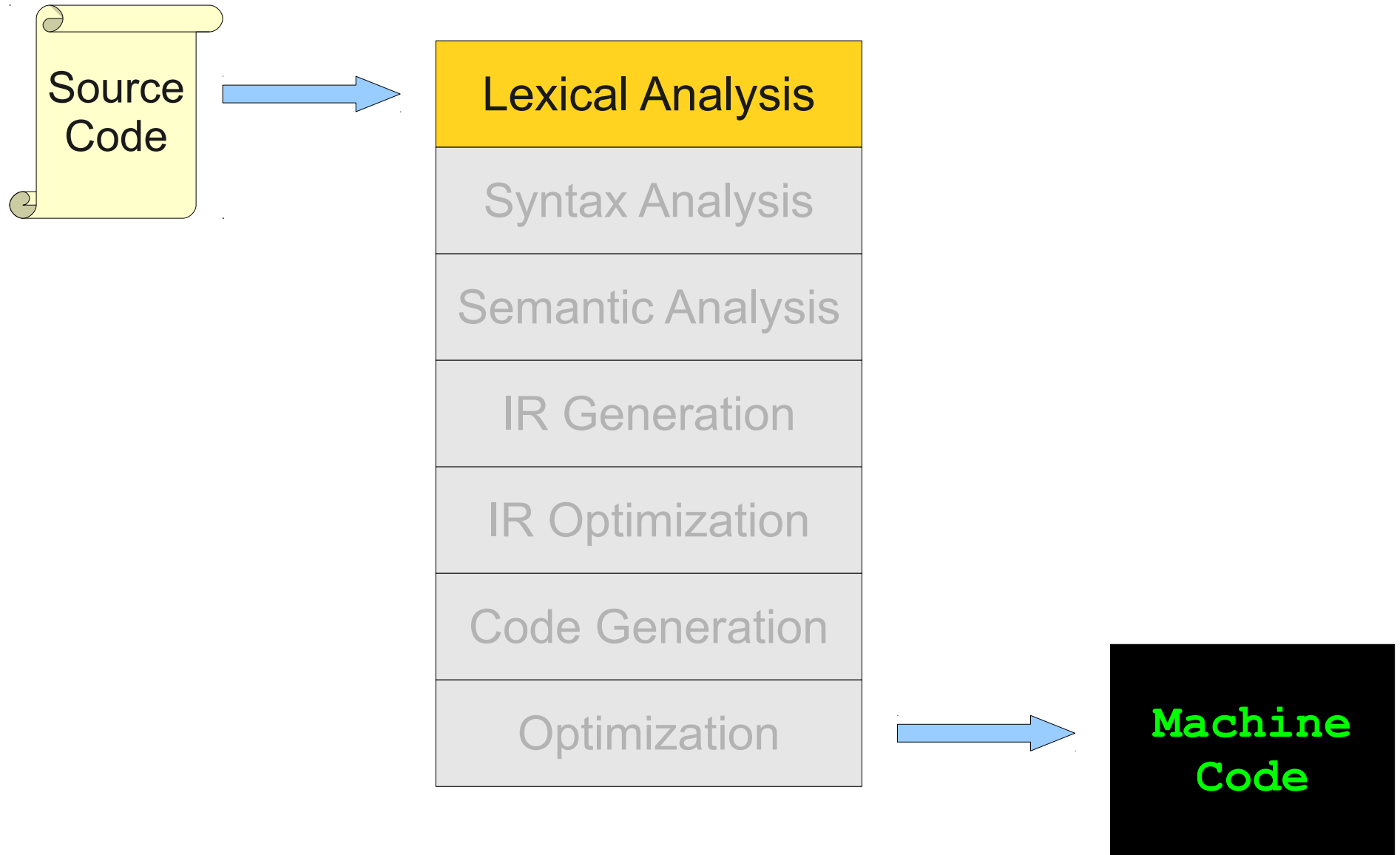
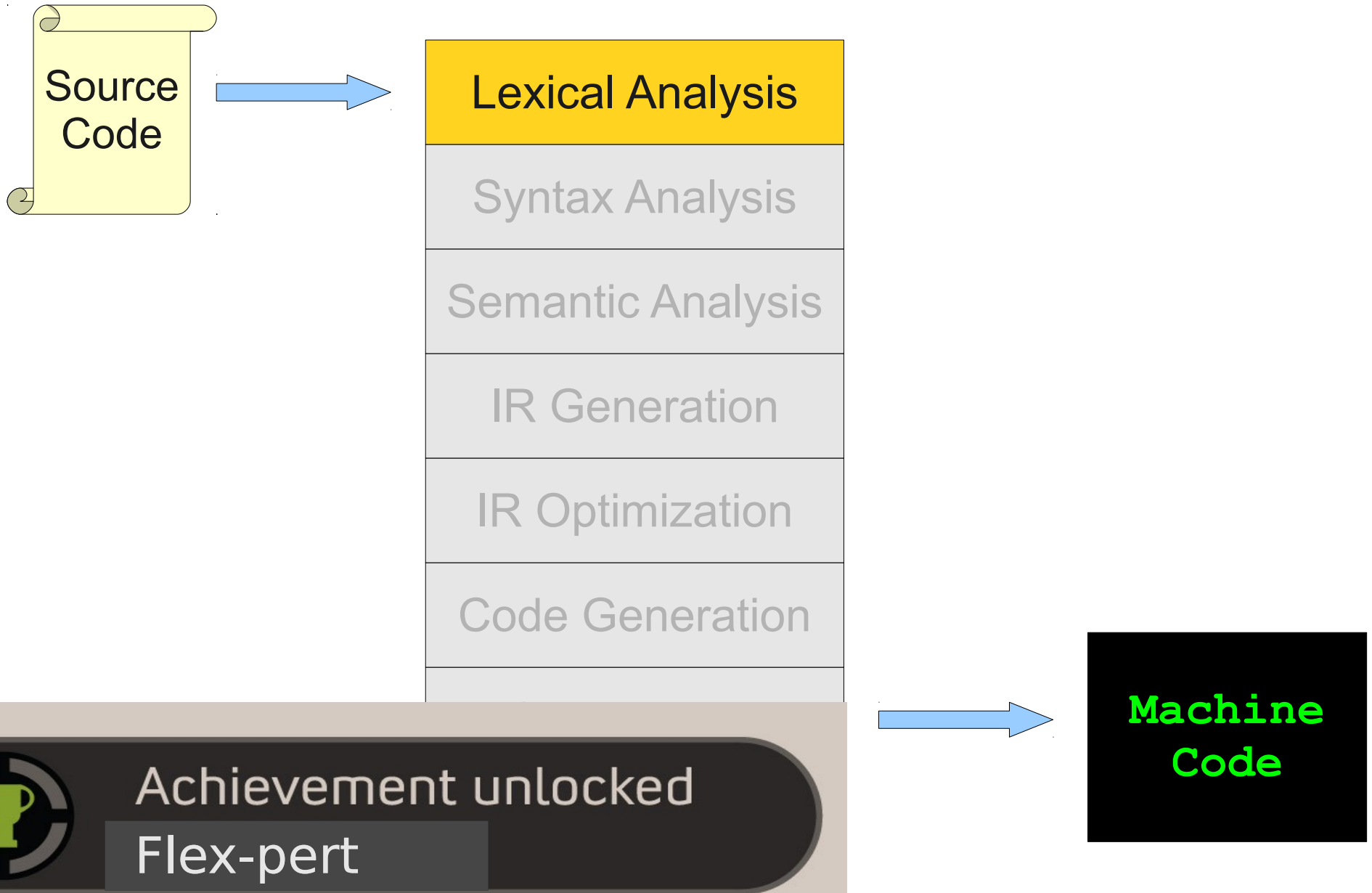


