

Synthesizing Transducers from Complex Specifications

Anvay Grover, Rüdiger Ehlers, Loris D'Antoni

Verifying string manipulating programs:

A Decision Procedure for Path Feasibility of String Manipulating Programs with Integer Data Types

String Constraints with Concatenation and Transducers
Solved Efficiently

Taolue Chen¹, Matthew Hague

Antonín Štěrba

High-Level Abstractions for Simplifying Extended String Constraints in SMT

Technology, Czech Republic

Technology, Czech Republic

ord, United Kingdom

ty, Sweden

Technology, Czech Republic

Andrew Reynolds¹, Andres Nötzli²,
Clark Barrett², and Cesare Tinelli

CertiStr: A Certified String Solver

Shuanglong Kan

Department of Computer Science
Technische Universität Kaiserslautern
Kaiserslautern, Germany
shuanglong@cs.uni-kl.de

Philipp Rümmer

Department of Information Technology
Uppsala University
Uppsala, Sweden
philipp.ruemmer@it.uu.se

Anthony Widjaja Lin

Department of Computer Science
Technische Universität Kaiserslautern & MPI-SWS
Kaiserslautern, Germany
lin@cs.uni-kl.de

Micha Schrader

Department of Computer Science
Technische Universität Kaiserslautern
Kaiserslautern, Germany
schrader@rhrk.uni-kl.de

```
var x = goog.string.htmlEscape(name);  
var y = goog.string.escapeString(x);  
nameElem.innerHTML = '<button onclick= "viewPerson(\'' + y + '\')">' + x + '</button>';
```



Sanitize inputs

Prevent XSS attacks

```
function toUpperCase(inputStr) {  
    ...  
}
```

`toUpperCase(Aba)` → ABA

`toUpperCase(ABA)` → ABA

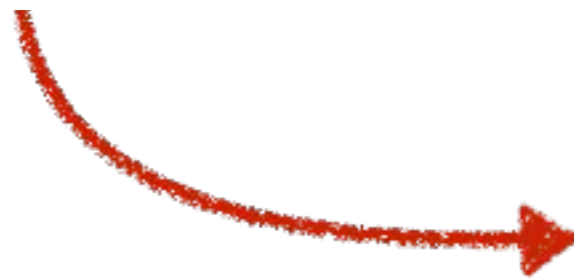
The `toUpperCase` function is idempotent

Suppose we want to prove:

$y = \text{toUpperCase}(x)$

A Decision Procedure for Path Feasibility
of String Manipulating Programs
with Integer Data Type

Taolue Chen¹, Matthew Hague², Jinlong He^{3,6}, Denghang Hu^{3,6},
Anthony Widjaja Lin⁴, Philipp Rümmer⁵, and Zhilin Wu^{3,7,8}(✉)

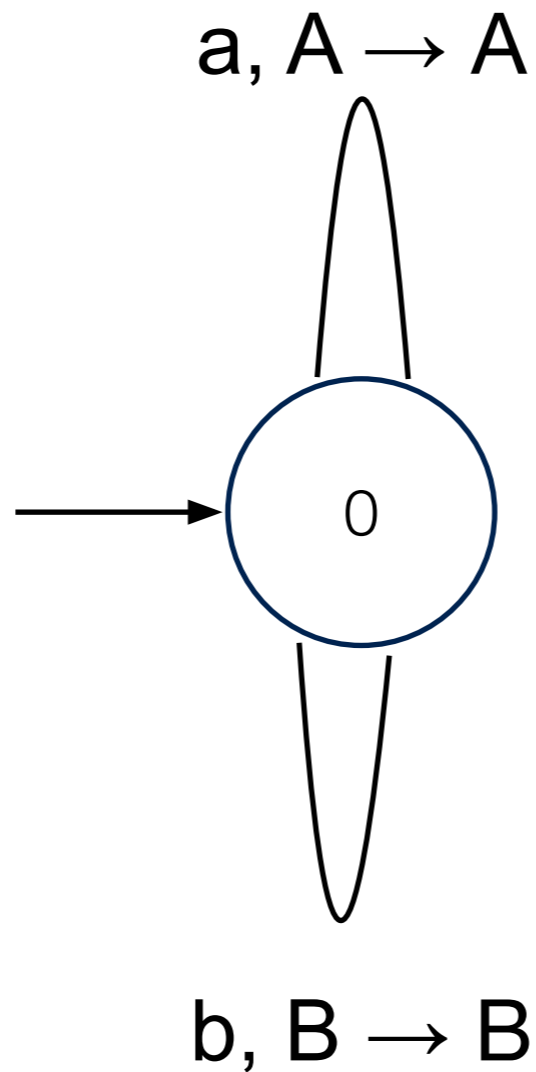


Ostrich: An SMT Solver
for string constraints

But...

Ostrich models string to string functions as transducers!

toUpperCase



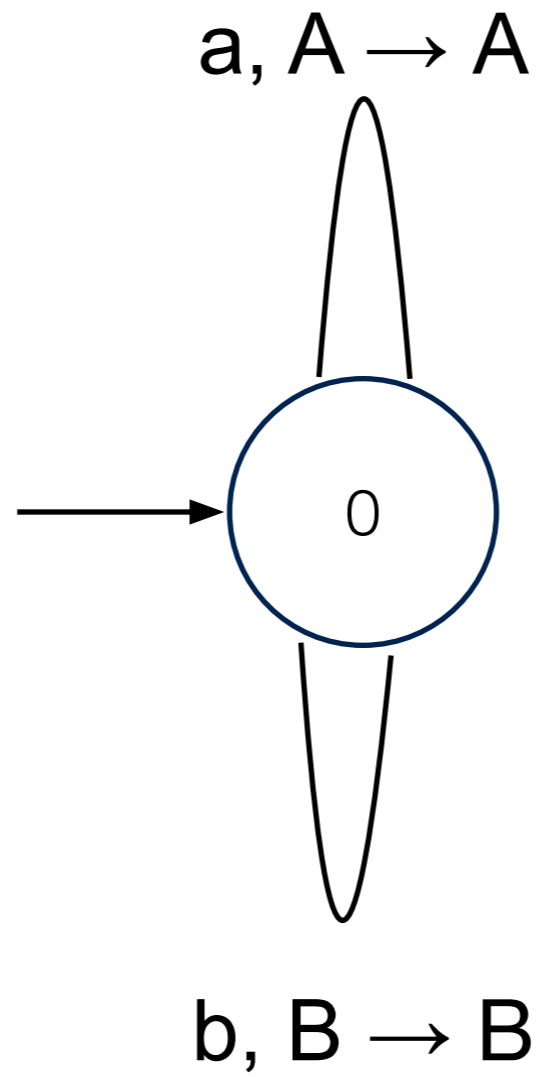
Why Transducers?

