# Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.
- Here is one possible CFG:

**E** → **int**

**E** → **E Op E**

**E** → **(E)**

**Op** → **+**

**Op** → **-**

**Op** → **\***

**Op** → **/**

**E**  
⇒ **E Op E**  
⇒ **E Op (E)**  
⇒ **E Op (E Op E)**  
⇒ **E \* (E Op E)**  
⇒ **int \* (E Op E)**  
⇒ **int \* (int Op E)**  
⇒ **int \* (int Op int)**  
⇒ **int \* (int + int)**

# Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.
- Here is one possible CFG:

**E** → **int**

**E** → **E Op E**

**E** → **(E)**

**Op** → **+**

**Op** → **-**

**Op** → **\***

**Op** → **/**

**E**  
⇒ **E Op E**  
⇒ **E Op int**  
⇒ **int Op int**  
⇒ **int / int**

# Context-Free Grammars

- Formally, a context-free grammar is a collection of four objects:
  - A set of **nonterminal symbols** (also called **variables**),
  - A set of **terminal symbols** (the **alphabet** of the CFG)
  - A set of **production rules** saying how each nonterminal can be converted by a string of terminals and nonterminals, and
  - A **start symbol** (which must be a nonterminal) that begins the derivation.

$$E \rightarrow \text{int}$$

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

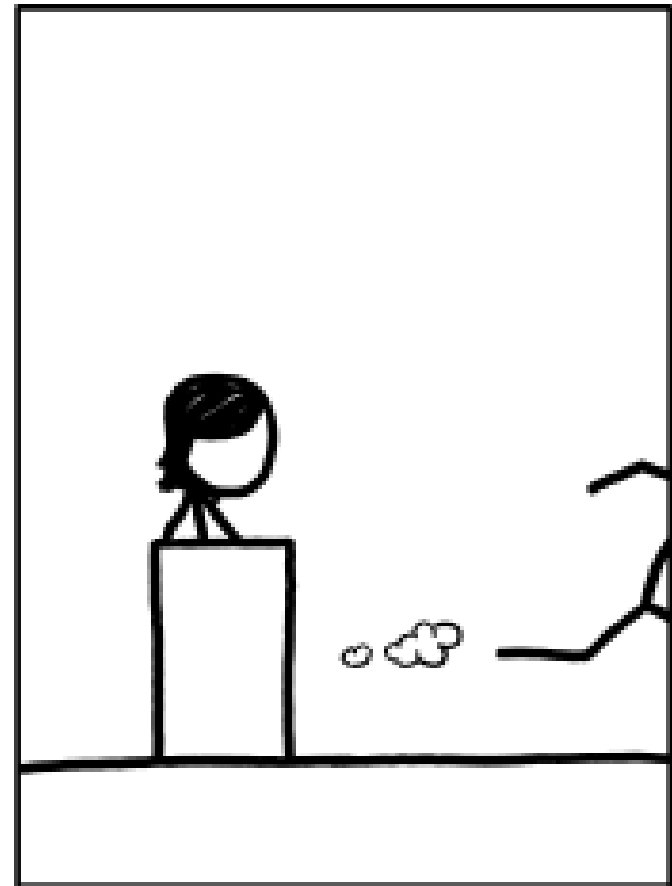
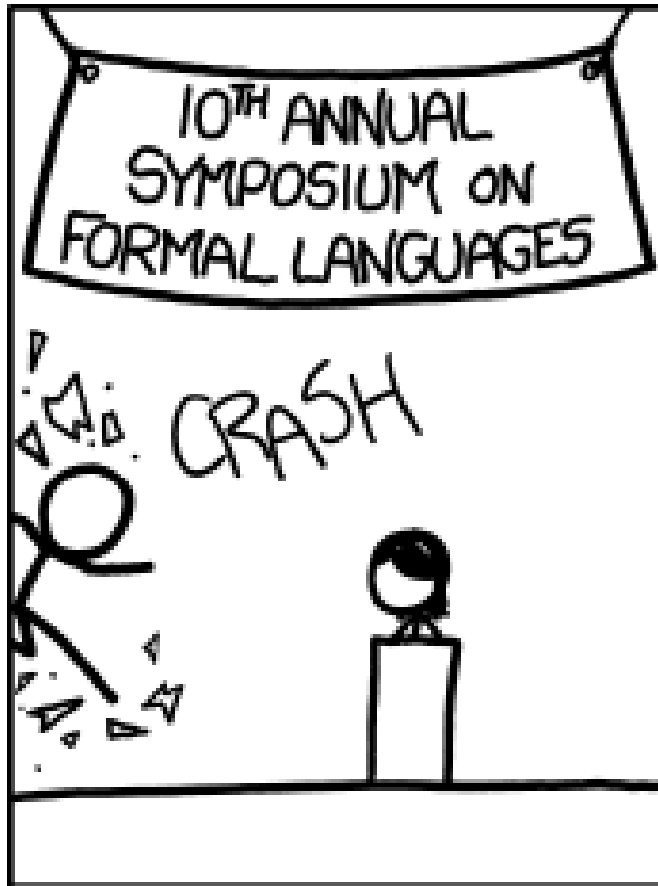
$$E \rightarrow (E)$$

$$\text{Op} \rightarrow +$$

$$\text{Op} \rightarrow -$$

$$\text{Op} \rightarrow *$$

$$\text{Op} \rightarrow /$$



# A Notational Shorthand

**E** → int

**E** → **E Op E**

**E** → (**E**)

**Op** → +

**Op** → -

**Op** → \*

**Op** → /



# A Notational Shorthand

**E** → *int* | **E Op E** | (**E**)

**Op** → + | - | \* | /

# Not Notational Shorthand

- The syntax for regular expressions does not carry over to CFGs.
- Cannot use \*, |, or parentheses.

**S** → **a\*b**

# Not Notational Shorthand

- The syntax for regular expressions does not carry over to CFGs.
- Cannot use \*, |, or parentheses.

$$\begin{aligned} S &\rightarrow Ab \\ A &\rightarrow Aa \mid \epsilon \end{aligned}$$

# Not Notational Shorthand

- The syntax for regular expressions does not carry over to CFGs.
- Cannot use \*, |, or parentheses.

**S** → **a (b | c\*)**

# Not Notational Shorthand

- The syntax for regular expressions does not carry over to CFGs.
- Cannot use \*, |, or parentheses.

**S** → **aX**

**X** → **b** | **C**

**C** → **Cc** | **ε**

# More Context-Free Grammars

- Chemicals!



**Form**  $\rightarrow$  **Cmp** | **Cmp Ion**

**Cmp**  $\rightarrow$  **Term** | **Term Num** | **Cmp Cmp**

**Term**  $\rightarrow$  **Elem** | **(Cmp)**

**Elem**  $\rightarrow$  **H** | **He** | **Li** | **Be** | **B** | **C** | ...

**Ion**  $\rightarrow$  **+** | **-** | **IonNum +** | **IonNum -**

**IonNum**  $\rightarrow$  **2** | **3** | **4** | ...

**Num**  $\rightarrow$  **1** | **IonNum**

# CFGs for Chemistry

**Form** → **Cmp** | **Cmp Ion**

**Cmp** → **Term** | **Term Num** | **Cmp Cmp**

**Term** → **Elem** | **(Cmp)**

**Elem** → **H** | **He** | **Li** | **Be** | **B** | **C** | ...

**Ion** → **+** | **-** | **IonNum +** | **IonNum -**

**IonNum** → **2** | **3** | **4** | ...

**Num** → **1** | **IonNum**

**Form**

⇒ **Cmp Ion**

⇒ **Cmp Cmp Ion**

⇒ **Cmp Term Num Ion**

⇒ **Term Term Num Ion**

⇒ **Elem Term Num Ion**

⇒ **Mn Term Num Ion**

⇒ **Mn Elem Num Ion**

⇒ **MnO Num Ion**

⇒ **MnO IonNum Ion**

⇒ **MnO<sub>4</sub> Ion**

⇒ **MnO<sub>4</sub><sup>-</sup>**

# CFGs for Programming Languages

**BLOCK** → **STMT**  
| **{ STMTS }**

**STMTS** →  $\epsilon$   
| **STMT STMTS**

**STMT** → **EXPR;**  
| **if (EXPR) BLOCK**  
| **while (EXPR) BLOCK**  
| **do BLOCK while (EXPR);**  
| **BLOCK**  
| ...

