



PLDI 2016 Tutorial

Automata-Based String Analysis

Tevfik Bultan, Abdulbaki Aydin, Lucas Bang
Verification Laboratory
<http://vlab.cs.ucsb.edu>
Department of Computer Science

Common Usages of Strings

- ▶ Input validation and sanitization



Google Search I'm Feeling Lucky

- ▶ Database query generation



- ▶ Formatted data generation



- ▶ Dynamic code generation

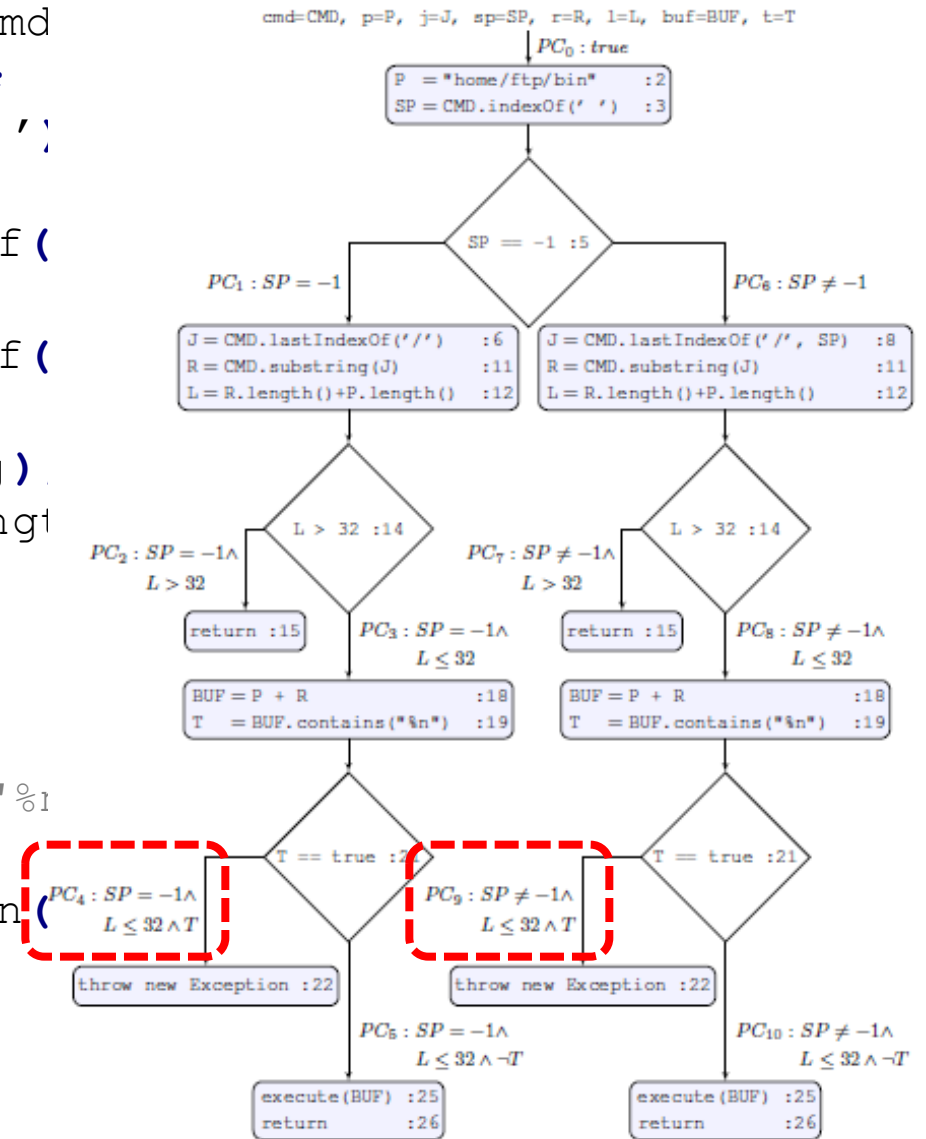


Symbolic Execution Example

```

1  public void site_exec(String cmd
2      String p = "home/ftp/bin";
3      int j, sp = cmd.indexOf(' ');
5      if (sp == -1) {
6          j = cmd.lastIndexOf(
7      } else {
8          j = cmd.lastIndexOf(
9      }
11     String r = cmd.substring(j);
12     int l = r.length() + p.length
14     if (l > 32) {
15         return;
16     }
18     String buf = p + r;
19     boolean t = buf.contains("%1
21     if (t == true) {
22         throw new Exception(
23     }
25     execute(buf);
26     return;
27 }

```



Model Counting String Constraint Solver

INPUT

string
constraint:

C

**Automata-Based
model Counting
string constraint
solver
(ABC)**

OUTPUT

counting
function:

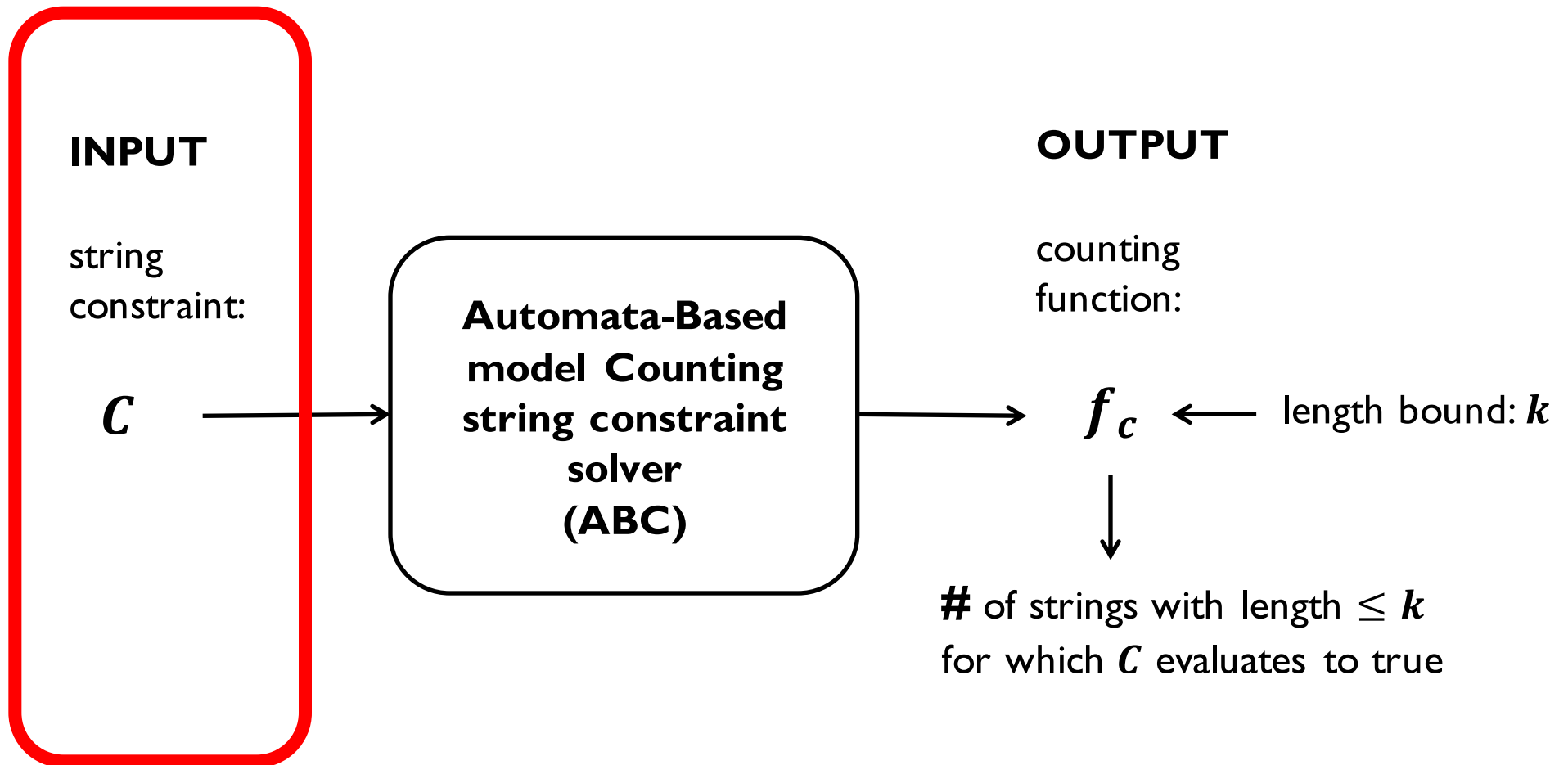
f_c

← length bound: k



of strings with length $\leq k$
for which C evaluates to true

Model Counting String Constraint Solver



String Constraint Language

C → *bterm*

bterm → *v* | true | false
| \neg *bterm* | *bterm* ∧ *bterm* | *bterm* ∨ *bterm* | (*bterm*)
| *sterm* = *sterm*
| match(*sterm*, *sterm*)
| contains(*sterm*, *sterm*)
| begins(*sterm*, *sterm*)
| ends(*sterm*, *sterm*)
| *iterm* = *iterm* | *iterm* < *iterm* | *iterm* > *iterm*

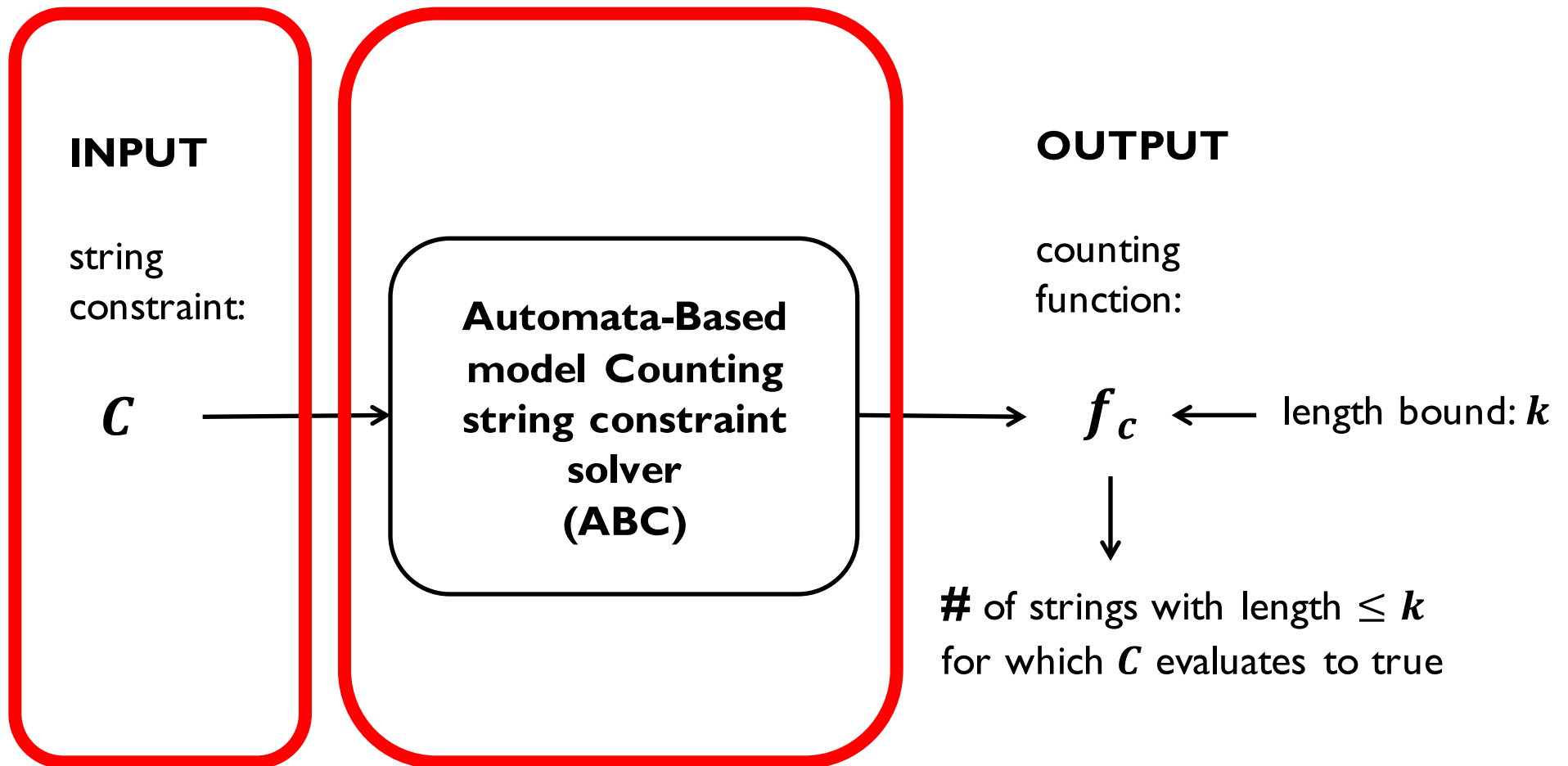
iterm → *v* | *n*
| *iterm* + *iterm* | *iterm* - *iterm* | *iterm* × *n* | (*iterm*)
| length(*sterm*) | toint(*sterm*)
| indexof(*sterm*, *sterm*)
| lastindexof(*sterm*, *sterm*)

sterm → *v* | ε | *s*
| *sterm*.*sterm* | *sterm*|*sterm* | *sterm** | (*sterm*)
| charat(*sterm*, *iterm*) | tostring(*iterm*)
| toupper(*sterm*) | tolower(*sterm*)
| substring(*sterm*, *iterm*, *iterm*)
| replacefirst(*sterm*, *sterm*, *sterm*)
| replacelast(*sterm*, *sterm*, *sterm*)
| replaceall(*sterm*, *sterm*, *sterm*)

Example String Expressions

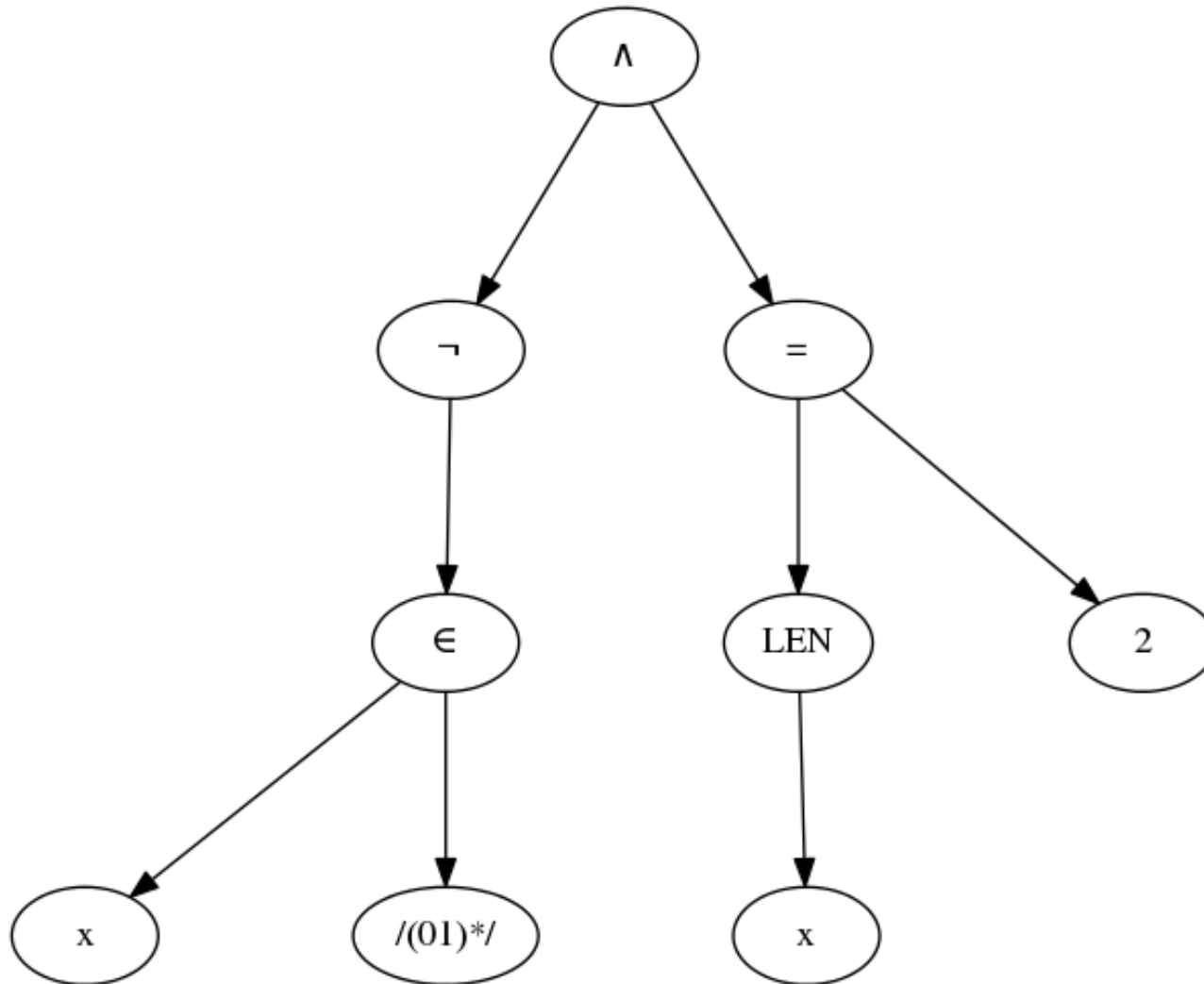
	String Expression	Constraint Language
Java	<code>s.length()</code>	<code>length(s)</code>
	<code>s.isEmpty()</code>	<code>length(s) == 0</code>
	<code>s.startsWith(t,n)</code>	$0 \leq n \wedge n \leq s \wedge$ <code>begins(substring(s,n, s),t)</code>
	<code>s.indexOf(t,n)</code>	<code>indexof(substring(s,n, s),t)</code>
	<code>s.replaceAll(p,r)</code>	<code>replaceall(s,p,r)</code>
PHP	<code>strrpos(s, t)</code>	<code>lastindexof(s,t)</code>
	<code>substr_replace(s, t,i,j)</code>	<code>substring(s,0,i).t.substring(s,j, s)</code>
	<code>strip_tags(s)</code>	<code>replaceall(s,("<a>" "<p>" ...), "")</code>
	<code>mysql_real_escape_string(s)</code>	<code>...replaceall(s, replaceall(s, "\\ ", "\\ \\ ") , "' ", "\\' ") ...</code>