Syntax Analysis

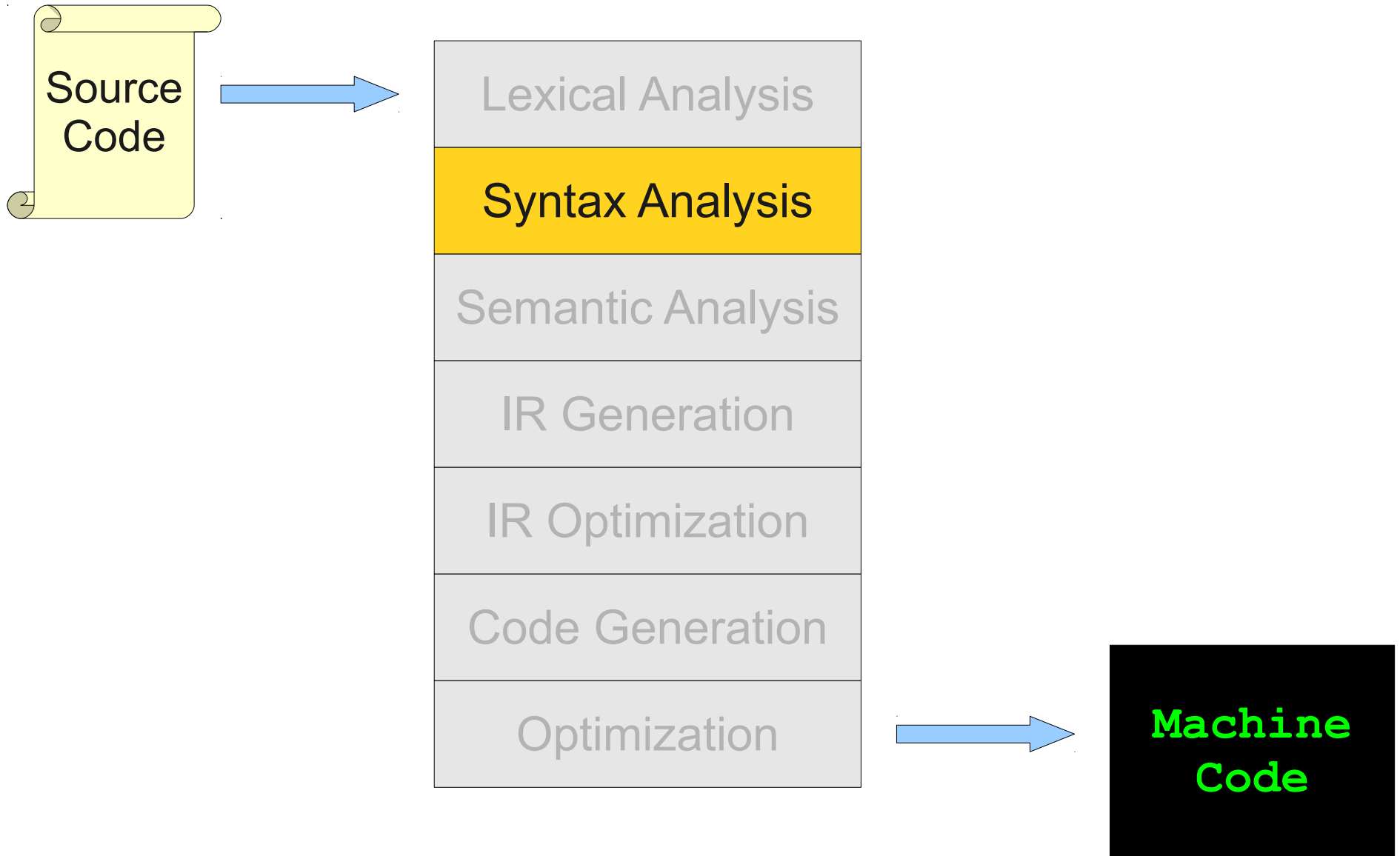
Where We Are



Where We Are



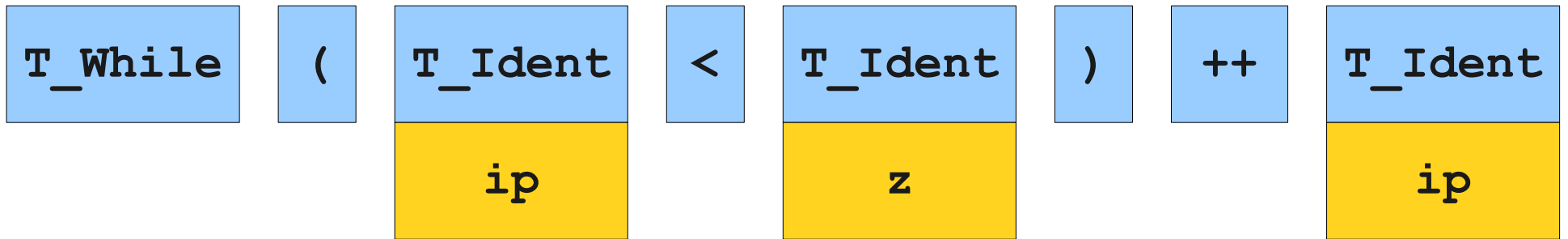
Where We Are



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

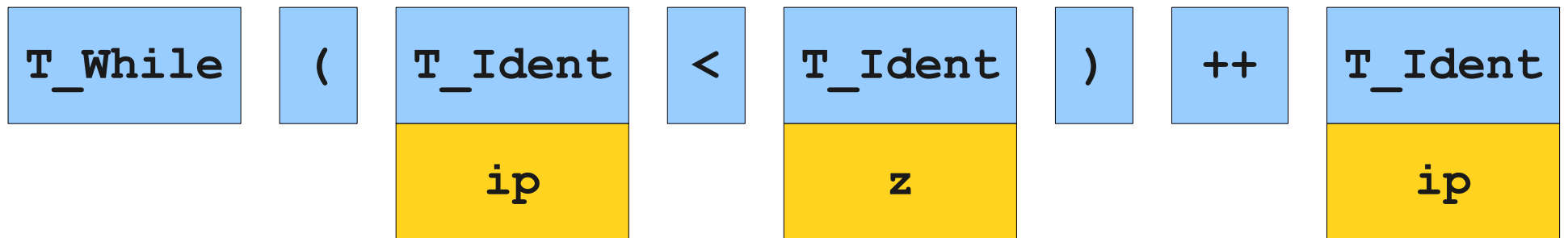
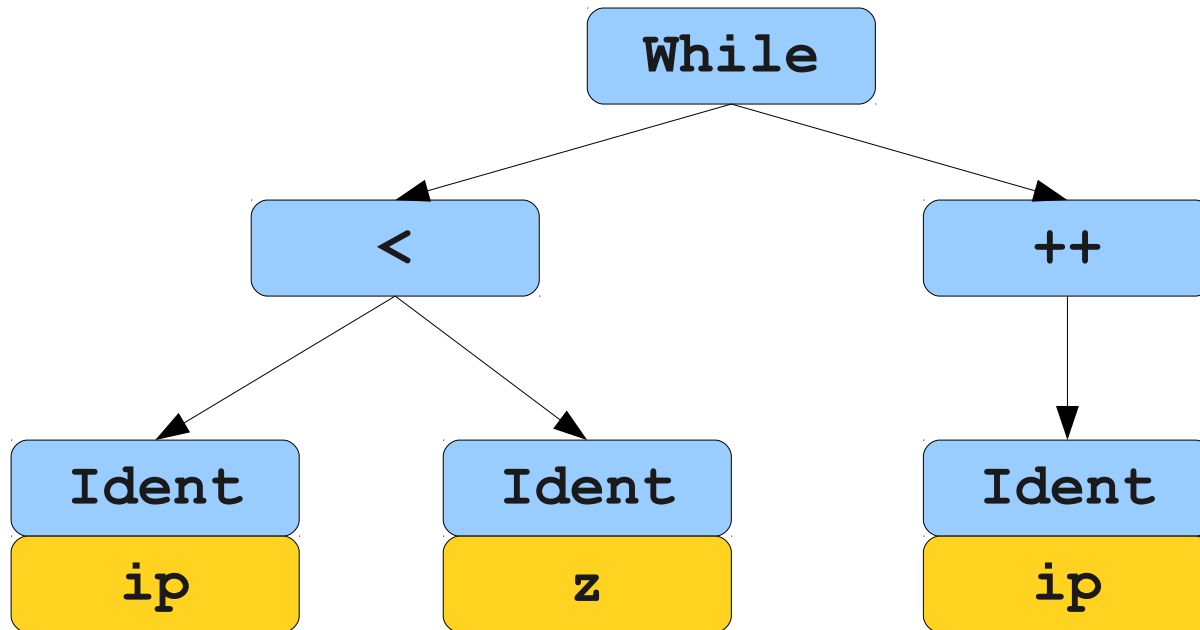
w	h	i	l	e		(i	p		<		z)	\n	\t	+	+	i	p	;
---	---	---	---	---	--	---	---	---	--	---	--	---	---	----	----	---	---	---	---	---

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



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

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



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

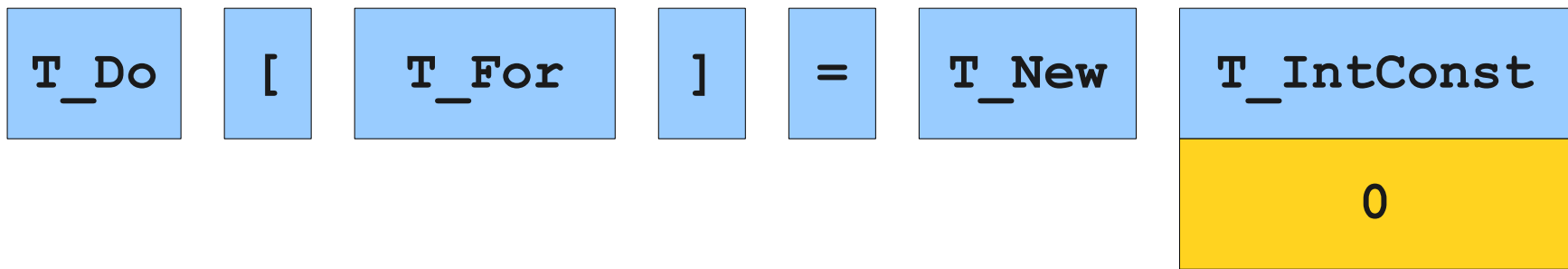
```
while (ip < z)
    ++ip;
```



```
do[for] = new 0;
```

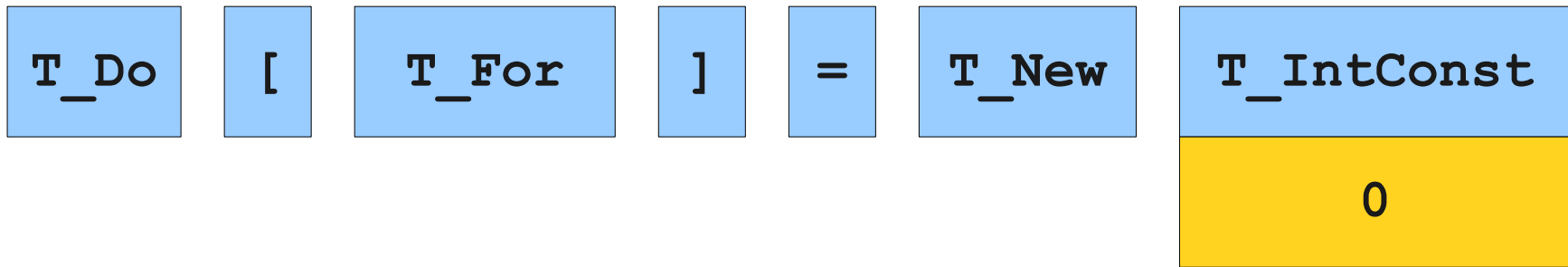
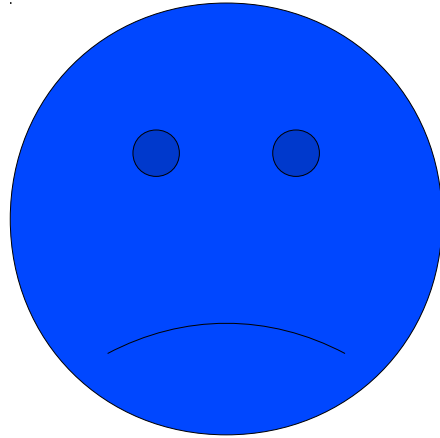
d	o	[f	o	r]		=		n	e	w		0	;
---	---	---	---	---	---	---	--	---	--	---	---	---	--	---	---

`do[for] = new 0;`



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

```
do[for] = new 0;
```



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

`do[for] = new 0;`

What is Syntax Analysis?

- After lexical analysis (scanning), we have a series of tokens.
- Goal: Recover the **structure** described by that series of tokens.
- Goal: Report **errors** if those tokens do not properly encode a structure.

Outline

- Today:
 - Formalisms for syntax analysis.
- Next Week:
 - Parsing algorithms.
 - Parser generators.

Recall: Languages

- Given an alphabet Σ , a **language** over Σ is a set of strings over Σ .
- During scanning, we were concerned with languages defining **tokens**:
 - Identifiers: {a, b, c, ..., aa, ab, ... }
 - Numbers: {0, 1, 2, ... }
 - etc.
- However, we can think of a programming language as a language as well.
 - Each string is a complete program.

Describing Programming Languages

- When dealing with tokens, we used regular expressions.
- Can we do the same for programming languages?
- In general, **no**.
- Intuition: **Automata have finite memory**.
 - The only information retained at any point is what state(s) the automaton is in.
 - No way of remembering unbounded paths through the automaton (**pumping lemma**).

Non-Regular Languages

- **Balanced parentheses.**
 - $\{ \varepsilon, (), (()), ()(), ((())), (())(), (())(), \text{etc.} \}$
- **Scheme programs.**
 - $\{ 1, 2, 3, \dots, (\text{lambda } (x) (+ x 1)), \text{etc.} \}$
- **Regular expressions.**
 - $\{ \varepsilon, a, aa, a \mid a, a \mid \varepsilon, a^*, \varepsilon^*, (a), (\varepsilon), \text{etc.} \}$

Context-Free Grammars

- A tool for describing languages that is strictly more powerful than regular languages.
- Excellent way of capturing recursive structure.
- Example: Regular expressions
 - A regular expression can be
 - Any letter
 - ϵ
 - The concatenation of regular expressions.
 - The disjunction of regular expressions.
 - The Kleene closure of a regular expression.
 - A parenthesized regular expression.

Context-Free Grammars

- A tool for describing languages that is strictly more powerful than regular languages.
- Excellent way of capturing recursive structure.
- Example: Regular expressions

$$R \rightarrow a$$
$$R \rightarrow \varepsilon$$
$$R \rightarrow RR$$
$$R \rightarrow R \mid R$$
$$R \rightarrow R^*$$
$$R \rightarrow (R)$$

CFG for Regular Expressions

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$ R

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$ **R**

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid R$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | R

$R \rightarrow RR$

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | R

$R \rightarrow RR$

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | **R**

R \rightarrow **RR**

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | R

$R \rightarrow RR$

R | RR

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | R

$R \rightarrow RR$

R | RR

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | R

$R \rightarrow RR$

R | RR

$R \rightarrow R | R$

R → R*

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid R\mathbf{R}$

$R \rightarrow R \mid R$

$R \mid RR^*$

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid$ **R**

$R \rightarrow RR$

$R \mid$ **RR**

$R \rightarrow R \mid R$

$R \mid$ **RR***

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid$ **R**

$R \rightarrow RR$

$R \mid$ **RR**

$R \rightarrow R \mid R$

$R \mid$ **RR***

$R \rightarrow R^*$

$R \rightarrow (R)$

CFG for Regular Expressions

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

R

$R \mid R$

$R \mid RR$

$R \mid RR^*$

CFG for Regular Expressions

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

R

$R \mid R$

$R \mid RR$

$R \mid RR^*$

$R \mid cR^*$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

R | **R**

$R \rightarrow RR$

R | **RR**

$R \rightarrow R | R$

R | **RR***

$R \rightarrow R^*$

$R \rightarrow (R)$

R | **cR***

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid \mathbf{RR}$

$R \rightarrow R \mid R$

$R \mid \mathbf{RR}^*$

$R \rightarrow R^*$

$R \rightarrow (R)$

R $\mid cR^*$

CFG for Regular Expressions

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

R

$R \mid R$

$R \mid RR$

$R \mid RR^*$

$R \mid cR^*$

CFG for Regular Expressions

R \rightarrow **a**

R \rightarrow ϵ

R \rightarrow RR

R \rightarrow R | R

R \rightarrow R*

R \rightarrow (R)

R

R | **R**

R | **RR**

R | **RR***

R | cR*

a | cR*

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid R\mathbf{R}$

$R \rightarrow R \mid R$

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

$R \rightarrow R^*$

$R \rightarrow (R)$

$\mathbf{R} \mid cR^*$

$a \mid cR^*$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid \mathbf{RR}$

$R \rightarrow R \mid R$

$R \mid \mathbf{RR}^*$

$R \rightarrow R^*$

$R \rightarrow (R)$

$\mathbf{R} \mid cR^*$

$a \mid c\mathbf{R}^*$

CFG for Regular Expressions

R \rightarrow **a**

R \rightarrow ϵ

R \rightarrow RR

R \rightarrow R | R

R \rightarrow R*

R \rightarrow (R)

R

R | **R**

R | **RR**

R | **RR***

R | cR*

a | c**R***

CFG for Regular Expressions

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

R

$R \mid R$

$R \mid RR$

$R \mid RR^*$

$R \mid cR^*$

$a \mid cR^*$

$a \mid cb^*$

CFG for Regular Expressions

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \mid \mathbf{R}$

$R \rightarrow RR$

$R \mid R\mathbf{R}$

$R \rightarrow R \mid R$

$R \mid \mathbf{RR}^*$

$R \rightarrow R^*$

R $\mid cR^*$

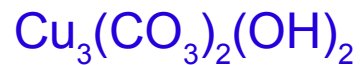
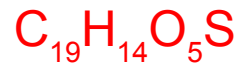
$R \rightarrow (R)$

$a \mid c\mathbf{R}^*$

$a \mid cb^*$

More Context-Free Grammars

- Example: Chemical Formulae



- $\text{Form} \rightarrow \text{Cmp} \mid \text{Cmp Ion}$
- $\text{Cmp} \rightarrow \text{Term} \mid \text{Term Num} \mid \text{Cmp Cmp}$
- $\text{Term} \rightarrow \text{Elem} \mid (\text{Cmp})$
- $\text{Elem} \rightarrow \text{H} \mid \text{He} \mid \text{Li} \mid \text{Be} \mid \text{B} \mid \text{C} \mid \text{N} \mid \text{O} \mid \dots$
- $\text{Ion} \rightarrow + \mid - \mid \text{IonNum} + \mid \text{IonNum} -$
- $\text{IonNum} \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $\text{Num} \rightarrow 1 \mid \text{IonNum}$

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$ $Form$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$ $Form$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$ $Cmp \ Ion$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$ $Form$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$ $Cmp \ Ion$
- $Term \rightarrow Elem \mid (Cmp)$ $Cmp \ Cmp \ Ion$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$ *Form*
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$ *Cmp Ion*
- $Term \rightarrow Elem \mid (Cmp)$ *Cmp Cmp Ion*
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$ *Cmp Term Num Ion*
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

Form

Cmp Ion

Cmp Cmp Ion

Cmp Term Num Ion

Term Term Num Ion

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

Form

Cmp Ion

Cmp Cmp Ion

Cmp Term Num Ion

Term Term Num Ion

Elem Term Num Ion

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

Form

Cmp Ion

Cmp Cmp Ion

Cmp Term Num Ion

Term Term Num Ion

Elem Term Num Ion

Mn *Term Num Ion*

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

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*

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

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*

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \text{ Ion}$
- $Cmp \rightarrow Term \mid Term \text{ Num} \mid Cmp \text{ Cmp}$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum + \mid IonNum -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

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*

MnO₄ *Ion*

CFGs for Chemistry

- $Form \rightarrow Cmp \mid Cmp \ Ion$
- $Cmp \rightarrow Term \mid Term \ Num \mid Cmp \ Cmp$
- $Term \rightarrow Elem \mid (Cmp)$
- $Elem \rightarrow H \mid He \mid Li \mid Be \mid B \mid C \mid N \mid \dots$
- $Ion \rightarrow + \mid - \mid IonNum \ + \mid IonNum \ -$
- $IonNum \rightarrow 2 \mid 3 \mid 4 \mid \dots$
- $Num \rightarrow 1 \mid IonNum$

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*

MnO₄ *Ion*

MnO₄⁻

CFGs in Java

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
| ...

Formal Definition of CFGs

- A **context-free grammar (CFG)** is

- A set of terminals **T**
- A set of nonterminals **N**
- A start symbol **S** in **N**
- A set of productions **P**

- A production is a rule of the form

$$A \rightarrow B_1 B_2 \dots B_n$$

where A is a nonterminal and the B_i 's are either terminals or nonterminals.