**EXPR** → **identifier**  
| **constant**  
| **EXPR + EXPR**  
| **EXPR - EXPR**  
| **EXPR \* EXPR**  
| ...



# Some CFG Notation

- Capital letters in **Bold Red Uppercase** will represent nonterminals.
  - i.e. **A, B, C, D**
- Lowercase letters in **blue monospace** will represent terminals.
  - i.e. **t, u, v, w**
- Lowercase Greek letters in *gray italics* will represent arbitrary strings of terminals and nonterminals.
  - i.e.  *$\alpha, \gamma, \omega$*

# Examples

- We might write an arbitrary production as

$$A \rightarrow \omega$$

- We might write a string of a nonterminal followed by a terminal as

$$At$$

- We might write an arbitrary production containing a nonterminal followed by a terminal as

$$B \rightarrow \alpha At \omega$$

# Derivations

$E$   
 $\Rightarrow E \text{ Op } E$   
 $\Rightarrow E \text{ Op } (E)$   
 $\Rightarrow E \text{ Op } (E \text{ Op } E)$   
 $\Rightarrow E * (E \text{ Op } E)$   
 $\Rightarrow \text{int} * (E \text{ Op } E)$   
 $\Rightarrow \text{int} * (\text{int} \text{ Op } E)$   
 $\Rightarrow \text{int} * (\text{int} \text{ Op } \text{int})$   
 $\Rightarrow \text{int} * (\text{int} + \text{int})$

- This sequence of steps is called a **derivation**.
- A string  $\alpha A \omega$  **yields** string  $\alpha \gamma \omega$  iff  $A \rightarrow \gamma$  is a production.
- If  $\alpha$  yields  $\beta$ , we write  $\alpha \Rightarrow \beta$ .
- We say that  $\alpha$  **derives**  $\beta$  iff there is a sequence of strings where
$$\alpha \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \beta$$
- If  $\alpha$  derives  $\beta$ , we write  $\alpha \Rightarrow^* \beta$ .

# The Language of a Grammar

- If  $G$  is a CFG with alphabet  $\Sigma$  and start symbol  $S$ , then the **language of  $G$**  is the set

$$\mathcal{L}(G) = \{ \omega \in \Sigma^* \mid S \Rightarrow^* \omega \}$$

- That is, the set of strings derivable from the start symbol.
- If  $L$  is a language and there is some CFG  $G$  such that  $L = \mathcal{L}(G)$ , then we say that  $L$  is a **context-free language** (or **CFL**).

# Regular and Context-Free Languages

- Consider the following CFG  $G$ :

$$S \rightarrow 0S1 \mid \epsilon$$

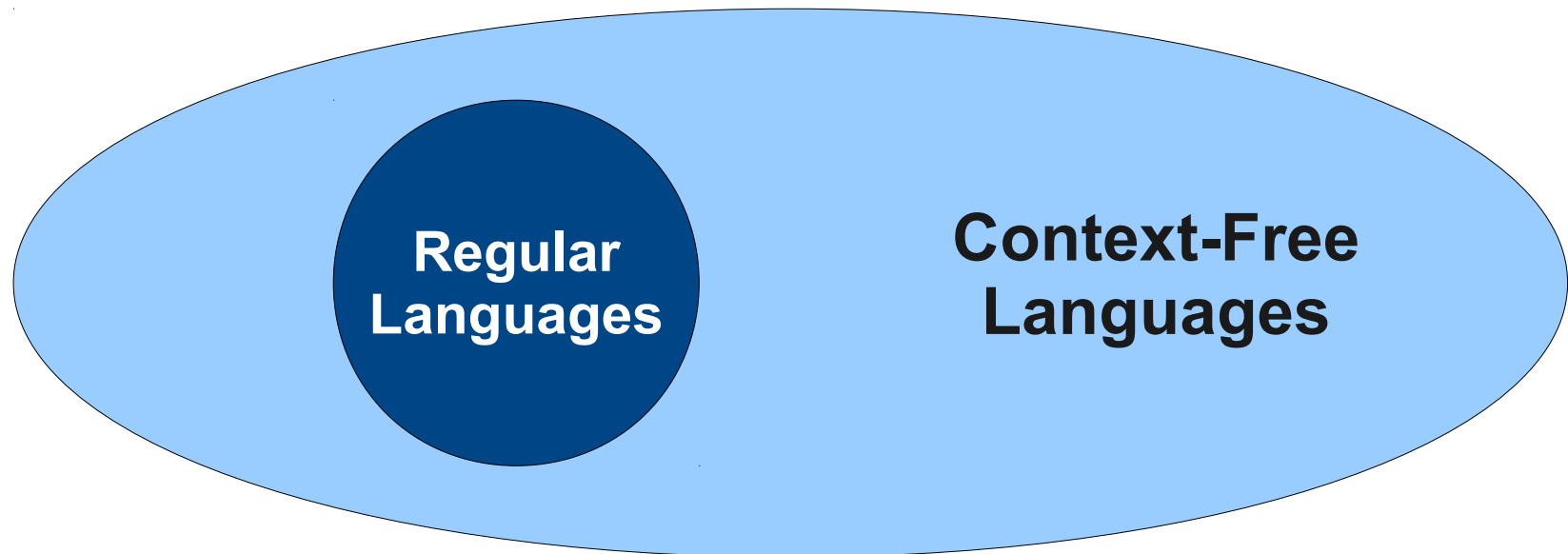
- What strings can this generate?

|   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|

$$\mathcal{L}(G) = \{ 0^n 1^n \mid n \in \mathbb{N} \}$$

# Regular and Context-Free Languages

- Context-free languages are a **strict superset** of the regular languages.
- Every regular language is context-free, but not necessarily the other way around.
- We'll see a proof of this next time.



# Leftmost Derivations

|              |   |  |
|--------------|---|--|
| <b>BLOCK</b> | → <b>STMT</b><br>  <b>{ STMTS }</b>   |  |
| <b>STMTS</b> | → $\epsilon$<br>  <b>STMT STMTS</b>   | <b>STMTS</b><br>⇒ <b>STMT STMTS</b>  |
| <b>STMT</b>  | → <b>EXPR;</b><br>  <b>if (EXPR) BLOCK</b><br>  <b>while (EXPR) BLOCK</b><br>  <b>do BLOCK while (EXPR);</b><br>  <b>BLOCK</b><br>  ...           | ⇒ <b>EXPR; STMTS</b><br>⇒ <b>EXPR = EXPR; STMTS</b><br>⇒ <b>id = EXPR; STMTS</b><br>⇒ <b>id = EXPR + EXPR; STMTS</b> |
| <b>EXPR</b>  | → <b>identifier</b><br>  <b>constant</b><br>  <b>EXPR + EXPR</b><br>  <b>EXPR - EXPR</b><br>  <b>EXPR * EXPR</b><br>  <b>EXPR = EXPR</b><br>  ... | ⇒ <b>id = id + EXPR; STMTS</b><br>⇒ <b>id = id + constant; STMTS</b><br>⇒ <b>id = id + constant;</b>                 |

# Leftmost Derivations

- A **leftmost derivation** is a derivation in which each step expands the leftmost nonterminal.
- A **rightmost derivation** is a derivation in which each step expands the rightmost nonterminal.
- These will be of great importance next lecture when we discuss *pushdown automata*.