Model Counting String Constraint Solver



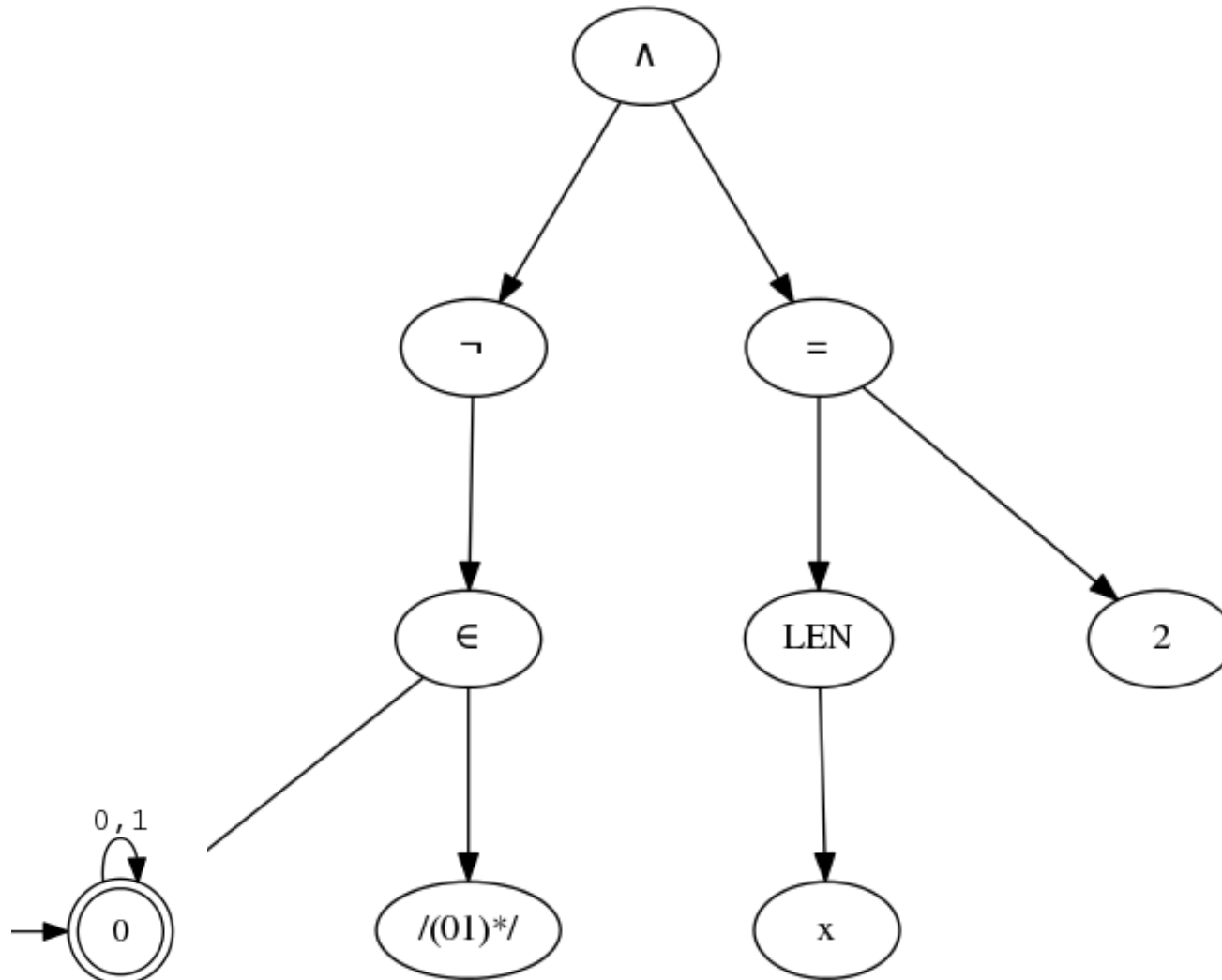
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



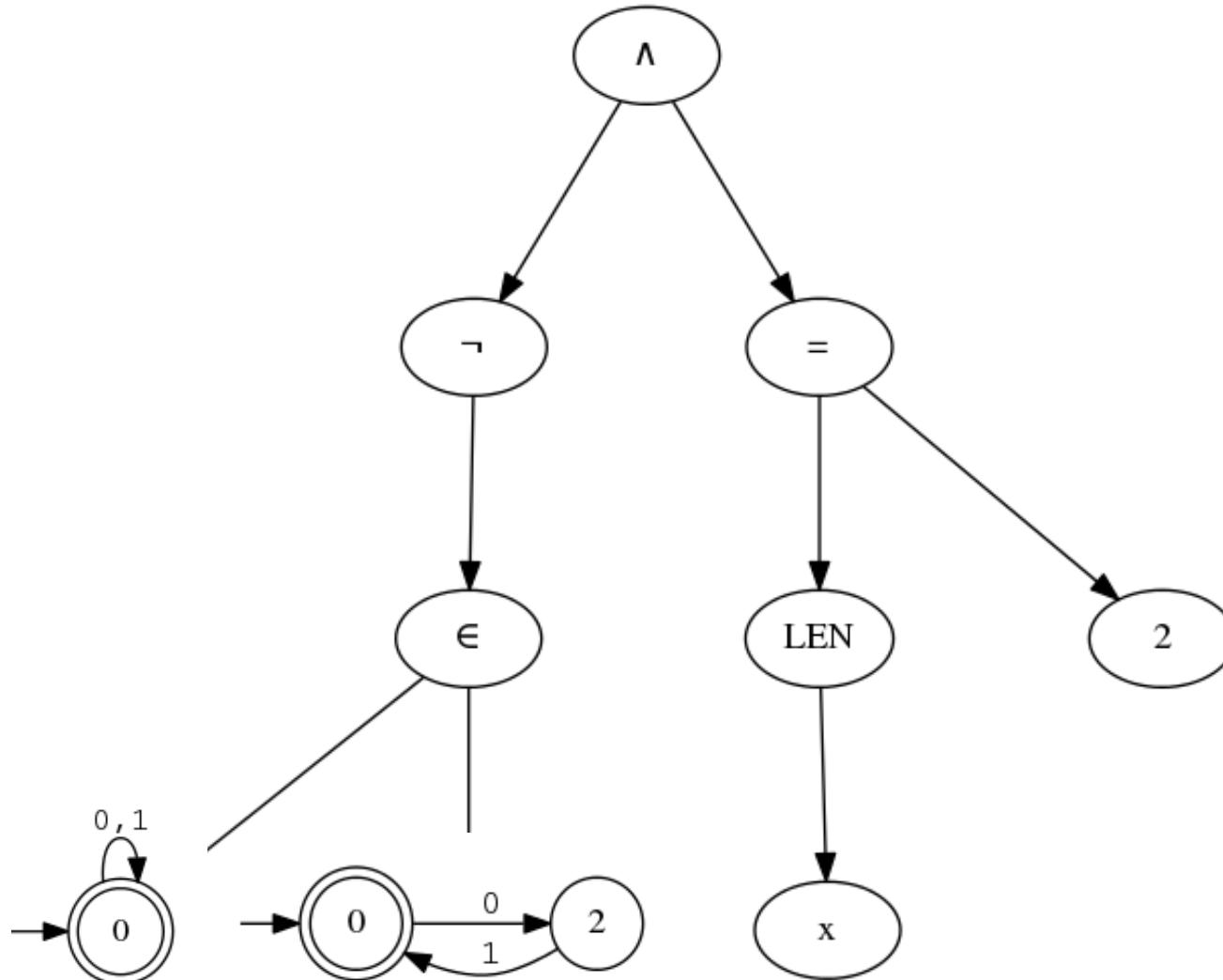
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



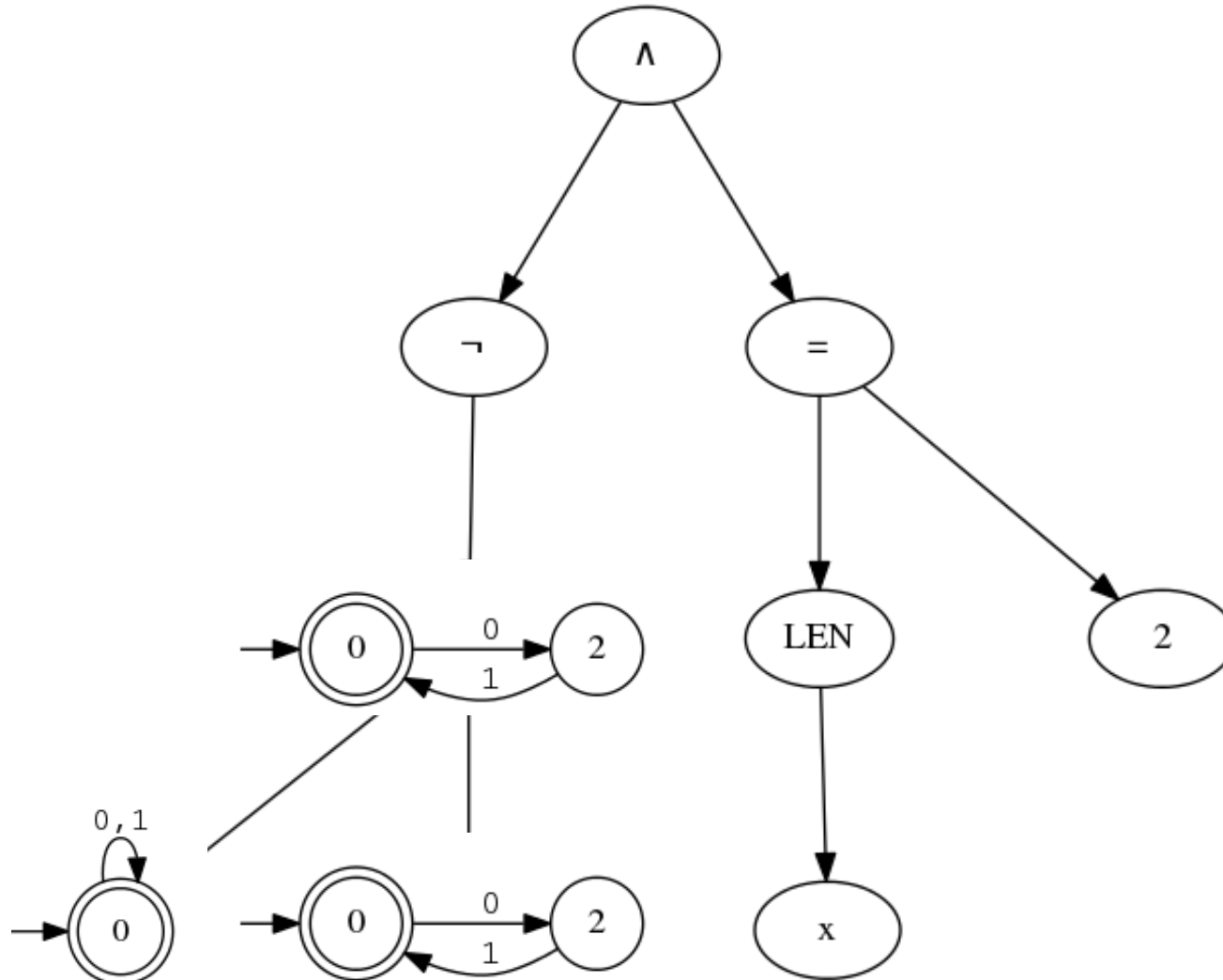
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



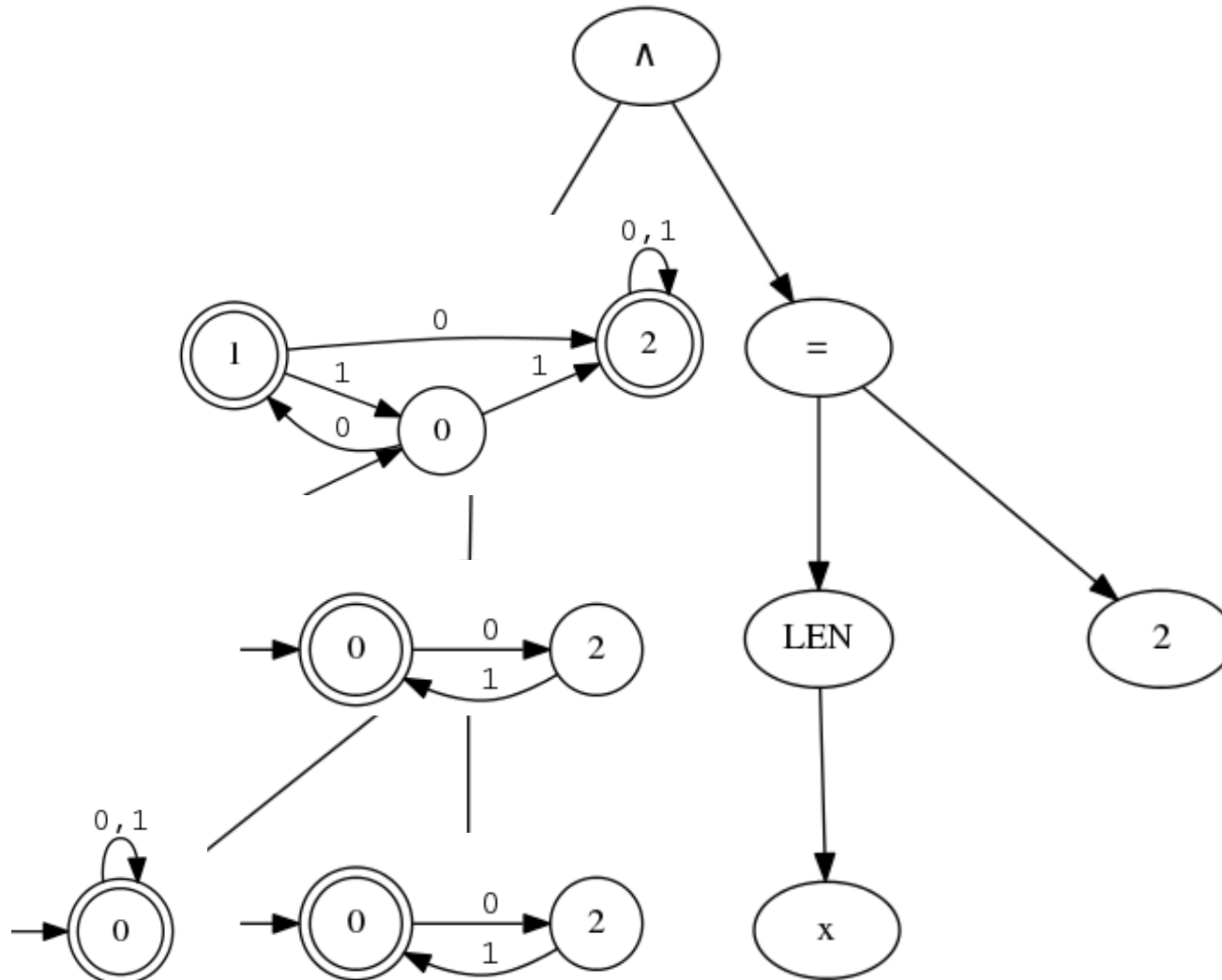
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