- We write $A \rightarrow \varepsilon$ if A can be replaced by no symbols.

Formal Definition of Derivations

- Given a CFG, a **derivation** is a transformation

$$A_1 A_2 A_3 \dots A_n \rightarrow A_1 A_2 \dots A_{j-1} B_1 B_2 \dots B_m A_{j+1} \dots A_n$$

if there is a production

$$A_j \rightarrow B_1 B_2 \dots B_m$$

- We say that

$$A_1 A_2 \dots A_n \rightarrow^* B_1 B_2 \dots B_m$$

if

$$A_1 A_2 \dots A_n \rightarrow \dots \rightarrow \dots \rightarrow B_1 B_2 \dots B_m$$

using zero or more derivations.

Leftmost Derivations

BLOCK → STMT
| { STMTS }

STMTS

STMTS → ε
| STMT STMTS

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

EXPR → **id**
| **constant**
| EXPR = EXPR
| EXPR + EXPR
| EXPR - EXPR
| EXPR * EXPR
| ...

Leftmost Derivations

BLOCK → STMT
| { STMTS }

STMTS

STMTS → ε
| STMT STMTS

→ *STMT STMTS*

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

EXPR → **id**
| **constant**
| EXPR = EXPR
| EXPR + EXPR
| EXPR - EXPR
| EXPR * EXPR
| ...

Leftmost Derivations

BLOCK → STMT
| { STMTS }

STMTS

STMTS → ε
| STMT STMTS

→ *STMT STMTS*

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

→ *EXPR; STMTS*

EXPR → **id**
| **constant**
| EXPR = EXPR
| EXPR + EXPR
| EXPR - EXPR
| EXPR * EXPR
| ...

Leftmost Derivations

BLOCK → STMT
| { STMTS }

STMTS

STMTS → ε
| STMT STMTS

→ *STMT STMTS*

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

→ *EXPR; STMTS*

→ *EXPR = EXPR; STMTS*

EXPR → **id**
| **constant**
| EXPR = EXPR
| EXPR + EXPR
| EXPR - EXPR
| EXPR * EXPR
| ...

Leftmost Derivations

BLOCK	→	STMT	
		{ STMTS }	
			<i>STMTS</i>
STMTS	→	ϵ	
		STMT STMTS	→ <i>STMT STMTS</i>
STMT	→	EXPR;	→ <i>EXPR; STMTS</i>
		if (EXPR) BLOCK	→ <i>EXPR = EXPR; STMTS</i>
		while (EXPR) BLOCK	
		do BLOCK while (EXPR);	→ <i>id = EXPR; STMTS</i>
		BLOCK	
		...	
EXPR	→	id	
		constant	
		EXPR = EXPR	
		EXPR + EXPR	
		EXPR - EXPR	
		EXPR * EXPR	
		...	

Leftmost Derivations

BLOCK → STMT
| { STMTS }

STMTS

STMTS → ε
| STMT STMTS

→ *STMT STMTS*

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

→ *EXPR; STMTS*

→ *EXPR = EXPR; STMTS*

→ *id = EXPR; STMTS*

→ *id = EXPR + EXPR; STMTS*

EXPR → **id**
| **constant**
| EXPR = EXPR
| EXPR + EXPR
| EXPR - EXPR
| EXPR * EXPR
| ...

Leftmost Derivations

BLOCK	→	STMT	
		{ STMTS }	
			<i>STMTS</i>
STMTS	→	ϵ	
		STMT STMTS	→ <i>STMT STMTS</i>
STMT	→	EXPR;	→ <i>EXPR; STMTS</i>
		if (EXPR) BLOCK	→ <i>EXPR = EXPR; STMTS</i>
		while (EXPR) BLOCK	
		do BLOCK while (EXPR);	→ <i>id = EXPR; STMTS</i>
		BLOCK	→ <i>id = EXPR + EXPR; STMTS</i>
		...	→ <i>id = id + EXPR; STMTS</i>
EXPR	→	id	
		constant	
		EXPR = EXPR	
		EXPR + EXPR	
		EXPR - EXPR	
		EXPR * EXPR	
		...	

Leftmost Derivations

BLOCK	→	STMT	
		{ STMTS }	
			<i>STMTS</i>
STMTS	→	ϵ	
		STMT STMTS	→ <i>STMT STMTS</i>
STMT	→	EXPR;	→ <i>EXPR; STMTS</i>
		if (EXPR) BLOCK	→ <i>EXPR = EXPR; STMTS</i>
		while (EXPR) BLOCK	
		do BLOCK while (EXPR);	→ <i>id = EXPR; STMTS</i>
		BLOCK	→ <i>id = EXPR + EXPR; STMTS</i>
		...	→ <i>id = id + EXPR; STMTS</i>
EXPR	→	id	→ <i>id = id + constant; STMTS</i>
		constant	
		EXPR = EXPR	
		EXPR + EXPR	
		EXPR - EXPR	
		EXPR * EXPR	
		...	

Leftmost Derivations

BLOCK	→	STMT	
		{ STMTS }	
			<i>STMTS</i>
STMTS	→	ϵ	
		STMT STMTS	→ <i>STMT STMTS</i>
STMT	→	EXPR;	→ <i>EXPR; STMTS</i>
		if (EXPR) BLOCK	→ <i>EXPR = EXPR; STMTS</i>
		while (EXPR) BLOCK	
		do BLOCK while (EXPR);	→ <i>id = EXPR; STMTS</i>
		BLOCK	→ <i>id = EXPR + EXPR; STMTS</i>
		...	→ <i>id = id + EXPR; STMTS</i>
EXPR	→	id	→ <i>id = id + constant; STMTS</i>
		constant	→ <i>id = id + constant;</i>
		EXPR = EXPR	
		EXPR + EXPR	
		EXPR - EXPR	
		EXPR * EXPR	
		...	

Leftmost Derivations

- A **leftmost derivation** expands the leftmost nonterminal first.
- A **rightmost derivation** expands the rightmost nonterminal first.
- These will be of great importance when we talk about parsing next week.

Parse Trees

$R \rightarrow a$

R

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$

Parse Trees

$R \rightarrow a$

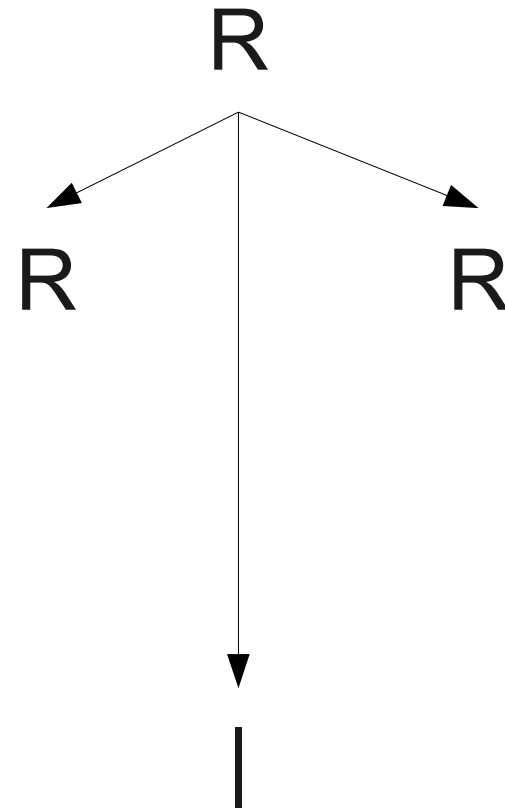
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

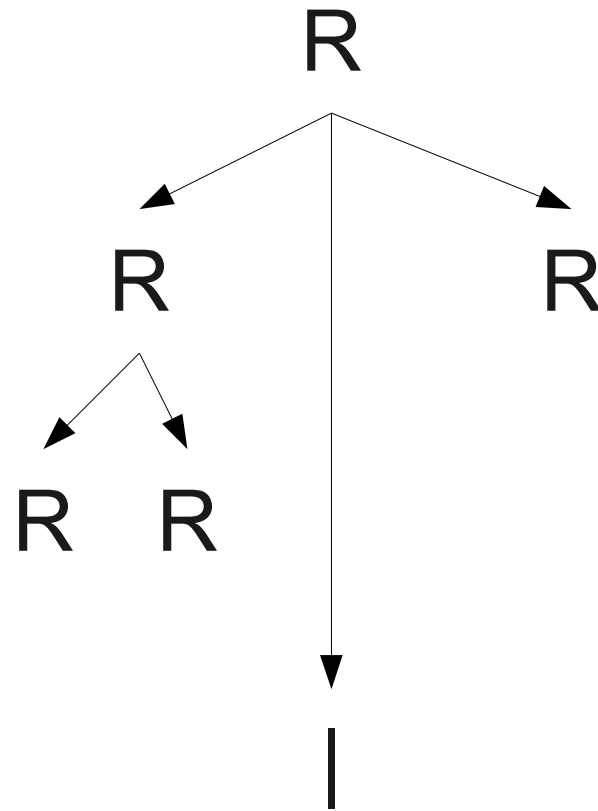
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

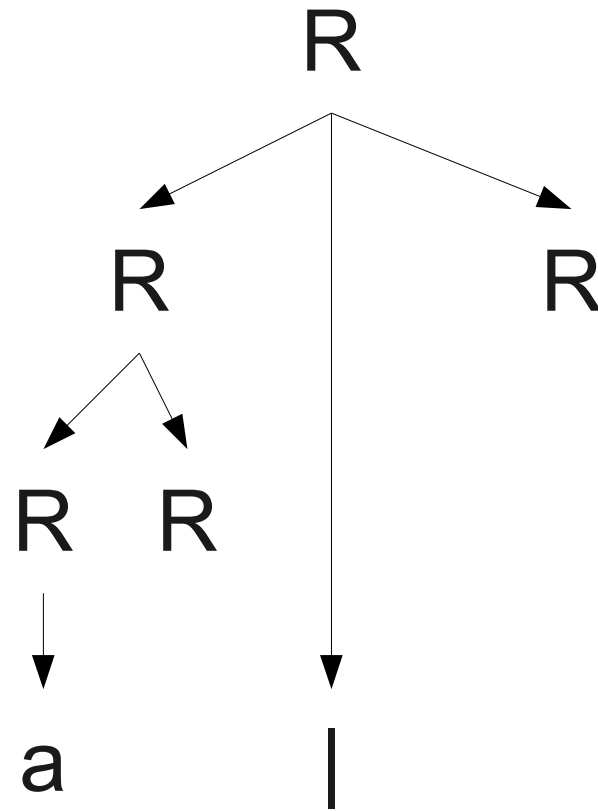
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

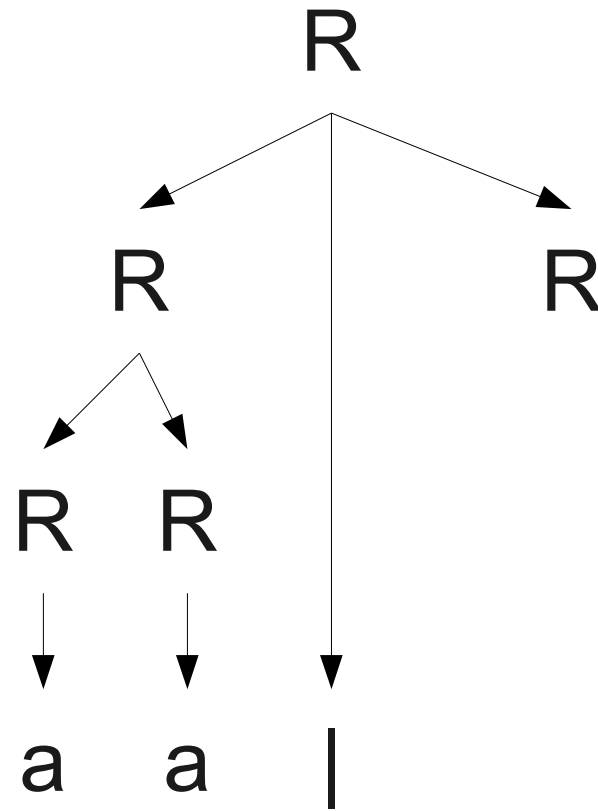
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

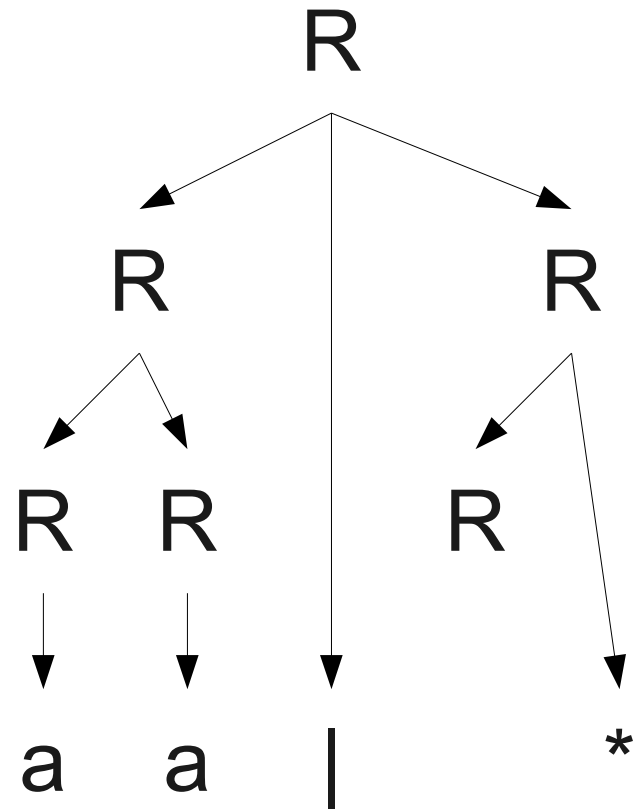
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

