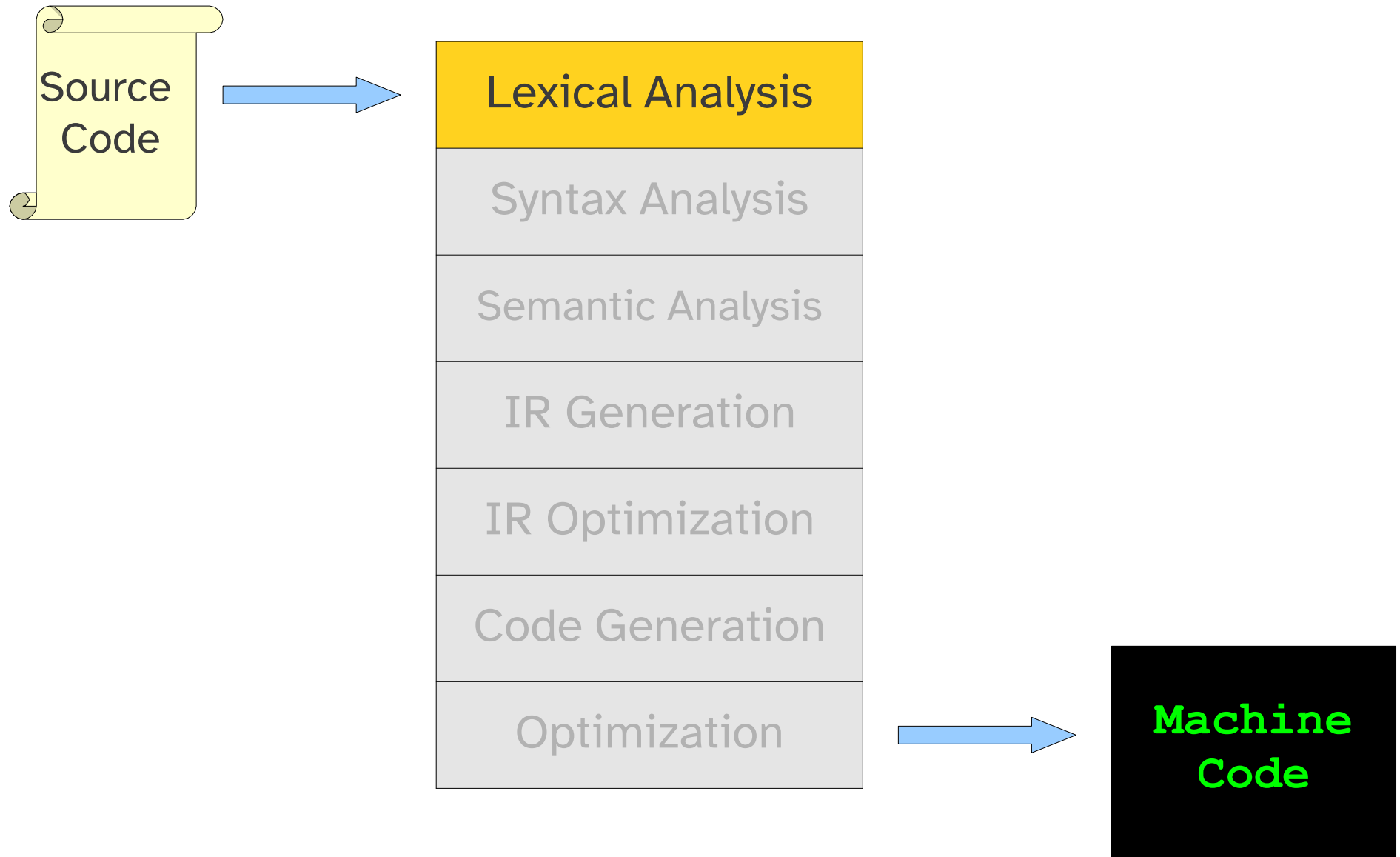
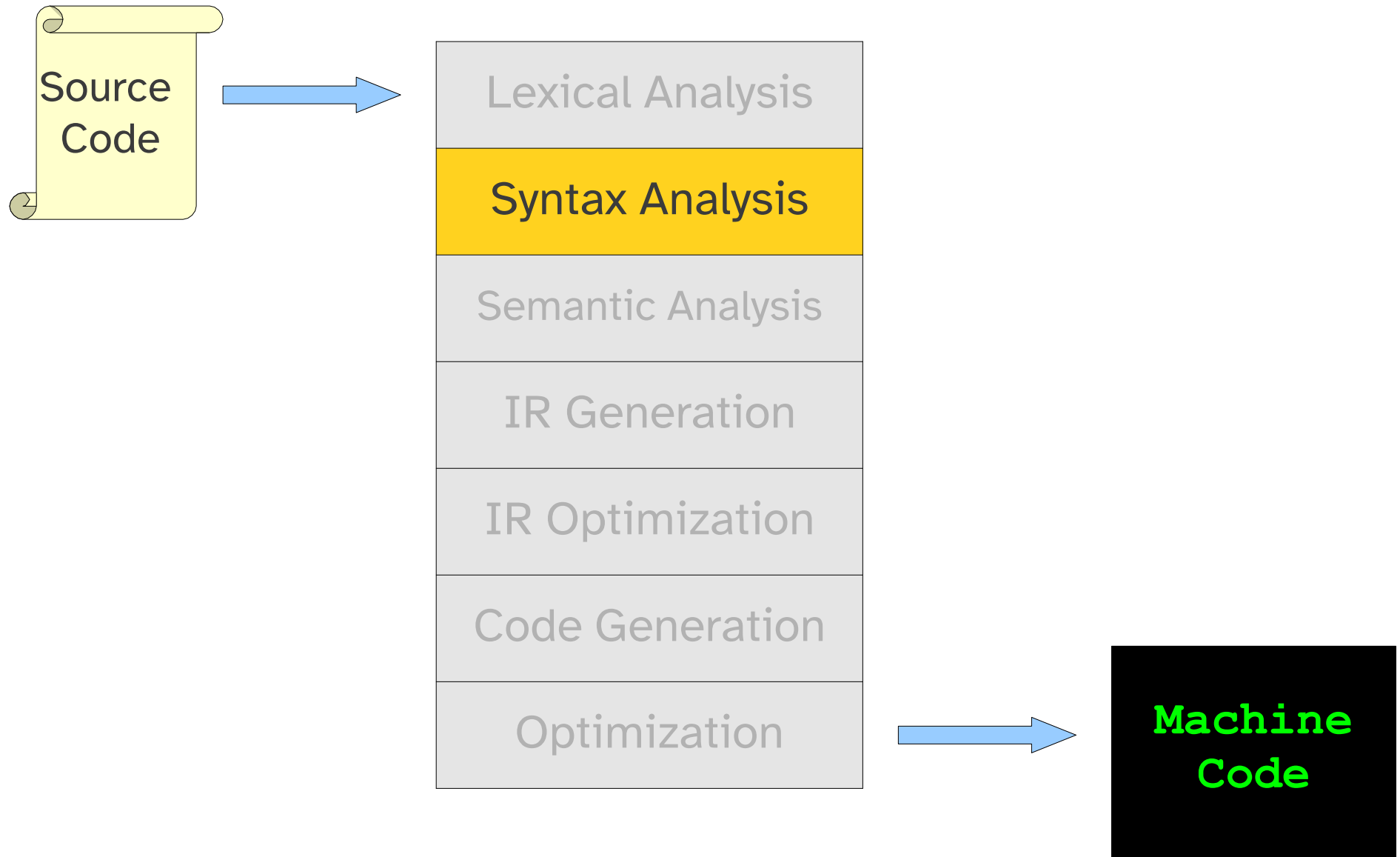