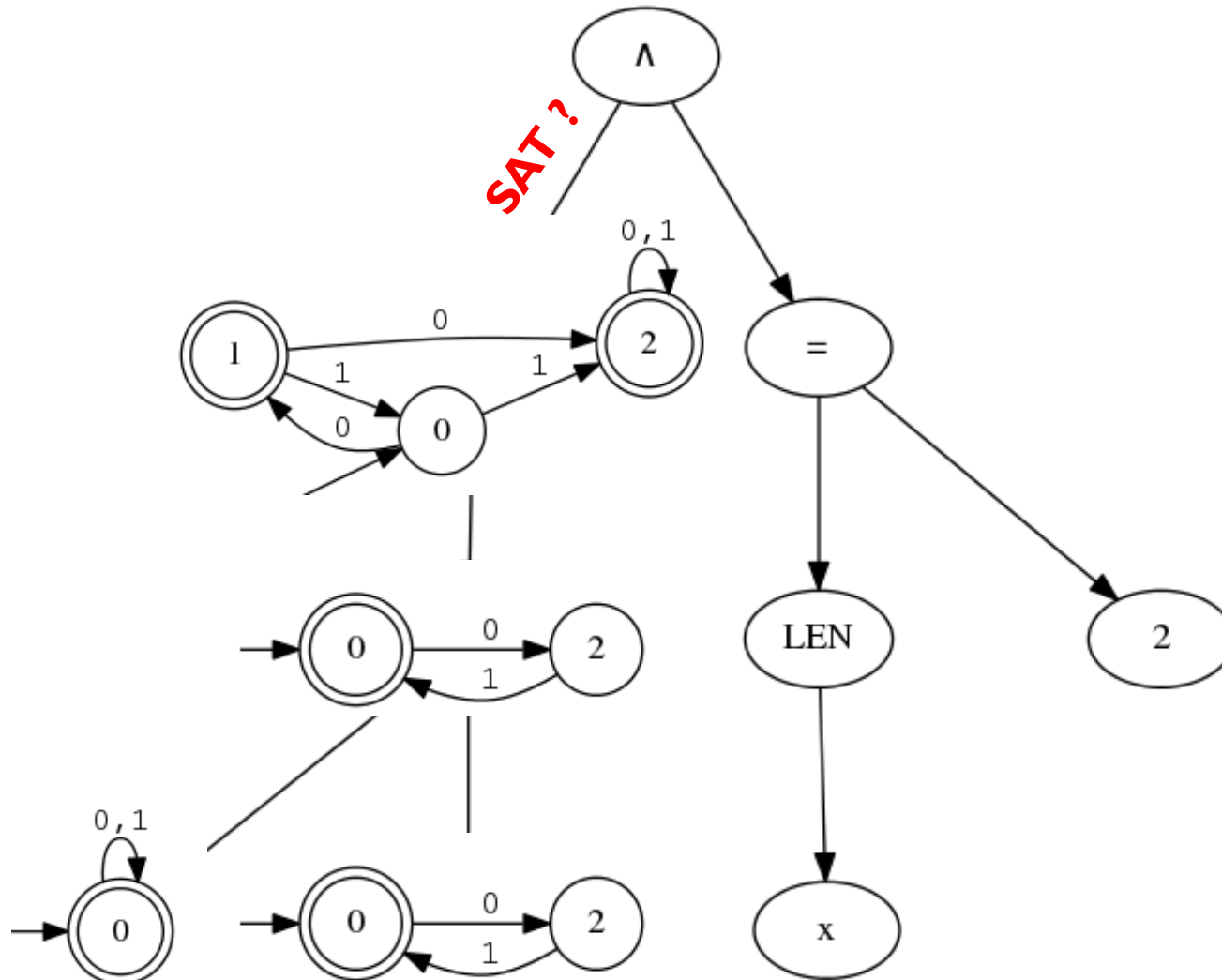
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



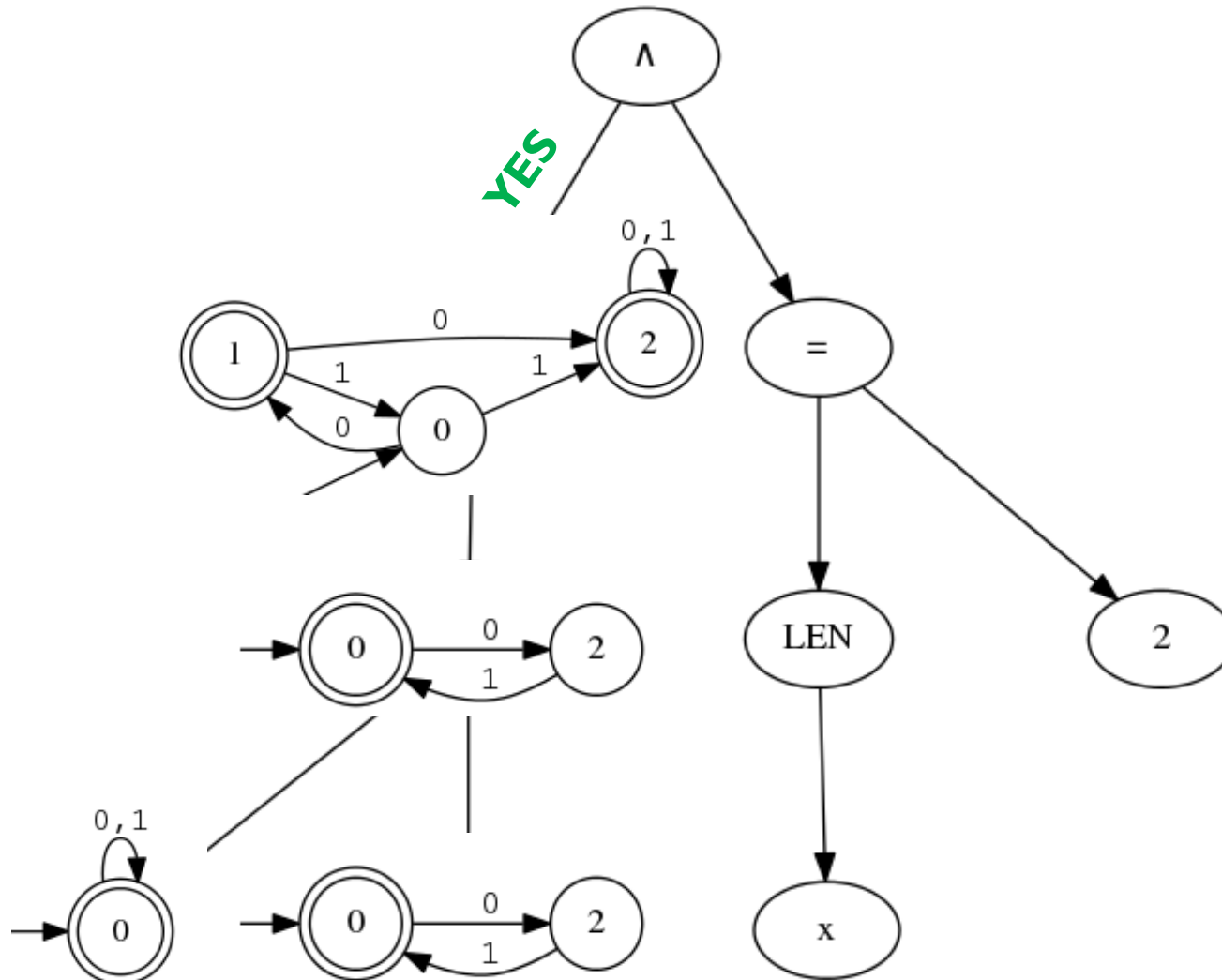
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



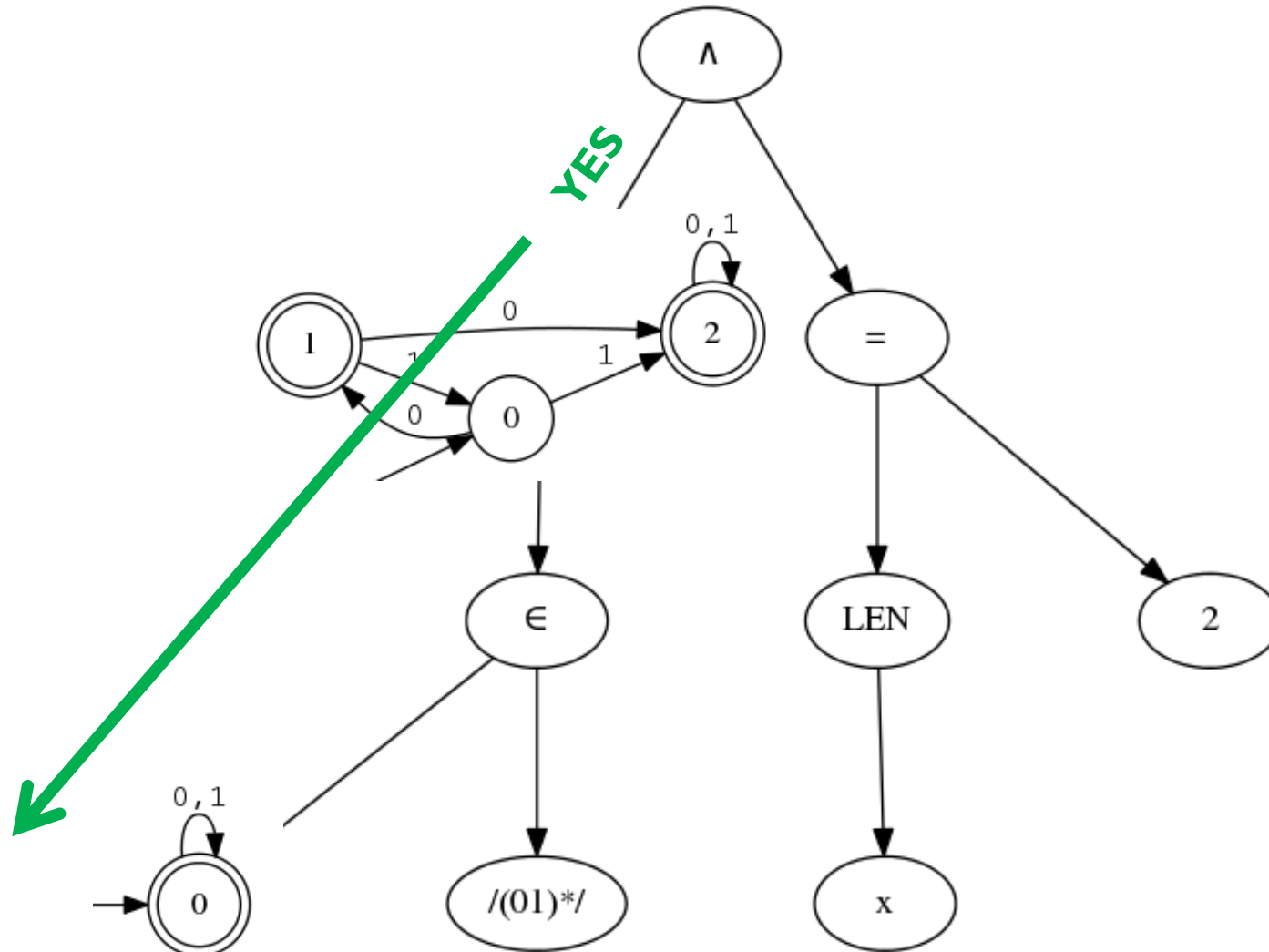
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



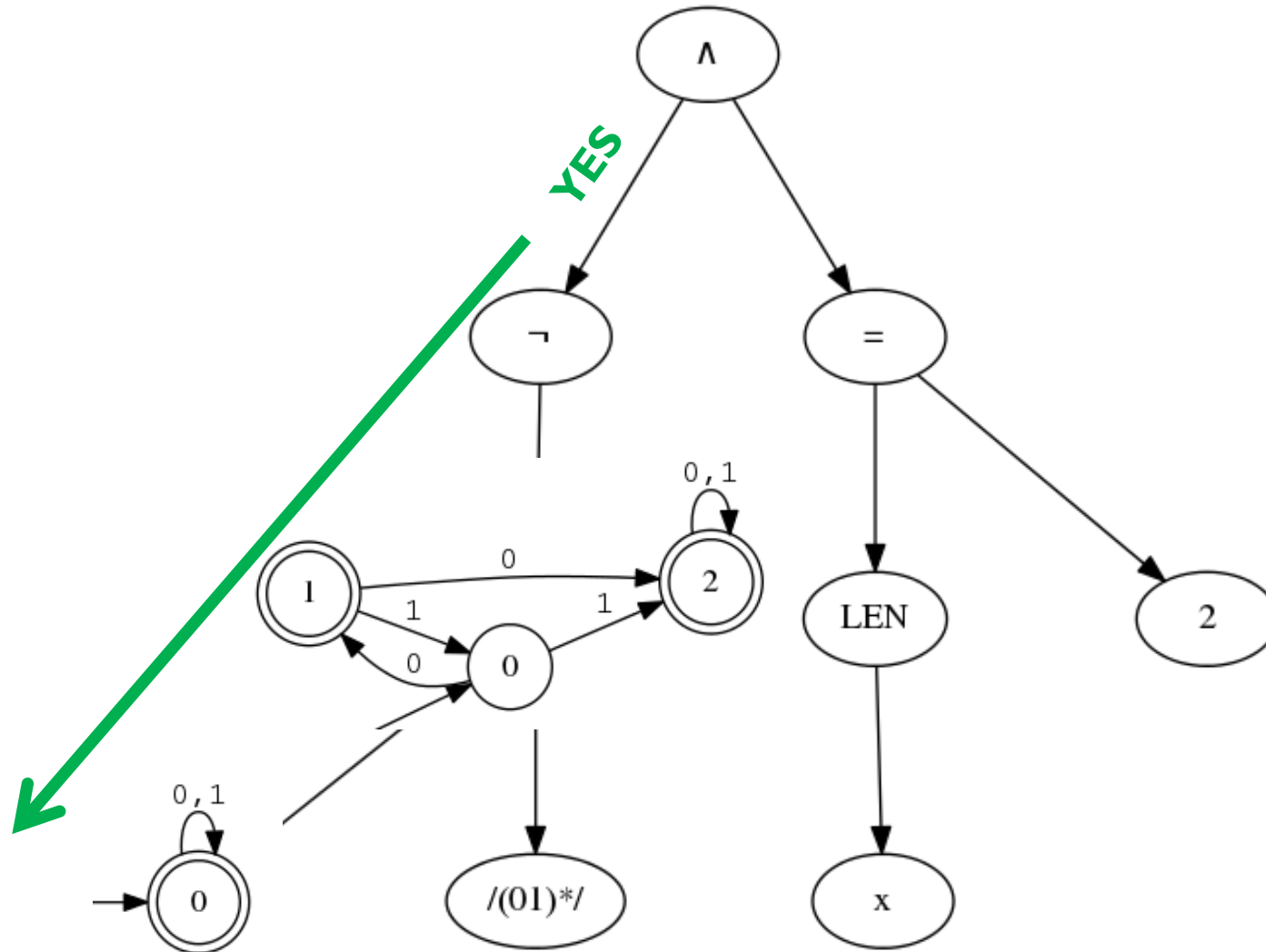
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



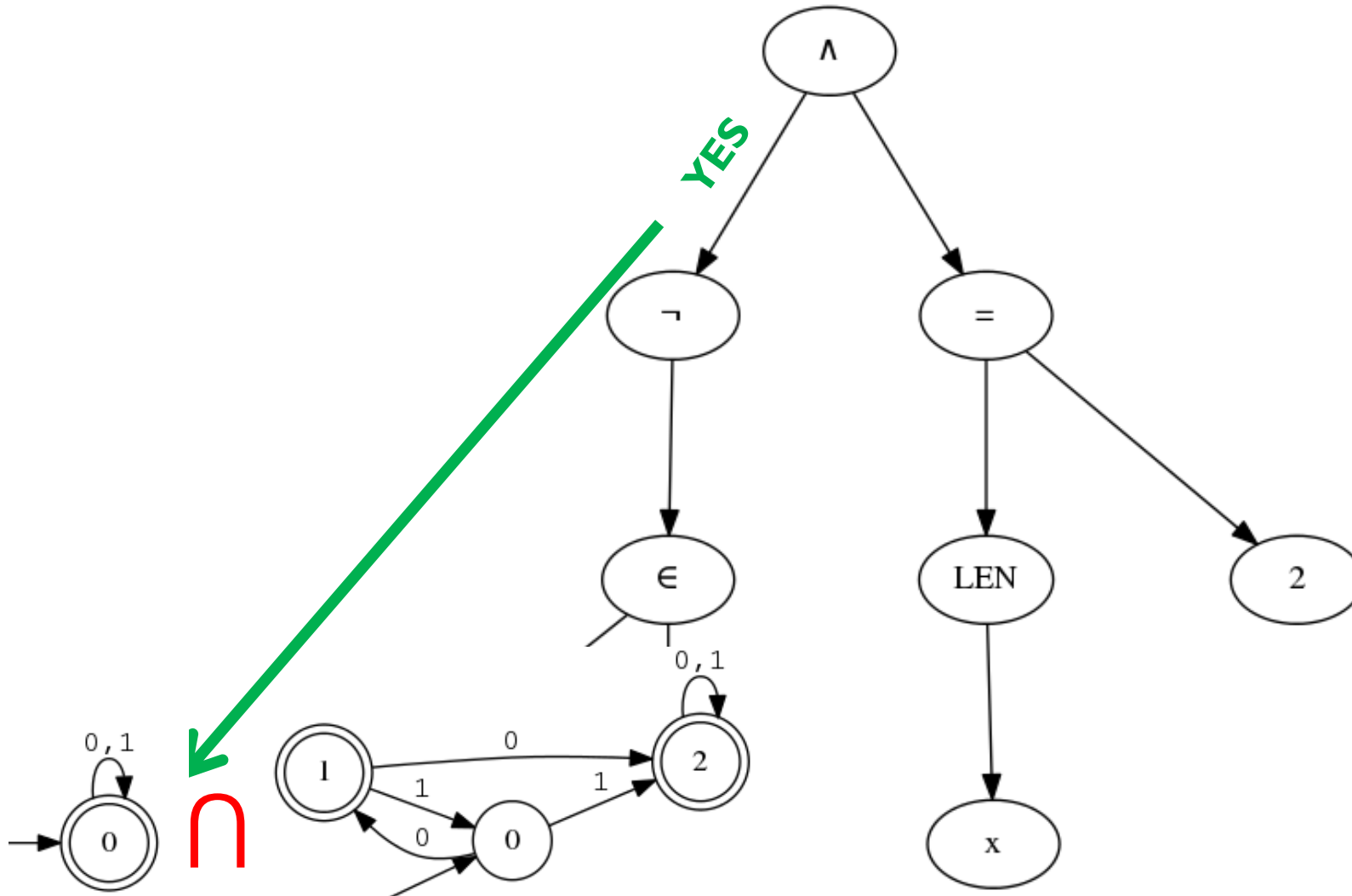
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