Decidability Properties

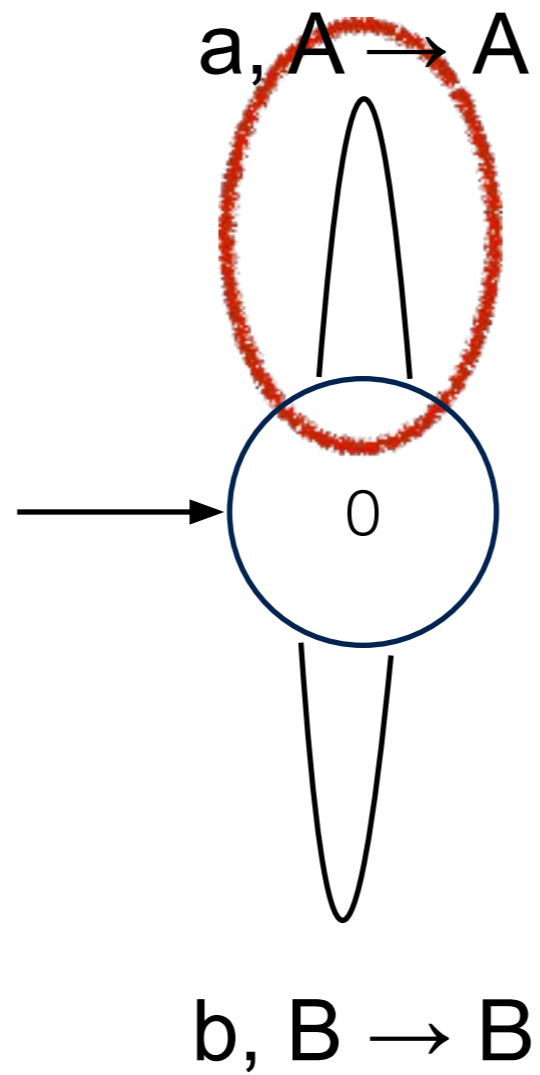
Closure under Composition

Equivalence Checking

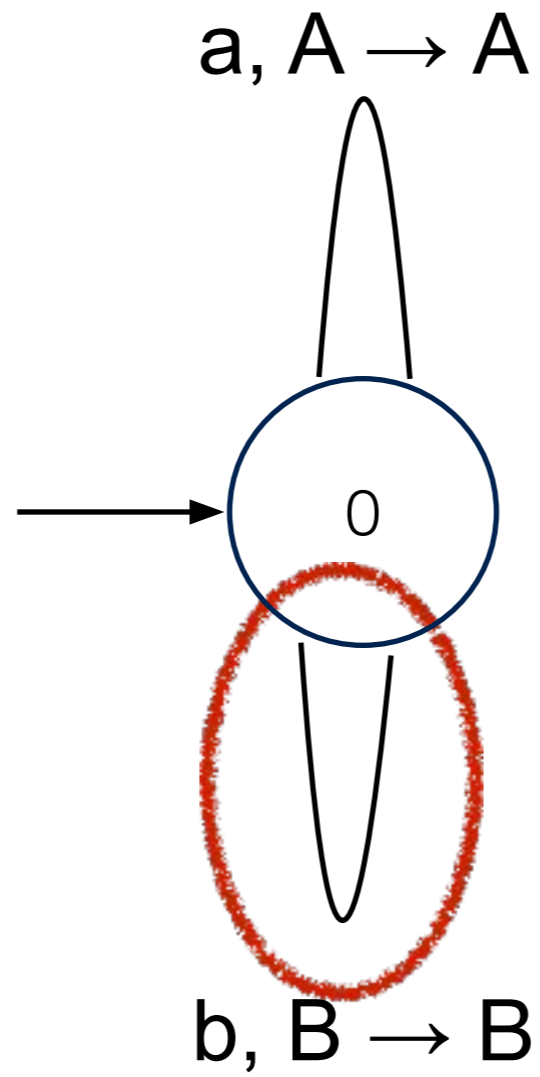
`toUpperCase(toUpperCase()) ≡ toUpperCase()`