# Related Derivations

**E**  
⇒ **E Op E**  
⇒ int **Op E**  
⇒ int \* **E**  
⇒ int \* (**E**)  
⇒ int \* (**E Op E**)  
⇒ int \* (int **Op E**)  
⇒ int \* (int + **E**)  
⇒ int \* (int + int)

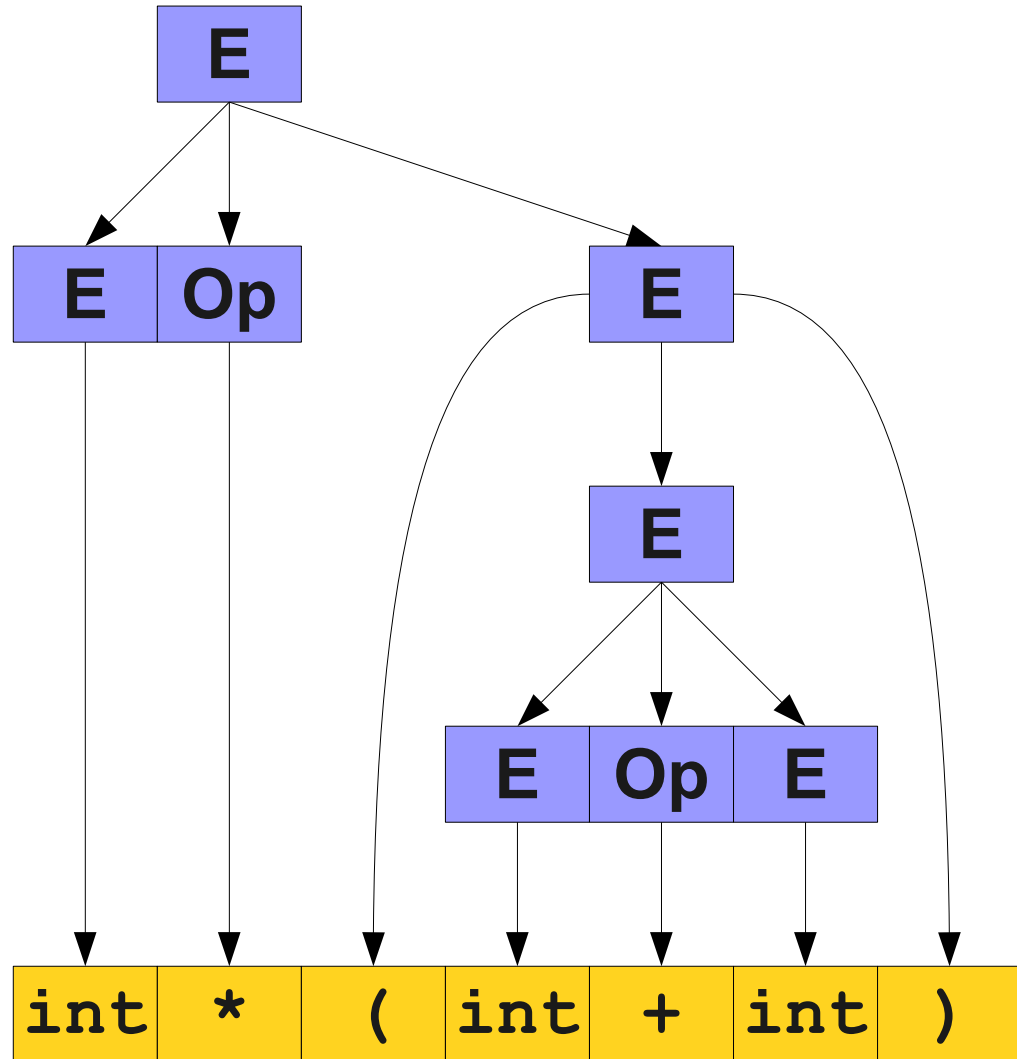
**E**  
⇒ **E Op E**  
⇒ **E Op (E)**  
⇒ **E Op (E Op E)**  
⇒ **E Op (E Op int)**  
⇒ **E Op (E + int)**  
⇒ **E Op (int + int)**  
⇒ **E \* (int + int)**  
⇒ int \* (int + int)

# Derivations Revisited

- A derivation encodes two pieces of information:
  - What productions were applied produce the resulting string from the start symbol?
  - In what order were they applied?
- Multiple derivations might use the same productions, but apply them in a different order.

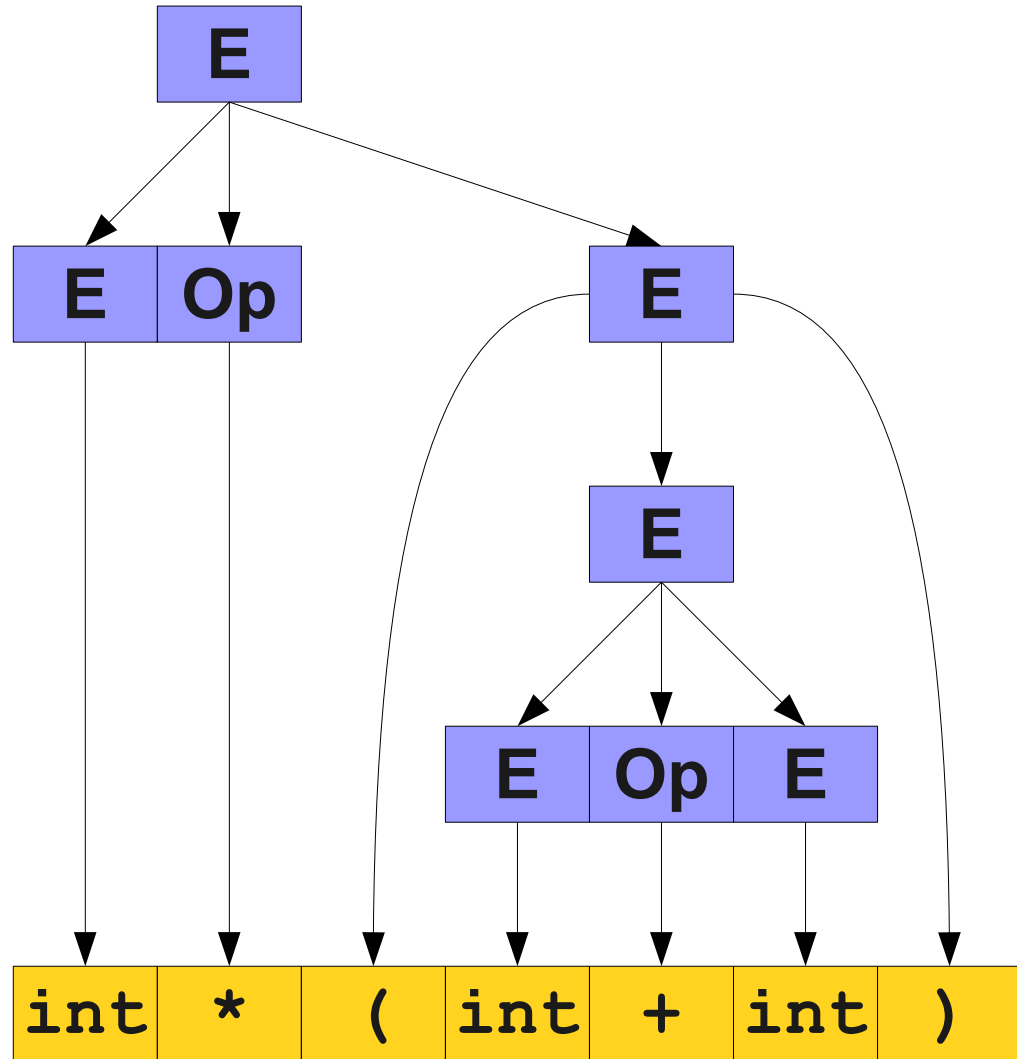
# Parse Trees

**E**  
⇒ **E Op E**  
⇒ **int Op E**  
⇒ **int \* E**  
⇒ **int \* (E)**  
⇒ **int \* (E Op E)**  
⇒ **int \* (int Op E)**  
⇒ **int \* (int + E)**  
⇒ **int \* (int + int)**