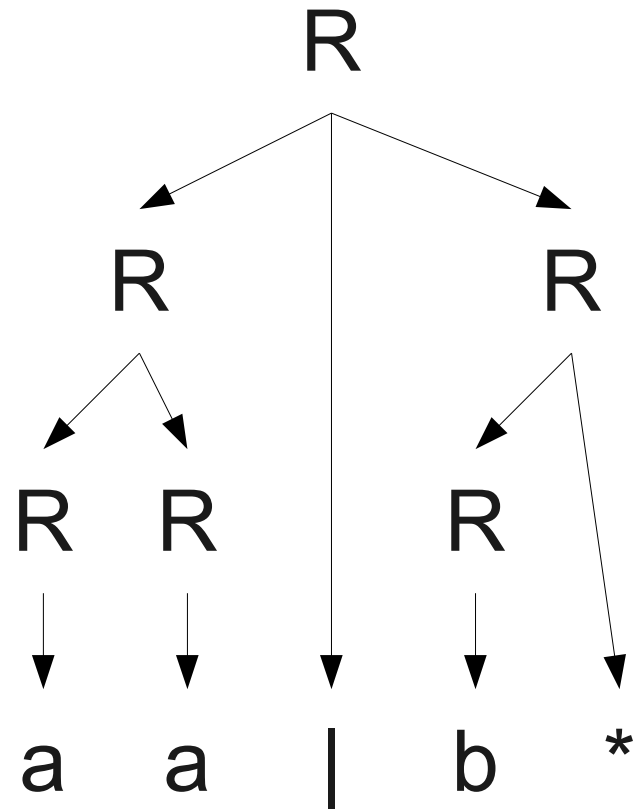
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

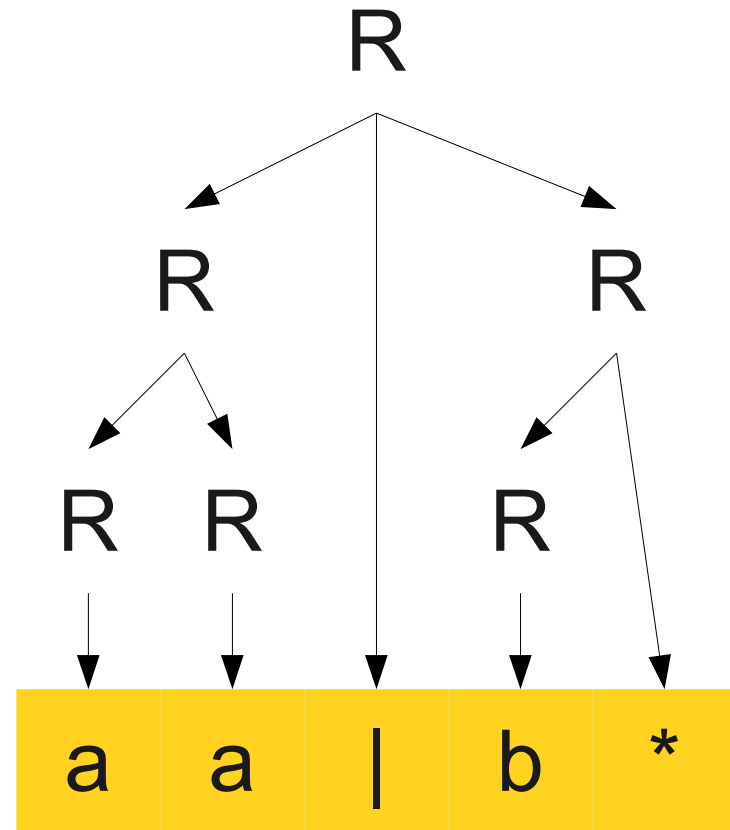
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

- Tree structure encoding the steps in a derivation.
- Inorder walk of the leaves contains the generated string.
- A derivation encodes **how** to produce the input.
- A parse tree encodes the **structure** of the input.

Syntax Analysis

- Goal of syntax analysis: Recover the **structure** described by a series of tokens.
- Use a context-free grammar:
 - Define a CFG for the language whose terminals are the tokens emitted by the scanner.
 - Obtain a parse tree for those tokens to recover the underlying structure.

Parse Trees

$R \rightarrow a$

$R \rightarrow \varepsilon$

$R \rightarrow RR$

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$

R

Parse Trees

$R \rightarrow a$

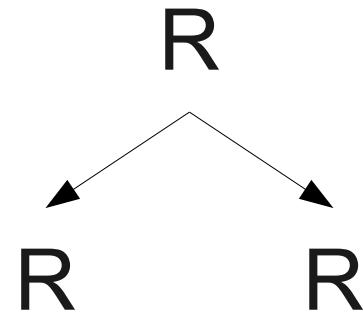
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

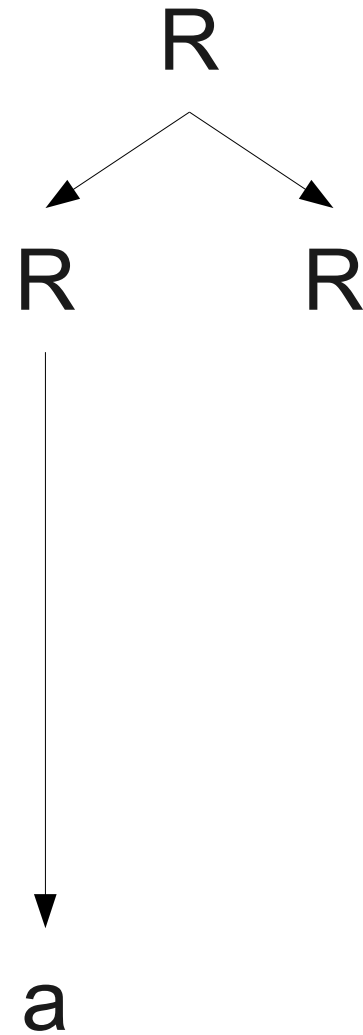
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

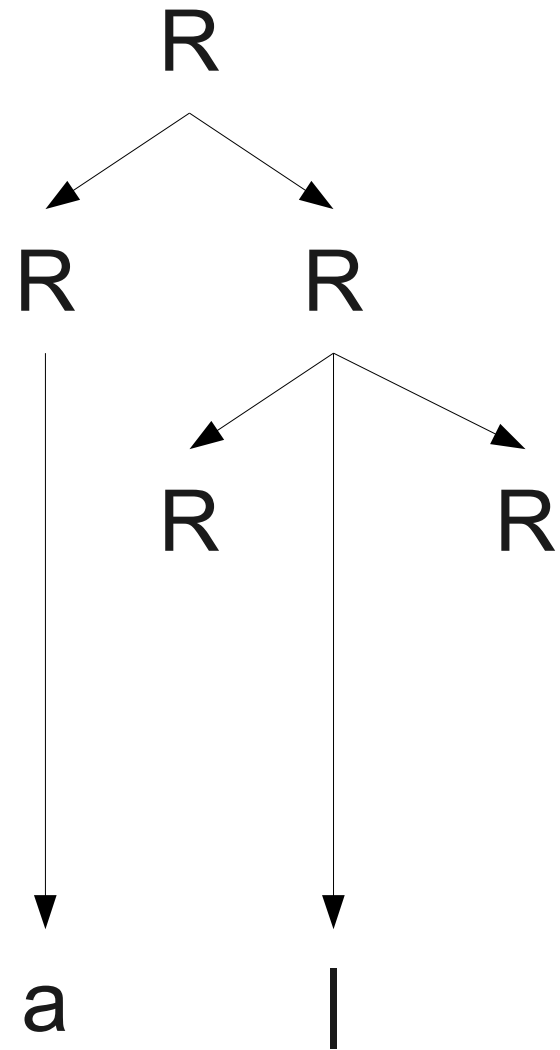
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

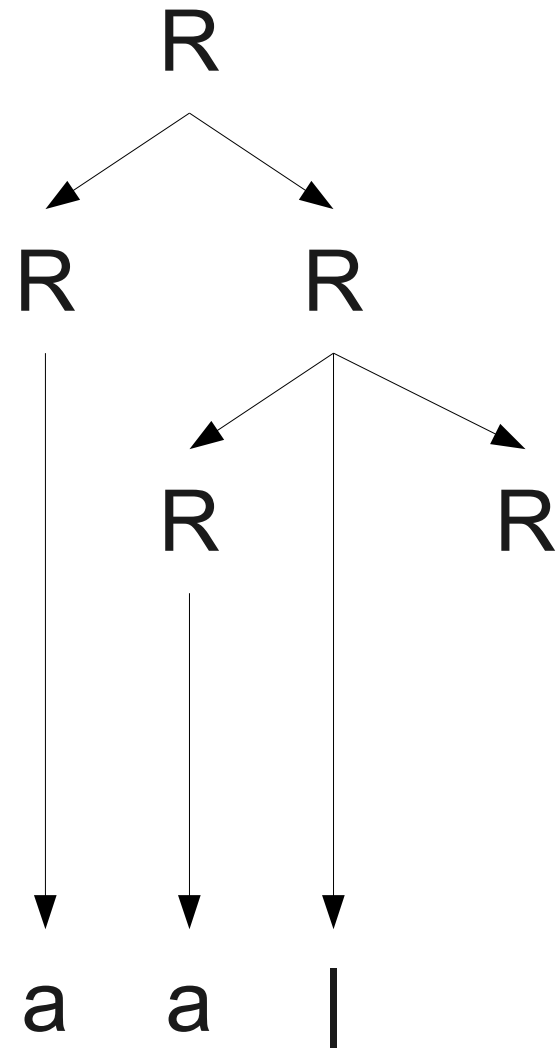
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

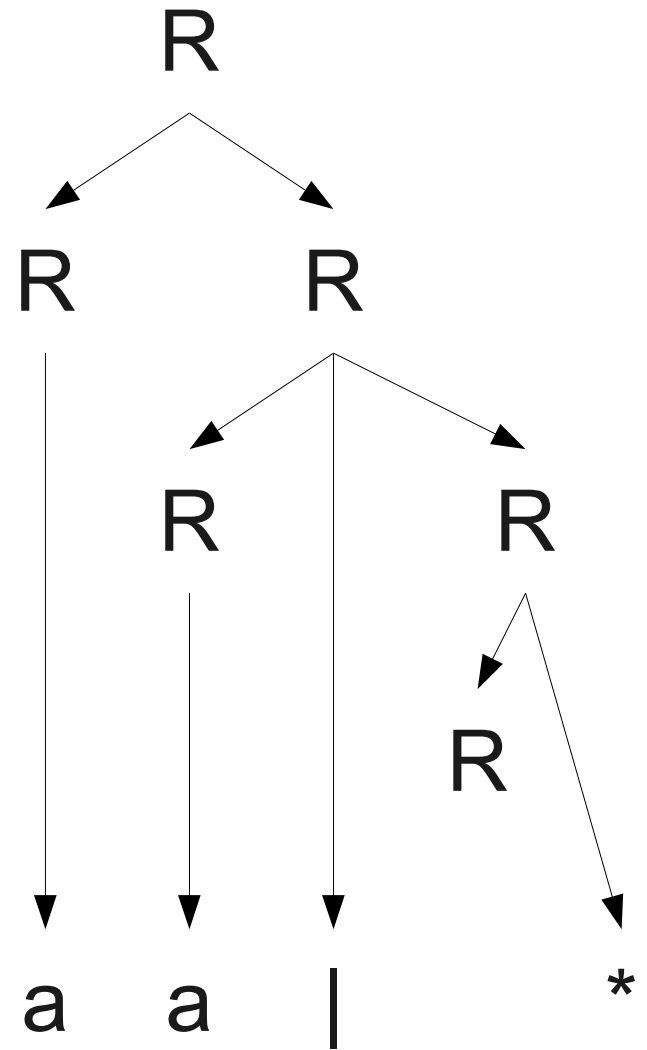
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

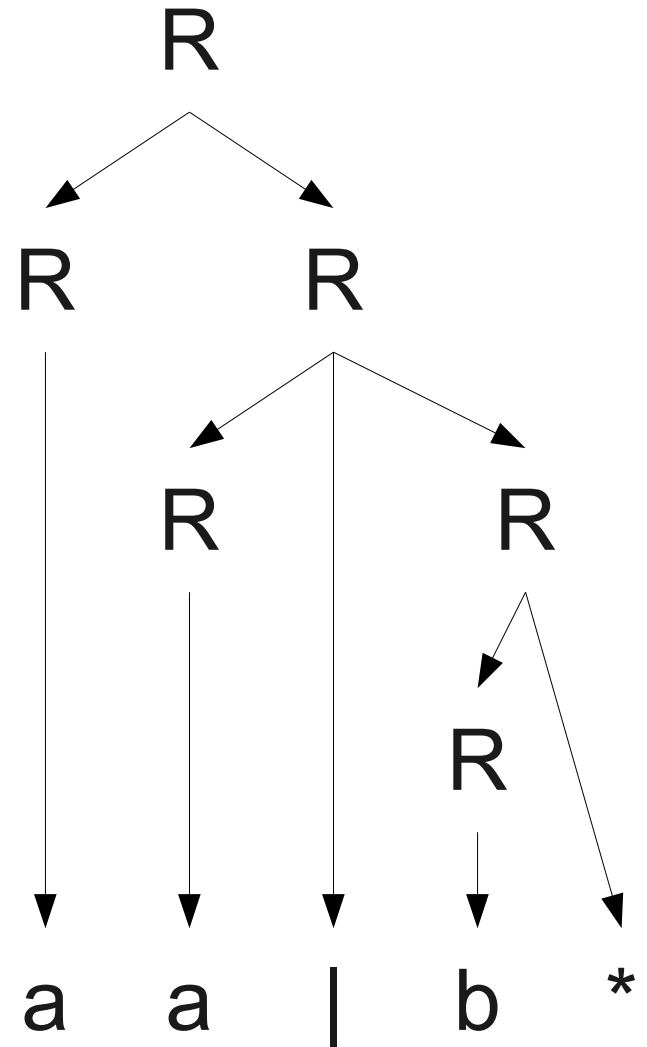
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R | R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Parse Trees