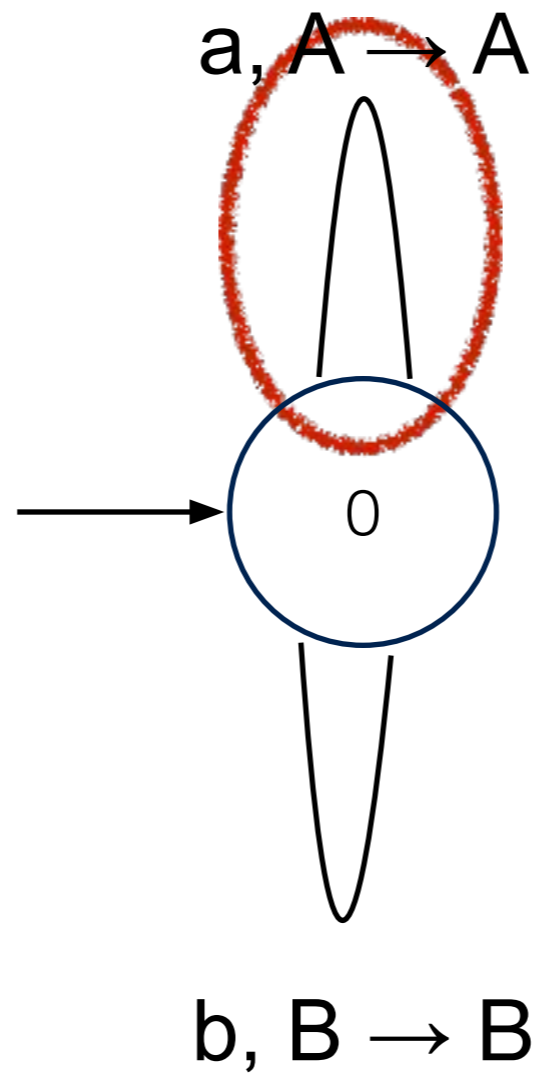
toUpperCase



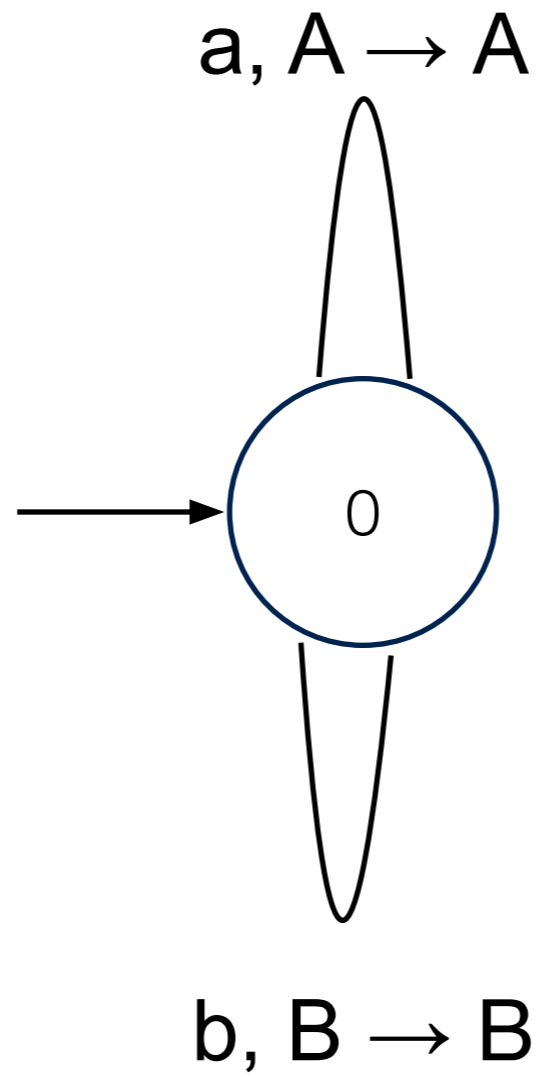
toUpperCase(aBa) → ABA



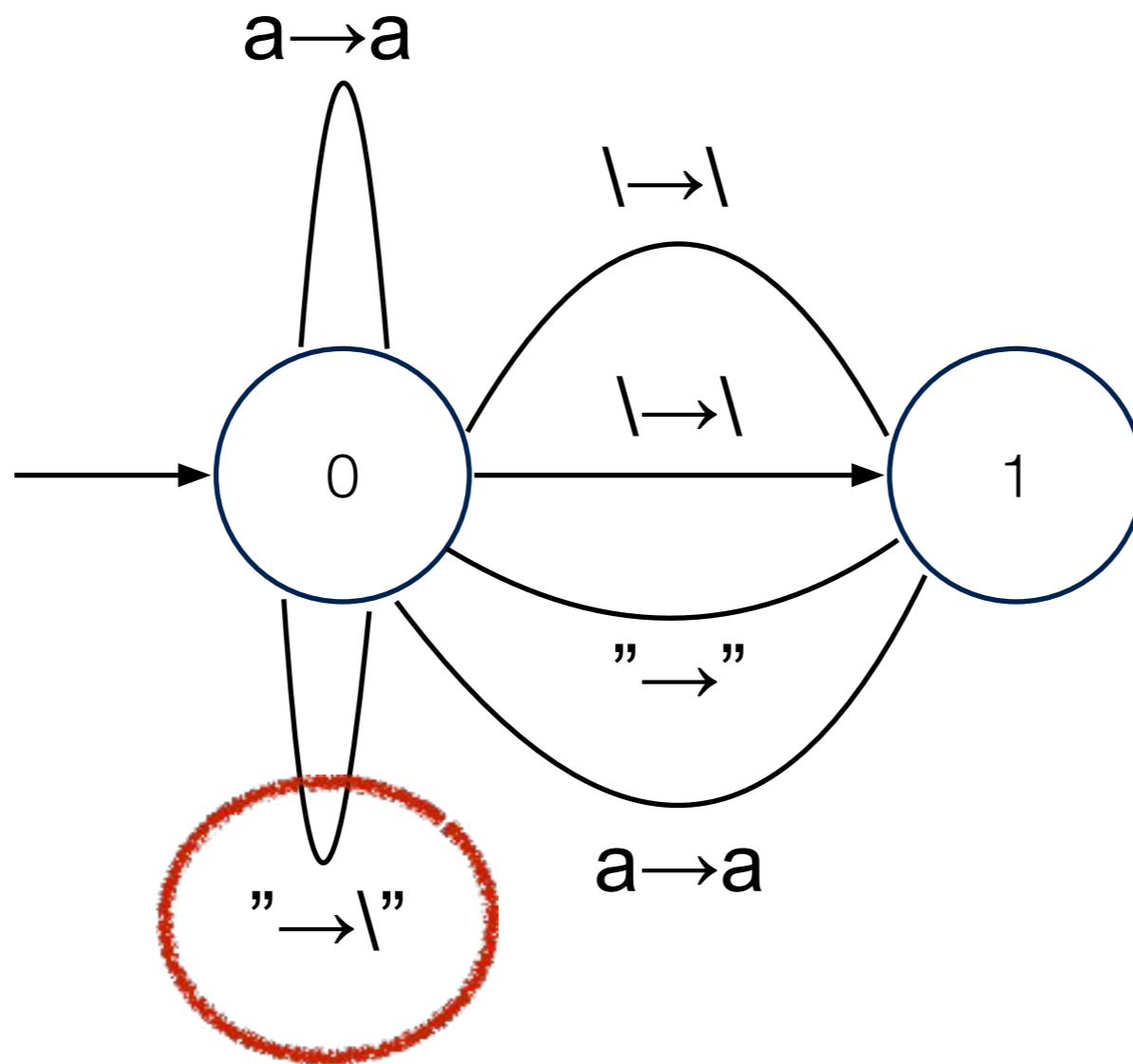
toUpperCase(aBa) \rightarrow ABA



toUpperCase(aBa) \rightarrow AB**A**



toUpperCase



escapeQuotes

We automate writing
transducers **by**
synthesizing them

<i>Input v₁</i>	<i>Output</i>
<i>International Business Machines</i>	<i>IBM</i>
<i>Principles Of Programming Languages</i>	<i>POPL</i>
<i>International Conference on Software Engineering</i>	<i>ICSE</i>

generate abbreviation

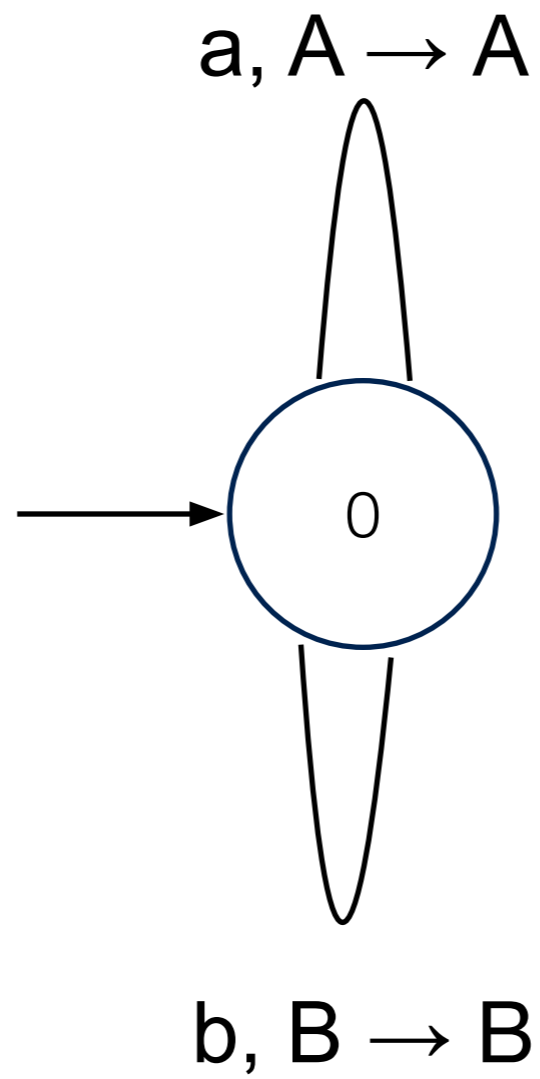
<i>Input v₁</i>	<i>Output</i>
<i>BTR KRNL WK CORN 15Z</i>	<i>15Z</i>
<i>CAMP DRY DBL NDL 3.6 OZ</i>	<i>3.6 OZ</i>
<i>CHORE BOY HD SC SPNG 1 PK</i>	<i>1 PK</i>
<i>FRENCH WORCESTERSHIRE 5 Z</i>	<i>5 Z</i>
<i>O F TOMATO PASTE 6 OZ</i>	<i>6 OZ</i>

extract quantities

[Automating String Processing in Spreadsheets using Input-Output Examples] POPL '11

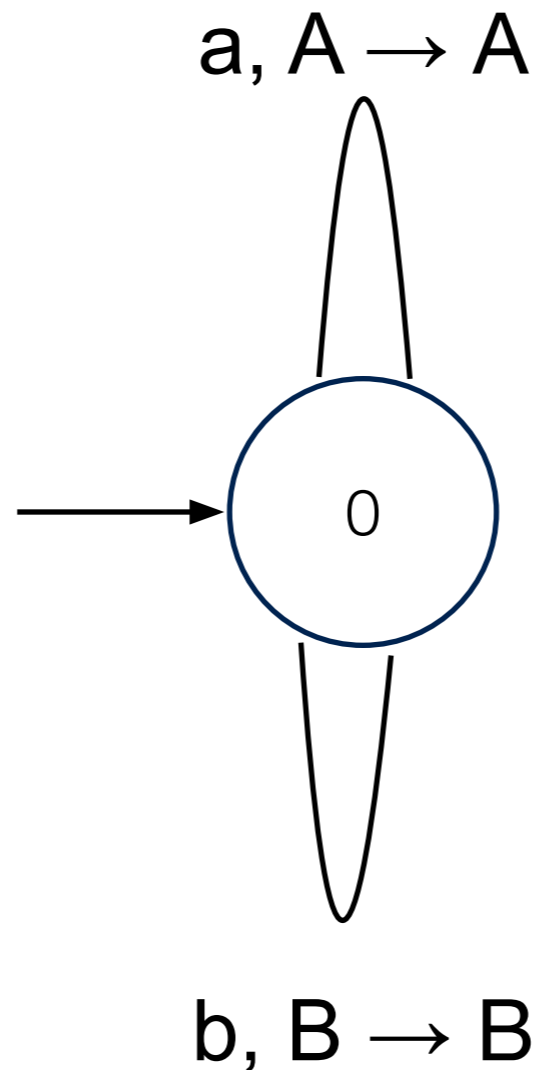
How do we do this?