# Parse Trees

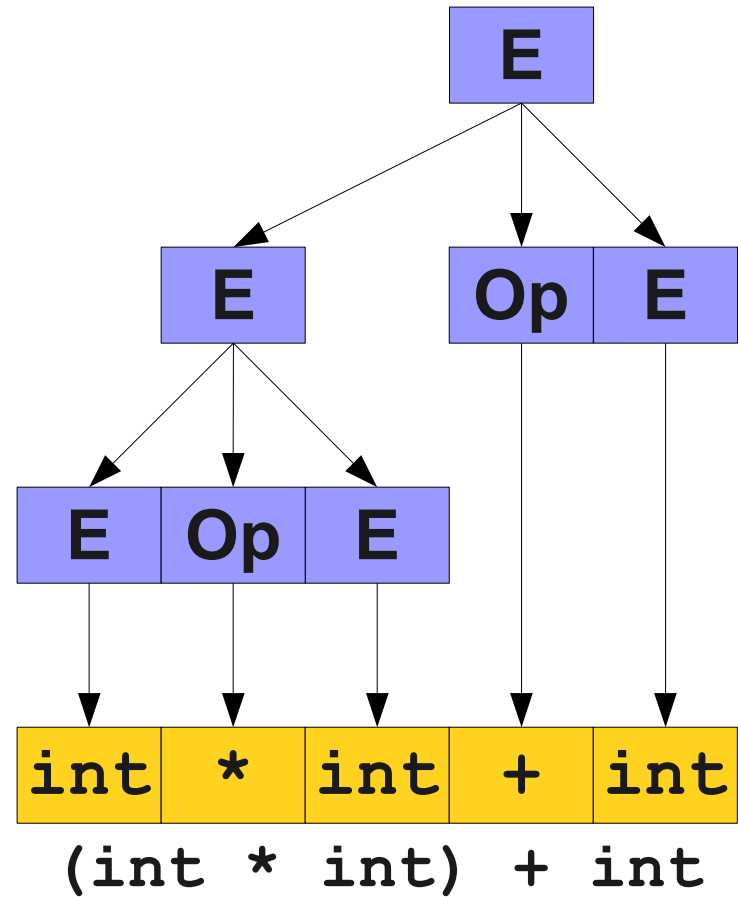
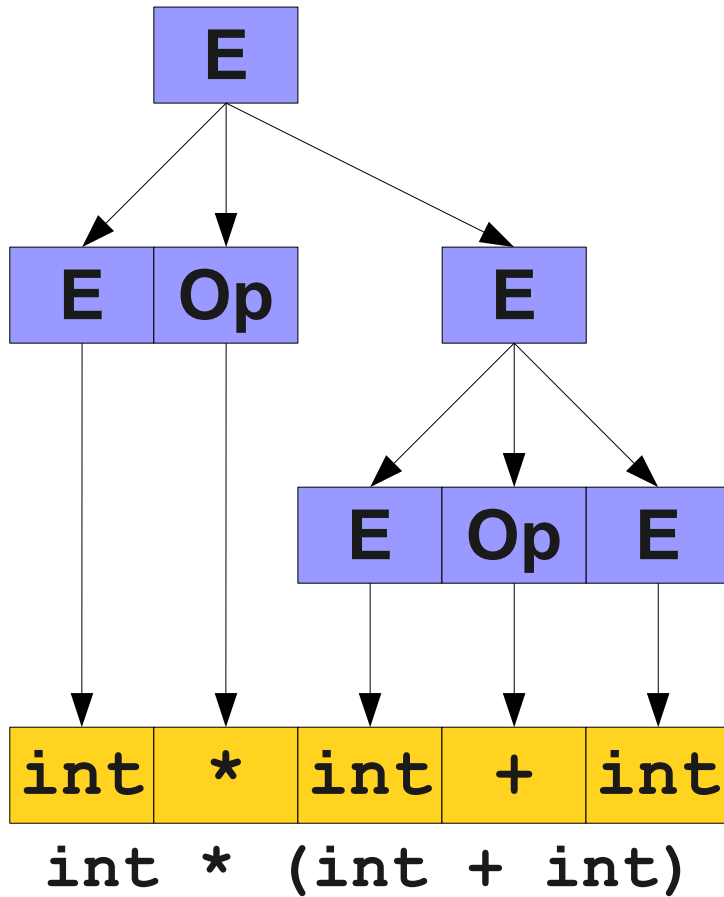
**E**  
⇒ **E Op E**  
⇒ **E Op (E)**  
⇒ **E Op (E Op E)**  
⇒ **E Op (E Op int)**  
⇒ **E Op (E + int)**  
⇒ **E Op (int + int)**  
⇒ **E \* (int + int)**  
⇒ **int \* (int + int)**



# Parse Trees

- A **parse tree** is a tree encoding the steps in a derivation.
- Internal nodes represent nonterminal symbols used in the production.
- Walking the leaves in order gives the generated string.
- Encodes what productions are used, not the order in which those productions are applied.

# A Serious Problem



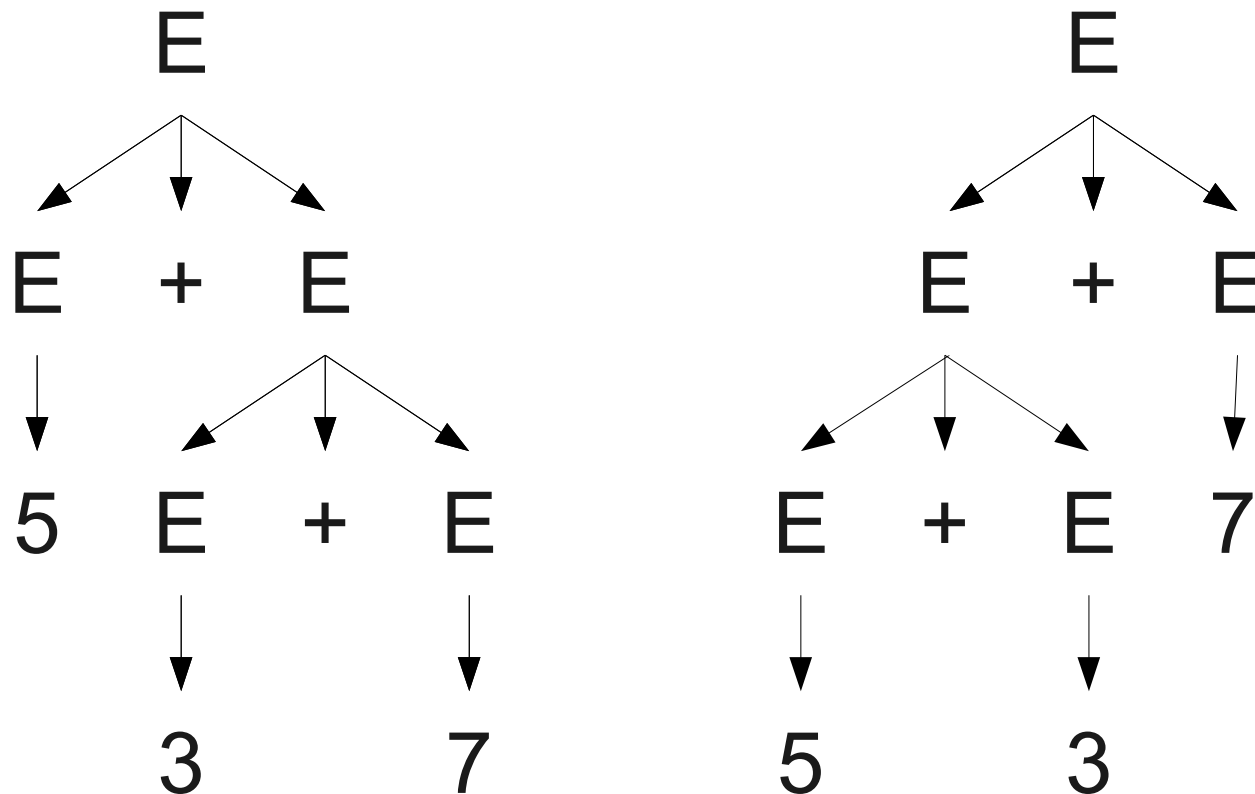
# Ambiguity

- A CFG is said to be **ambiguous** if there is at least one string with two or more parse trees.
- There is no algorithm for converting an arbitrary ambiguous grammar into an unambiguous one.
  - Some languages are **inherently ambiguous**, meaning that no unambiguous grammar exists for them.
- There is no algorithm for detecting whether an arbitrary grammar is ambiguous.

# Is Ambiguity a Problem?

- Depends on **semantics**.

