

بسم الله الرحمن الرحيم

Parsing:

Context-Free Grammars, Parsing and Programming Languages (2)

Comparison with Lexical Analysis

| Phase | Input | Output |
|--------|----------------------|---------------------------------|
| Lexer | String of characters | String of tokens |
| Parser | String of tokens | Parse tree (may be implicit) |

Example

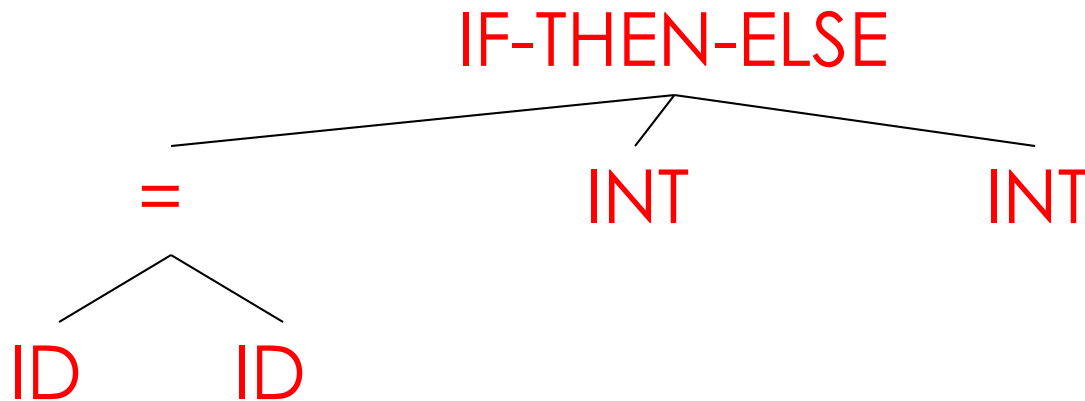
- Cool

if x = y then 1 else 2 fi

- Parser input

IF ID = ID THEN INT ELSE INT FI

- Parser output

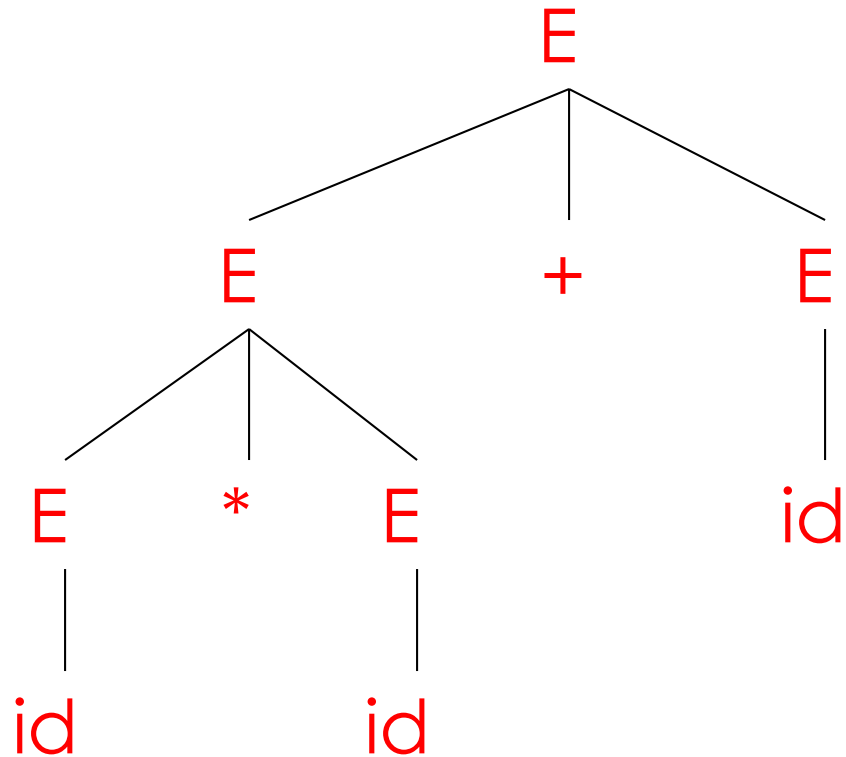


Ambiguity

- Grammar $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- String $id * id + id$

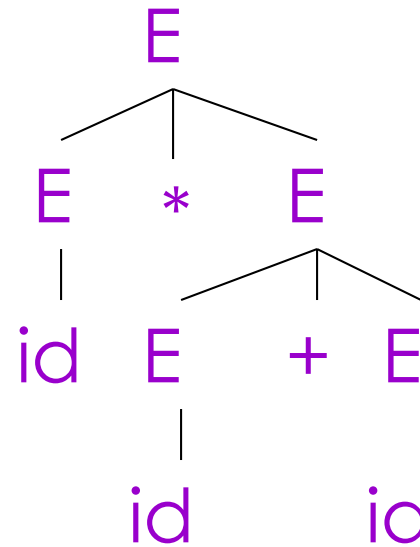
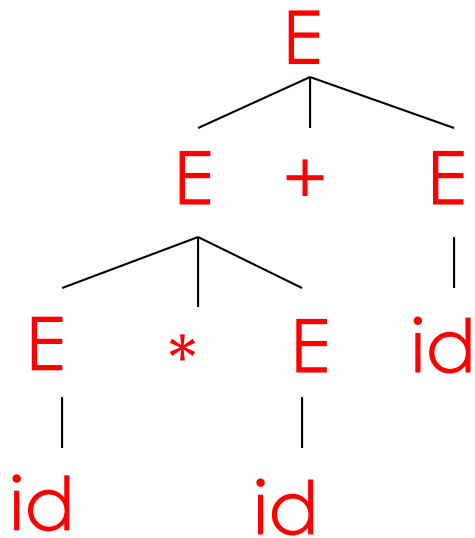
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$
 $\rightarrow id * id + id$



Ambiguity (Cont.)

This string has two parse trees

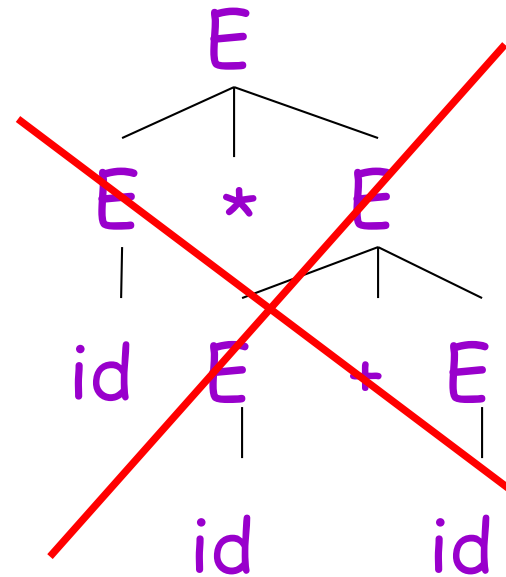
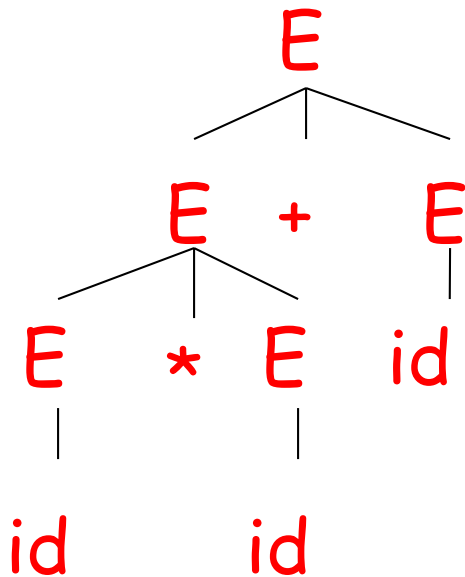


Ambiguity in Arithmetic Expressions

- Recall the grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

- The string $id * id + id$ has two parse trees:



Dealing with Ambiguity

- There are several ways to handle ambiguity
- Most direct method is to rewrite grammar unambiguously

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \text{id}$$

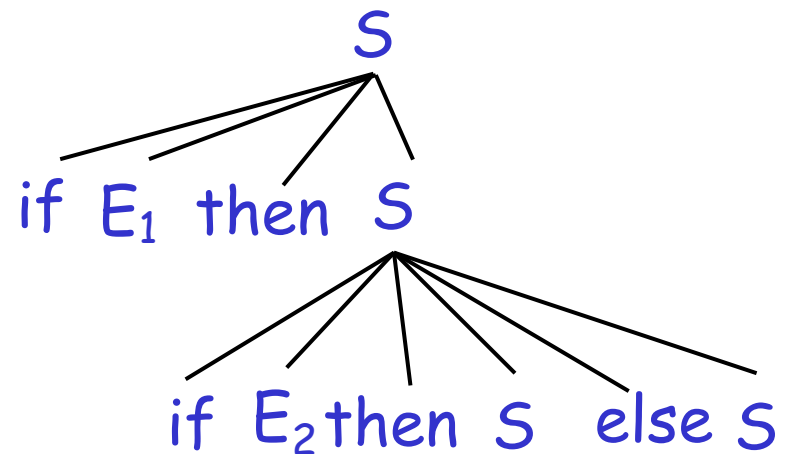
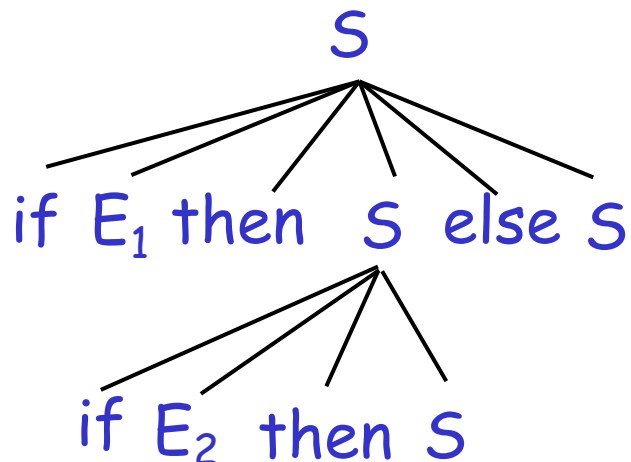
- Enforces precedence of $*$ over $+$

Ambiguity: The Dangling Else

- Consider the grammar
$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S \\ &\quad | \text{if } E \text{ then } S \text{ else } S \\ &\quad | \text{OTHER} \end{aligned}$$
- This grammar is also ambiguous

The Dangling Else: Example

- The expression
if E_1 then if E_2 then S else S
has two parse trees



- Typically we want the second form

The Dangling Else: A Fix

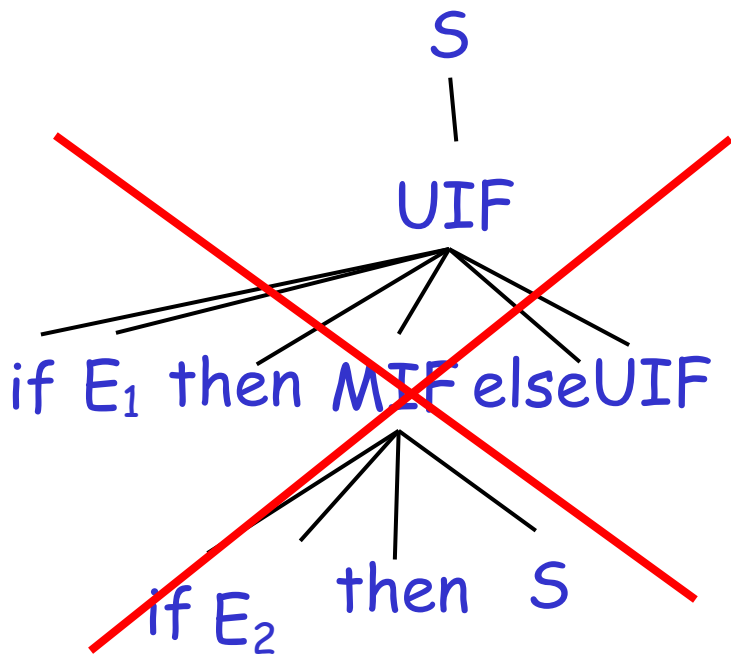
- `else` matches the closest unmatched `then`
- We can describe this in the grammar

```
S → MIF          /* all then are matched */  
   | UIF          /* some then is unmatched */  
MIF → if E then MIF else MIF  
     | OTHER  
UIF → if E then S  
     | if E then MIF else UIF
```

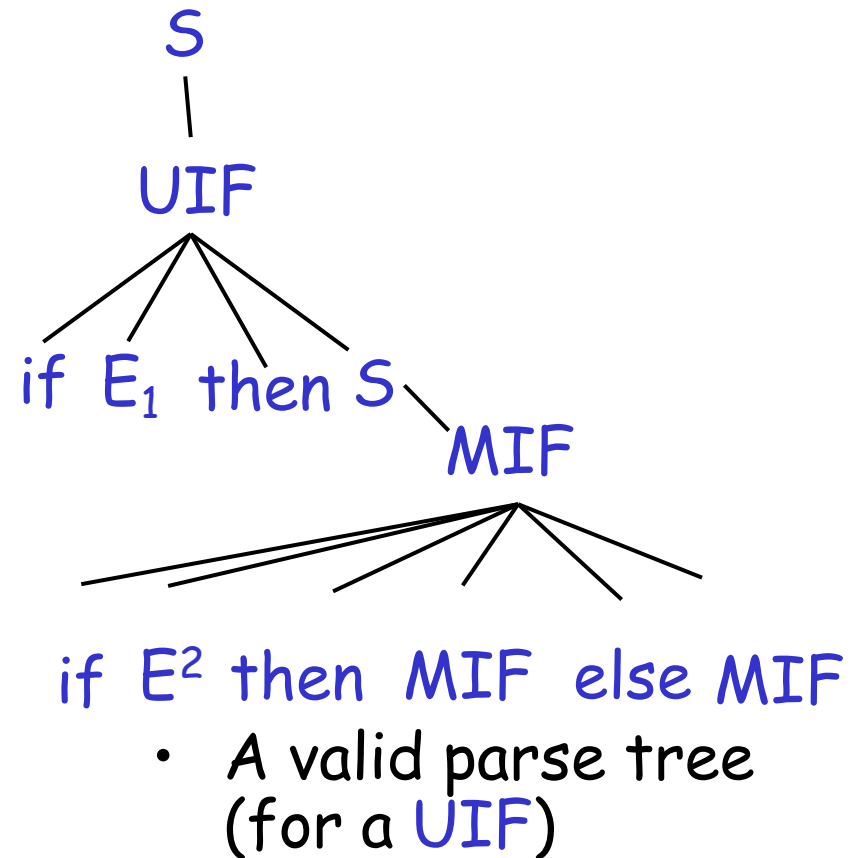
- Describes the same set of strings

The Dangling Else: Example Revisited

- The expression $\text{if } E_1 \text{ then if } E_2 \text{ then } S \text{ else } S$



- Not valid because the then expression is not a MIF



Ambiguity

- No general techniques for handling ambiguity
- Impossible to convert automatically an ambiguous grammar to an unambiguous one
- Used with care, ambiguity can simplify the grammar
 - Sometimes allows more natural definitions
 - We need disambiguation mechanisms

Resolving Ambiguity

- If a grammar can be made unambiguous, 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:
 - ϵ
 - $()$
 - $((())())$
 - $(((())))((()))()$
 -
- Here is one possible grammar for balanced parentheses:

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

Balanced Parentheses

- Given the grammar $\mathbf{P} \rightarrow \epsilon \mid \mathbf{PP} \mid (\mathbf{P})$
- How might we generate the string $((())())$?

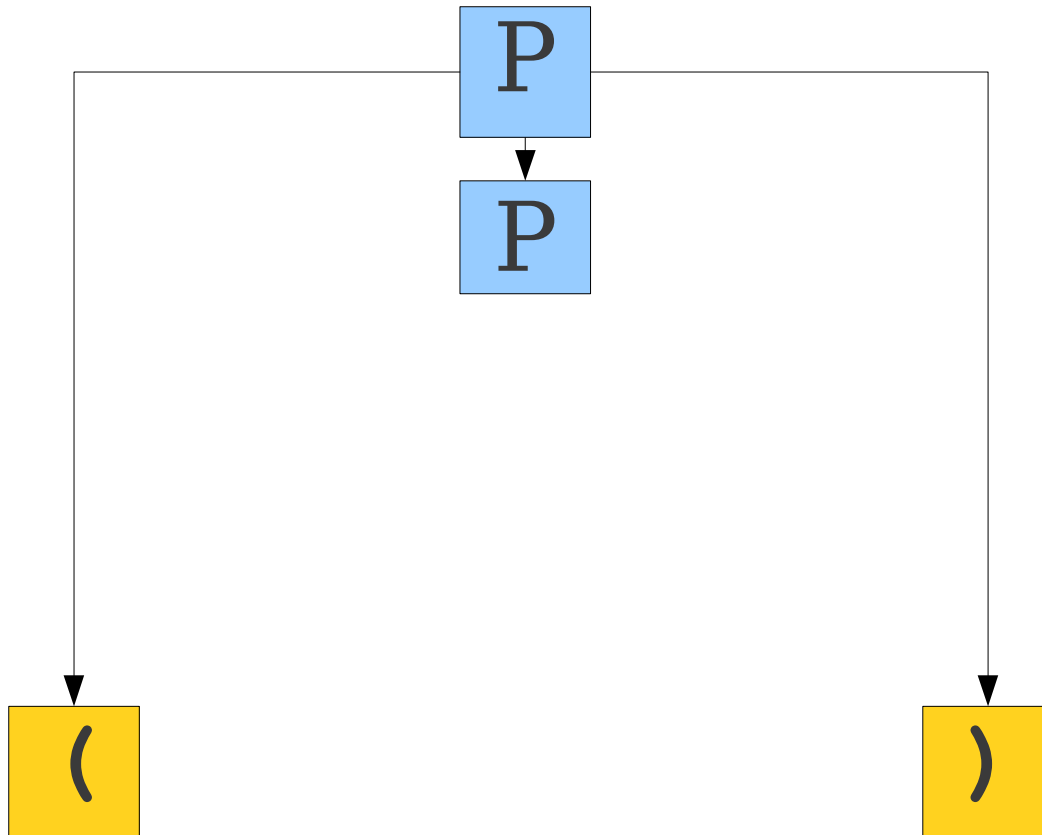
Balanced Parentheses

- Given the grammar $\mathbf{P} \rightarrow \epsilon \mid \mathbf{PP} \mid (\mathbf{P})$
- How might we generate the string $((())())$?

P

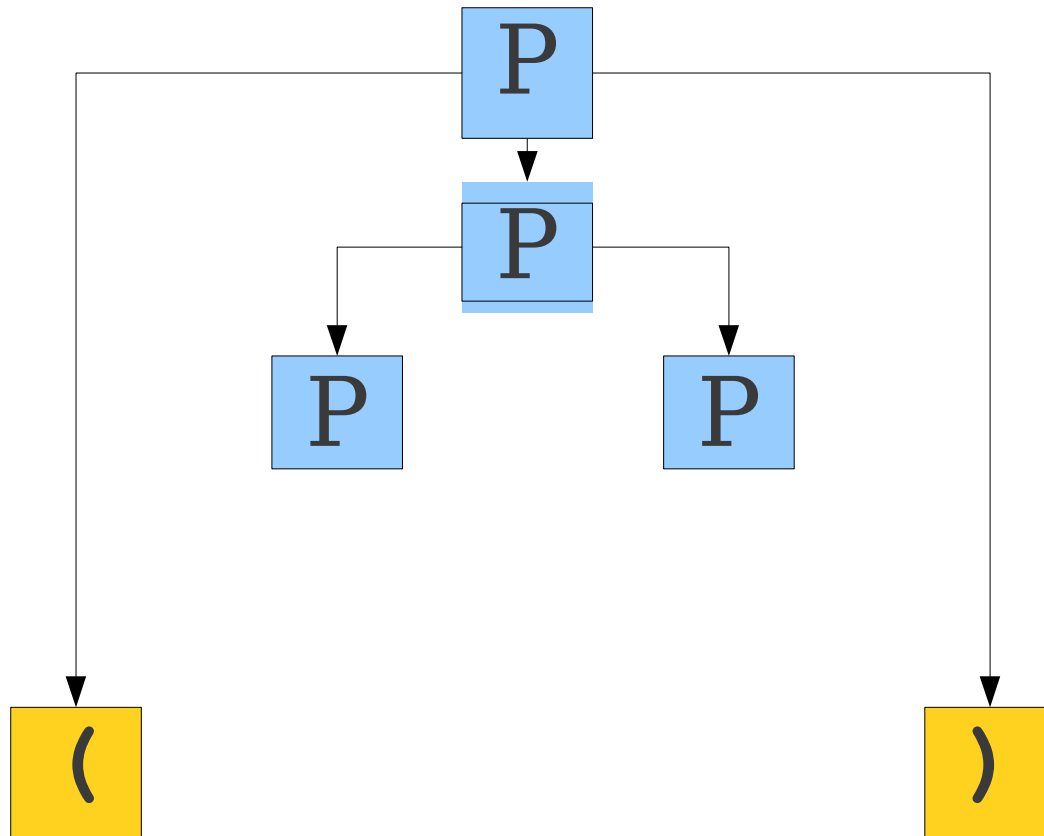
Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((())())$?



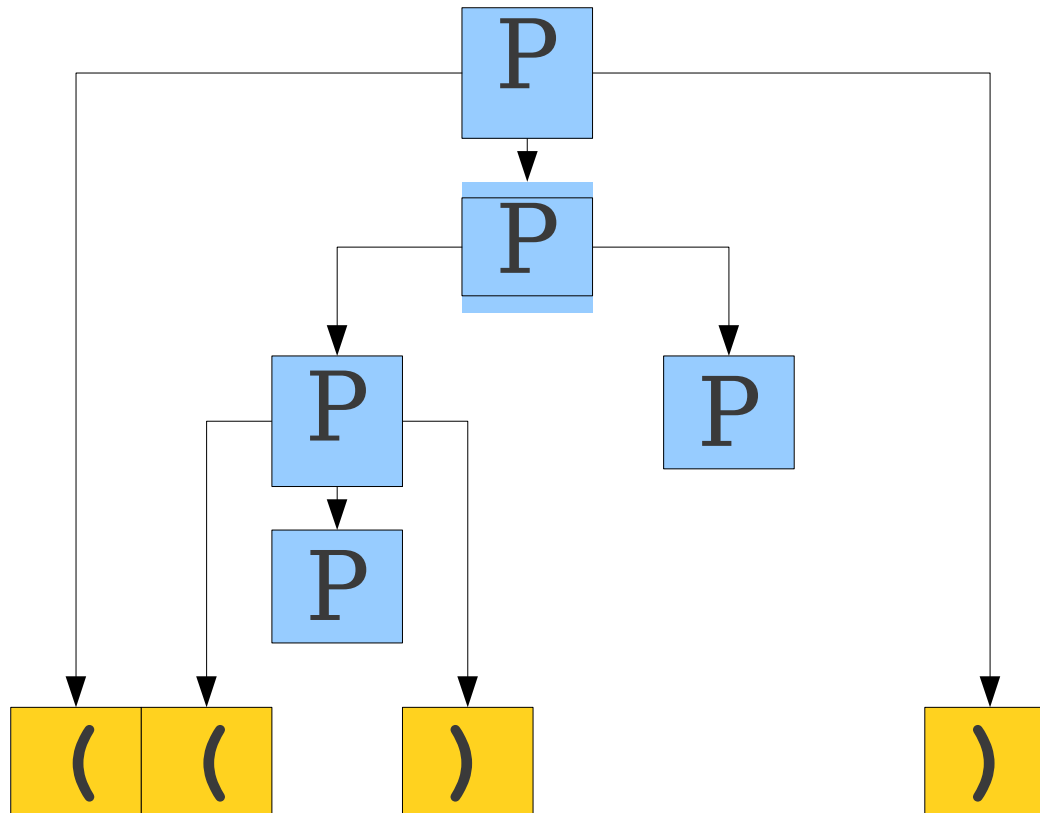
Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((()))$?



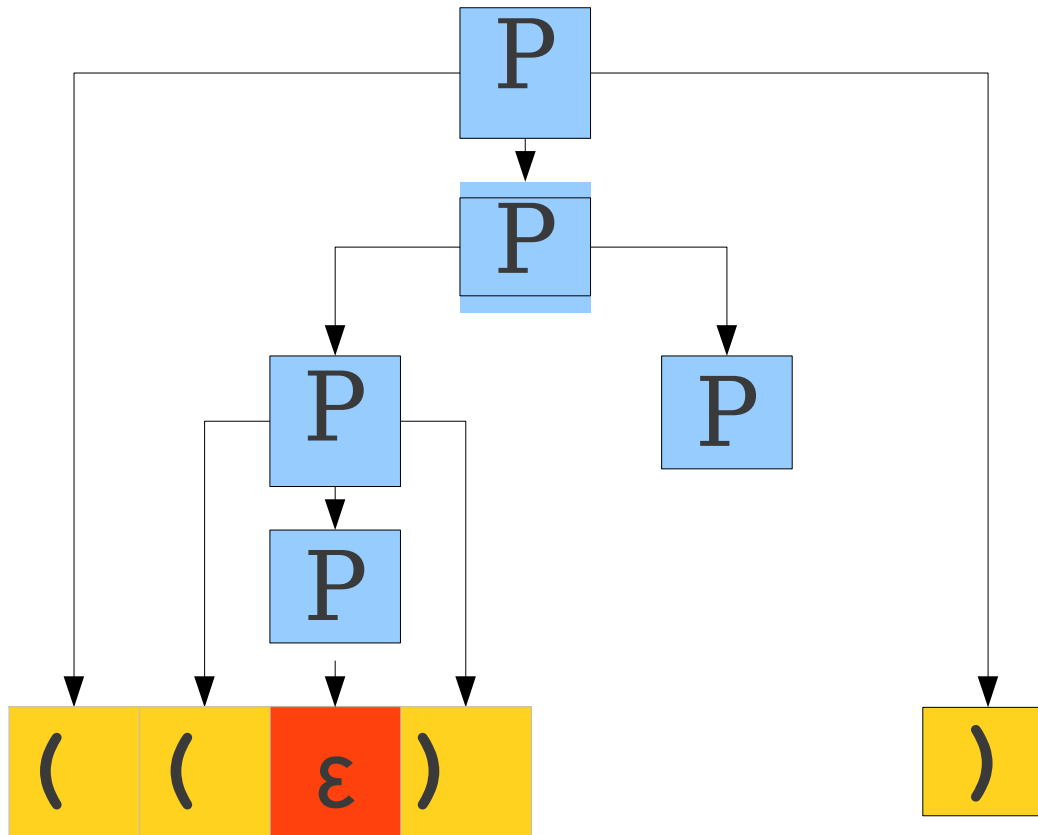
Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((()))$?