$R \rightarrow a$

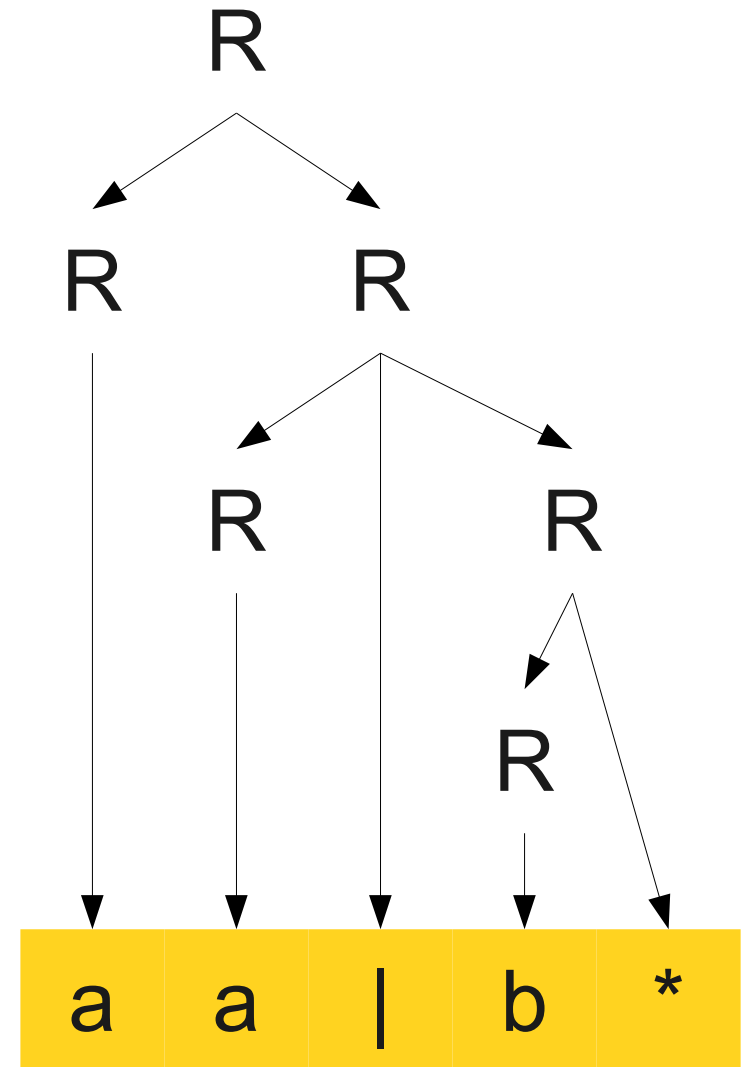
$R \rightarrow \epsilon$

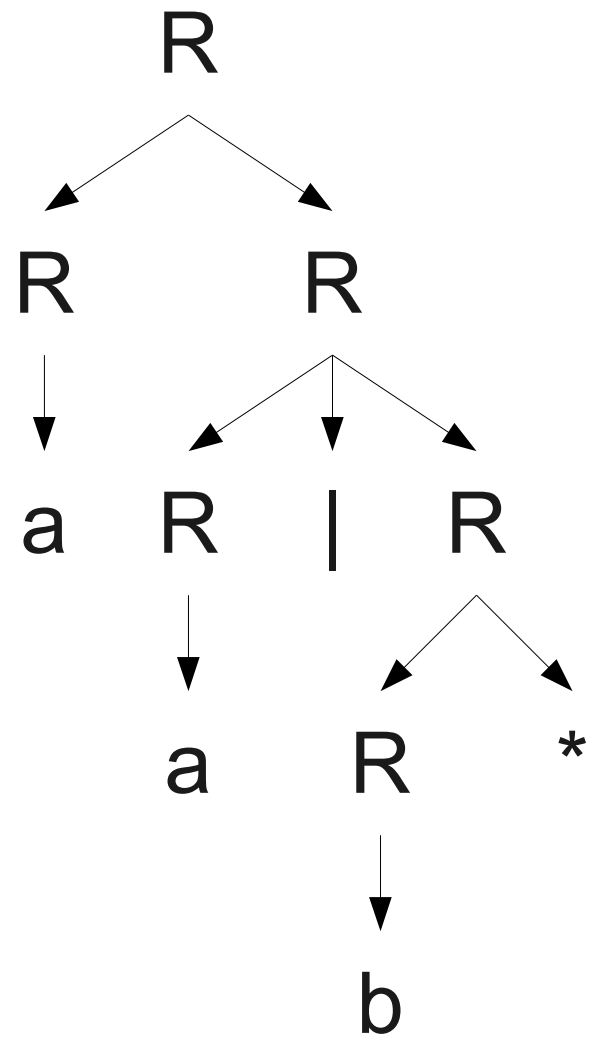
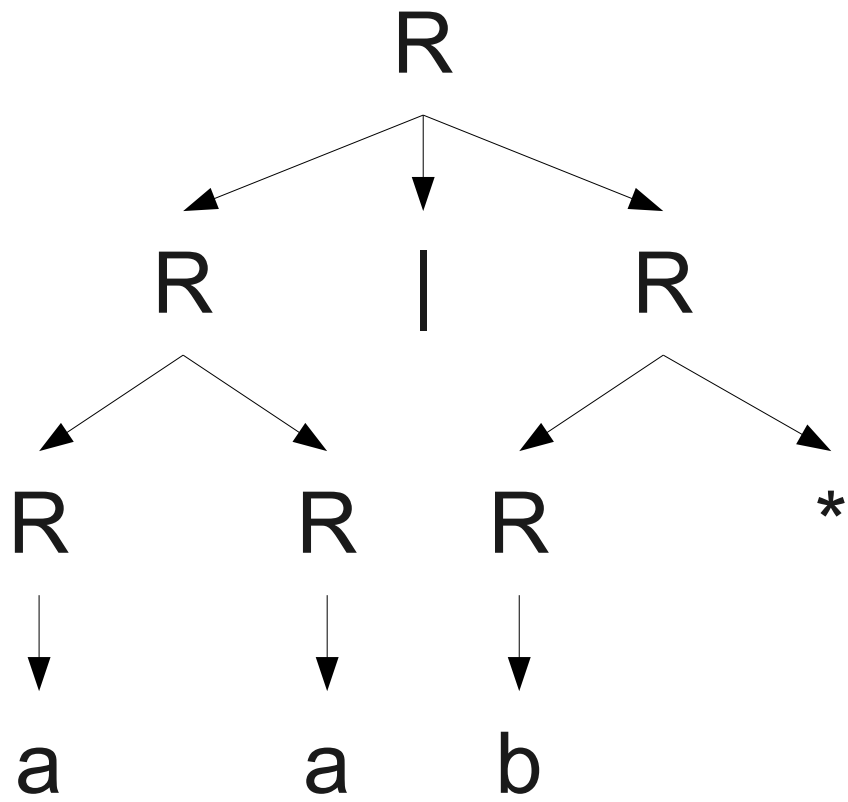
$R \rightarrow RR$

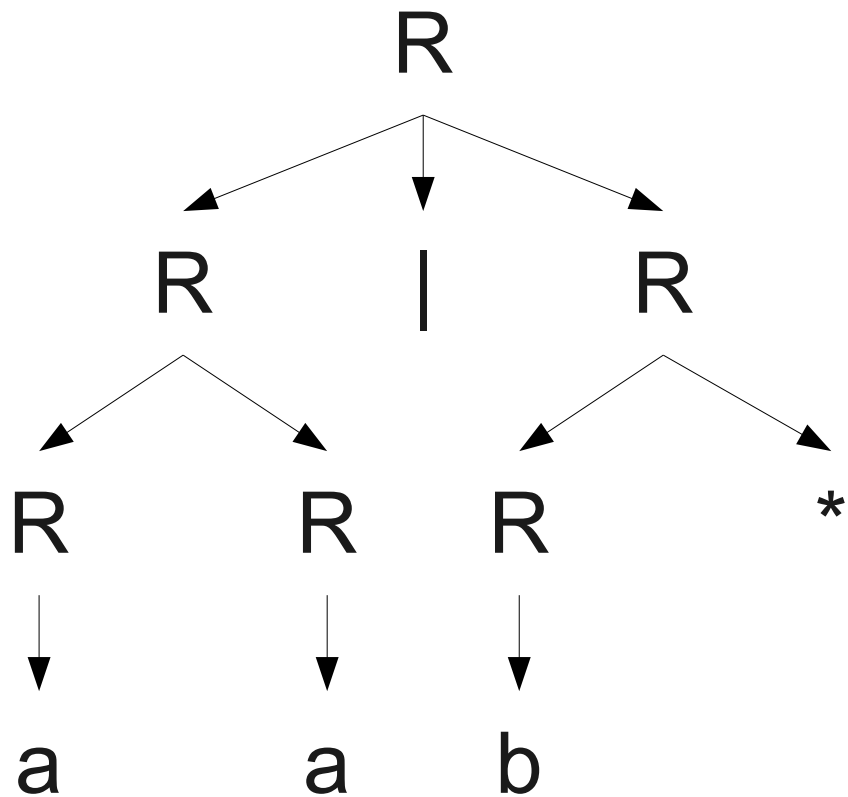
$R \rightarrow R | R$

$R \rightarrow R^*$

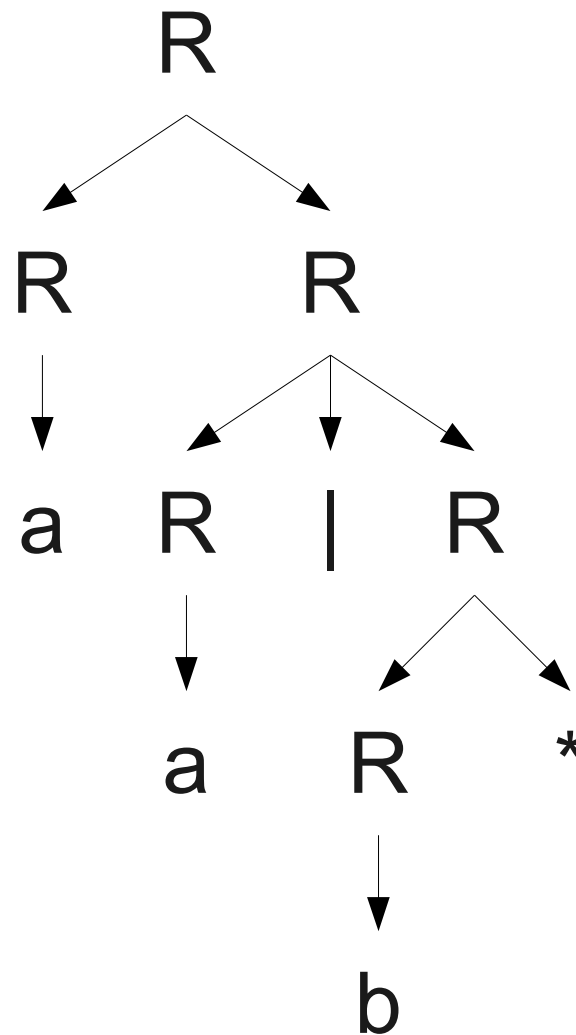
$R \rightarrow (R)$







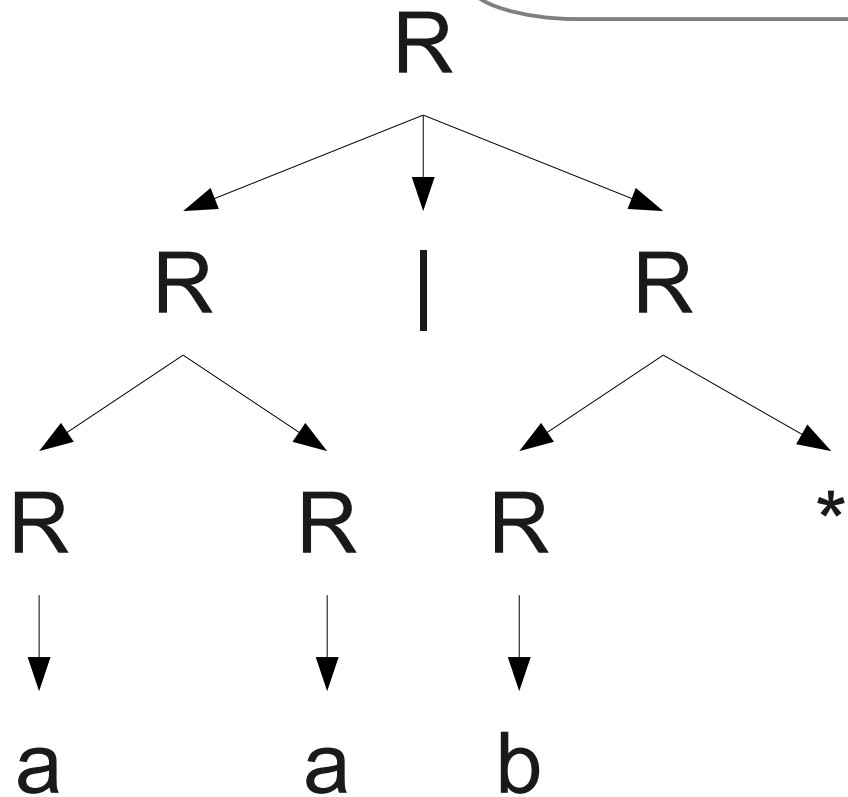
$(aa) \mid (b^*)$



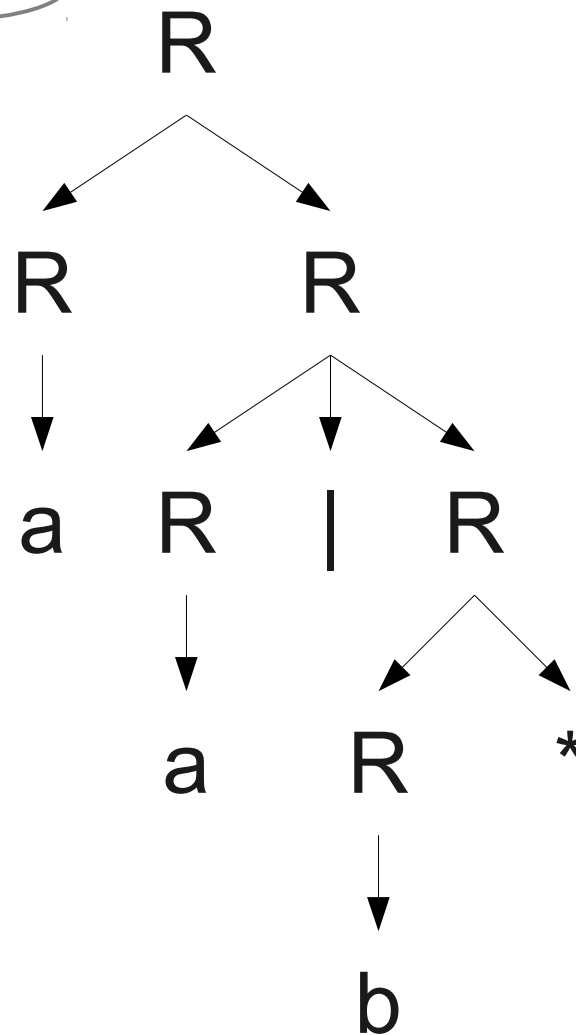
$a(a \mid (b^*))$



Problem?