**E**  $\rightarrow$  **int** | **E + E**

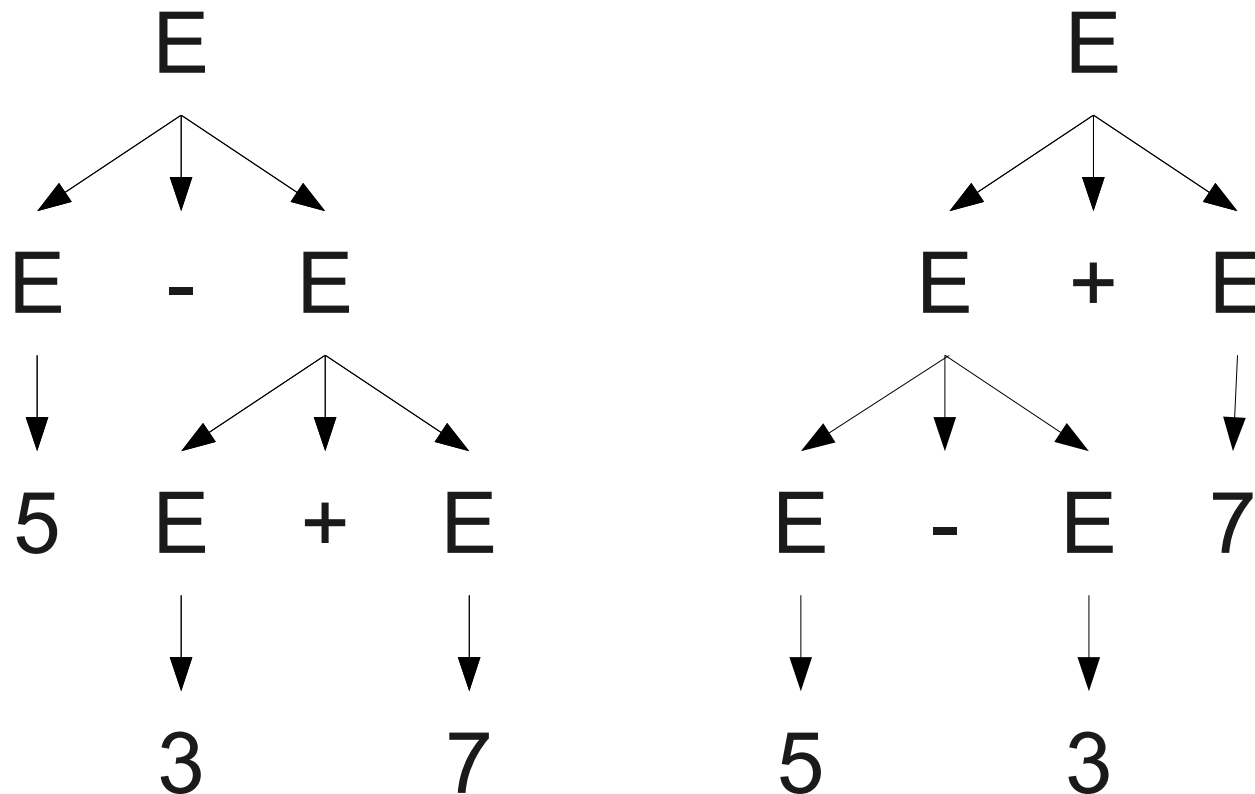




# Is Ambiguity a Problem?

- Depends on **semantics**.

$E \rightarrow \text{int} \mid E + E \mid E - E$



# Resolving Ambiguity

- If a grammar can be made unambiguous at all, it is usually made unambiguous through **layering**.
- Have exactly one way to build each piece of the string.
- Have exactly one way of combining those pieces back together.

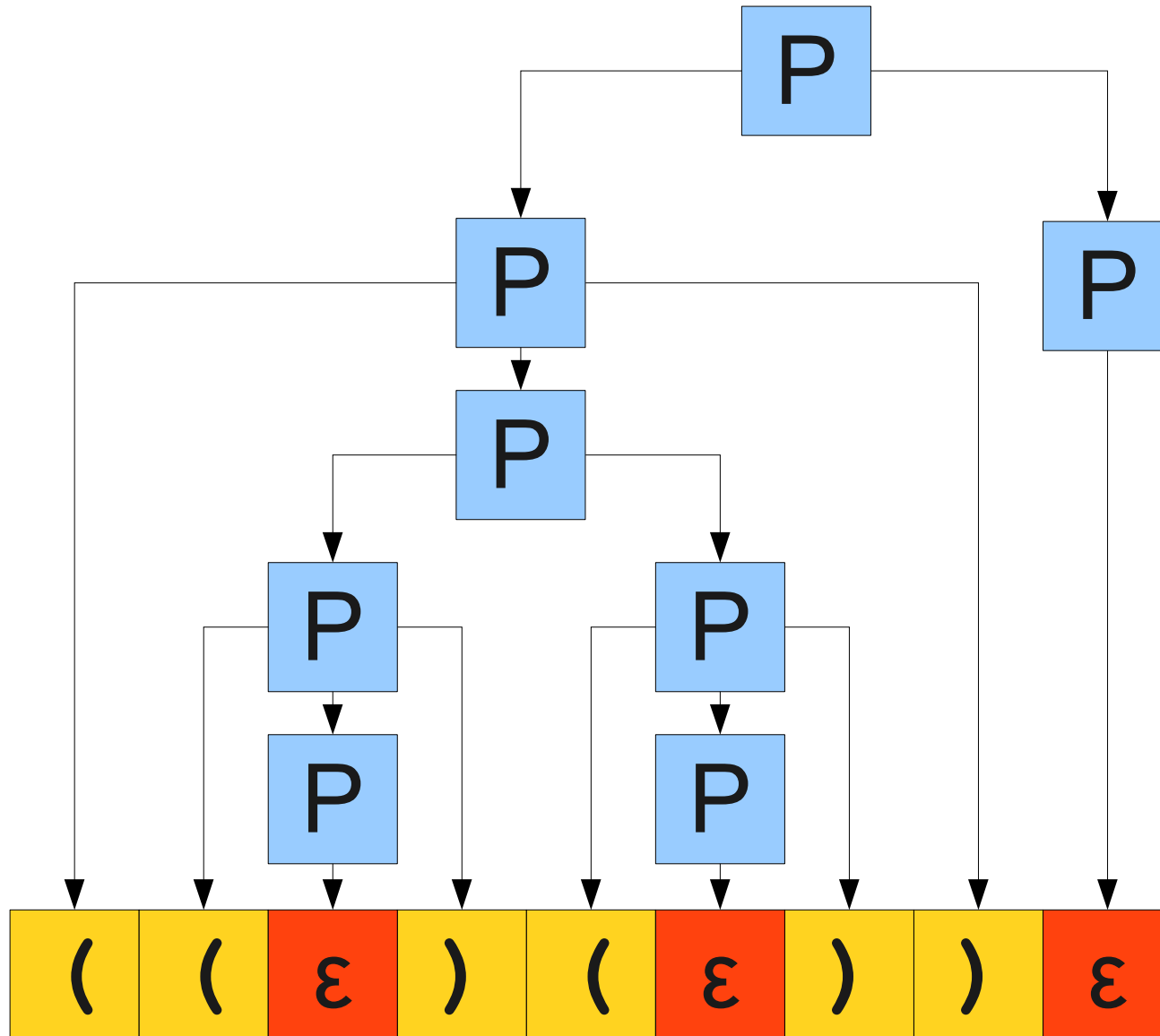
# Example: Balanced Parentheses

- Consider the language of all strings of balanced parentheses.
- Examples:
  - $\epsilon$
  - $()$
  - $((())())$
  - $(((())))((())())$
- Here is one possible grammar for balanced parentheses:

$$\mathbf{P} \rightarrow \epsilon \mid \mathbf{PP} \mid (\mathbf{P})$$



# Balanced Parentheses



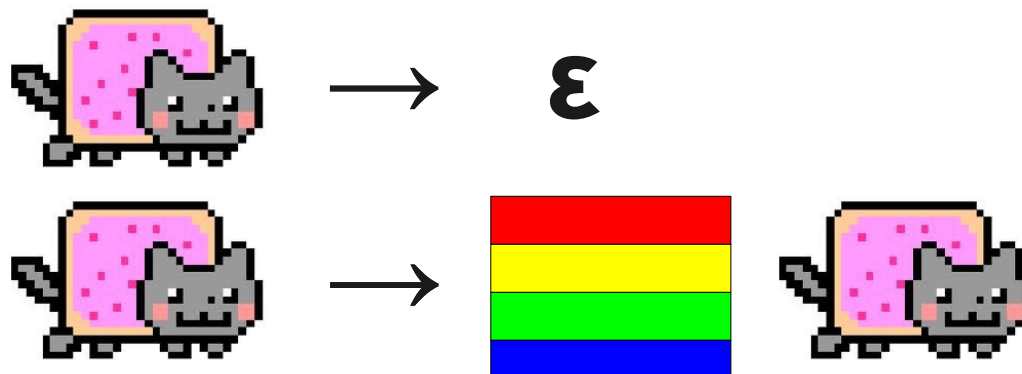
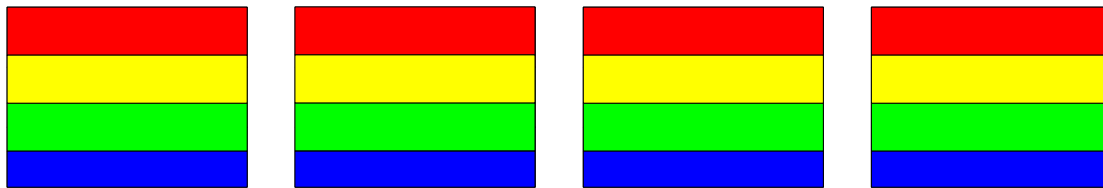
How to resolve this ambiguity?

