

Languages, Grammars, Regular Expressions

Synopsis

- ◇ **Languages:** how to define
- ◇ **Grammars**
 - Specification of languages
 - Language generators
 - Derivations
 - Analysis
 - Language Structure
- ◇ **Regular Languages**
 - Regular grammars
 - Deterministic Finite Automata
 - Regular expressions
 - Use of REs in Ruby

Languages

- ◇ Symbols: elements of an alphabet, typically ASCII symbols
- ◇ Strings: sequences of symbols.
- ◇ Language: *set of strings*

Examples of languages:

$$\{a, ab, aba, bab, bbb\}$$
$$\{0, 1, \dots, 9, 10, 11, \dots, 19, 21, \dots\}$$
$$\left\{ \text{int main()}\{\text{printf("Hello world");}\} , \dots \right\}$$

Languages

- ◇ Symbols: elements of an alphabet, typically ASCII symbols
- ◇ Strings: sequences of symbols.
- ◇ Language: *set of strings*

The language C is the set of all possible C programs!



Examples of languages:

$\{a, ab, aba, bab, bbb\}$

$\{0, 1, \dots, 9, 10, 11, \dots, 19, 21, \dots\}$

$\left\{ \text{int main()}\{\text{printf("Hello world");}\} , \dots \right\}$

Specification of Languages

- ◇ Enumeration is not an efficient way of specifying a language.
- ◇ Must be finitary.
- ◇ Must provide an efficient way of deciding whether a string belongs to the language or not.
- ◇ Must capture the structure of the language.
 - ◇ For PL, execution must take structure into account.
- ◇ **Grammars:** standard mechanism of specifying programming languages

Grammars

- ◇ Mathematical concept in the area of *formal languages*
- ◇ Language generator: specifies a mechanism for generating all elements of the language.
 - ◇ The language is the set of strings that can be generated.
 - ◇ Can be converted into an *acceptor*: a procedure that tests whether a string is in the language or not.
- ◇ Practitioner's approach: introduce concepts step-by-step

Production Rules

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$

Production Rules

Production rule

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

Production Rules

$$\begin{array}{lll} E & \rightarrow & (E) \\ E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & D \\ D & \rightarrow & 0D \\ D & \rightarrow & 1D \\ D & \rightarrow & 0 \\ D & \rightarrow & 1 \end{array}$$

Set of production rules



Production Rules

Left hand side

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

Right hand side

Left hand side can be *rewritten* as the right hand side.

Production Rules

$$\begin{array}{lcl} E & \rightarrow & (E) \\ E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & D \\ D & \rightarrow & 0D \\ D & \rightarrow & 1D \\ D & \rightarrow & 0 \\ D & \rightarrow & 1 \end{array}$$

Non-terminals

- ◇ Denoted by **capitals**
- ◇ Not part of the generated language.
- ◇ An aid to generating the language.

Production Rules

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$

Terminals:

- ◇ The rest of the symbols.
- ◇ Part of the generated language.

Production Rules

Derivation:

$E \rightarrow$

$$\begin{array}{lcl} E & \rightarrow & (E) \\ E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & D \\ D & \rightarrow & 0D \\ D & \rightarrow & 1D \\ D & \rightarrow & 0 \\ D & \rightarrow & 1 \end{array}$$

Production Rules

Derivation:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

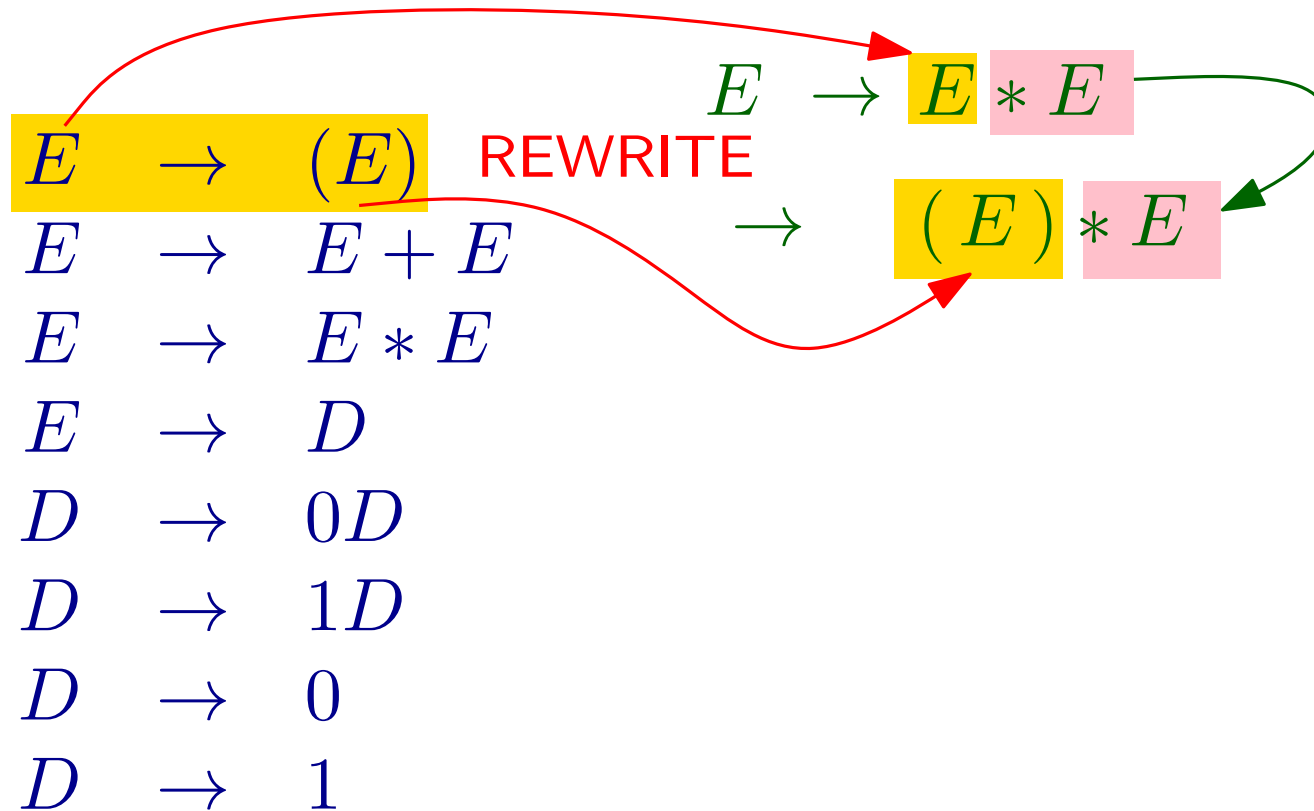
$$D \rightarrow 0$$

$$D \rightarrow 1$$

$$E \rightarrow E * E$$

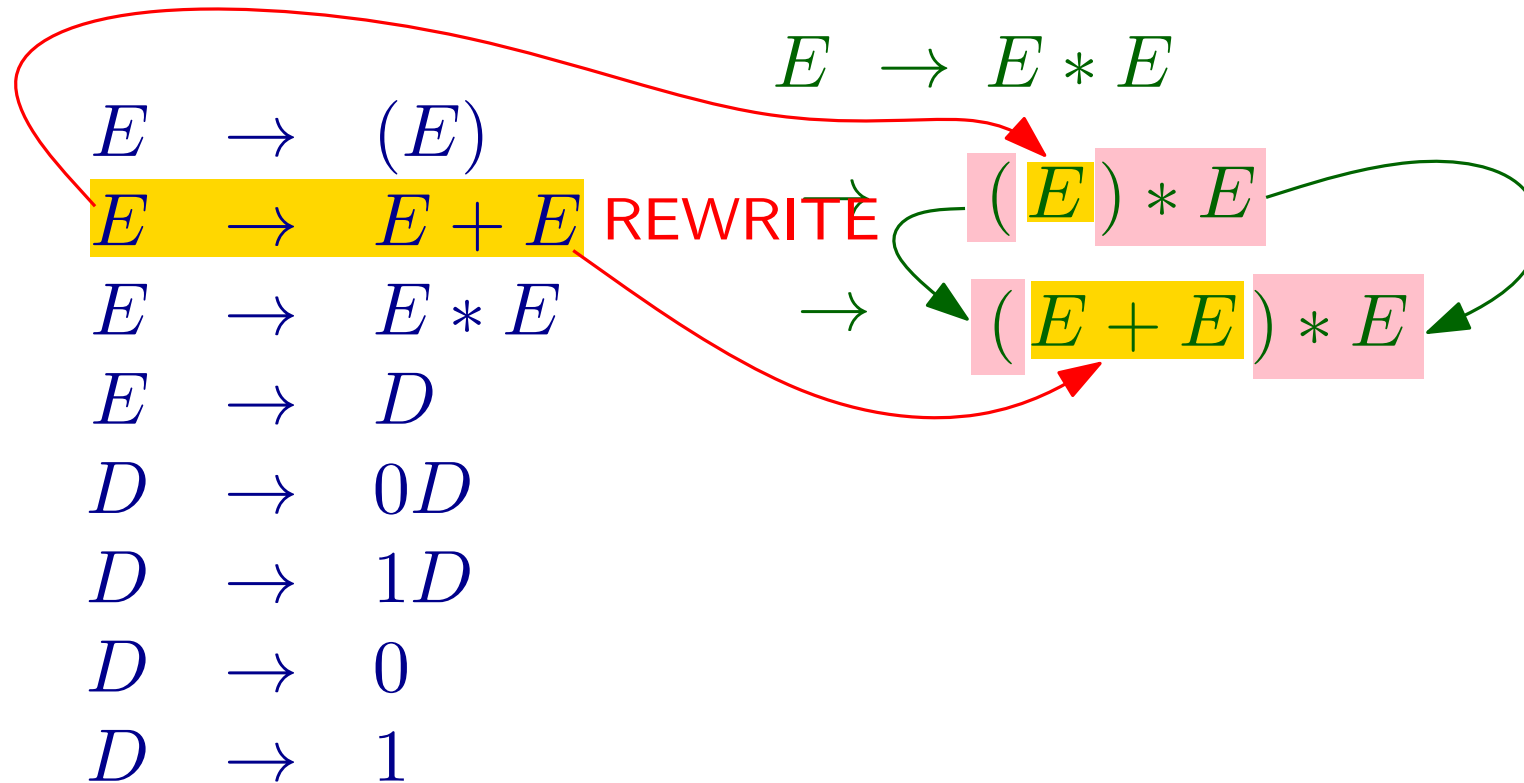

Production Rules

Derivation:



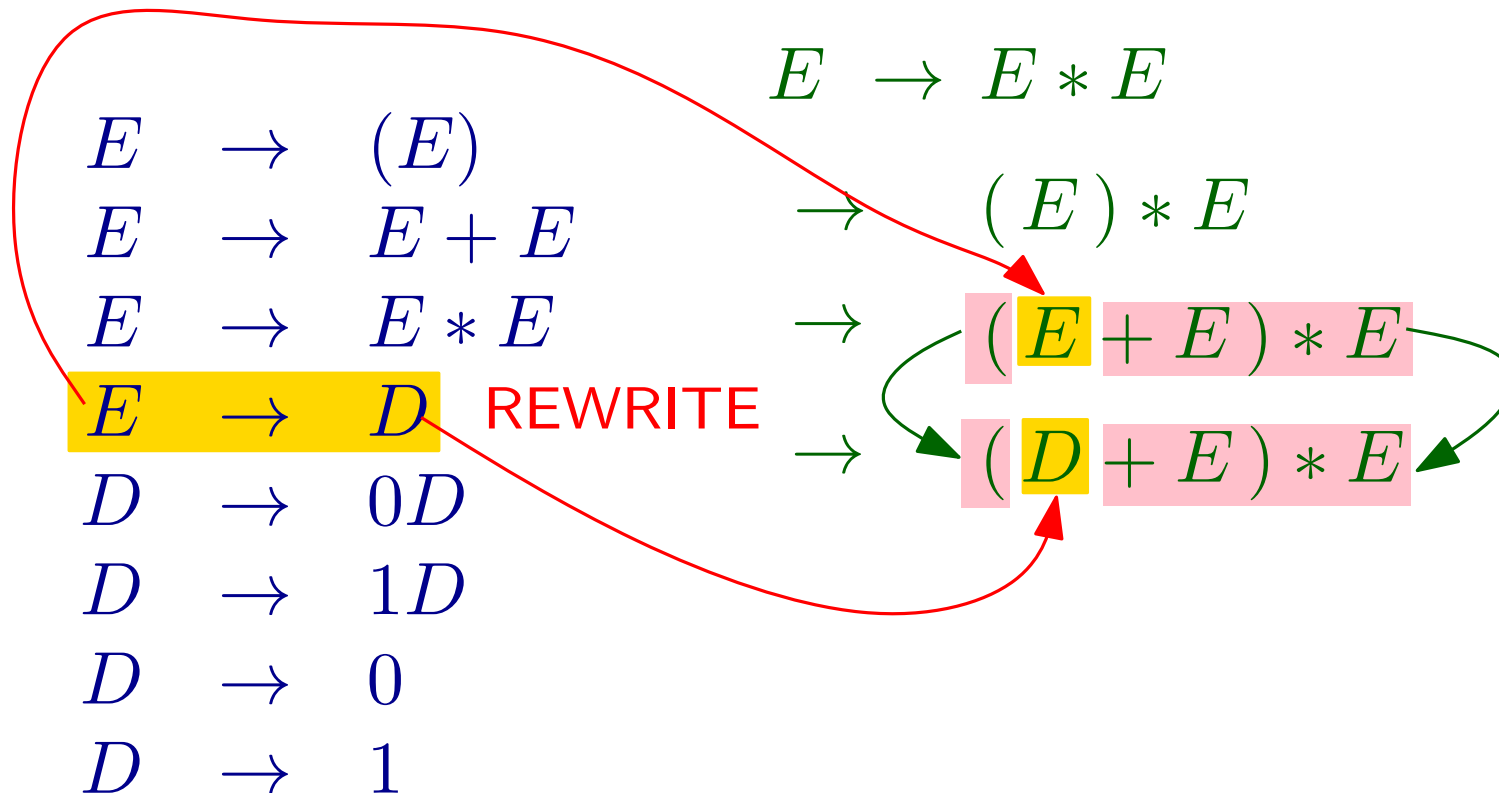
Production Rules

Derivation:



Production Rules

Derivation:



Production Rules

Derivation:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1 \text{ REWRITE}$$

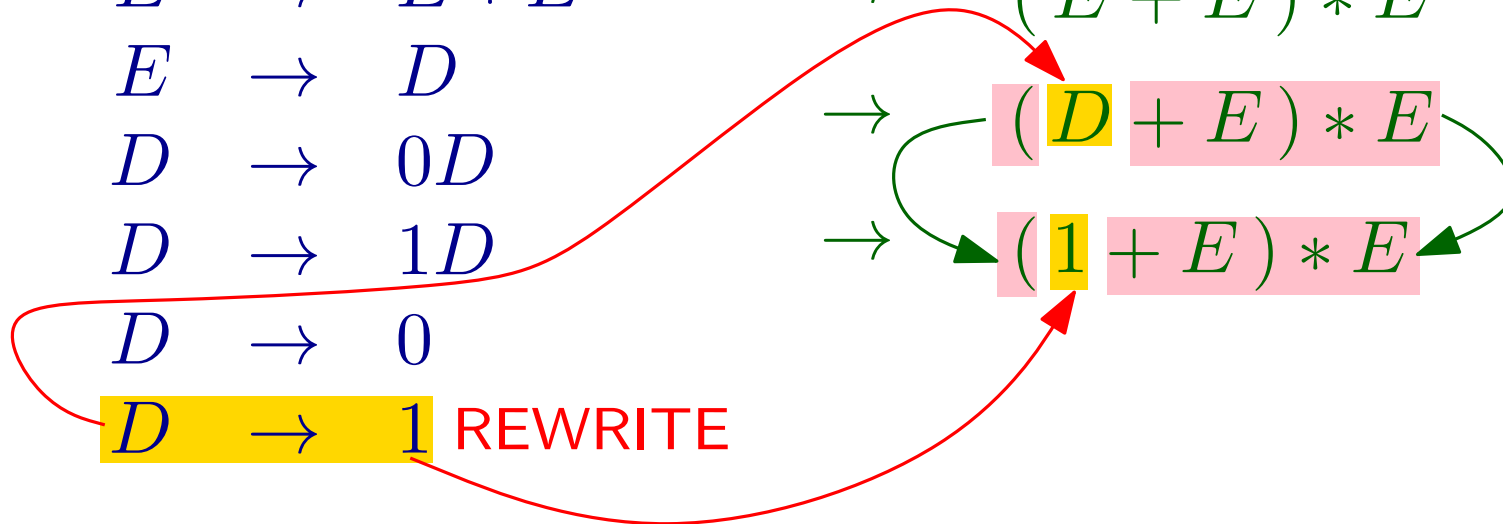
$$E \rightarrow E * E$$

$$\rightarrow (E) * E$$

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

$$\rightarrow (D + E) * E$$

$$\rightarrow (1 + E) * E$$



Production Rules

Derivation:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D \text{ REWRITE}$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

$$E \rightarrow E * E$$

$$\rightarrow (E) * E$$

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

$$\rightarrow (D + E) * E$$

$$\rightarrow (1 + E) * E$$

$$\rightarrow (1 + D) * E$$

Production Rules

Derivation:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D \text{ REWRITE}$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

$$E \rightarrow E * E$$

$$\rightarrow (E) * E$$

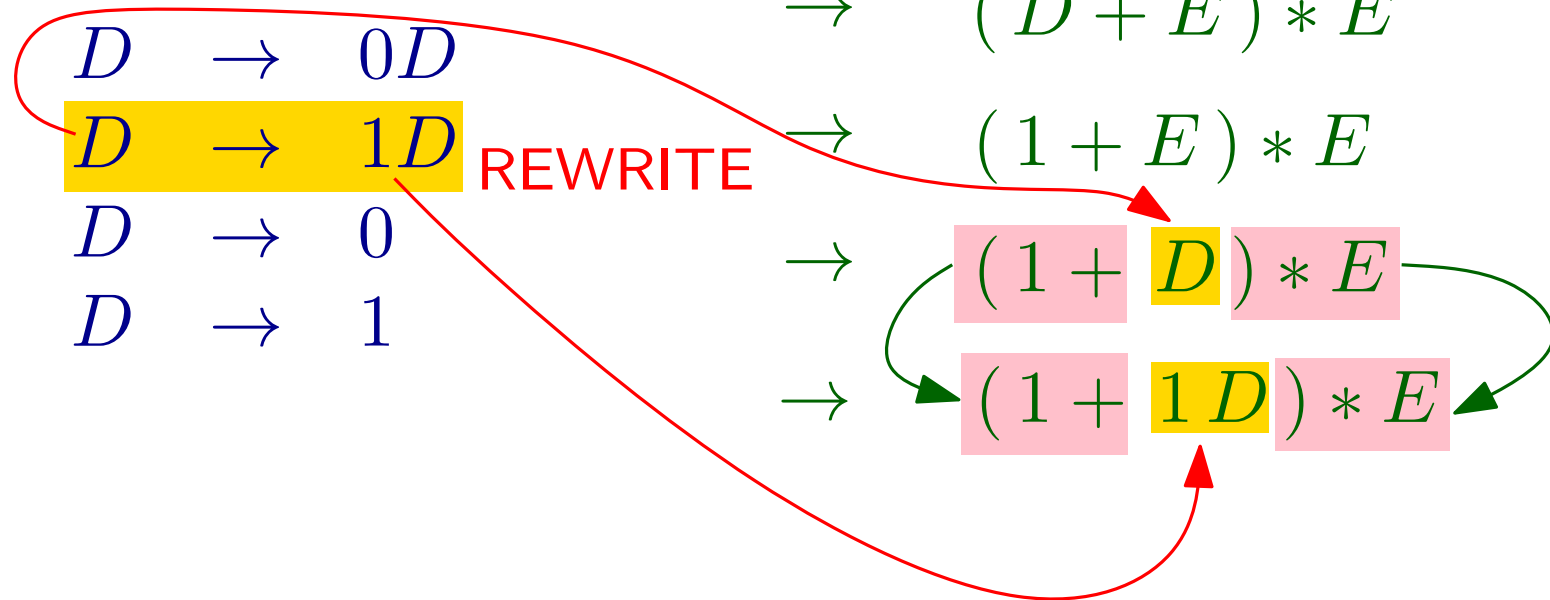
$$\rightarrow (E + E) * E$$

$$\rightarrow (D + E) * E$$

$$\rightarrow (1 + E) * E$$

$$\rightarrow (1 + D) * E$$

$$\rightarrow (1 + 1D) * E$$



Production Rules

Derivation:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0 \text{ REWRITE}$$

$$D \rightarrow 1$$

$$E \rightarrow E * E$$

$$\rightarrow (E) * E$$

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

$$\rightarrow (D + E) * E$$

$$\rightarrow (1 + E) * E$$

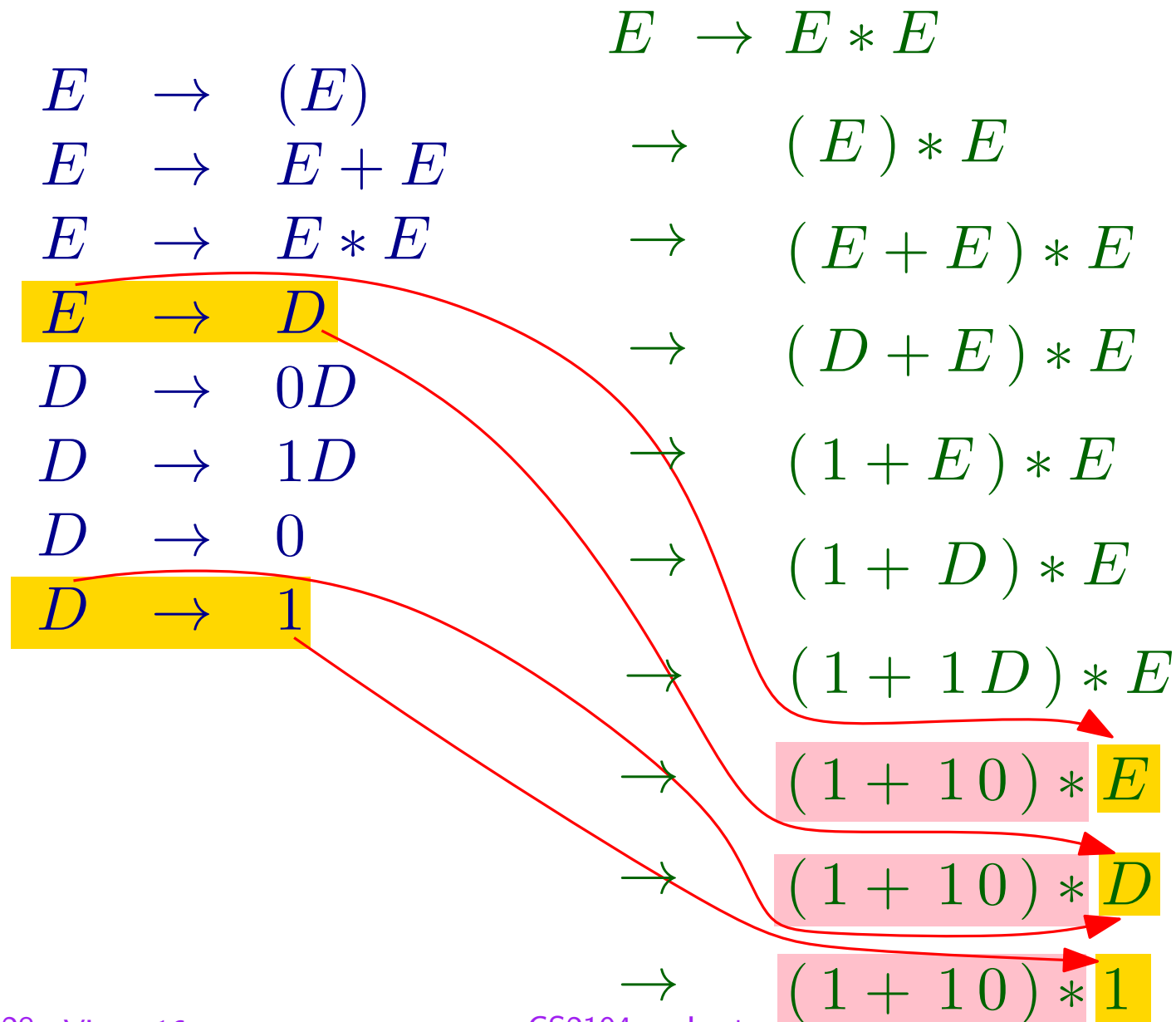
$$\rightarrow (1 + D) * E$$

$$\rightarrow (1 + 1D) * E$$

$$\rightarrow (1 + 10) * E$$

Production Rules

Derivation:



Production Rules

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$

Derivation:

$$(1 + 10) * 1 \in \mathcal{L}(E)$$

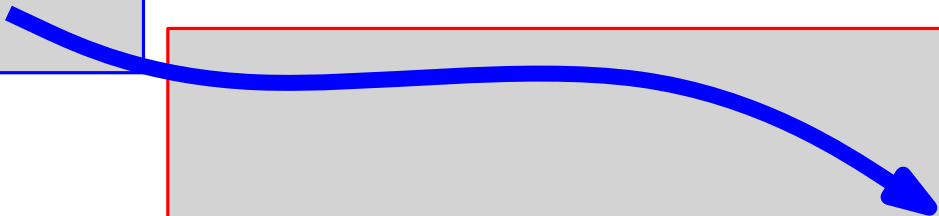
$$\begin{aligned} &\rightarrow (1 + E) * E \\ &\rightarrow (1 + D) * E \\ &\rightarrow (1 + 1D) * E \\ &\rightarrow (1 + 10) * E \\ &\rightarrow (1 + 10) * D \\ &\rightarrow (1 + 10) * 1 \end{aligned}$$

Production Rules

Language generated by
non-terminal E

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$

Derivation:


$$\underbrace{(1 + 10) * 1}_{\text{terminals only}} \in \mathcal{L}(E)$$

Similarly:

$$1001 \in \mathcal{L}(D)$$

$$\rightarrow (1 + 10) * E$$

$$\rightarrow (1 + 10) * D$$

$$\rightarrow (1 + 10) * 1$$

Formal Definition of Grammar

Grammar: tuple $G \equiv \langle \Sigma, N, \Pi, S \rangle$

Σ - alphabet of *terminal symbols*

N - set of *non-terminal symbols*

Π - set of production rules

S - *start non-terminal*, $S \in N$

$\mathcal{L}(G) \equiv \mathcal{L}(S)$ (contains only strings of terminal symbols)

In practice...

Rules are enough in practice.

Grammar can be derived from the rules:

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

In practice...

Rules are enough in practice.

Grammar can be derived from the rules:

$$\Sigma = \{ (,), +, *, 0, 1 \}$$

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

In practice...

Rules are enough in practice.

Grammar can be derived from the rules:

$$\Sigma = \{ (,), +, *, 0, 1 \}$$

$$N = \{ E, D \}$$

Diagram illustrating the derivation of grammar rules from a set of rules. Red arrows connect the rules to the sets Σ and N .