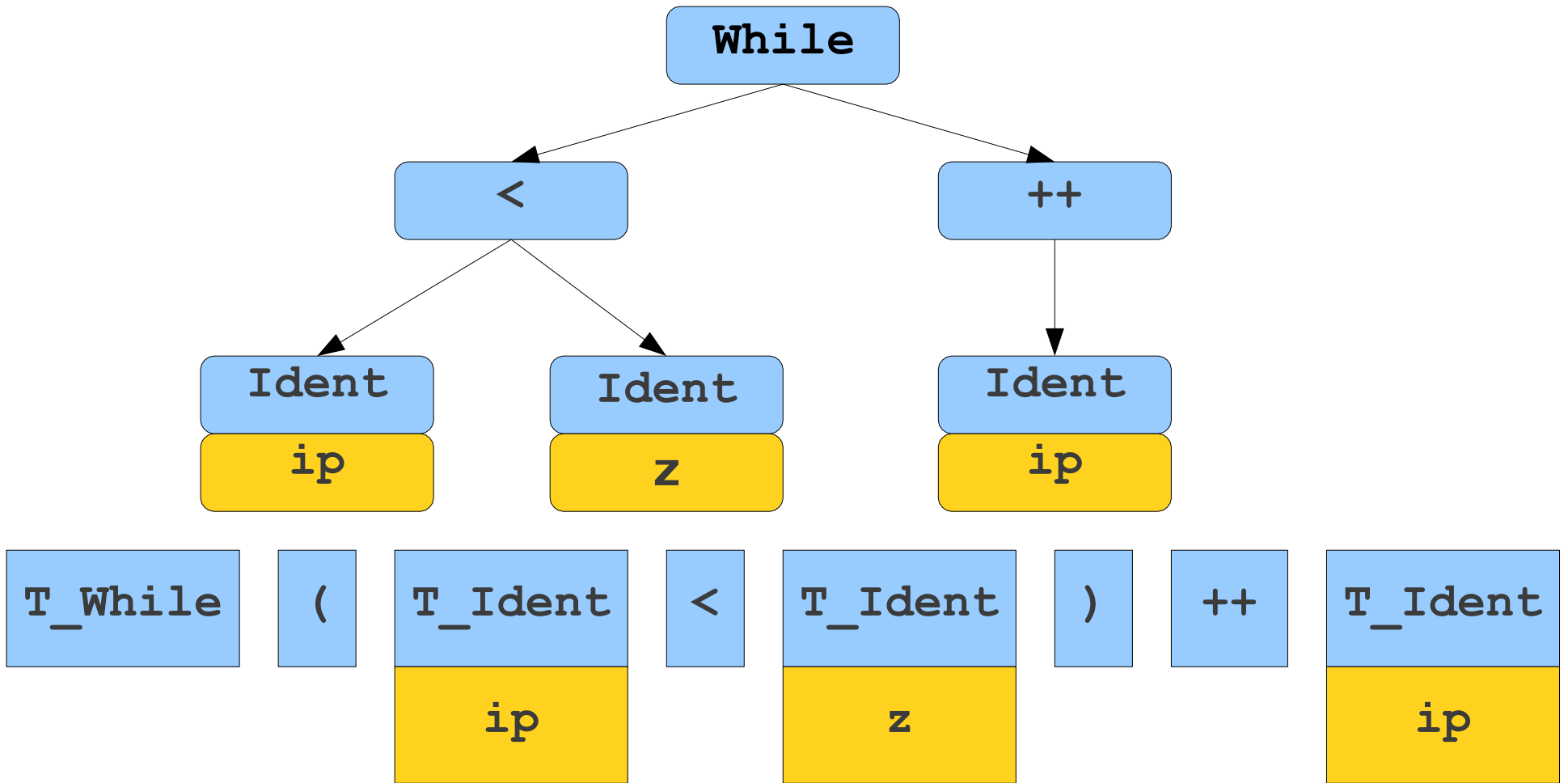
Syntax Analysis

Where We Were



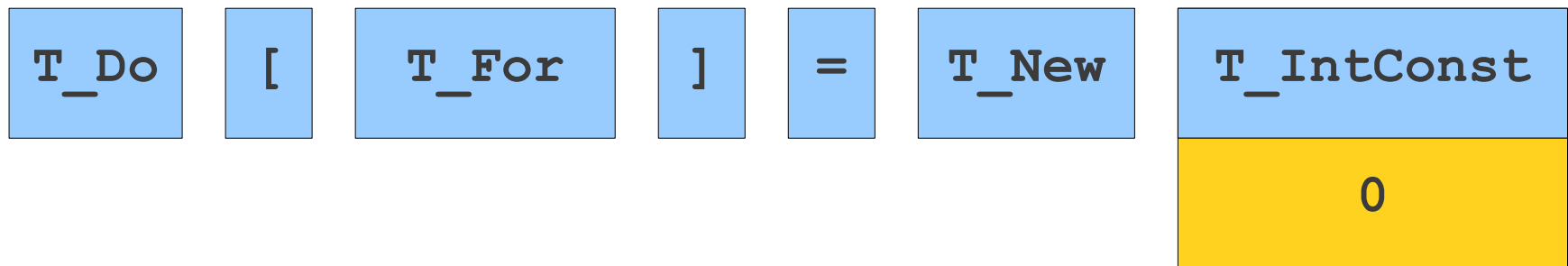
Where We Are





w	h	i	l	e		(i	p		<		z)	\n	\t	+	+	i	p	;
---	---	---	---	---	--	---	---	---	--	---	--	---	---	----	----	---	---	---	---	---

```
while (ip < z)
    ++ip;
```



d	o	[f	o	r]		=		n	e	w		0	;
---	---	---	---	---	---	---	--	---	--	---	---	---	--	---	---

`do[for] = new 0;`

Outline

- Today: Formalisms for syntax analysis
 - Context-Free Grammars
 - Derivations
 - Concrete and Abstract Syntax Trees
 - Ambiguity
- Later: Parsing algorithms
 - Top-Down Parsing
 - Bottom-Up Parsing

Context-Free Grammars

- A **context-free grammar** (or **CFG**) is a formalism for defining languages.
- Can define the **context-free languages**, a strict superset of the regular 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 \rightarrow **int**

E \rightarrow **E Op E**

E \rightarrow **(E)**

Op \rightarrow **+**

Op \rightarrow **-**

Op \rightarrow *****

Op \rightarrow **/**

Context-Free Grammars (CFGs)

A CFG is a collection of four objects:

- A set of **nonterminal symbols** (or **variables**),
- A set of **terminal symbols**,
- A set of **production rules** saying how each nonterminal can be converted by a string of terminals and nonterminals, and
- A **start symbol** that begins the derivation.