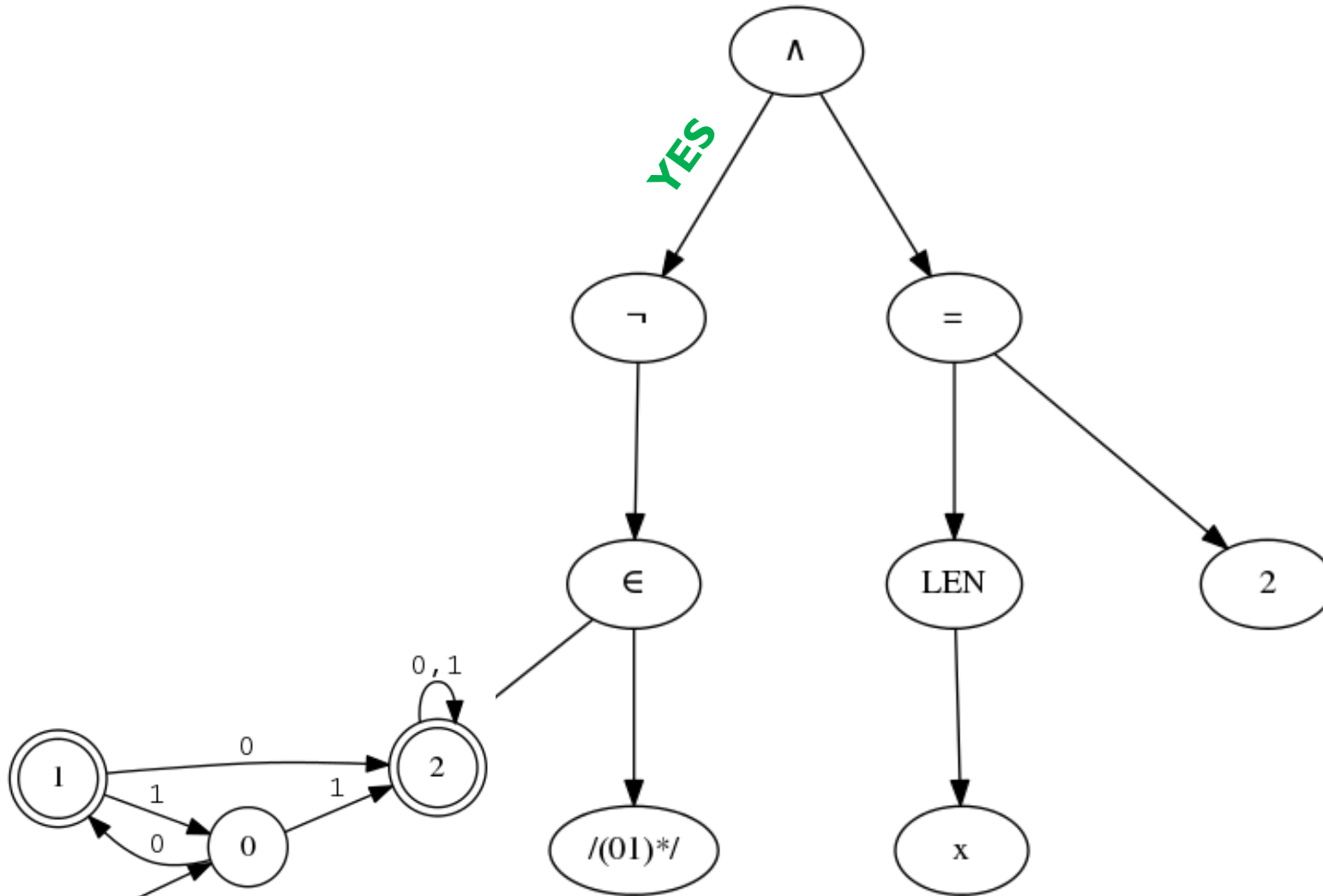
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



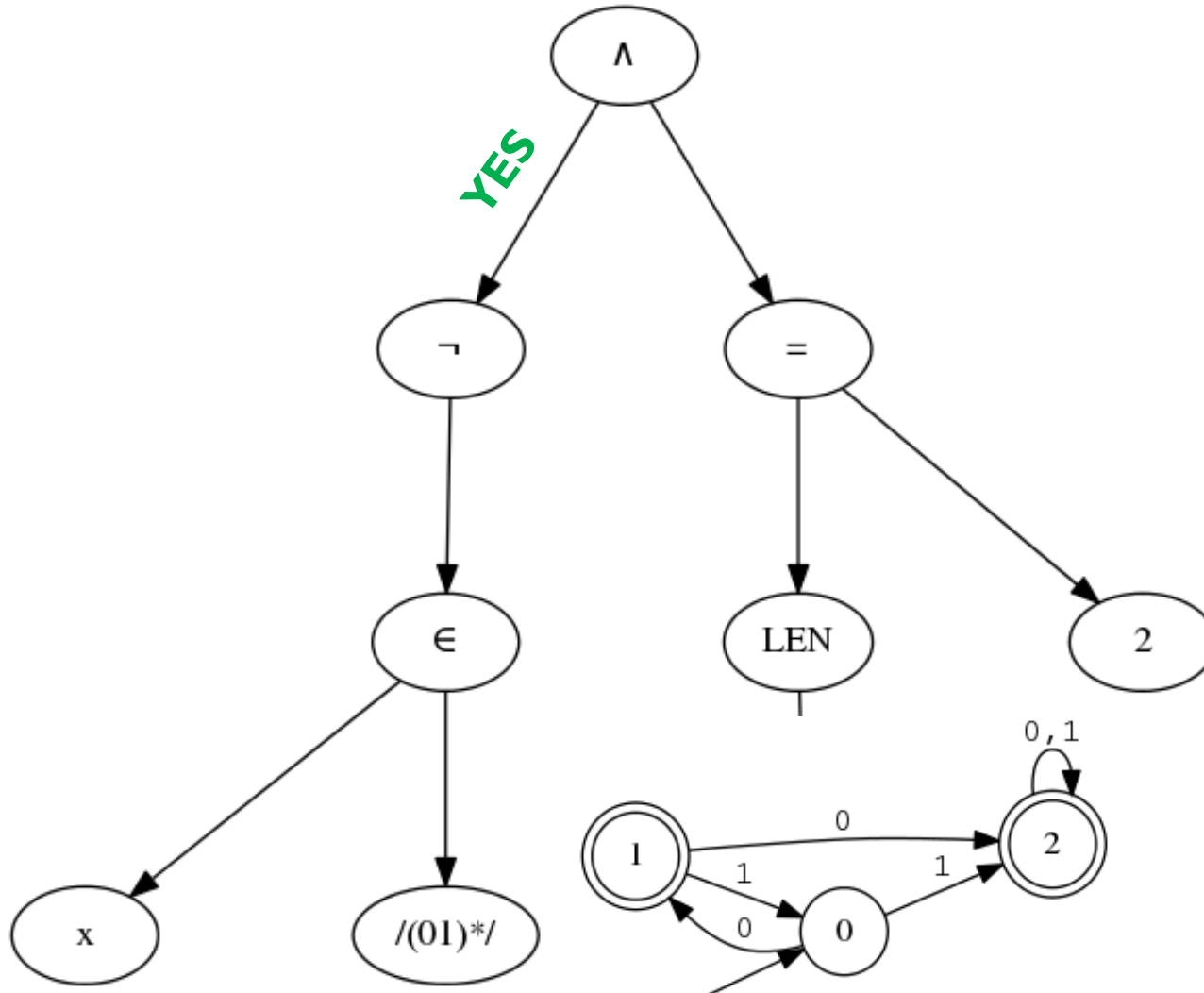
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



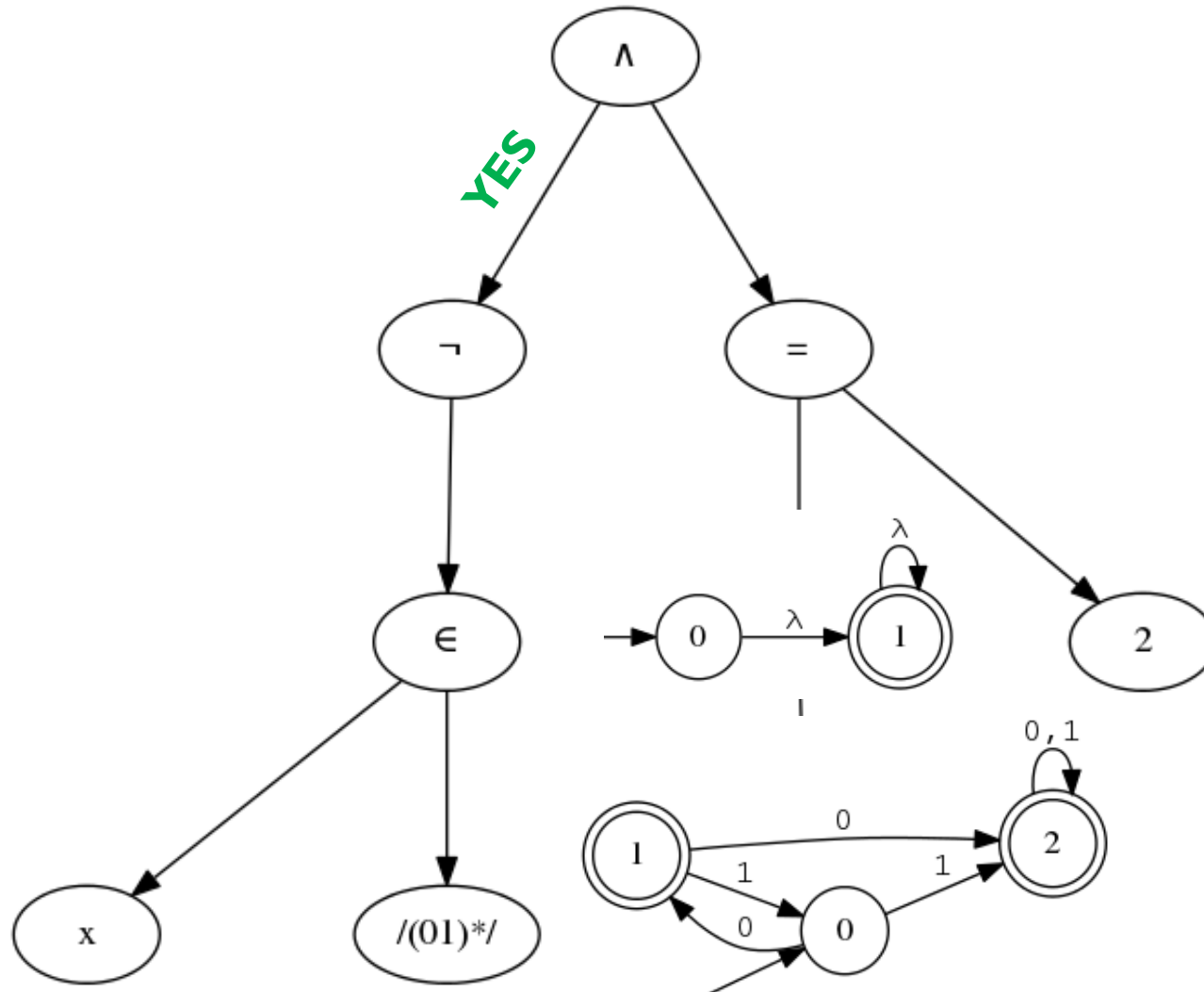
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



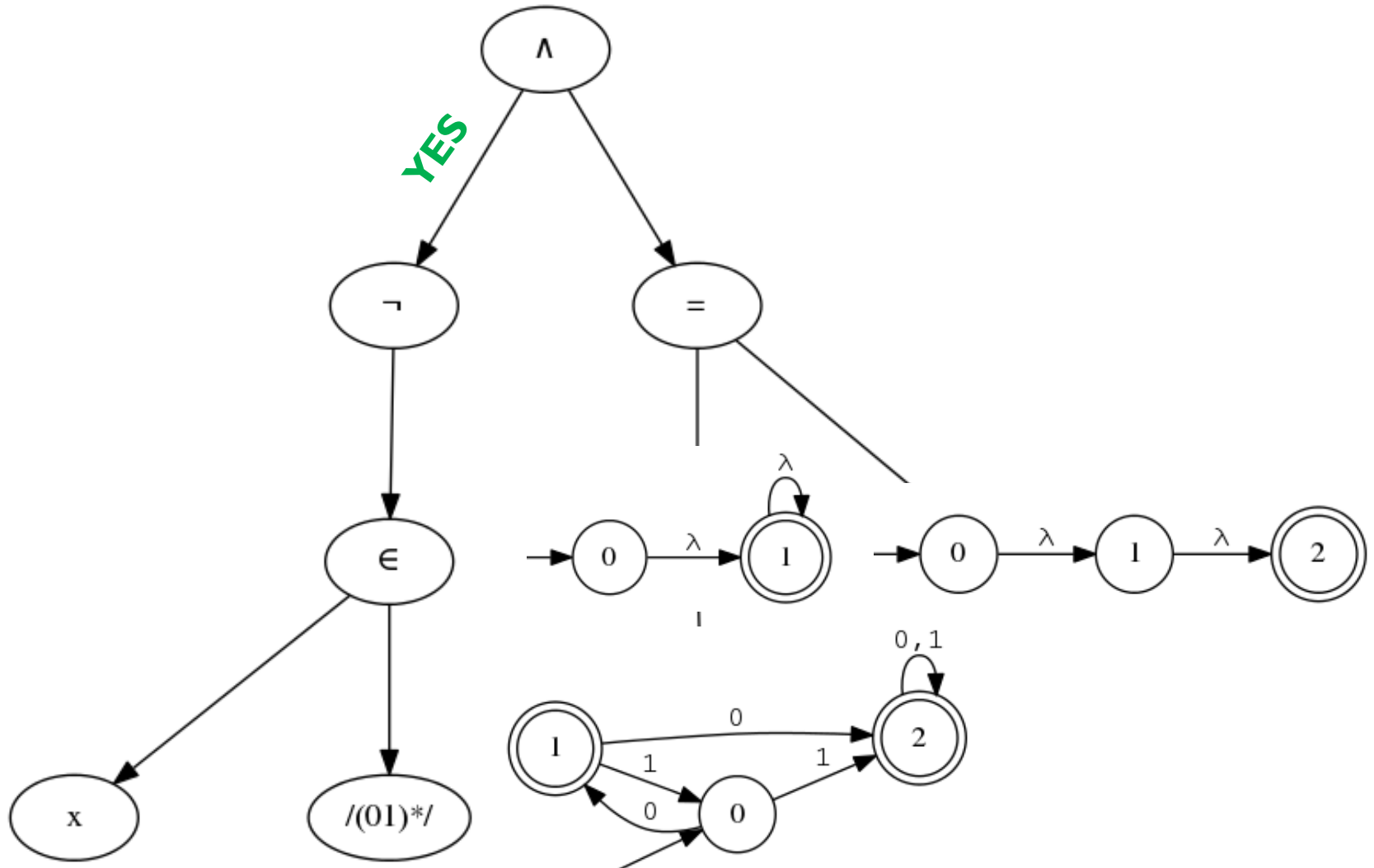
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