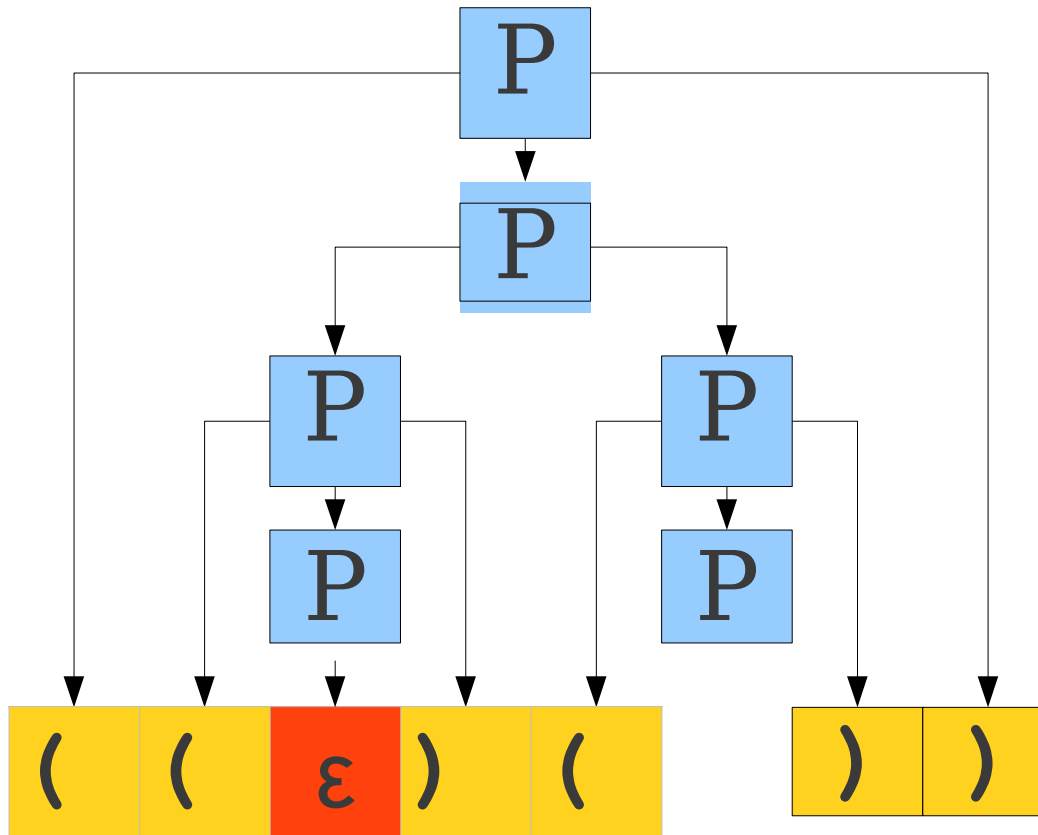
Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((\epsilon))$?



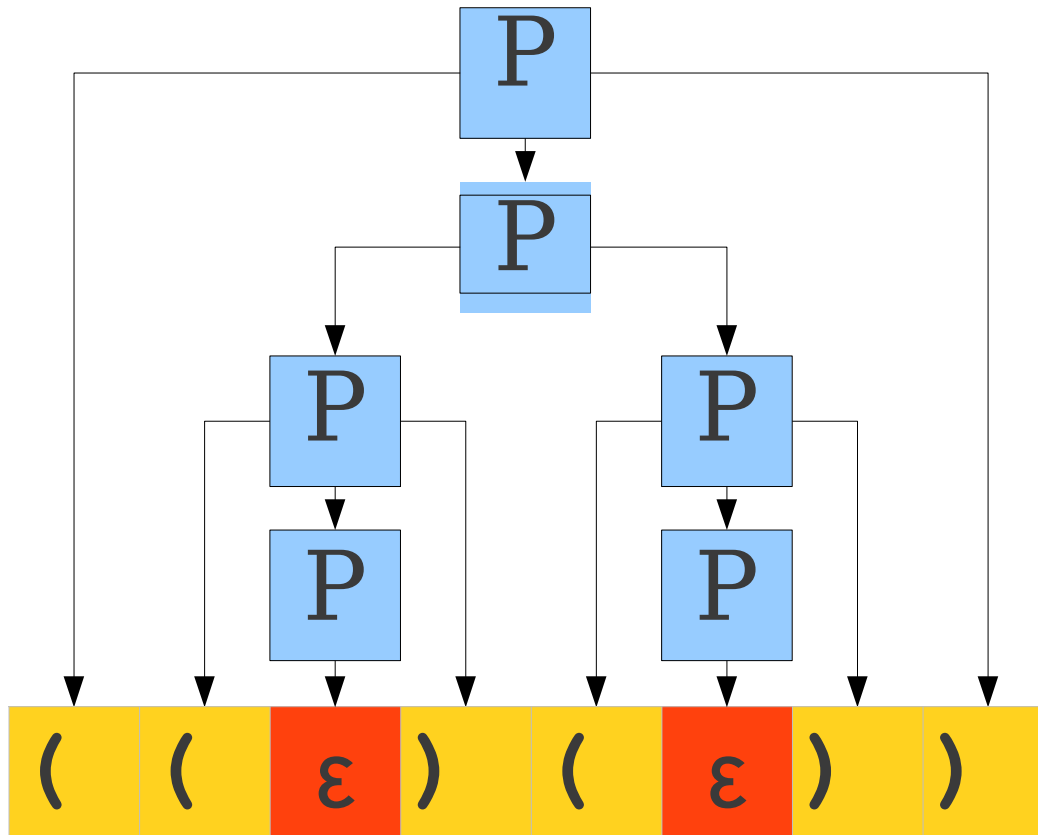
Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((\epsilon))$?

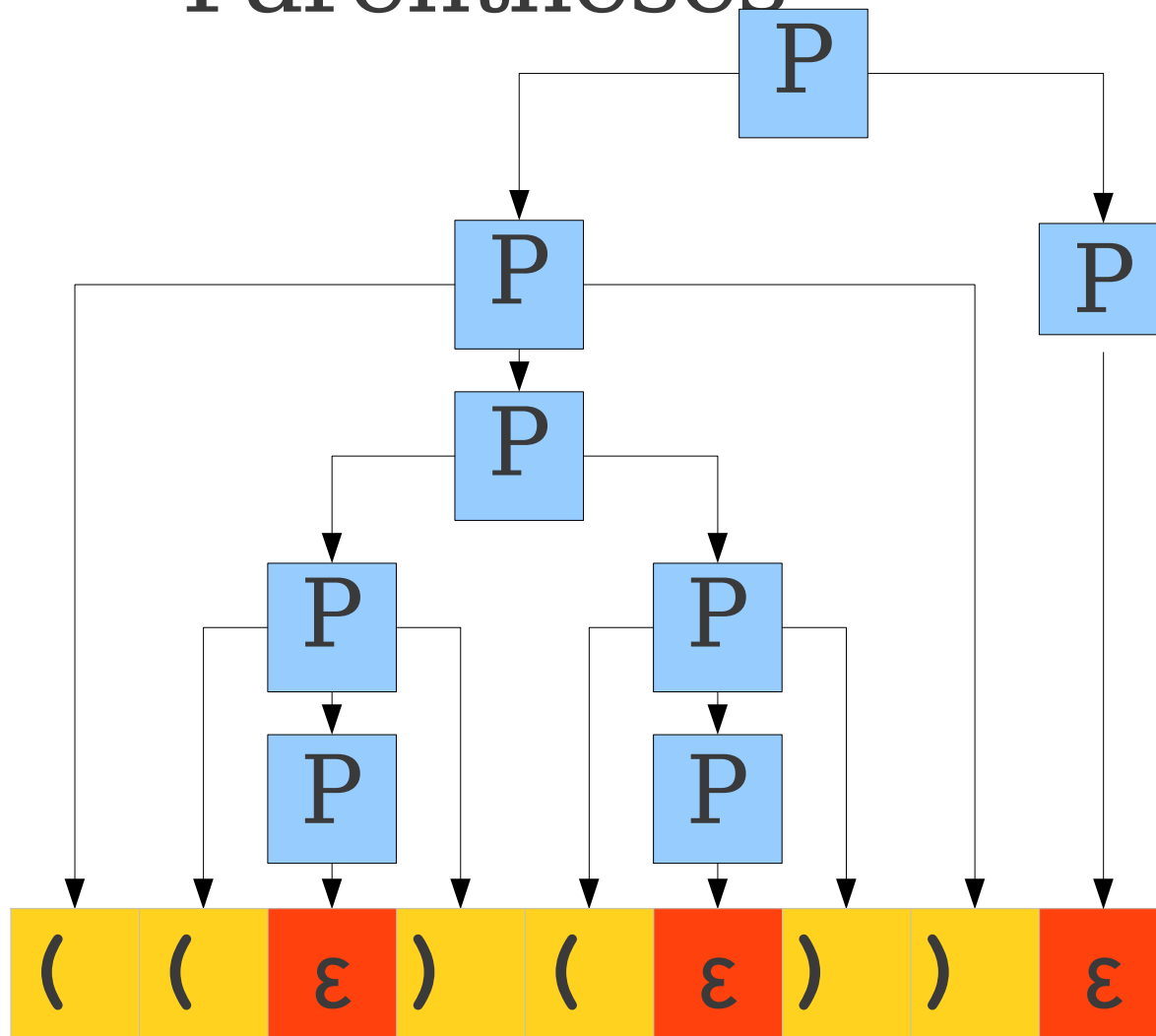


Balanced Parentheses

- Given the grammar $P \rightarrow \epsilon \mid PP \mid (P)$
- How might we generate the string $((\epsilon)(\epsilon))$?



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.

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** | ϵ

P \rightarrow (**S**)

Building Parentheses

S \rightarrow **P** **S** | ε

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
 - ε
 - 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:

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

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

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

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

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

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

An Ambiguous Grammar

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

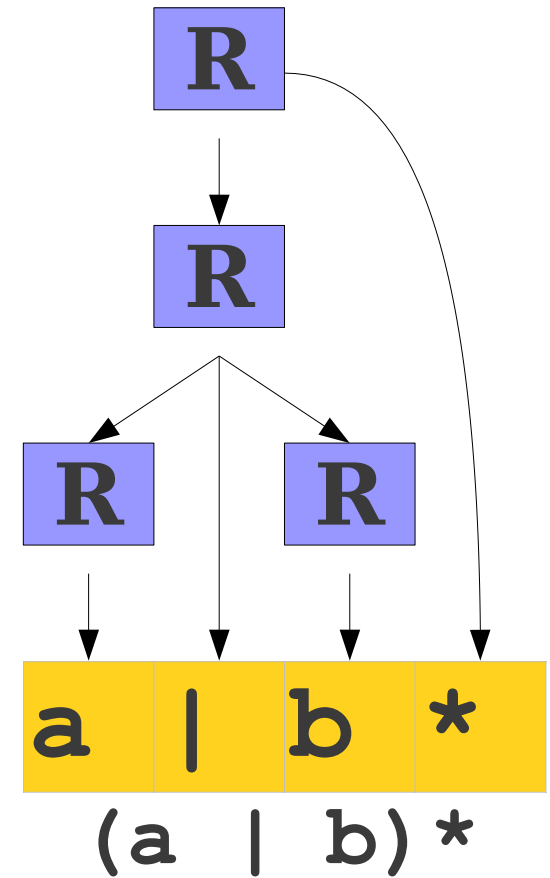
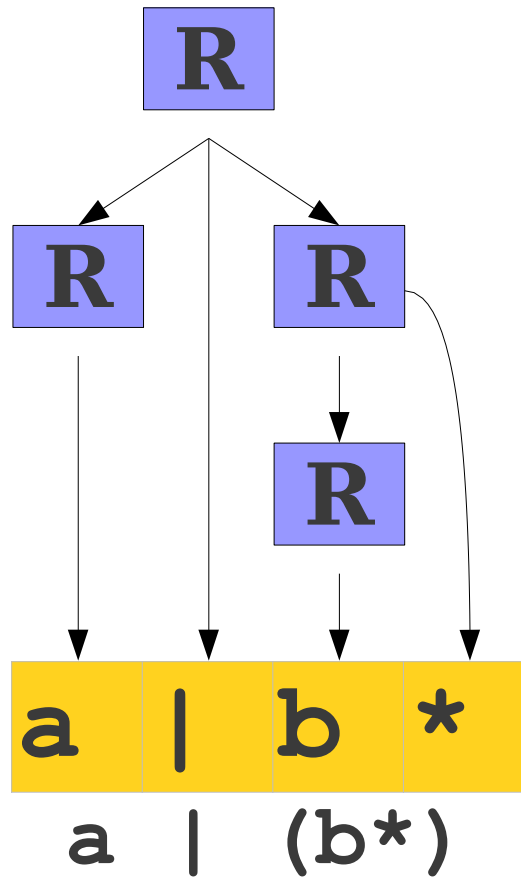
$R \rightarrow \epsilon$

$R \rightarrow RR$

$R \rightarrow R \mid R$

$R \rightarrow R^*$

$R \rightarrow (R)$



Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

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

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

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

$$\mathbf{R} \rightarrow \mathbf{R} \mid \mathbf{R}$$

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

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

| | | | | |
|---|---|--|---|---|
| a | a | | b | * |
|---|---|--|---|---|

Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

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

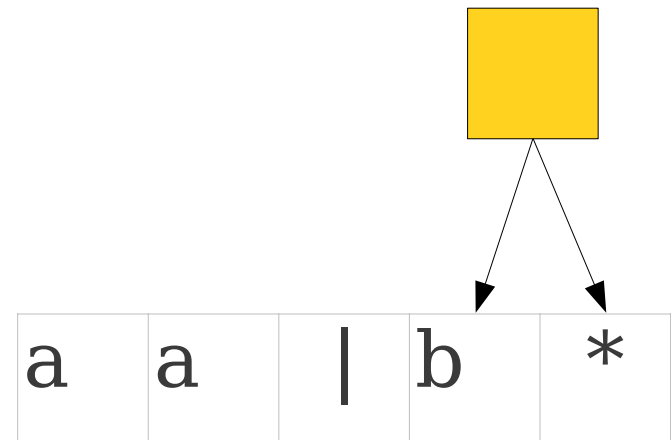
R \rightarrow " **ϵ** "

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** | **b** | **c** | ...

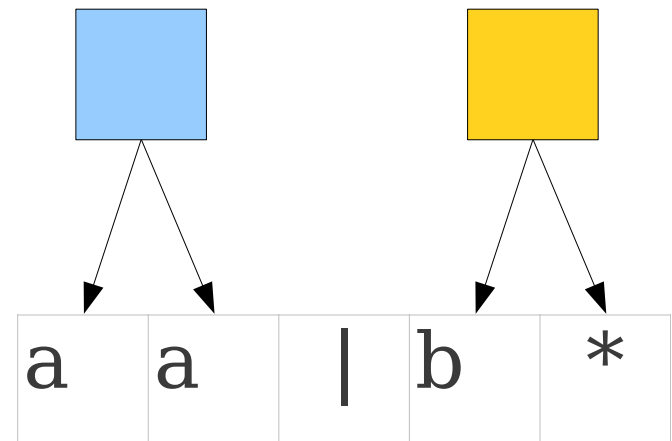
R \rightarrow " **ϵ** "

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** | **b** | **c** | ...

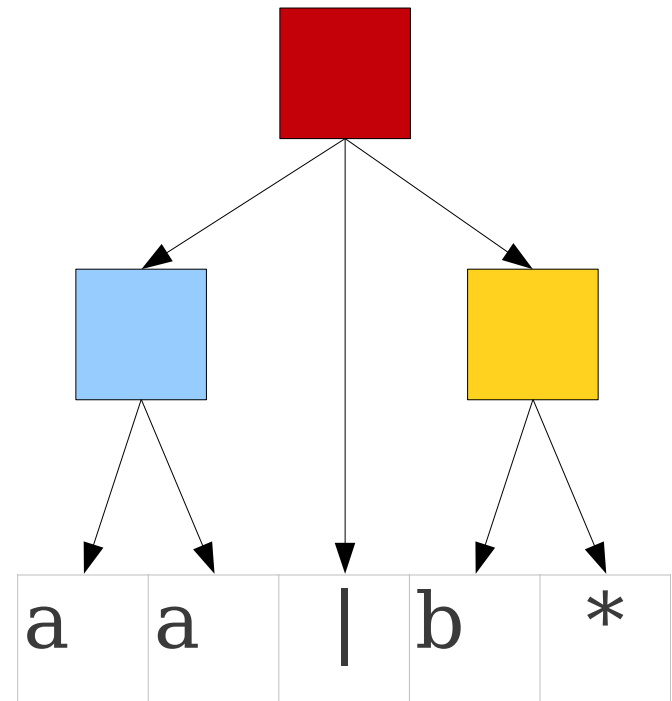
R \rightarrow " **ϵ** "

R \rightarrow **RR**

R \rightarrow **R** "**|**" **R**

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

R \rightarrow (**R**)



Resolving Ambiguity

- We can try to resolve the ambiguity via layering:

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

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

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

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

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

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

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

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

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

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

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

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

Why is this unambiguous?

$$R \rightarrow S \mid R \text{ “|” } S$$
$$S \rightarrow T \mid ST$$
$$T \rightarrow U \mid T^*$$
$$U \rightarrow a \mid b \mid \dots$$
$$U \rightarrow \text{“}\epsilon\text{”}$$
$$U \rightarrow (R)$$

Why is this unambiguous?

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

$S \rightarrow T \mid ST$

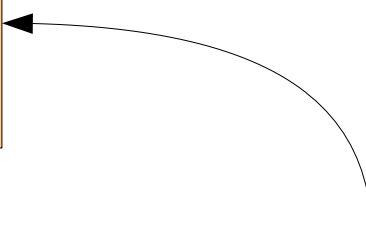
$T \rightarrow U \mid T^*$

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

$U \rightarrow \text{“}\epsilon\text{”}$

$U \rightarrow (R)$

Only generates
“atomic” expressions



Why is this unambiguous?

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

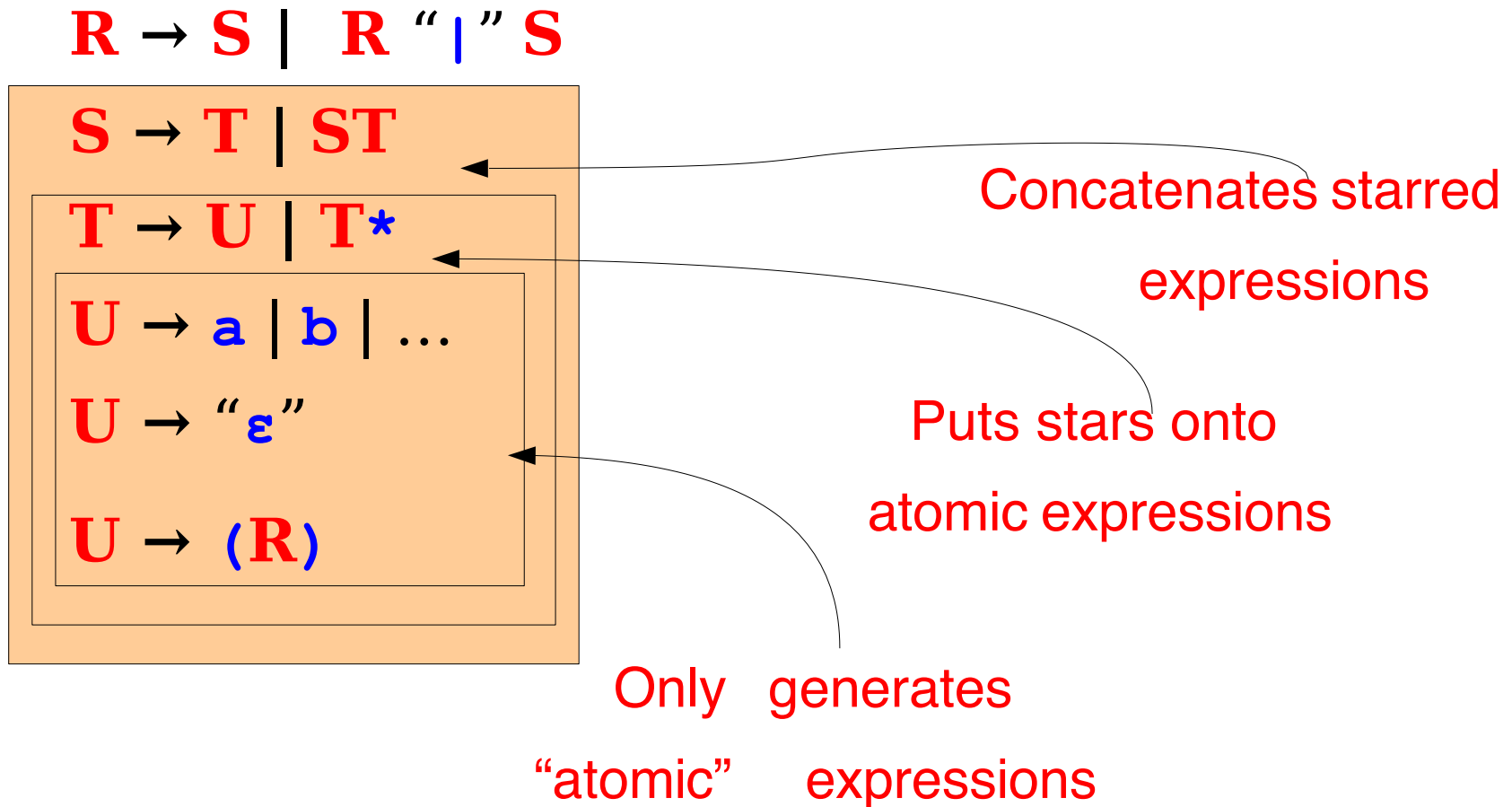
$U \rightarrow \epsilon$

$U \rightarrow (R)$

Puts stars onto
atomic expressions

Only generates
“atomic” expressions

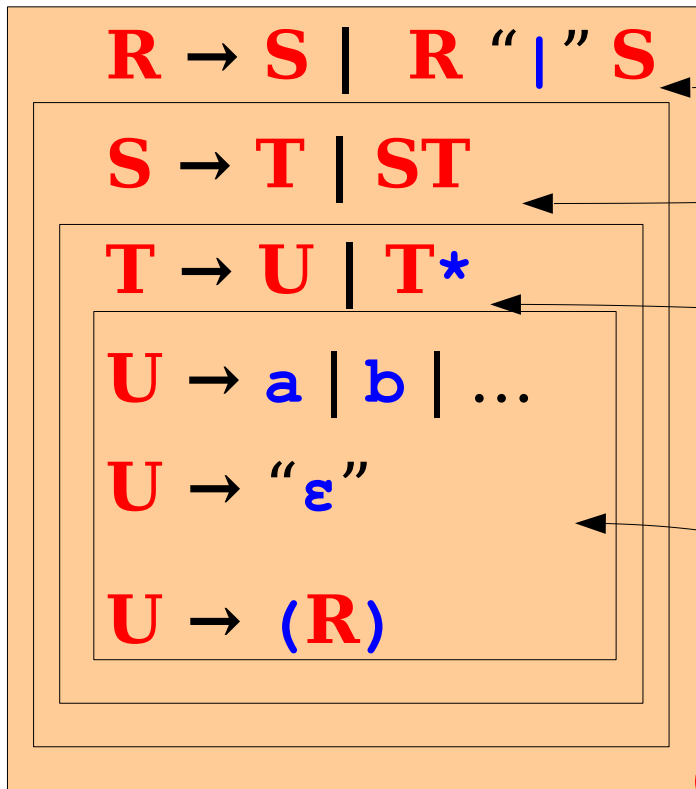
Why is this unambiguous?



Why is this unambiguous?

Unions

concatenated
expressions



Concatenates starred
expressions

Puts stars onto
atomic expressions

Only generates
“atomic” expressions

R

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \text{“}\epsilon\text{”}$

$U \rightarrow (R)$

a

b

|

c

|

a

*

R \rightarrow **S** | **R** “|” **S**

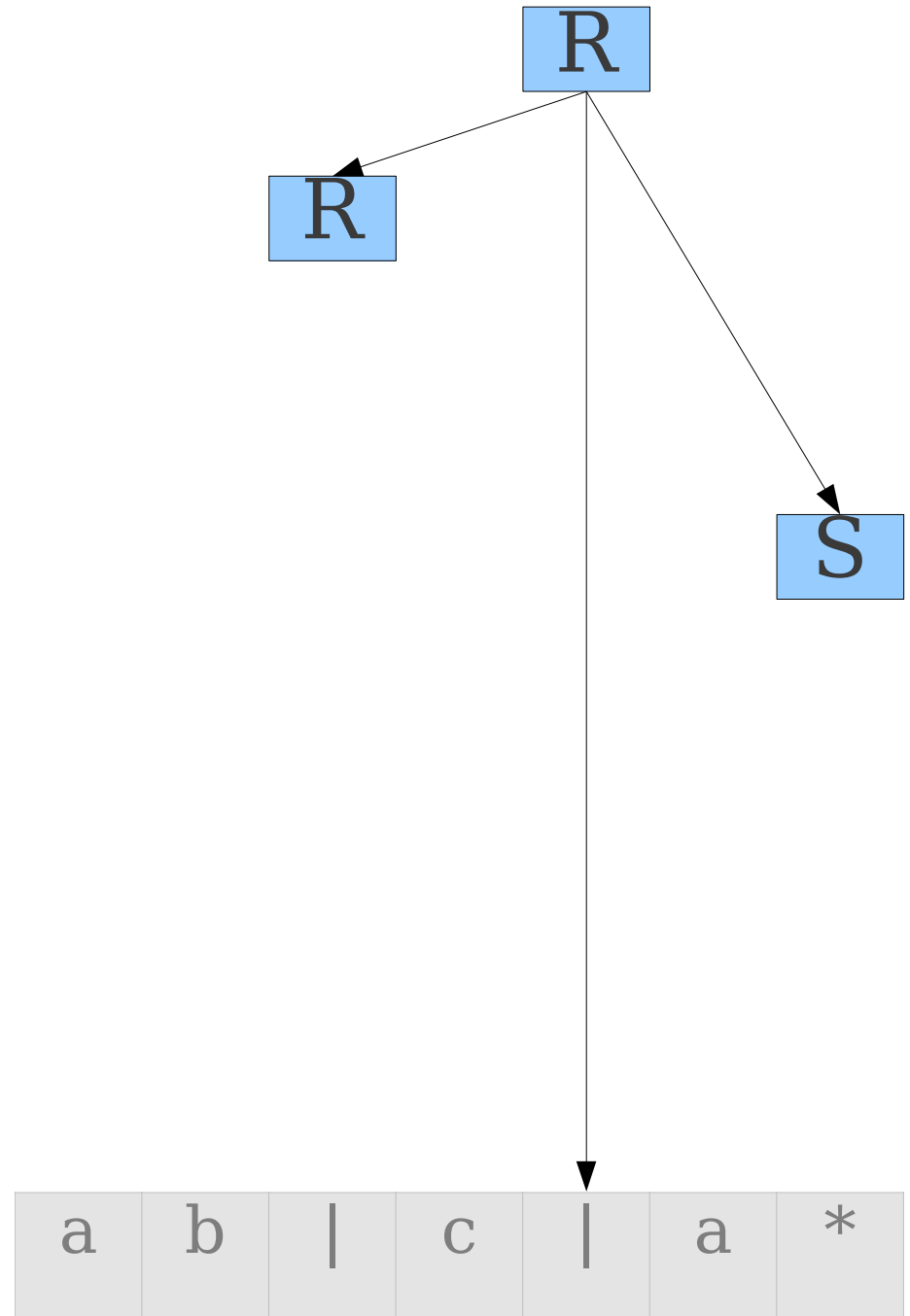
S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

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

U \rightarrow “ **ϵ** ”

U \rightarrow (**R**)