**Supporting various
specification
mechanisms**



$aBa \rightarrow ABA$

$ABA \rightarrow ABA$

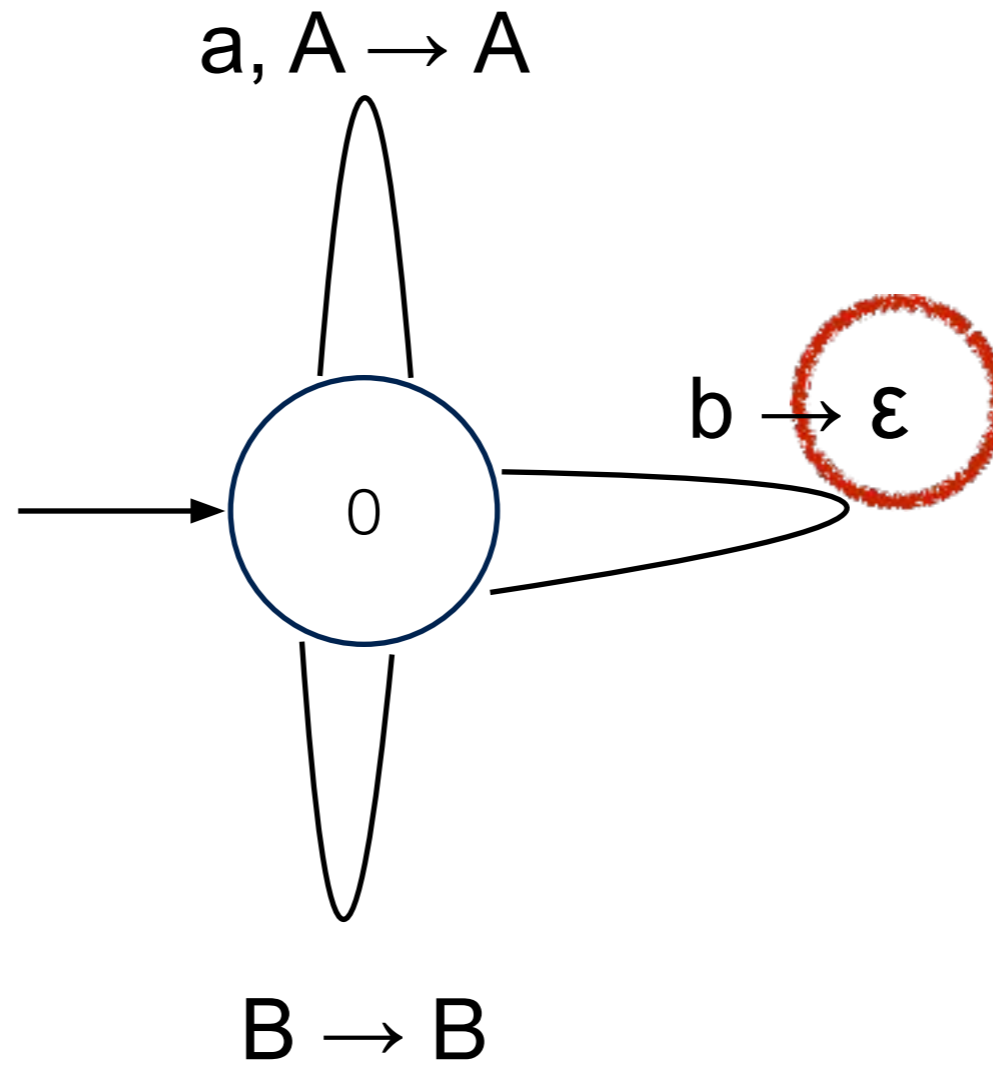


1. Input-Output Examples

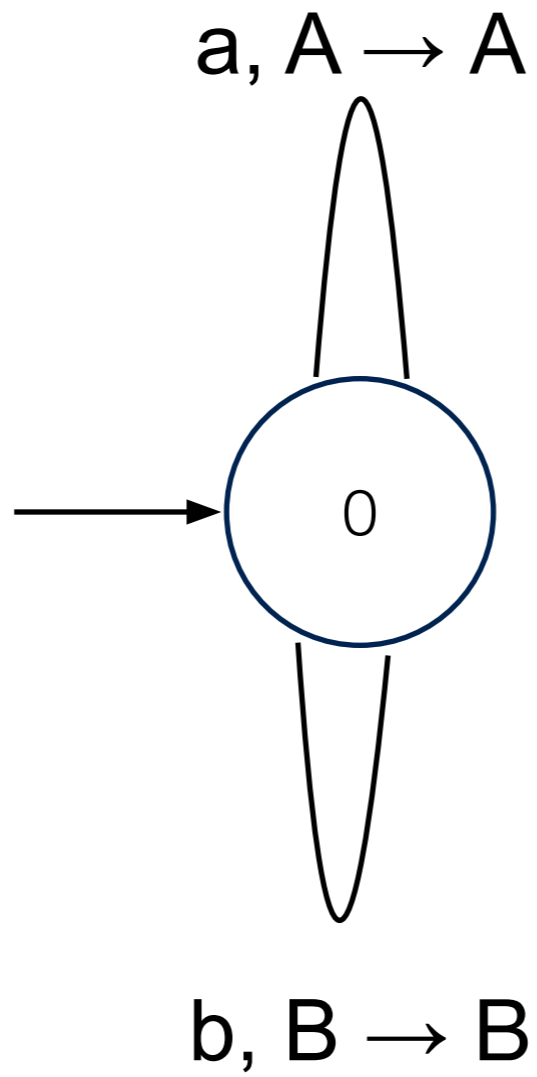
For an example $i \rightarrow o$, we write $T(i) = o$

aBa \rightarrow ABA

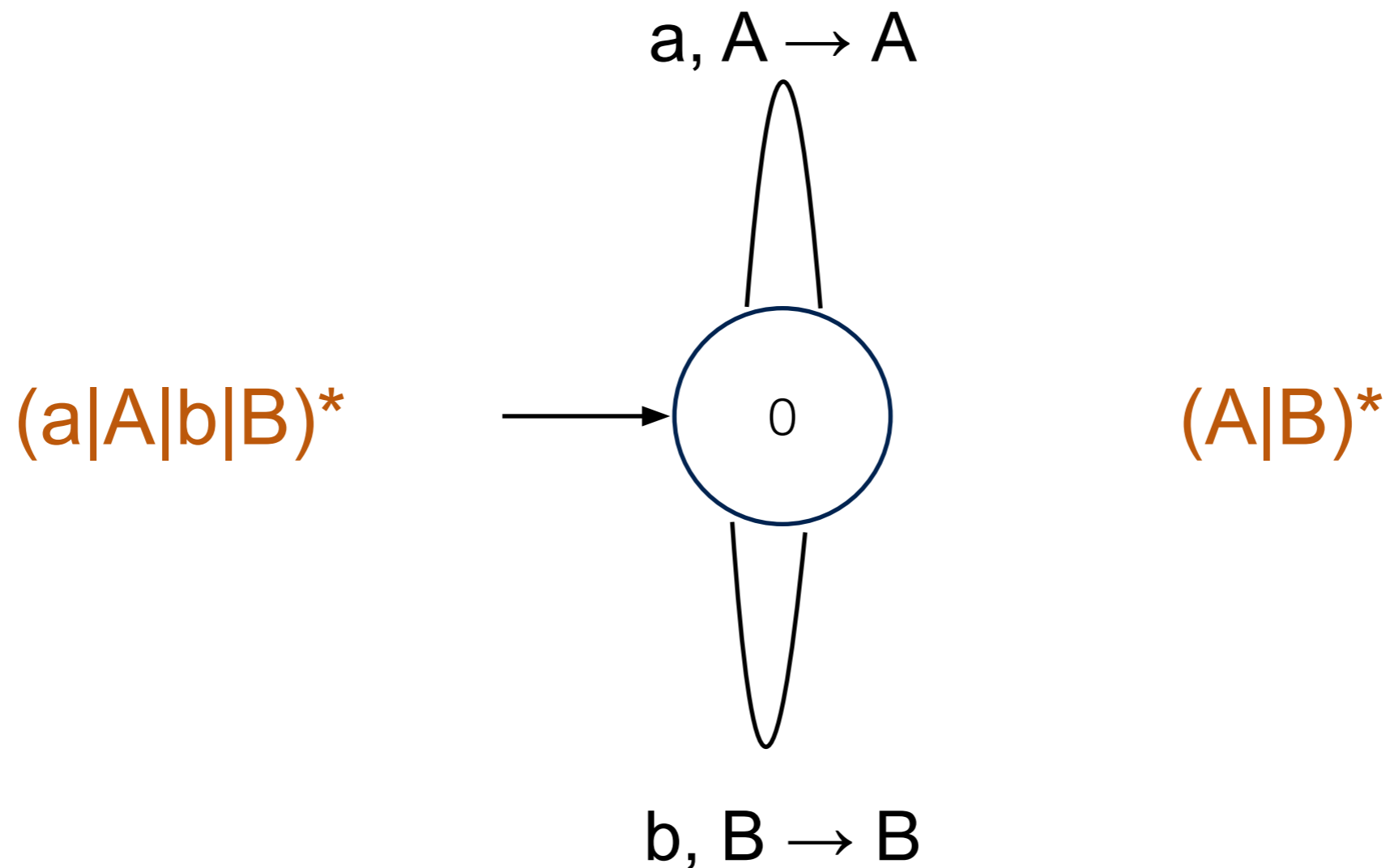
ABA \rightarrow ABA



A valid transducer can overfit!