# Rethinking Parentheses

- A string of balanced parentheses is a sequence of strings that are themselves balanced parentheses.
- To avoid ambiguity, we can build the string in two steps:
  - Decide how many different substrings we will glue together.
  - Build each substring independently.

# Um... what?

- The way the nyancat flies across the sky is similar to how we can build up strings of balanced parentheses.





# Building Parentheses

- Spread a string of parentheses across the string. There is exactly one way to do this for any number of parentheses.
- Expand out each substring by adding in parentheses and repeating.

**S**  $\rightarrow$  **P S** |  $\epsilon$

**P**  $\rightarrow$  ( **S** )

# Building Parentheses

**S**  $\rightarrow$  **P S** |  $\epsilon$

**P**  $\rightarrow$  ( **S** )

**S**  
 $\Rightarrow$  **PS**  
 $\Rightarrow$  **PPS**  
 $\Rightarrow$  **PP**  
 $\Rightarrow$  ( **S** ) **P**  
 $\Rightarrow$  ( **S** ) ( **S** )  
 $\Rightarrow$  ( **PS** ) ( **S** )  
 $\Rightarrow$  ( **P** ) ( **S** )  
 $\Rightarrow$  ( ( **S** ) ) ( **S** )  
 $\Rightarrow$  ( ( ) ) ( **S** )  
 $\Rightarrow$  ( ( ) ) ( )

# Context-Free Grammars

- A regular expression can be
  - Any letter
  - $\epsilon$
  - $\emptyset$
  - The concatenation of regular expressions.
  - The union of regular expressions.
  - The Kleene closure of a regular expression.
  - A parenthesized regular expression.

# Context-Free Grammars

- This gives us the following CFG:

$$\mathbf{R} \rightarrow \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots$$

$$\mathbf{R} \rightarrow \text{“}\varepsilon\text{”}$$

$$\mathbf{R} \rightarrow \emptyset$$

$$\mathbf{R} \rightarrow \mathbf{RR}$$

$$\mathbf{R} \rightarrow \mathbf{R} \text{ “} \mid \text{” } \mathbf{R}$$

$$\mathbf{R} \rightarrow \mathbf{R}^*$$

$$\mathbf{R} \rightarrow (\mathbf{R})$$

# An Ambiguous Grammar

$R \rightarrow a \mid b \mid c \mid \dots$

$R \rightarrow \epsilon$

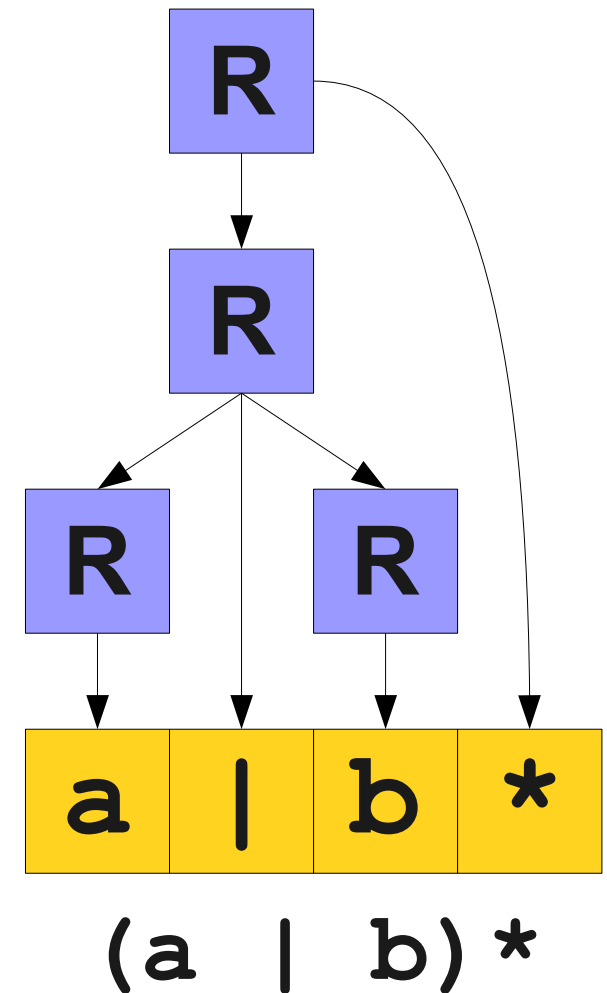
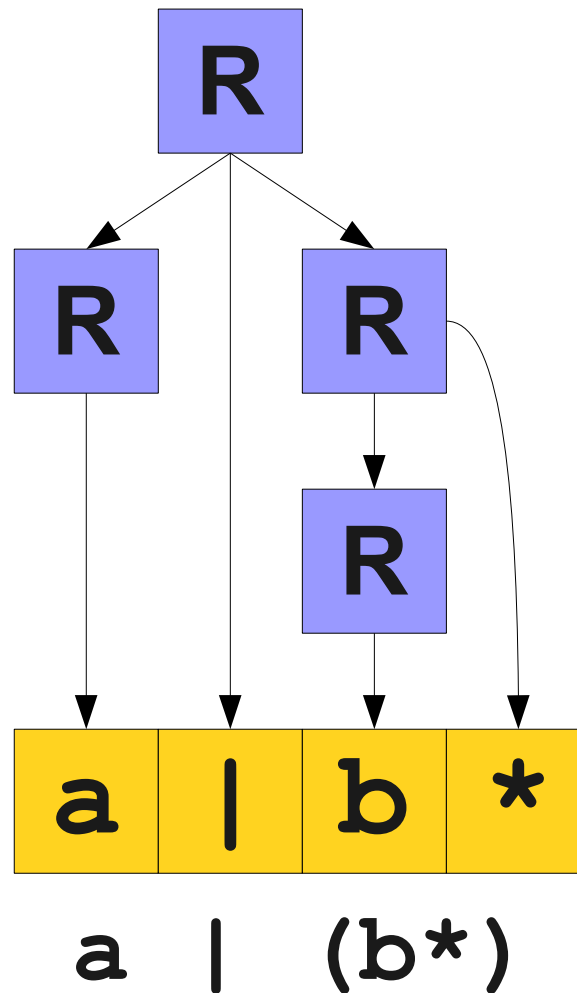
$R \rightarrow \emptyset$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



# Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

**R** → **a** | **b** | **c** | ...

**R** → “**ε**”

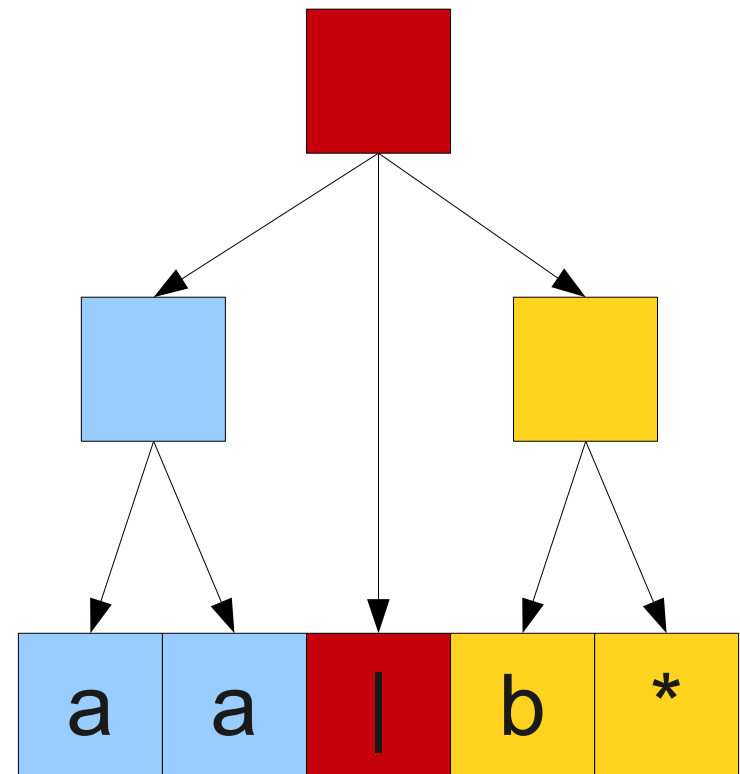
**R** → **∅**

**R** → **RR**

**R** → **R** “|” **R**

**R** → **R\***

**R** → **(R)**



# Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

$$\mathbf{R} \rightarrow \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots$$