- $E \rightarrow (E)$ connects to Σ (parentheses) and N (E).
- $E \rightarrow E + E$ connects to Σ ($+$) and N (E).
- $E \rightarrow E * E$ connects to Σ ($*$) and N (E).
- $E \rightarrow D$ connects to N (D).
- $D \rightarrow 0D$ connects to Σ (0) and N (D).
- $D \rightarrow 1D$ connects to Σ (1) and N (D).
- $D \rightarrow 0$ connects to Σ (0).
- $D \rightarrow 1$ connects to Σ (1).

In practice...

Rules are enough in practice.

Grammar can be derived from the rules:

$$\Sigma = \{ (,), +, *, 0, 1 \}$$

$$N = \{ E, D \}$$

$$\Pi = \{ E \rightarrow (E), \dots \}$$

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$

In practice...

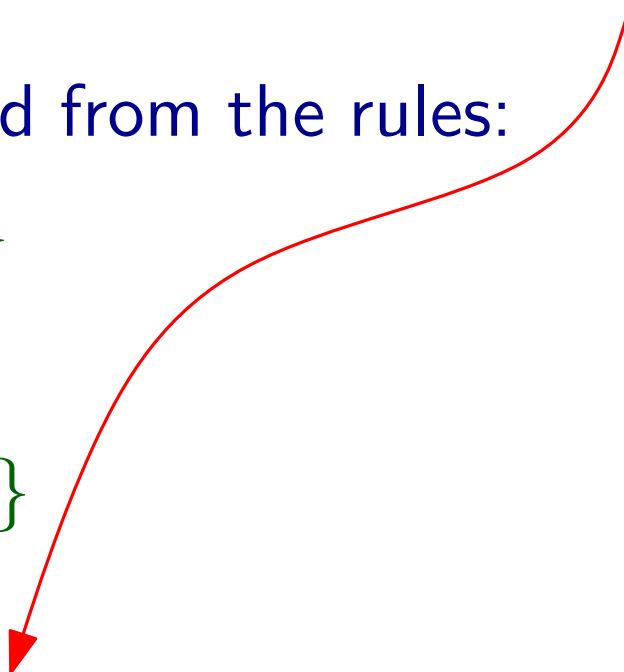
Rules are enough in practice.

Grammar can be derived from the rules:

$$\Sigma = \{ (,), +, *, 0, 1 \}$$

$$N = \{ E, D \}$$

$$\Pi = \{ E \rightarrow (E), \dots \}$$


$$\begin{array}{lcl} E & \rightarrow & (E) \\ E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & D \\ D & \rightarrow & 0D \\ D & \rightarrow & 1D \\ D & \rightarrow & 0 \\ D & \rightarrow & 1 \end{array}$$

Start non-terminal: E (the LHS of the first rule)

In practice...

Rules are enough in practice.

Grammar can be derived from the rules:

$$\Sigma = \{ (,), +, *, 0, 1 \}$$

$$N = \{ E, D \}$$

$$\Pi = \{ E \rightarrow (E), \dots \}$$

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$

Start non-terminal: E (the LHS of the first rule)

$$G = \langle \Sigma, N, \Pi, E \rangle$$

Derivation Tree

$$\begin{array}{lll} E & \rightarrow & (E) \\ E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & D \\ D & \rightarrow & 0D \\ D & \rightarrow & 1D \\ D & \rightarrow & 0 \\ D & \rightarrow & 1 \end{array} \quad E$$

$$(\quad 1 \quad + \quad 1 \quad 0 \quad) \quad * \quad 1$$

Derivation Tree

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

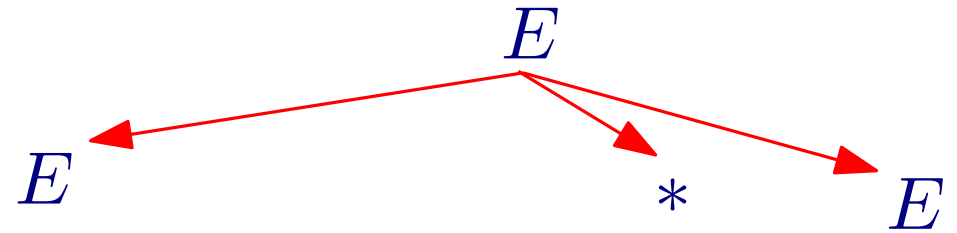
$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$



(1 + 1 0) * 1

Derivation Tree

$E \rightarrow (E)$

$E \rightarrow E + E$

$E \rightarrow E * E$

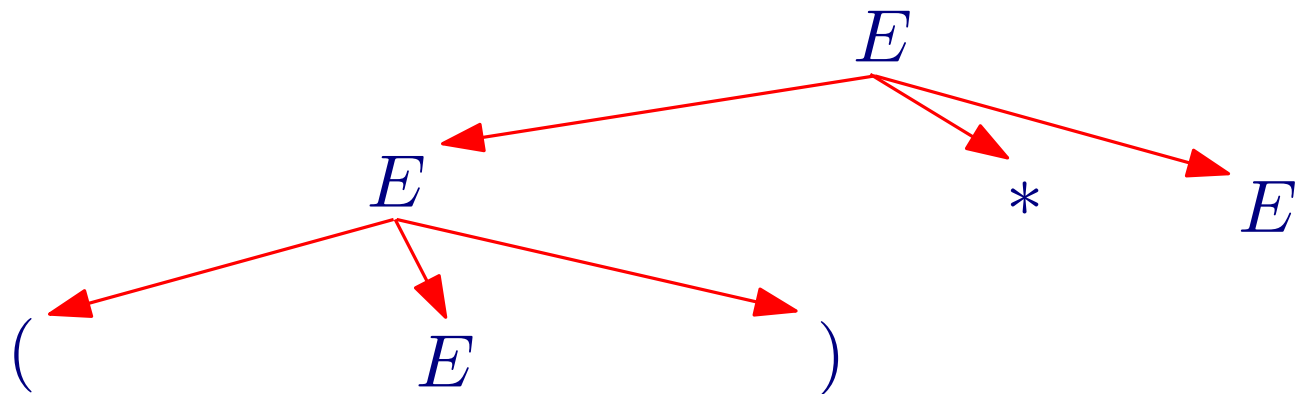
$E \rightarrow D$

$D \rightarrow 0D$

$D \rightarrow 1D$

$D \rightarrow 0$

$D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

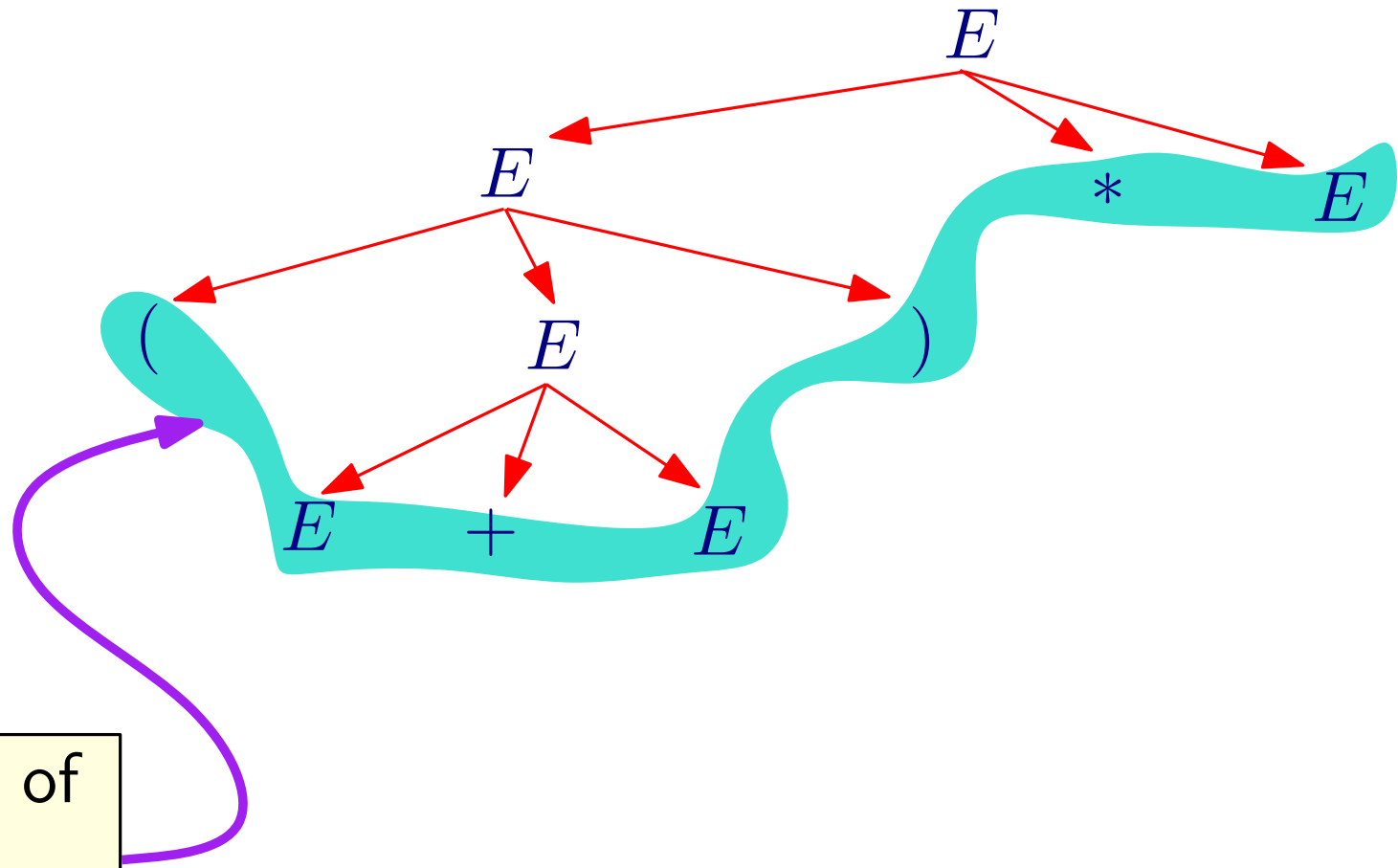
$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

$$D \rightarrow 1$$



Current frontier of
derivation tree:
sentential form

matched: (1 + 1 0) * 1

Derivation Tree

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

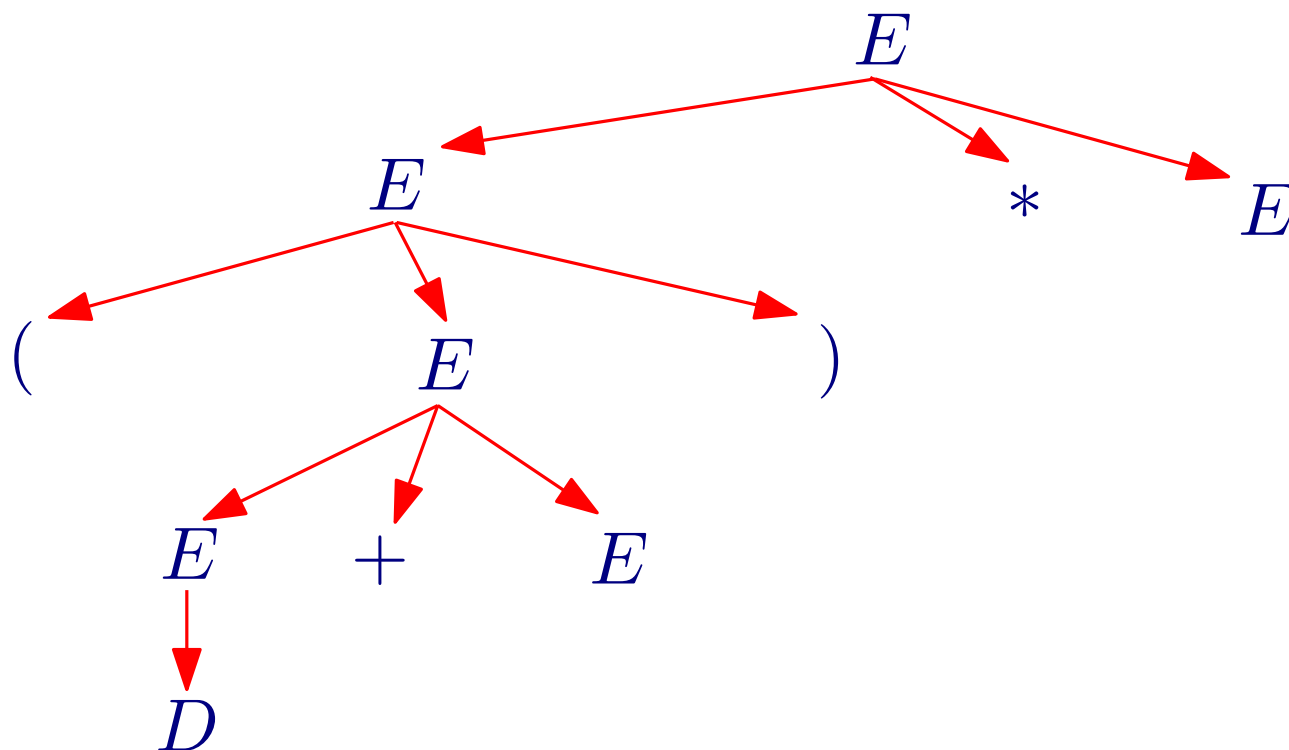
$$E \rightarrow D$$

$$D \rightarrow 0D$$

$$D \rightarrow 1D$$

$$D \rightarrow 0$$

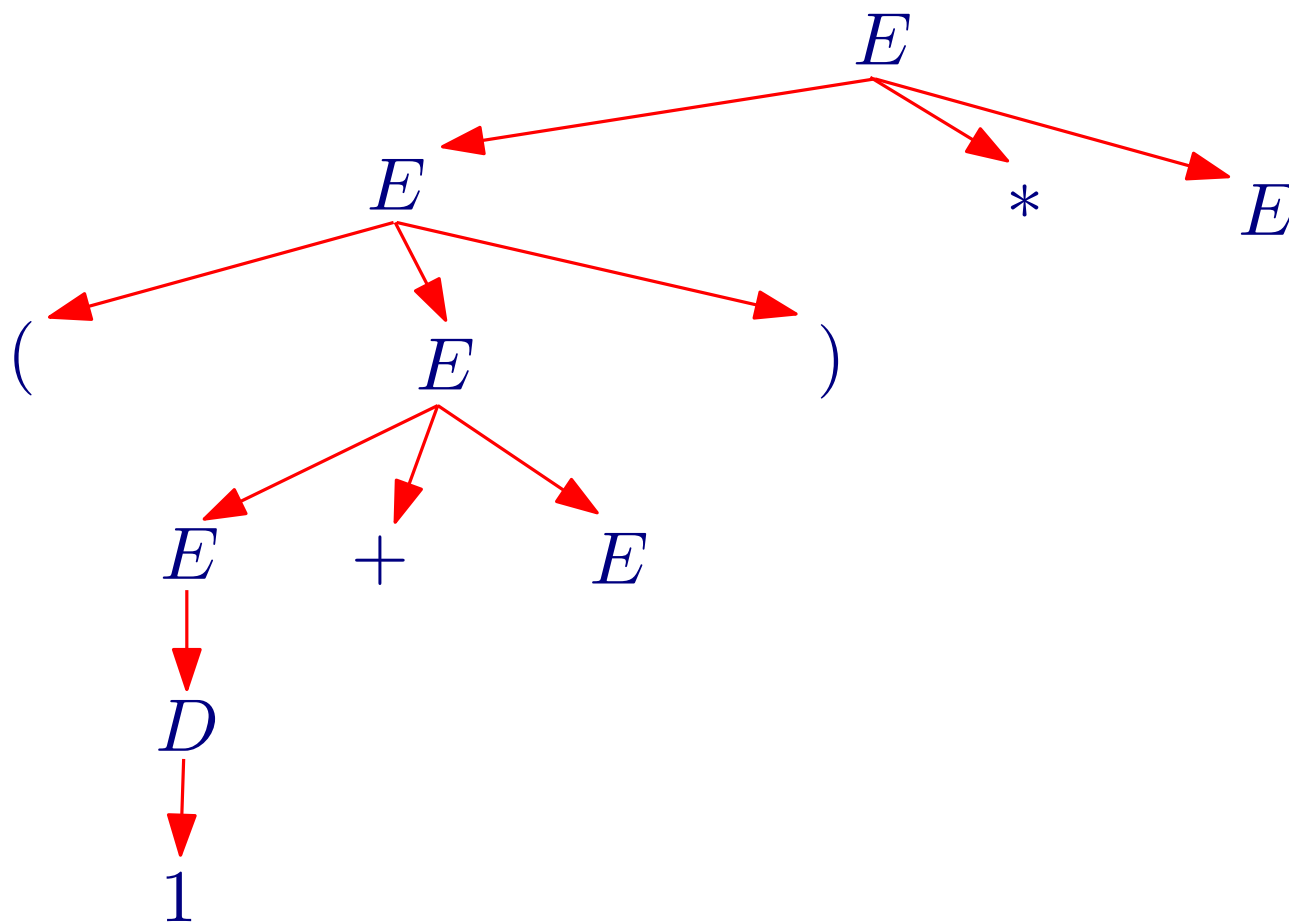
$$D \rightarrow 1$$



matched: (1 + 1 0) * 1

Derivation Tree

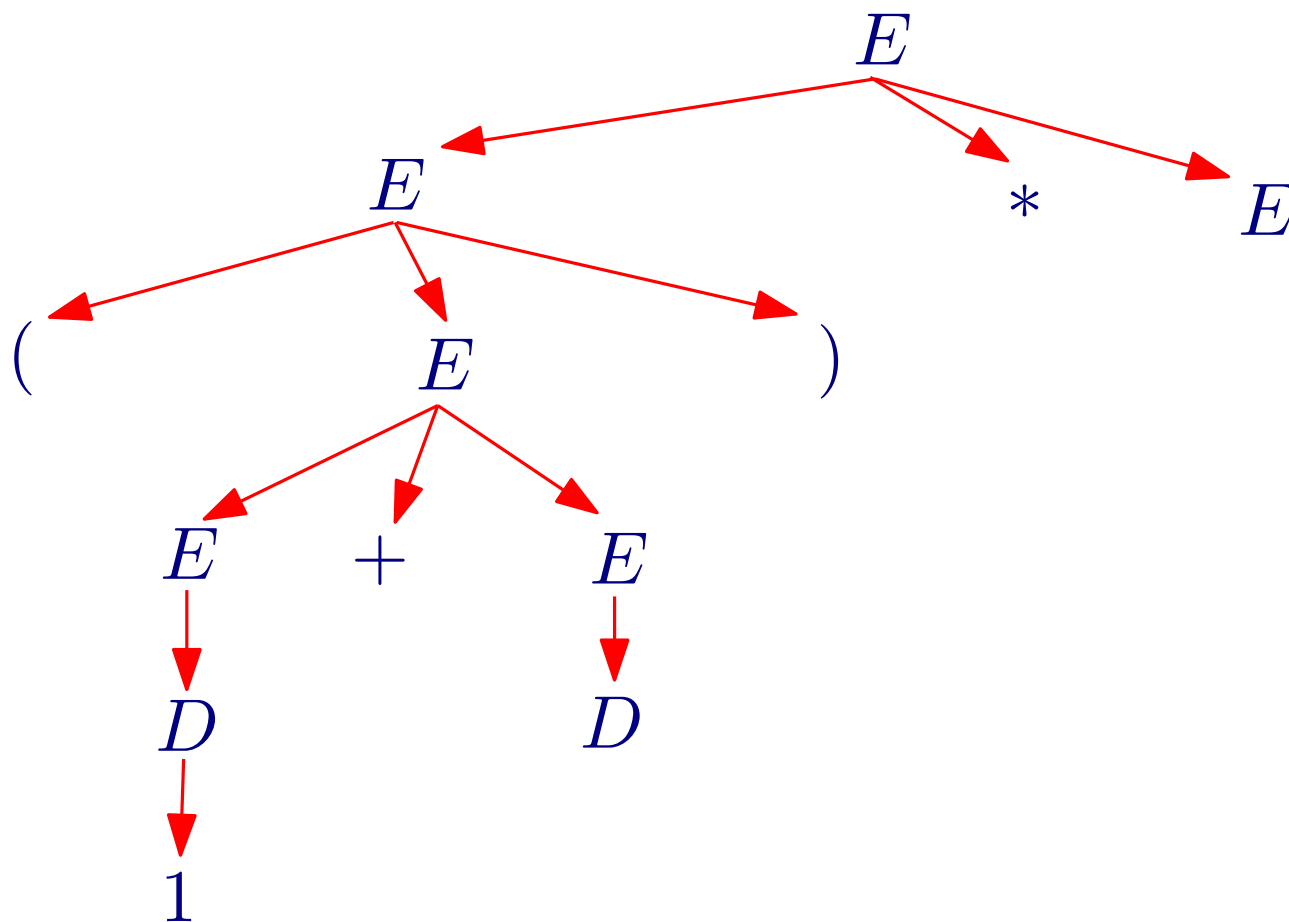
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

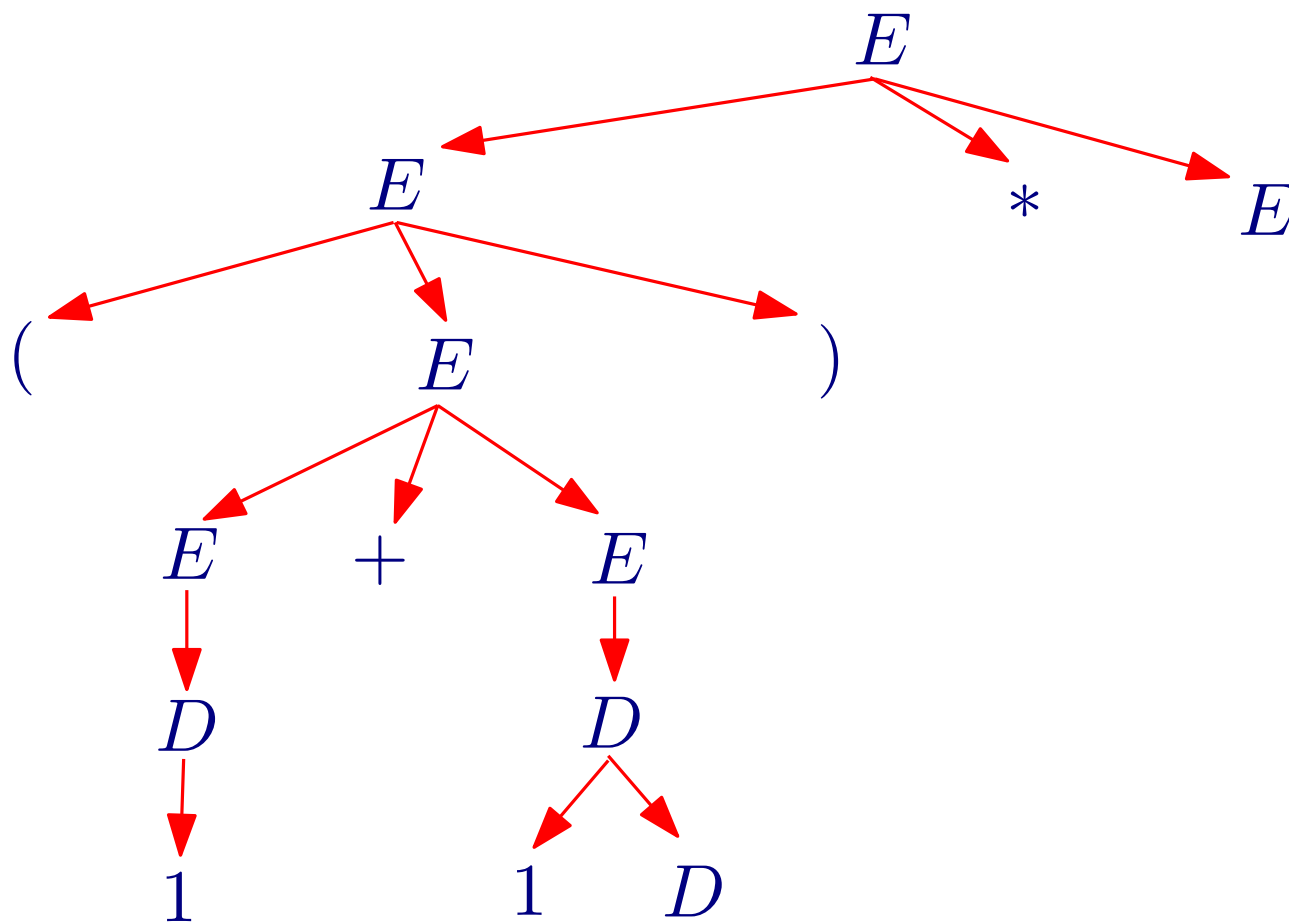
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

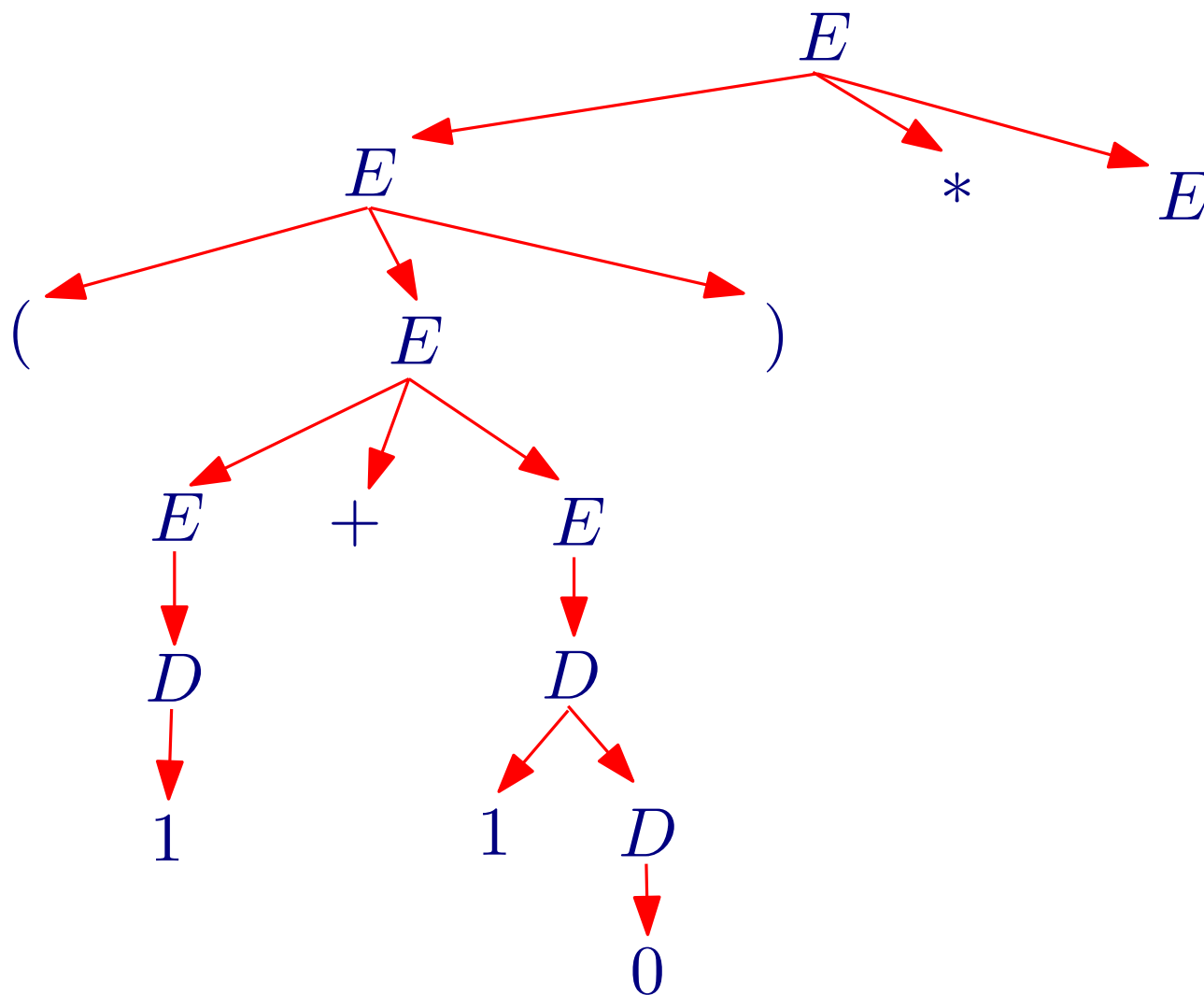
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