R \rightarrow **S** | **R** " | " **S**

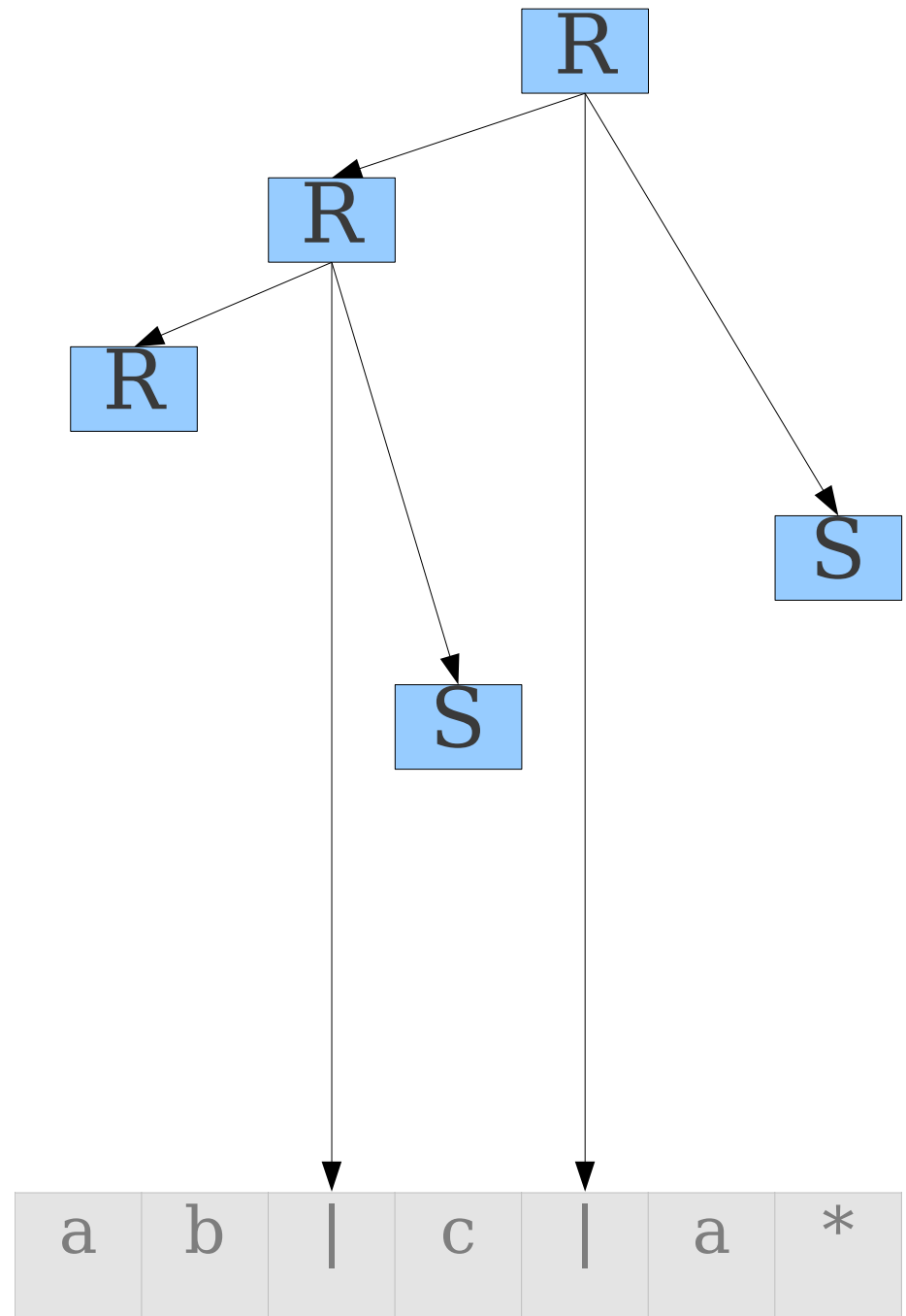
S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

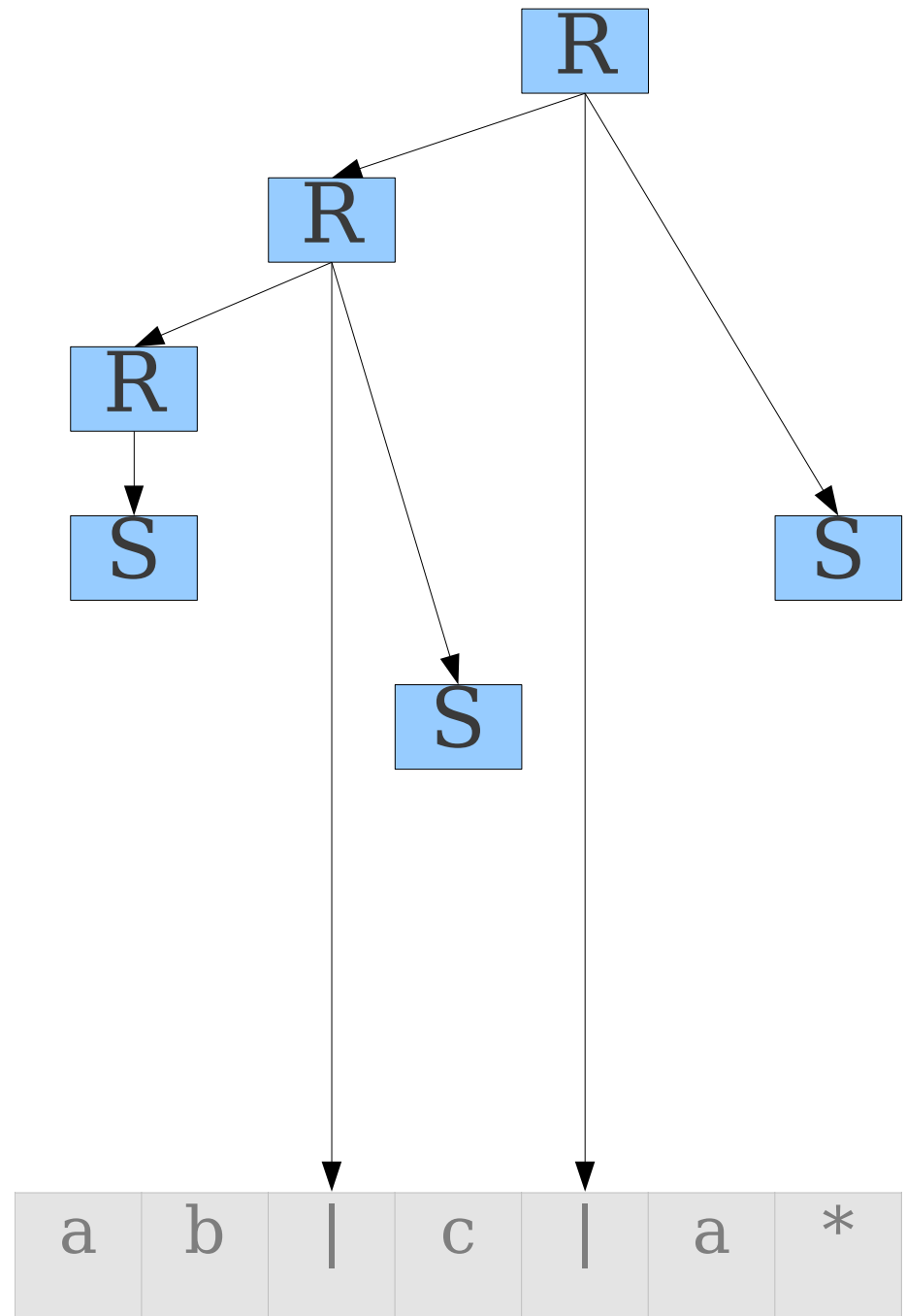
U \rightarrow **a** | **b** | **c** | ...

U \rightarrow " **ϵ** "

U \rightarrow (**R**)



U \rightarrow **(R)**



R \rightarrow **S** | **R** " | " **S**

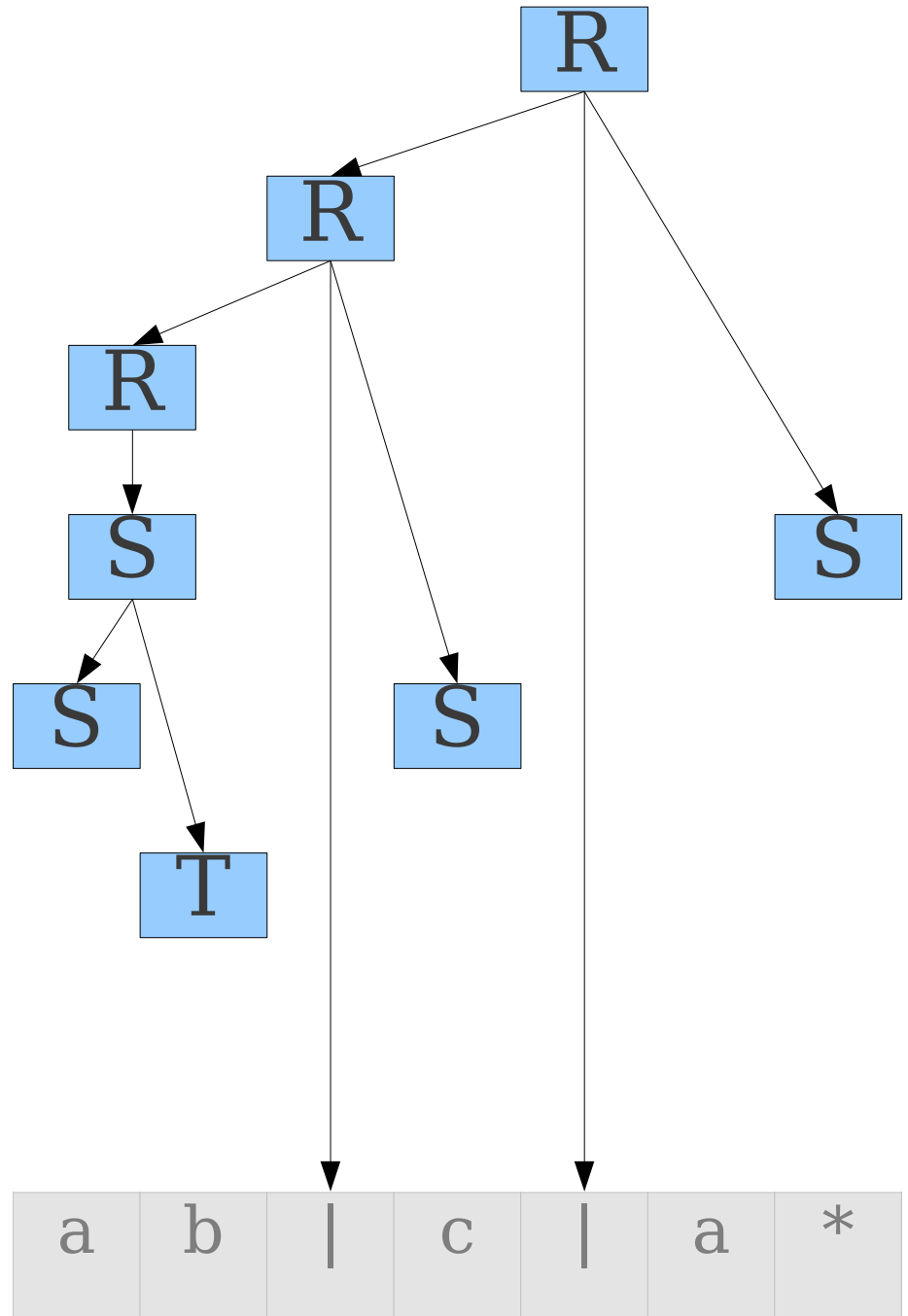
S \rightarrow **T** | **ST**

T → **U** | **T***

U \rightarrow a | b | c | ...

U → “ε”

U \rightarrow **(R)**



R \rightarrow **S** | **R** " | " **S**

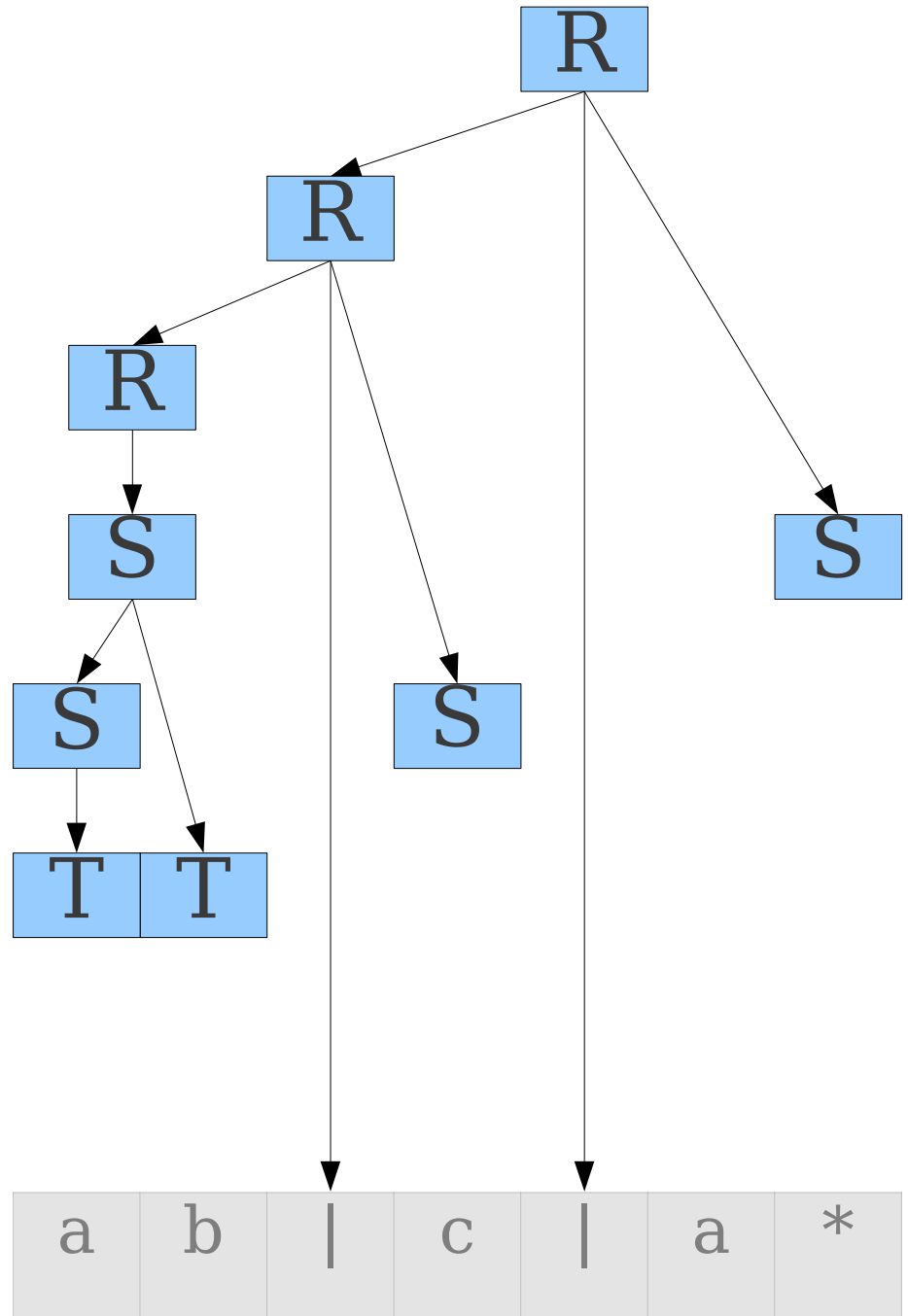
S \rightarrow **T** | **ST**

T → **U** | **T***

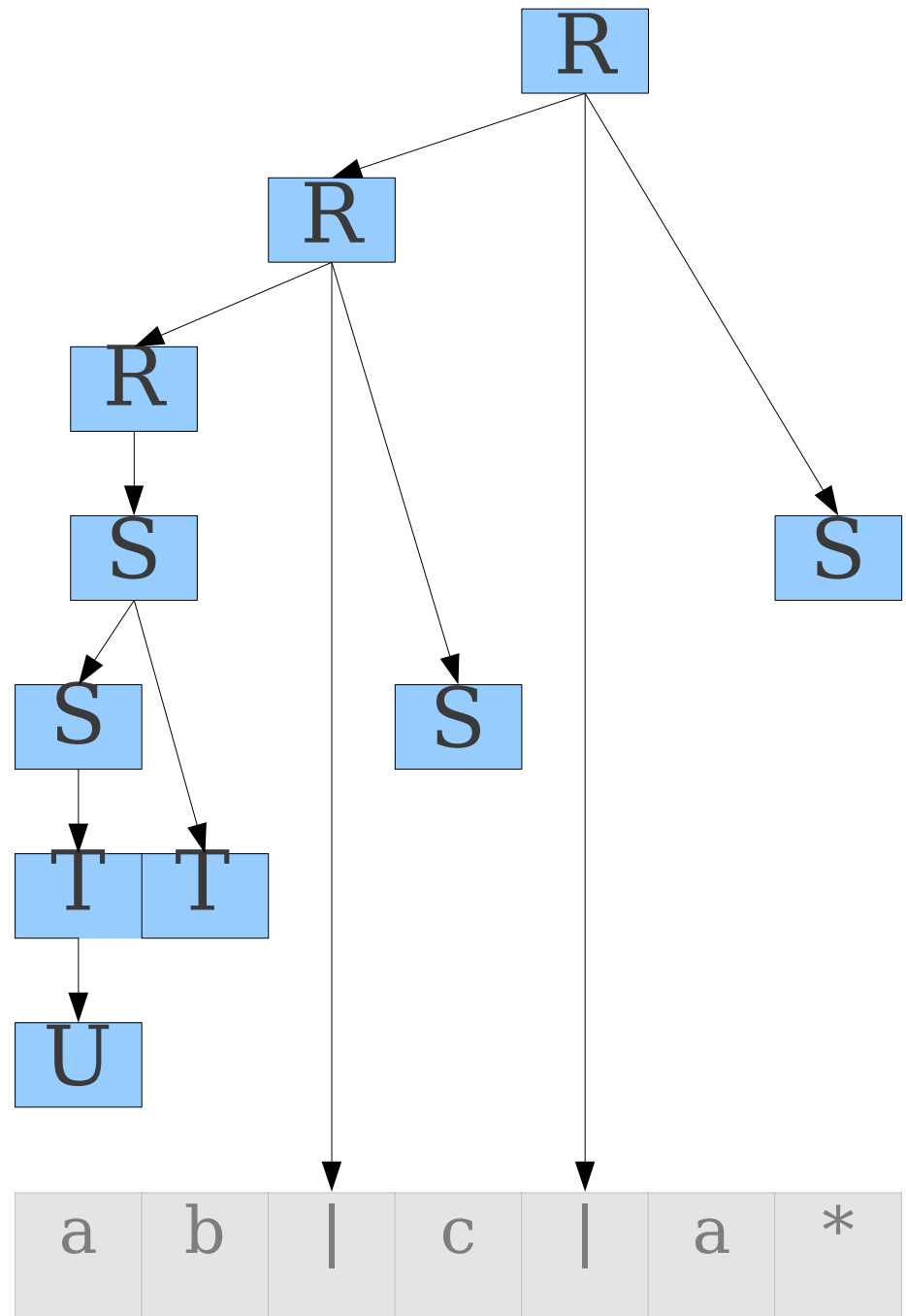
U \rightarrow **a** | **b** | **c** | ...

U → “**ε**”

U \rightarrow **(R)**



U \rightarrow **(R)**



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

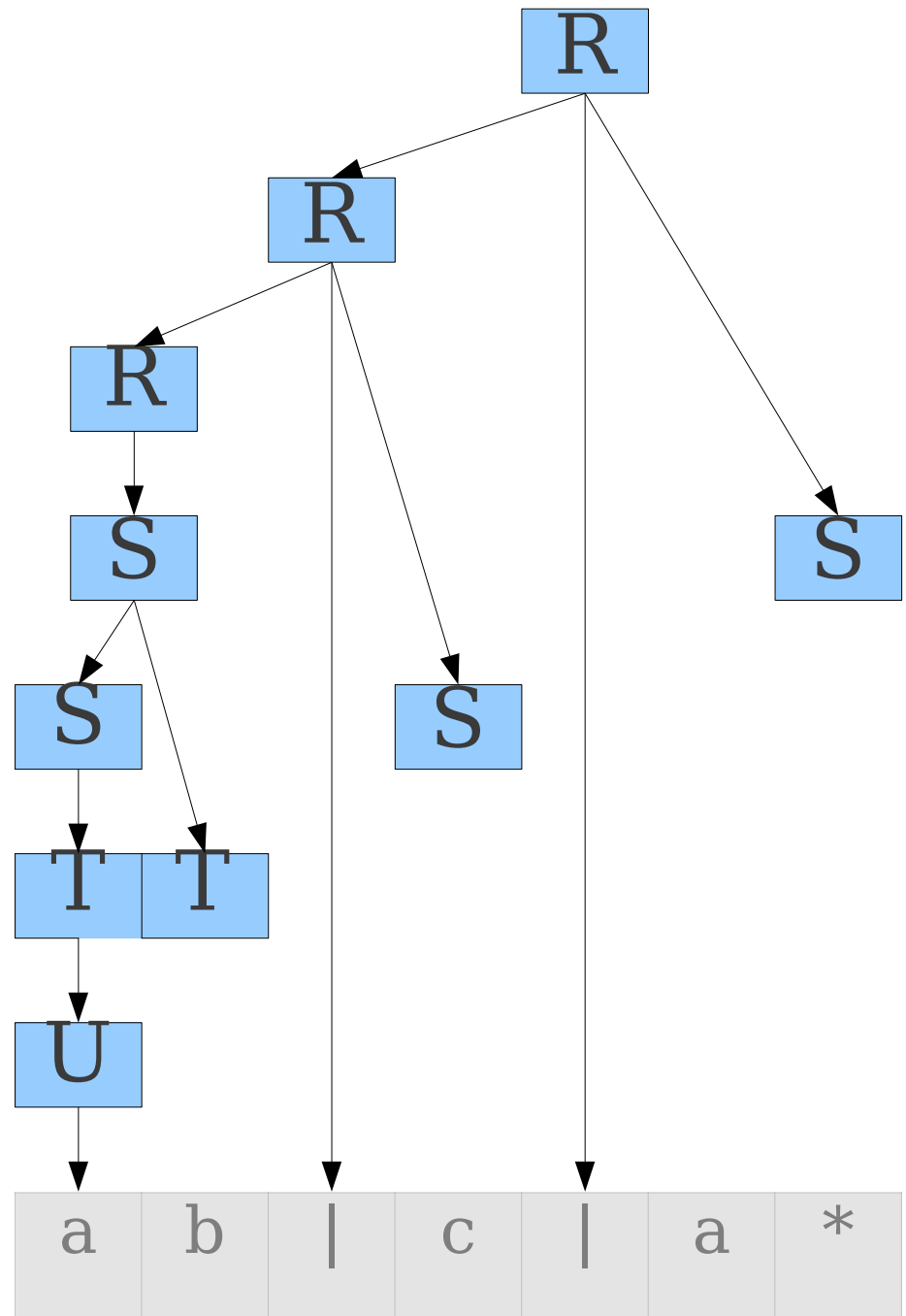
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

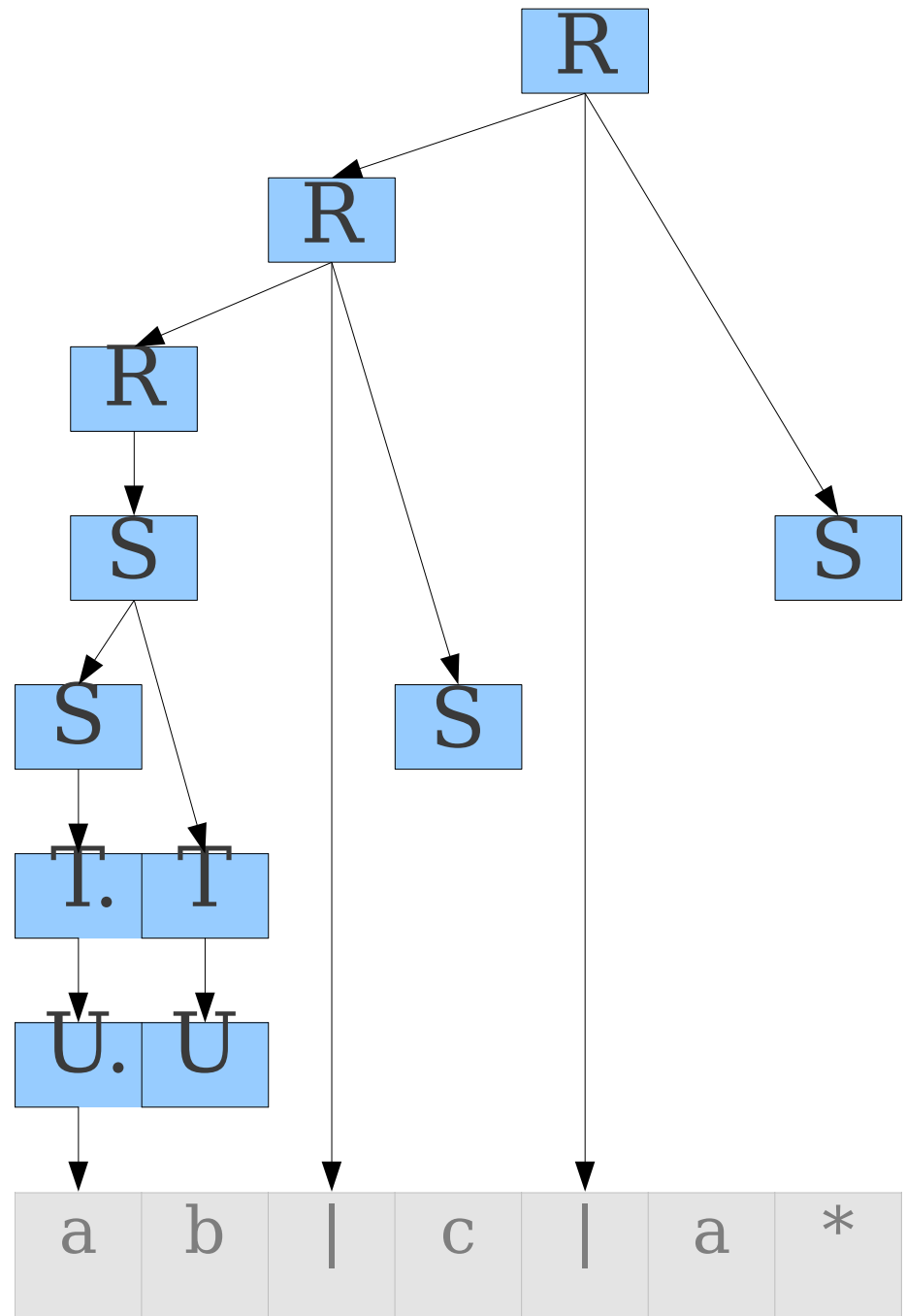
$U \rightarrow a \mid b \mid c \mid \dots$