E \rightarrow **int**

E \rightarrow **E Op E**

E \rightarrow **(E)**

Op \rightarrow **+**

Op \rightarrow **-**

Op \rightarrow *****

Op \rightarrow **/**

CFGs for Programming Languages

BLOCK	→	STMT
		{ STMTS }
STMTS	→	ϵ
		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 at the beginning of the alphabet will represent nonterminals.
 - i.e., **A**, **B**, **C**, **D**
- Lowercase letters at the end of the alphabet will represent terminals.
 - i.e., **t**, **u**, **v**, **w**
- Lowercase Greek letters will represent arbitrary strings of terminals and nonterminals.
 - i.e., α , γ , ω

Examples

- We might write an arbitrary production as

$$\mathbf{A} \rightarrow \omega$$

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

$$\mathbf{A}t$$

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

$$\mathbf{B} \rightarrow \alpha \mathbf{A}t\omega$$

Derivations

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

- 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 α yields β , we write $\alpha \Rightarrow \beta$.
- We say that α **derives** β iff there is a sequence of strings where
$$\alpha \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \beta$$
- If α derives β , we write $\alpha \Rightarrow^* \beta$.

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 when we talk about parsing later.

Derivations Examples

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

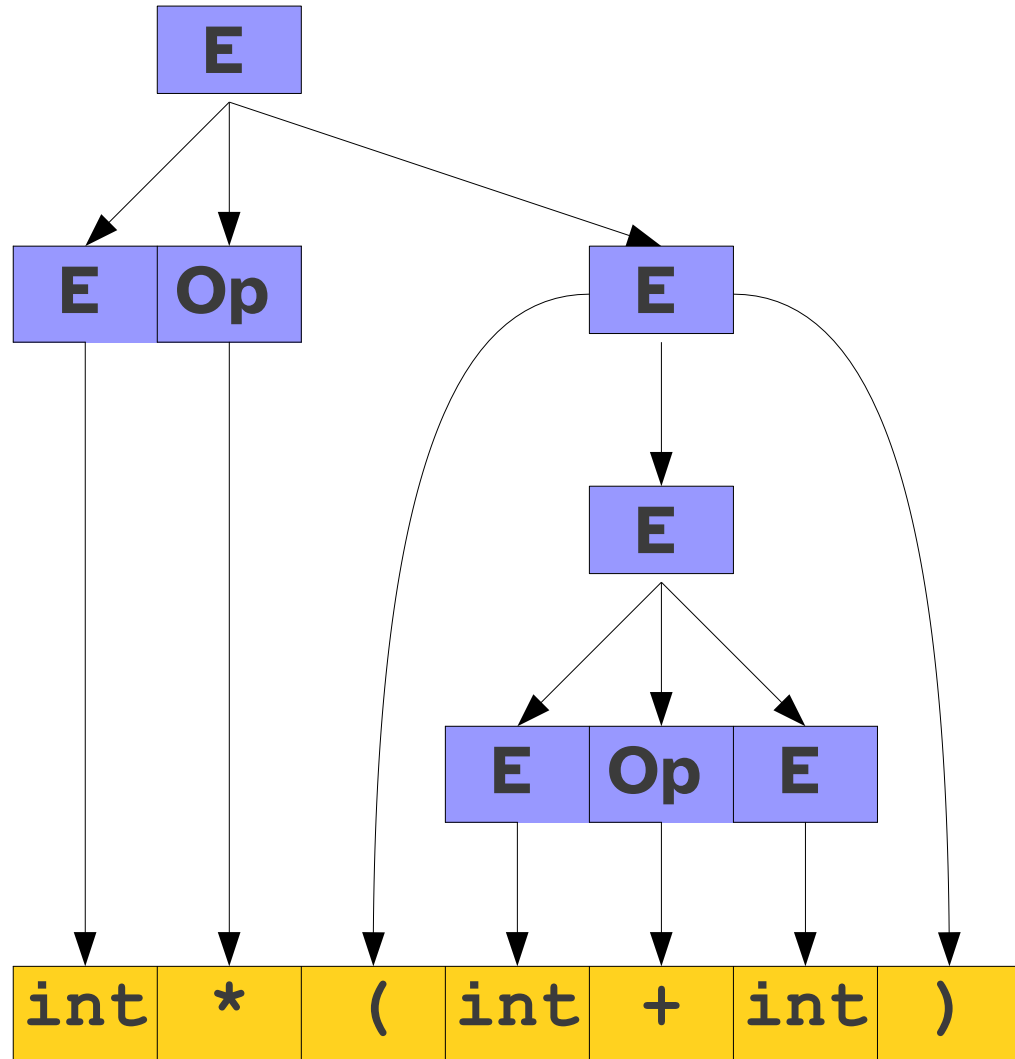
- A derivation encodes two pieces of information:
 - What productions were applied produce the resulting string from the start symbol?
 - In what order were they applied?
- Multiple derivations might use the same productions, but apply them in a different order.

Parse Trees

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

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)