$$(a|A|b|B)^* \rightarrow (A|B)^*$$



2. Input-Output Types

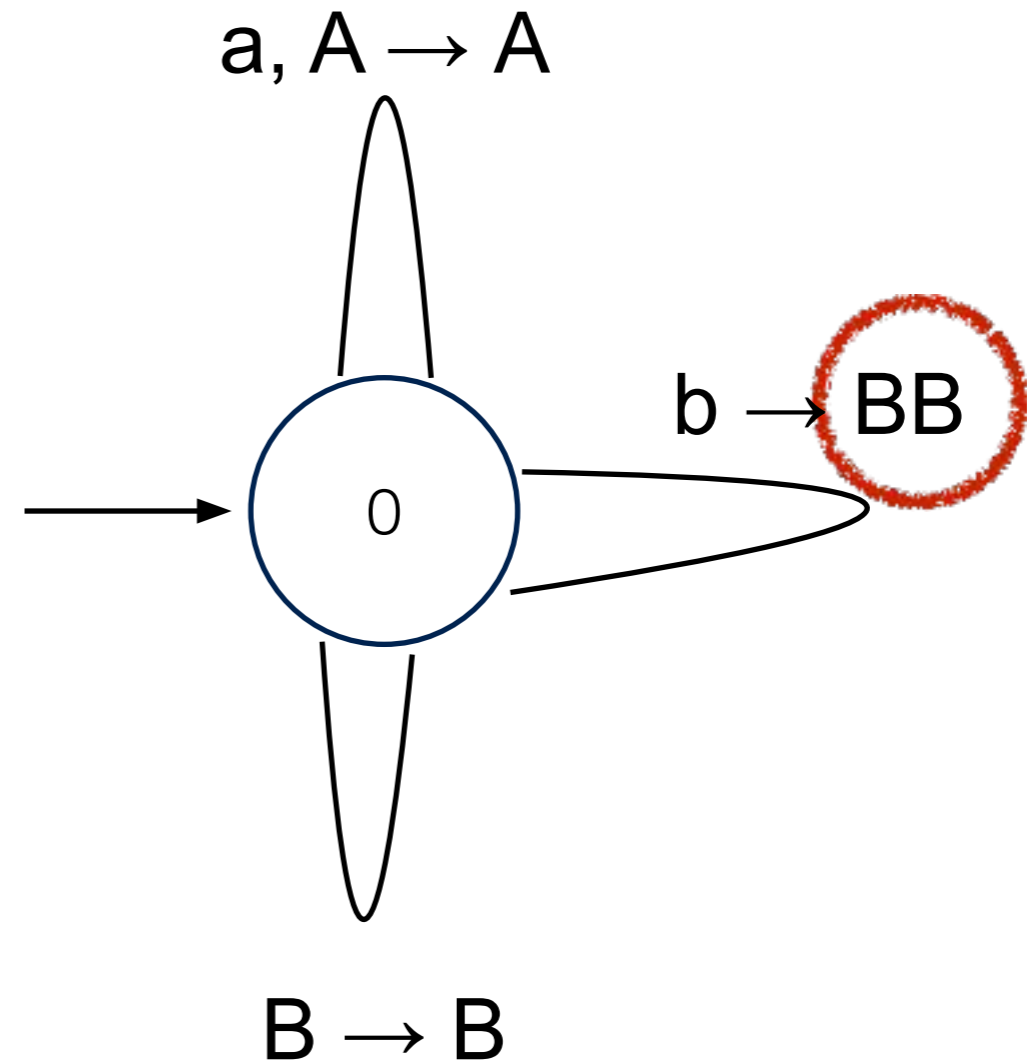
Specify pre and post conditions

For input-output regexes I , O , we write $\{I\} T \{O\}$

$aBa \rightarrow ABA$

$ABA \rightarrow ABA$

$(a|A|b|B)^* \rightarrow (A|B)^*$



Still not enough:

Possible to satisfy types in a way that is unintended

aBa → ABA

ABA → ABA

toUpperCase

aBa → ABA

ABA → ABA

toUpperCase

How 'far apart' can input-output strings be?

aBa → ABA

ABA → ABA

toUpperCase

3. Input-Output Distance

aBa → ABA

We made 2 edits over a string of length 3



mean edit distance = edits/length = 2/3

Why normalize by length?

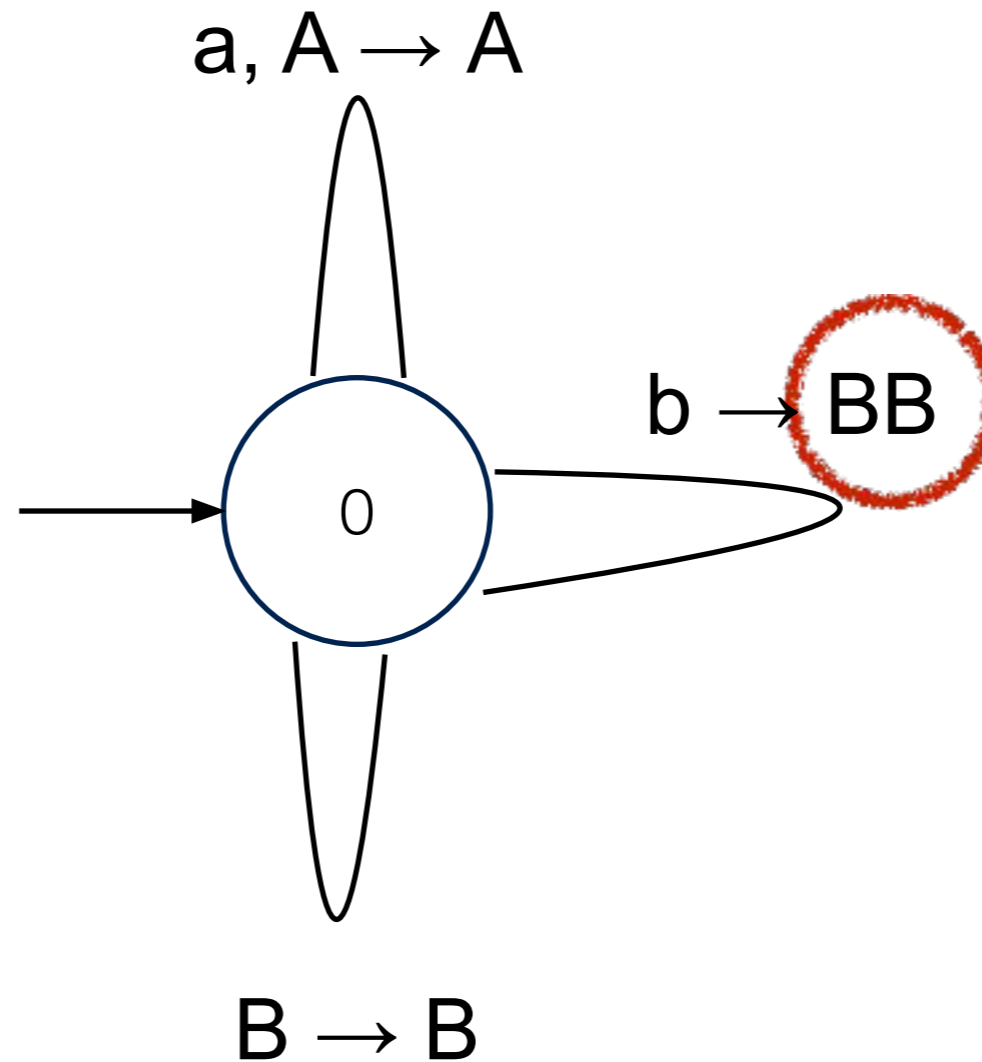
Distance can be bounded even when making an unbounded number of edits

aBa → ABA

For two strings s , t , we write ***dist(s, t)*** for their mean edit distance

Transducer T satisfies mean edit distance d if

$$\forall i \in I. \text{dist}(i, T(i)) \leq d$$



To prevent extra characters from being written:

$$d = 1/1$$

Inputs: regexes I, O ; examples $\{\langle i_1, o_1 \rangle, \dots, \langle i_n, o_n \rangle\}$; distance d

Output: T such that

$\{ I \} T \{ O \}$

Types

$T(i_1) = o_1$

Examples

...