$U \rightarrow \epsilon$

$U \rightarrow (R)$



U \rightarrow **(R)**



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

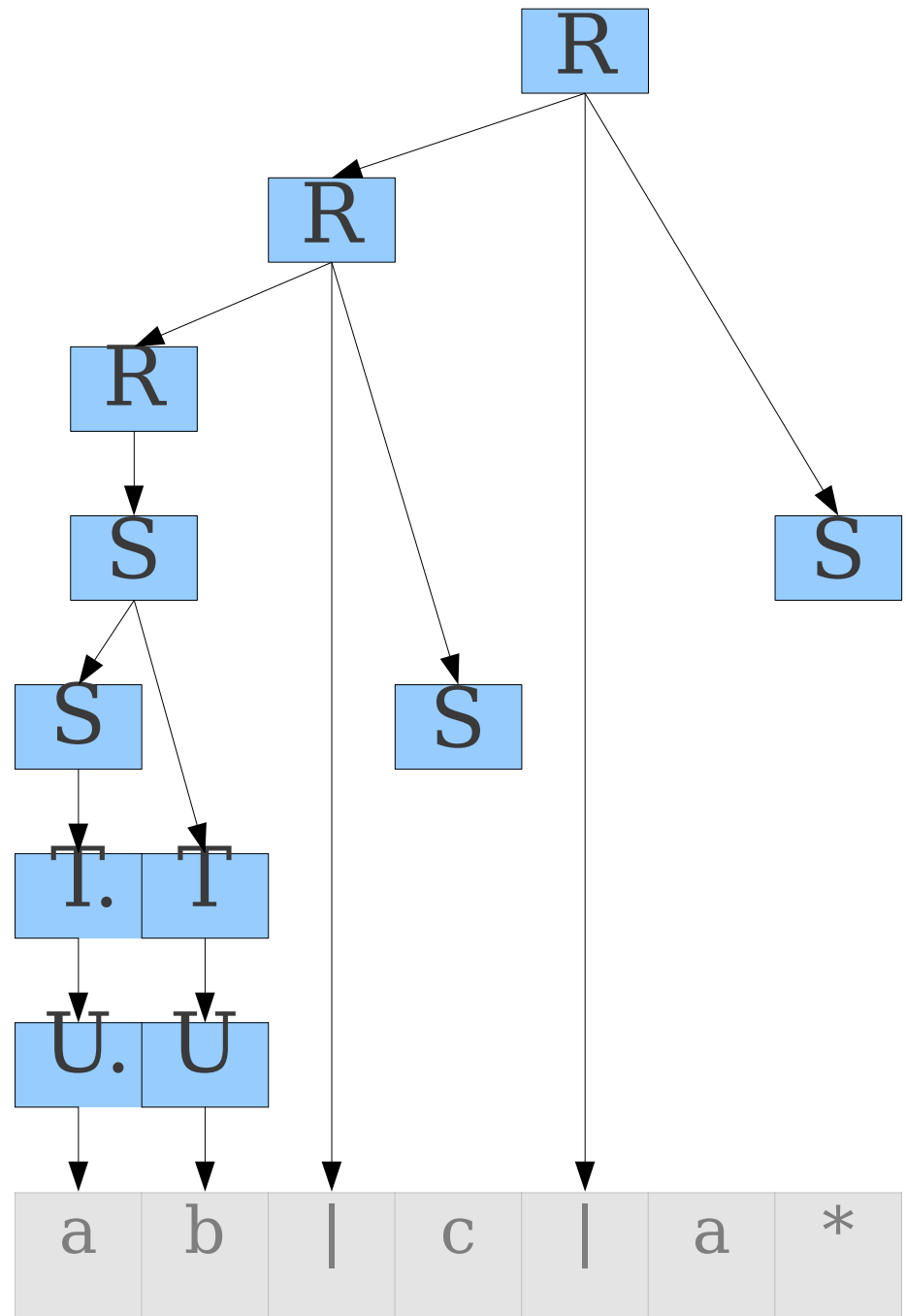
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$



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

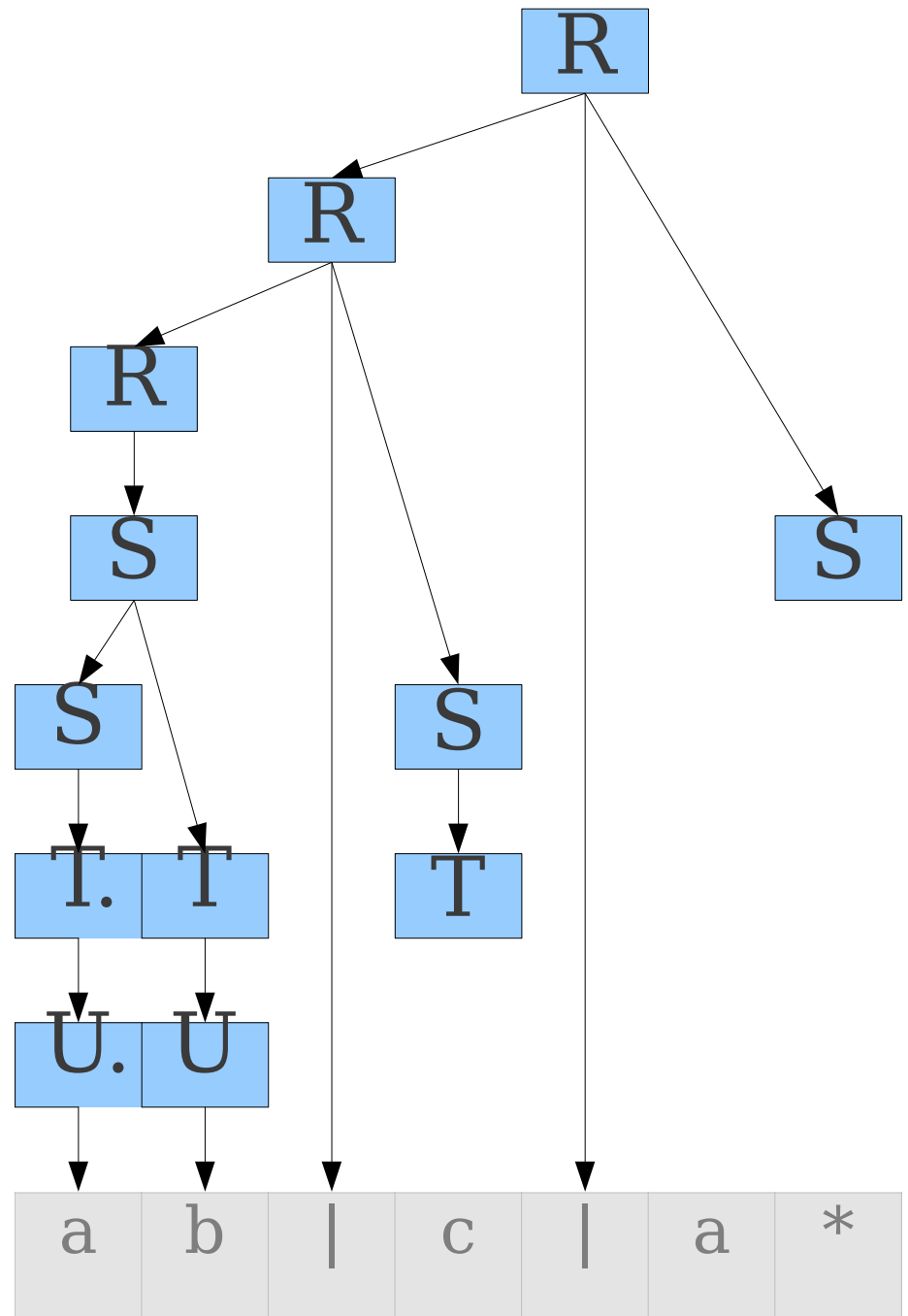
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$



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

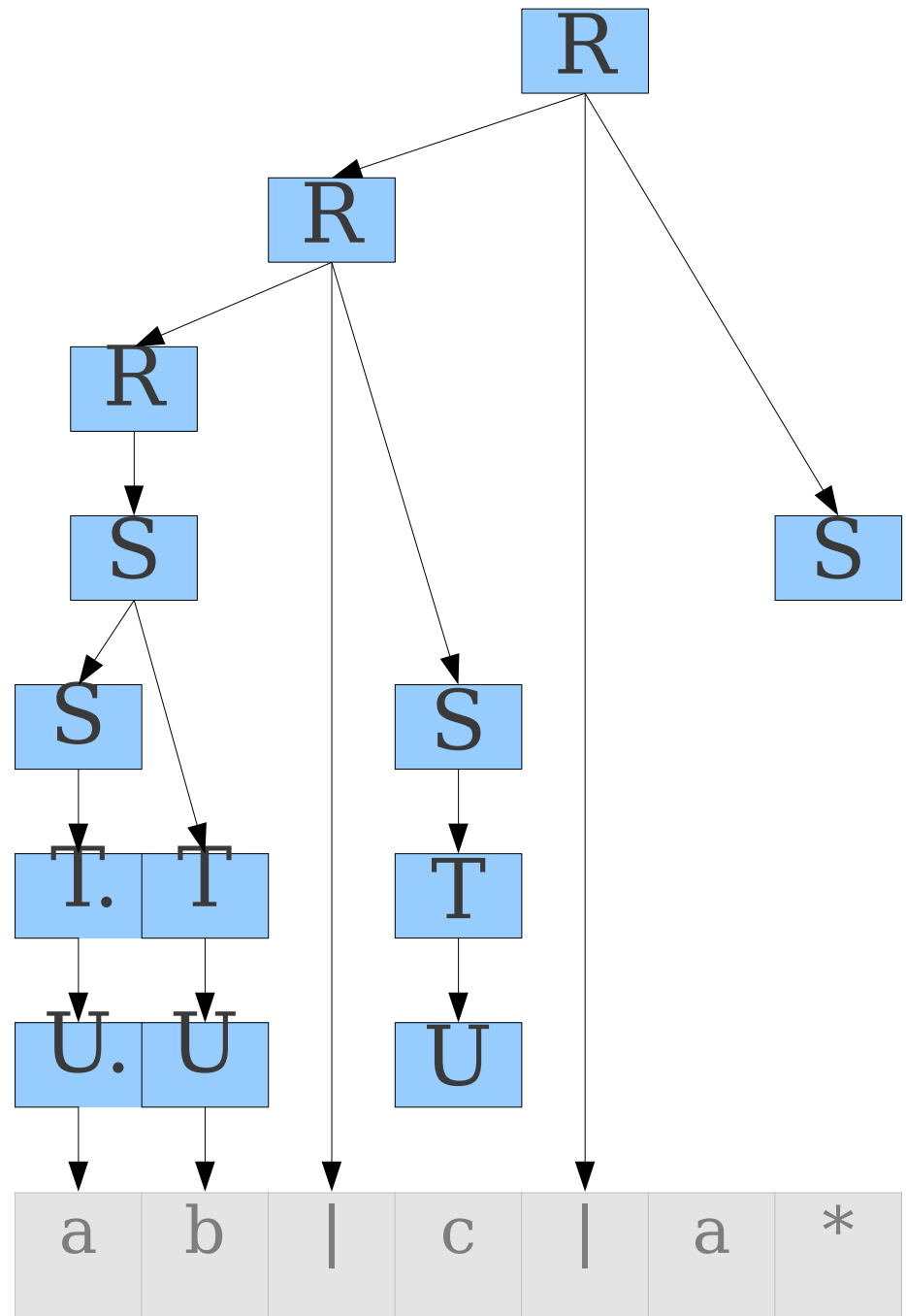
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$



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

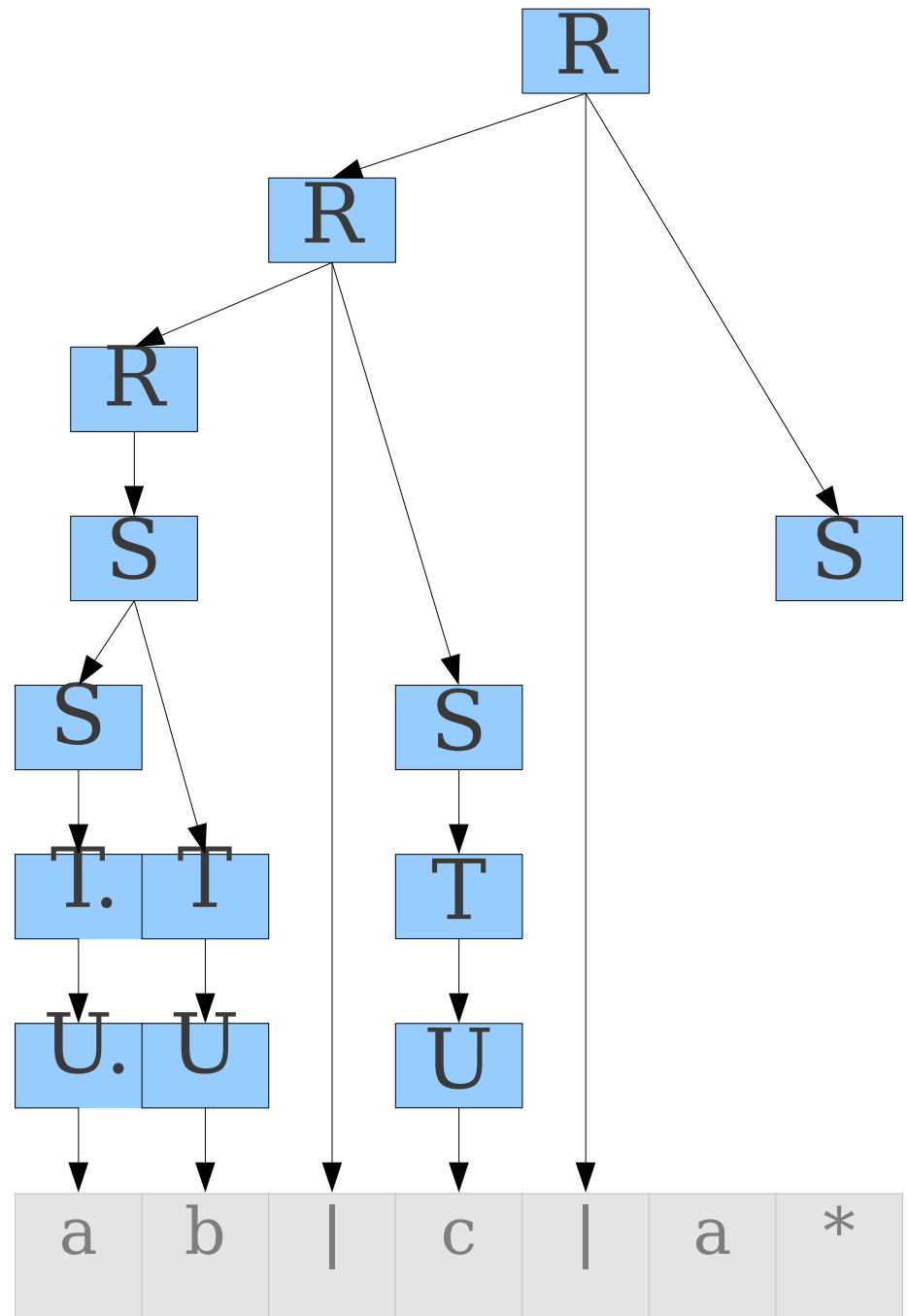
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$



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

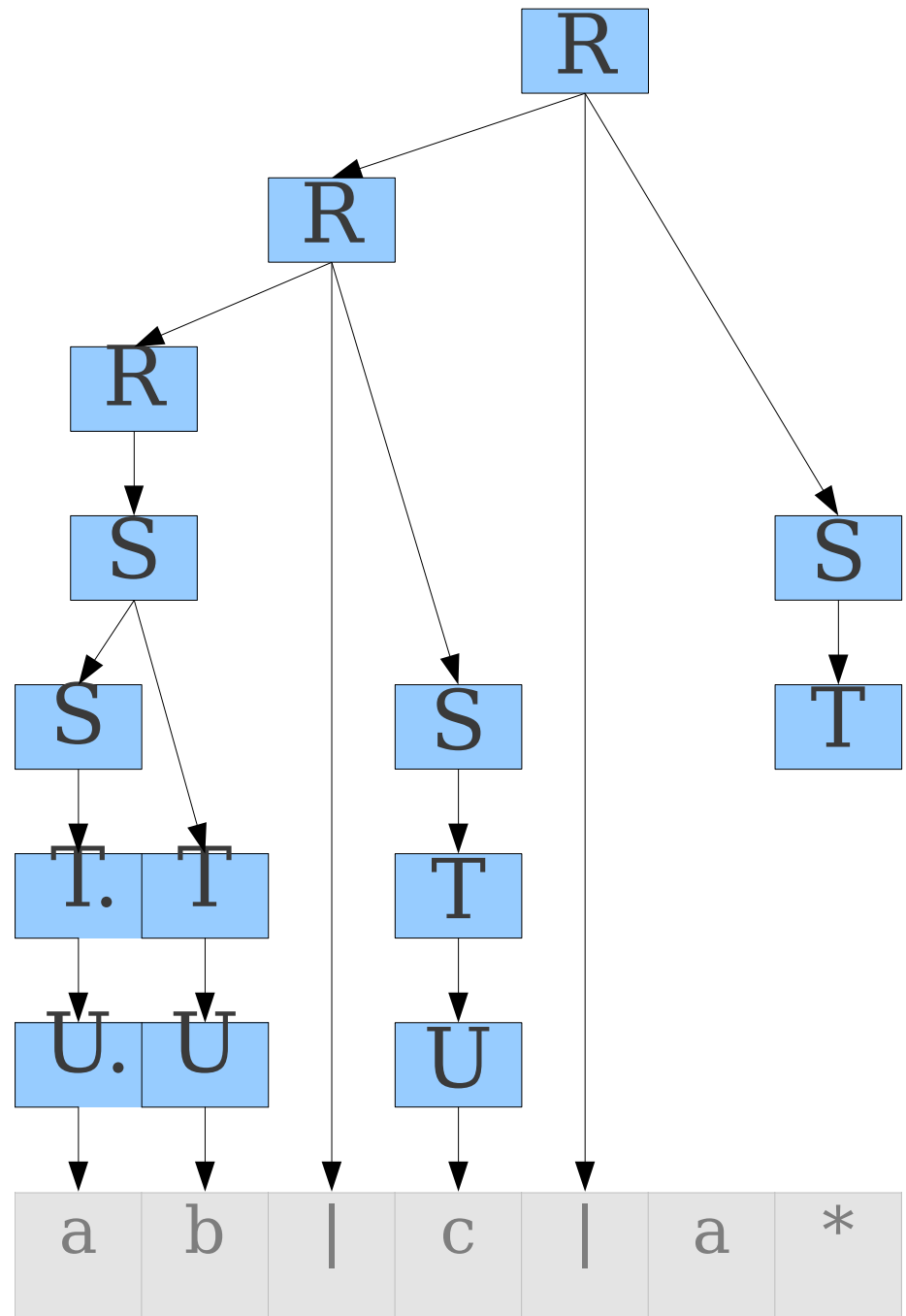
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$



R \rightarrow **S** | **R** “|” **S**

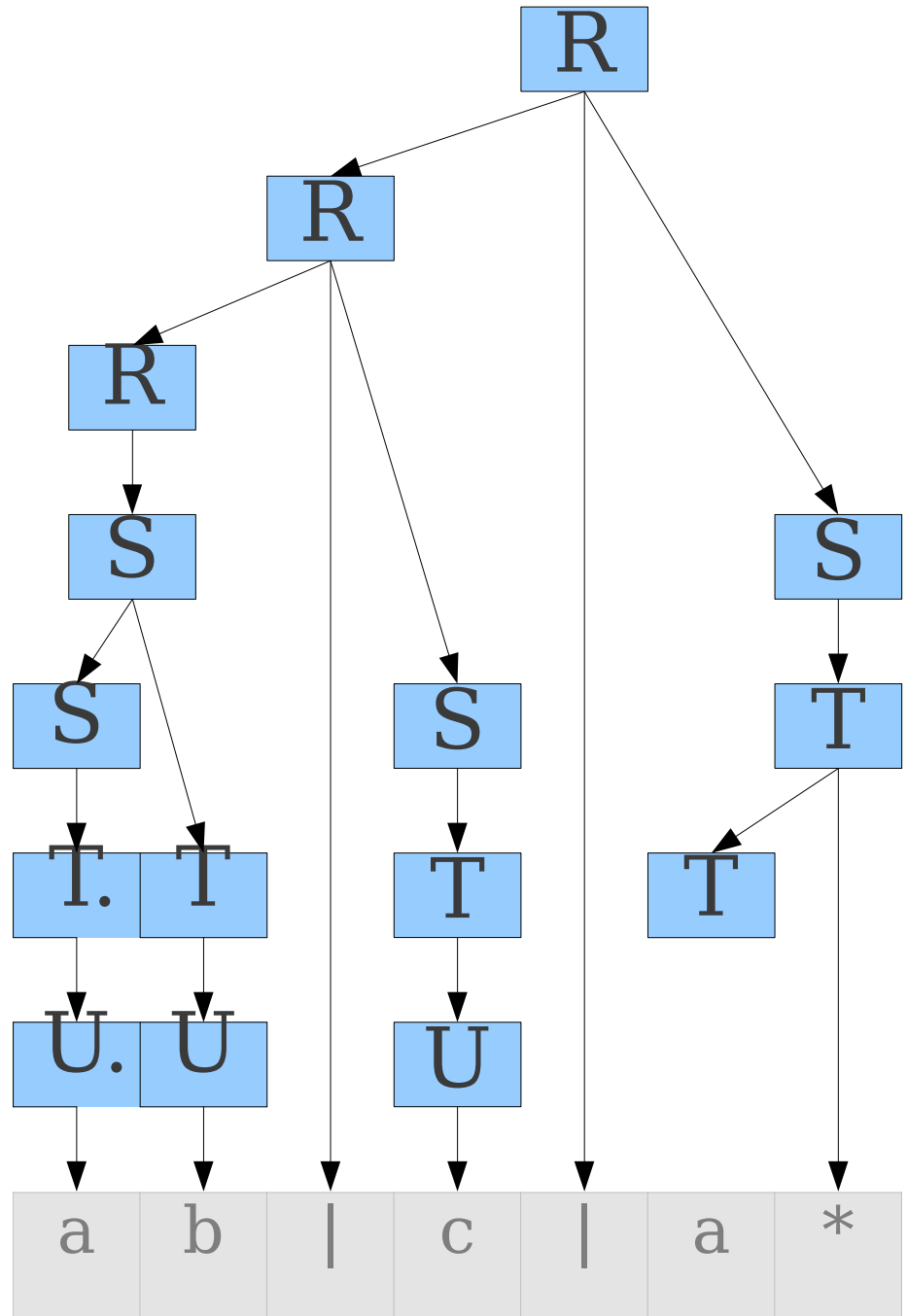
S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

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

U → “**ε**”

U \rightarrow **(R)**



R \rightarrow **S** | **R** “ | ” **S**

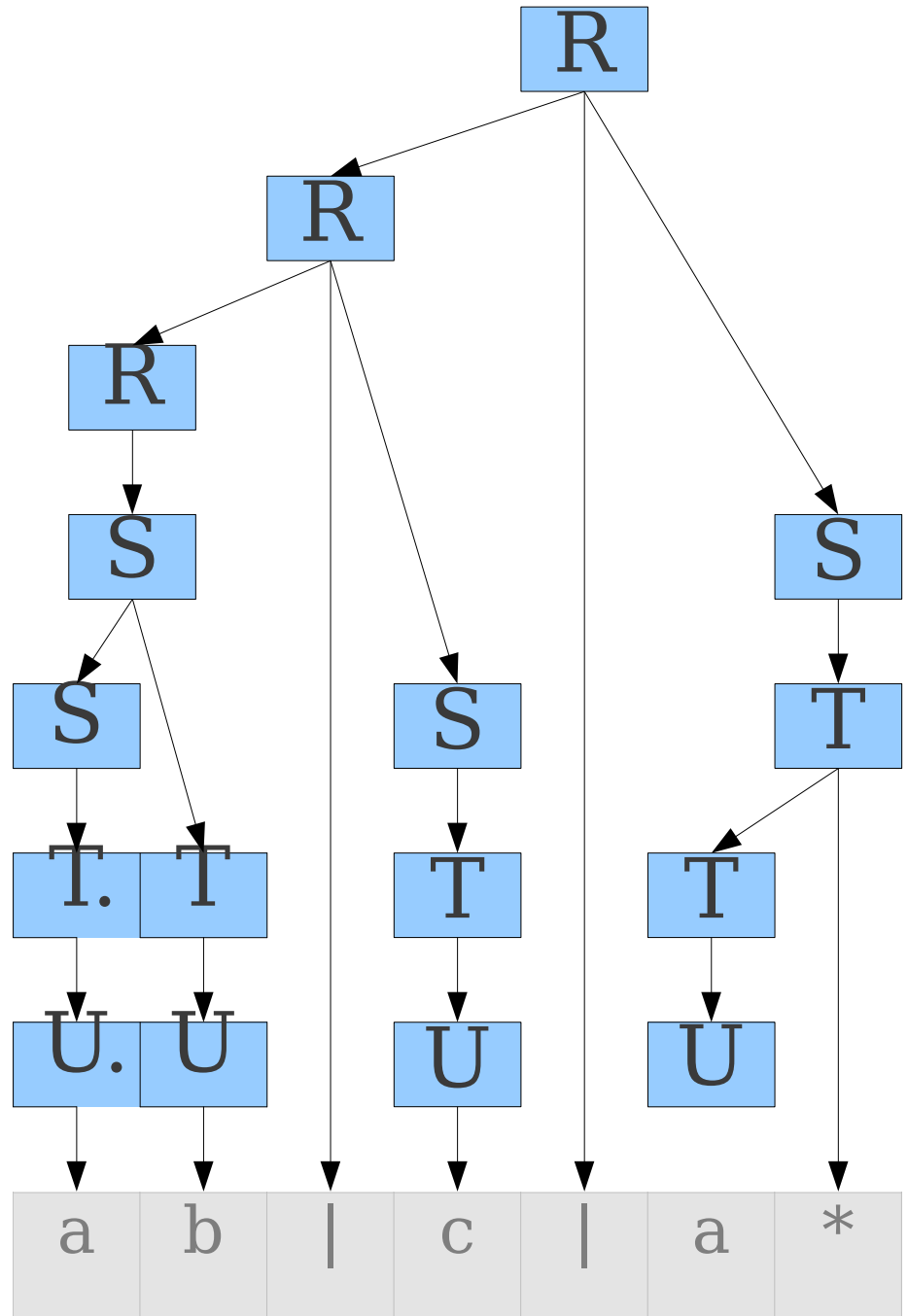
S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

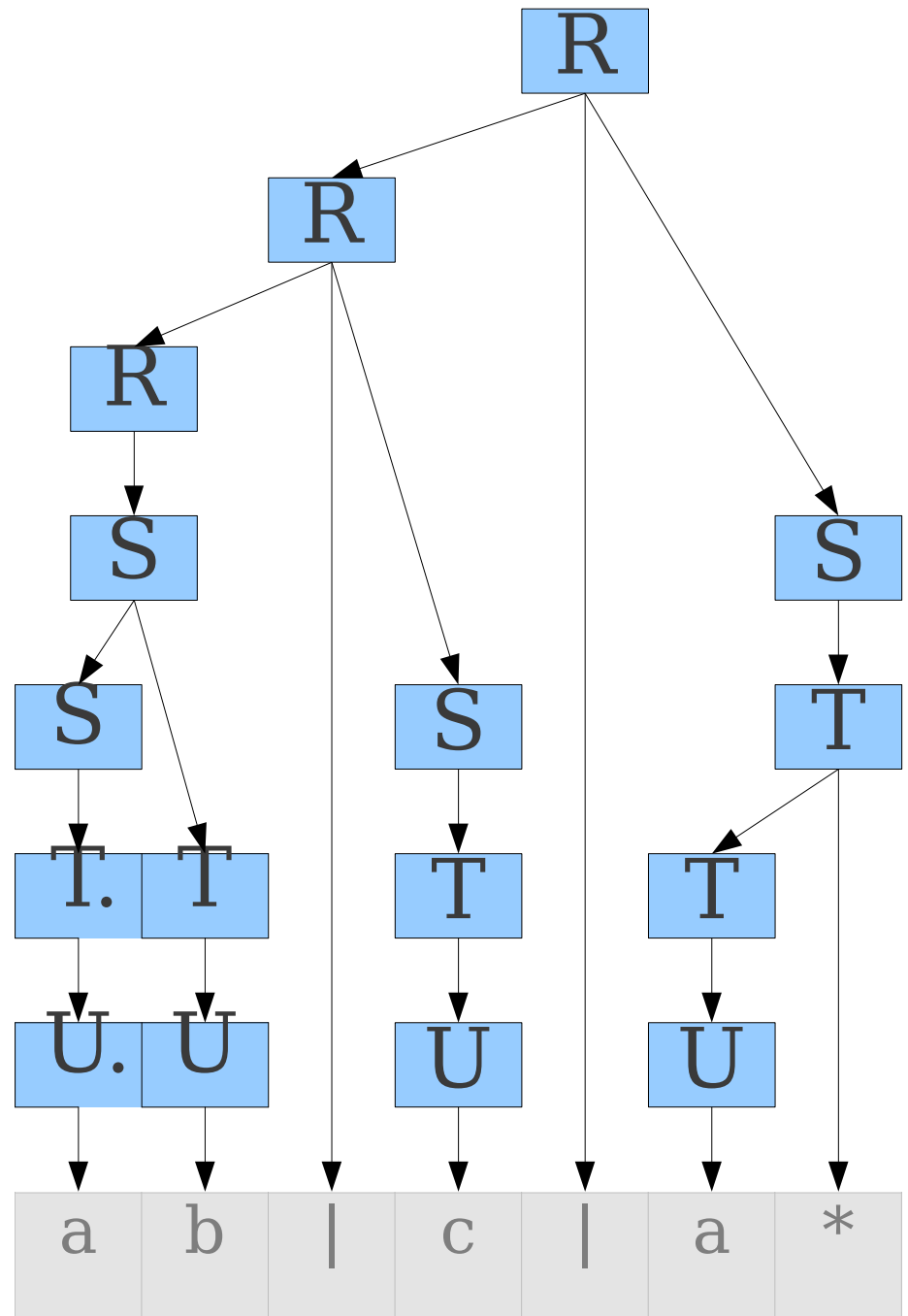
U \rightarrow **a** | **b** | **c** | ...