$(aa) \mid (b^*)$



$a(a \mid (b^*))$

Ambiguity

- A CFG is said to be **ambiguous** if there is at least one string with two or more derivations.
- Unfortunately:
 - It is **undecidable** whether a given grammar is ambiguous.
 - There are some languages that are **inherently** ambiguous.
- However, in most cases, the ambiguity is in the **grammar** and not the **language**.

Is Ambiguity a Problem?

Is Ambiguity a Problem?

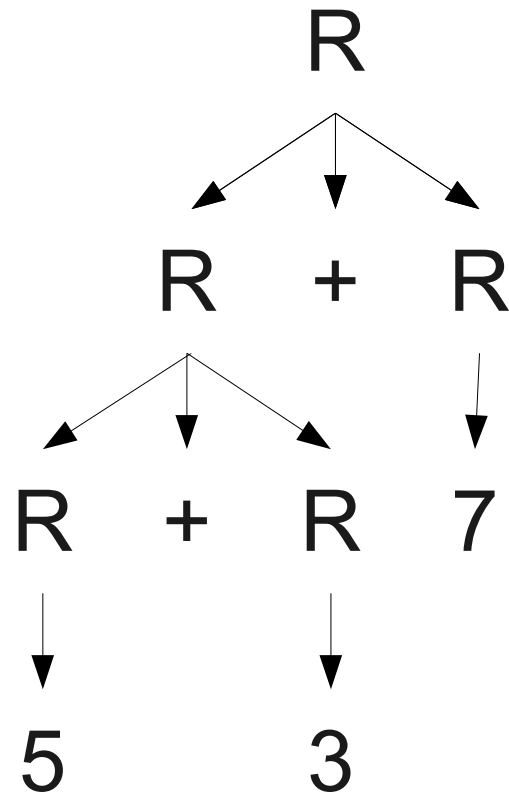
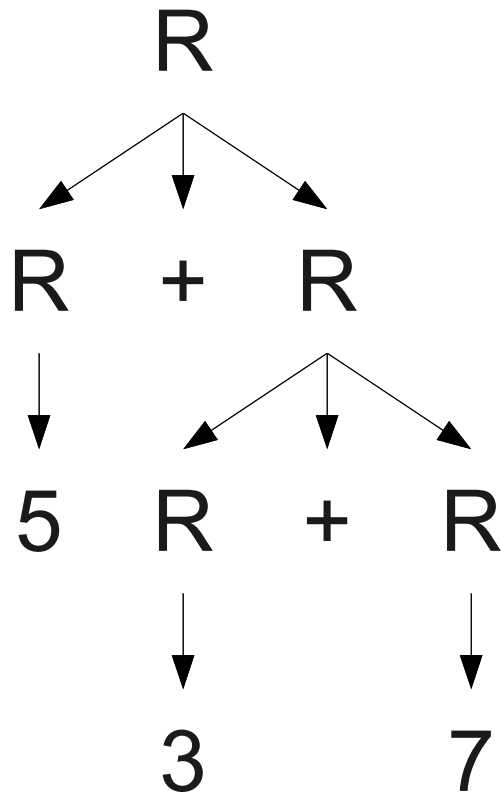
- Depends on **semantics**.

Is Ambiguity a Problem?

- Depends on **semantics**.
- $E \rightarrow \mathbf{int} \mid E + E$

Is Ambiguity a Problem?

- Depends on **semantics**.
- $E \rightarrow \text{int} \mid E + E$

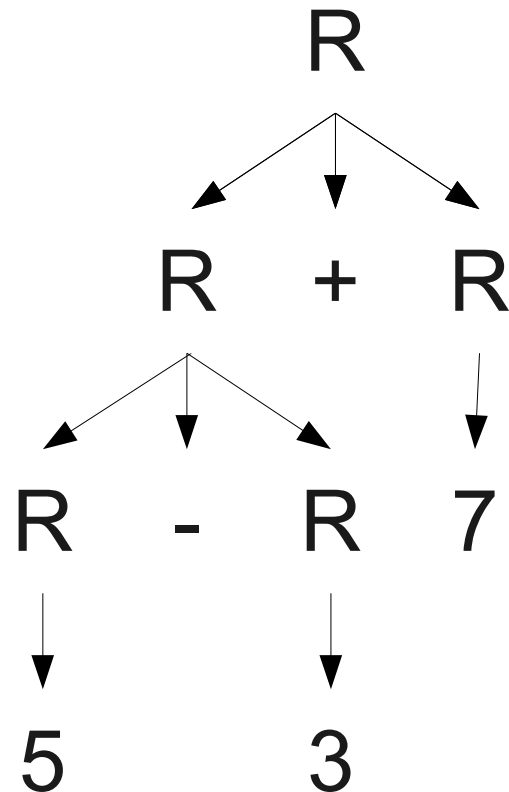
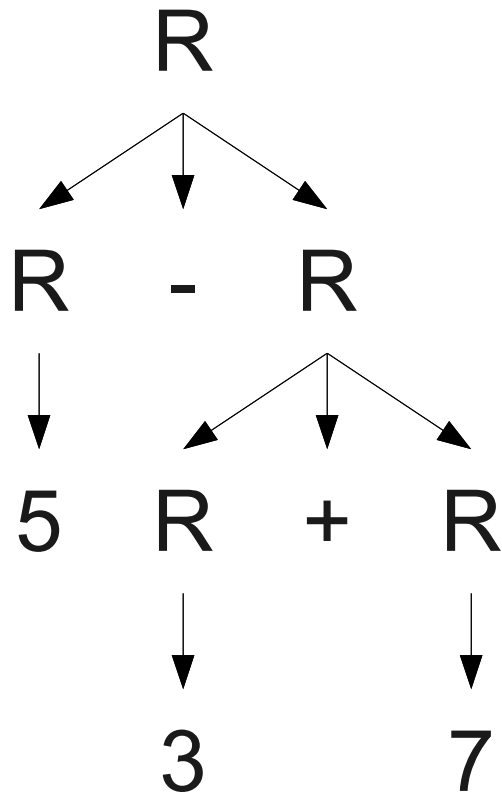


Is Ambiguity a Problem?

- Depends on **semantics**.
- $E \rightarrow \mathbf{int} \mid E + E \mid E - E$

Is Ambiguity a Problem?

- Depends on **semantics**.
- $E \rightarrow \mathbf{int} \mid E + E \mid E - E$



Resolving Ambiguity

- Return to our CFG for regular expressions:

$$R \rightarrow R \mid R$$

$$R \rightarrow RR$$

$$R \rightarrow R^*$$

$$R \rightarrow a$$

$$R \rightarrow \varepsilon$$

$$R \rightarrow (R)$$

- Kleene closure should have highest precedence, then concatenation, then disjunction.

Resolving Ambiguity

- Return to our CFG for regular expressions:

$R \rightarrow R \mid R$

$R \rightarrow RR$

$R \rightarrow R^*$

$R \rightarrow "a"$

$R \rightarrow "\epsilon"$

$R \rightarrow (R)$

- Kleene closure should have highest precedence, then concatenation, then disjunction.

Resolving Ambiguity

- Return to our CFG for regular expressions:

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

$$R \rightarrow RR$$

$$R \rightarrow R^*$$

$$R \rightarrow \text{"a"}$$

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

$$R \rightarrow (R)$$

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

$$S \rightarrow T \mid ST$$

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

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

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

$$U \rightarrow (R)$$

- Kleene closure should have highest precedence, then concatenation, then disjunction.

Why is this unambiguous?

$$R \rightarrow S \mid R \mid S$$
$$S \rightarrow T \mid ST$$
$$T \rightarrow U \mid T^*$$
$$U \rightarrow "a"$$
$$U \rightarrow "\epsilon"$$
$$U \rightarrow (R)$$

Why is this unambiguous?

$R \rightarrow S \mid R \mid S$
 $S \rightarrow T \mid ST$
 $T \rightarrow U \mid T^*$

$U \rightarrow "a"$
 $U \rightarrow "\epsilon"$
 $U \rightarrow (R)$

Only generates
"atomic" expressions

Why is this unambiguous?

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

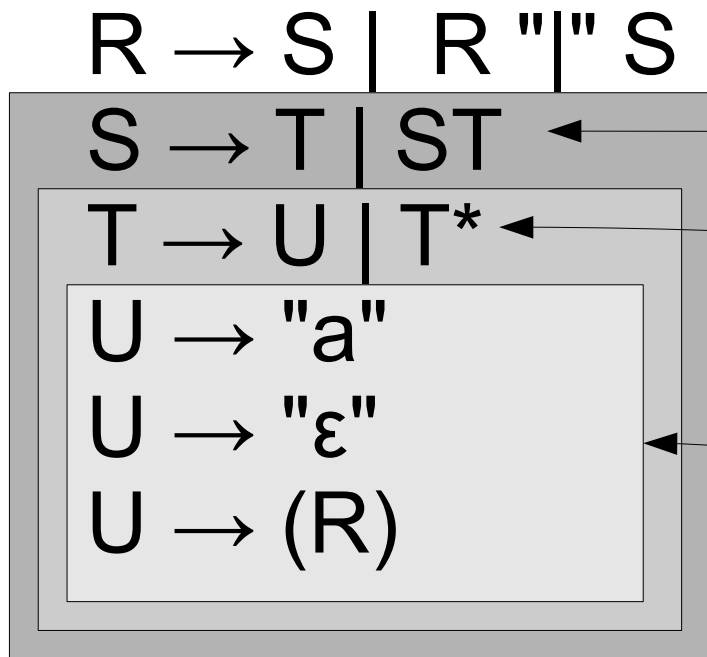
$U \rightarrow \text{"\epsilon"}$

$U \rightarrow (R)$

Puts stars onto
atomic expressions

Only generates
"atomic" expressions

Why is this unambiguous?

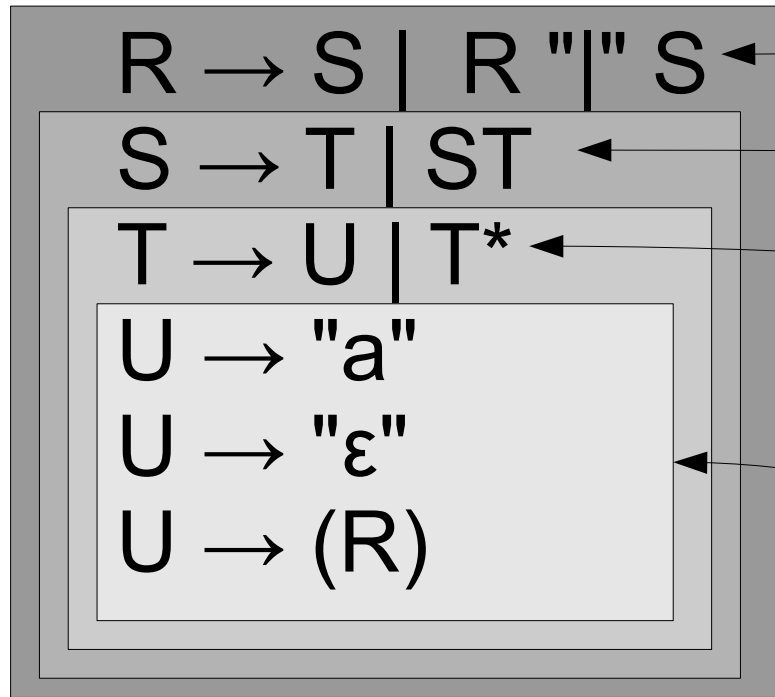


Concatenates starred expressions

Puts stars onto atomic expressions

Only generates "atomic" expressions

Why is this unambiguous?



Ors concatenated expressions

Concatenates starred expressions

Puts stars onto atomic expressions

Only generates "atomic" expressions

Resolving Ambiguity: Precedence

- If we leave the world of pure CFGs, we can often resolve ambiguities through **precedence declarations**.
- Used in many parser generators, including **bison**.
- We'll see how these are implemented later on.