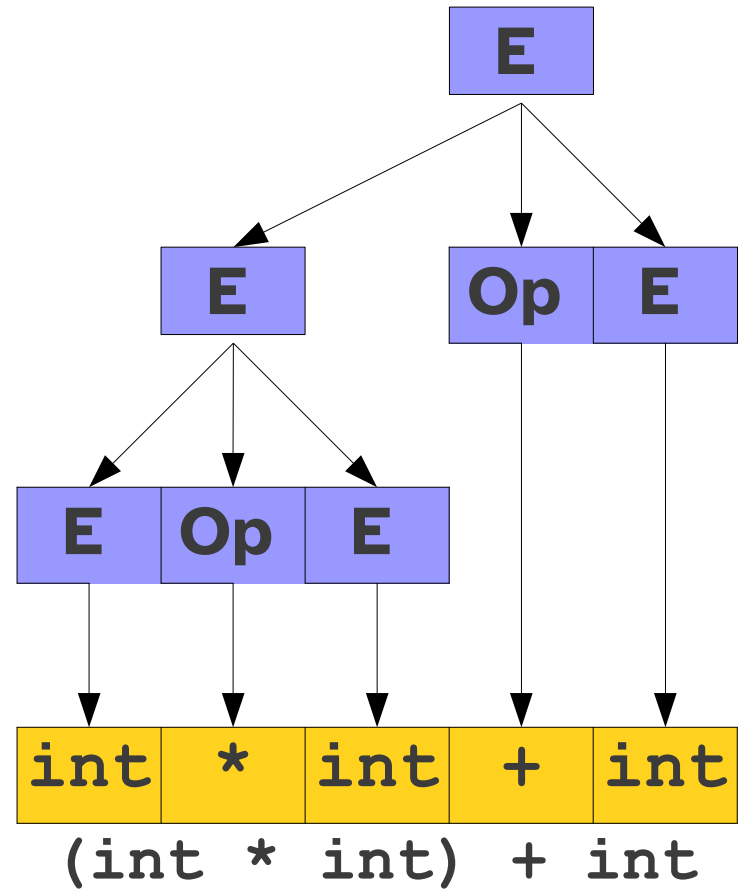
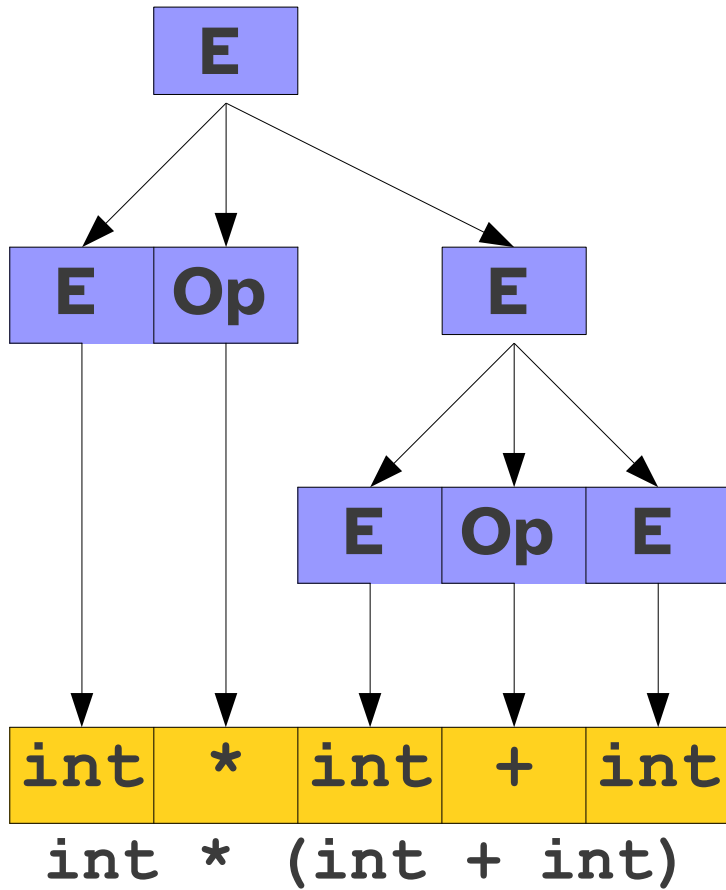


The Goal of Parsing

- Goal of syntax analysis: Recover the **structure** described by a series of tokens.
- If language is described as a CFG, goal is to recover a parse tree for the the input string.

Challenges in Parsing

A Serious Problem



Ambiguity

- A CFG is said to be **ambiguous** if there is at least one string with two or more parse trees.
 - Using the same leftmost/rightmost derivation, one string can have at least 2 ways to be derived.
- There is no algorithm for converting an arbitrary ambiguous grammar into an unambiguous one.
- There is no algorithm for detecting whether an arbitrary grammar is ambiguous.

Example: With string `int * (int + int)`

E \rightarrow `int` |
 E Op E |
 (E)
Op \rightarrow `+` | `-` |
 `*` | `/`

E

Example: With string `int * (int + int)`

E \rightarrow `int` |

E Op E |

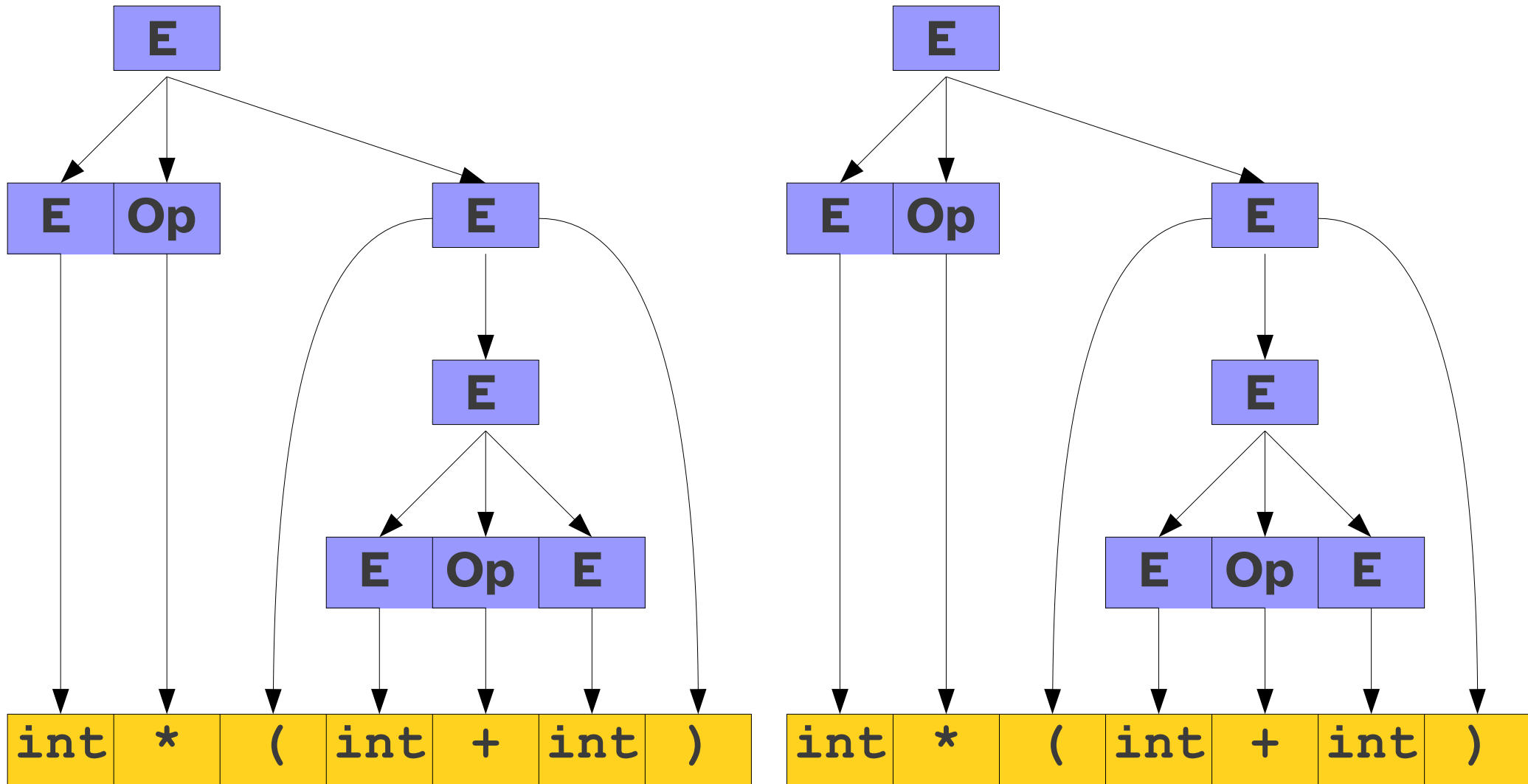
(E)