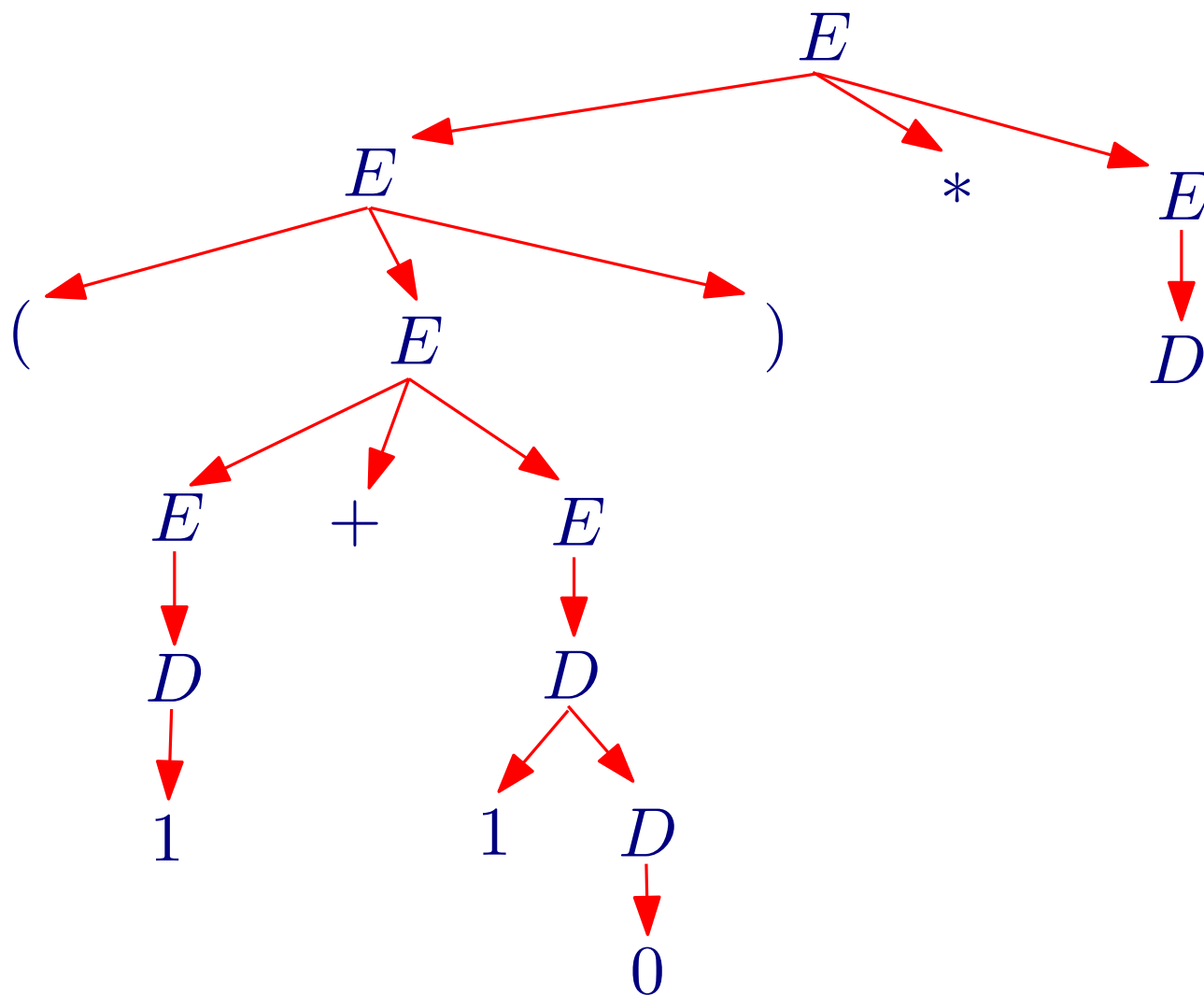
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

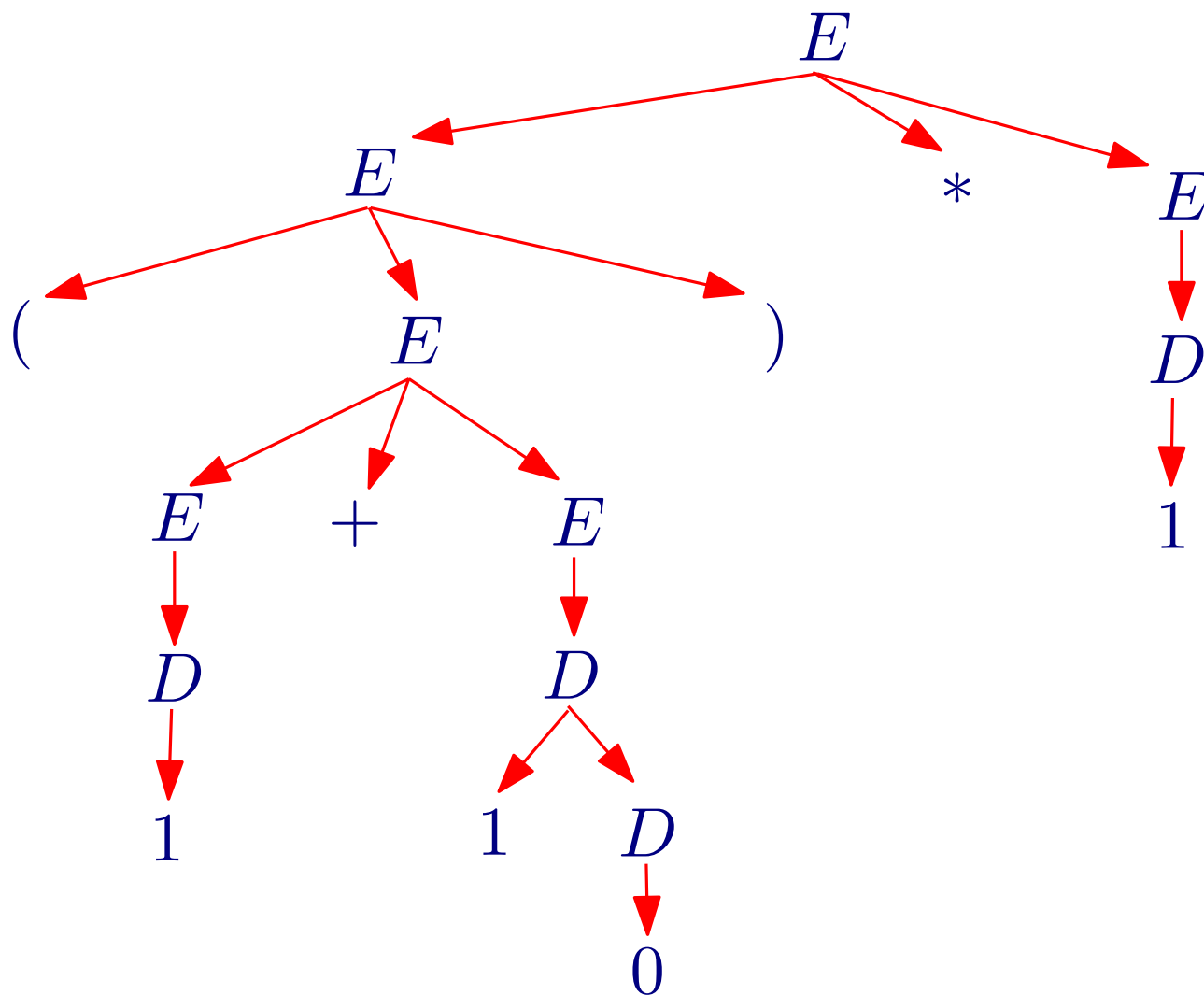
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

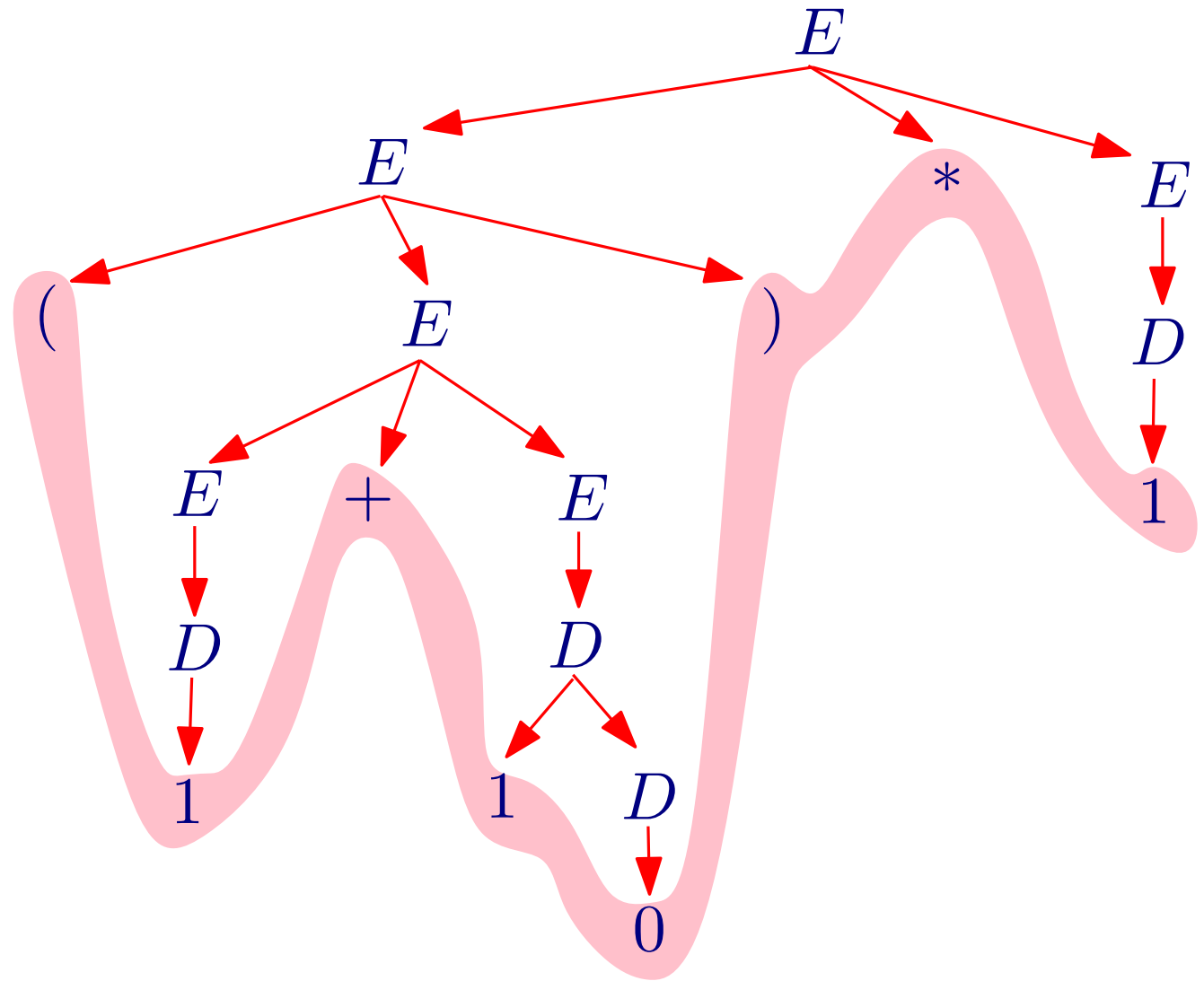
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



matched: (1 + 1 0) * 1

Derivation Tree

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$

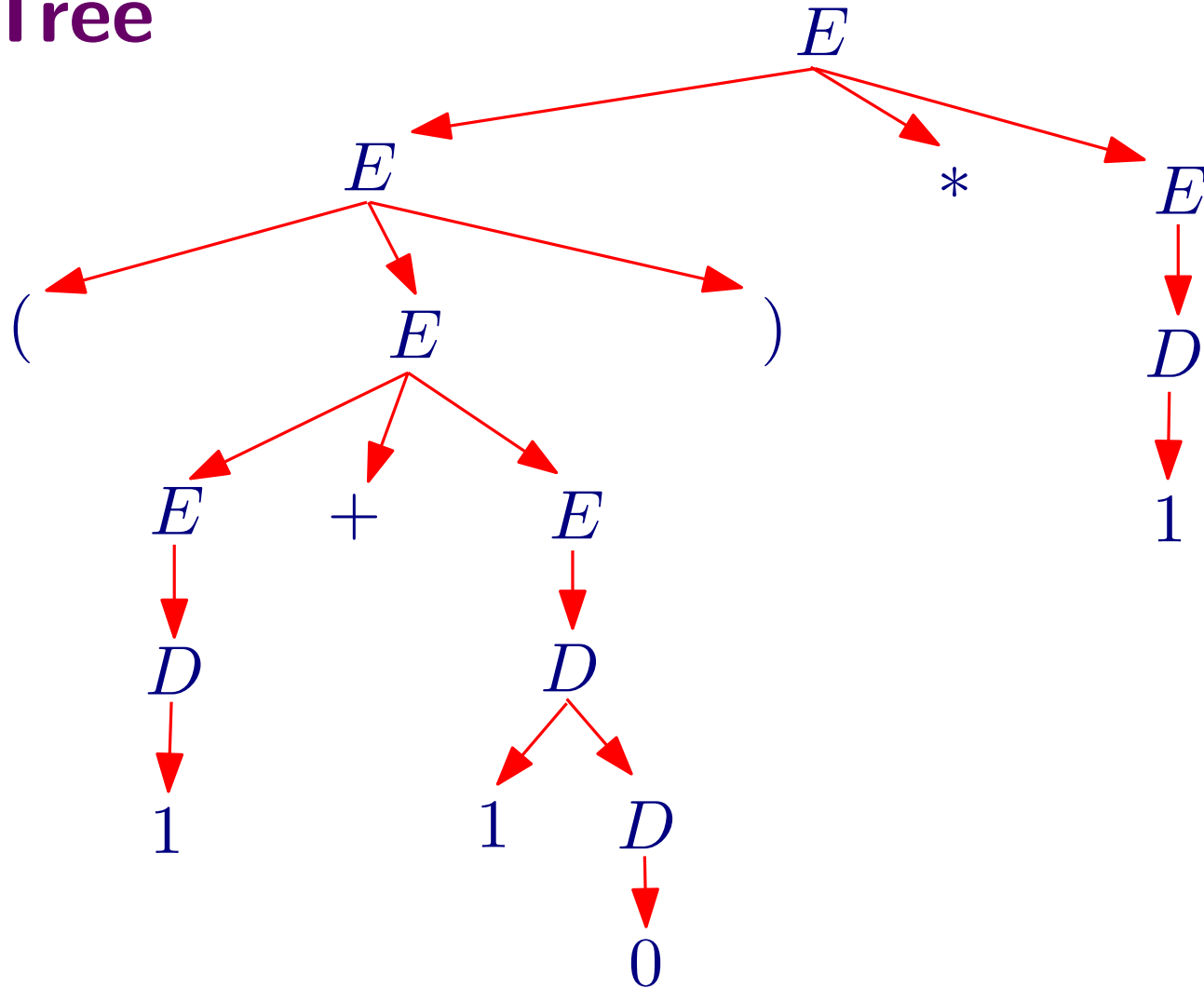


matched: (1 + 1 0) * 1

Abstract Syntax Tree

Simplified version of the parse tree:

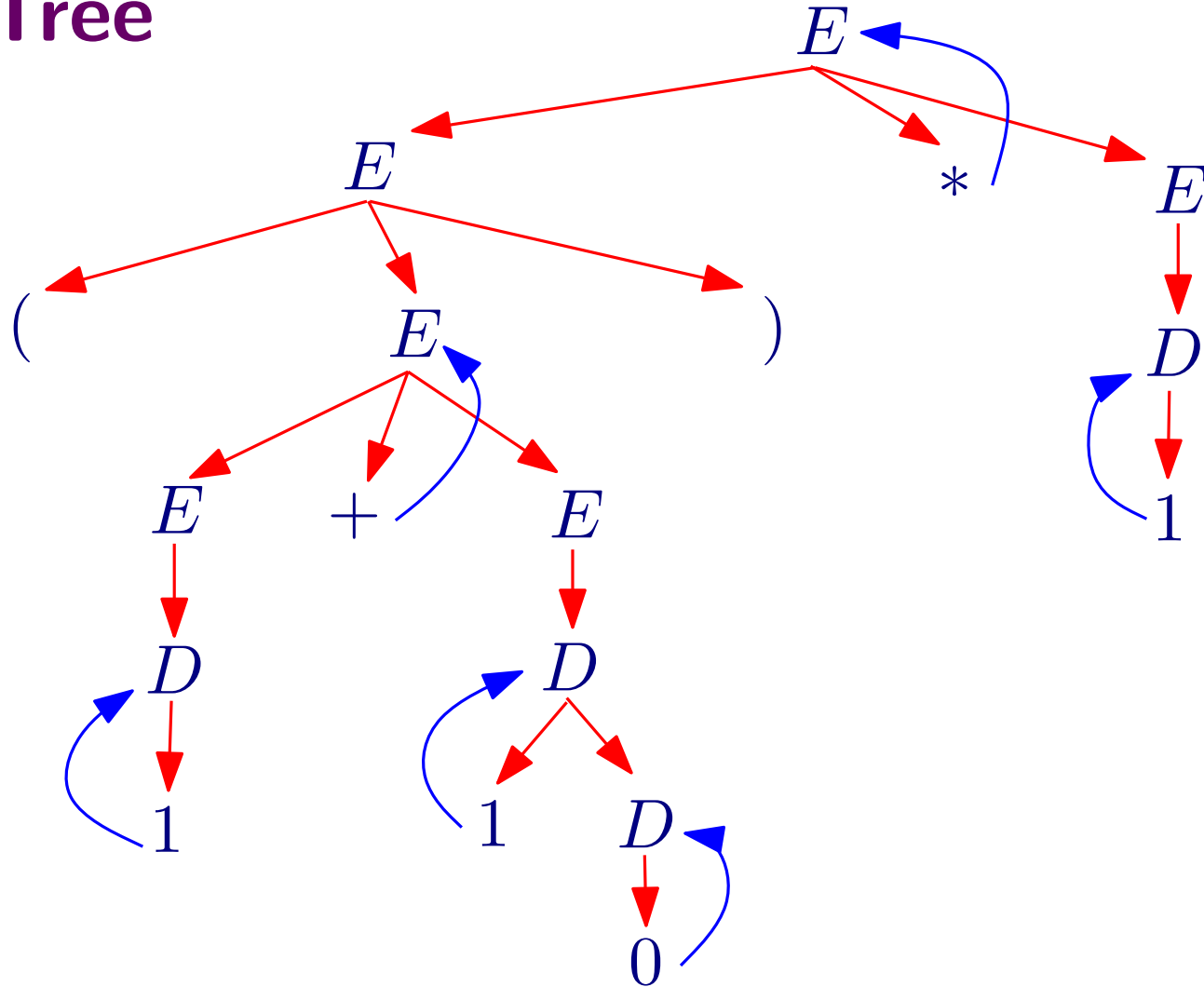
- ◇ if node has only one terminal child, make the terminal the label of current node and remove the child
- ◇ if node has only one child, shortcut the node
- ◇ many other customized rules, discussed as they are needed



Abstract Syntax Tree

Simplified version of the parse tree:

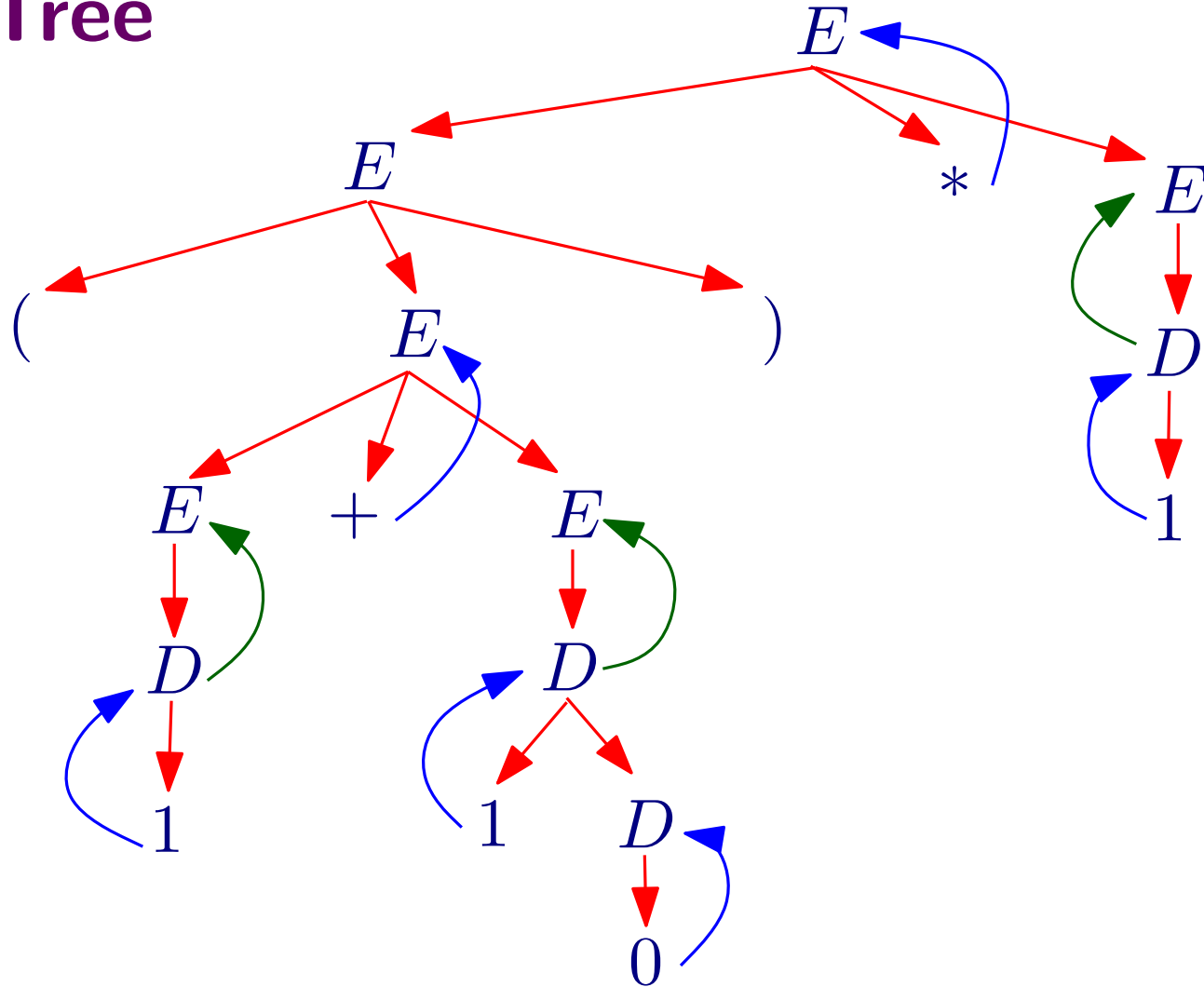
- ◇ if node has only one terminal child, make the terminal the label of current node and remove the child
- ◇ if node has only one child, shortcut the node
- ◇ many other customized rules, discussed as they are needed



Abstract Syntax Tree

Simplified version of the parse tree:

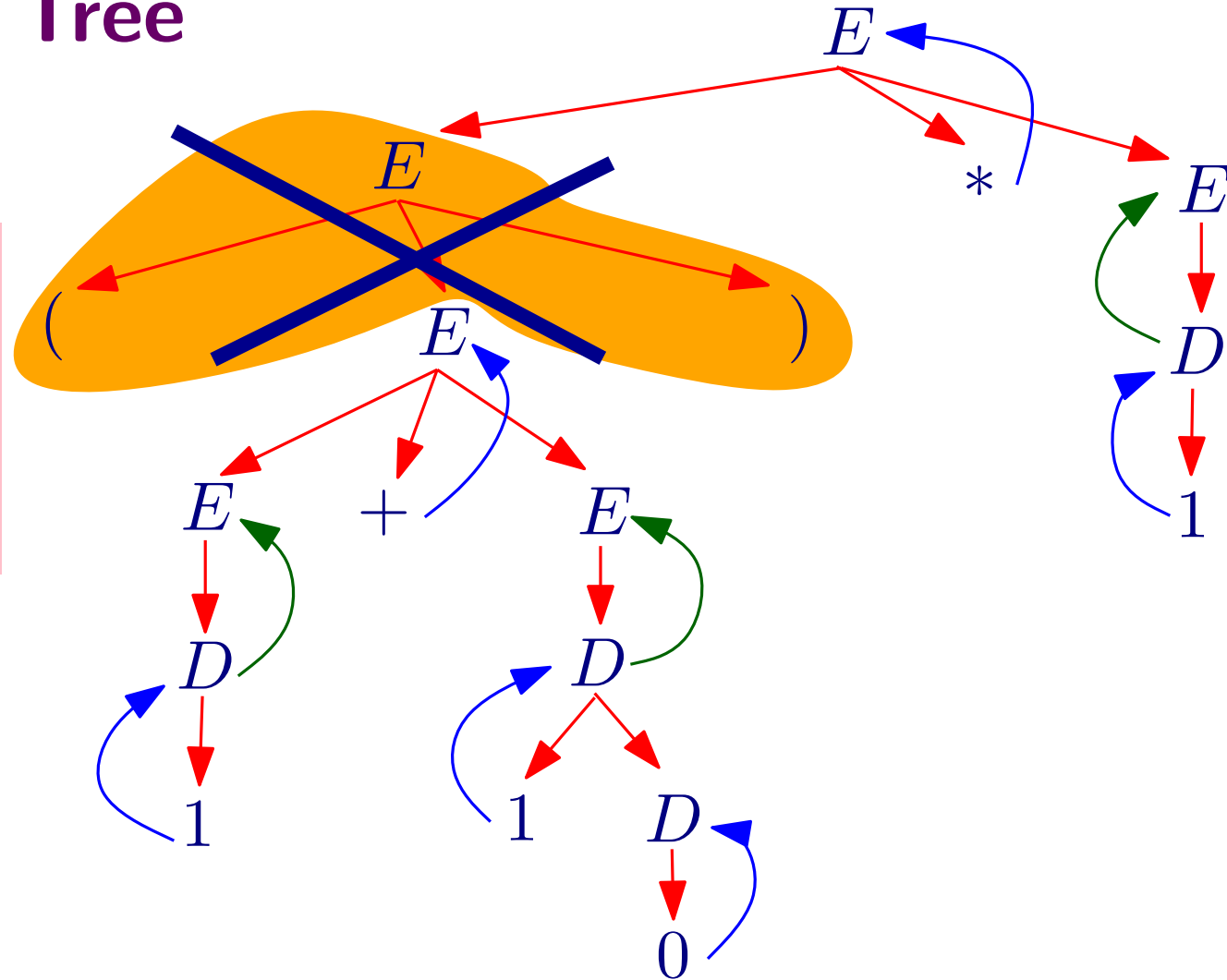
- ◇ if node has only one terminal child, make the terminal the label of current node and remove the child
- ◇ if node has only one child, shortcut the node
- ◇ many other customized rules, discussed as they are needed