U → “ε”

U \rightarrow **(R)**



U \rightarrow **(R)**



Precedence Declarations

- If we leave the world of pure CFGs, we can often resolve ambiguities through **precedence declarations**.
 - e.g. multiplication has higher precedence than addition, but lower precedence than exponentiation.
- Allows for unambiguous parsing of ambiguous grammars.
- We'll see how this is implemented later on.

The Structure of a Parse Tree

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

$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

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

$U \rightarrow \epsilon$

$U \rightarrow (R)$

The Structure of a Parse Tree

R \rightarrow **S** | **R** “|” **S**

S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

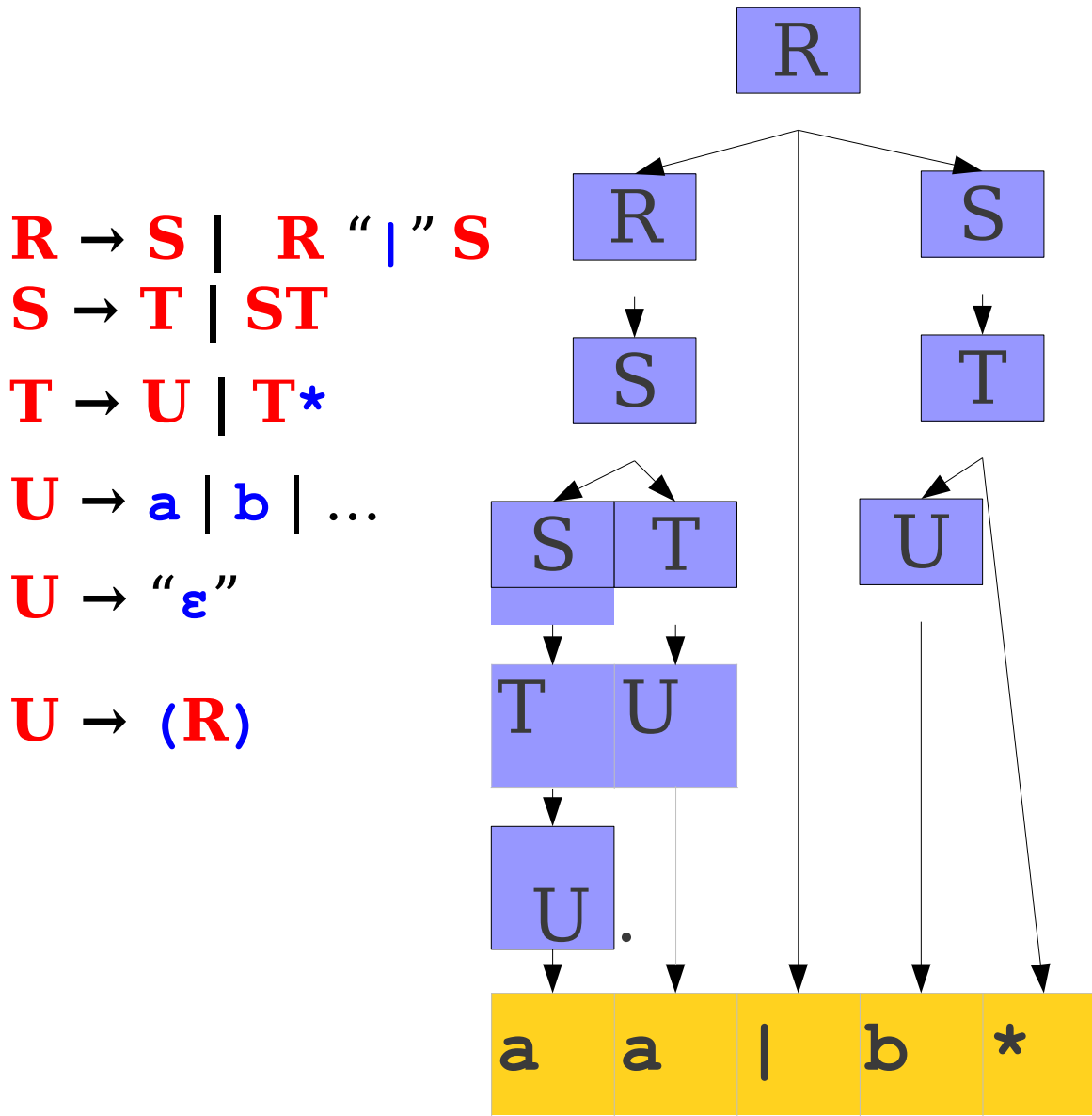
U \rightarrow **a** | **b** | ...

U \rightarrow “ ϵ ”

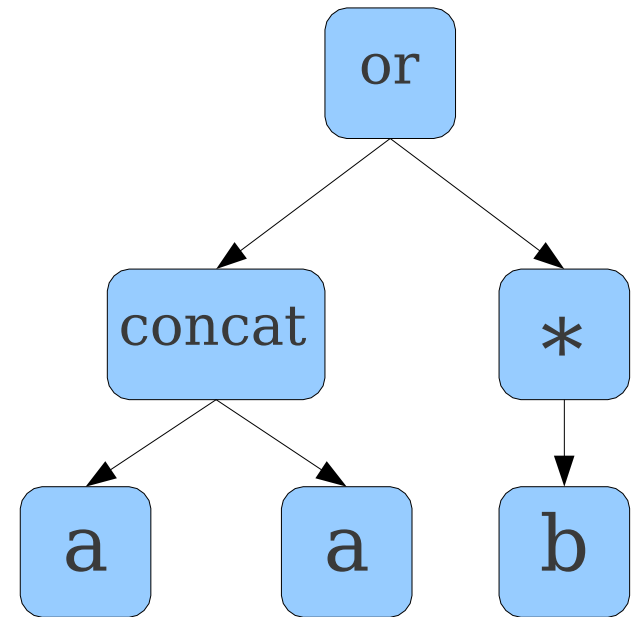
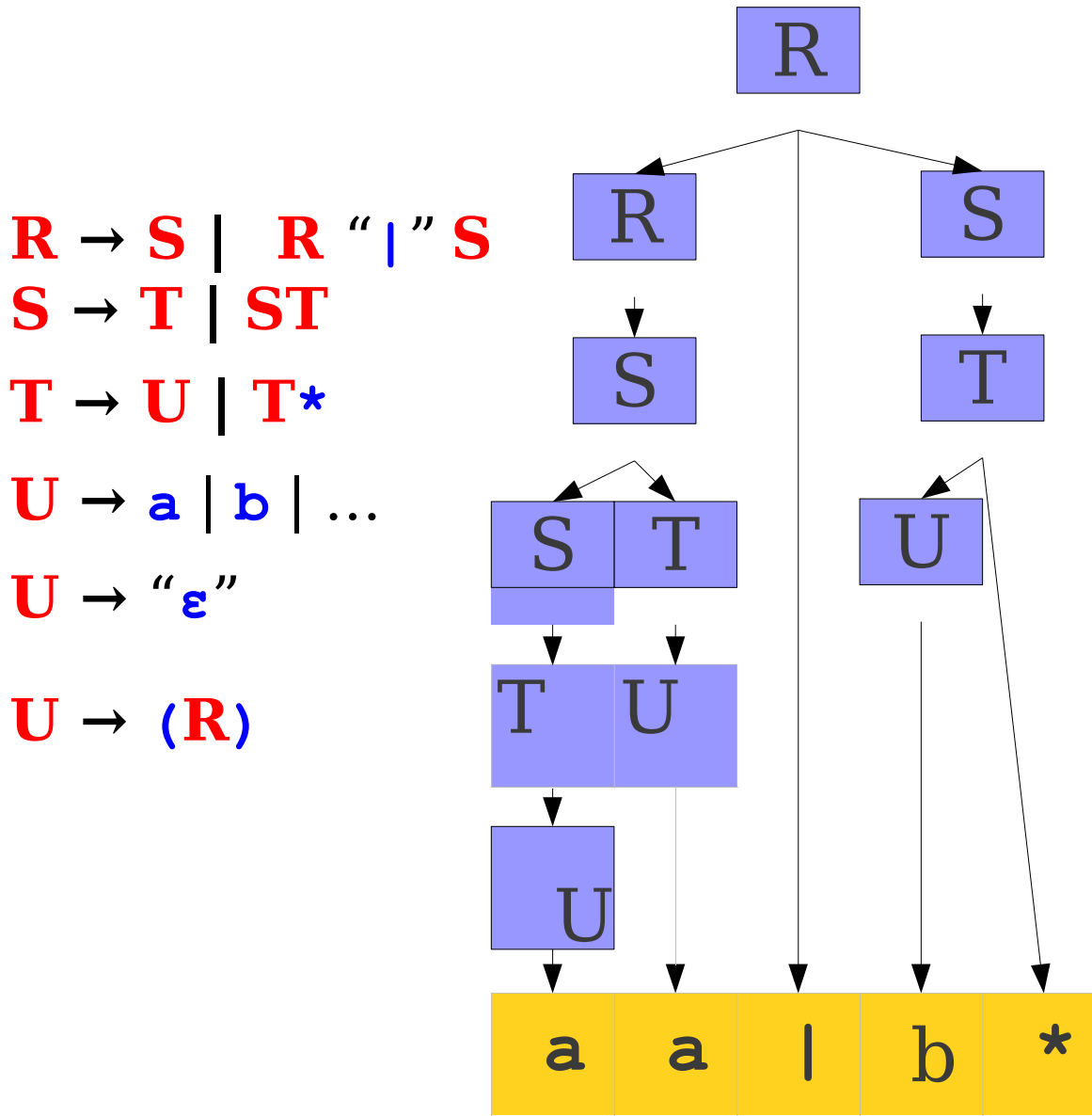
U \rightarrow (**R**)

| | | | | |
|---|---|--|---|---|
| a | a | | b | * |
|---|---|--|---|---|

The Structure of a Parse Tree



The Structure of a Parse Tree



R \rightarrow **S** | **R** “|” **S**

S \rightarrow **T** | **ST**

T \rightarrow **U** | **T***

U \rightarrow **a** | **b** | ...

U \rightarrow “ ϵ ”

U \rightarrow (**R**)

| | | | | | |
|---|---|---|--|---|---|
| a | (| b | | c |) |
|---|---|---|--|---|---|

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

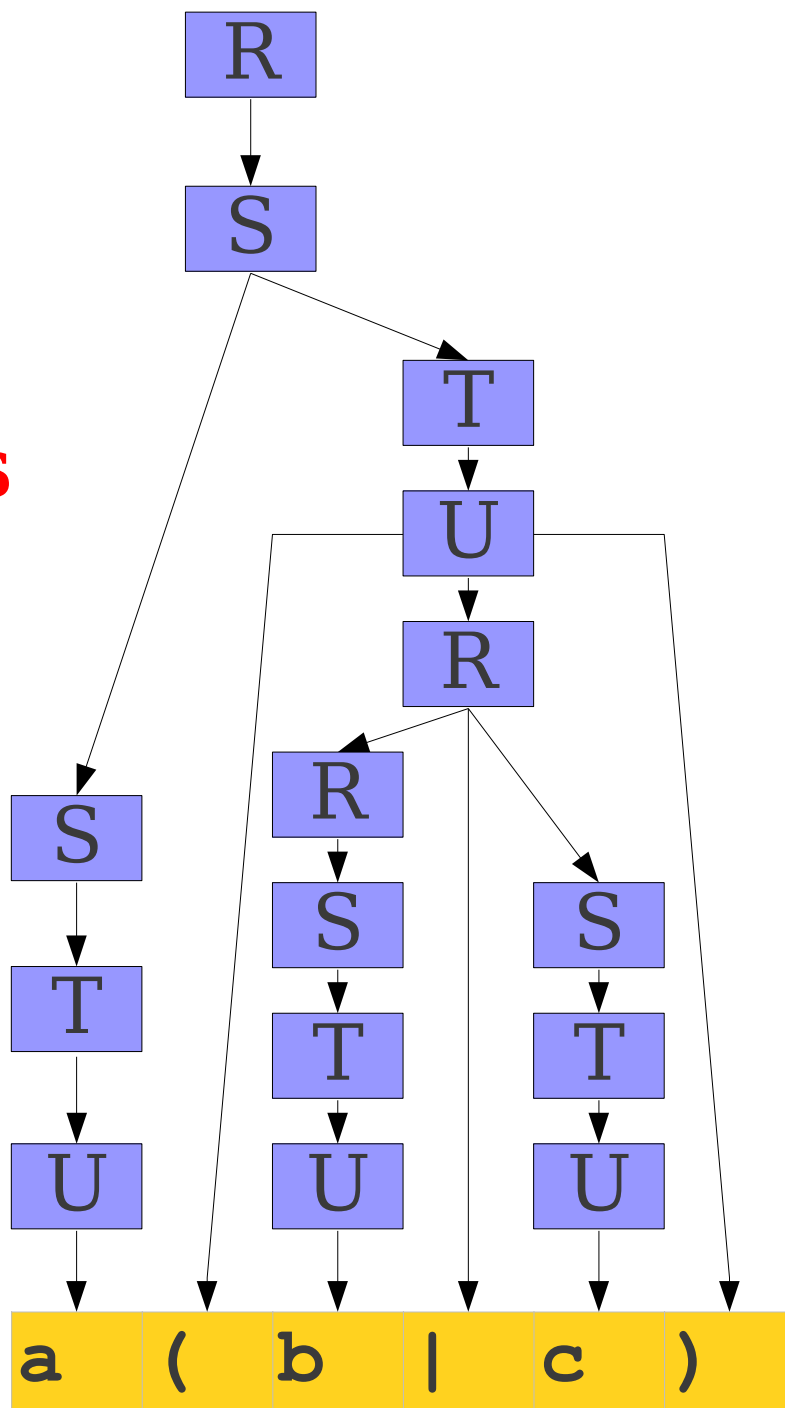
$S \rightarrow T \mid ST$

$T \rightarrow U \mid T^*$

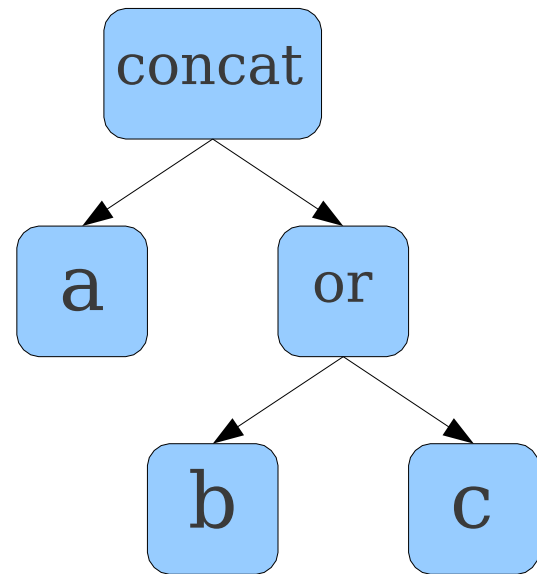
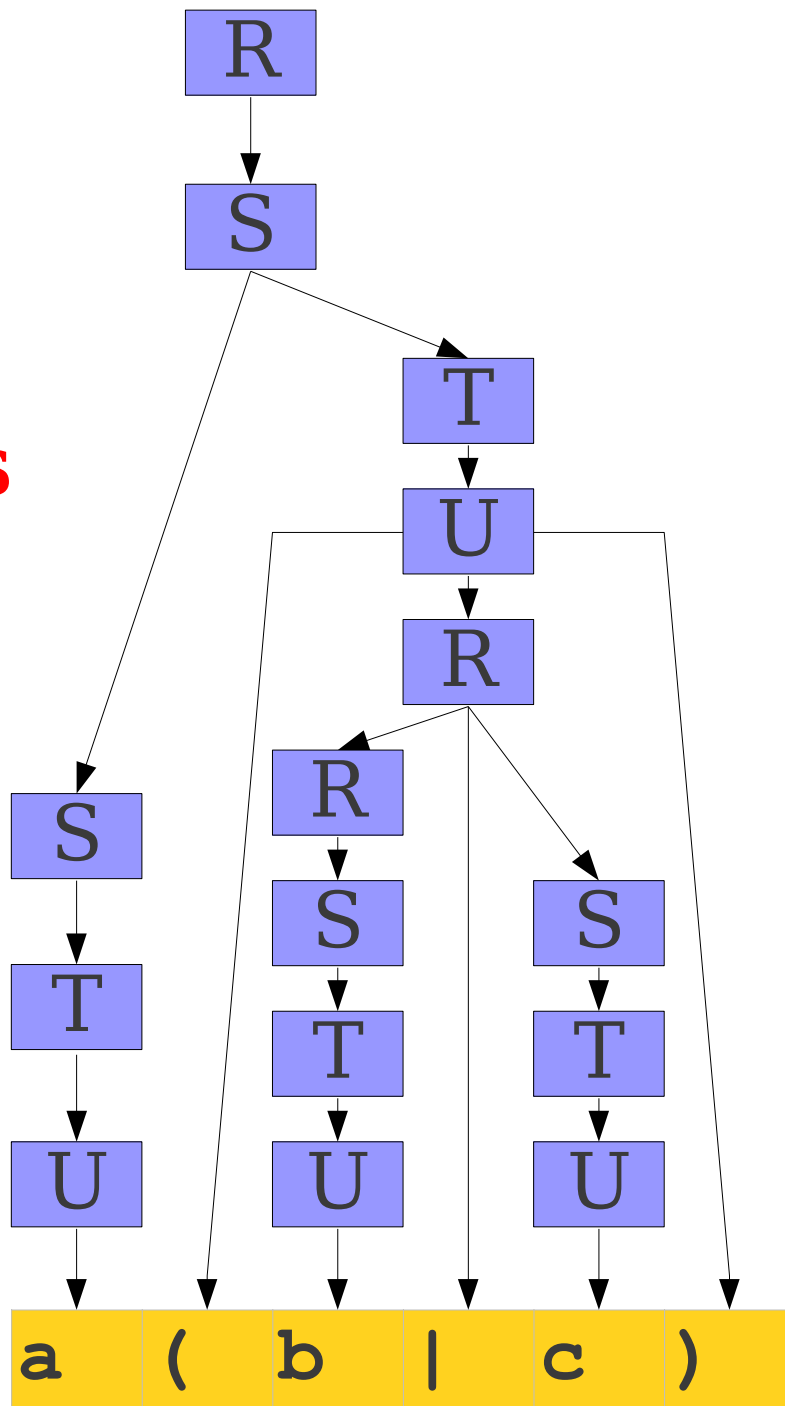
$U \rightarrow a \mid b \mid \dots$

$U \rightarrow \epsilon$

$U \rightarrow (R)$



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



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.
 - Exact order of code execution depends on the parsing method used.
- This is called a **syntax-directed translation**.

Different Types of Parsing

- **Top-Down Parsing**
 - Beginning with the start symbol, try to guess the productions to apply to end up at the user's program.
- **Bottom-Up Parsing**
 - Beginning with the user's program, try to apply productions in reverse to convert the program back into the start symbol.

Challenges in Top-Down Parsing

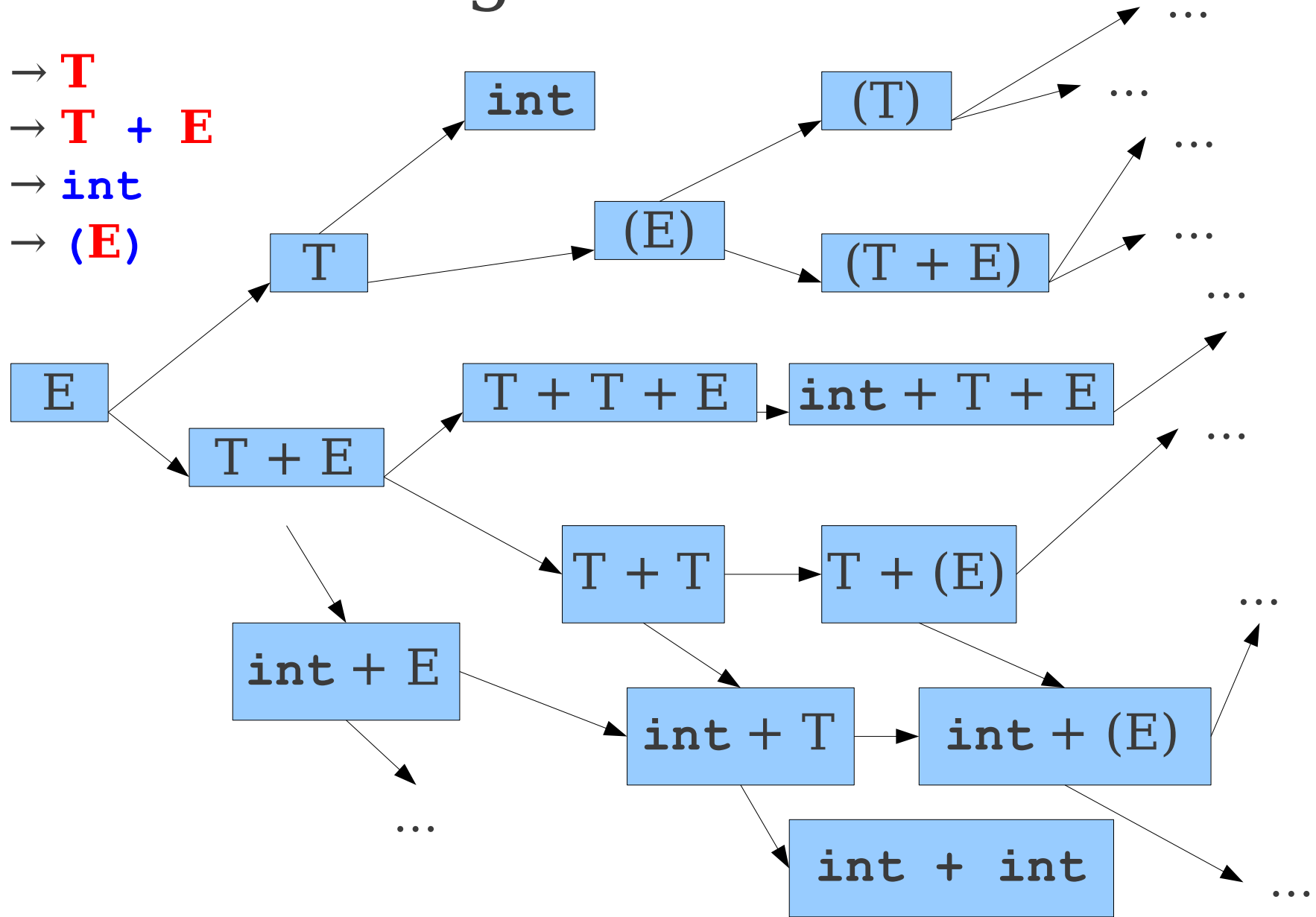
- Top-down parsing begins with virtually no information.
 - Begins with just the **start symbol**, which matches *every* program.
- How can we know which productions to apply?
- In general, we can't.
 - There are some grammars for which the best we can do is guess and backtrack if we're wrong.
 - If we have to guess, how do we do it?

Parsing as a Search

- An idea: **treat parsing as a graph search**.
- Each node is a **sentential form** (a string of terminals and nonterminals derivable from the start symbol).
- There is an edge from node α to node β iff $\alpha \Rightarrow \beta$.

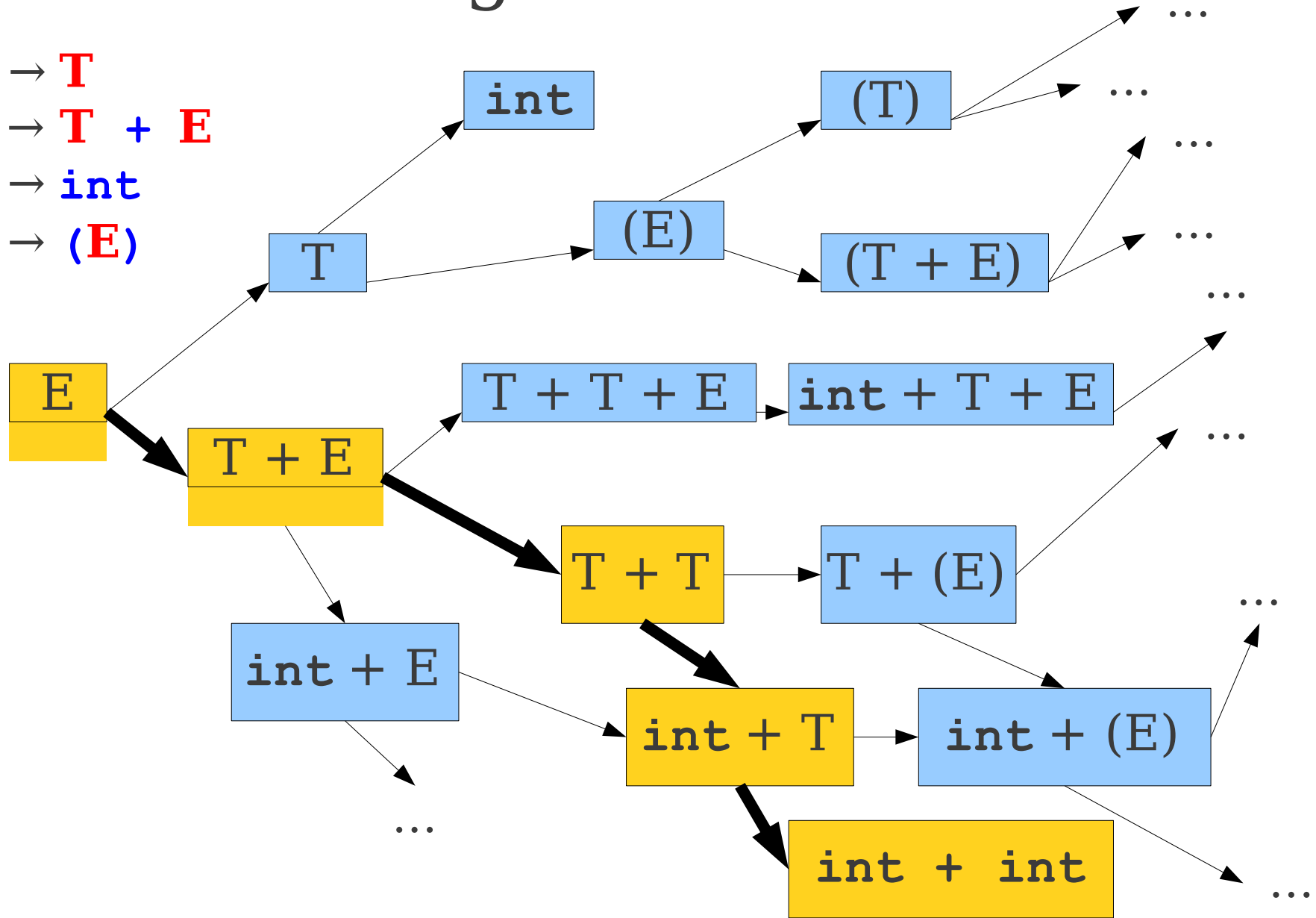
Parsing as a Search

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



Parsing as a Search

T \rightarrow **(E)**



Our First Top-Down Algorithm

- **Breadth-First Search**
- Maintain a worklist of sentential forms, initially just the start symbol **S**.
- While the worklist isn't empty:
 - Remove an element from the worklist.
 - If it matches the target string, you're done.
 - Otherwise, for each possible string that can be derived in one step, add that string to the worklist.
- Can recover a parse tree by tracking what productions we applied at each step.