Op \rightarrow `+` | `-` |

`*` | `/`

E

This is ambiguous



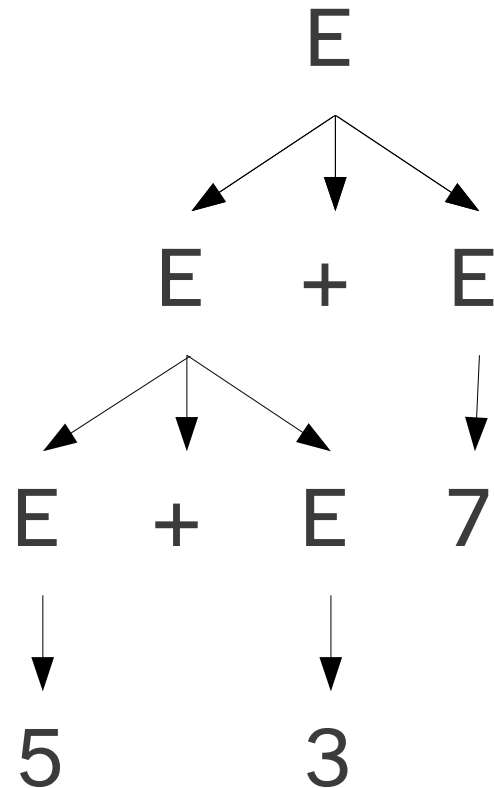
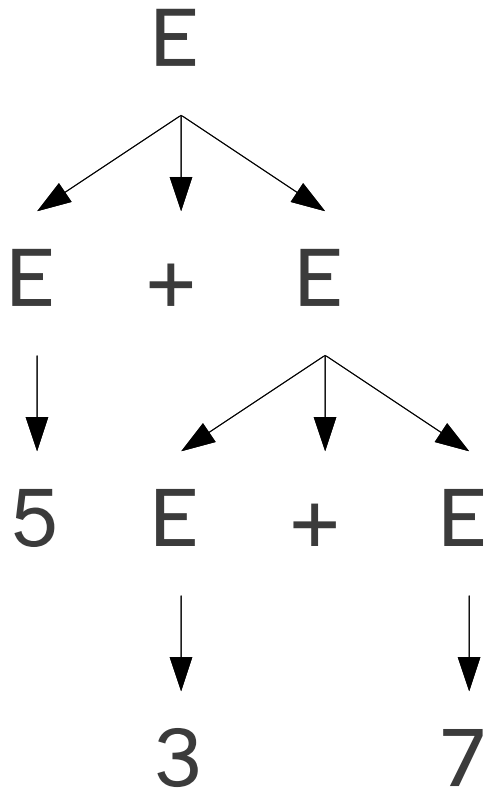
Because 1 strings created by 2 ways

Wrong! The parse trees are exactly the same!

Is Ambiguity a Problem?

- Depends on **semantics**.

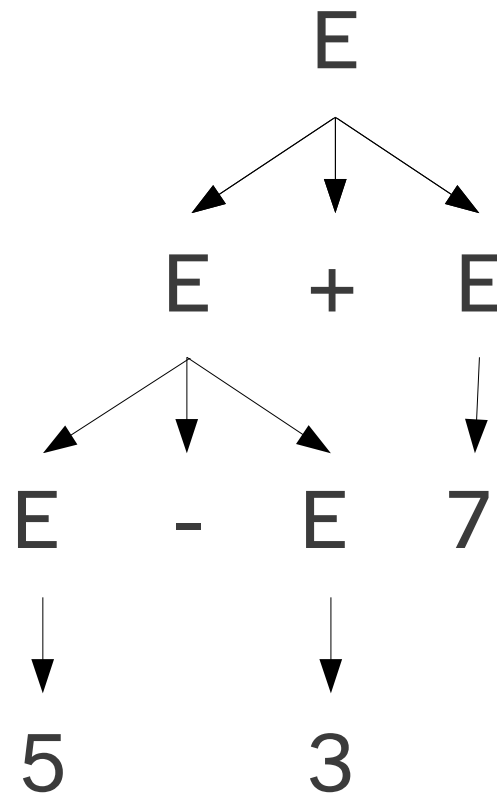
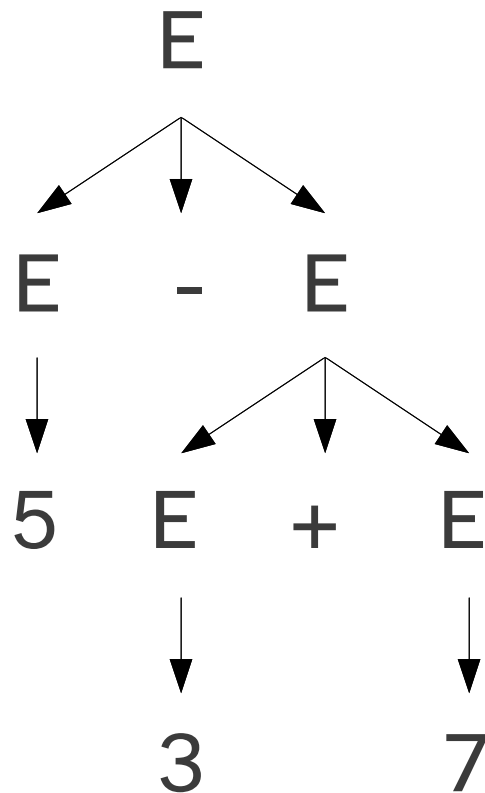
E \rightarrow **int** | **E + E**



Is Ambiguity a Problem?

- Depends on **semantics**.

E \rightarrow **int** | **E + E** | **E - E**



Drawbacks of ambiguous grammars

- Ambiguous semantics
- Parsing complexity
- May affect other phases
- Solutions
 - Allow only non-ambiguous grammars
 - Transform grammar into non-ambiguous
 - Handle as part of parsing method
 - Using special form of “precedence”

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.

CFG for RegEx

- Recall: A regular expression can be
 - Any letter
 - ϵ
 - The concatenation of regular expressions.
 - The union of regular expressions.
 - The Kleene closure of a regular expression.
 - A parenthesized regular expression.

CFG for RegEx

- This gives us the following CFG:

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

$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

An Ambiguous Grammar

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

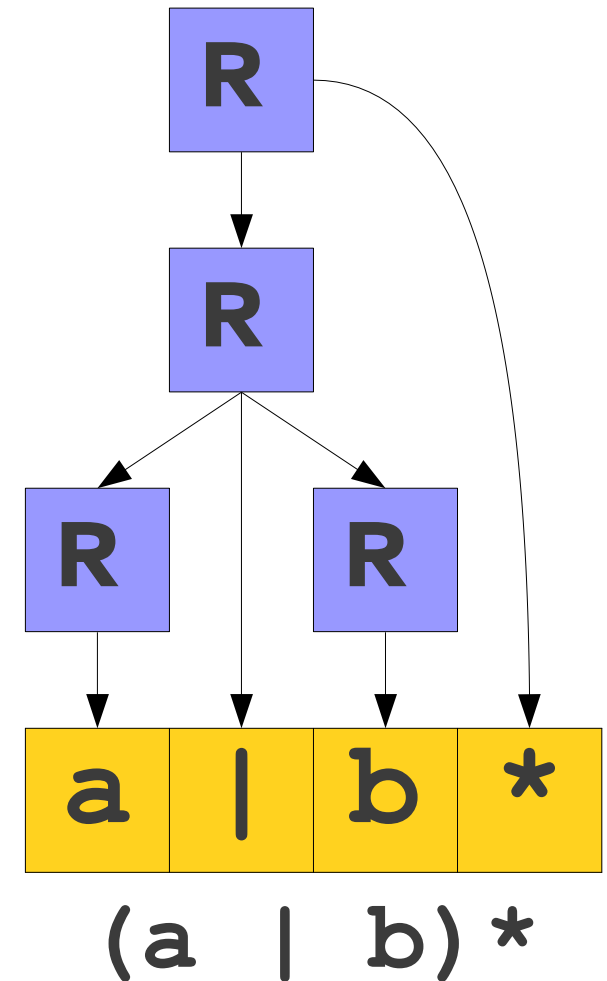
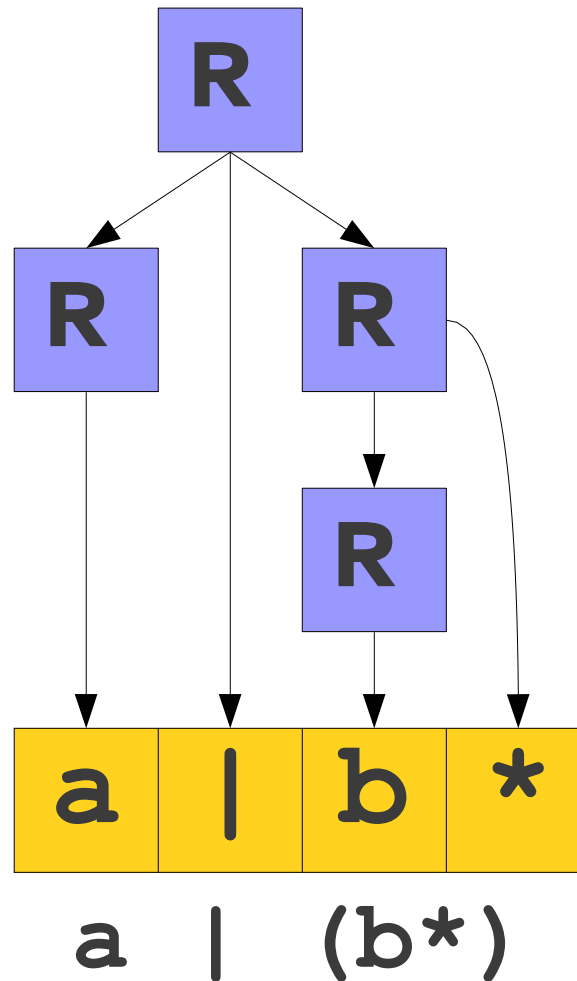
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

R \rightarrow **a** | **b** | **c** | ...

R \rightarrow “ **ϵ** ”

R \rightarrow **RR**

R \rightarrow **R** “**|**” **R**

R \rightarrow **R*******

R \rightarrow (**R**)

Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

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

$$R \rightarrow \epsilon$$

$$R \rightarrow RR$$

$$R \rightarrow R \mid R$$

$$R \rightarrow R^*$$

$$R \rightarrow (R)$$

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

$$S \rightarrow T \mid ST$$

$$T \rightarrow U \mid T^*$$

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

$$U \rightarrow \epsilon$$

$$U \rightarrow (R)$$

Why is this unambiguous?

$R \rightarrow S \mid R \text{ " | " } S$

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \text{"}\epsilon\text{"}$

$U \rightarrow (R)$

Infix Expression Example

Given grammar G:

E \rightarrow **E** "+" **E** |
E "-" **E** |
"-" **E** |
E "*" **E** |
E "/" **E** |
E "^" **E** |
"(" **E** ")" |
var

G is ambiguous

What about G1? (Wirth notation)

E \rightarrow **P** { **B** **P** }
P \rightarrow **var** | "(" **E** ")" | **U** **P**
B \rightarrow "+" | "-" | "*" | "/" | "^"
U \rightarrow "-"

Also, is $L(G1) = L(G)$?

Another solution: G2 grammar

E \rightarrow **T** { ("+" | "-") **T** }

T \rightarrow **F** { ("*" | "/") **F** }

F \rightarrow **P** ["^" **F**]

P \rightarrow **var** | "(" **E** ")" | "-" **T**

- Is G2 ambiguous?
- Is $L(G2)=L(G1)$?
- Is $L(G2)=L(G)$?

What can we tell about G?

E \rightarrow **T**

E \rightarrow **E** + **T**

T \rightarrow **F**

T \rightarrow **T** * **F**

F \rightarrow **var**

F \rightarrow (**E**)

- Unambiguous
- () has highest precedence
- * has precedence over +
- * and +: left associative: operations to be done from left to right
- e.g., $n+n+n$, $n+n*n$

If G is changed to G'

E \rightarrow **T**

E \rightarrow **T** + **E**

T \rightarrow **F**

T \rightarrow **T** * **F**

F \rightarrow **var**

F \rightarrow (**E**)

- Unambiguous
- * has precedence over +
- * is left associative
 - left recursive
- + is right associative
 - right recursive

Note:

- while $L(G)$ is equivalent to $L(G')$,
- G is **not** equivalent to G' : + is now right associative