Abstract Syntax Tree

Simplified version of the parse tree:

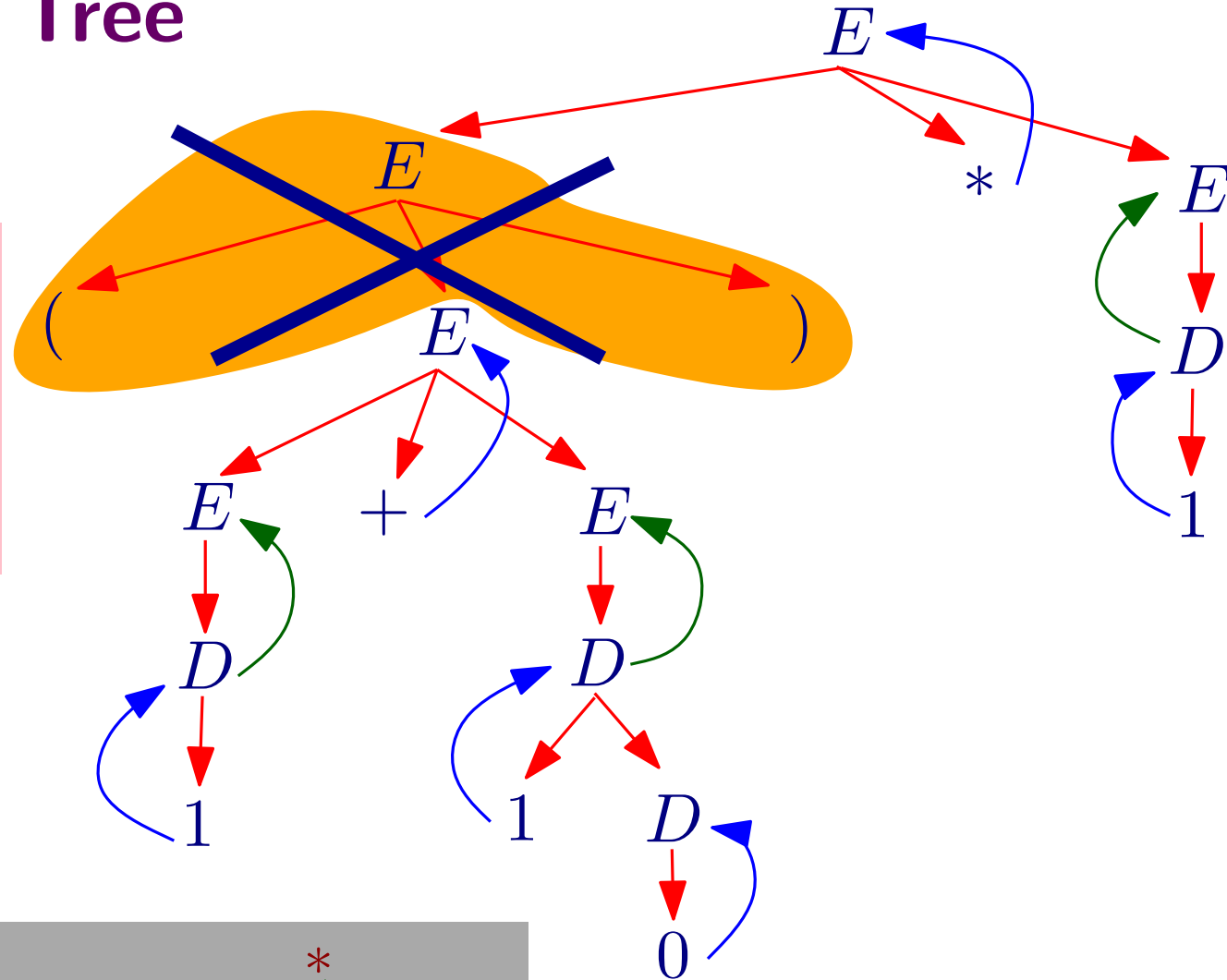
- ◇ if node has only one terminal child, make the terminal the label of current node and remove the child
- ◇ if node has only one child, shortcut the node
- ◇ many other customized rules, discussed as they are needed



Abstract Syntax Tree

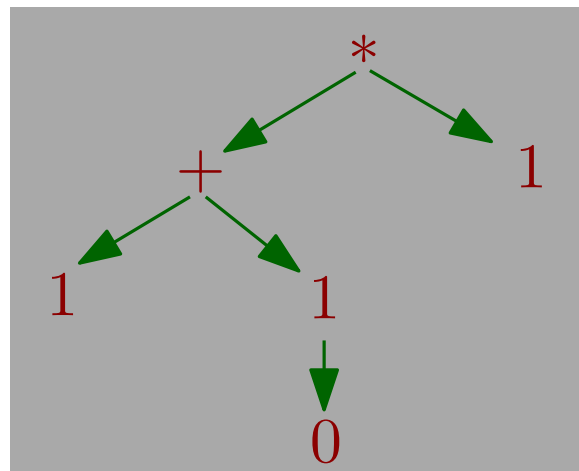
Simplified version of the parse tree:

- ◇ if node has only one terminal child, make the terminal the label of current node and remove the child
- ◇ if node has only one child, shortcut the node
- ◇ many other customized rules, discussed as they are needed



Resulting tree:

Simple, but still captures the structure of the expression.



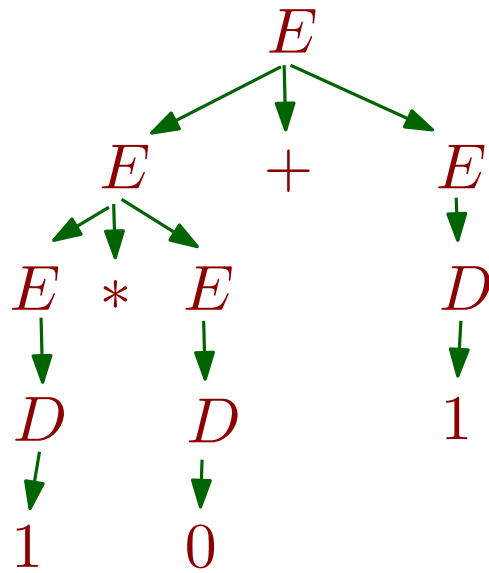
What Are Grammars Good For?

- ◇ Allow specification of (programming) languages – *generation*.
- ◇ Allow deciding whether a string belongs to the language or not – *acceptance*.
- ◇ Allow *syntactic analysis*.
- ◇ **Syntactic analysis:** building the Abstract Syntax Tree from a string or program.
- ◇ The AST can be used by compilers, program analyzers, and other code manipulators.

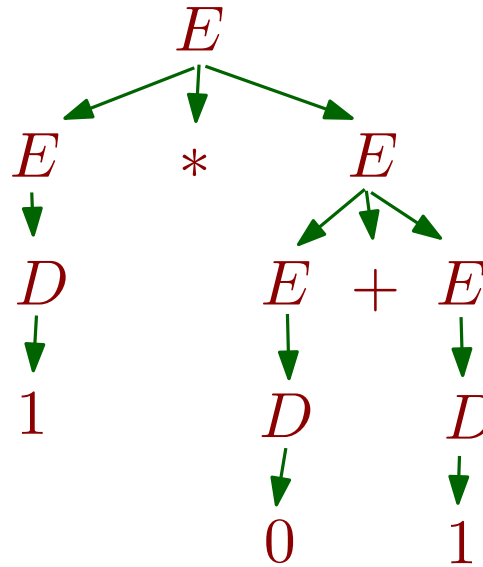
Ambiguity

Non-unique parse trees: *ambiguous grammar*

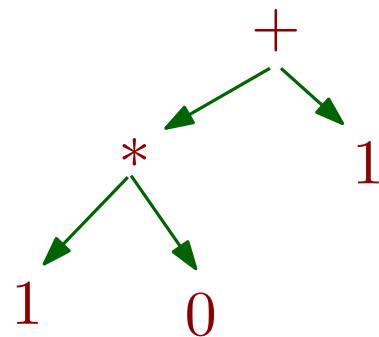
$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 0$
 $D \rightarrow 1$



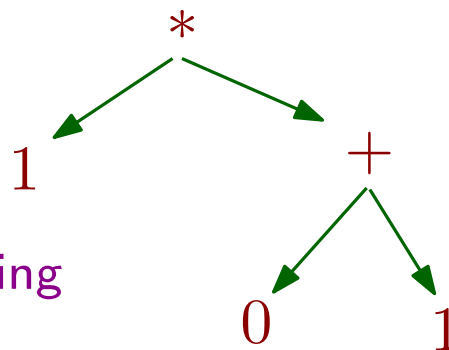
1 * 0 + 1



1 * 0 + 1

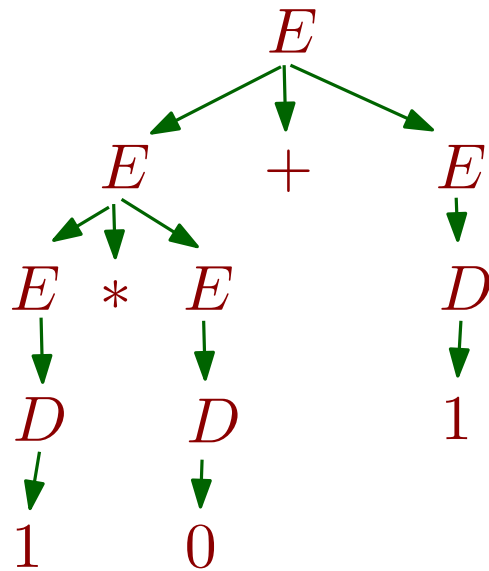


Corresponding
ASTs

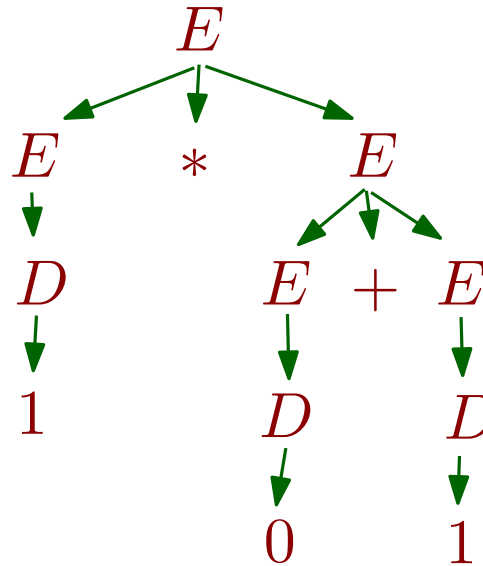


Ambiguity

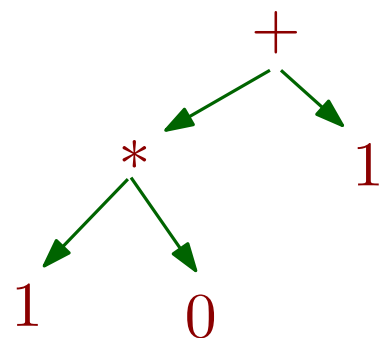
Non-unique parse trees: *ambiguous grammar*

$$\begin{aligned} E &\rightarrow (E) \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow D \\ D &\rightarrow 0D \\ D &\rightarrow 1D \\ D &\rightarrow 0 \\ D &\rightarrow 1 \end{aligned}$$


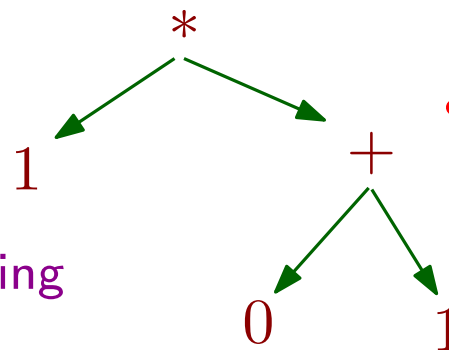
1 * 0 + 1



1 * 0 + 1



Corresponding
ASTs



Ambiguous
precedence of
operators!

Ambiguity Should Be Avoided!

- ◇ An ambiguous grammar can always be replaced by a non-ambiguous one.
- ◇ Ambiguous grammars have fewer rules, but tend to capture less of the language's structure.
- ◇ Precedence and associativity of operators is crucial to structure of languages, and should be captured in the grammar.
- ◇ Languages with non-ambiguous grammars can be parsed more efficiently.

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

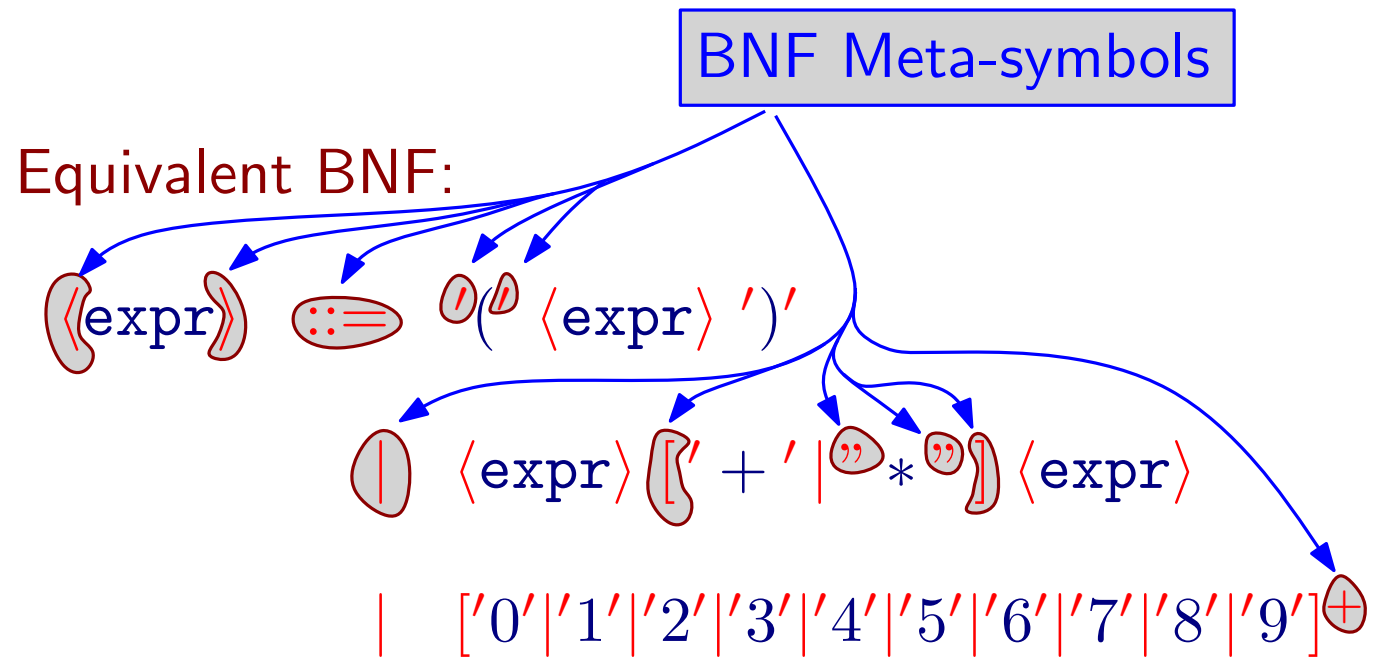
Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$
 $\quad | \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$
 $\quad | [' 0' | ' 1' | ' 2' | ' 3' | ' 4' | ' 5' | ' 6' | ' 7' | ' 8' | ' 9']^+$

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$



Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Non-terminals
are enclosed in
angle brackets

Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$
 $\quad | \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$
 $\quad | [' 0 ' | ' 1 ' | ' 2 ' | ' 3 ' | ' 4 ' | ' 5 ' | ' 6 ' | ' 7 ' | ' 8 ' | ' 9 ']^+$

Backus-Naur Form

Terminals are enclosed in simple or double quotes.

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= (' \langle \text{expr} \rangle ')$
 $| \langle \text{expr} \rangle ['+' | '*'] \langle \text{expr} \rangle$
 $| ['0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9']^+$

Single quote terminal: `' '`
 Double quote terminal: `' ' ' '`

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= (' \langle \text{expr} \rangle ')$
 $| \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$
 $| ['0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9']^+$

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Grouping

Equivalent BNF:

$\langle \text{expr} \rangle ::= (' \langle \text{expr} \rangle ')$
 $| \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$
 $| ['0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9']^+$

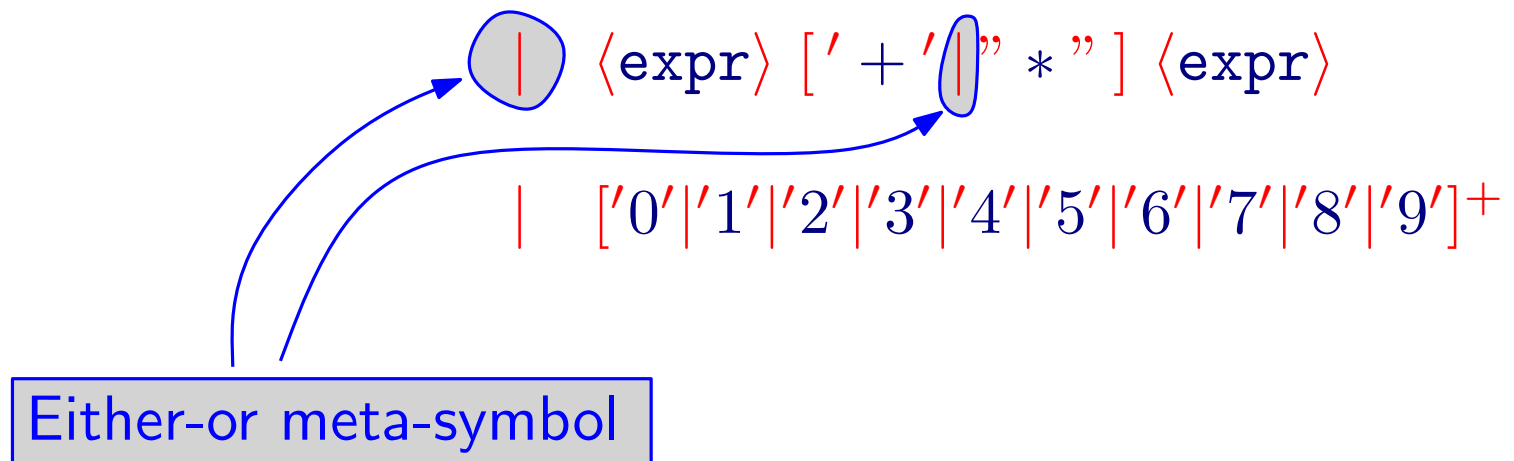
Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$



Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$

$| \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$

$| [' 0 ' | ' 1 ' | ' 2 ' | ' 3 ' | ' 4 ' | ' 5 ' | ' 6 ' | ' 7 ' | ' 8 ' | ' 9 ']^+$

Factorization

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$

$| \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$

$| [' 0 ' | ' 1 ' | ' 2 ' | ' 3 ' | ' 4 ' | ' 5 ' | ' 6 ' | ' 7 ' | ' 8 ' | ' 9 ']^+$

Iteration

One or more repetitions

Backus-Naur Form

Original grammar:

$E \rightarrow (E)$
 $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow D$
 $D \rightarrow 0D$
 $D \rightarrow 1D$
 $D \rightarrow 2D$
 $D \rightarrow 3D$
 $D \rightarrow 4D$
 $D \rightarrow 5D$
 $D \rightarrow 6D$
 $D \rightarrow 7D$
 $D \rightarrow 8D$
 $D \rightarrow 9D$
 $D \rightarrow 0$
 $D \rightarrow 1$
 $D \rightarrow 2$
 $D \rightarrow 3$
 $D \rightarrow 4$
 $D \rightarrow 5$
 $D \rightarrow 6$
 $D \rightarrow 7$
 $D \rightarrow 8$
 $D \rightarrow 9$

Equivalent BNF:

$\langle \text{expr} \rangle ::= '(' \langle \text{expr} \rangle ')'$

$| \langle \text{expr} \rangle [' + ' | ' * '] \langle \text{expr} \rangle$

$| [' 0 ' | ' 1 ' | ' 2 ' | ' 3 ' | ' 4 ' | ' 5 ' | ' 6 ' | ' 7 ' | ' 8 ' | ' 9 ']^+$

Iteration

One or more repetitions

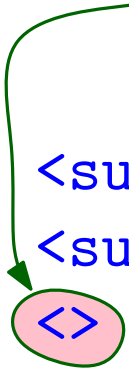
For 0 or more repetitions, use *

A Non-Ambiguous Grammar for Expressions

```
<expr>      ::= <subexpr> <term>
<subexpr>   ::= <subexpr> <term> ['+' | '-']
              | <>
<term>      ::= <subterm> <factor>
<subterm>   ::= <subterm> <factor> ['*' | '/']
              | <>
<factor>    ::= <base> <restexp>
<restexp>   ::= '^' <base> <restexp>
              | <>
<base>      ::= '(' <expr> ')',
              | a | b | c | d
```

A Non-Ambiguous Grammar for Expressions

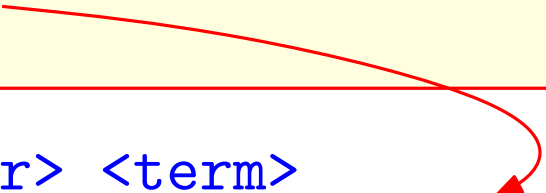
Empty string: can appear
inbetween any two terminals



```
<expr>      ::= <subexpr> <term>
<subexpr>   ::= <subexpr> <term> ['+' | '-' ]
              | <>
<term>      ::= <subterm> <factor>
<subterm>   ::= <subterm> <factor> ['*' | '/' ]
              | <>
<factor>    ::= <base> <restexp>
<restexp>   ::= '^' <base> <restexp>
              | <>
<base>      ::= '(' <expr> ') '
              | a | b | c | d
```

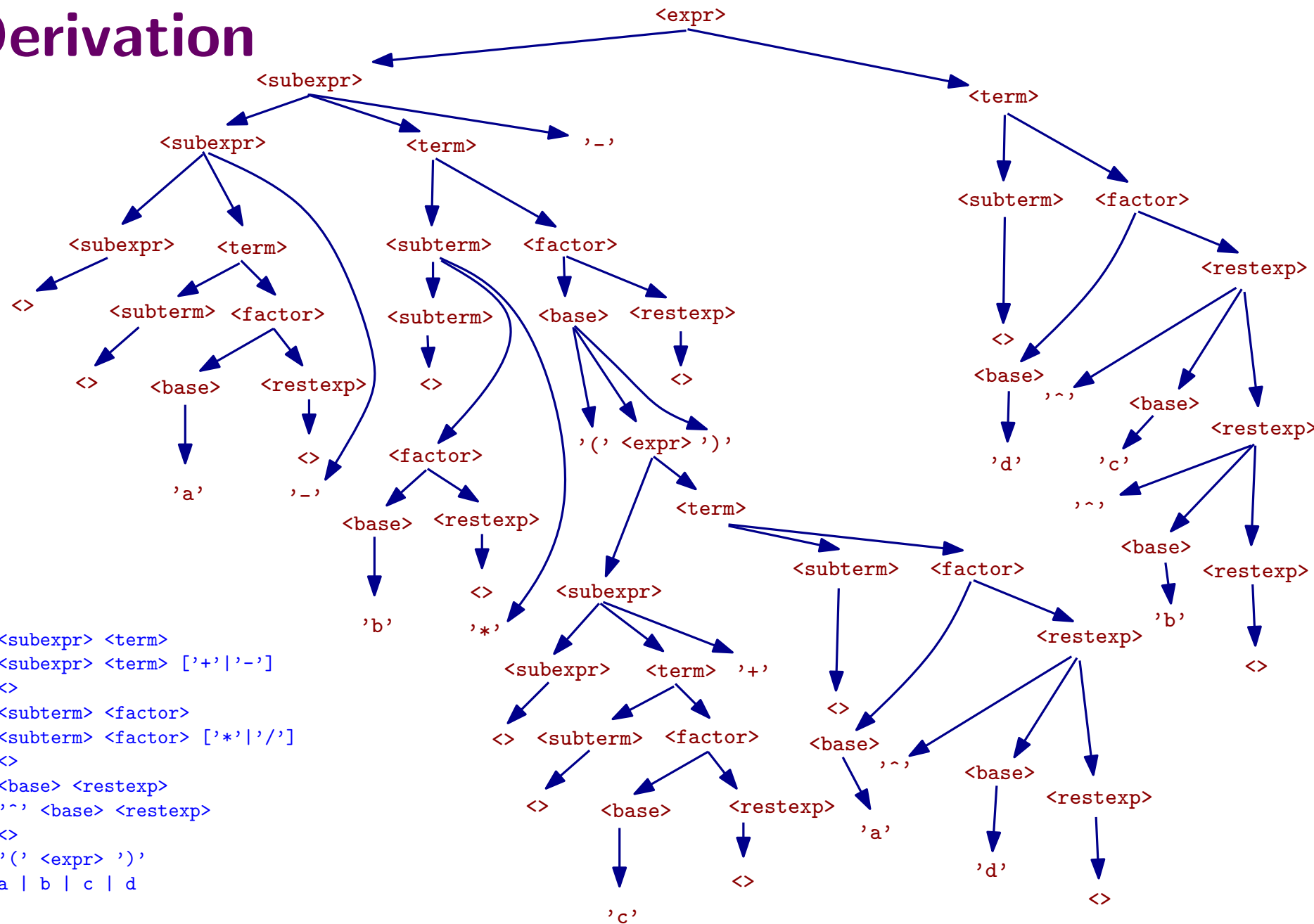

A Non-Ambiguous Grammar for Expressions

Compress the spec, but information about associativity is lost!