Breadth-First Search Parsing

Worklist

E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Breadth-First Search Parsing



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Breadth-First Search Parsing

Worklist

E

E → **T**

E → **T** + **E**

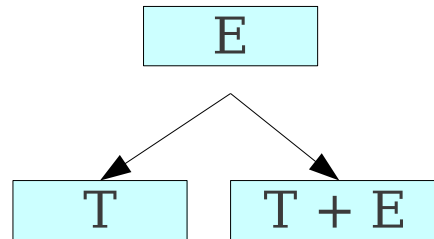
T → **int**

T → (**E**)

int + int

Breadth-First Search Parsing

Worklist



E → **T**

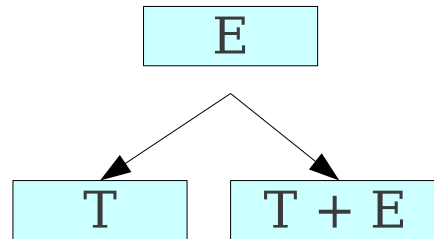
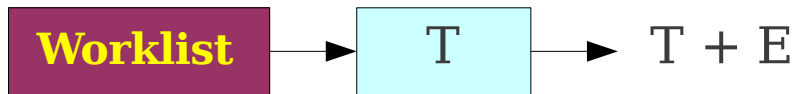
E → **T** + **E**

T → **int**

T → **(E)**

`int + int`

Breadth-First Search Parsing



E → **T**

E → **T** + **E**

T → **int**

T → **(E)**

`int + int`

Breadth-First Search Parsing



E → **T**

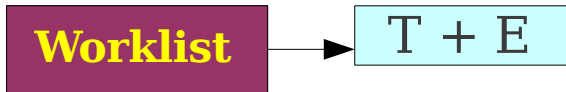
E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Breadth-First Search Parsing

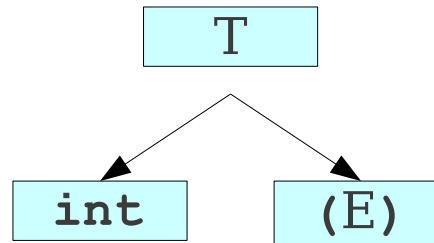
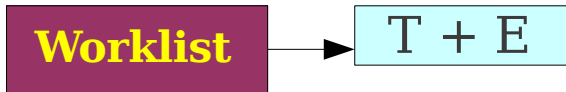


T

E → **T**
E → **T** + **E**
T → **int**
T → **(E)**

int + int

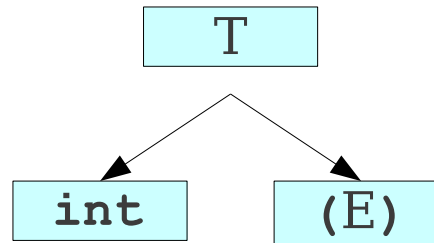
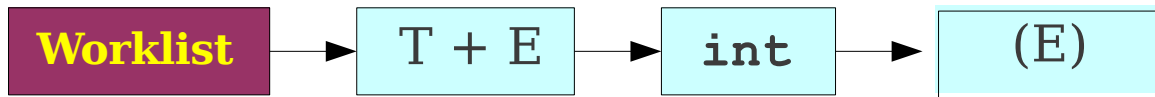
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

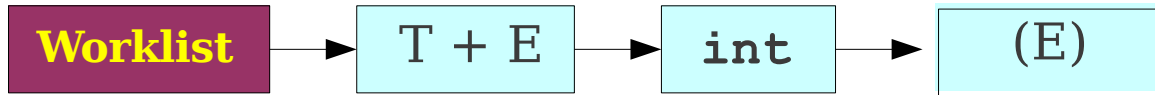
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Breadth-First Search Parsing



E → **T**

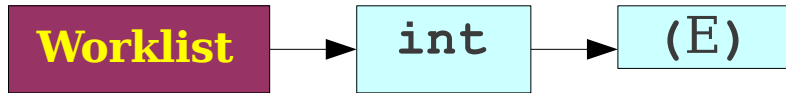
E → **T** + **E**

T → **int**

T → **(E)**

`int + int`

Breadth-First Search Parsing



T + E

E → **T**

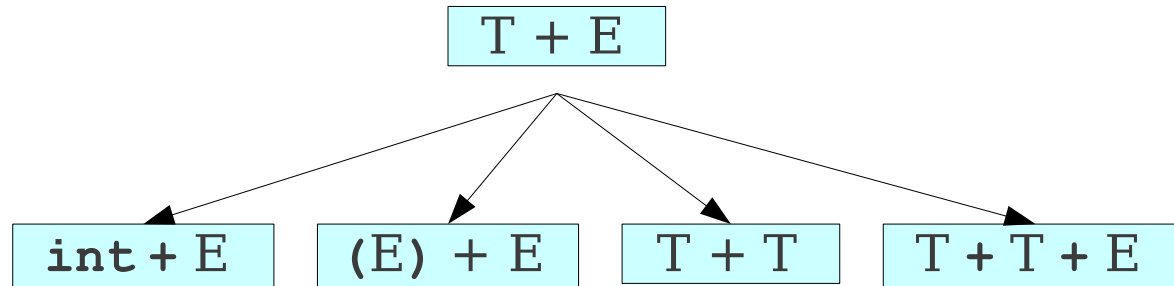
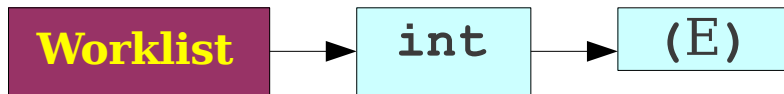
E → **T** + **E**

T → **int**

T → **(E)**

int + int

Breadth-First Search Parsing



E → **T**

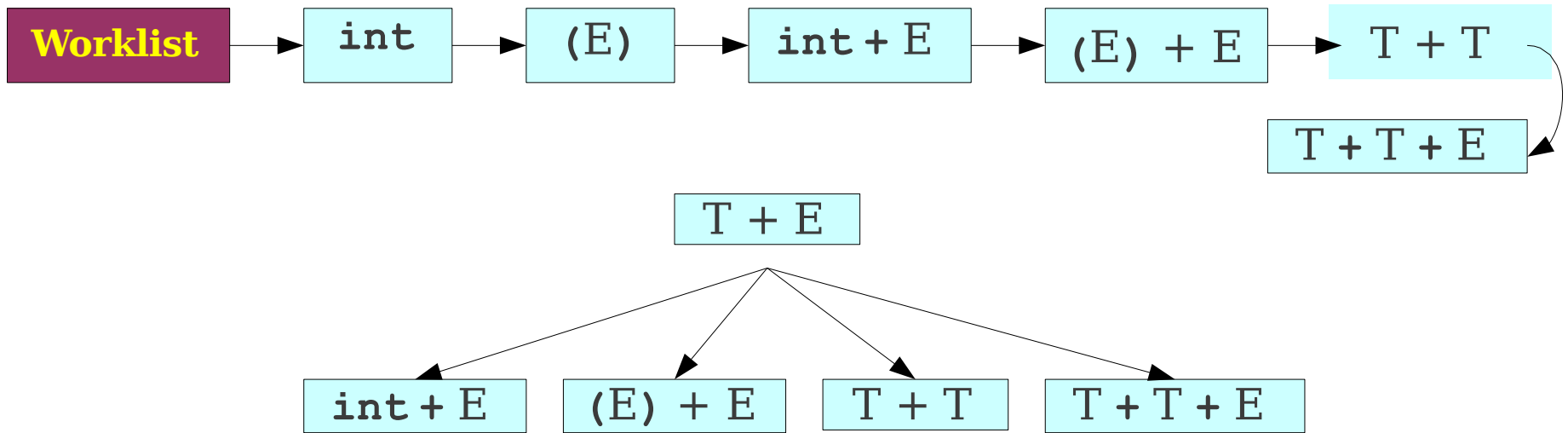
E → **T** + **E**

T → int

T → (**E**)

int + int

Breadth-First Search Parsing



E → **T**

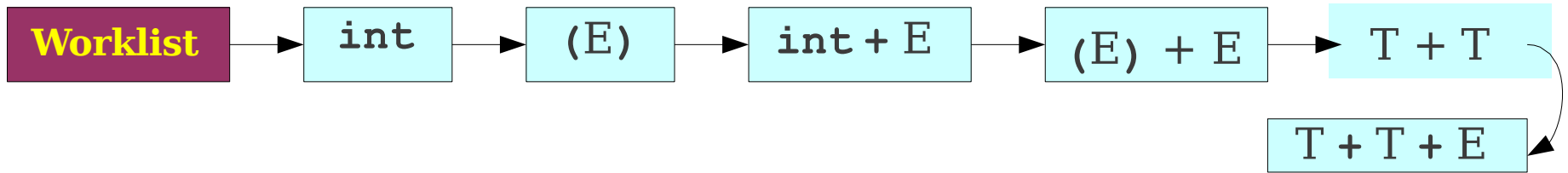
E → **T** + **E**

T → int

T → (**E**)

int + int

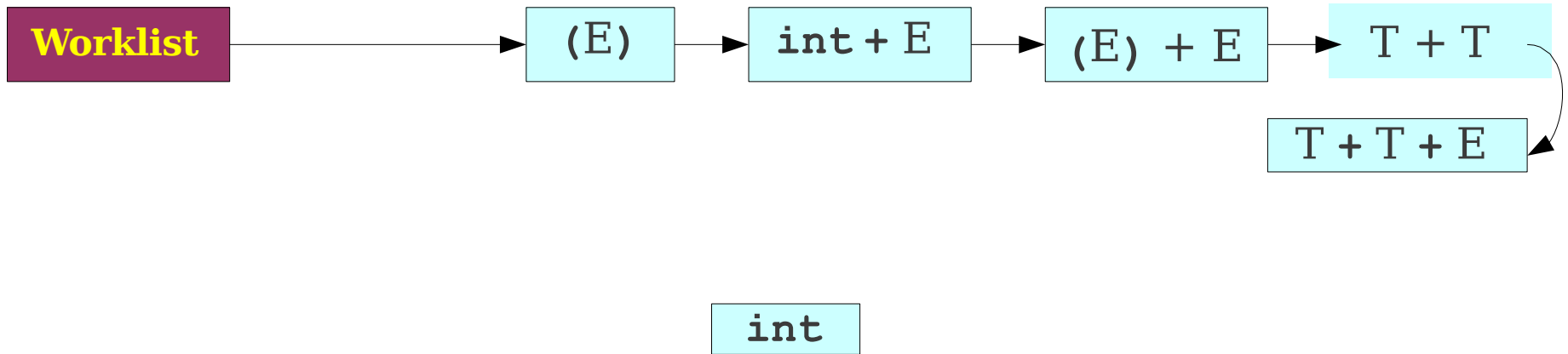
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

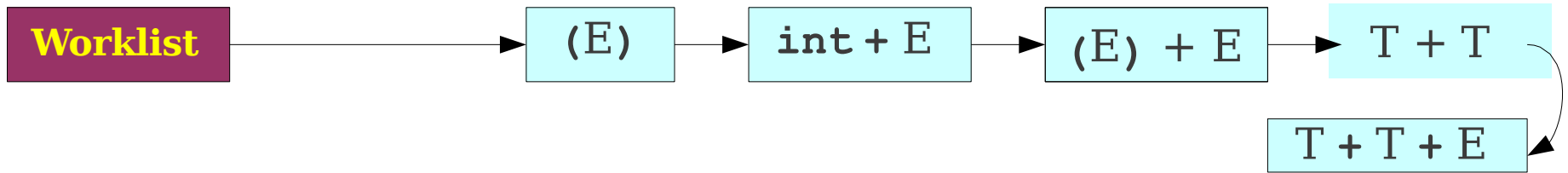
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

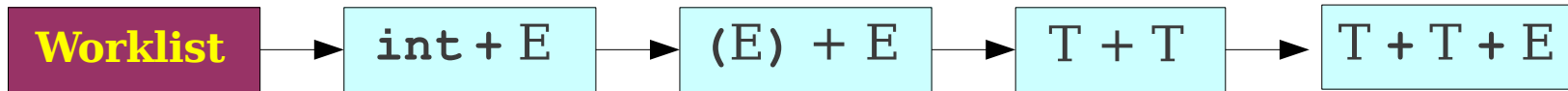
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

Breadth-First Search Parsing



(E)

E → **T**

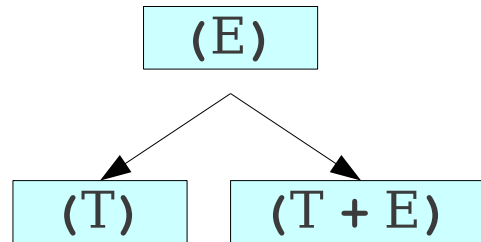
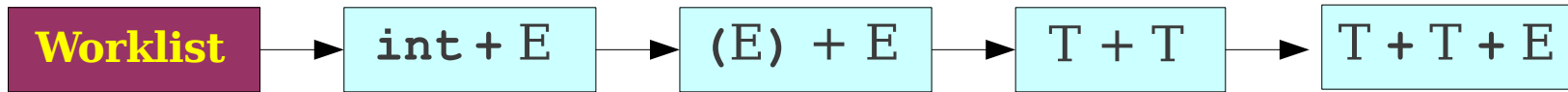
E → **T** + **E**

T → **int**

T → **(E)**

int + int

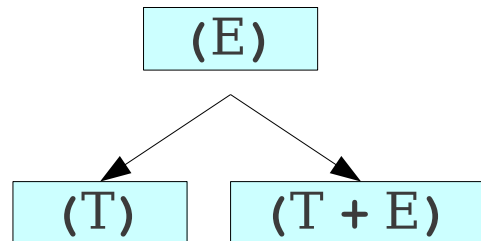
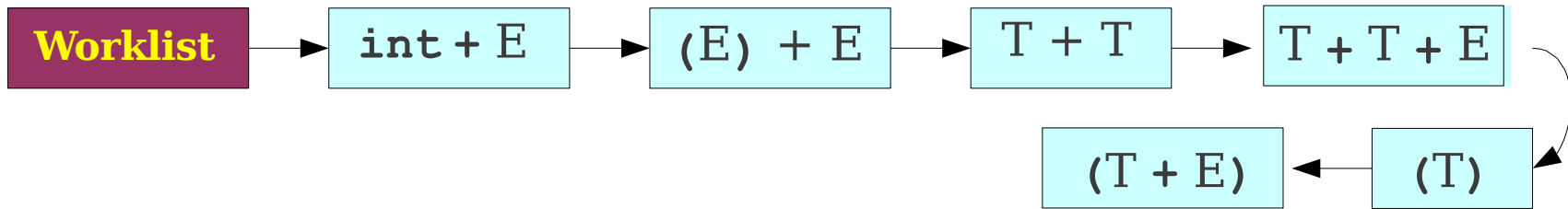
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

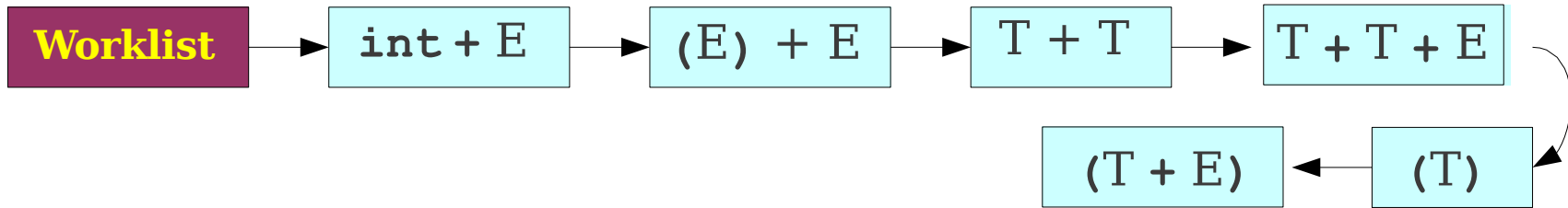
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Breadth-First Search Parsing



E → **T**

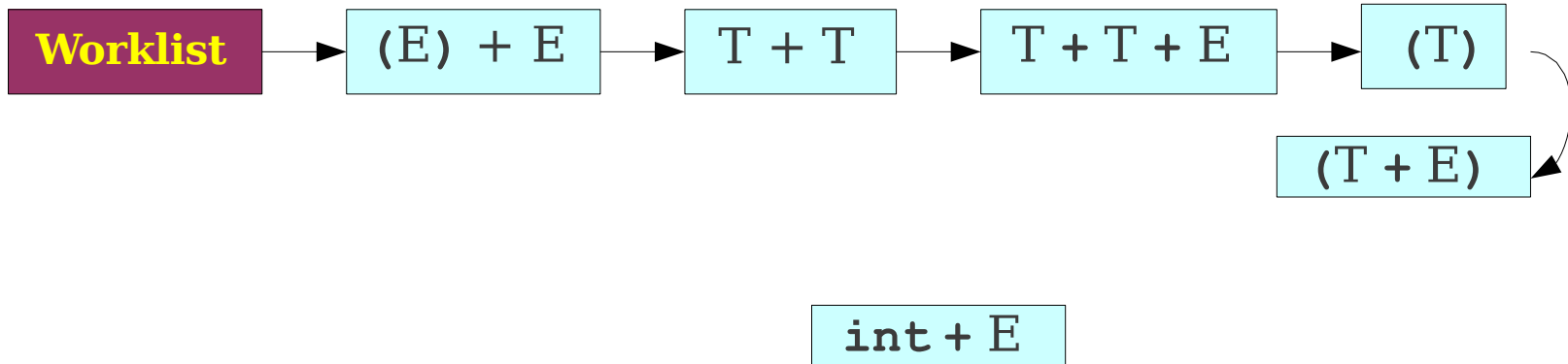
E → **T** + **E**

T → **int**

T → **(E)**

int + int

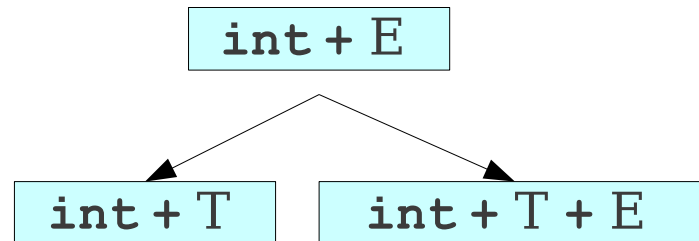
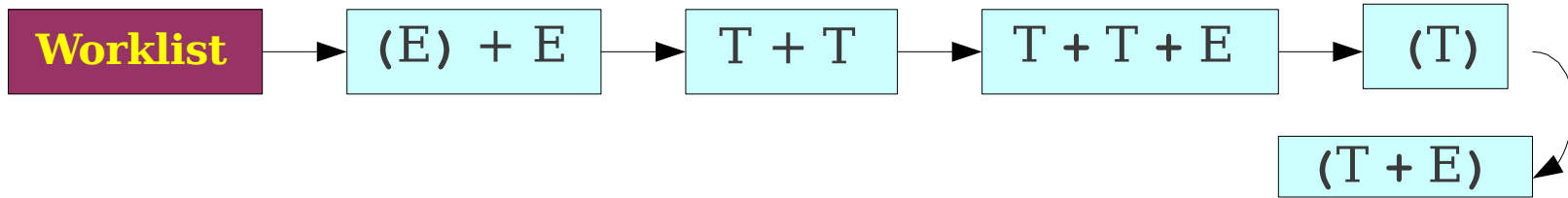
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow int$
 $T \rightarrow (E)$

`int + int`

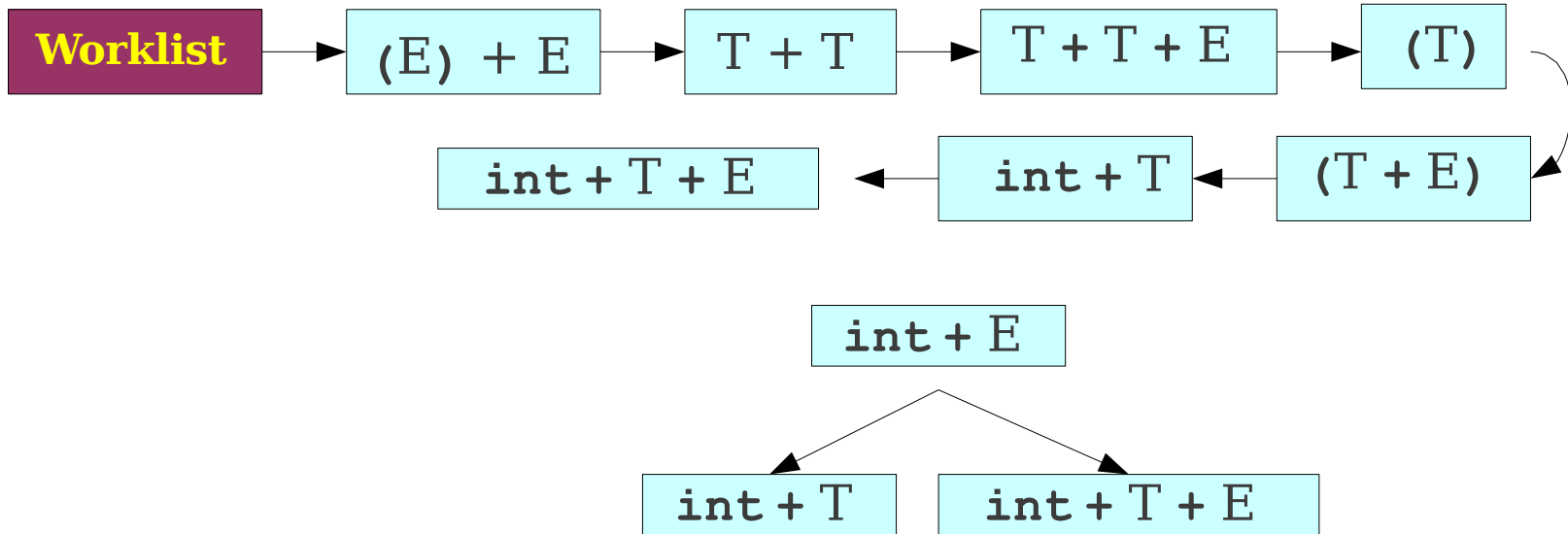
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

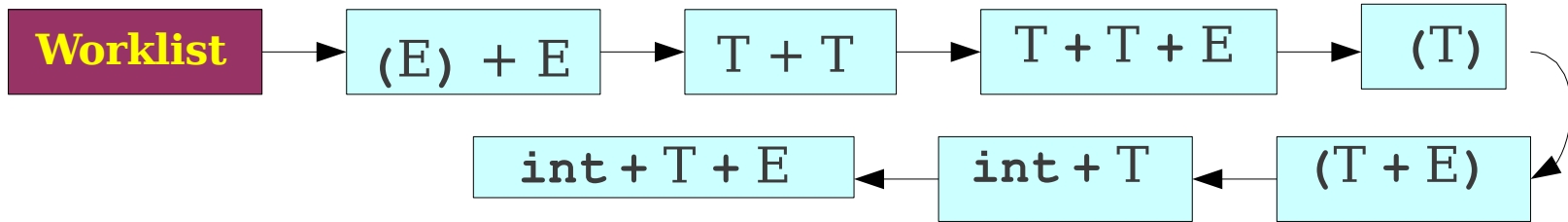
Breadth-First Search Parsing



E → **T**
E → **T** + **E**
T → **int**
T → **(E)**

int + int

Breadth-First Search Parsing



E → **T**

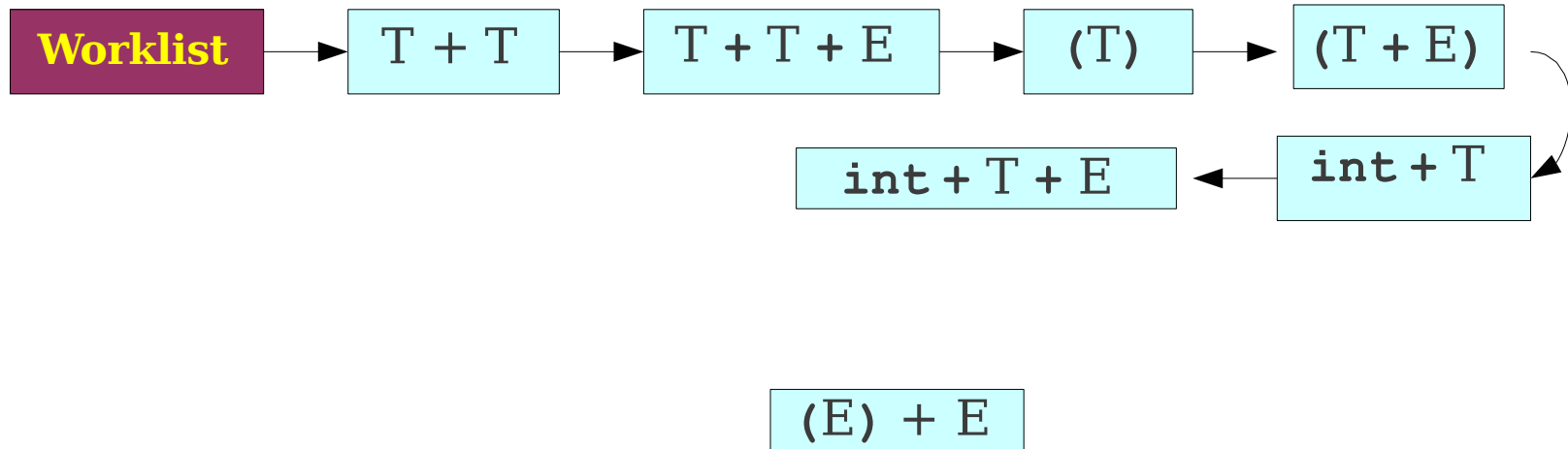
E → **T** + **E**

T → **int**

T → **(E)**

`int + int`

Breadth-First Search Parsing



E → **T**

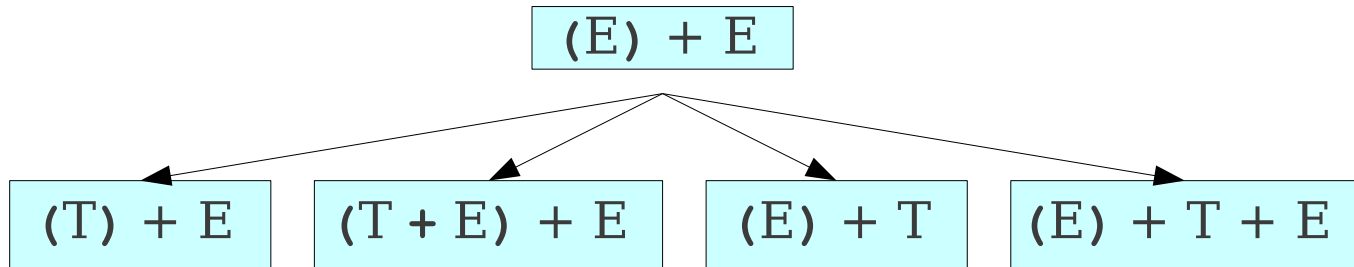
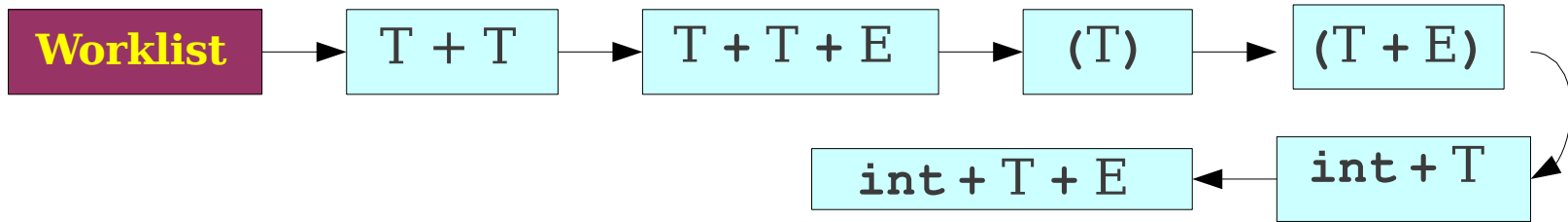
E → **T + E**