$$\mathbf{R} \rightarrow \text{“}\varepsilon\text{”}$$

$$\mathbf{R} \rightarrow \emptyset$$

$$\mathbf{R} \rightarrow \mathbf{RR}$$

$$\mathbf{R} \rightarrow \mathbf{R} \text{ “} \mid \text{” } \mathbf{R}$$

$$\mathbf{R} \rightarrow \mathbf{R}^*$$

$$\mathbf{R} \rightarrow (\mathbf{R})$$

$$\mathbf{R} \rightarrow \mathbf{S} \mid \mathbf{R} \text{ “} \mid \text{” } \mathbf{S}$$

$$\mathbf{S} \rightarrow \mathbf{T} \mid \mathbf{ST}$$

$$\mathbf{T} \rightarrow \mathbf{U} \mid \mathbf{T}^*$$

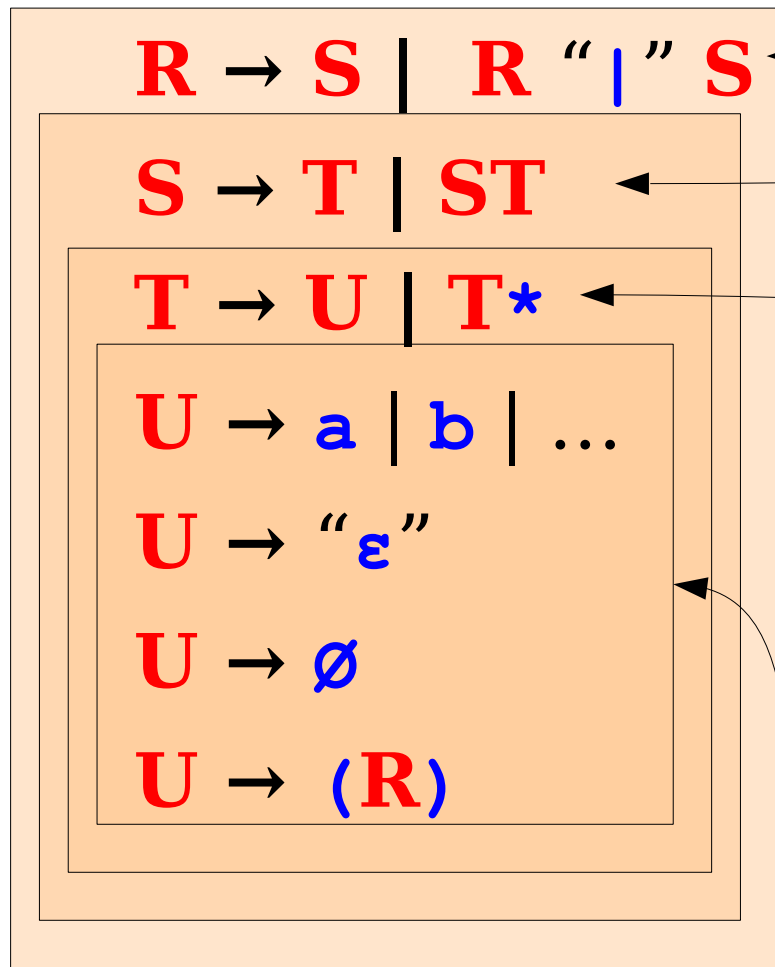
$$\mathbf{U} \rightarrow \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots$$

$$\mathbf{U} \rightarrow \text{“}\varepsilon\text{”}$$

$$\mathbf{U} \rightarrow \emptyset$$

$$\mathbf{U} \rightarrow (\mathbf{R})$$

# Why is this unambiguous?



Unions  
concatenated  
expressions

Concatenates starred  
expressions

Puts stars onto  
atomic expressions

Only generates  
"atomic" expressions



$R \rightarrow S \mid R \mid S$

$S \rightarrow T \mid ST$

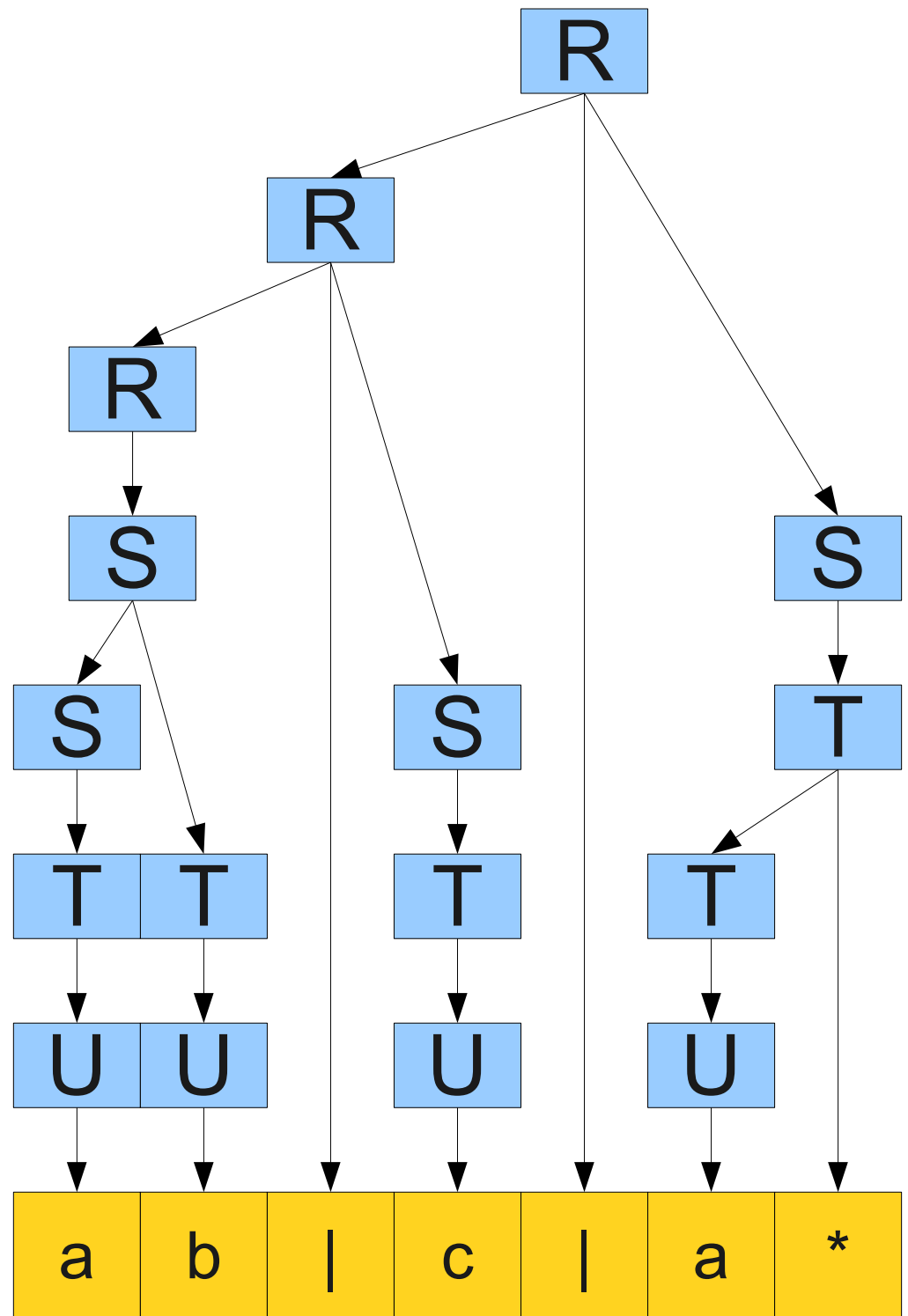
$T \rightarrow U \mid T^*$

$U \rightarrow a \mid b \mid c \mid \dots$

$U \rightarrow \epsilon$

$U \rightarrow \emptyset$

$U \rightarrow (R)$



# Summary

- **Context-free grammars** give a way to describe a class of formal languages (the **context-free languages**) that are strictly larger than the regular languages.
- A **parse tree** shows how a string can be **derived** from a grammar.
- A grammar is **ambiguous** if it can derive the same string multiple ways.
- There is no algorithm for eliminating ambiguity; it must be done by hand.