$T(i_n) = o_n$

$\forall i \in I. \text{dist}(i, T(i)) \leq d$

Distance

Inputs: regexes I, O ; examples $\{\langle i_1, o_1 \rangle, \dots, \langle i_n, o_n \rangle\}$; distance d

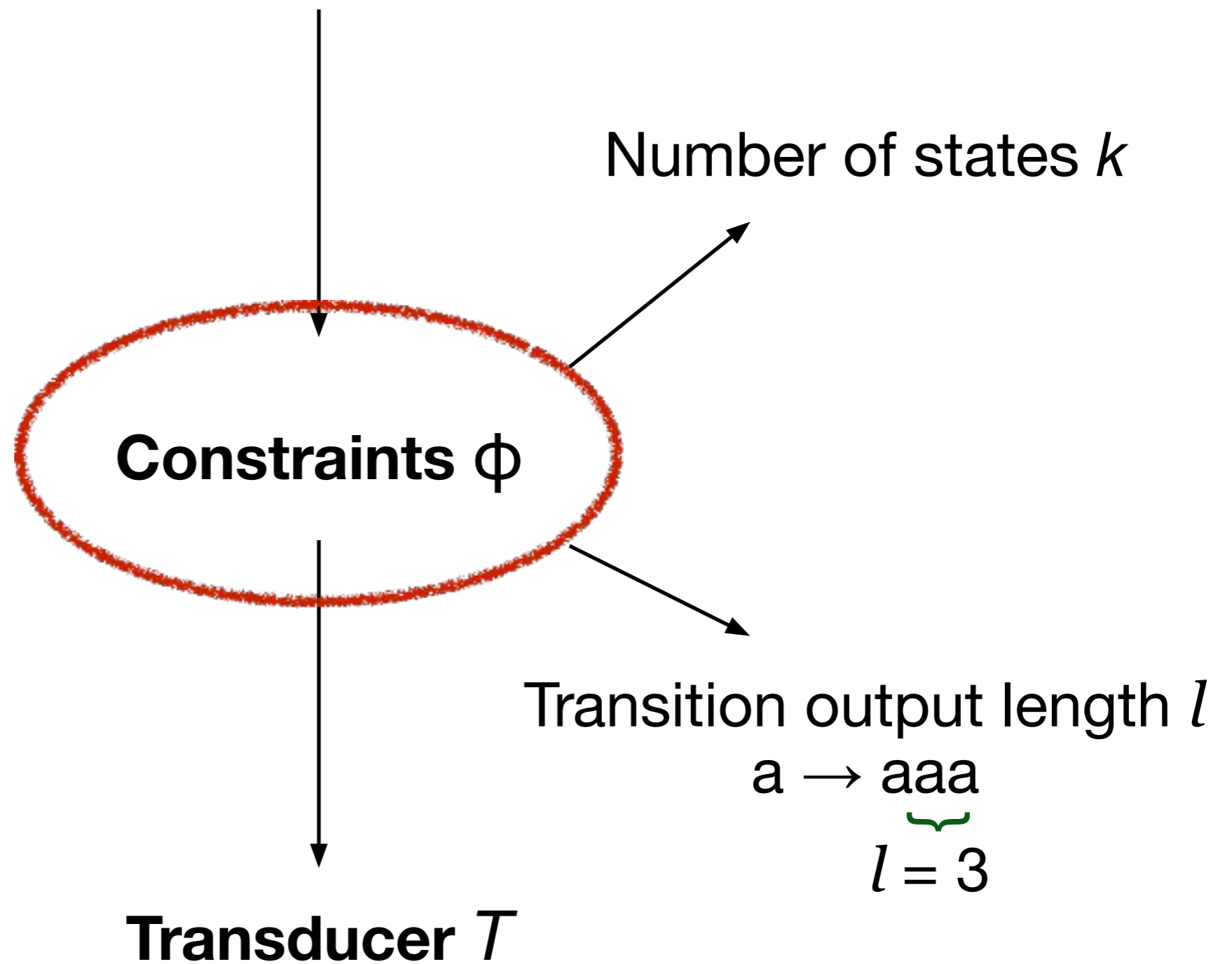


Constraints ϕ



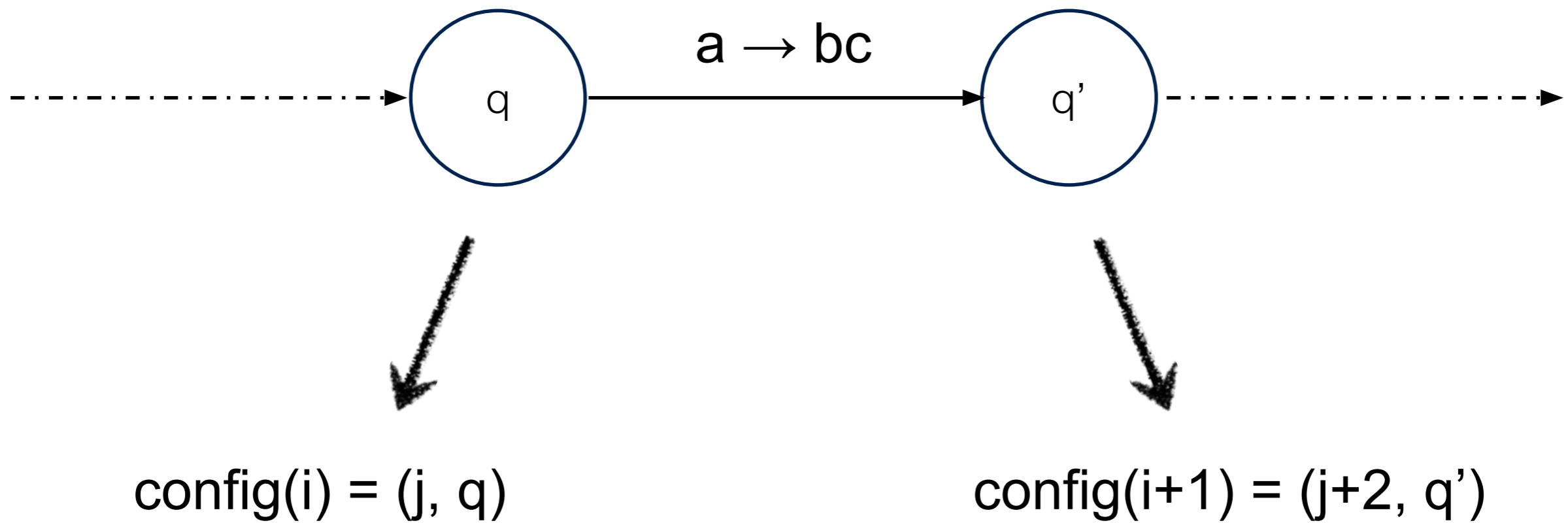
Transducer T

Inputs: regexes I, O ; examples $\{\langle i_1, o_1 \rangle, \dots, \langle i_n, o_n \rangle\}$; distance d