T → **int**

T → **(E)**

int + int

Breadth-First Search Parsing



E → **T**

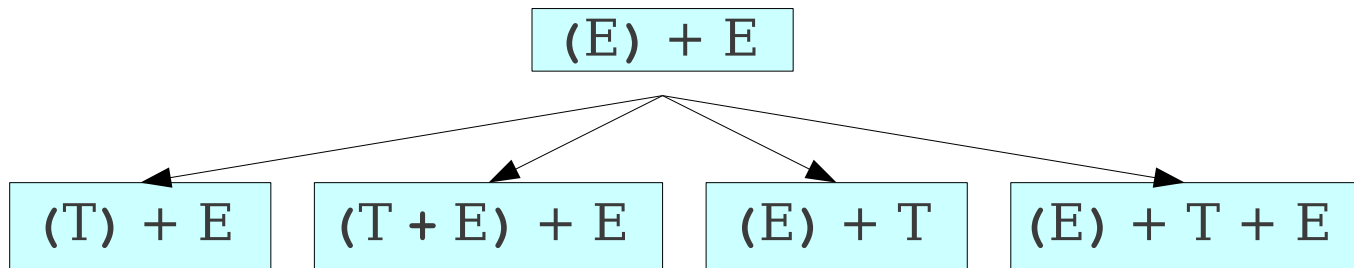
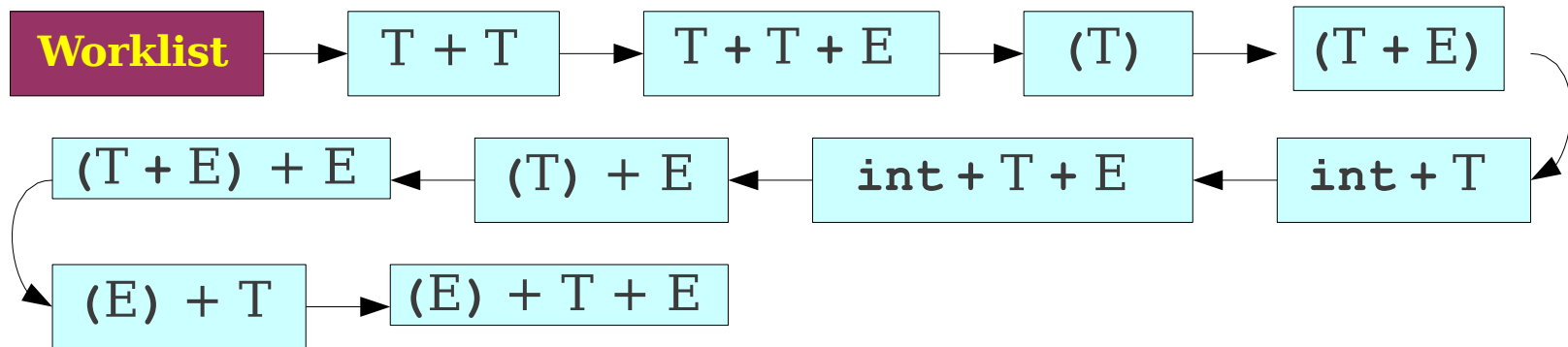
E → **T** + **E**

T → **int**

T → **(E)**

int + **int**

Breadth-First Search Parsing



E → **T**

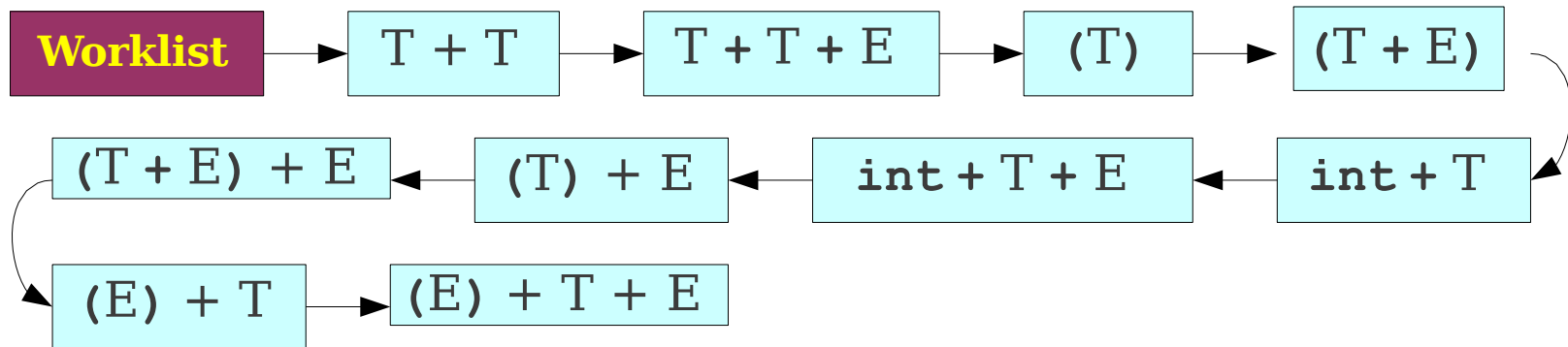
E → **T** + **E**

T → **int**

T → **(E)**

int + int

Breadth-First Search Parsing



E → **T**

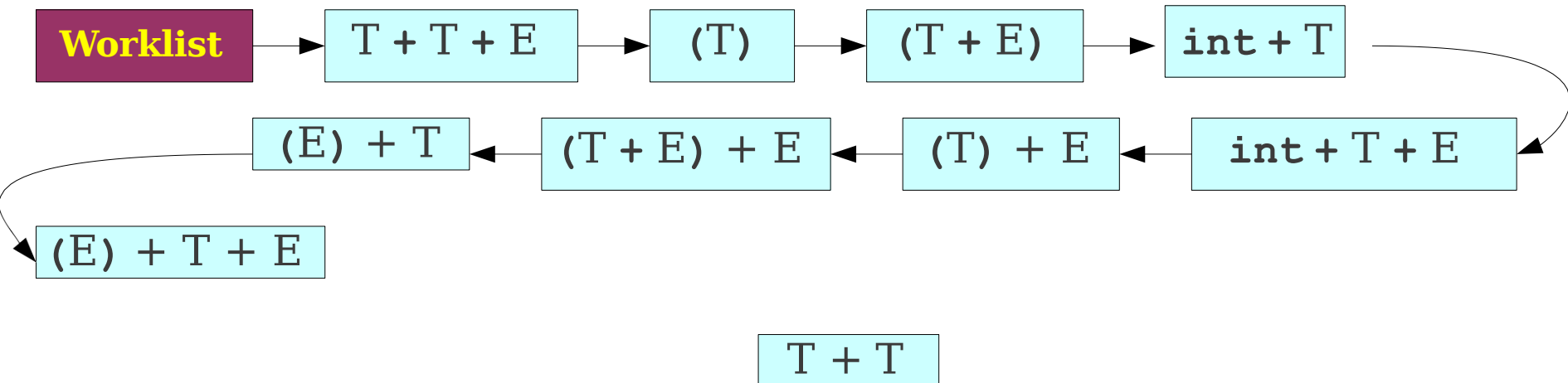
E → **T** + **E**

T → **int**

T → **(E)**

int + int

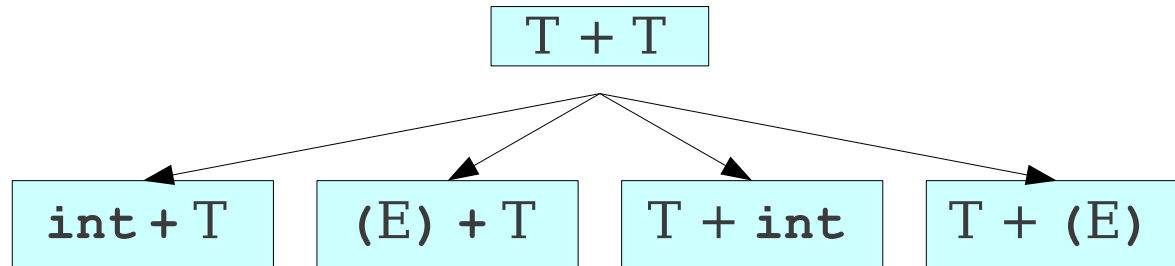
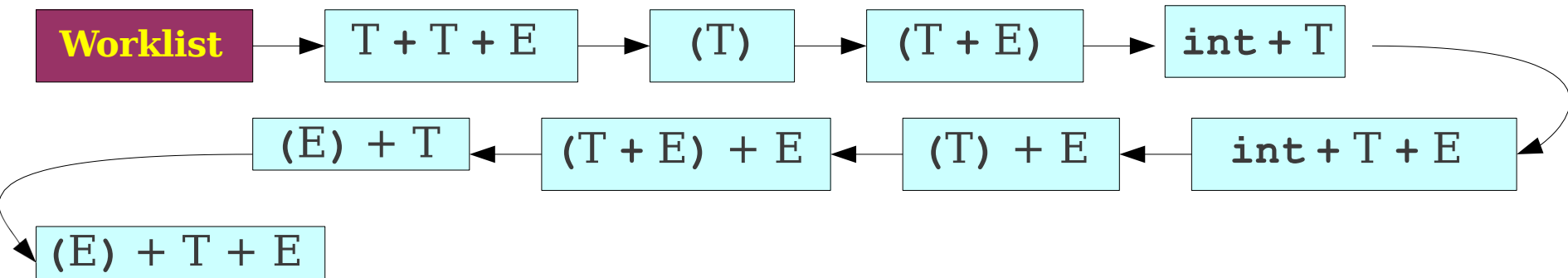
Breadth-First Search Parsing



E → **T**
E → **T** + **E**
T → **int**
T → **(E)**

int + int

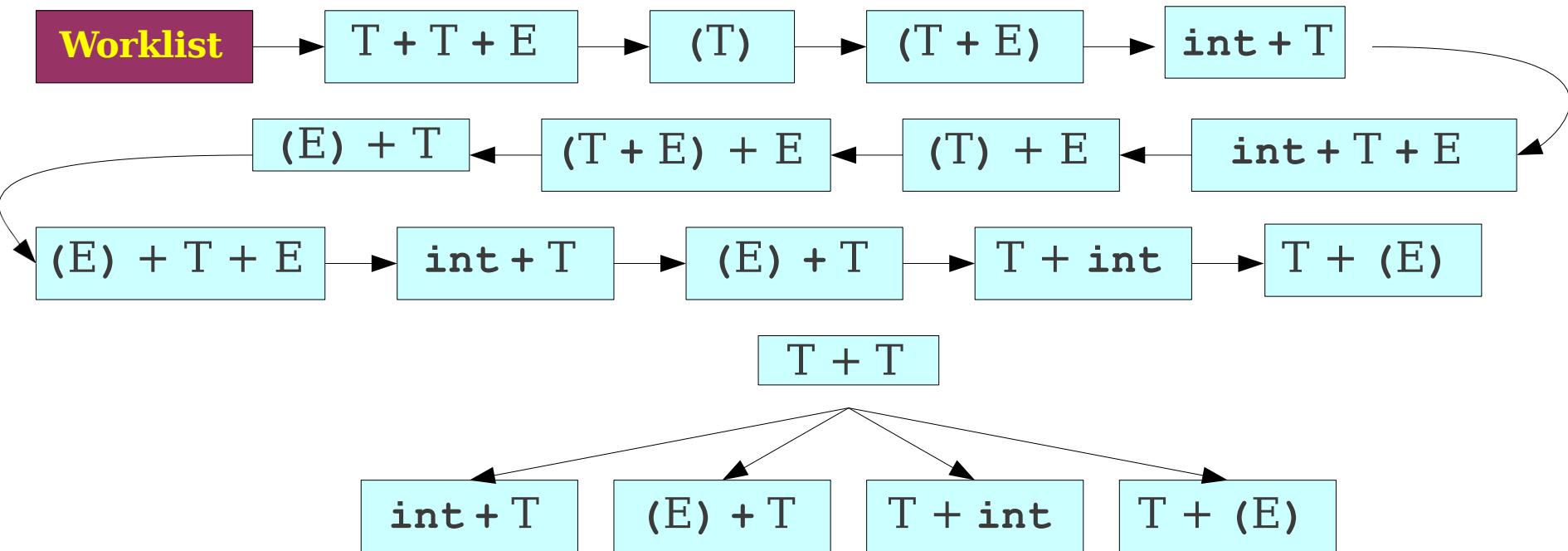
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

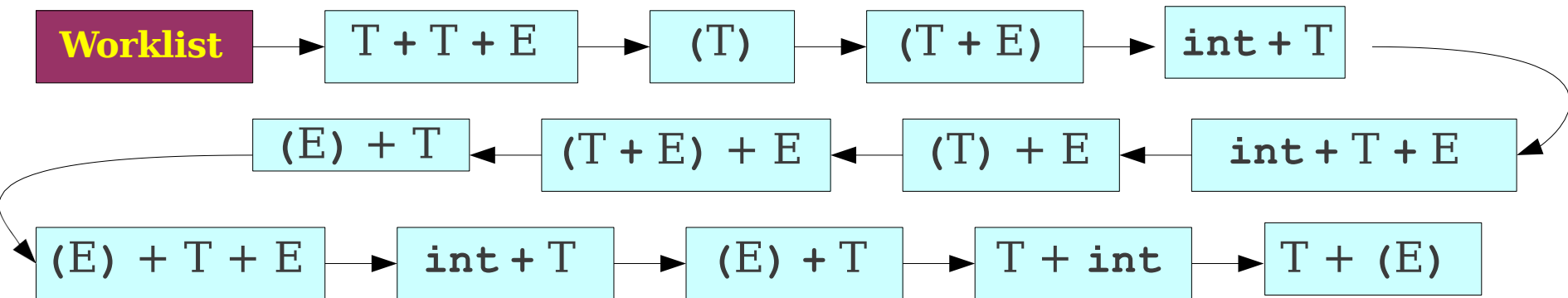
Breadth-First Search Parsing



E → **T**
E → **T + E**
T → **int**
T → **(E)**

int + int

Breadth-First Search Parsing



$E \rightarrow T$

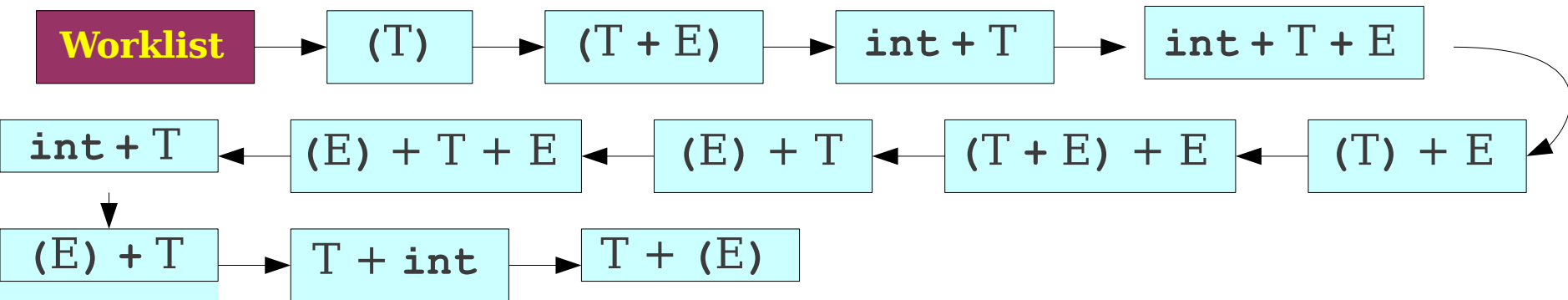
$E \rightarrow T + E$

$T \rightarrow int$

$T \rightarrow (E)$

$int + int$

Breadth-First Search Parsing

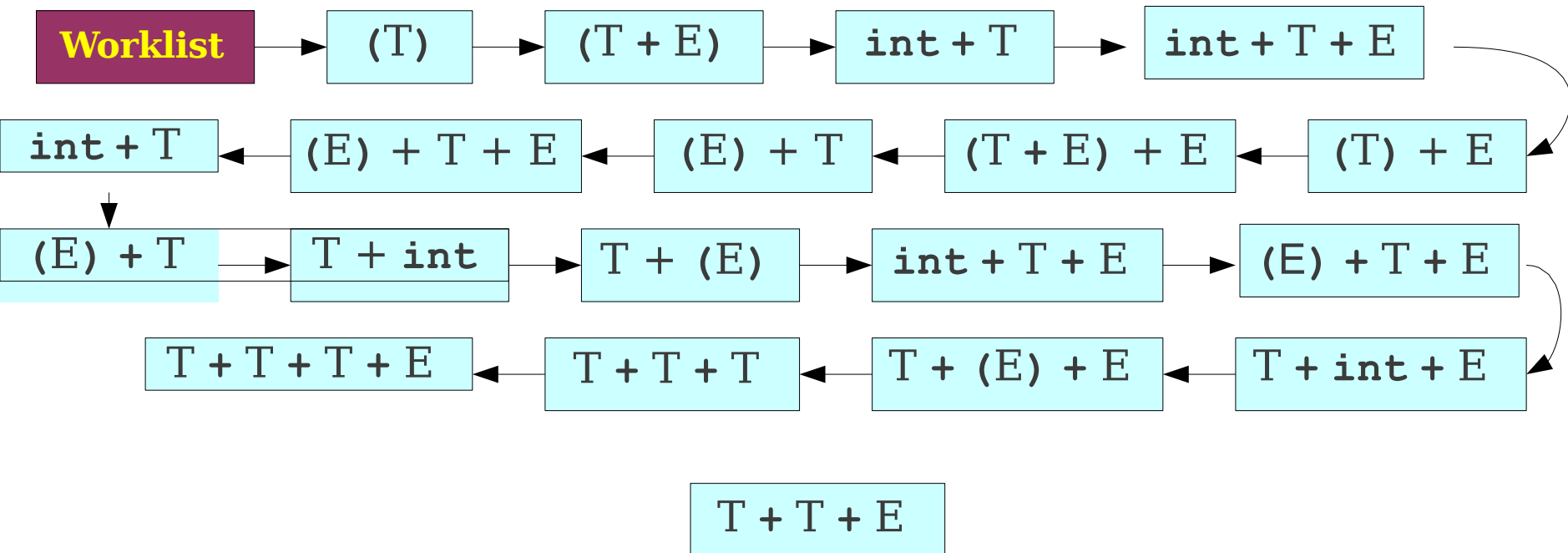


$T + T + E$

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

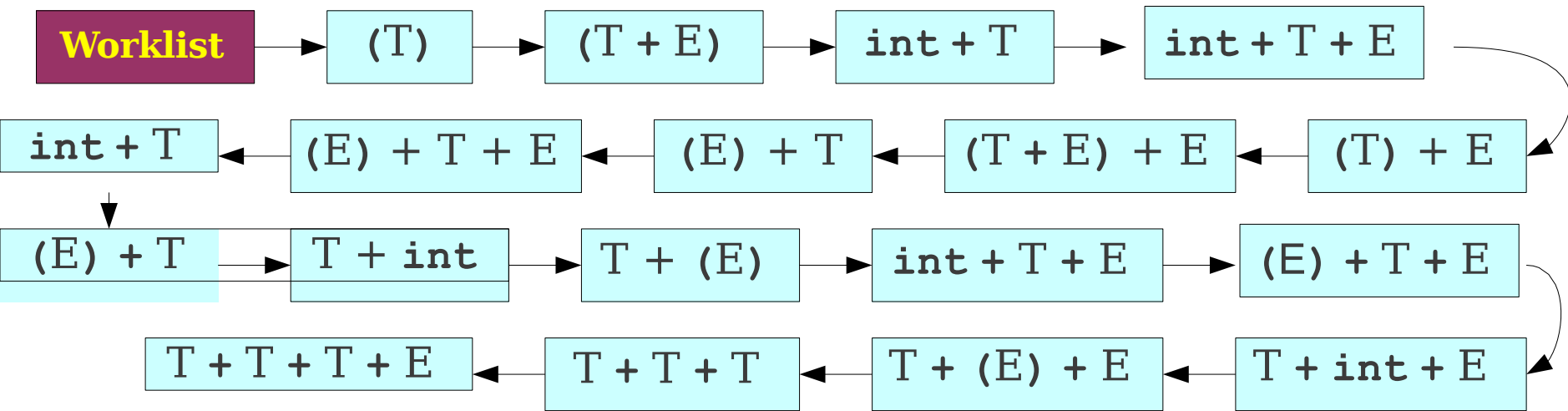
Breadth-First Search Parsing



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

Breadth-First Search Parsing



E → **T**

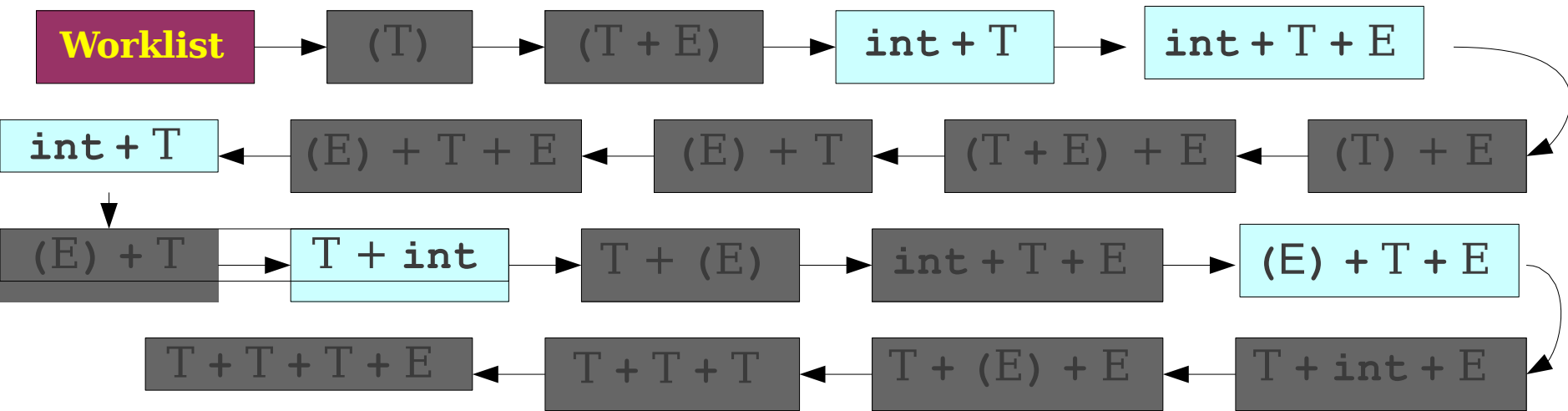
E → **T** + **E**

T → **int**

T → (**E**)

int + int

Breadth-First Search Parsing



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

int + int

BFS is Slow

- Enormous time and memory usage:
 - Lots of **wasted effort**:
 - Generates a lot of sentential forms that couldn't possibly match.
 - But in general, extremely hard to tell whether a sentential form can match – that's the job of parsing!
 - High **branching factor**:
 - Each sentential form can expand in (potentially) many ways for each nonterminal it contains.

Reducing Wasted Effort

- Suppose we're trying to match a string γ .
- Suppose we have a sentential form (derived from CFG) $\tau = \alpha\omega$, where α is a string of terminals and ω is a string of terminals and nonterminals.
- If α isn't a prefix of γ , then no string derived from τ can ever match γ .
- If we can find a way to try to get a prefix of terminals at the front of our sentential forms, then we can start pruning out impossible options.

Reducing the Branching Factor

- If a string has many nonterminals in it, the branching factor can be high.
 - Sum of the number of productions of each nonterminal involved.
- If we can restrict which productions we apply, we can keep the branching factor lower.

Leftmost Derivations

- Recall: A **leftmost derivation** is one where we always expand the leftmost symbol first.
- Updated algorithm:
 - Do a breadth-first search, **only considering leftmost derivations**.
 - Dramatically drops branching factor.
 - Increases likelihood that we get a prefix of nonterminals.
 - Prune sentential forms that can't possibly match.
 - Avoids wasted effort.