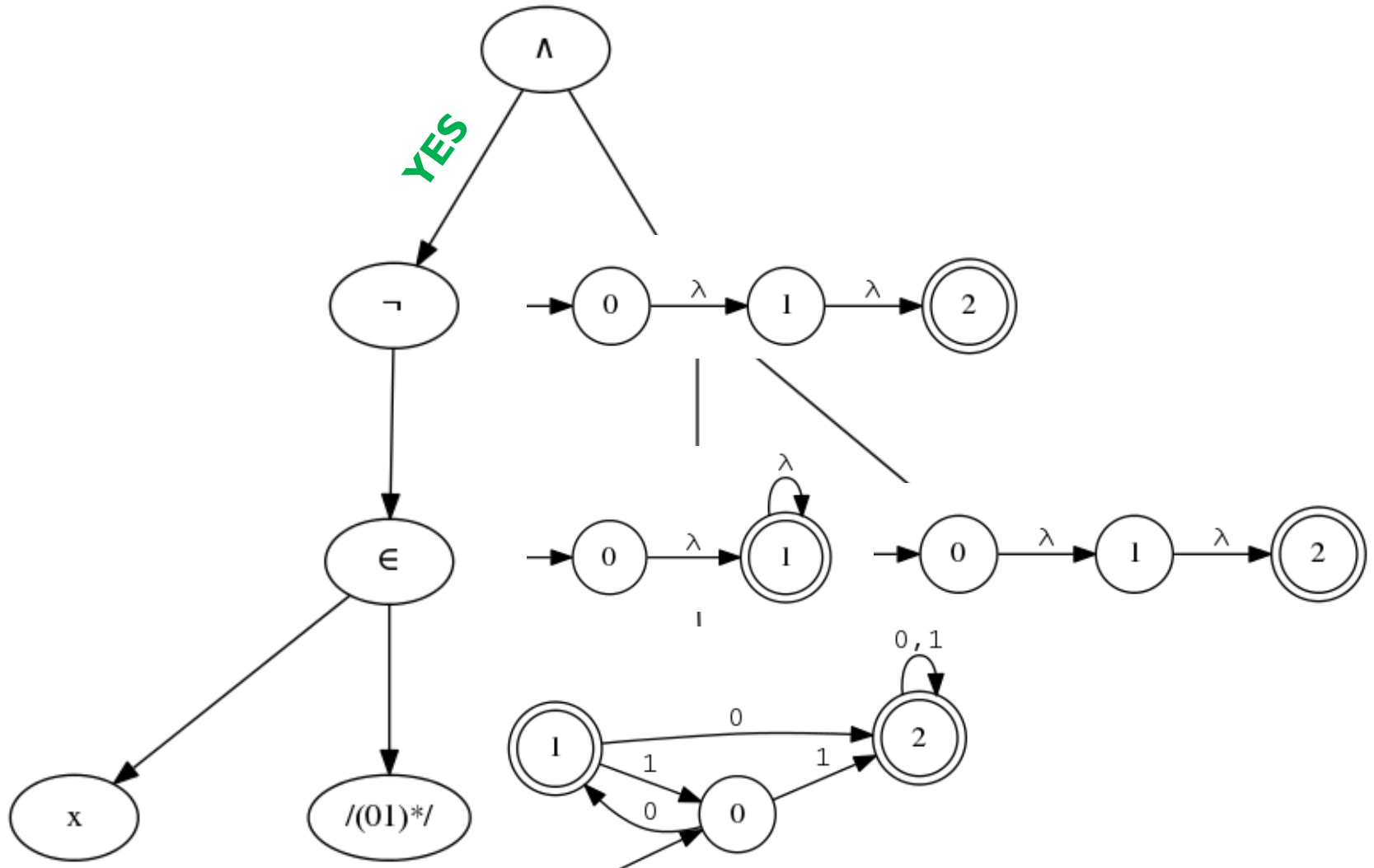
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



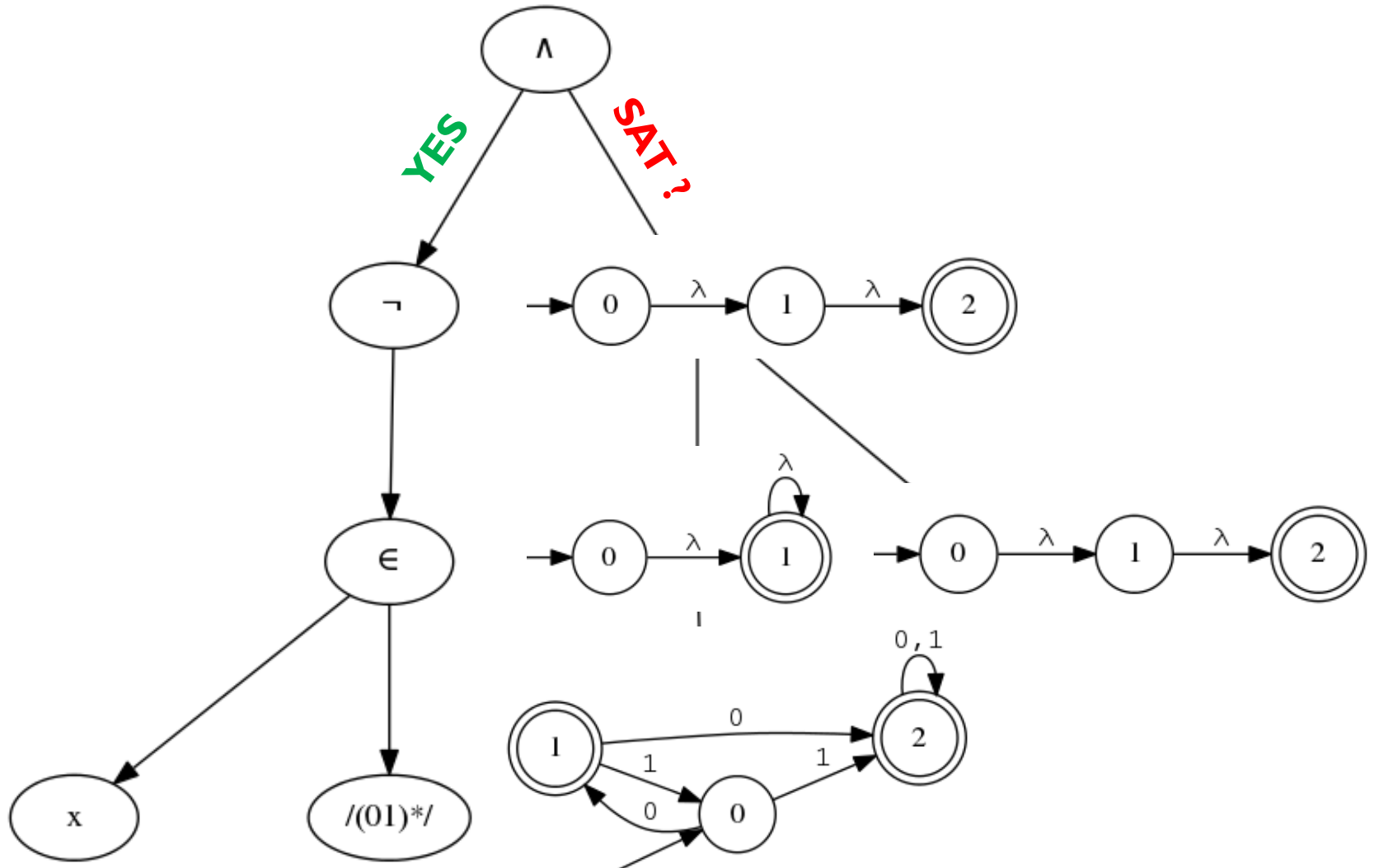
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



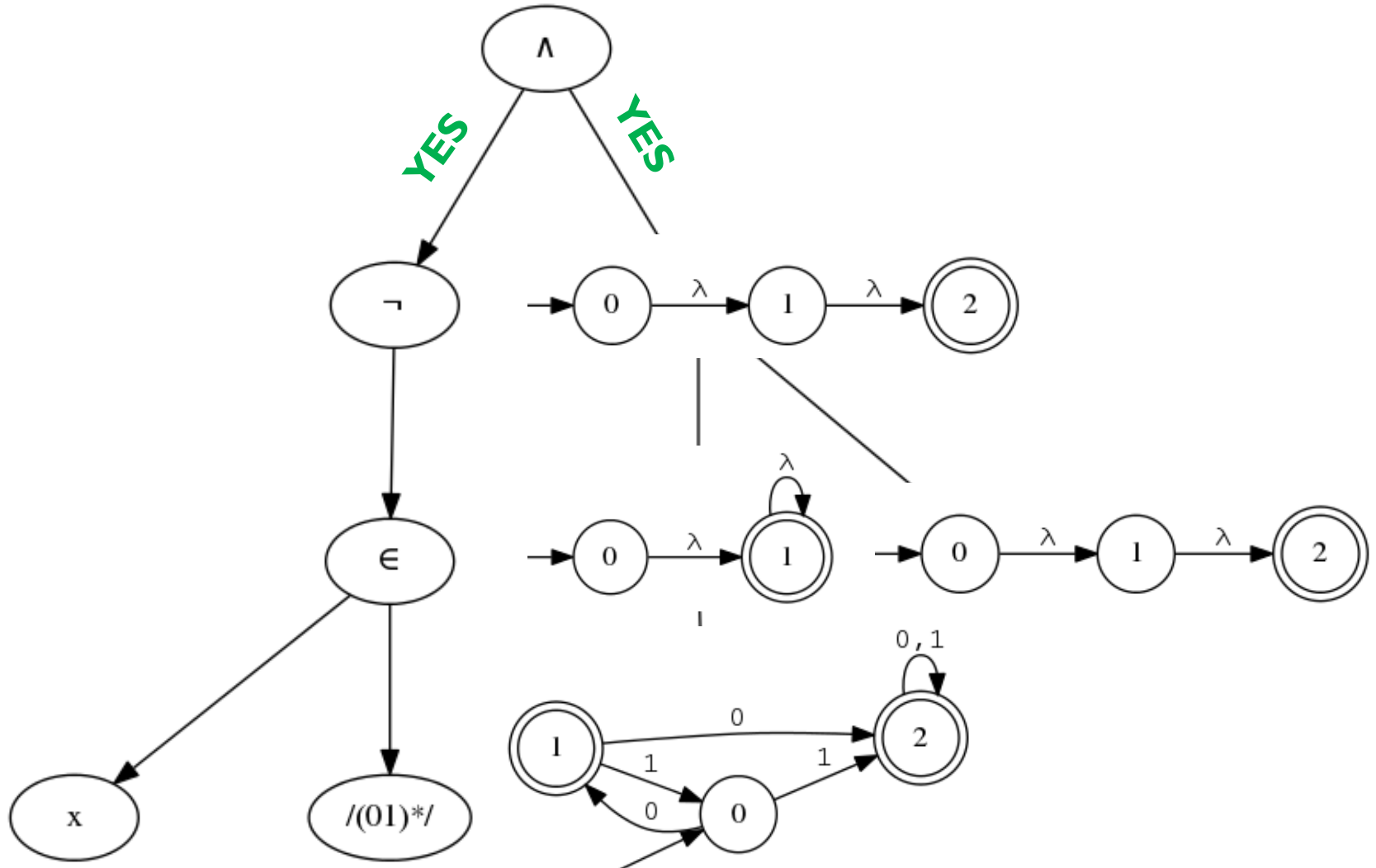
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



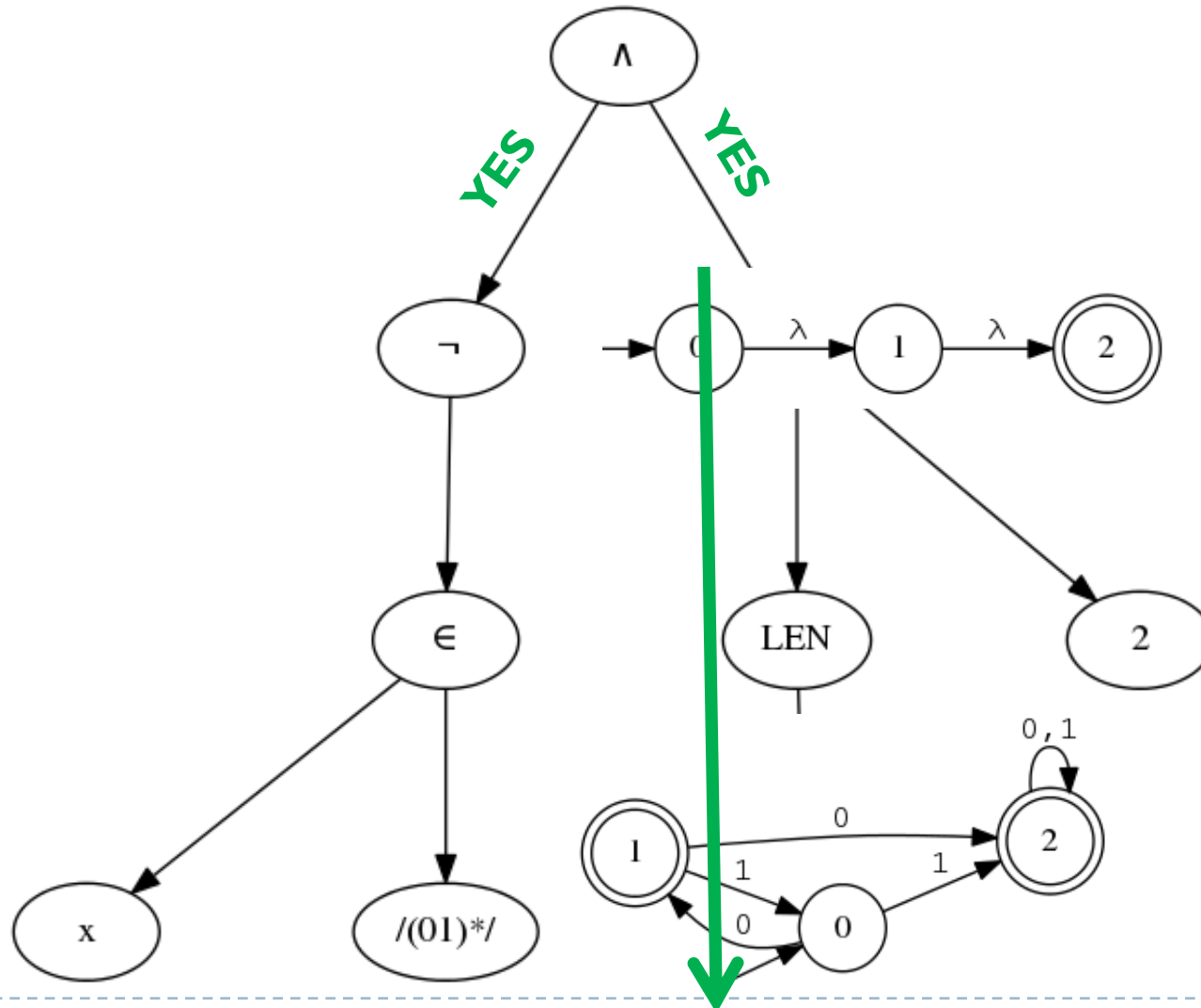
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



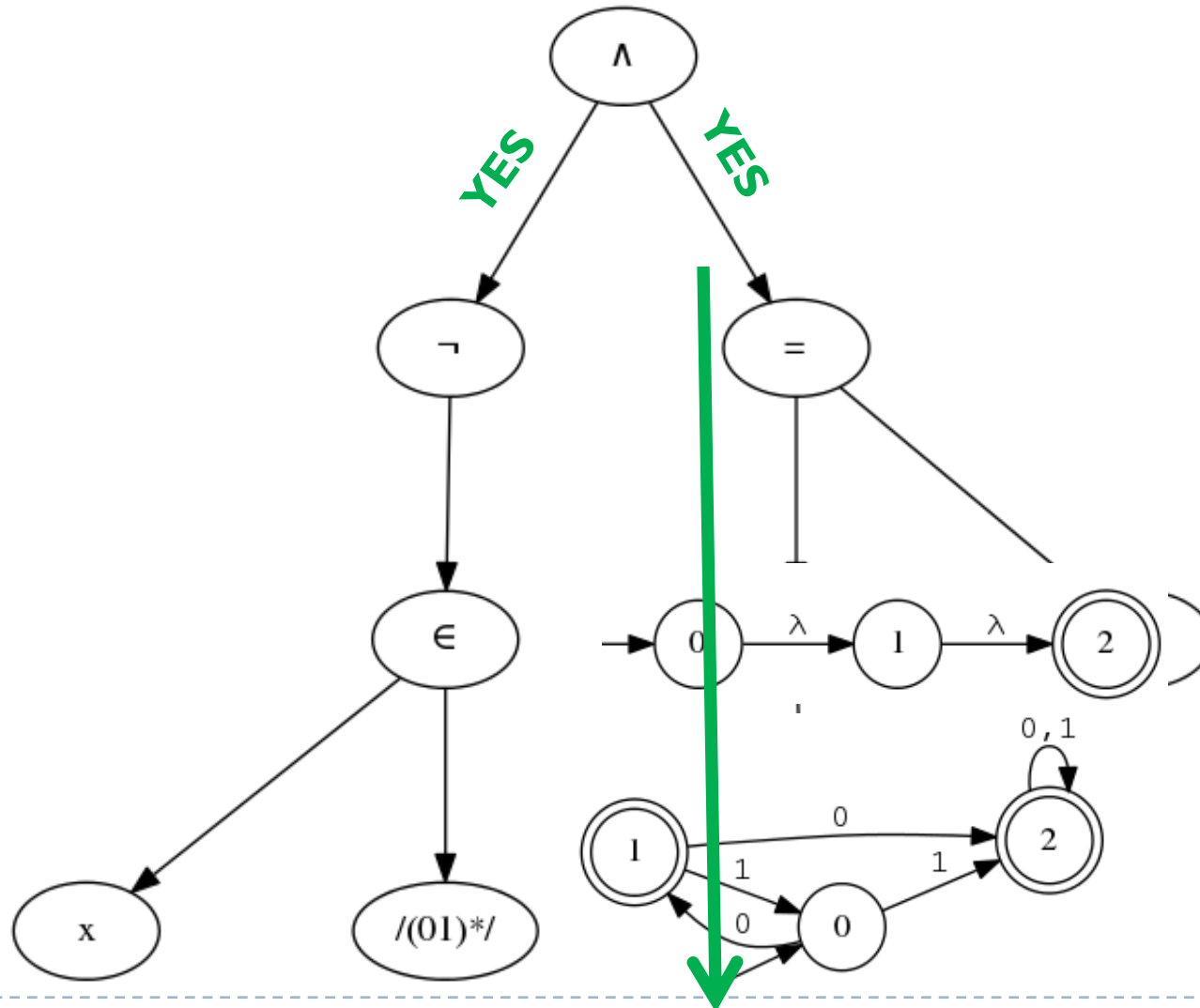
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