Input-Output Examples

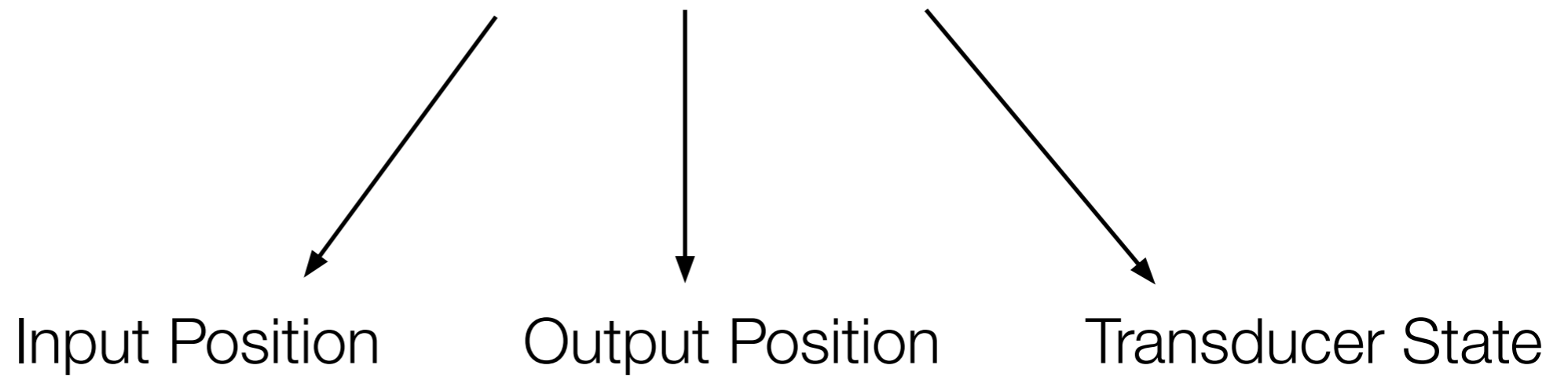
$$x_0 \dots x_i a \dots x_n \rightarrow y_0 \dots y_j bc \dots y_n$$



We have read i chars and written j chars

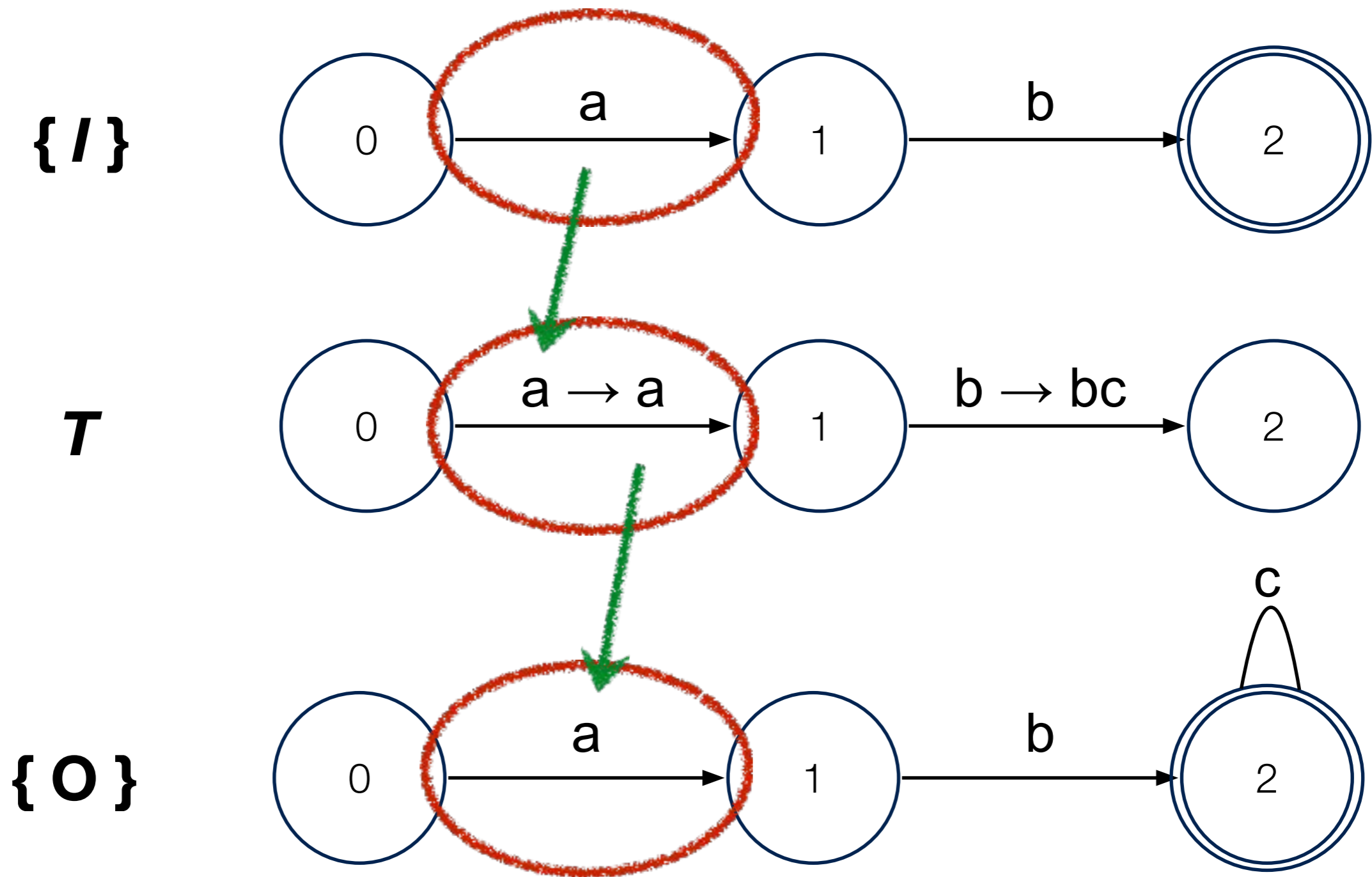
For each example:

$$\text{config: } Z \rightarrow Z \times Q_T$$

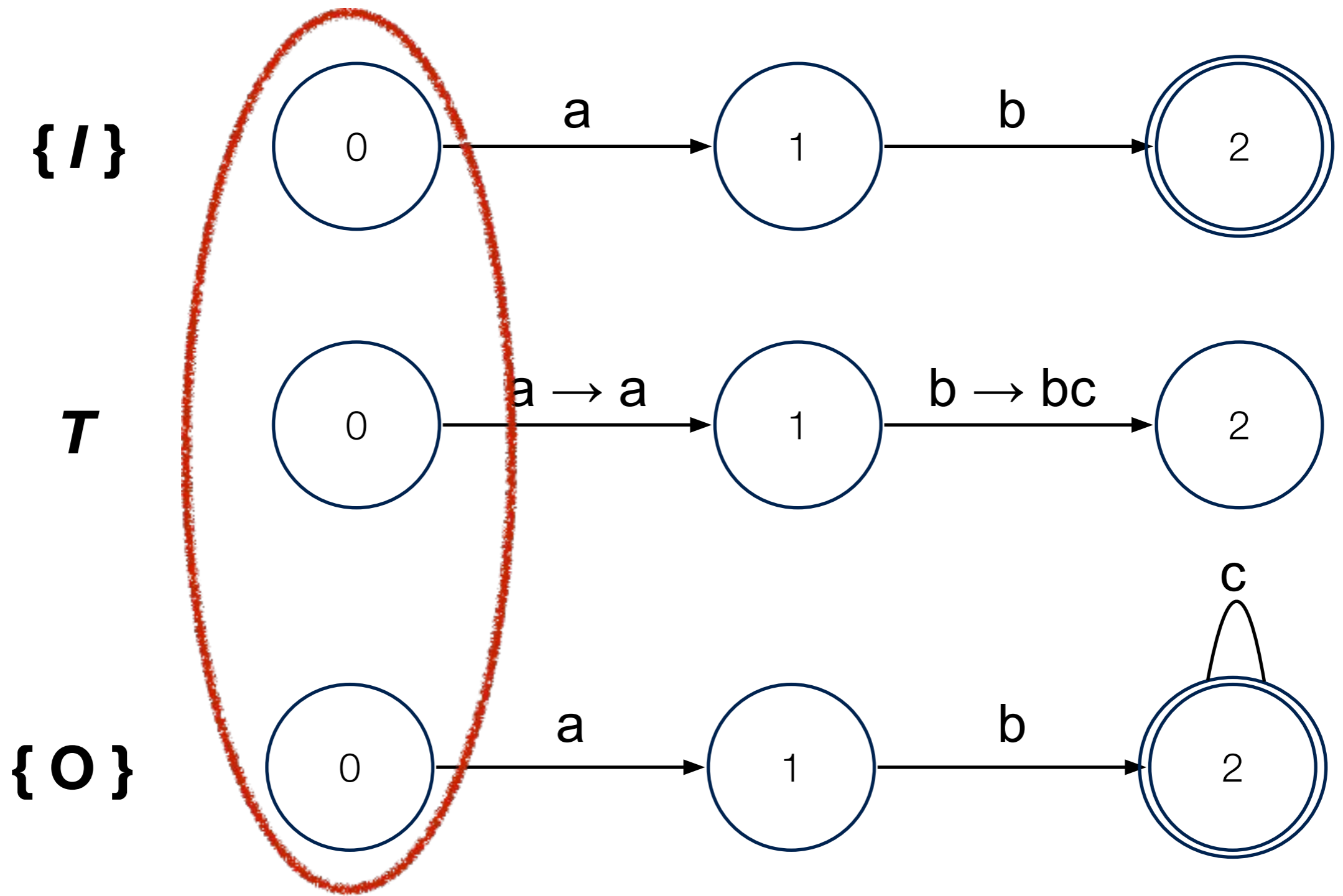


Unknowns!