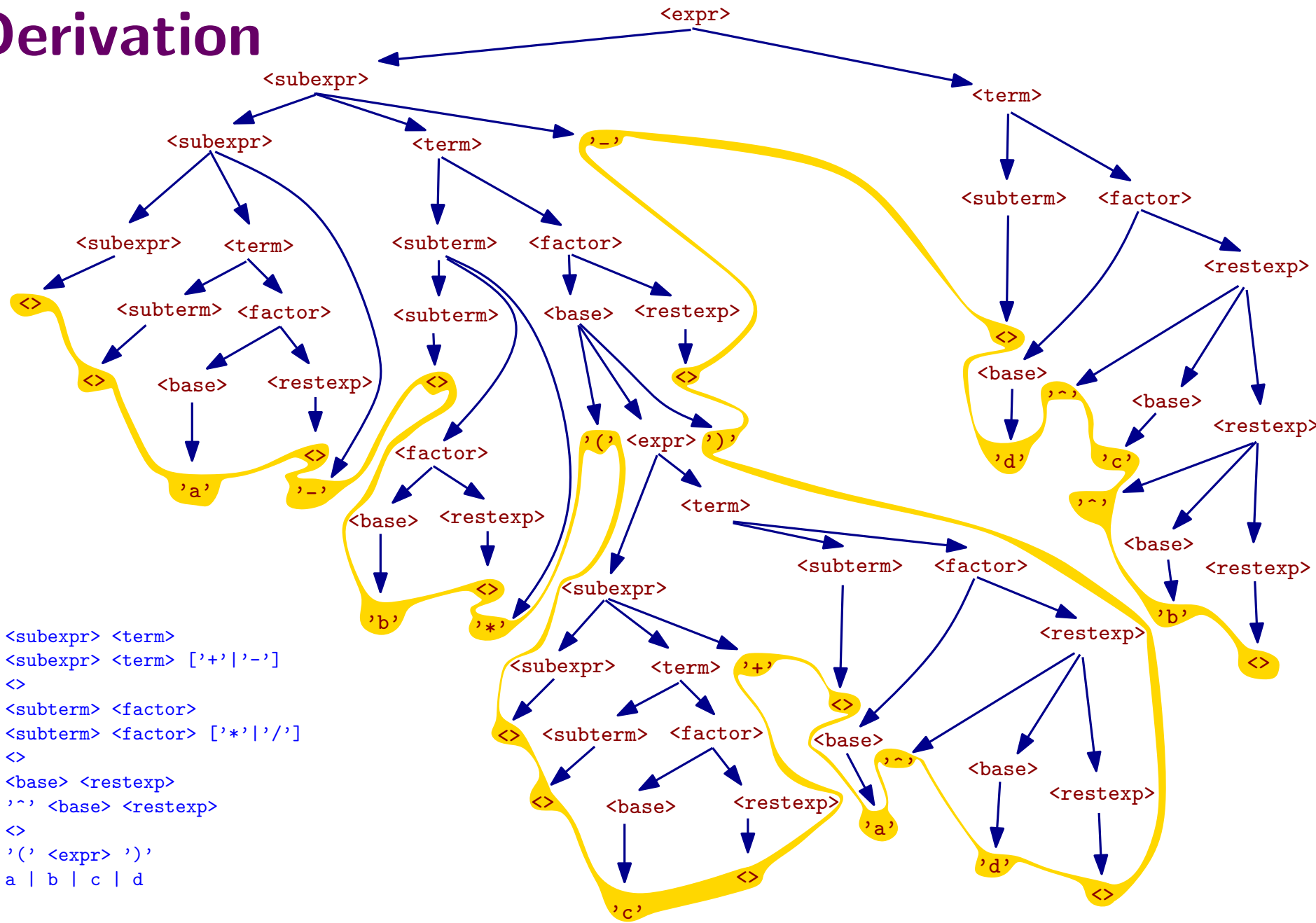
```
<expr>      ::= <subexpr> <term>
<subexpr>   ::= [ <term> ['+' | '-' ] ]*
<term>      ::= <subterm> <factor>
<subterm>   ::= <subterm> <factor> ['*' | '/']
              | <>
<factor>    ::= <base> <restexp>
<restexp>   ::= '^' <base> <restexp>
              | <>
<base>      ::= '(' <expr> ') '
              | a | b | c | d
```

A Derivation



a - b * (c + a ^ d) - d ^ c ^ b

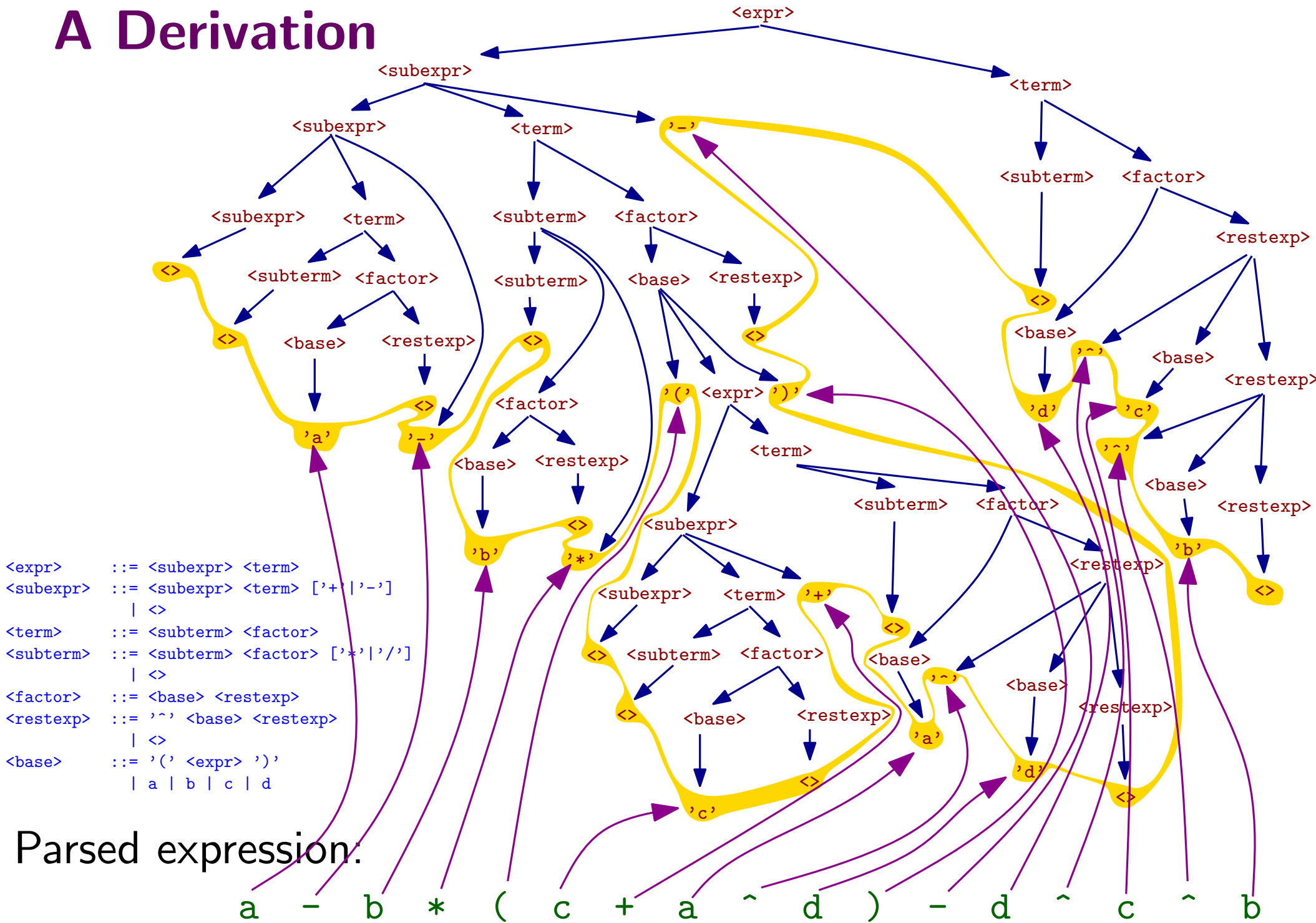
A Derivation



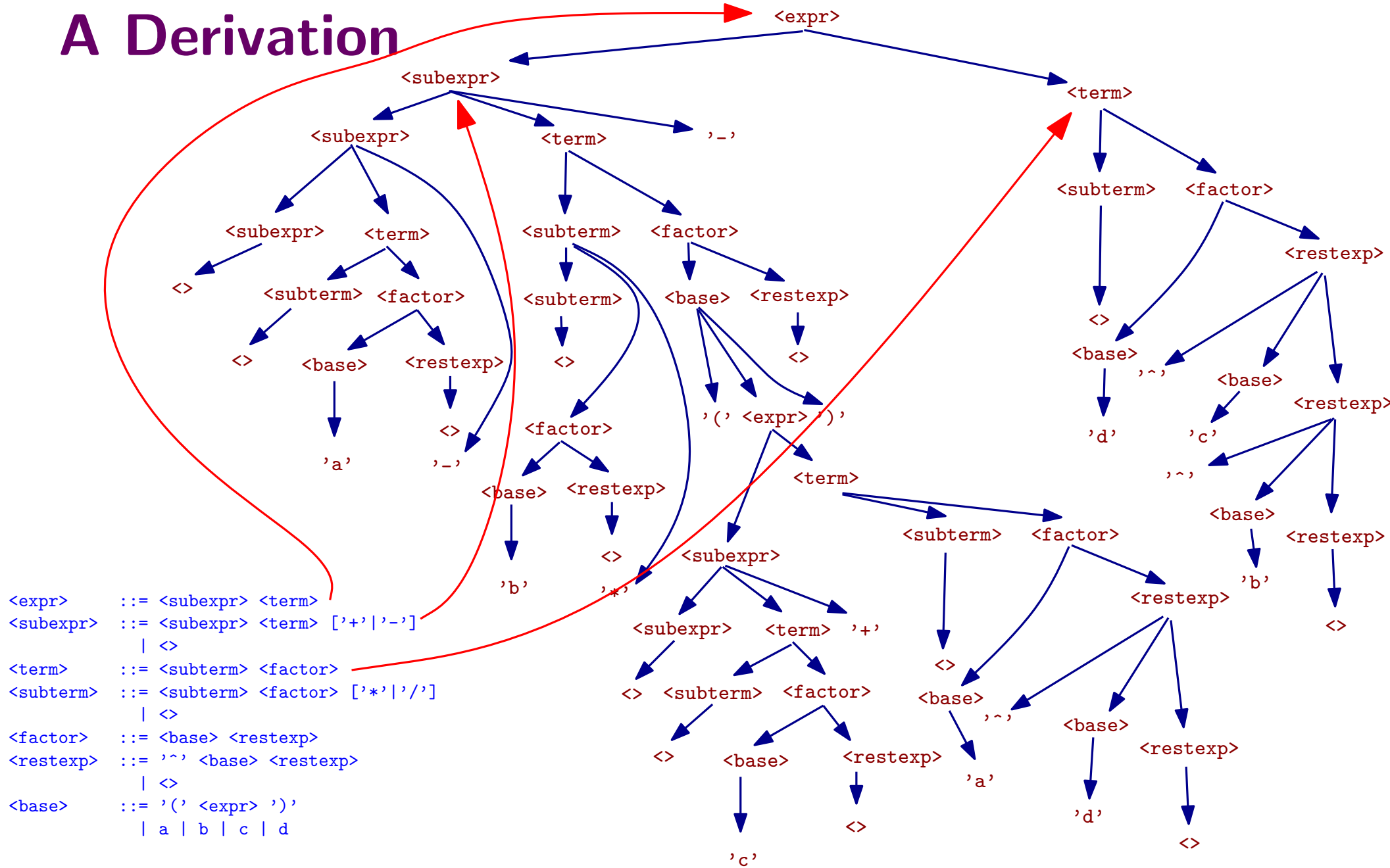
Parsed expression:

$$a - b * (c + a ^ d) - d ^ c ^ b$$

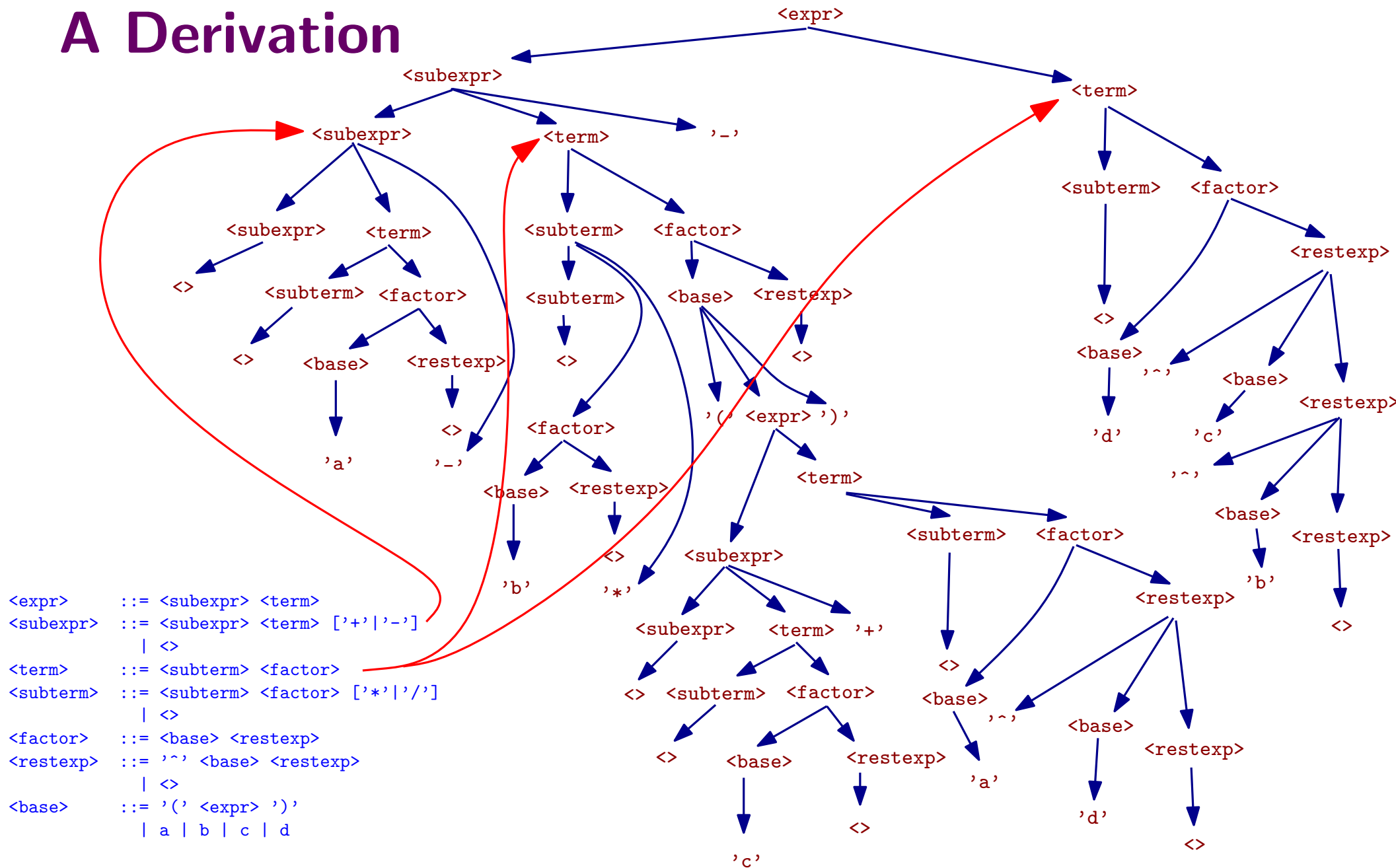
A Derivation



A Derivation

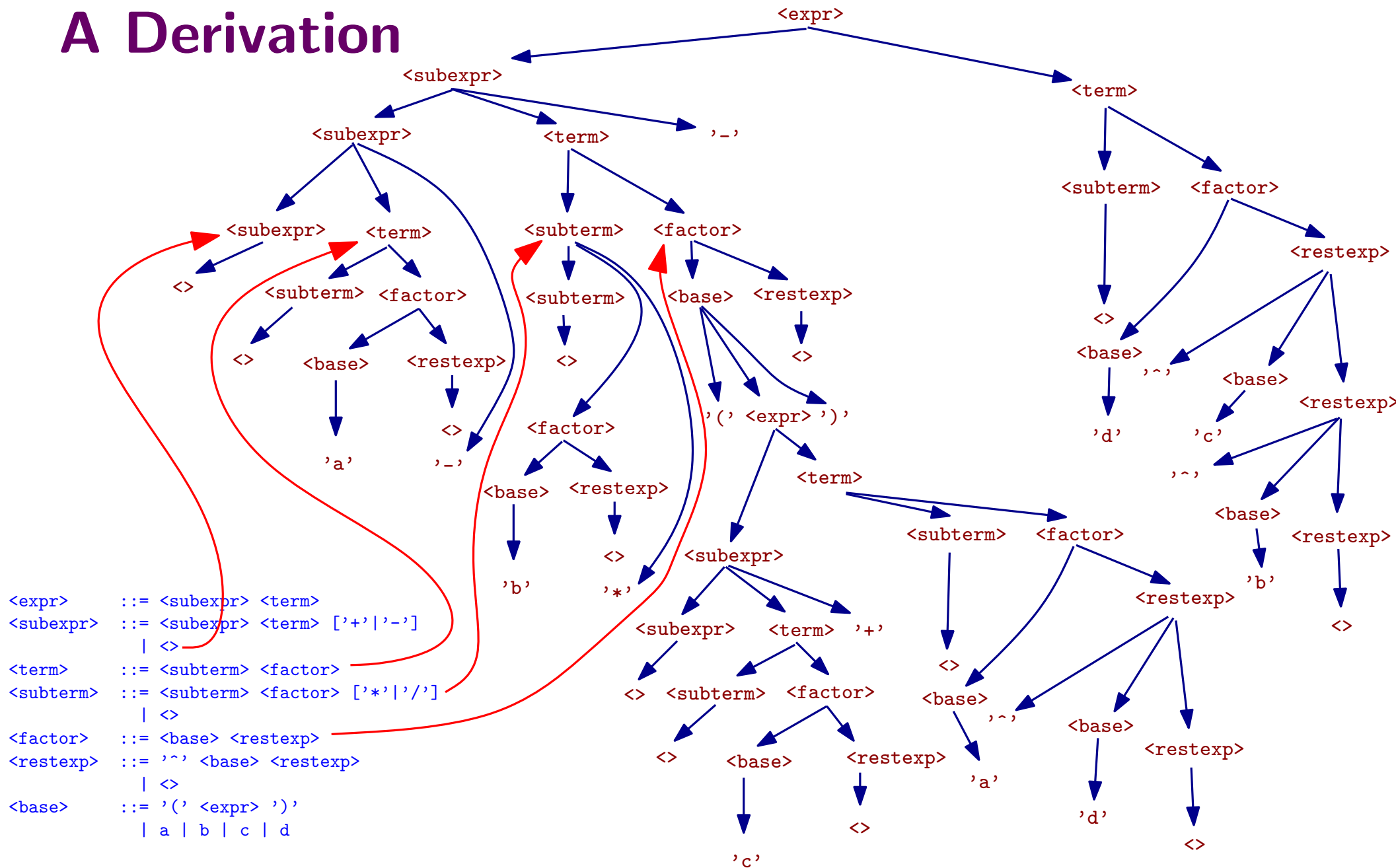

$$a - b * (c + a ^ d) - d ^ c ^ b$$

A Derivation

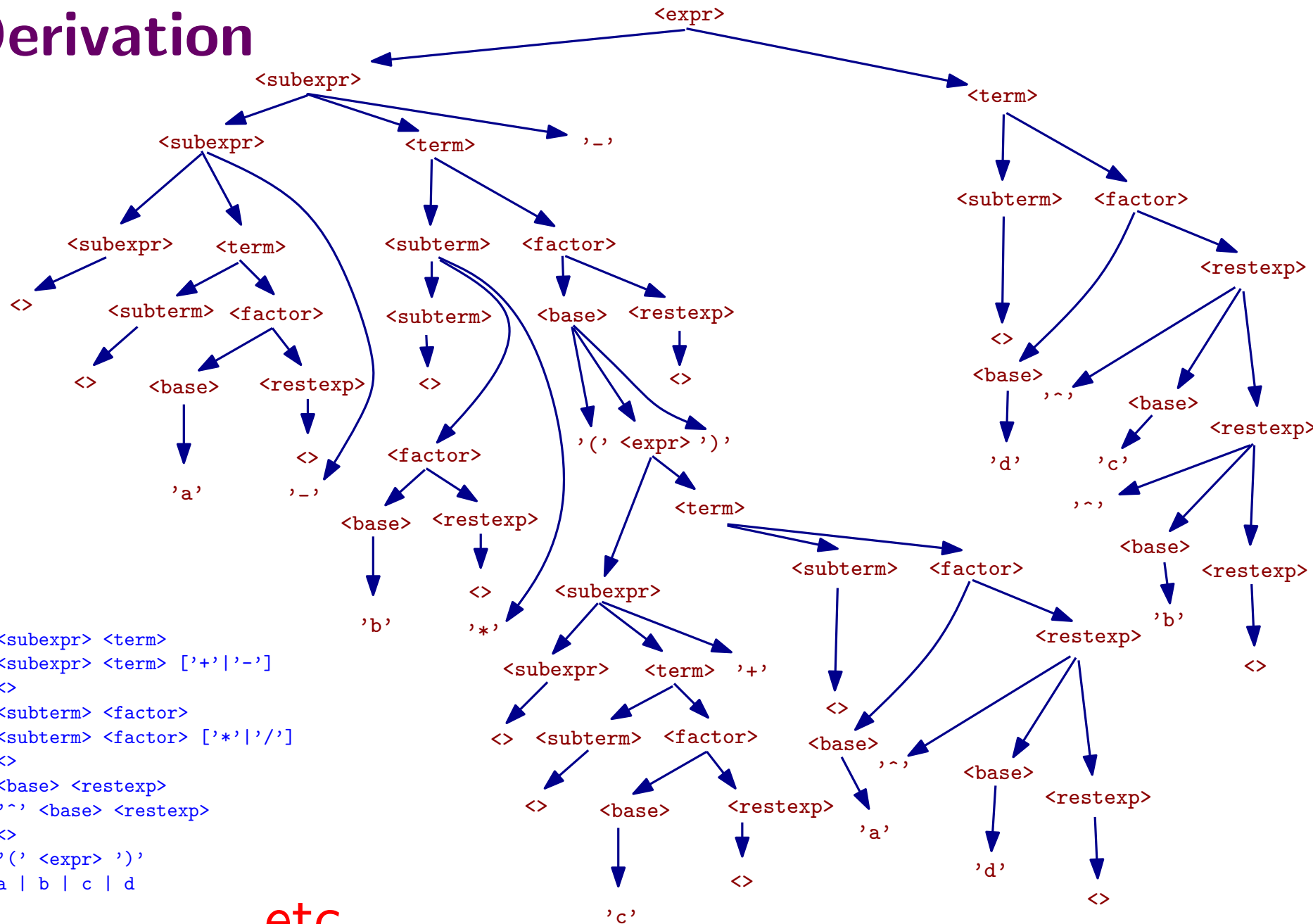


a - b * (c + a ^ d) - d ^ c ^ b

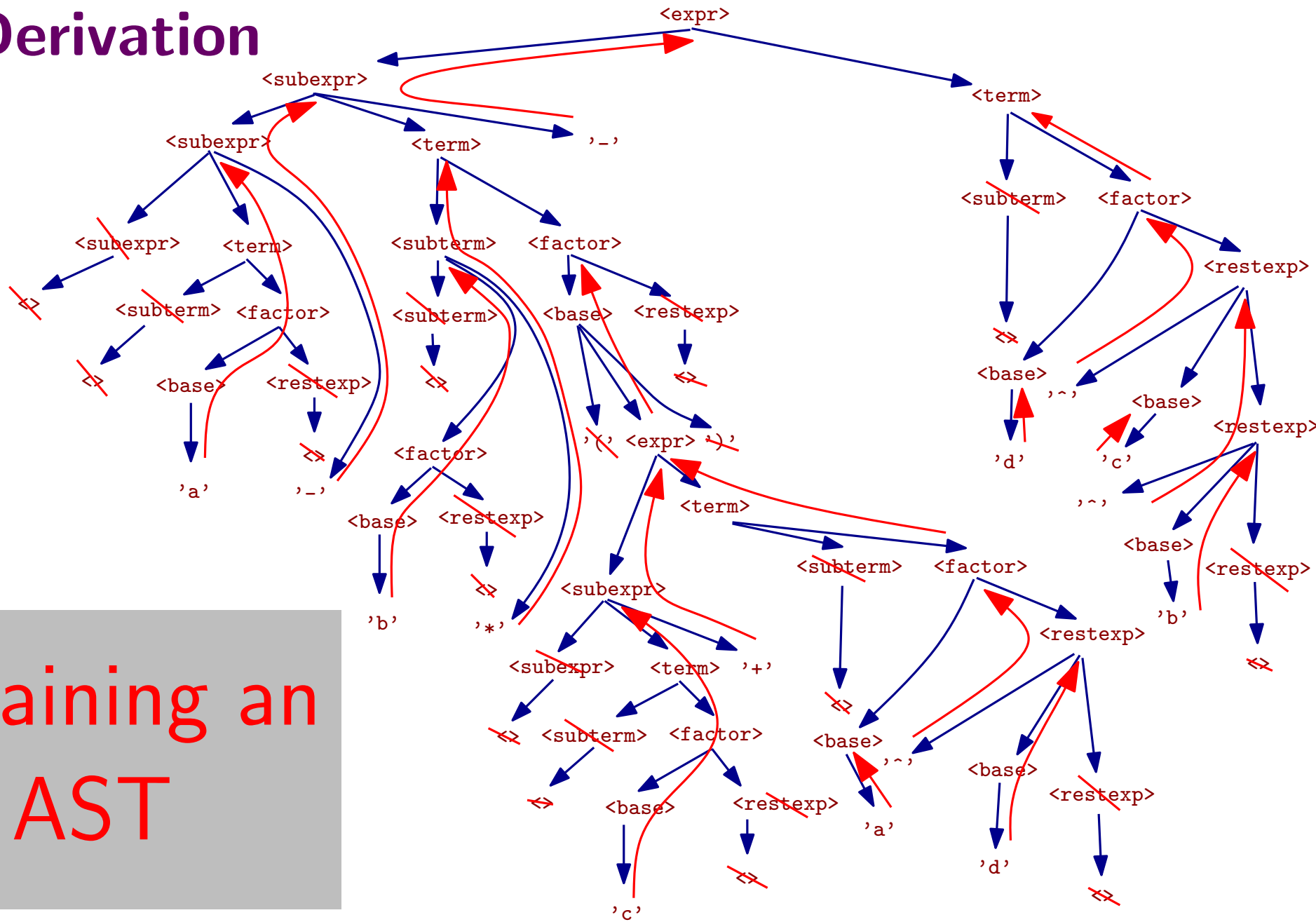
A Derivation


$$a - b * (c + a ^ d) - d ^ c ^ b$$

A Derivation


$$a - b * (c + a ^ d) - d ^ c ^ b$$

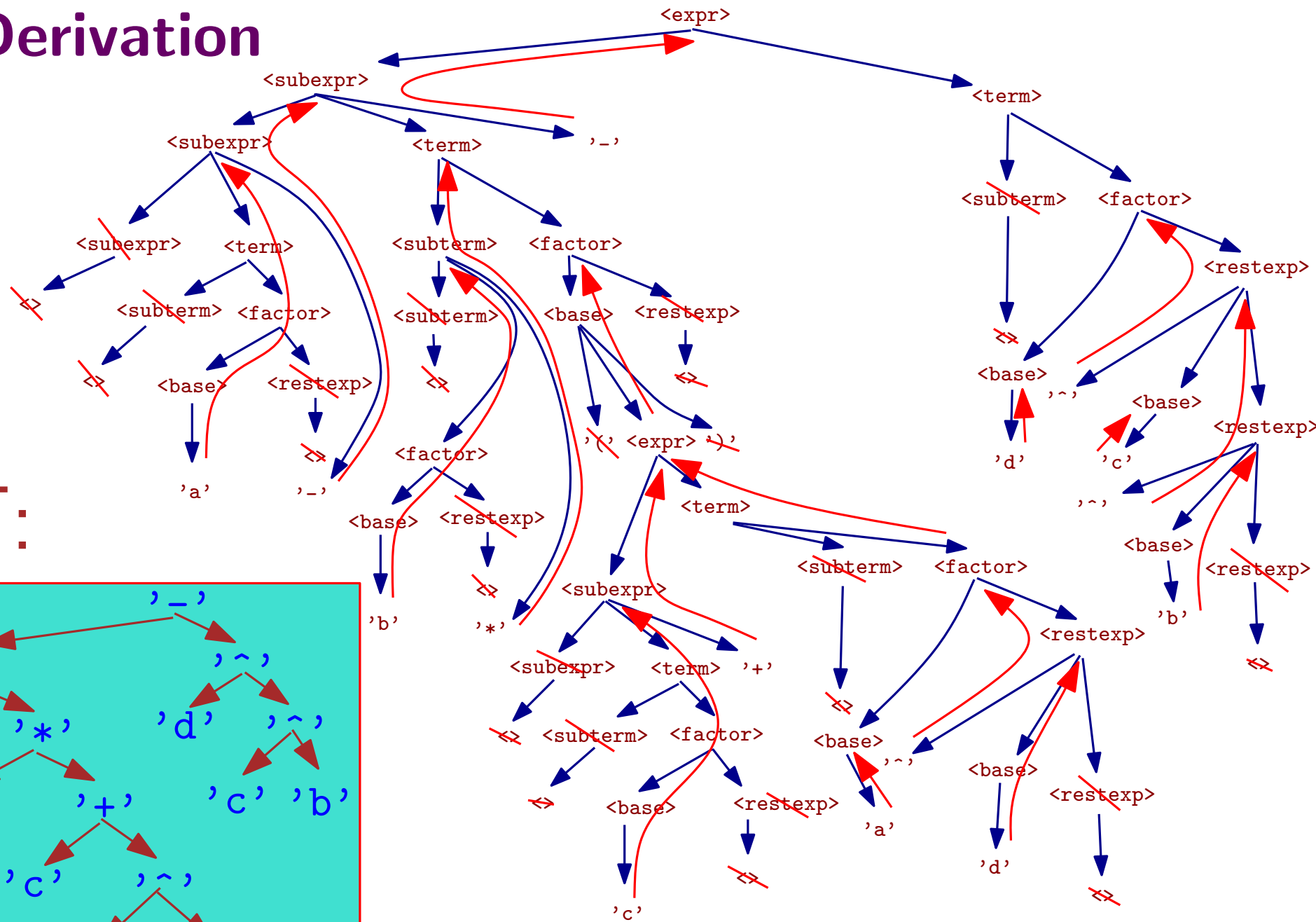
A Derivation



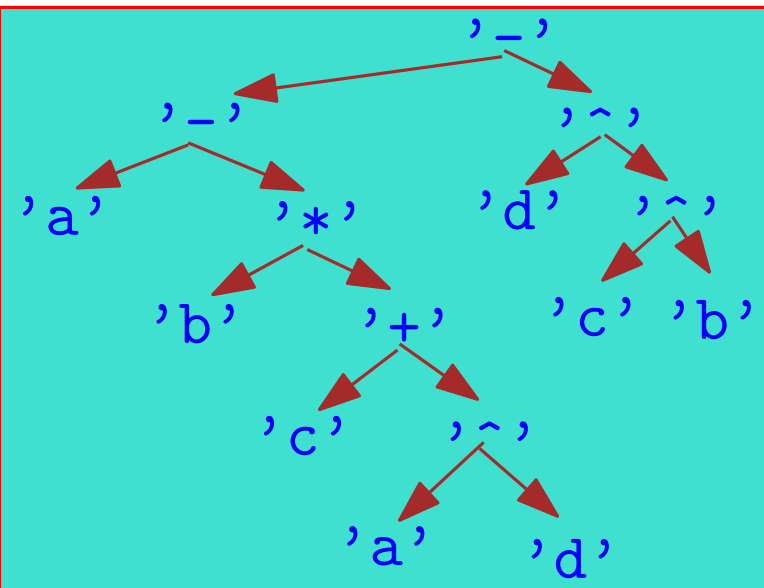
Obtaining an AST

$$a - b * (c + a ^ d) - d ^ c ^ b$$

A Derivation



AST:



a - b * (c + a ^ d) - d ^ c ^ b

Lessons Learned

- ◇ There are many grammars for the same language.
- ◇ *Parsing*: Building an AST for a given string, provided it is in the grammar's language.
- ◇ Unambiguous grammars are preferred for that purpose.
- ◇ Grammar should try to capture structural properties of the language, such as associativity of operators, and nesting of blocks.
- ◇ The AST can be built from the parse tree through an algorithmic process.
- ◇ Each tree segment is processed in a systematic way.

Stratification

- ◇ For syntactic analysis of programming languages, a 2-layered approach is used.
- ◇ The primitive elements of the language, such as *constants*, *identifiers*, *operators*, *keywords*, etc, *special purpose grammars are defined*.
- ◇ The start symbols of the special purpose grammars are then use as *terminals of a higher-level grammar*, which defines the language.
- ◇ Low-level parsing: *lexical analysis*
- ◇ High-level parsing: *syntactic analysis*
- ◇ Low-level parsing: the language has very little structure, grammars are in general very simple.
- ◇ High-level parsing: structure is important, and so is abstracting away from the lexical level; an identifier or a constant becomes a *terminal symbol of the higher-level grammar*.

Lexical Analysis Grammars

```
<PI> ::= <Digit> | <Digit> <PI>  
<Digit> ::= '0' | '1' | ... | '9'
```

```
<String> ::= '"' <Seq> '"'  
<Seq> ::= [ 'a' | 'b' | ... | 'z' |  
            '0' | ... | '9' |  
            '+' | '-' | ... ] <Seq>  
        | <>
```

```
<Id> ::= <Alph> <AlnumSeq>  
<Alph> ::= [ 'a' | ... | 'z' | 'A' | ... | 'Z' | '_' ]  
<Alnum> ::= <Alph> | <Digit>  
<AlnumSeq> ::= <Alnum> <AlnumSeq>  
             | <>
```

Lexical Analysis Grammars

Right
recursive!

$\langle \text{PI} \rangle ::= \langle \text{Digit} \rangle \mid \langle \text{Digit} \rangle \langle \text{PI} \rangle$
 $\langle \text{Digit} \rangle ::= '0' \mid '1' \mid \dots \mid '9'$

$\langle \text{String} \rangle ::= ' "' \langle \text{Seq} \rangle ' "'$
 $\langle \text{Seq} \rangle ::= ['a' \mid 'b' \mid \dots \mid 'z' \mid '0' \mid \dots \mid '9' \mid '+' \mid '-' \mid \dots] \mid \langle \rangle$

$\langle \text{Id} \rangle ::= \langle \text{Alph} \rangle \langle \text{AlnumSeq} \rangle$
 $\langle \text{Alph} \rangle ::= ['a' \mid \dots \mid 'z' \mid 'A' \mid \dots \mid 'Z' \mid '_']$
 $\langle \text{Alnum} \rangle ::= \langle \text{Alph} \rangle \mid \langle \text{Digit} \rangle$
 $\langle \text{AlnumSeq} \rangle ::= \langle \text{Alnum} \rangle \langle \text{AlnumSeq} \rangle \mid \langle \rangle$

Grammars whose recursion is always in the rightmost position in the body of the rule is called *regular*.

Lexical Analysis Grammars

Right
recursive!

$\langle \text{PI} \rangle ::= \langle \text{Digit} \rangle \mid \langle \text{Digit} \rangle \langle \text{PI} \rangle$
 $\langle \text{Digit} \rangle ::= '0' \mid '1' \mid \dots \mid '9'$

$\langle \text{String} \rangle ::= ' "' \langle \text{Seq} \rangle ' "'$
 $\langle \text{Seq} \rangle ::= ['a' \mid 'b' \mid \dots \mid 'z' \mid '0' \mid \dots \mid '9' \mid '+' \mid '-' \mid \dots] \mid \langle \rangle$

$\langle \text{Id} \rangle ::= \langle \text{Alph} \rangle \langle \text{AlnumSeq} \rangle$
 $\langle \text{Alph} \rangle ::= ['a' \mid \dots \mid 'z' \mid 'A' \mid \dots \mid 'Z' \mid '_']$
 $\langle \text{Alnum} \rangle ::= \langle \text{Alph} \rangle \mid \langle \text{Digit} \rangle$
 $\langle \text{AlnumSeq} \rangle ::= \langle \text{Alnum} \rangle \langle \text{AlnumSeq} \rangle \mid \langle \rangle$

Oversimplification,
doesn't consider
mutual recursion!

Grammars whose recursion is always in the rightmost position in the body of the rule is called *regular*.

Regular Languages

A *right regular grammar* (also called *right linear grammar*) is a grammar (N, Σ, P, S) such that all the production rules in P are one of the following forms:

- (a) $B \rightarrow a$, where $B \in N$ and $a \in \Sigma$.
- (b) $B \rightarrow aC$, where $B, C \in N$ and $a \in \Sigma$.
- (c) $B \rightarrow \epsilon$, where $B \in N$ and ϵ is the empty string.

Regular Languages

A *right regular grammar* (also called *right linear grammar*) is a grammar (N, Σ, P, S) such that all the production rules in P are one of the following forms:

- (a) $B \rightarrow a$, where $B \in N$ and $a \in \Sigma$.
- (b) $B \rightarrow aC$, where $B, C \in N$ and $a \in \Sigma$.
- (c) $B \rightarrow \epsilon$, where $B \in N$ and ϵ is the empty string.

A similar definition exists for *left regular grammars*. Given a left regular grammar G_L , there exists a right regular grammar G_R such that $\mathcal{L}(G_L) = \mathcal{L}(G_R)$.

Regular Languages

A *right regular grammar* (also called *right linear grammar*) is a grammar (N, Σ, P, S) such that all the production rules in P are one of the following forms:

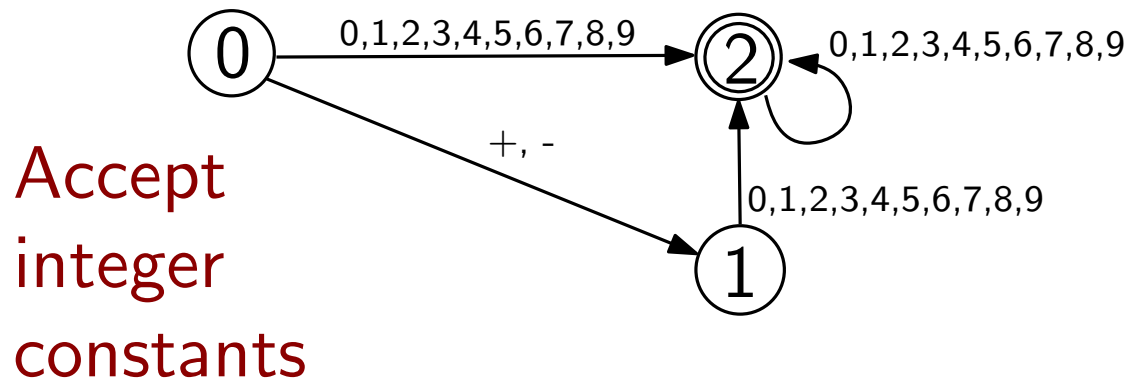
- (a) $B \rightarrow a$, where $B \in N$ and $a \in \Sigma$.
- (b) $B \rightarrow aC$, where $B, C \in N$ and $a \in \Sigma$.
- (c) $B \rightarrow \epsilon$, where $B \in N$ and ϵ is the empty string.

A similar definition exists for *left regular grammars*. Given a left regular grammar G_L , there exists a right regular grammar G_R such that $\mathcal{L}(G_L) = \mathcal{L}(G_R)$.

We shall only focus on *right regular grammars*. All languages generated by right or left regular grammars are called *regular languages*.

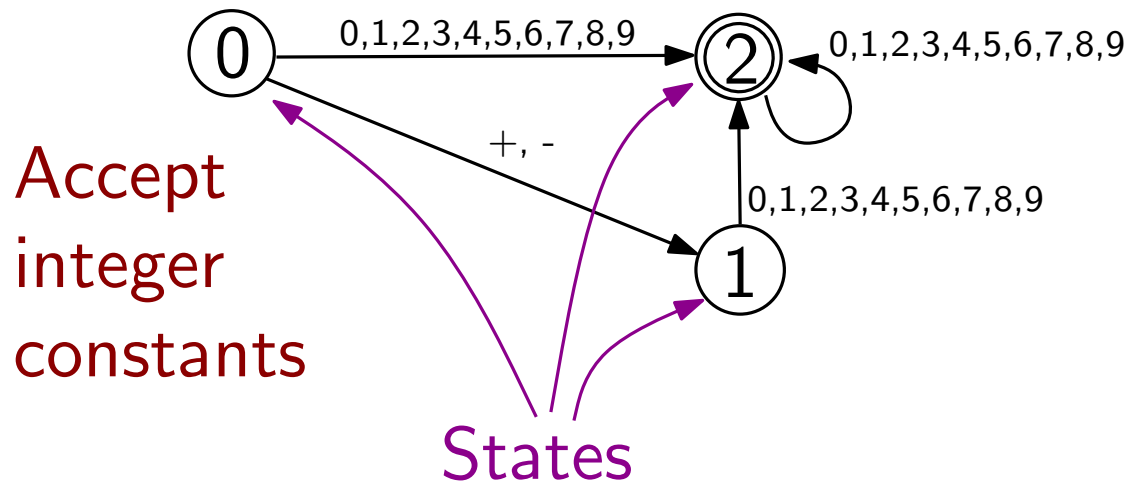
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



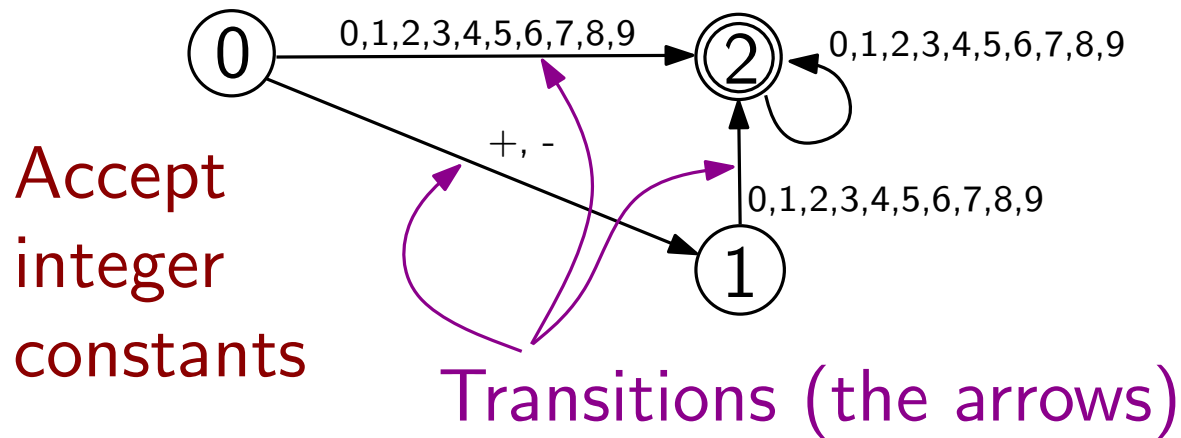
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



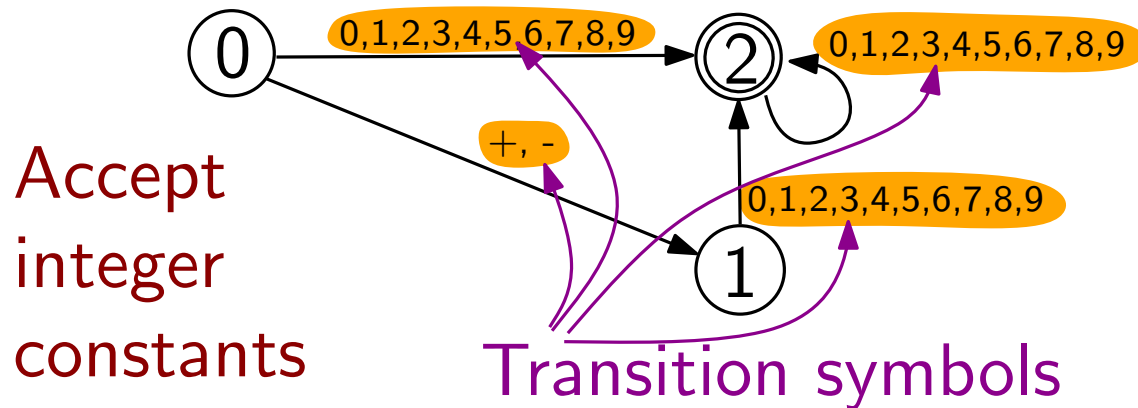
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



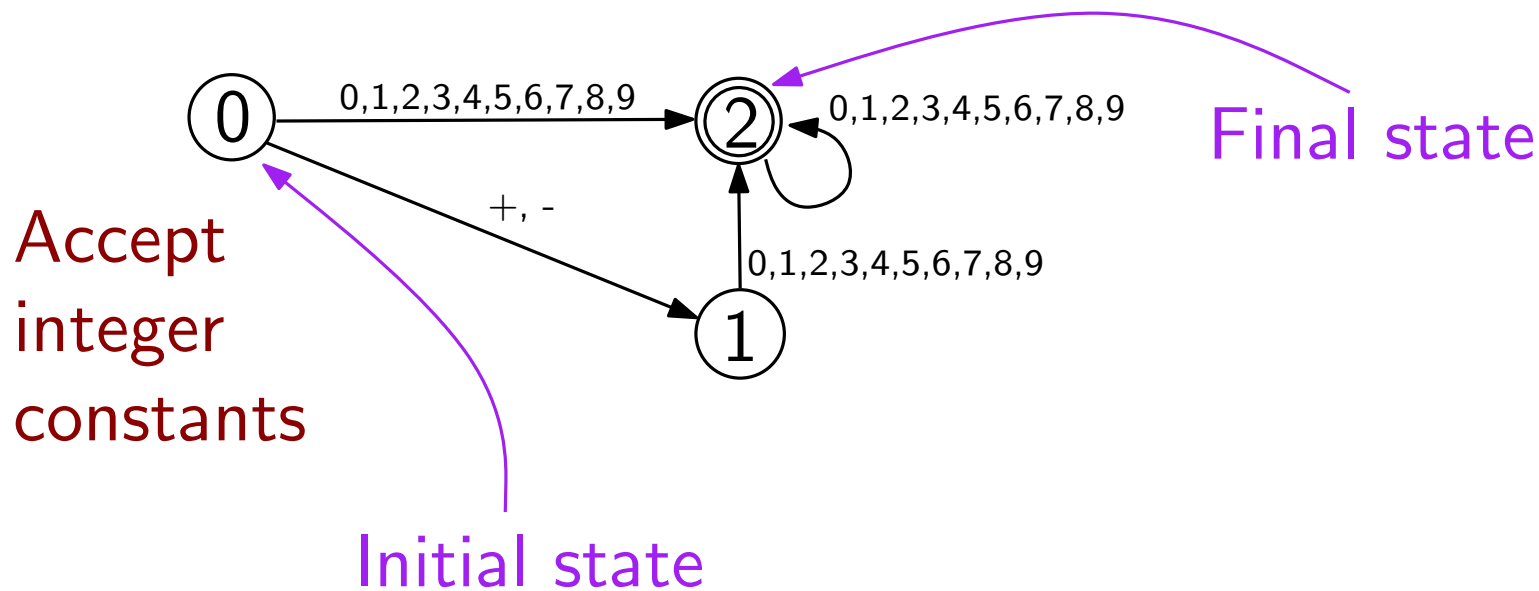
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Deterministic Finite Automata

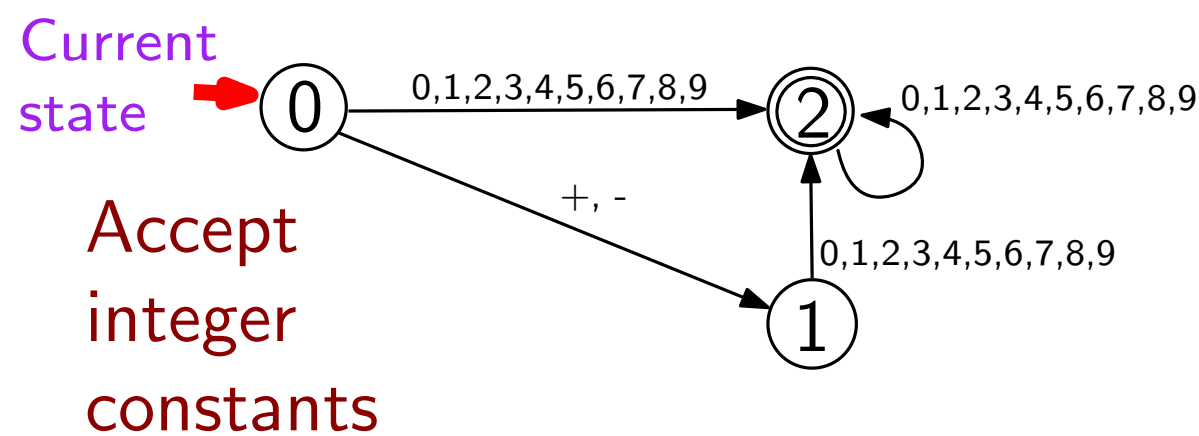
Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



(usually labelled 0 or q_0 , and appearing in the top left corner of diagram)

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



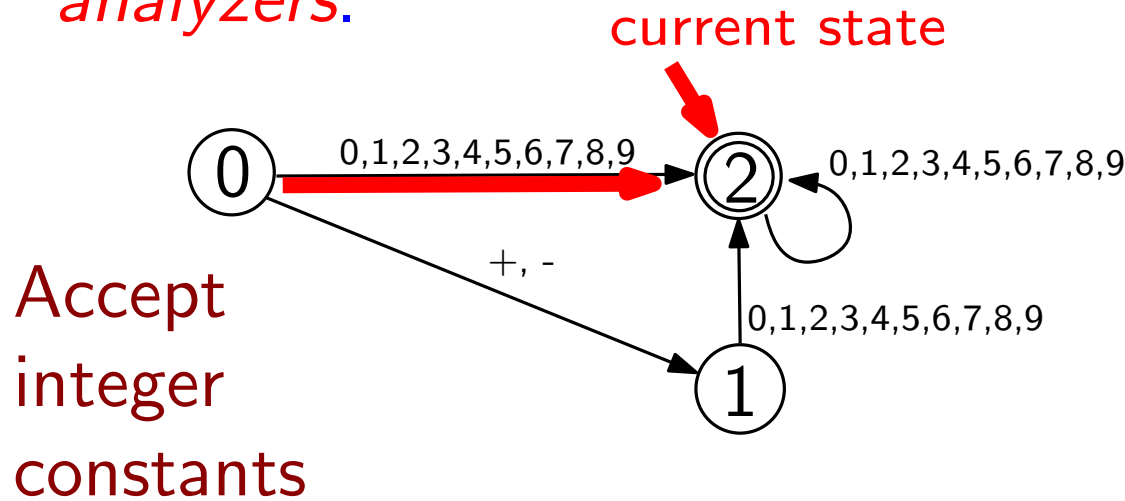
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



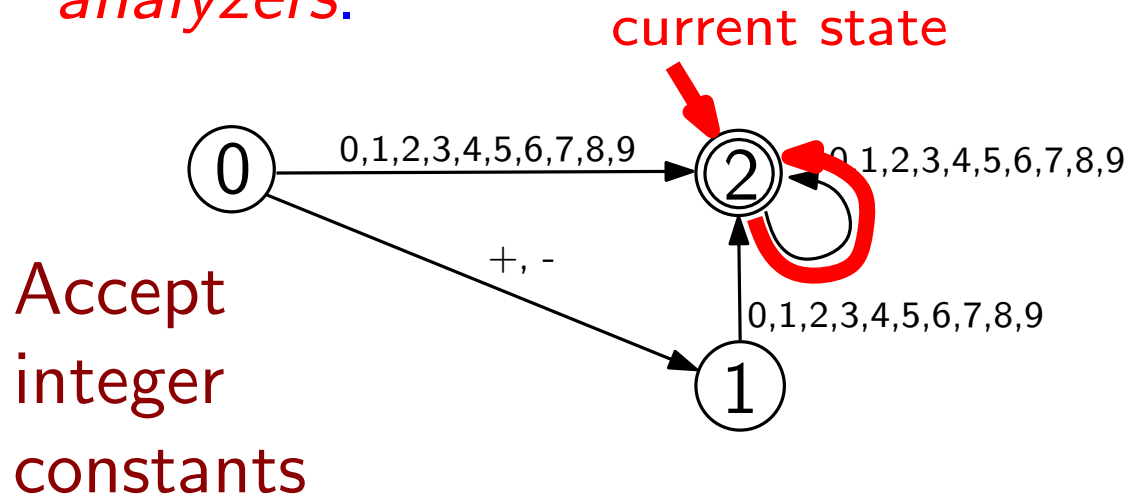
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



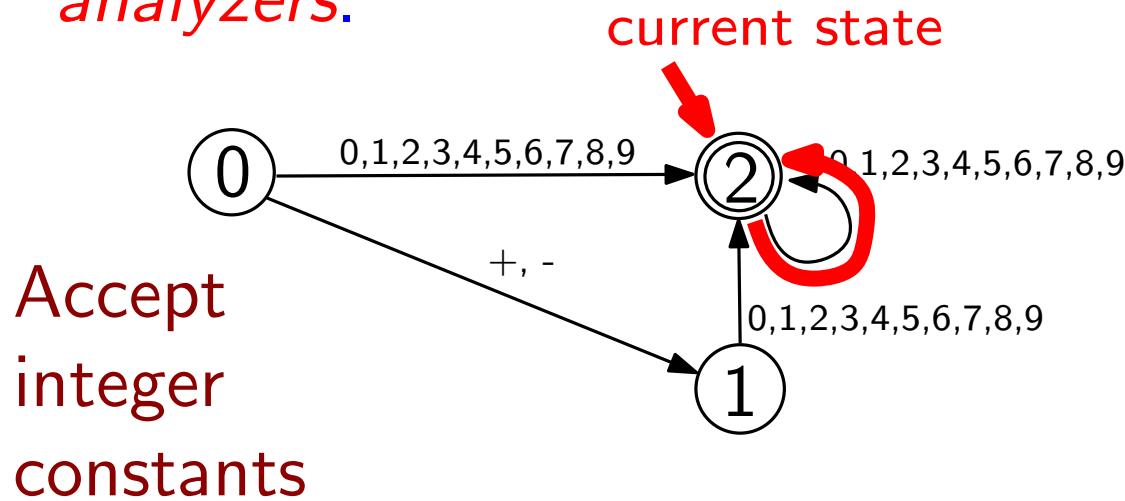
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