The Structure of a Parse Tree

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

$U \rightarrow "a"$

$U \rightarrow "\epsilon"$

$U \rightarrow (R)$

The Structure of a Parse Tree

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

$U \rightarrow "a"$

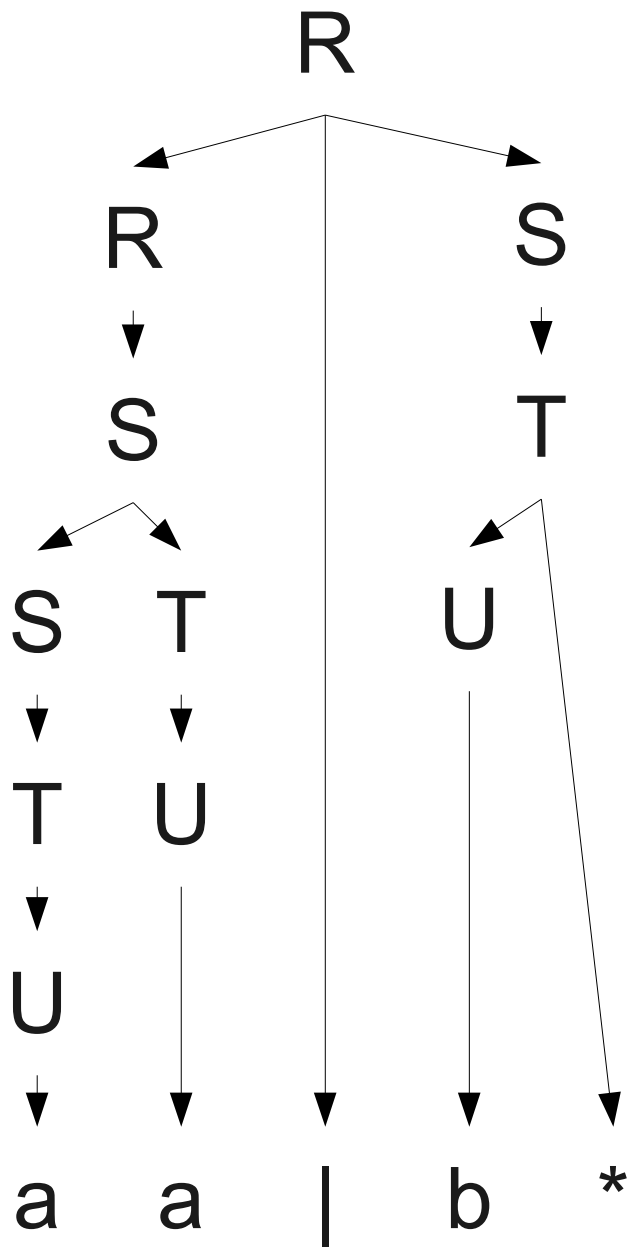
$U \rightarrow "\epsilon"$

$U \rightarrow (R)$

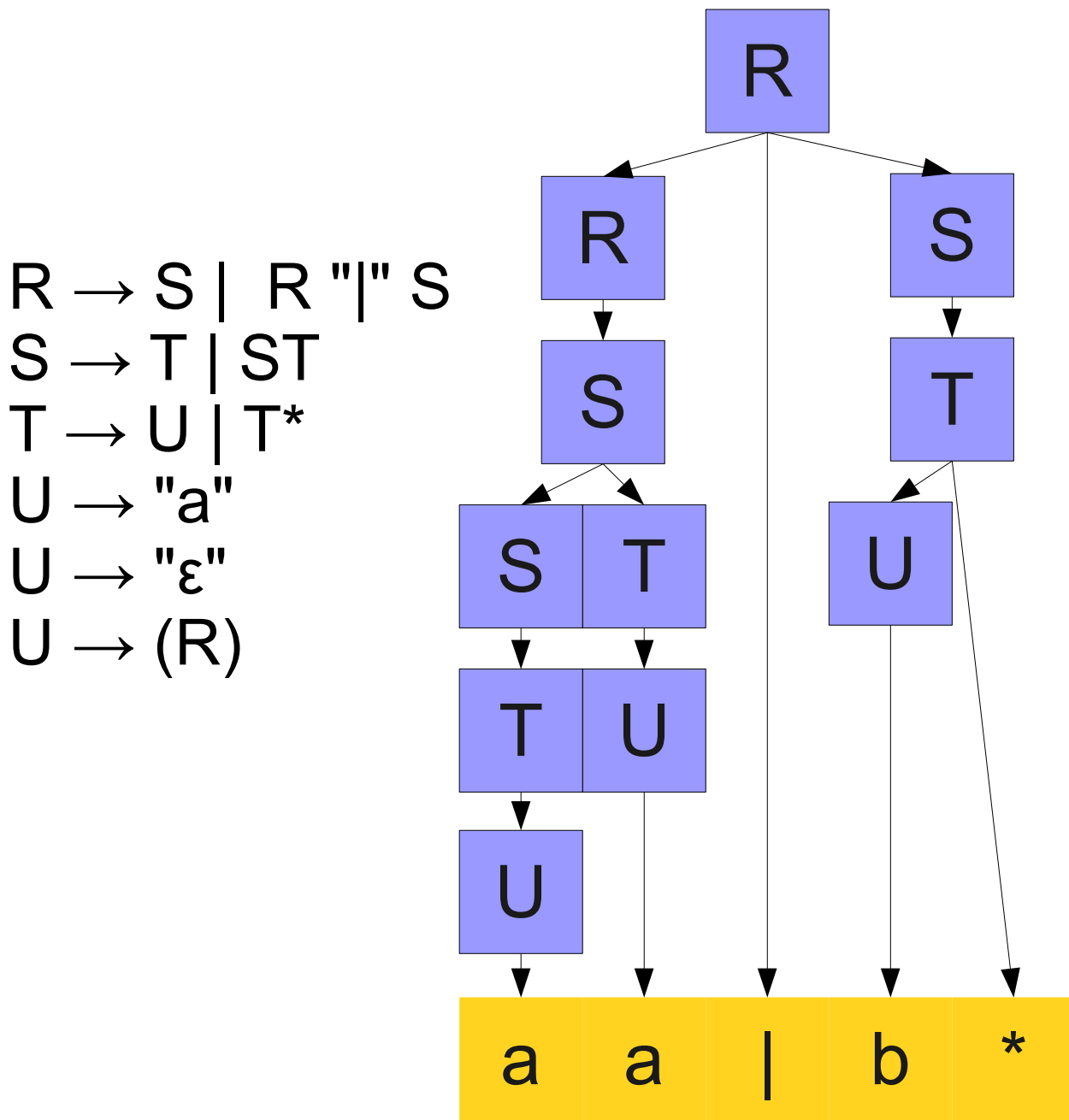
a a | b *

The Structure of a Parse Tree

$R \rightarrow S \mid R \mid S$
 $S \rightarrow T \mid ST$
 $T \rightarrow U \mid T^*$
 $U \rightarrow "a"$
 $U \rightarrow "\epsilon"$
 $U \rightarrow (R)$

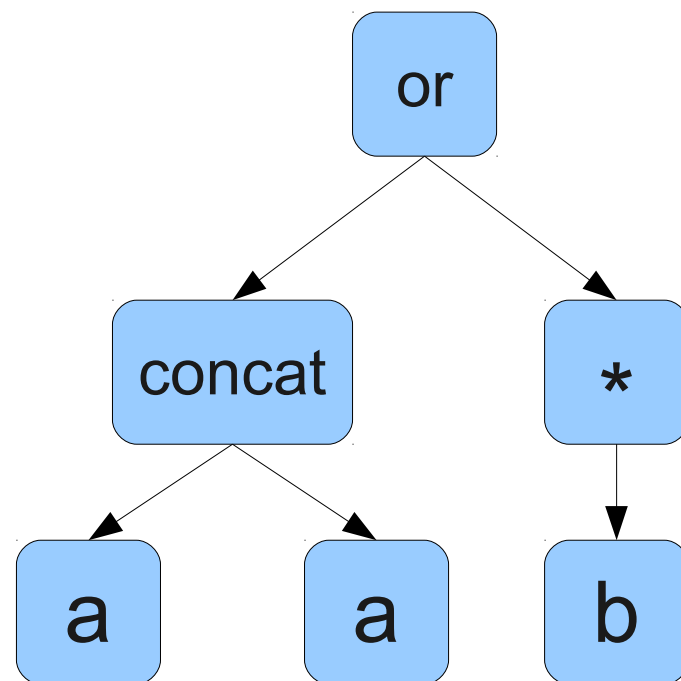
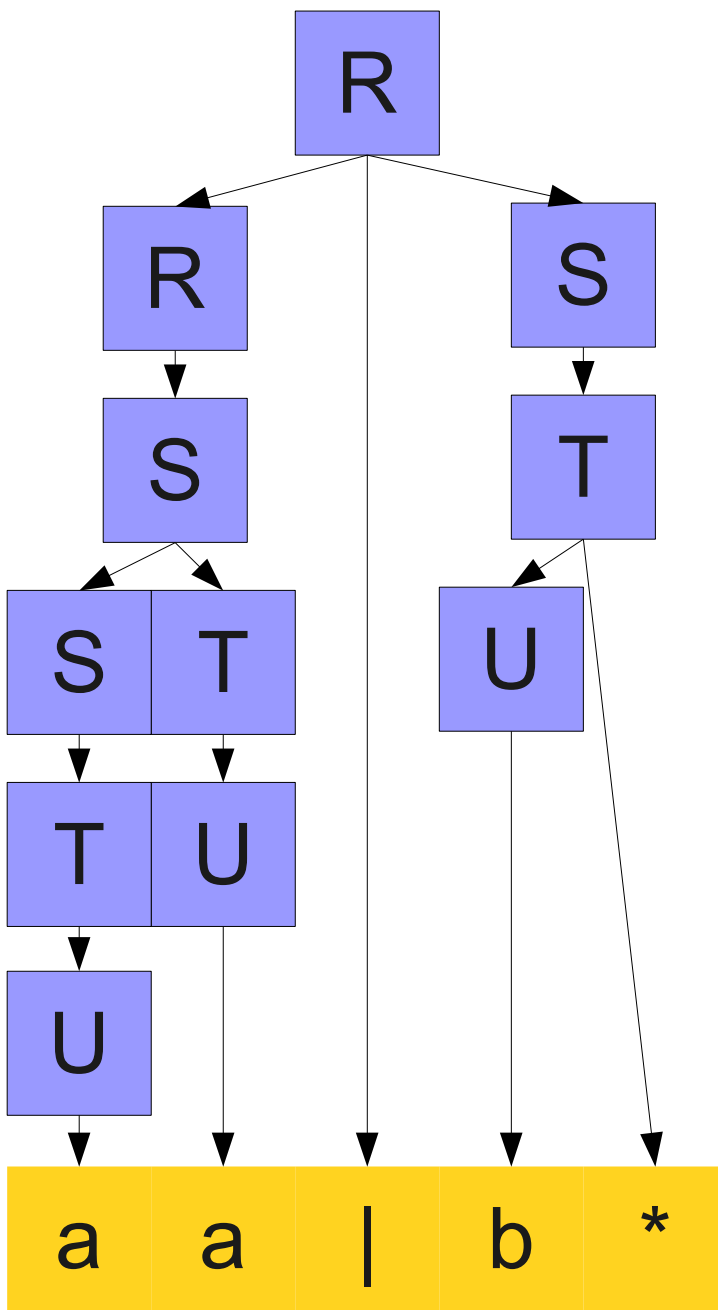


The Structure of a Parse Tree



The Structure of a Parse Tree

$R \rightarrow S \mid R \mid S$
 $S \rightarrow T \mid ST$
 $T \rightarrow U \mid T^*$
 $U \rightarrow "a"$
 $U \rightarrow "\epsilon"$
 $U \rightarrow (R)$



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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

$U \rightarrow "a"$

$U \rightarrow "\epsilon"$

$U \rightarrow (R)$

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

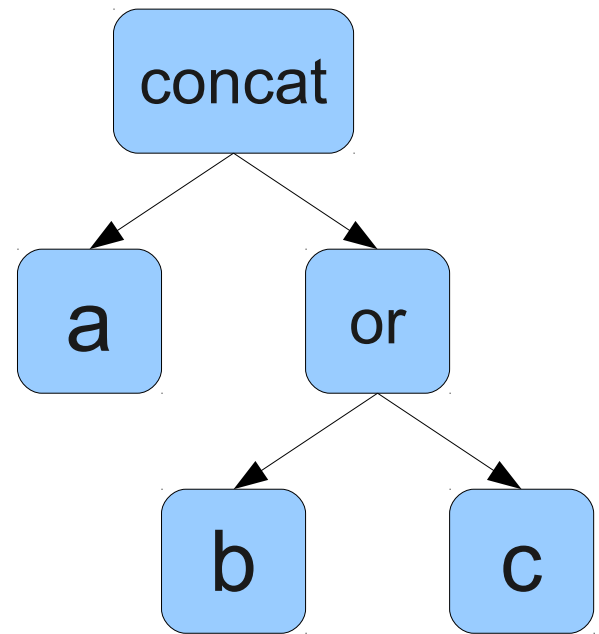
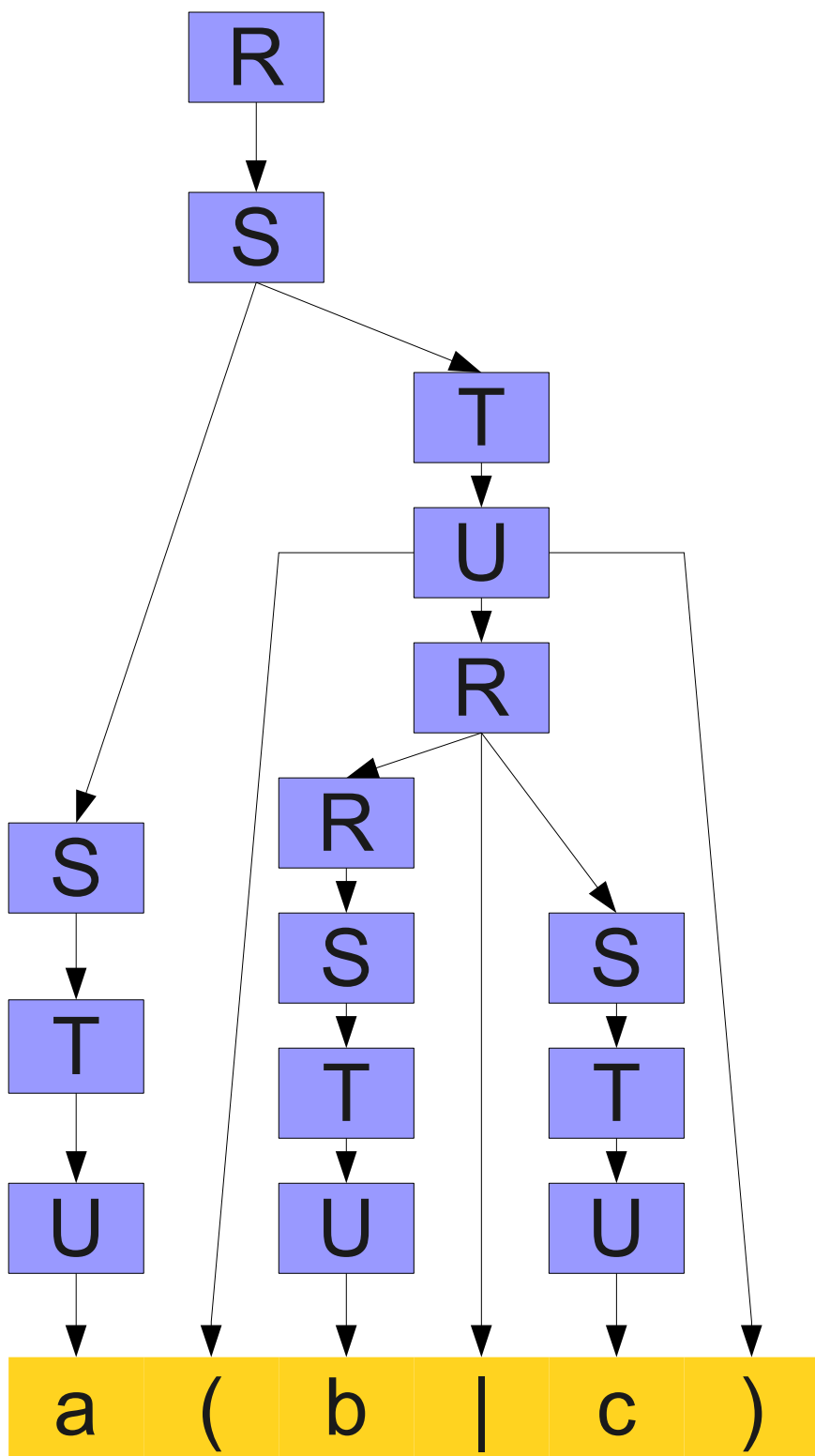
$U \rightarrow "a"$

$U \rightarrow "\epsilon"$

$U \rightarrow (R)$

a (b | c)

$R \rightarrow S \mid R \mid S$
 $S \rightarrow T \mid ST$
 $T \rightarrow U \mid T^*$
 $U \rightarrow \text{"a"}$
 $U \rightarrow \text{"\epsilon"}$
 $U \rightarrow (R)$



Abstract Syntax Trees (ASTs)

- Tree structure encoding the **logical structure** of a piece of code.
 - Abstracts from the particular syntax being used.
- Eliminate unnecessary grammar symbols.
- Eliminate nodes from grammar internals.
- Ultimate output of syntax analysis.