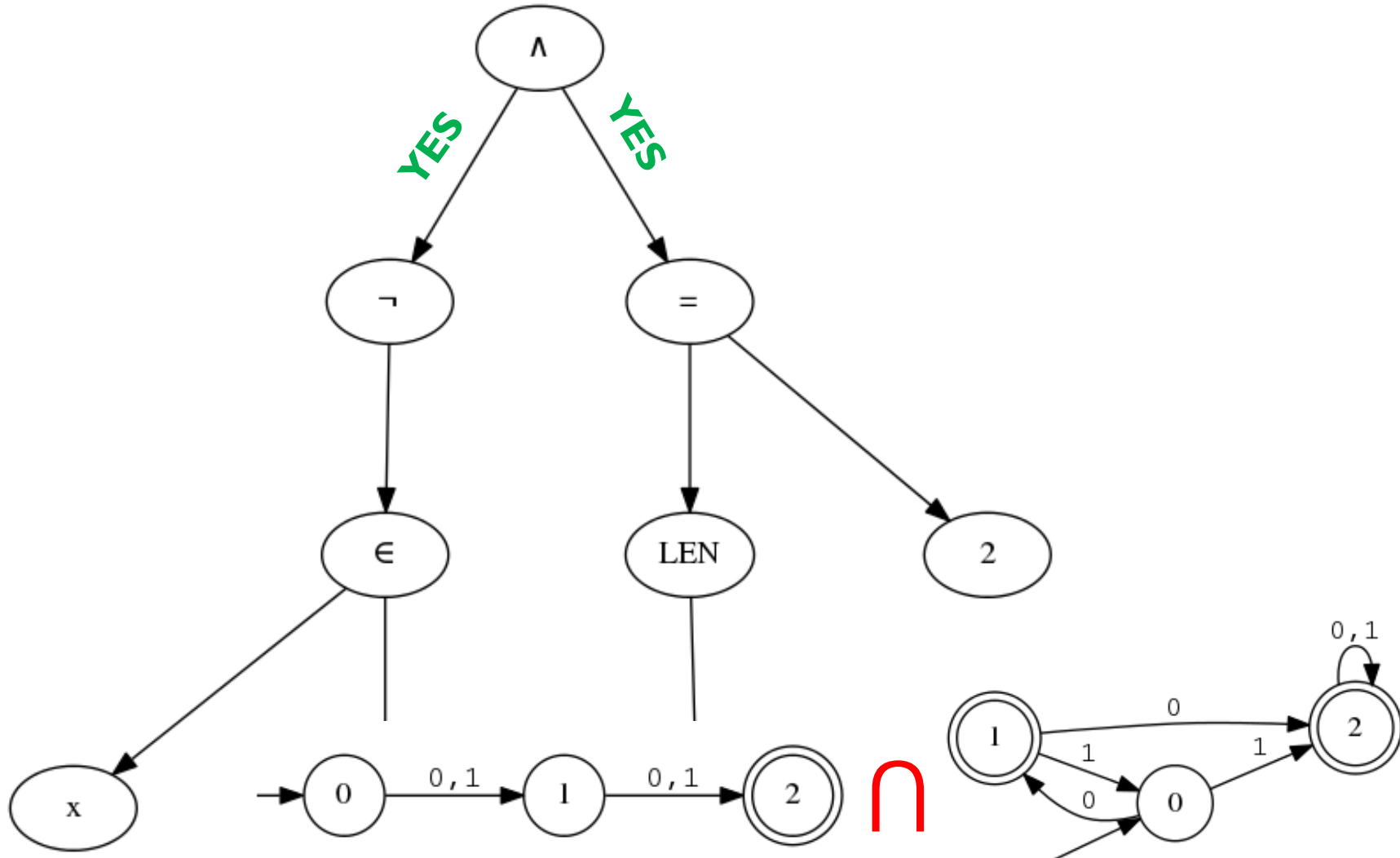
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



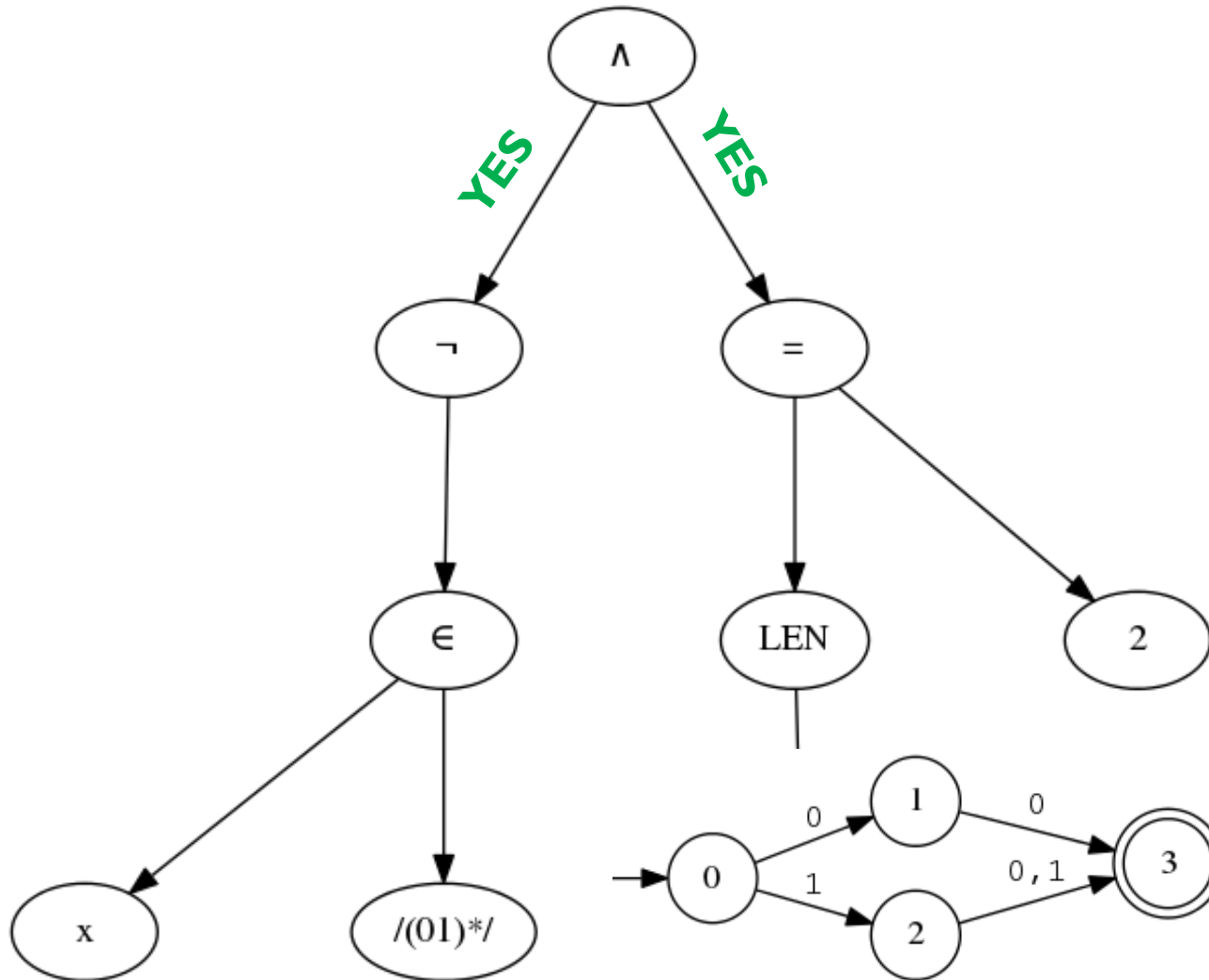
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



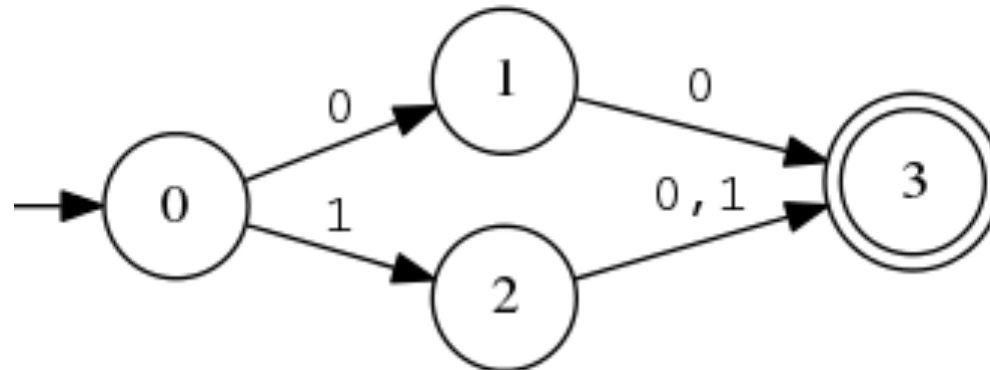
String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



String Automata Construction

$$C \equiv \neg(x \in (01)^*) \wedge \text{LEN}(x) = 2$$



00, 10, 11



Integer Constraints

$C \rightarrow bterm$

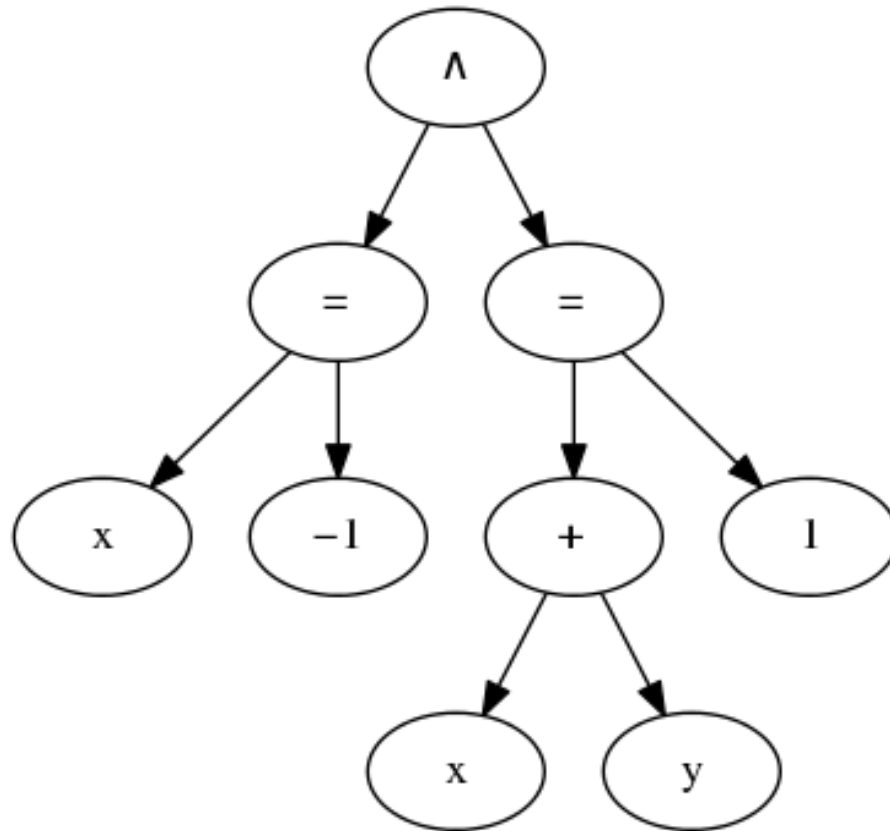
$bterm \rightarrow v \mid true \mid false$
| $\neg bterm \mid bterm \wedge bterm \mid bterm \vee bterm \mid (bterm)$
| $stern = stern$
| $match(stern, stern)$
| $contains(stern, stern)$
| $begins(stern, stern)$
| $ends(stern, stern)$
| $iterm = iterm \mid iterm < iterm \mid iterm > iterm$

$iterm \rightarrow v \mid n$
| $iterm + iterm \mid iterm - iterm \mid iterm \times n \mid (iterm)$
| $length(stern) \mid toint(stern)$
| $indexof(stern, stern)$
| $lastindexof(stern, stern)$

$stern \rightarrow v \mid \epsilon \mid s$
| $stern.stern \mid stern|stern \mid stern^* \mid (stern)$
| $charat(stern, iterm) \mid tostring(iterm)$
| $toupper(stern) \mid tolower(stern)$
| $substring(stern, iterm, iterm)$
| $replacefirst(stern, stern, stern)$
| $replacelast(stern, stern, stern)$
| $replaceall(stern, stern, stern)$

Integer Automata Construction

$$C \equiv x = -1 \wedge x + y = 1$$

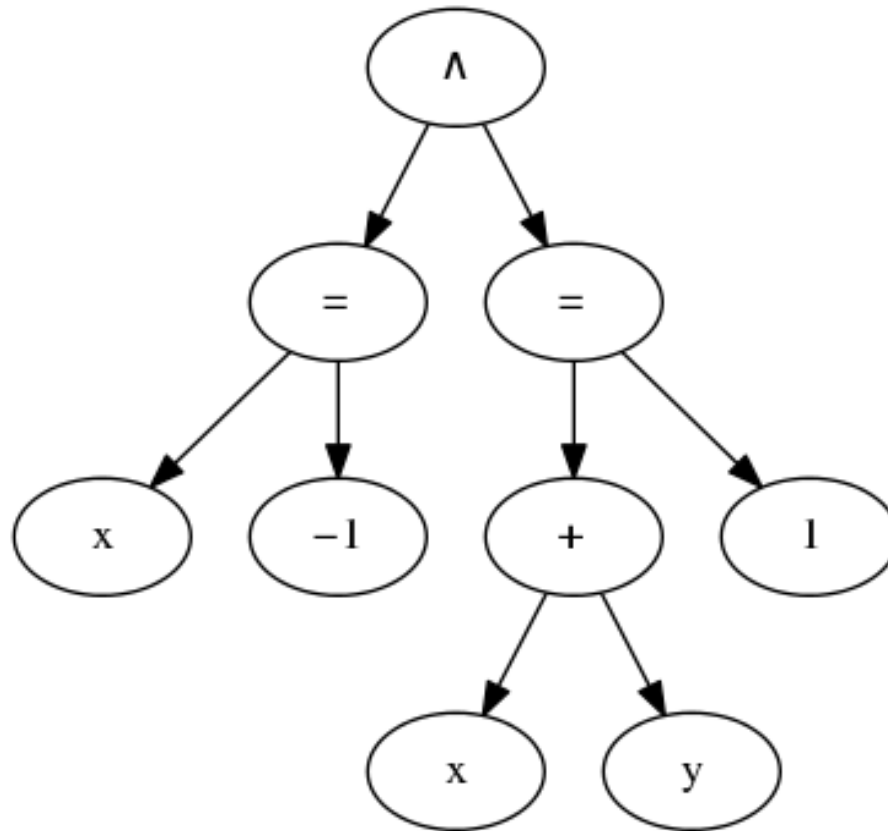


Integer Automata Construction

$$C \equiv x = -1 \wedge x + y = 1$$

$$C_1 \equiv x + 0 * y + 1 = 0 \Rightarrow [1 \ 0 \ 1]$$

$$C_2 \equiv x + y - 1 = 0 \Rightarrow [1 \ 1 \ -1]$$

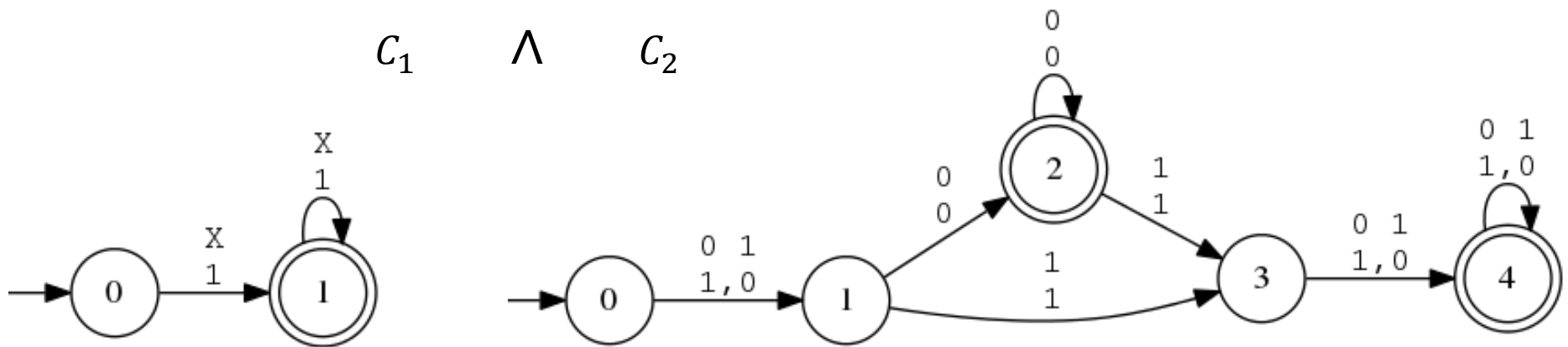


Integer Automata Construction

$$C \equiv x = -1 \wedge x + y = 1$$

$$C_1 \equiv x + 0 * y + 1 = 0 \Rightarrow [1 \ 0 \ 1]$$

$$C_2 \equiv x + y - 1 = 0 \Rightarrow [1 \ 1 \ -1]$$

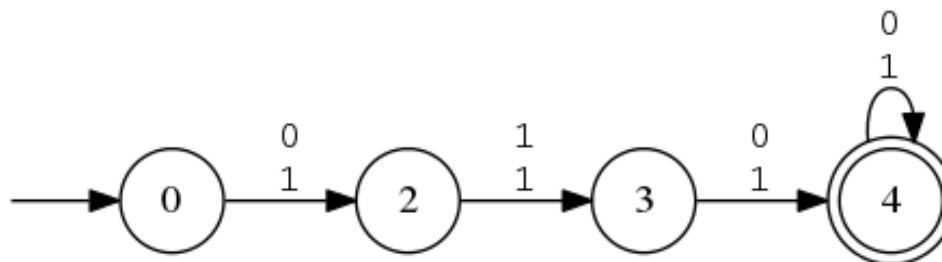


- ▶ Using automata construction techniques described in:

C. Bartzis and Tefvik Bultan. Efficient symbolic representations for arithmetic constraints in verification. *Int. J. Found. Comput. Sci.*, 2003

Integer Automata Construction

$$C \equiv x = -1 \wedge x + y = 1$$



$$(111, 010) = (-1, 2)$$

- ▶ Conjunction and disjunction is handled by automata product, negation is handled by automata complement

Mixing String – LIA Constraints

$C \rightarrow bterm$

$bterm \rightarrow v \mid true \mid false$
| $\neg bterm \mid bterm \wedge bterm \mid bterm \vee bterm \mid (bterm)$
| $stern = stern$
| $match(stern, stern)$
| $contains(stern, stern)$
| $begins(stern, stern)$
| $ends(stern, stern)$
| $iterm = iterm \mid iterm < iterm \mid iterm > iterm$

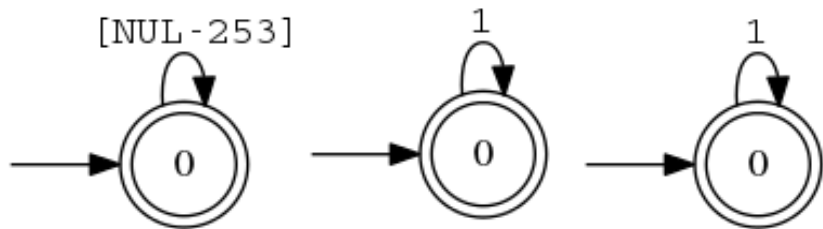
$iterm \rightarrow v \mid n$
| $iterm + iterm \mid iterm - iterm \mid iterm \times n \mid (iterm)$
| $length(stern) \mid toint(stern)$
| $indexof(stern, stern)$
| $lastindexof(stern, stern)$

$stern \rightarrow v \mid \epsilon \mid s$
| $stern.stern \mid stern|stern \mid stern^* \mid (stern)$
| $charat(stern, iterm) \mid tostring(iterm)$
| $toupper(stern) \mid tolower(stern)$
| $substring(stern, iterm, iterm)$
| $replacefirst(stern, stern, stern)$
| $replacelast(stern, stern, stern)$
| $replaceall(stern, stern, stern)$

Mixing String – LIA Constraints

$$\text{length}(s) = x \wedge x = 2y$$

- ▶ We can try to solve using string constraint solving
- ▶ Integer variables are represented with unary automaton
 - ▶ $A_s = L(\Sigma_s^*), A_x = L(\lambda^*), A_y = L(\lambda^*), A_{2y} = L(\lambda^*)$
 - ▶ $A_s \rightarrow A_{\text{length}}; L(\Sigma_s^*) \rightarrow L(\lambda^*)$
 - ▶ $L(\lambda^*) = L(\lambda^*) \wedge L(\lambda^*) = L(\lambda^*)$

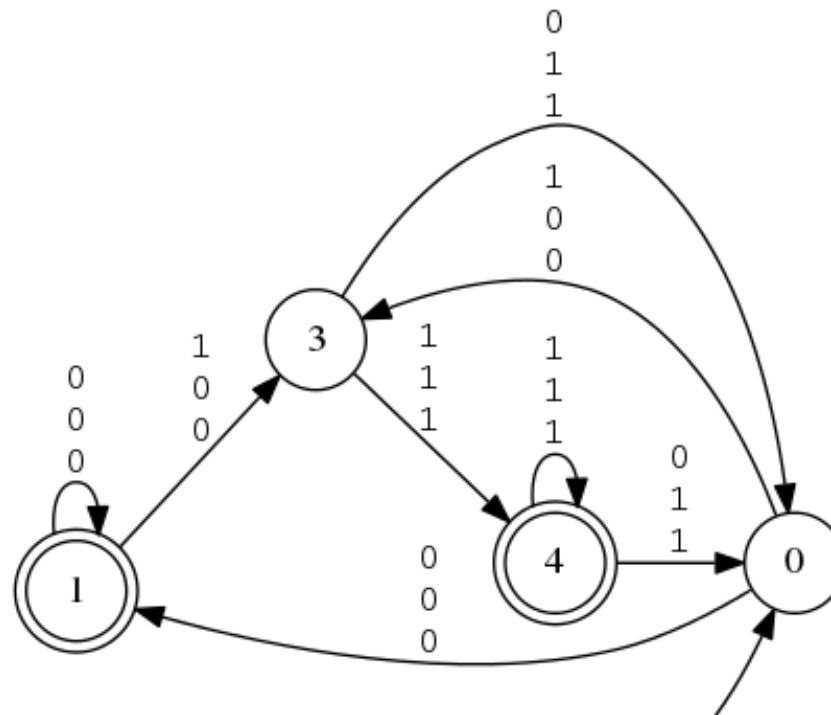


- ▶ $s = 'abc', x = 1, y = 2$

Mixing String – LIA Constraints

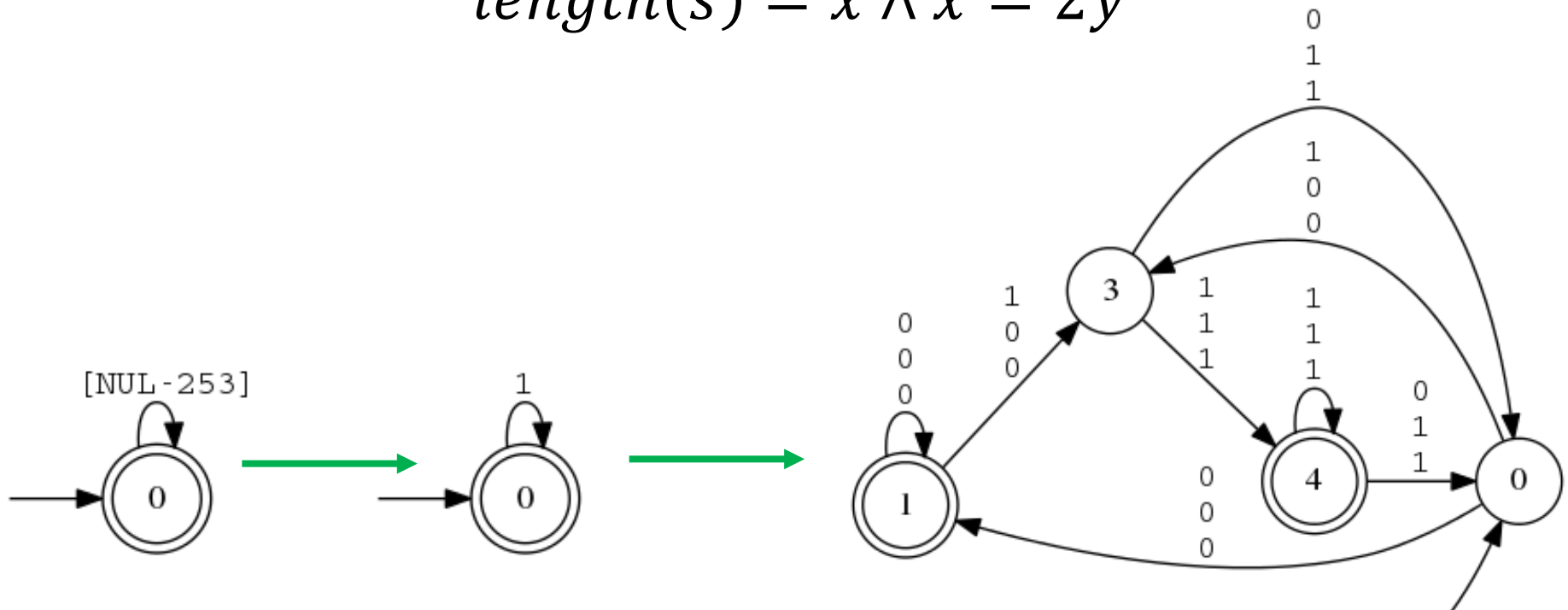
$$\text{length}(s) = x \wedge x = 2y$$

- ▶ We can do better by using multi-track binary integer automaton
- ▶ There will be three integer variables: $\text{length}(s), x, y$



Mixing String – LIA Constraints

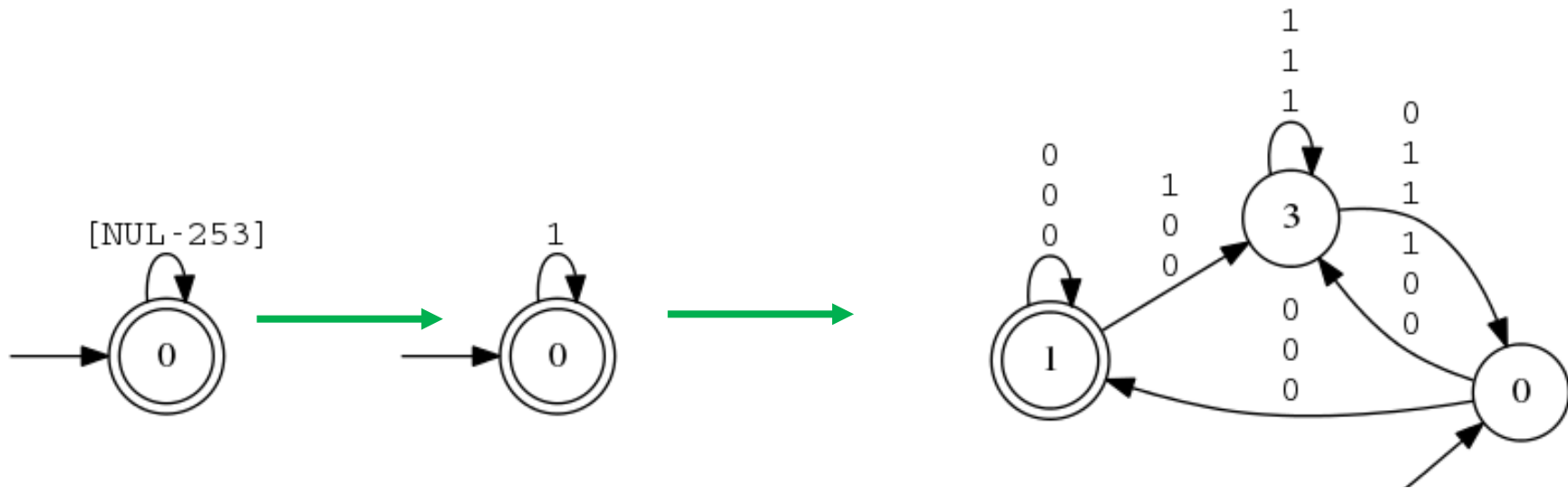
$$\text{length}(s) = x \wedge x = 2y$$



- ▶ Using automata techniques described in:
F. Yu, T. Bultan, O. H. Ibarra. Symbolic String Verification: Combining String Analysis and Size Analysis. *TACAS'09*

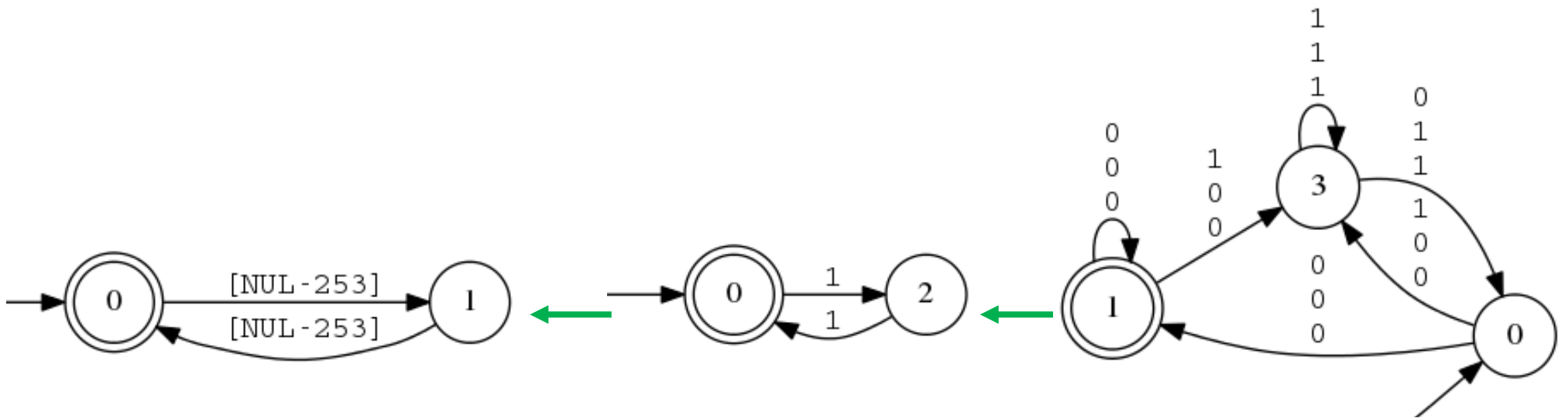
Mixing String – LIA Constraints

$$\text{length}(s) = x \wedge x = 2y$$



Mixing String – LIA Constraints

$$\text{length}(s) = x \wedge x = 2y$$



Relational String Constraints

$C \rightarrow bterm$

$bterm \rightarrow v \mid true \mid false$
| $\neg bterm \mid bterm \wedge bterm \mid bterm \vee bterm \mid (bterm)$
| $stern = stern$
| $match(stern, stern)$
| $contains(stern, stern)$
| $begins(stern, stern)$
| $ends(stern, stern)$
| $iterm = iterm \mid iterm < iterm \mid iterm > iterm$

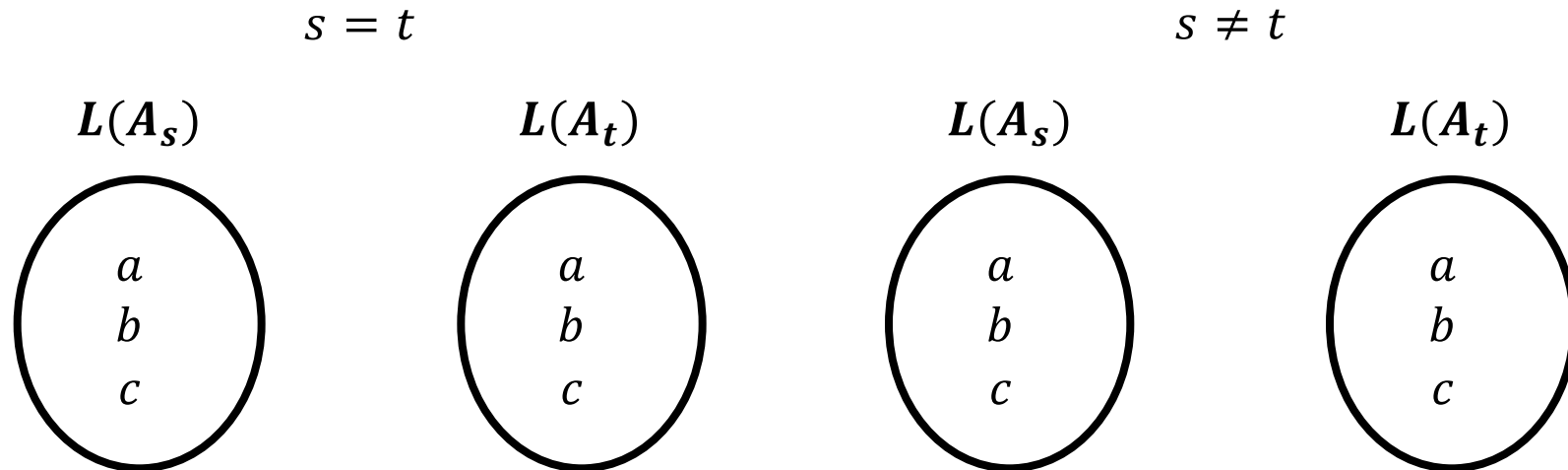
$iterm \rightarrow v \mid n$
| $iterm + iterm \mid iterm - iterm \mid iterm \times n \mid (iterm)$
| $length(stern) \mid toint(stern)$
| $indexof(stern, stern)$
| $lastindexof(stern, stern)$