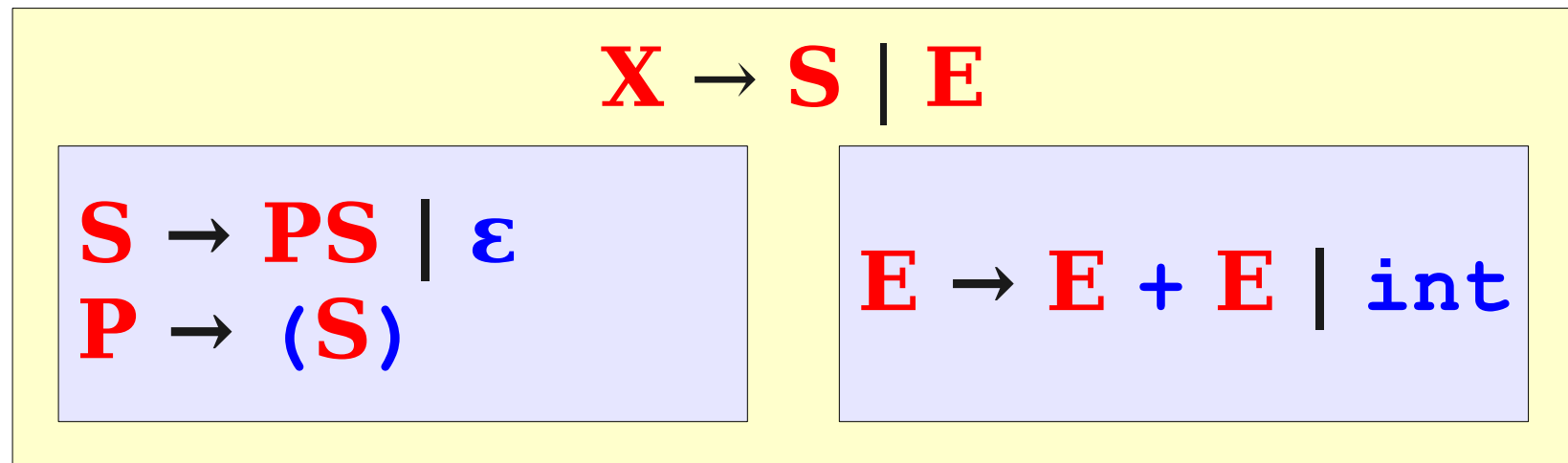
# Closure Properties of Context-Free Languages

# Closure Properties

- If  $L_1$  and  $L_2$  are regular, then
  - $\bar{L}_1$  is regular.
  - $L_1 \cup L_2$  is regular.
  - $L_1 \cap L_2$  is regular.
  - $L_1 L_2$  is regular.
  - $L_1^*$  is regular.
  - $h^*(L_1)$  is regular.
- How many of these properties still hold for context-free languages?

# The Union of CFLs

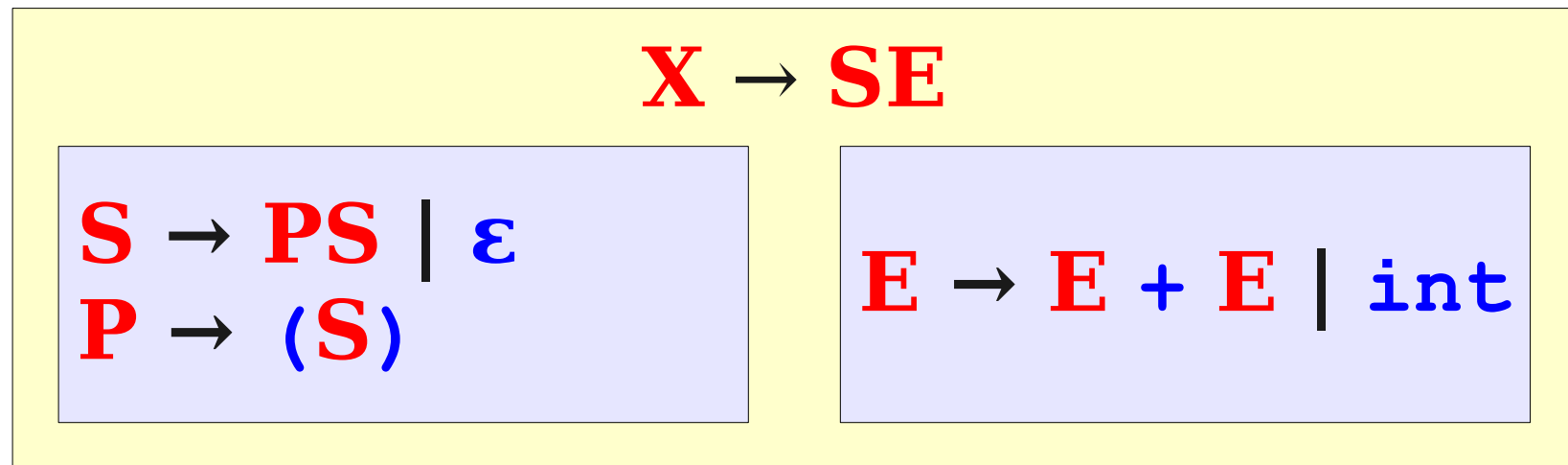
- Suppose that  $L_1$  and  $L_2$  are *context-free* languages.
- Is  $L_1 \cup L_2$  a context-free language?



- **Yes!** Use the above construction.
  - Rename nonterminals in the two grammars if necessary.

# The Concatenation of CFLs

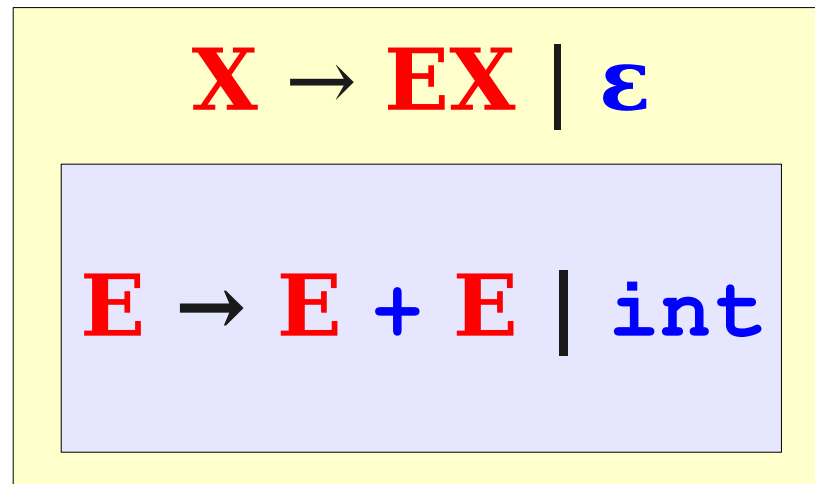
- Suppose that  $L_1$  and  $L_2$  are *context-free* languages.
- Is  $L_1 L_2$  a context-free language?



- **Yes!** Use the above construction.
  - Rename nonterminals in the two grammars if necessary.

# The Kleene Closure of CFLs

- Suppose that  $L$  is a *context-free* language.
- Is  $L^*$  a context-free language?



- **Yes!** Use the above construction.

# Closure Properties of CFLs

- If  $L_1$  and  $L_2$  are context-free languages, then
  - $L_1 \cup L_2$  is context-free.
  - $L_1 L_2$  is context-free.
  - $L_1^*$  is context-free.
  - $h^*(L_1)$  is context-free.
- Do the other properties still hold?
- We'll see early next week...



# Next Time

- **Pushdown Automata**
  - Automata for recognizing CFLs.
  - A beautiful generalization of DFAs and NFAs.
  - An easy proof that any regular language is context-free.