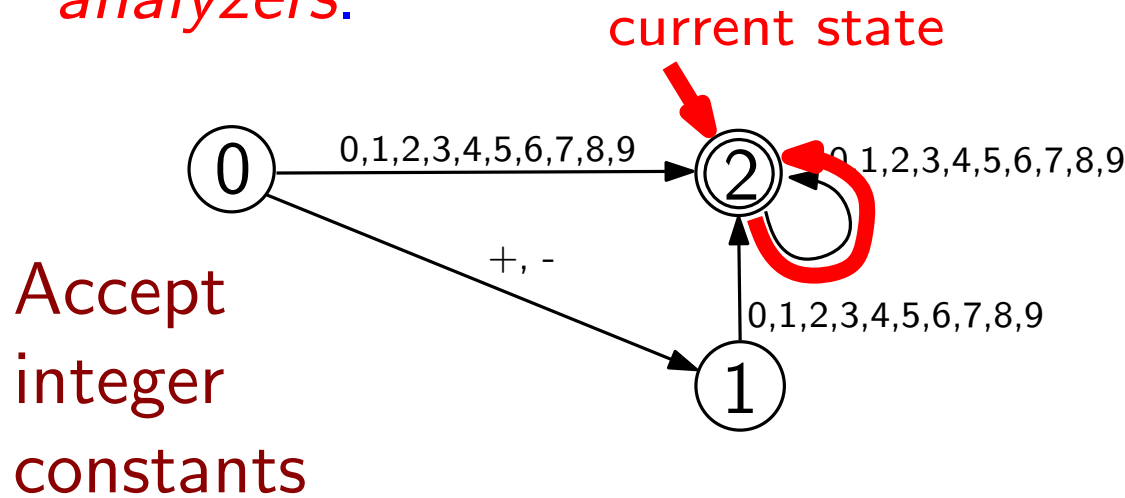
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



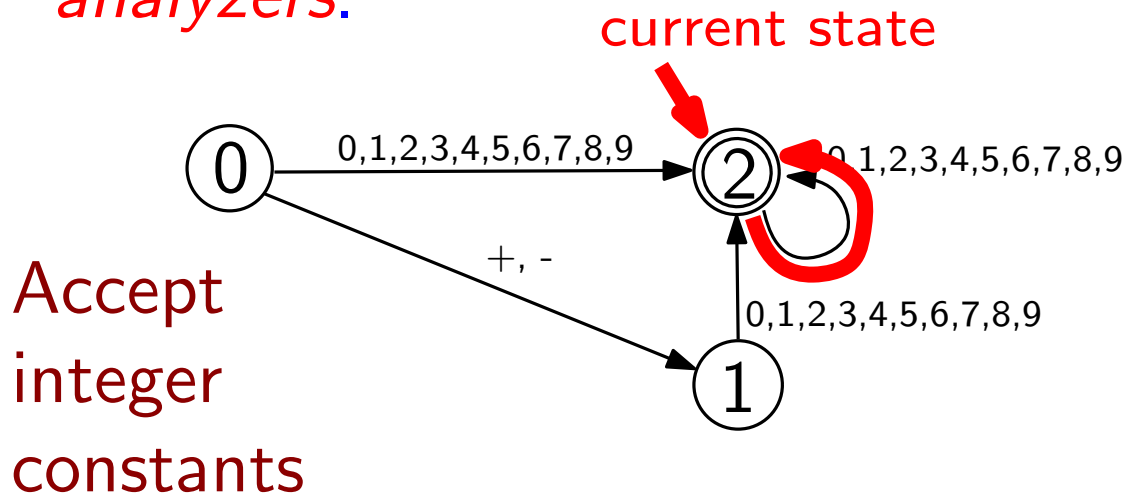
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



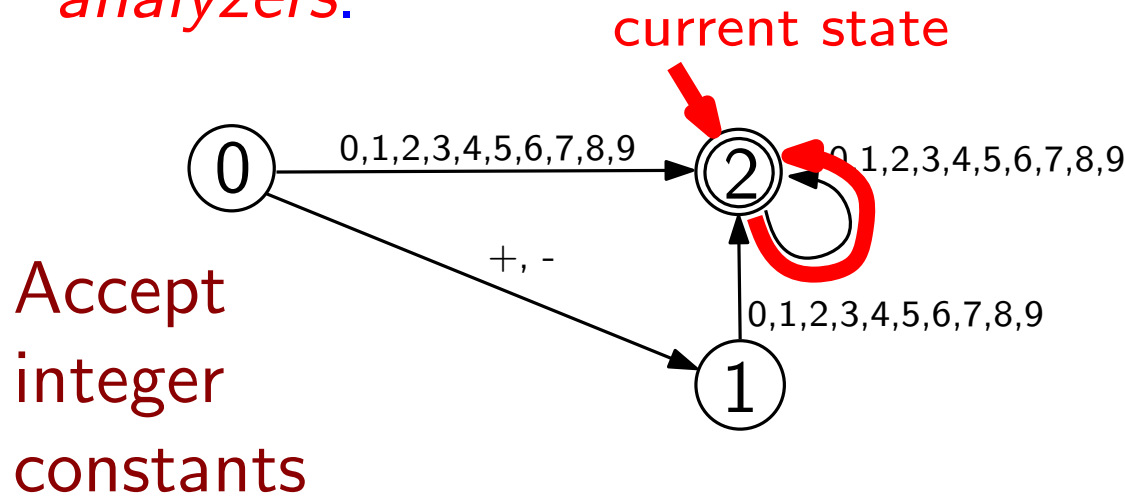
Simulation

Input: 0 1 2 3 4 5

Input "head" ↑

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Simulation

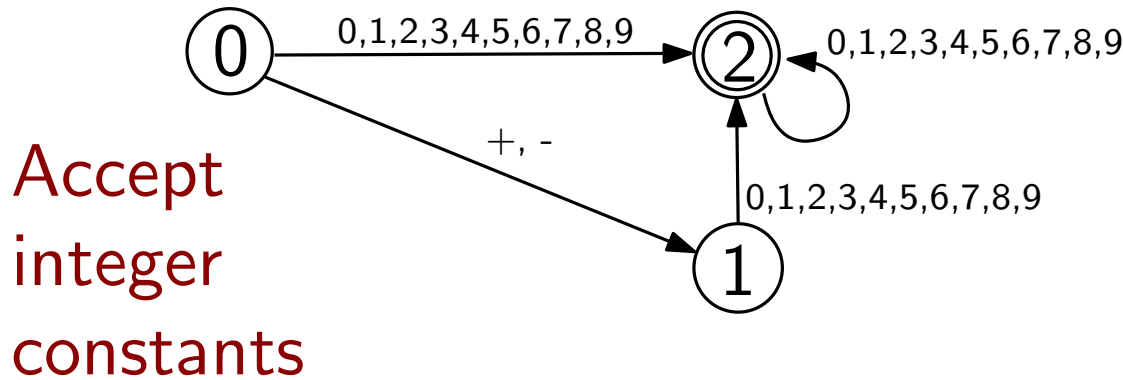
Input: 0 1 2 3 4 5

Input "head" ↑

If at the end of the input, or whenever the "head" encounters a symbol that is not a current transition symbol, the machine is in a final state, the string is "*accepted*". Otherwise, the string is "*rejected*".

Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



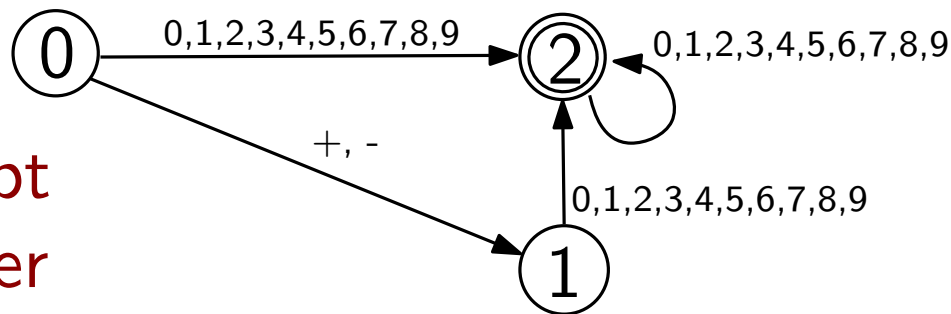
Equivalent regular grammar:

$$\begin{aligned} \langle 0 \rangle &::= [\text{'0'} \mid \dots \mid \text{'9'}] \langle 2 \rangle \\ &\quad \mid [\text{'+'} \mid \text{'-'}] \langle 1 \rangle \\ \langle 1 \rangle &::= [\text{'0'} \mid \dots \mid \text{'9'}] \langle 2 \rangle \\ \langle 2 \rangle &::= [\text{'0'} \mid \dots \mid \text{'9'}] \langle 2 \rangle \\ &\quad \mid \langle \rangle \end{aligned}$$

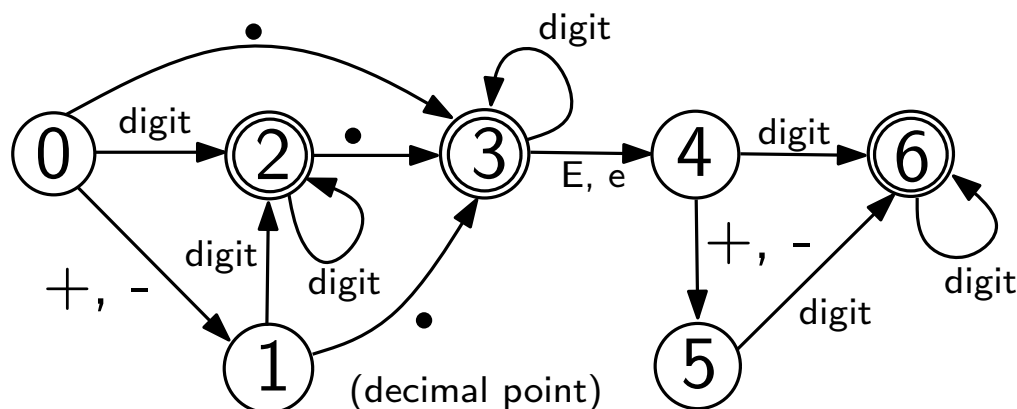
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.

Accept
integer
constants

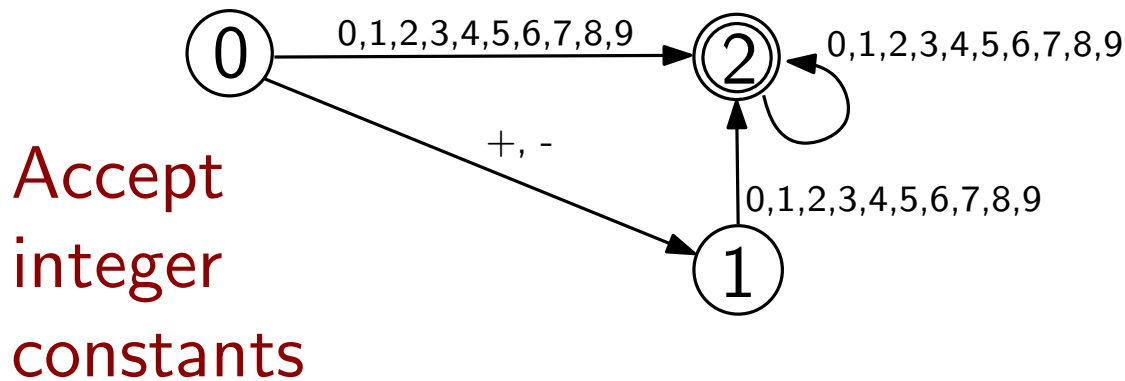


Accept real
constants



Deterministic Finite Automata

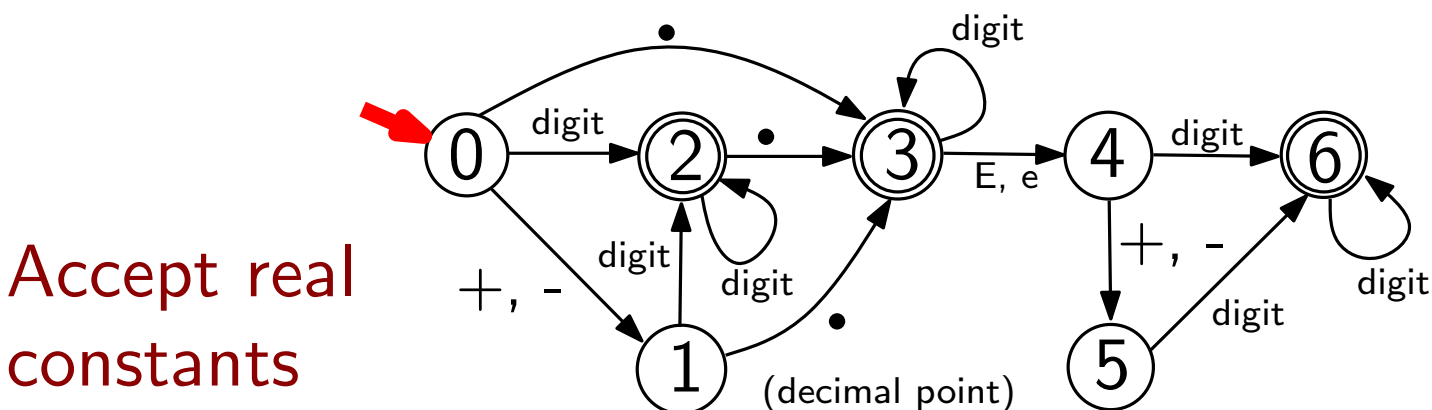
Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Simulation

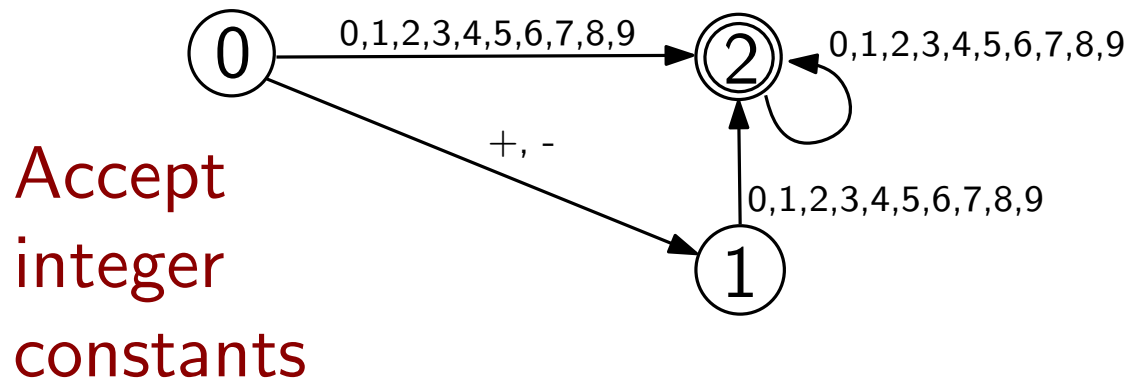
Input: - 1 . 2 E + 3

↑
Input "head"



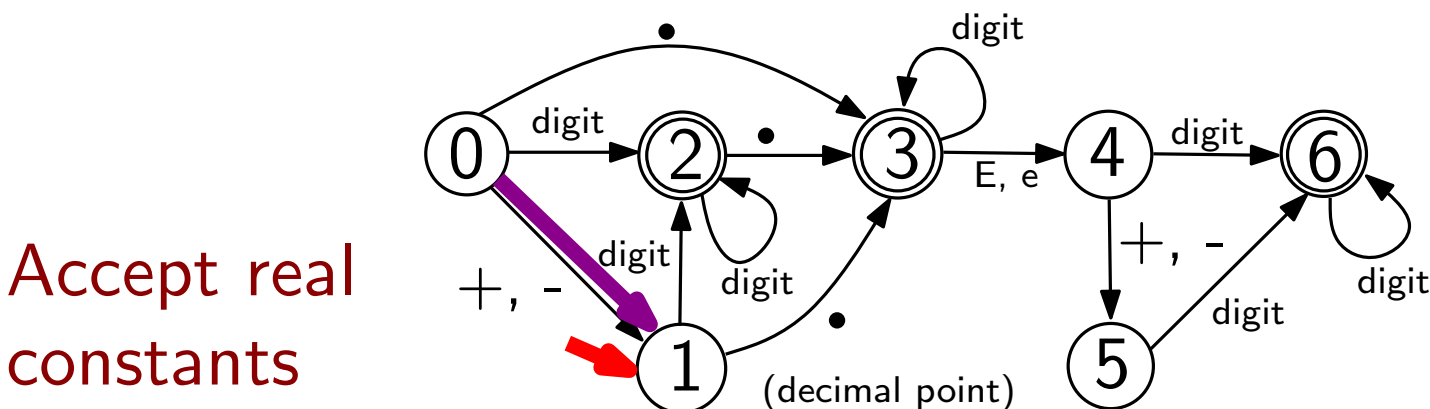
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



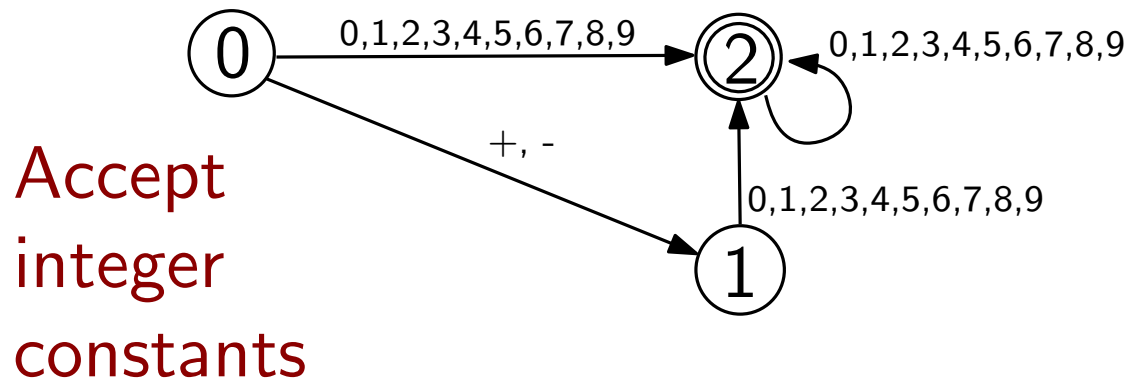
Simulation

Input: - 1 . 2 E + 3
 ↑
Input "head"



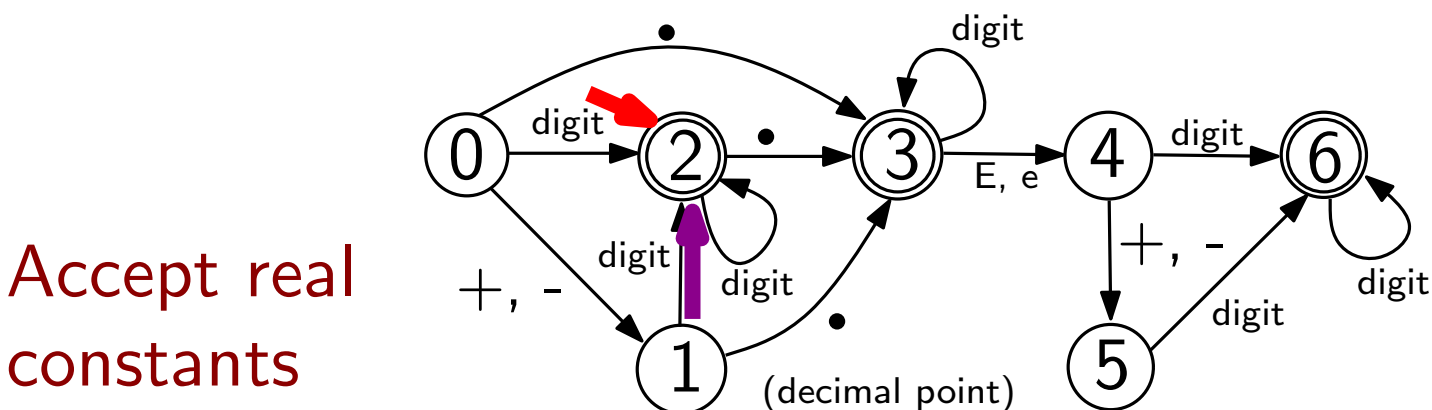
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



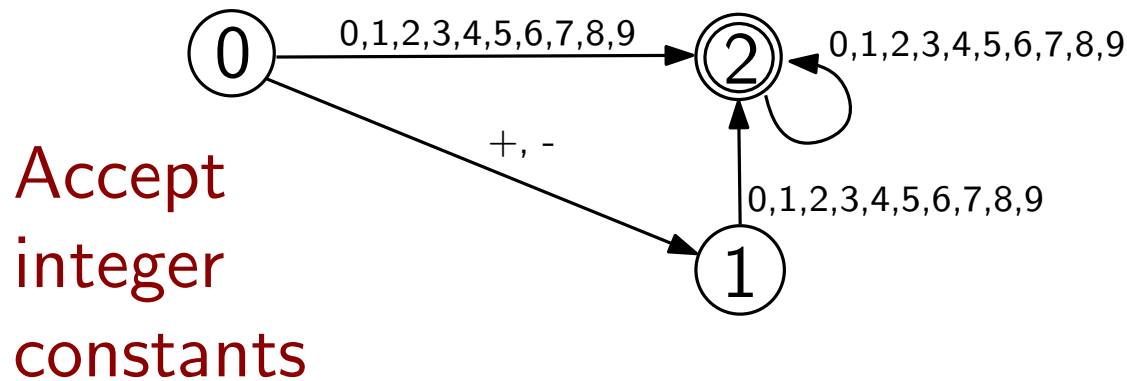
Simulation

Input: - 1 . 2 E + 3
Input "head" ↑



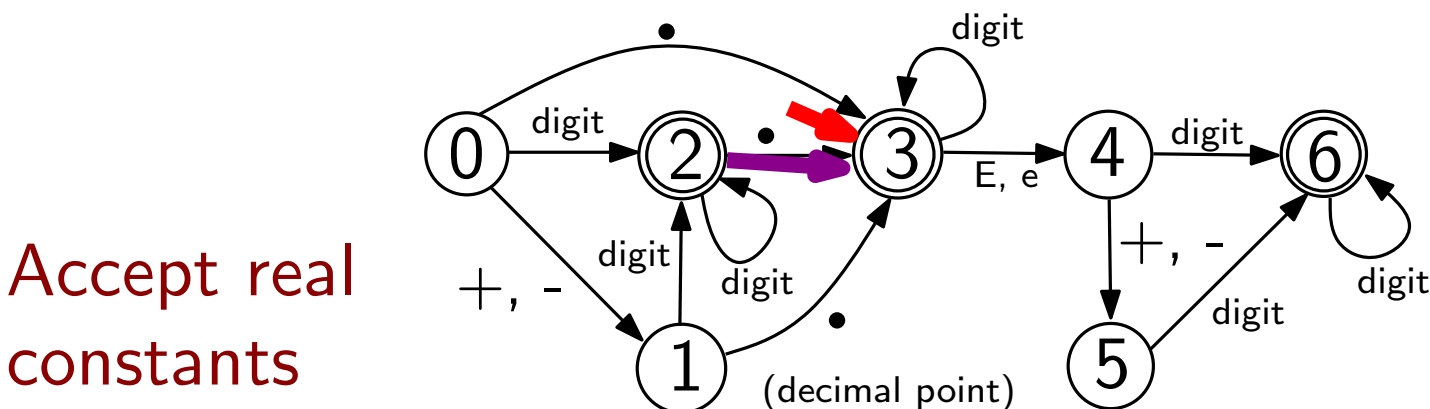
Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



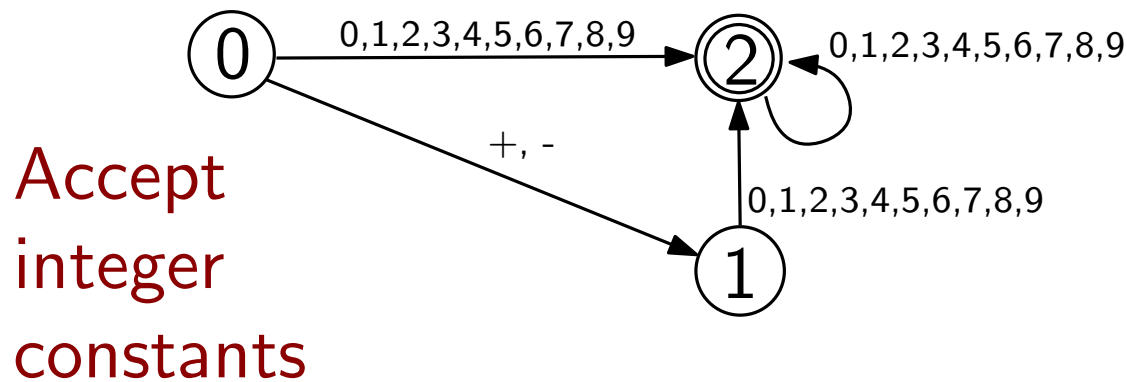
Simulation

Input: - 1 . 2 E + 3
Input "head" ↑



Deterministic Finite Automata

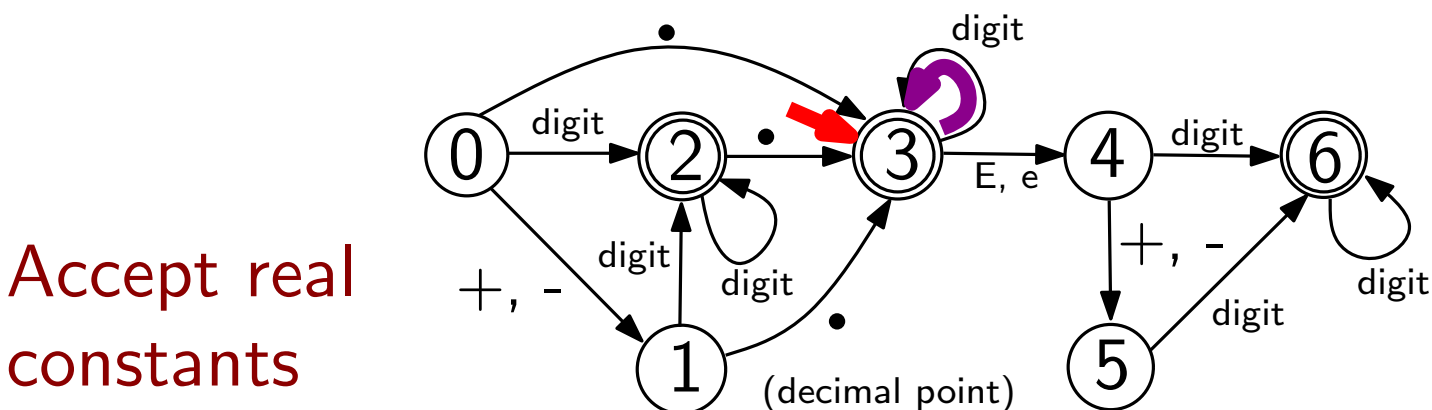
Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Simulation