The Structure of a Parse Tree

$R \rightarrow S \mid R \text{ " | " } S$

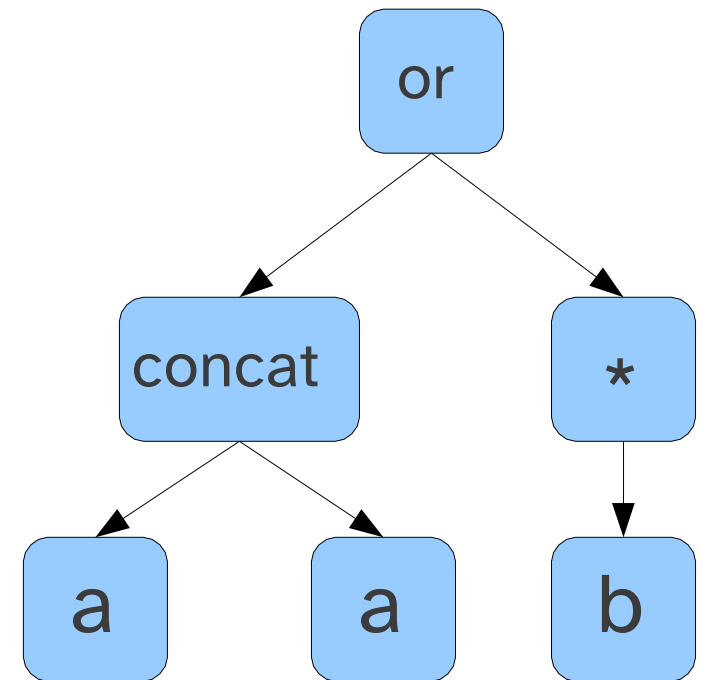
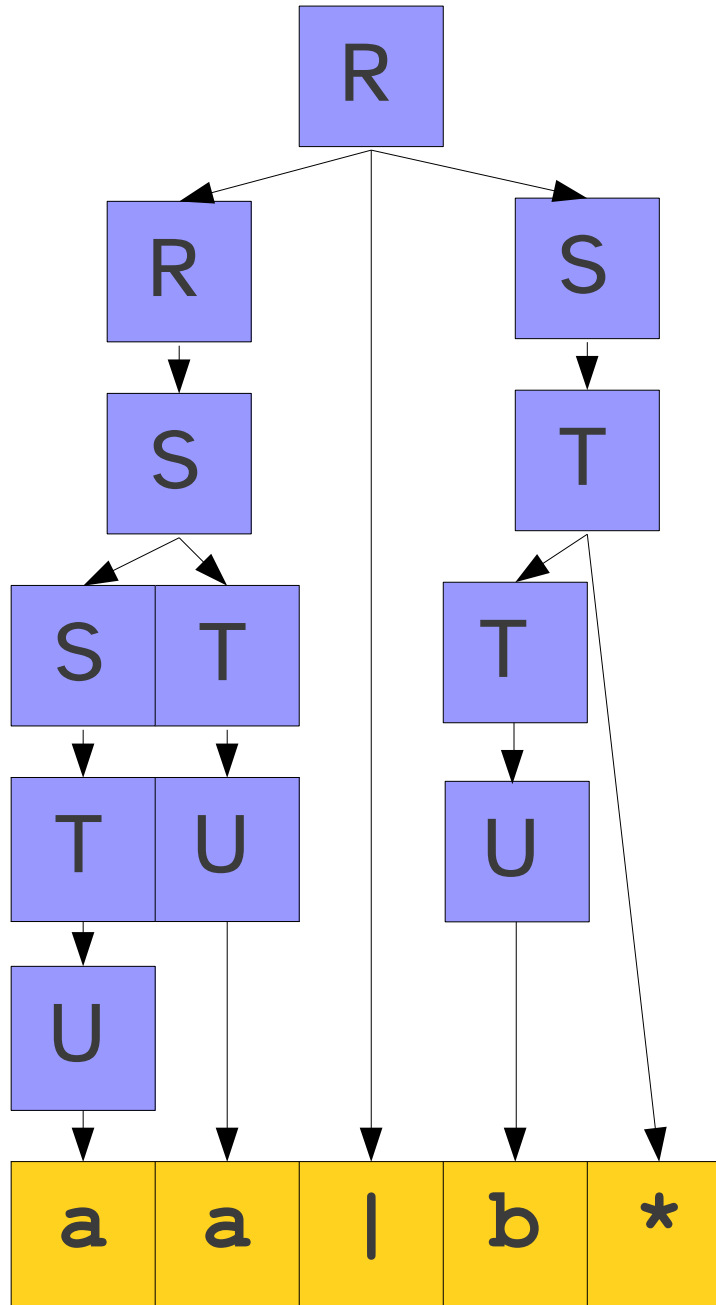
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \text{"}\epsilon\text{"}$

$U \rightarrow (R)$



$R \rightarrow S \mid R \text{ " | " } S$

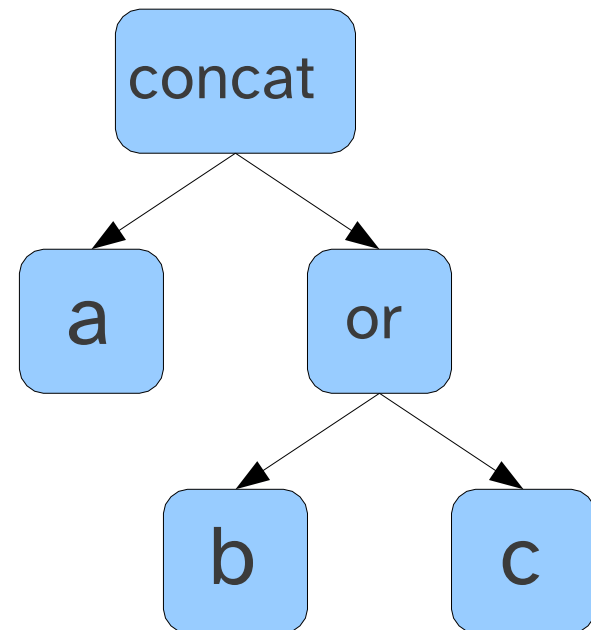
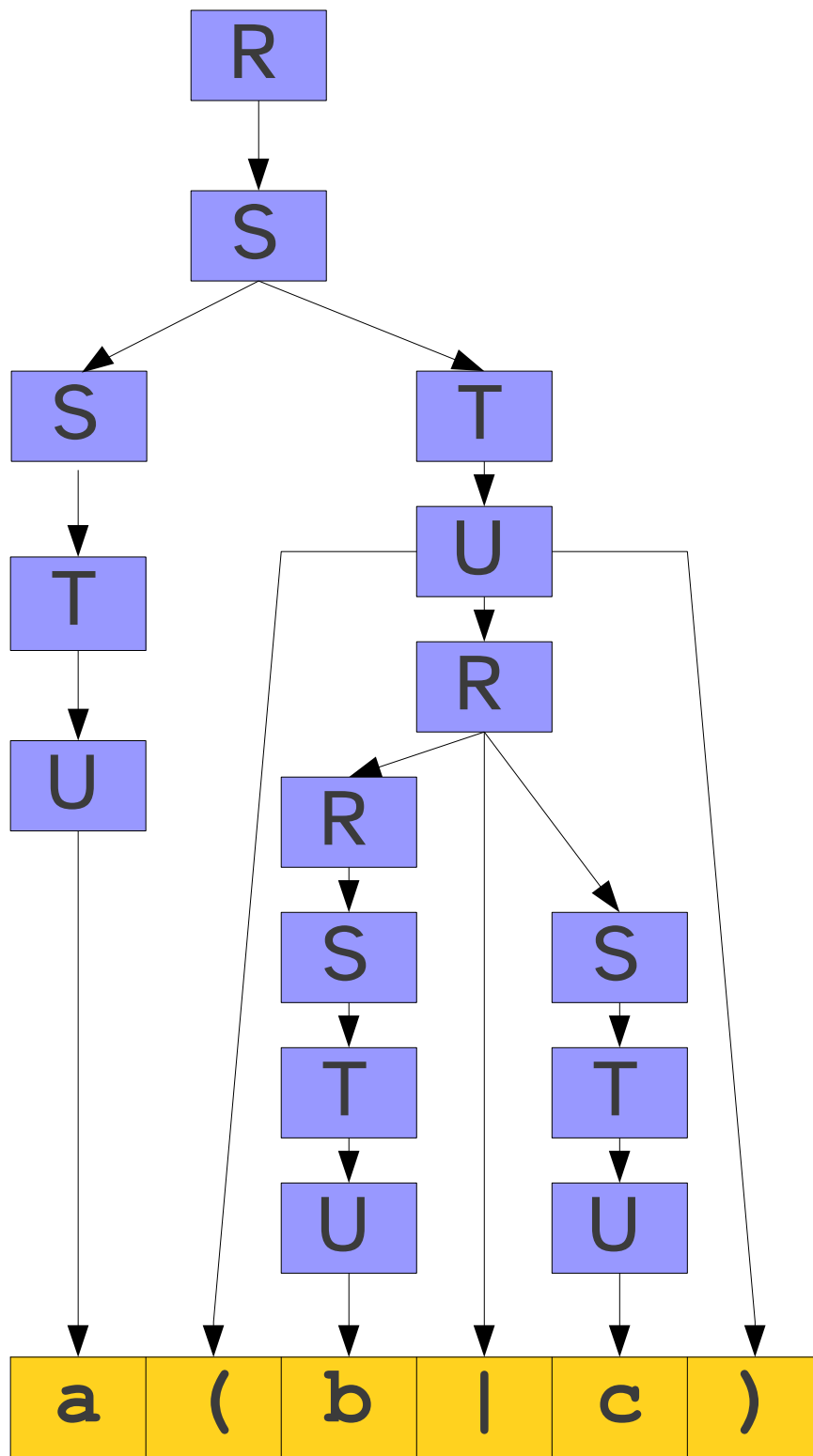
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \text{"}\epsilon\text{"}$