How to build an AST?

- Typically done through **semantic actions**.
- Idea: **Associate code with each production**.
- As the input is parsed, execute this code to build the AST.
 - Exact order of code execution depends on the parsing method used.
- This is called a **syntax-directed translation**.

Simple Semantic Actions

$E \rightarrow T + E$ $E_1.val = T.val + E_2.val$

$E \rightarrow T$ $E.val = T.val$

$T \rightarrow \mathbf{int}$ $T.val = \mathbf{int}.val$

$T \rightarrow \mathbf{int} * T$ $T.val = \mathbf{int}.val * T.val$

$T \rightarrow (E)$ $T.val = E.val$

Simple Semantic Actions

$E \rightarrow T + E$ $E_1.val = T.val + E_2.val$

$E \rightarrow T$ $E.val = T.val$

$T \rightarrow \mathbf{int}$ $T.val = \mathbf{int}.val$

$T \rightarrow \mathbf{int} * T$ $T.val = \mathbf{int}.val * T.val$

$T \rightarrow (E)$ $T.val = E.val$



Simple Semantic Actions

$E \rightarrow T + E$

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

$E \rightarrow T$

$E.val = T.val$

$T \rightarrow \mathbf{int}$

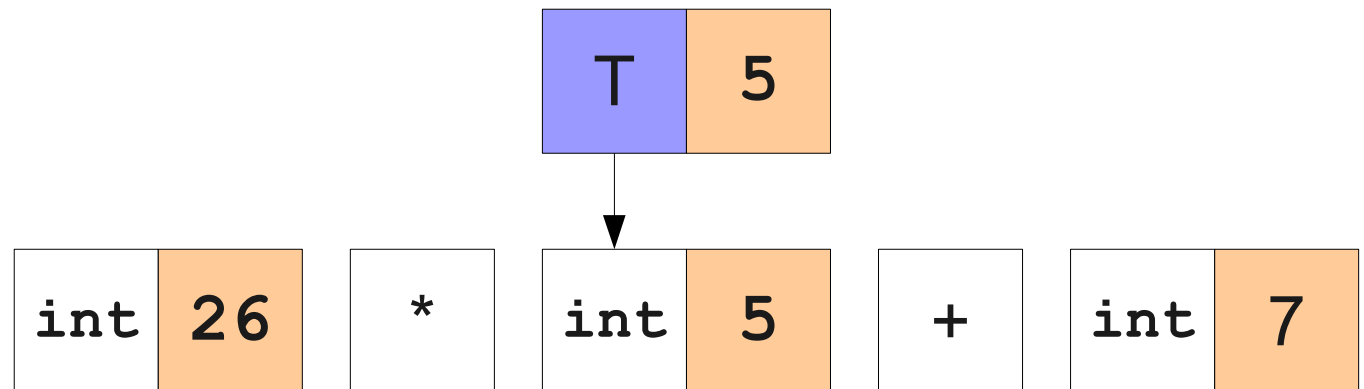
$T.val = \mathbf{int}.val$

$T \rightarrow \mathbf{int} * T$

$T.val = \mathbf{int}.val * T.val$

$T \rightarrow (E)$

$T.val = E.val$



Simple Semantic Actions

$E \rightarrow T + E$

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

$E \rightarrow T$

$E.val = T.val$

$T \rightarrow \text{int}$

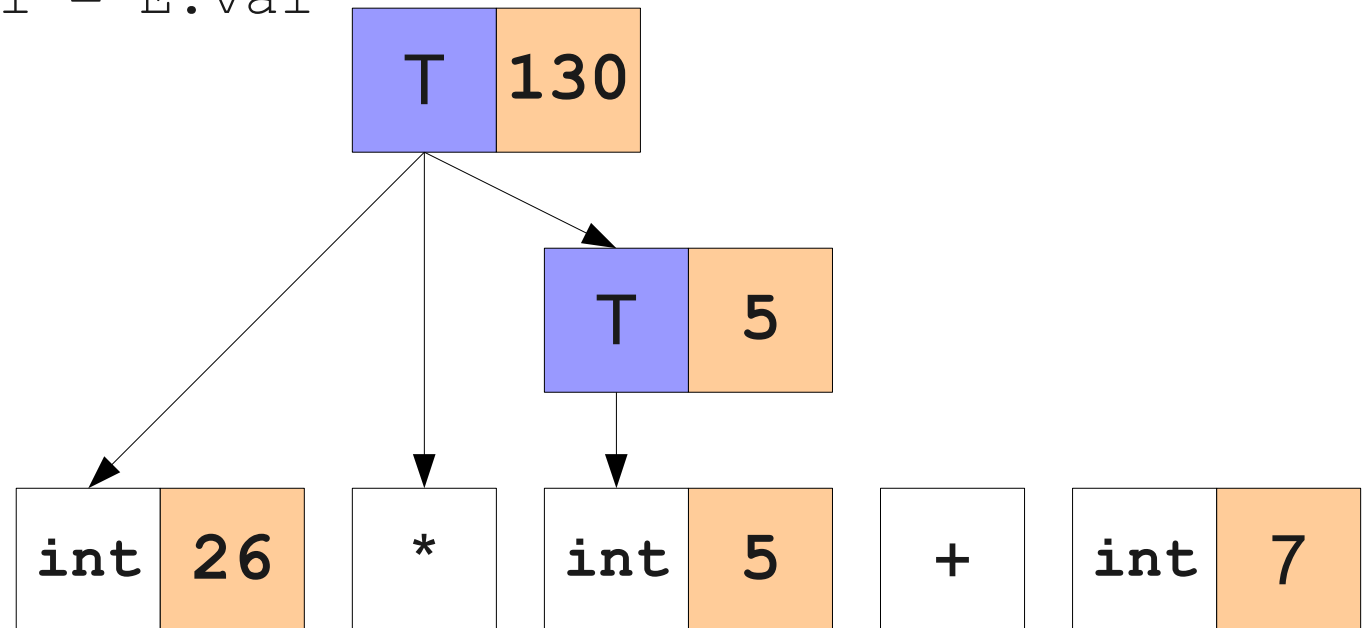
$T.val = \text{int}.val$

$T \rightarrow \text{int} * T$

$T.val = \text{int}.val * T.val$

$T \rightarrow (E)$

$T.val = E.val$



Simple Semantic Actions

$E \rightarrow T + E$

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

$E \rightarrow T$

$E.val = T.val$

$T \rightarrow \text{int}$

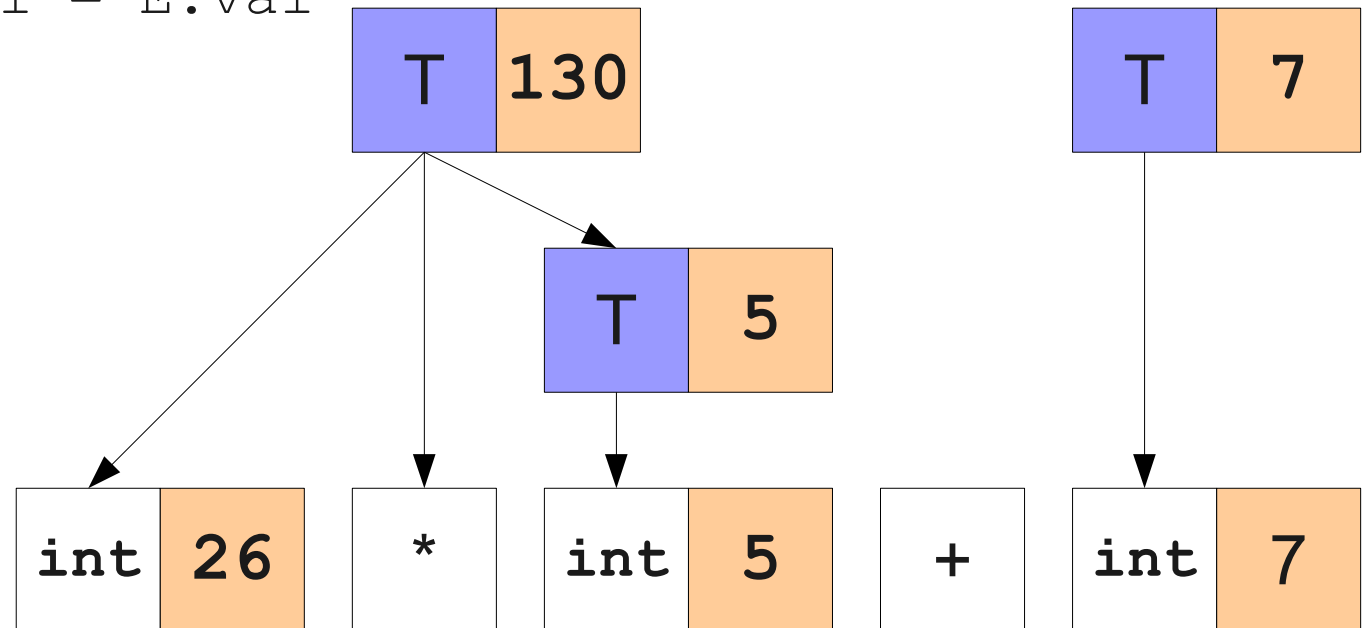
$T.val = \text{int}.val$

$T \rightarrow \text{int} * T$

$T.val = \text{int}.val * T.val$

$T \rightarrow (E)$

$T.val = E.val$



Simple Semantic Actions

$E \rightarrow T + E$

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

$E \rightarrow T$

$E.val = T.val$

$T \rightarrow \text{int}$

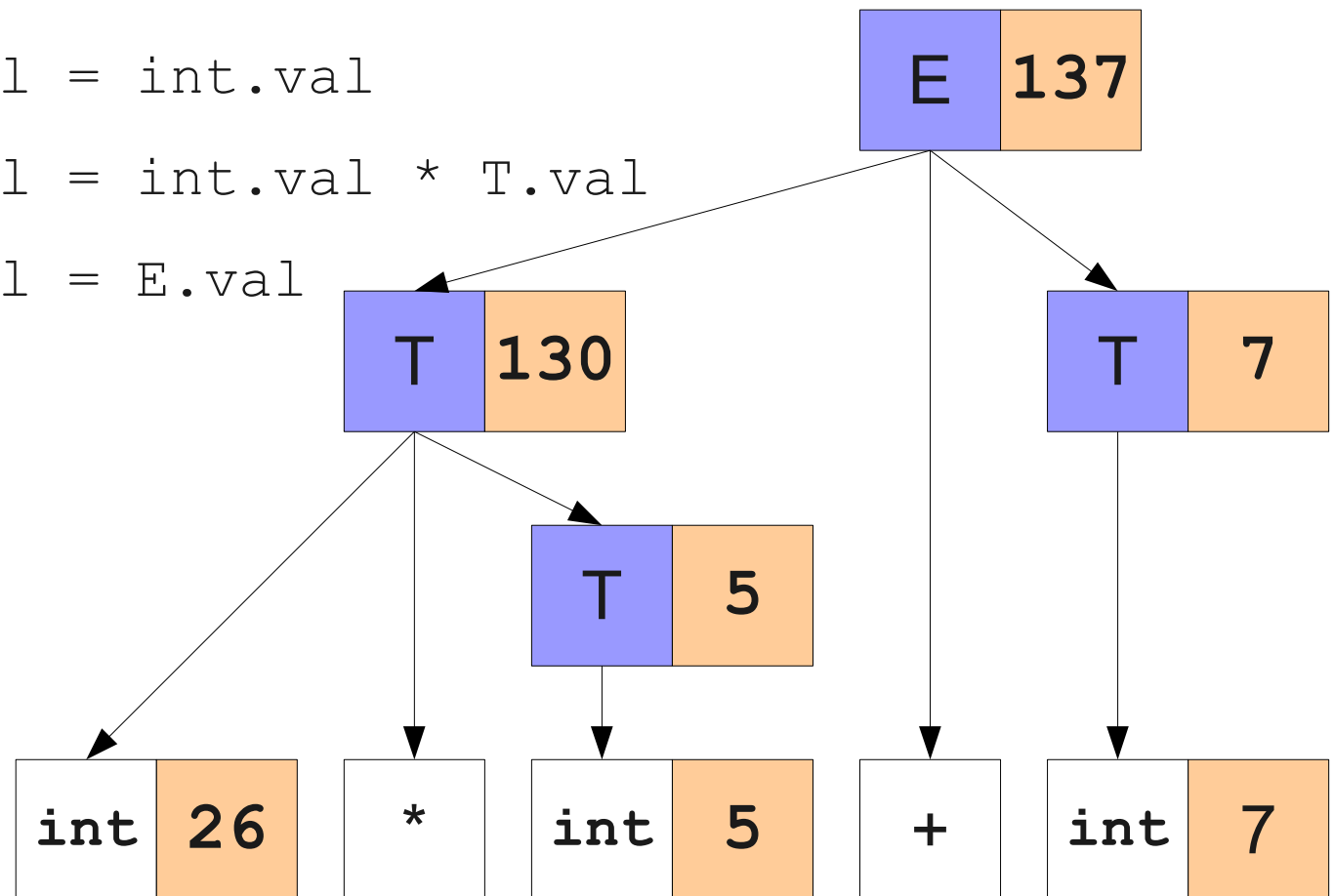
$T.val = \text{int}.val$

$T \rightarrow \text{int} * T$

$T.val = \text{int}.val * T.val$

$T \rightarrow (E)$

$T.val = E.val$



Semantic Actions to Build ASTs

```
R → S      R.ast = S.ast;  
R → R "|" S  R1.ast = new Or(R2.ast, S.ast);  
S → T      S.ast = T.ast;  
S → ST     S1.ast = new Concat(S2.ast, T.ast);  
T → U      T.ast = U.ast;  
T → T*     T1.ast = new Star(T2.ast);  
U → "a"    U.ast = SingleChar('a');  
U → "ε"    U.ast = new Epsilon();  
U → (R)    U.ast = R.ast;
```

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
 - Recursive Descent
 - LL(1)