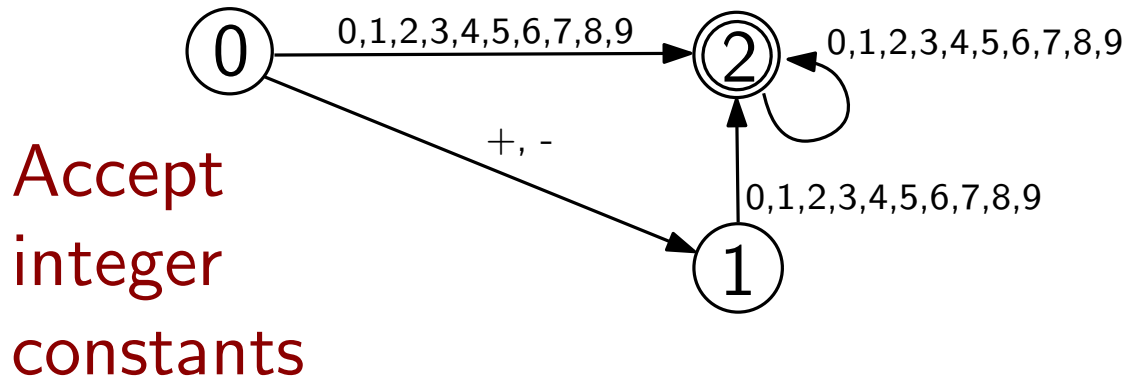
Input: - 1 . 2 E + 3

Input "head" ↑



Deterministic Finite Automata

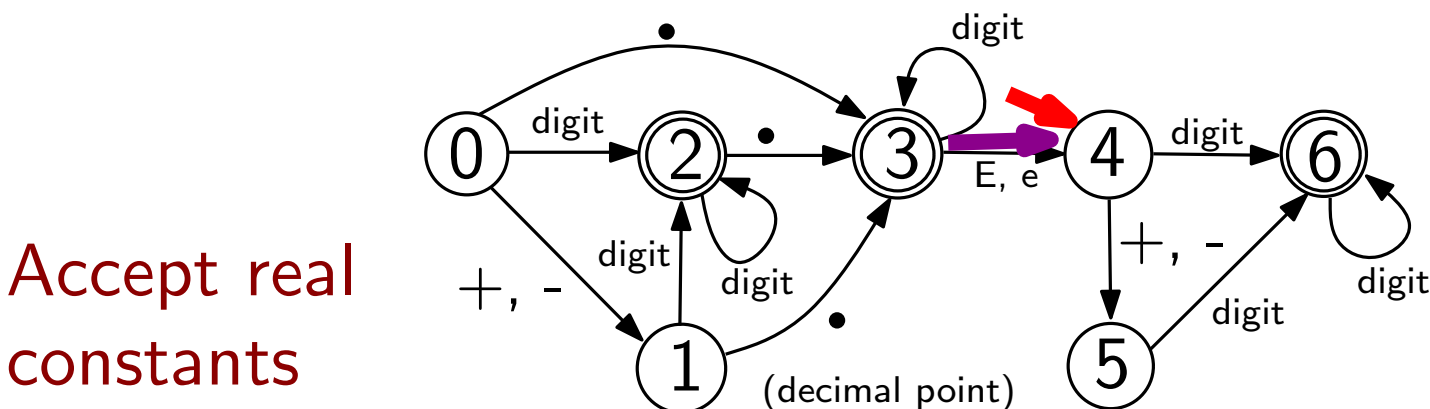
Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Simulation

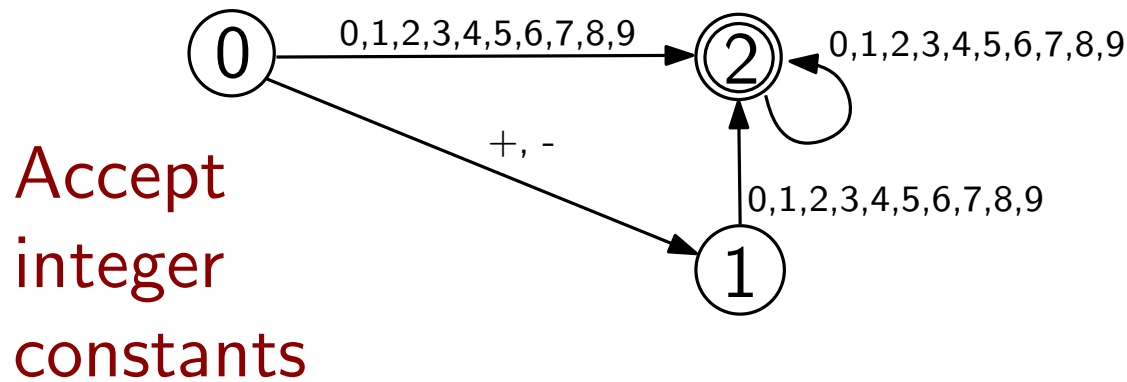
Input: - 1 . 2 E + 3

Input "head" ↑



Deterministic Finite Automata

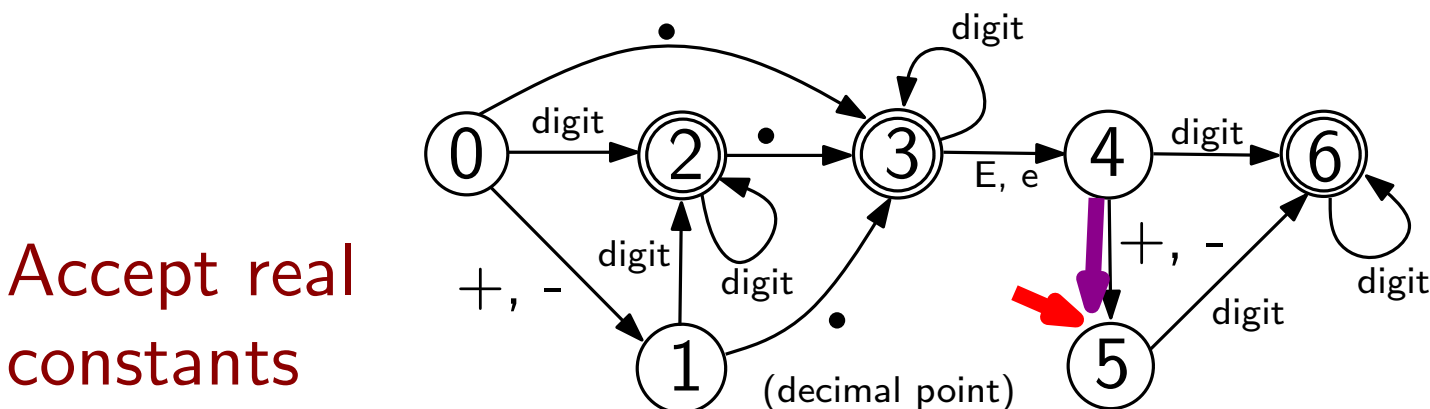
Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.



Simulation

Input: - 1 . 2 E + 3

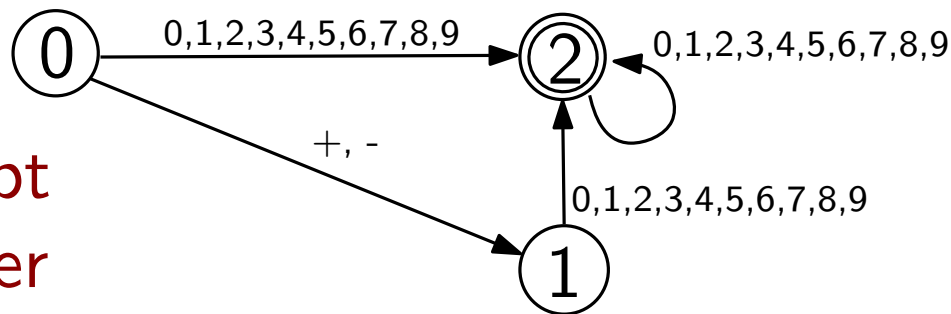
Input "head" ↑



Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.

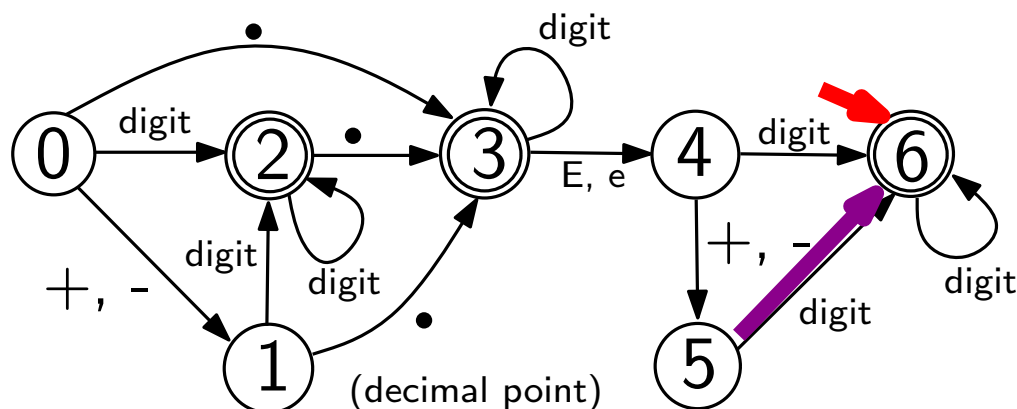
Accept
integer
constants



Simulation

Input: - 1 . 2 E + 3
Input "head" ↑

Accept real
constants



Deterministic Finite Automata

Regular languages can be accepted by *Deterministic Finite Automata*, which is the preferred way of implementing *lexical analyzers*.

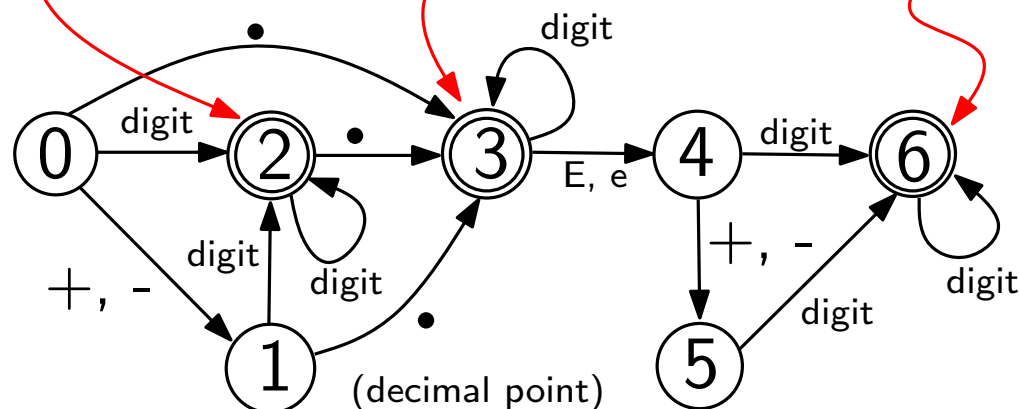
Automata can be bundled together!

Attach a *lexeme type* to each final state.

Integer constant

Fixed point real constant

Floating point real constant



Accept real constants

What have we learned

- ◇ DFAs are essentially graphs, where nodes are states, and arcs are transitions; they also have an input tape on which the input string is placed – a "head" reads one symbol at a time.
- ◇ The automaton starts in its initial state, and has its head on the first symbol of its "input tape".
- ◇ A transition is performed only if the current state has an outgoing arrow labeled with the current symbol.
- ◇ The result of the transition is a new current state, pointed to by the current arc; the "head" advances by one symbol during the transition.
- ◇ If there is no outgoing transition for the current symbol, the automaton stops.
- ◇ If the automaton stops in a final state, the part of the string that was analyzed so far is turned into a *lexeme* with the type indicated by the final state's label.
- ◇ If the automaton does not stop in a final state, that indicates an error that stops the lexical analysis process.

Regular Expressions

- ◇ Language that defines languages (like the BNF)
- ◇ Equivalent to regular grammars and DFAs
 - For every RE E , there exists an RG G and a DFA A s.t. $\mathcal{L}(E) = \mathcal{L}(G) = \mathcal{L}(A)$. The converse is also true.
- ◇ *Definition:* Let $\Sigma = \{a, b, c, \dots\}$ be an alphabet of symbols.
 - ϵ is an RE with $\mathcal{L}(\epsilon) = \{\epsilon\}$
 - For each $x \in \Sigma$, x is an RE with $\mathcal{L}(x) = \{x\}$
 - Given two REs r and s , $r|s$ is an RE with $\mathcal{L}(r|s) = \mathcal{L}(r) \cup \mathcal{L}(s)$.
 - Given two REs r and s , rs is an RE with $\mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s)$.
 - Given an RE r , r^* is an RE with $\mathcal{L}(r^*) = \bigcup_{i \geq 0} (\mathcal{L}(r))^i$.

Regular Expressions

- ◇ Language that defines languages (like the BNF)

- ◇ Equivalent to regular grammars and DFAs

- For e
s.t. $\mathcal{L}($

$$L_1 \cdot L_2 = \{s_1 s_2 \mid s_1 \in L_1 \text{ and } s_2 \in L_2\}$$

A
ue.

- ◇ *Definition*
symbols.

Concatenation of languages

- ϵ is a
- For each $x \in \Sigma$, x is an RE with $\mathcal{L}(x) = \{x\}$
- Given two REs r and s , $r|s$ is an RE with $\mathcal{L}(r|s) = \mathcal{L}(r) \cup \mathcal{L}(s)$.
- Given two REs r and s , rs is an RE with $\mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s)$.
- Given an RE r , r^* is an RE with $\mathcal{L}(r^*) = \bigcup_{i \geq 0} (\mathcal{L}(r))^i$.

Regular Expressions

- ◇ Language that defines languages (like the BNF)
- ◇ Equivalent to regular grammars and DFAs
 - For every RE E , there exists an RG G and a DFA A

◇ Def
syn

$$L^0 = \{\epsilon\}$$
$$L^i = \underbrace{L \cdot L \cdot \dots \cdot L}_{i \text{ times}, i \geq 1}$$

$$\mathcal{L}(r|s) = \mathcal{L}(r) \cup \mathcal{L}(s).$$

- Given two REs r and s , rs is an RE with

$$\mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s).$$

- Given an RE r , r^* is an RE with $\mathcal{L}(r^*) = \bigcup_{i \geq 0} (\mathcal{L}(r))^i$.

Regular Expressions

- ◇ Language that defines languages (like the BNF)
- ◇ Equivalent to regular grammars and DFAs
 - For every RE E , there exists an RG G and a DFA A s.t. $\mathcal{L}(E) = \mathcal{L}(G) = \mathcal{L}(A)$. The converse is also true.

◇ **Definition:** Let $\Sigma =$ symbols.

- ϵ is an RE with
- For each $x \in \Sigma$,
- Given two REs r and s , $r|s$ is an RE with $\mathcal{L}(r|s) = \mathcal{L}(r) \cup \mathcal{L}(s)$.
- Given two REs r and s , rs is an RE with $\mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s)$.
- Given an RE r , r^* is an RE with $\mathcal{L}(r^*) =$

$$\{s_1 \cdots s_k \mid k \geq 0, s_i \in \mathcal{L}(r), 1 \leq i \leq k\}$$

Any number of strings from $\mathcal{L}(r)$ concatenated.

$$\bigcup_{i \geq 0} (\mathcal{L}(r))^i$$

Regular Expressions

Examples

$$\mathcal{L}((a|b|c)d) = \{ad, bd, cd\}$$

$$\mathcal{L}((a|b)^*) = \{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$$

$$\mathcal{L}((ab)^*) = \{\epsilon, ab, abab, ababab, \dots\}$$

$$\mathcal{L}((ab)^*c) = \{c, abc, ababc, abababc, \dots\}$$

$$\mathcal{L}(a a^* b) = \{ab, aab, aaab, \dots\}$$

◇ *Definition:* Let $\Sigma = \{a, b, c, \dots\}$ be an alphabet of symbols.

- ϵ is an RE with $\mathcal{L}(\epsilon) = \{\epsilon\}$
- For each $x \in \Sigma$, x is an RE with $\mathcal{L}(x) = \{x\}$
- Given two REs r and s , $r|s$ is an RE with $\mathcal{L}(r|s) = \mathcal{L}(r) \cup \mathcal{L}(s)$.
- Given two REs r and s , rs is an RE with $\mathcal{L}(rs) = \mathcal{L}(r) \cdot \mathcal{L}(s)$.
- Given an RE r , r^* is an RE with $\mathcal{L}(r^*) = \bigcup_{i \geq 0} (\mathcal{L}(r))^i$.

The Language Ruby

- ◇ Developed as the pet project of Japanese programmer Yukihiro "Matz" Matsumoto.
- ◇ Multiparadigm: functional-imperative object-oriented reflective language.
- ◇ One of the few languages where regular expressions are *first-class objects*
 - can be assigned to variables
 - can be parameters to functions
 - its methods can be invoked
- ◇ Available as either interactive interpreter (command: `irb`) with *read-print-eval* loop, or as non-interactive script interpretation utility (command: `ruby`).
- ◇ Great sandbox for practicing regular expressions.
 - Watch the video for demo

Ruby Regular Expressions

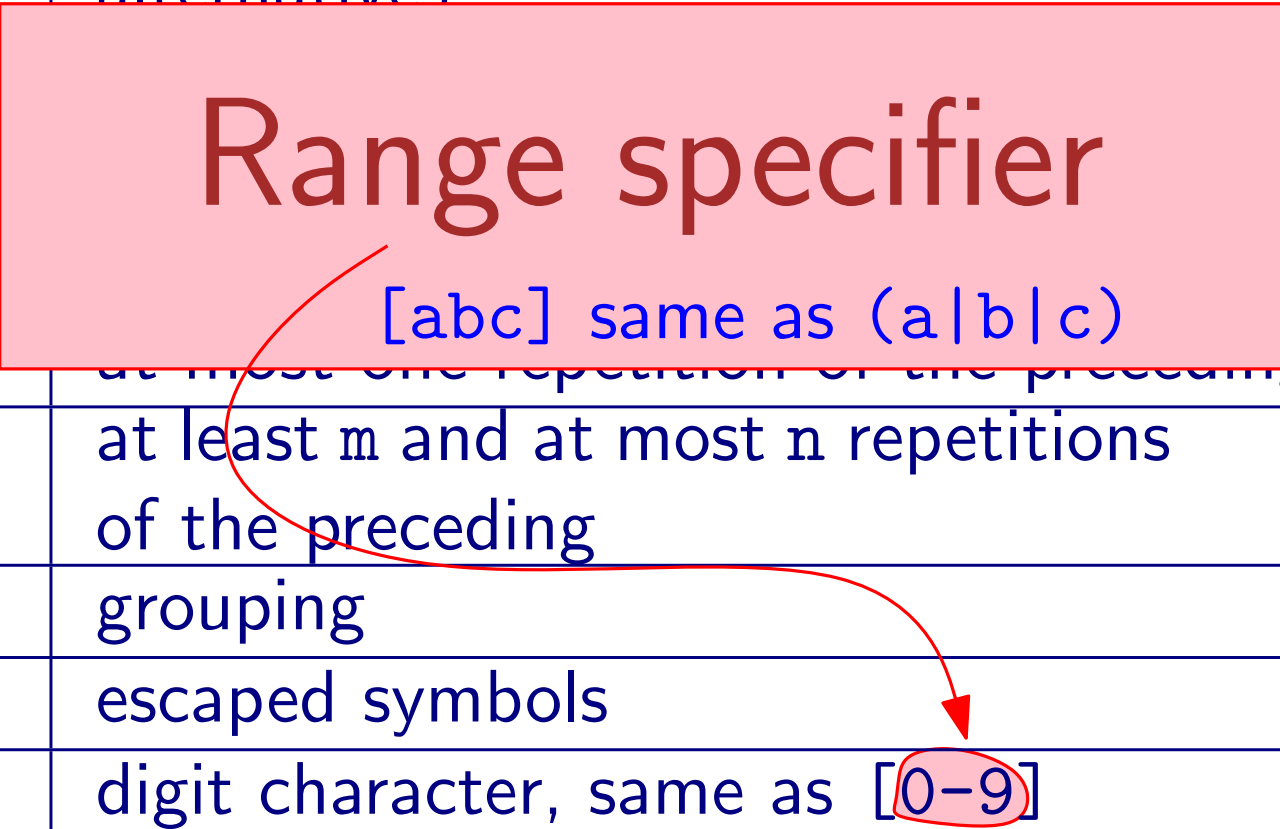
<code>()</code>	the empty string
<code> </code>	alternatives
<code>[]</code>	one-symbol alternatives
<code>*</code>	zero or more repetitions of the preceding
<code>+</code>	one or more repetitions of the preceding
<code>?</code>	at most one repetition of the preceding
<code>{m,n}</code>	at least m and at most n repetitions of the preceding
<code>(...)</code>	grouping
<code>\(, \), \ , ...</code>	escaped symbols
<code>\d</code>	digit character, same as <code>[0-9]</code>
<code>\D</code>	non-digit character, same as <code>[^0-9]</code>
<code>\w</code>	word character, same as <code>[0-9A-Za-z_]</code>
<code>\W</code>	non-word character, same as <code>[^0-9A-Za-z_]</code>

Ruby Regular Expressions

()	the empty string
	alternatives
[]	
*	
+	
?	
{m,n}	
(...)	
\(, \), \ , ...	
\d	
\D	
\w	
\W	

Range specifier

[abc] same as (a|b|c)



Ruby Regular Expressions

()	the empty string
	alternatives
[]	one-symbol alternatives
*	zero or more repetitions of the preceding
+	one or more repetitions of the preceding
?	zero or one repetition of the preceding
{m}	m repetitions
(.	capture
\(, \), \ , ...	escaped symbols
\d	digit character, same as [0-9]
\D	non-digit character, same as [^0-9]
\w	word character, same as [0-9A-Za-z_]
\W	non-word character, same as [^0-9A-Za-z_]

Exclusion operator

must appear at beginning of [...]

Ruby Regular Expressions

()	the empty string
	alternatives
[]	one-symbol alternatives
*	zero or more
+	one or more
?	at most one
{m,n}	at least m, at most n of the preceding
(...)	grouping
\(, \), \ , ...	escaped
\d	digit character, same as [0-9]
\D	non-digit character, same as [^0-9]
\w	word character, same as [0-9A-Za-z_]
\W	non-word character, same as [^0-9A-Za-z_]

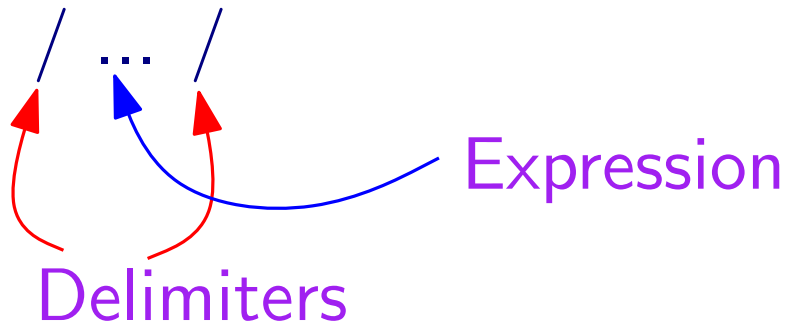
Multiple ranges

Ruby's RE usage

Syntax: `/ ... /`

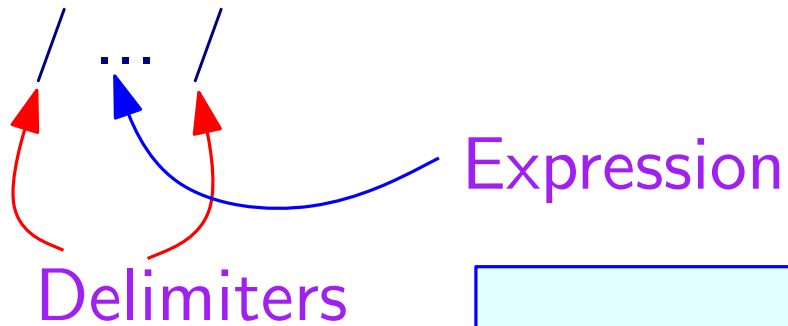
Ruby's RE usage

Syntax:



Ruby's RE usage

Syntax:



Example: `/ab*([cd]+|e?)/`

Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

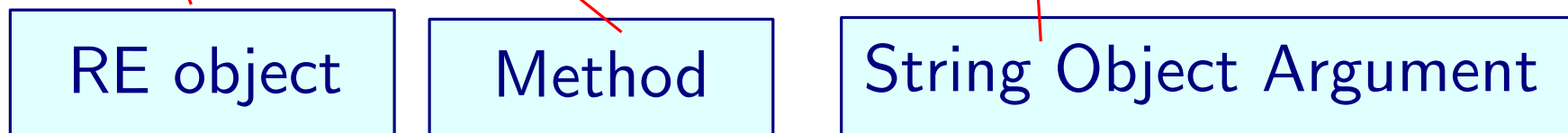
RE matching: `x.match "1abbc2"`

Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

RE matching: `x.match "1abbc2"`



Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

RE matching: `x.match "1abbc2"`

Evaluates to: `"abbc"`

Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

RE matching: `x.match "1abbc2"`

Evaluates to: ~~"abbc"~~

The result is of type `MatchData`

Must be converted to `String` to be visualized

Ruby's RE usage

Syntax: `/ ... /`

First-class value: `x = /ab*c/`

RE matching: `(x.match "1abbc2").to_s`

Evaluates to: `"abbc"`



The result is of type `MatchData`

Must be converted to `String` to be visualized

What Have We Learned?

- ◇ Ruby is an object oriented language
 - Expressions have the form
`obj1.method(obj2)` (brackets optional)
instead of
`method(obj1,obj2)`
- ◇ Regular expressions are first class objects.
- ◇ Method `match` can be invoked to match an RE with a string.
- ◇ Result is the first occurrence of a matched substring (must be converted using `to_s`).
- ◇ REs are specified by RE-specific language.
- ◇ RE expressions must be enclosed between forward slashes.

Conclusion

- ◇ Programming languages are specified by *grammars*, in a stratified manner.
- ◇ The lower level, that of *lexical analysis*, uses *regular grammars*, and their counterpart, *regular expressions*.
 - convert a program into a sequence of *lexemes* — more in tutorial
- ◇ The higher level, called *syntactic analysis* uses more sophisticated grammars.
 - capture the structure of the language;
 - use *lexemes* as terminals
 - shall be covered in more detail next time
- ◇ *Regular expressions* are a basic data type in Ruby.
 - Ruby can be used to build a *toy lexer*.
 - Examples in the tutorial.