$stern \rightarrow v \mid \epsilon \mid s$
| $stern.stern \mid stern|stern \mid stern^* \mid (stern)$
| $charat(stern, iterm) \mid tostring(iterm)$
| $toupper(stern) \mid tolower(stern)$
| $substring(stern, iterm, iterm)$
| $replacefirst(stern, stern, stern)$
| $replacelast(stern, stern, stern)$
| $replaceall(stern, stern, stern)$

Relational String Constraint Example

$$s = t \wedge s \neq t$$

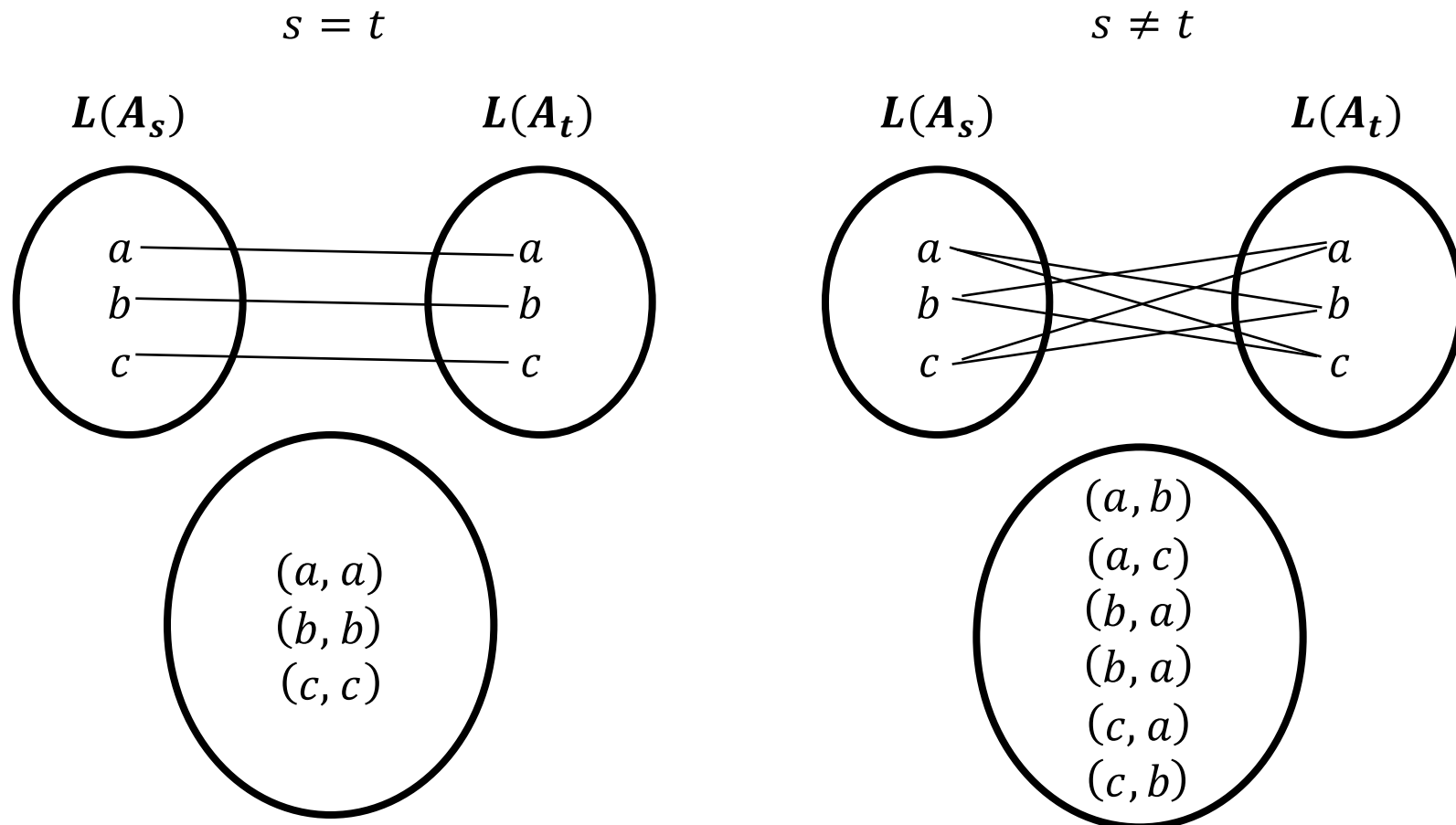
- ▶ Single track string automata results in SAT (w/ optimizations)
 - ▶ Assume $\Sigma = \{a, b, c\}$



Relational String Constraint Example

$$s = t \wedge s \neq t$$

- ▶ Multi-track string automata can keep relations

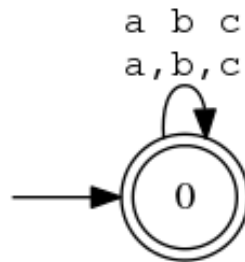


Relational String Constraint Example

$$s = t \wedge s \neq t$$

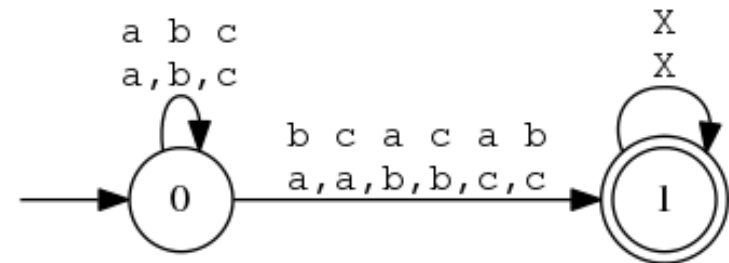
- ▶ Each track represents values of one string variable

$$s = t$$



\wedge

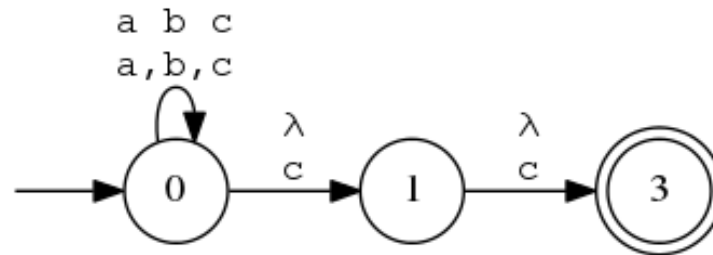
$$s \neq t$$



Word Equations

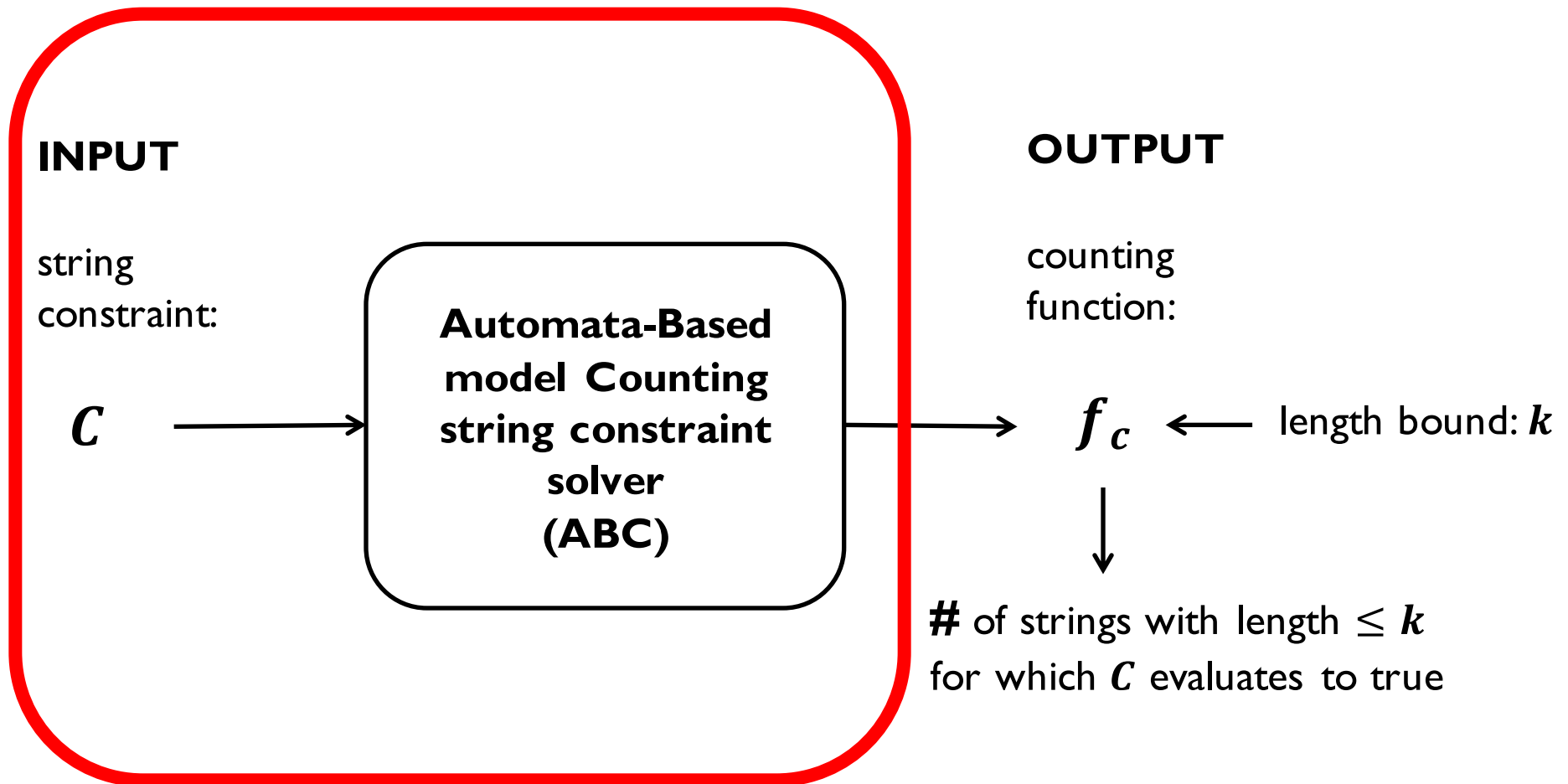
- ▶ Let X (the first track), Y (the second track), be two string variables
- ▶ λ is the padding symbol

$$X = Y.c c$$

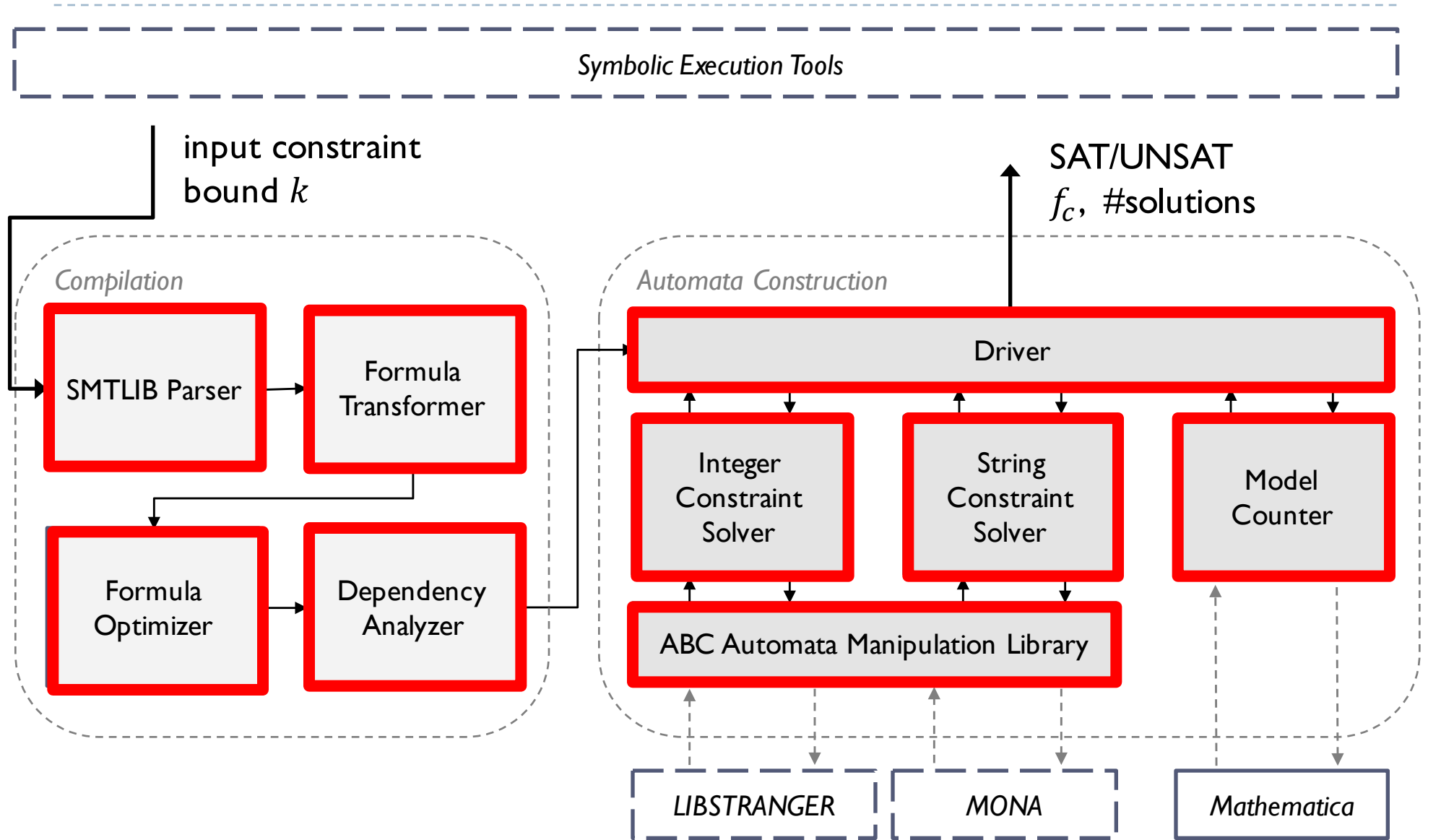


- ▶ $X = Y.C$, $X = C.Y$, $C = X.Y$ can be precisely represented with multi-track automaton
- ▶ $X = Y.Z$ is non-linear, can be represented with over- or under-approximation

Model Counting String Constraints Solver



ABC Tool



Experimental Evaluation

- ▶ We conducted experiments on 3 benchmark sets
 - ▶ A java benchmark with wide range of string operations
 - ▶ Kaluza Small & Kaluza Big benchmarks for satisfiability check

	Frequency of Operations Per 1000 Formulas				
	€	.	=	LENGTH	REPLACEALL, BEGINS SUBSTRING, ENDS CONTAINS, INDEXOF
Kausler	0.42	386.10	129.39	382.54	716.5
Kaluza Small	30.29	93.89	224.87	46.84	0
Kaluza Big	38.12	129.53	164.64	60.46	0

Satisfiability Check Comparison

- ▶ Compared with CVC4
 - ▶ Used SMT-lib format of Kaluza benchmarks from CVC4

	ABC-CVC4 sat-sat	ABC-CVC4 unsat-unsat	ABC-CVC4 sat-unsat	ABC-CVC4 unsat-sat	ABC-CVC4 sat-timeout
sat/small	19728	3	0	0	0
sat/big	1587	0	0	0	0
unsat/small	8139	3013	74	0	0
unsat/big	3419	5904	2385	0	2359

- ▶ Constraint solver performance for Kaluza benchmarks

	ABC Avg.Time (seconds)	CVC4 Avg.Time (seconds)
big	0.44	0.01
small	0.01	0.015

Kauser's Benchmark

- ▶ Extracted from 7 real-world server-side Java applications
 - ▶ Constraints were generated by extracting program path constraints through dynamic symbolic execution
 - ▶ CVC4 1.4 stable version was not able to handle Kauser's benchmark

# of satisfiable path constraints	Avg. # of BDD Nodes (each 16 bytes)	Avg. Running Time (seconds)
66236	69.51	0.0015

- ▶ Automata representation we use encodes the transition relation as a binary decision diagram

Model Counting String Constraints Solver