Leftmost BFS



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Leftmost BFS

Worklist

E

E → **T**

E → **T** + **E**

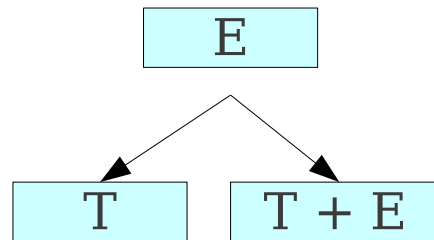
T → **int**

T → (**E**)

`int + int`

Leftmost BFS

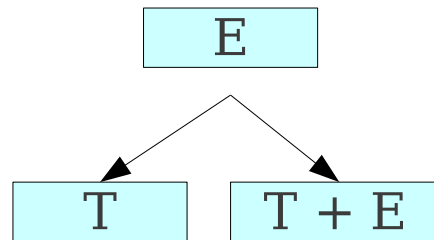
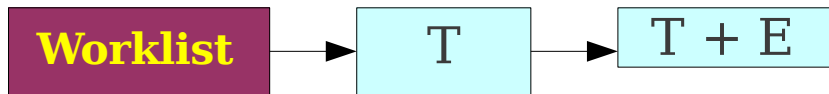
Worklist



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

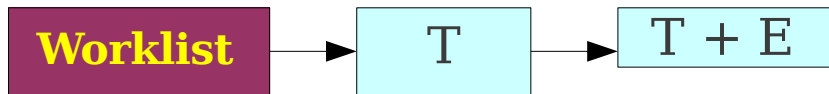
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Leftmost BFS



E → **T**

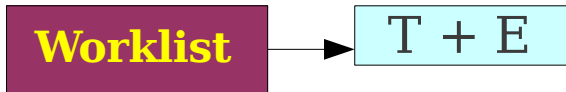
E → **T** + **E**

T → **int**

T → **(E)**

`int + int`

Leftmost BFS

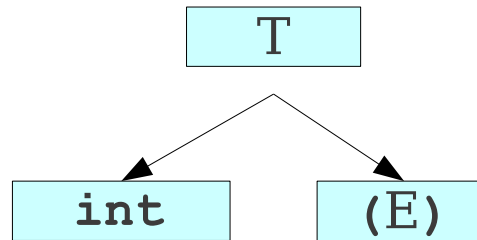
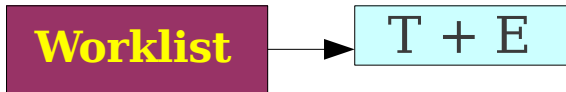


T

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

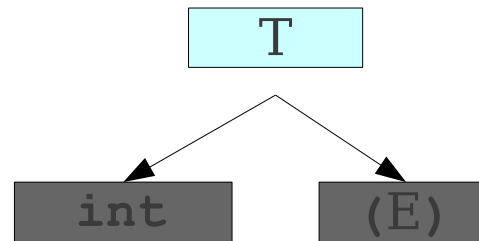
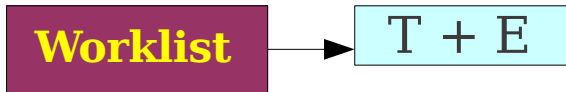
Leftmost BFS



E → **T**
E → **T** + **E**
T → int
T → (**E**)

int + int

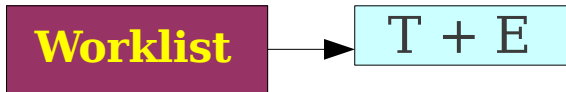
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Leftmost BFS



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Leftmost BFS

Worklist

T + E

E → **T**

E → **T** + **E**

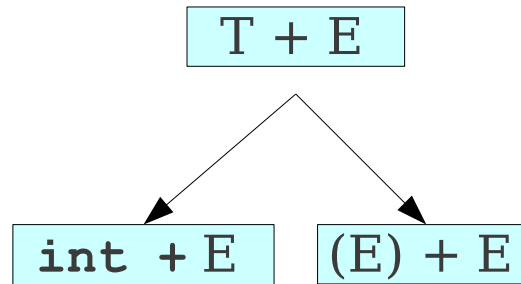
T → **int**

T → (**E**)

int + int

Leftmost BFS

Worklist

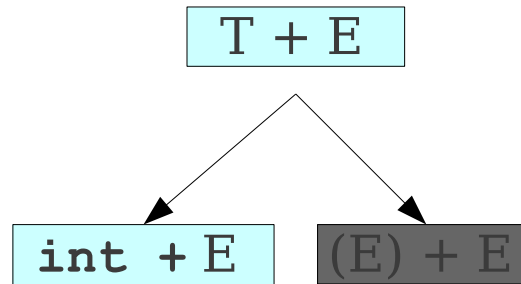


$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Leftmost BFS

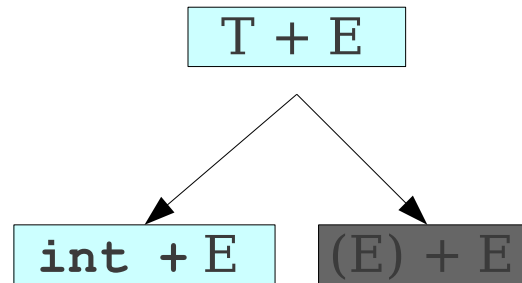
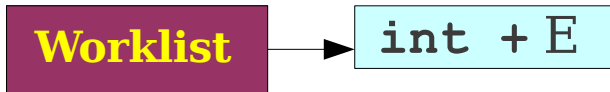
Worklist



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

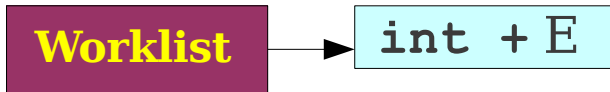
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

Leftmost BFS



E → **T**

E → **T** + **E**

T → int

T → (**E**)

int + int

Leftmost BFS

Worklist

int + E

E → **T**

E → **T** + **E**

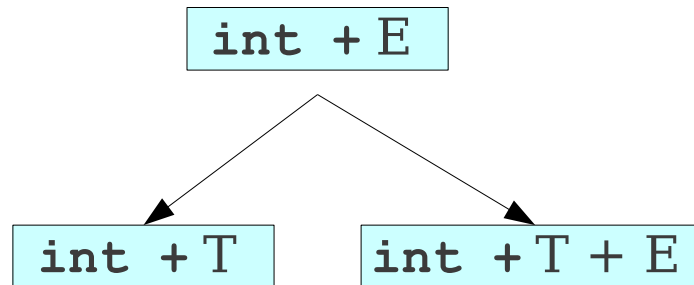
T → int

T → (**E**)

int + int

Leftmost BFS

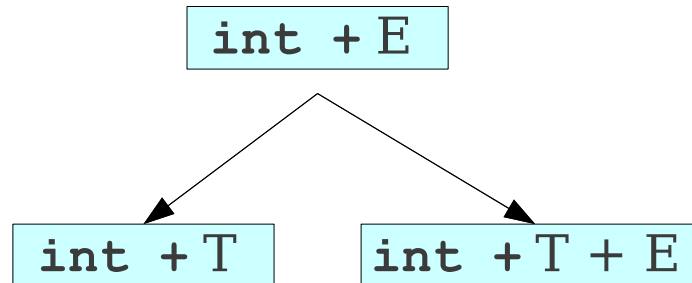
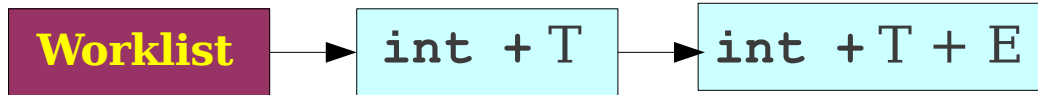
Worklist



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

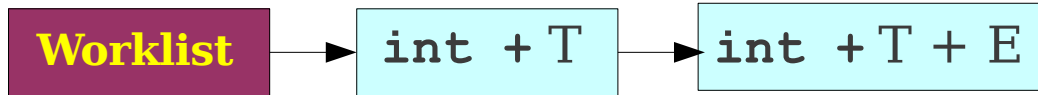
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

`int + int`

Leftmost BFS



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

int + int

Leftmost BFS

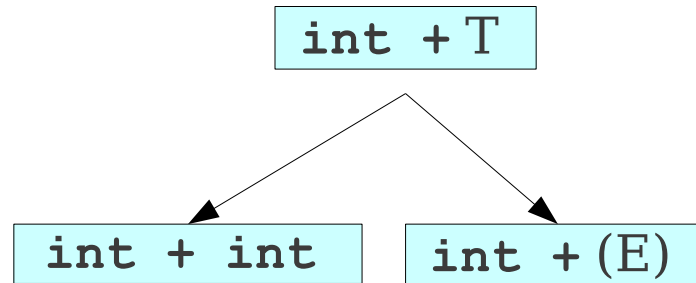
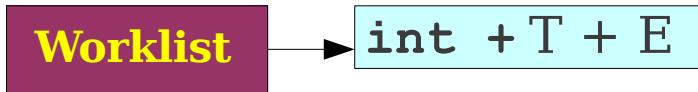
Worklist → int + T + E

int + T

E → **T**
E → **T** + **E**
T → int
T → (**E**)

int + int

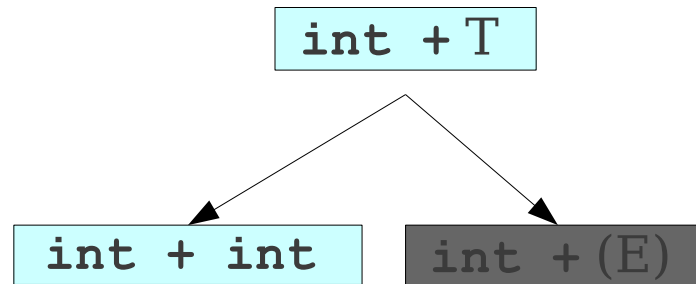
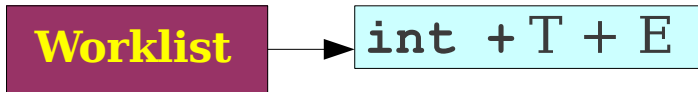
Leftmost BFS



$E \rightarrow T'$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

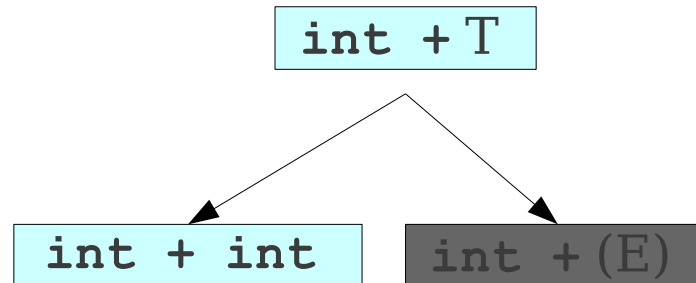
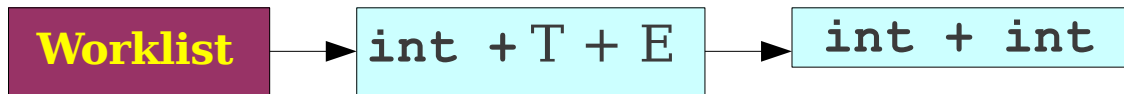
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

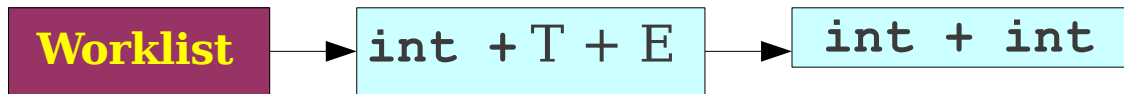
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

$\text{int} + \text{int}$

Leftmost BFS



E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

int + int

Leftmost BFS

Worklist → int + int

int + T + E

E → **T**

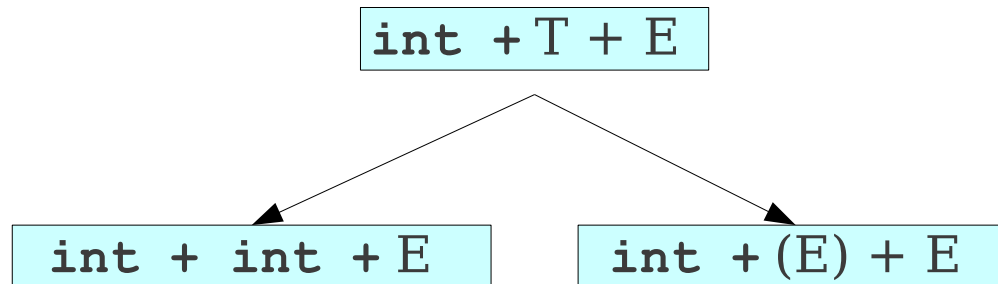
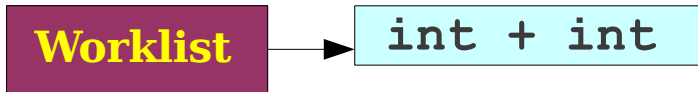
E → **T** + **E**

T → int

T → (**E**)

int + int

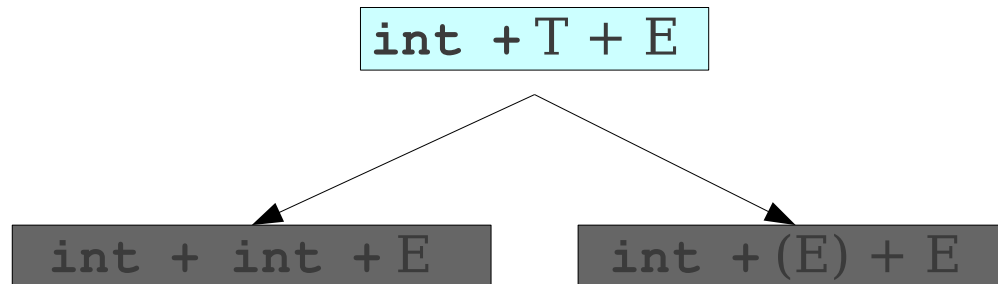
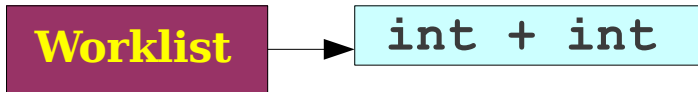
Leftmost BFS



$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

int + int

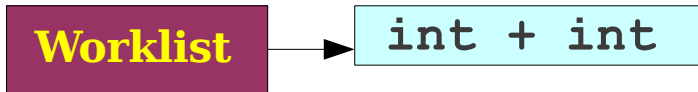
Leftmost BFS



E → **T**
E → **T** + **E**
T → int
T → (**E**)

int + int

Leftmost BFS



E → **T**

E → **T** + **E**

T → int

T → (**E**)

int + int

Leftmost BFS

Worklist

int + int

E → **T**

E → **T** + **E**

T → int

T → (**E**)

int + int

Leftmost BFS

Worklist



int + int

E → **T**

E → **T** + **E**

T → int

T → (**E**)

int + int

Leftmost BFS

- Substantial improvement over naïve algorithm.
- Will always find a valid parse of a program if one exists.
- Can easily be modified to find if a program can't be parsed.
- But, there are still problems.

Leftmost BFS Has Problems

Worklist

$A \rightarrow Aa \mid Ab \mid c$

Leftmost BFS Has Problems

Worklist

A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems

Worklist

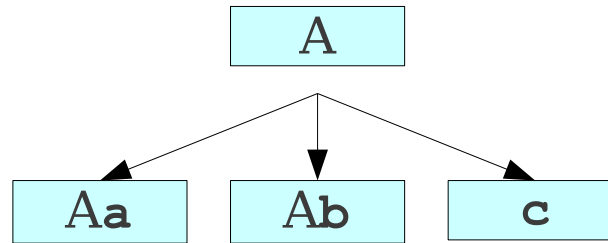
A

A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems

Worklist

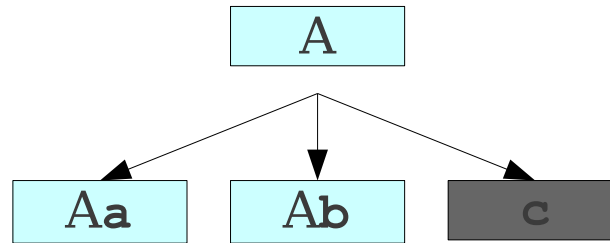


A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems

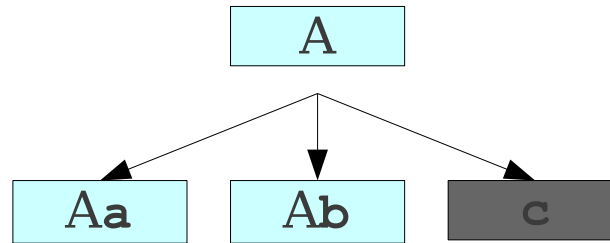
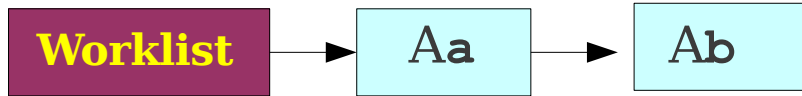
Worklist



A → **A**a | **A**b | c

caaaaaaaaaa

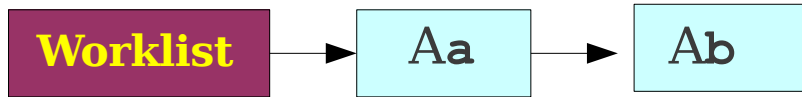
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaaa

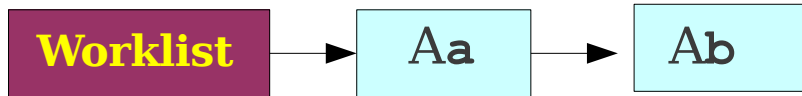
Leftmost BFS Has Problems



A → **A****a** | **A****b** | **c**

caaaaaaaaaa

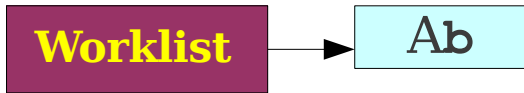
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaaa

Leftmost BFS Has Problems

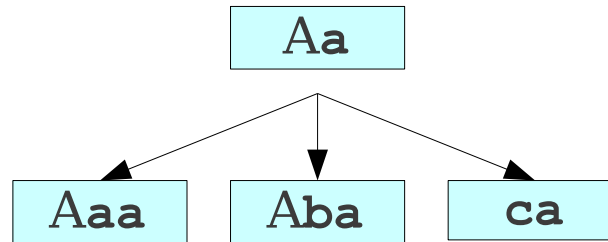


Aa

A → **A**a | **A**b | c

caaaaaaaaaaa

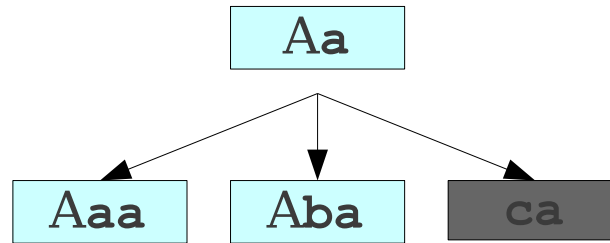
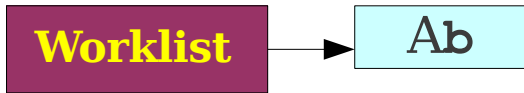
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

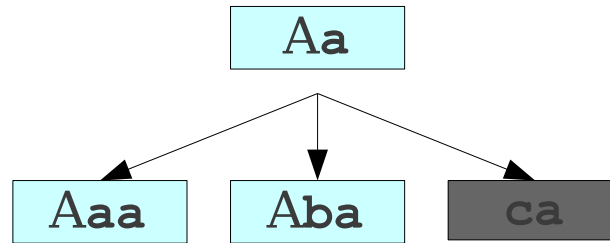
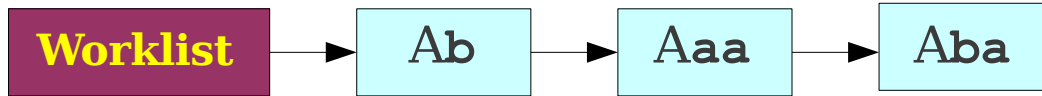
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaaa

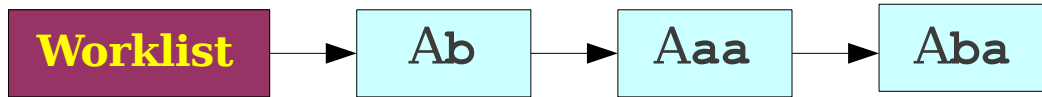
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

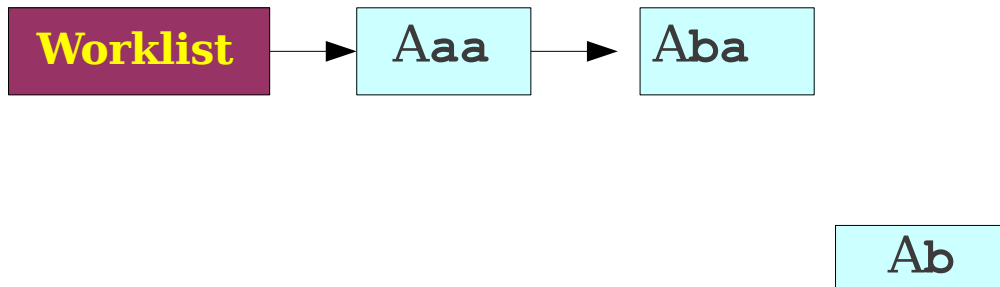
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

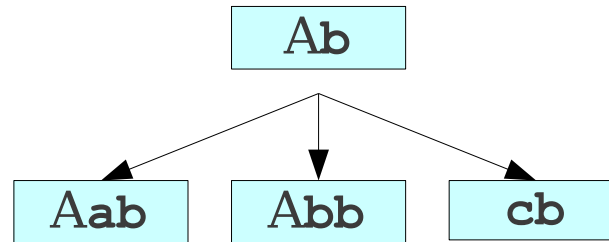
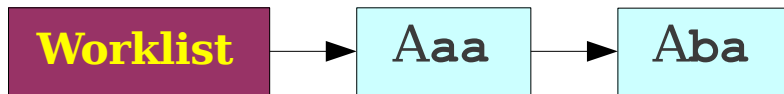
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

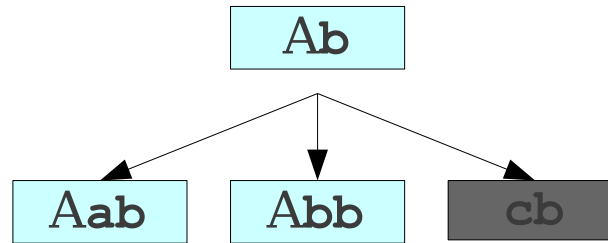
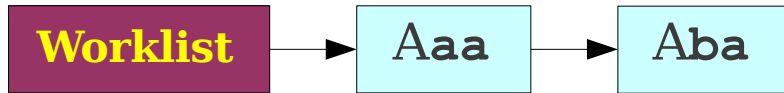
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

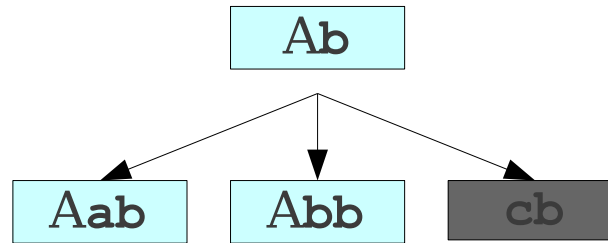
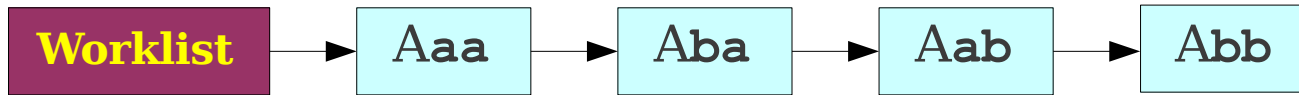
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

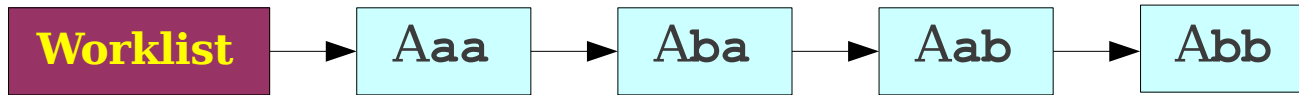
Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

Leftmost BFS Has Problems



A → **A**a | **A**b | c

caaaaaaaaaa

Problems with Leftmost BFS

- Grammars like this can make parsing take **exponential time**.
- Also uses **exponential memory**.
- What if we search the graph with a different algorithm?

Leftmost DFS

- Idea: Use **depth-first** search.
- Advantages:
 - **Lower memory** usage: Only considers one branch at a time.
 - **High performance**: On many grammars, runs very quickly.
 - **Easy to implement**: Can be written as a set of mutually recursive functions.

Leftmost DFS

E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

Leftmost DFS

E → **T**

E → **T** + **E**

T → **int**

T → (**E**)

`int + int`

Leftmost DFS

E

E → **T**

E → **T** + **E**

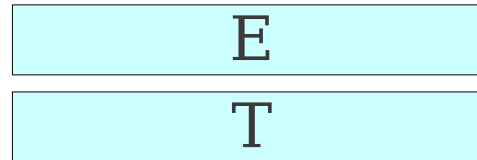
T → **int**

T → (**E**)

`int + int`

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

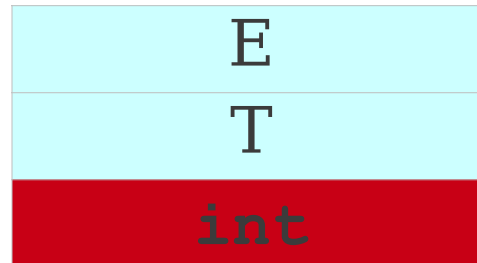
E → **T**
E → **T** + **E**
T → **int**
T → (**E**)

| |
|-----|
| E |
| T |
| int |

int + int

Leftmost DFS

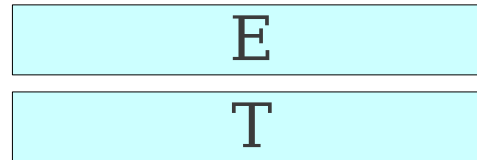
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



int + int

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

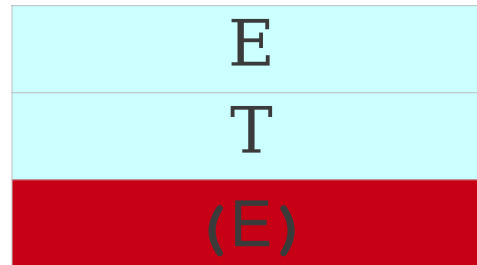
E → **T**
E → **T** + **E**
T → **int**
T → (**E**)

| |
|-----|
| E |
| T |
| (E) |

`int + int`

Leftmost DFS

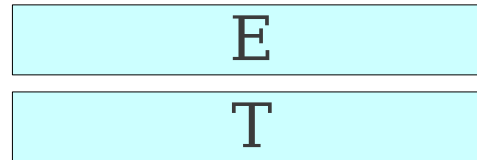
E → **T**
E → **T** + **E**
T → **int**
T → (**E**)



int + int

Leftmost DFS

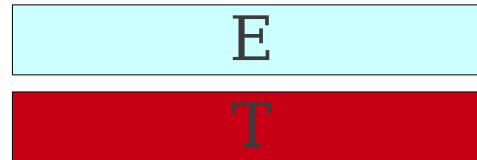
$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

E

E → **T**

E → **T** + **E**

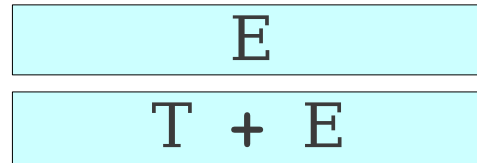
T → **int**

T → (**E**)

`int + int`

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$



`int + int`

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

| |
|---------|
| E |
| T + E |
| int + E |

int + int

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

| |
|---------|
| E |
| T + E |
| int + E |
| int + T |

int + int

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

| |
|-----------|
| E |
| T + E |
| int + E |
| int + T |
| int + int |

int + int

Leftmost DFS

$E \rightarrow T$
 $E \rightarrow T + E$
 $T \rightarrow \text{int}$
 $T \rightarrow (E)$

| |
|-----------|
| E |
| T + E |
| int + E |
| int + T |
| int + int |



int + int

Problems with Leftmost DFS

A → **A**a | c

| |
|-------|
| A |
| Aa |
| Aaa |
| Aaaa |
| Aaaaa |



c