$U \rightarrow (R)$



Abstract Syntax Trees (ASTs)

- A parse tree is a **concrete syntax tree**; it shows exactly how the text was derived.
- A more useful structure is an **abstract syntax tree**, which retains only the essential structure of the input.

How to build an AST?

- Typically done through **semantic actions**.
- Associate a piece of code to execute with each production.
- As the input is parsed, execute this code to build the AST.

Semantic Actions to Build ASTs

$R \rightarrow S$	<code>R.ast = S.ast;</code>
$R \rightarrow R \mid S$	<code>R₁.ast = new Or(R₂.ast, S.ast);</code>
$S \rightarrow T$	<code>S.ast = T.ast;</code>
$S \rightarrow ST$	<code>S₁.ast = new Concat(S₂.ast, T.ast);</code>
$T \rightarrow U$	<code>T.ast = U.ast;</code>
$T \rightarrow T^*$	<code>T₁.ast = new Star(T₂.ast);</code>
$U \rightarrow a$	<code>U.ast = new SingleChar('a');</code>
$U \rightarrow \epsilon$	<code>U.ast = new Epsilon();</code>
$U \rightarrow (R)$	<code>U.ast = R.ast;</code>

Another example: Semantic Actions

E → **T** + **E**

$E_1.val = T.val + E_2.val$

E → **T**

$E.val = T.val$

T → **int**

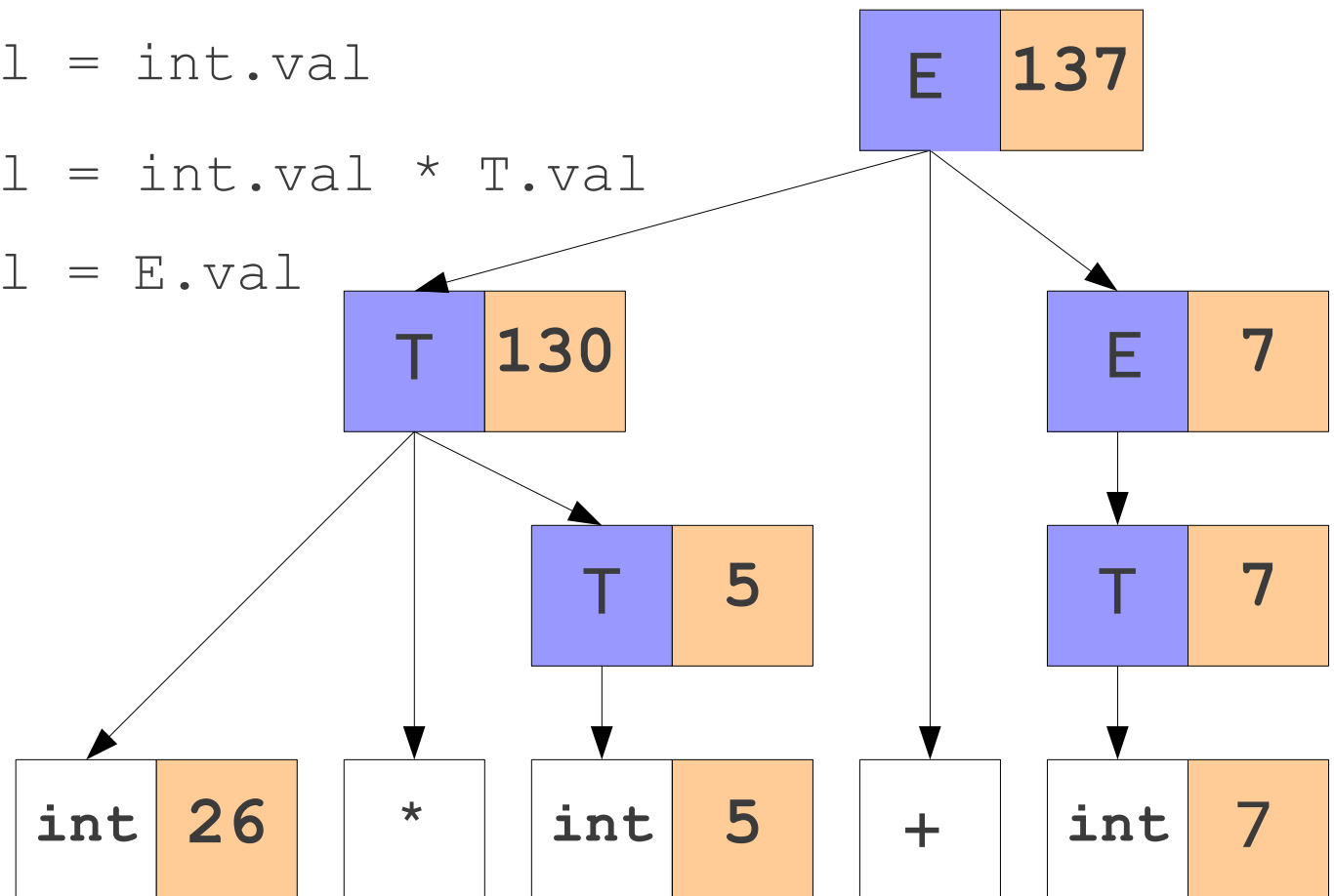
$T.val = int.val$

T → **int** * **T**

$T.val = int.val * T.val$

T → (**E**)

$T.val = E.val$



Summary

- Syntax analysis (**parsing**) extracts the structure from the tokens produced by the scanner.
- Languages are usually specified by **context-free grammars (CFGs)**.
- 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.
- **Abstract syntax trees (ASTs)** contain an abstract representation of a program's syntax.
- **Semantic actions** associated with productions can be used to build ASTs.

Next Time

- **Top-Down Parsing**
 - Parsing as a Search
 - Backtracking Parsers
 - Predictive Parsers LL(1)