Input-Output Types + Distance

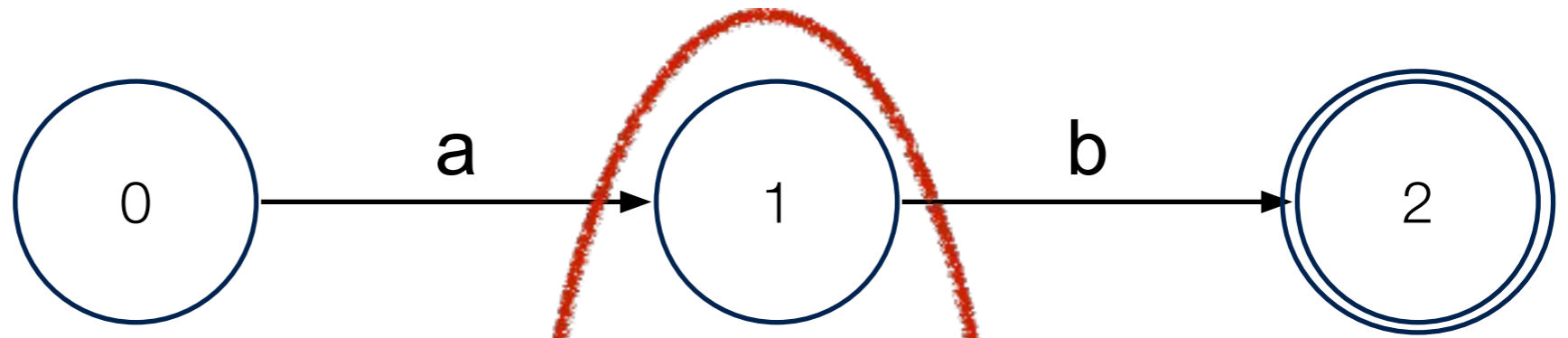


How do we run all 3 automata together?
 A simulation relation

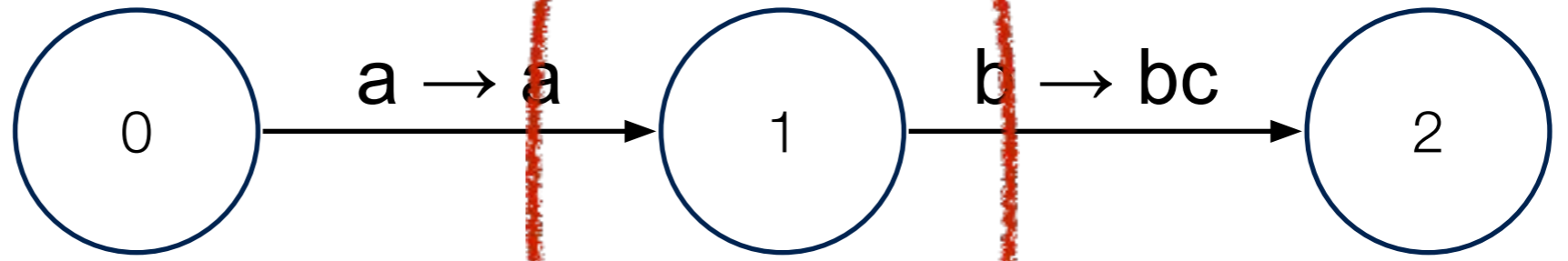


$$\text{sim}(q^I_{\text{init}}, q^T_{\text{init}}, q^O_{\text{init}})$$

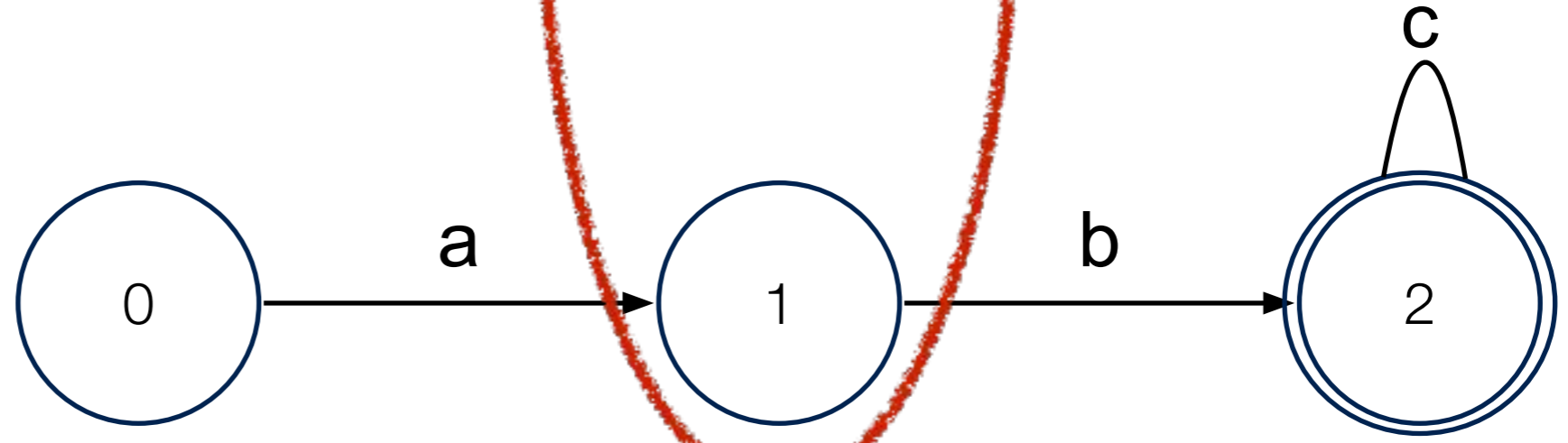
{ / }



T

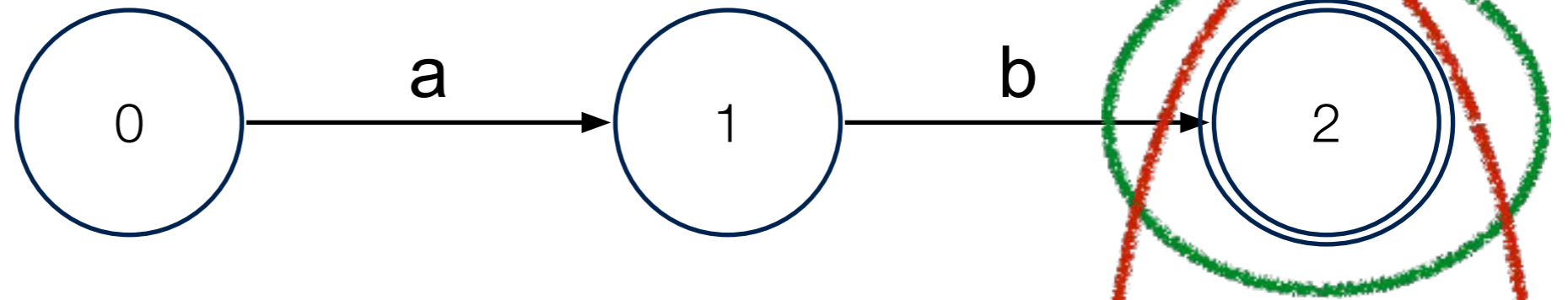


{ 0 }

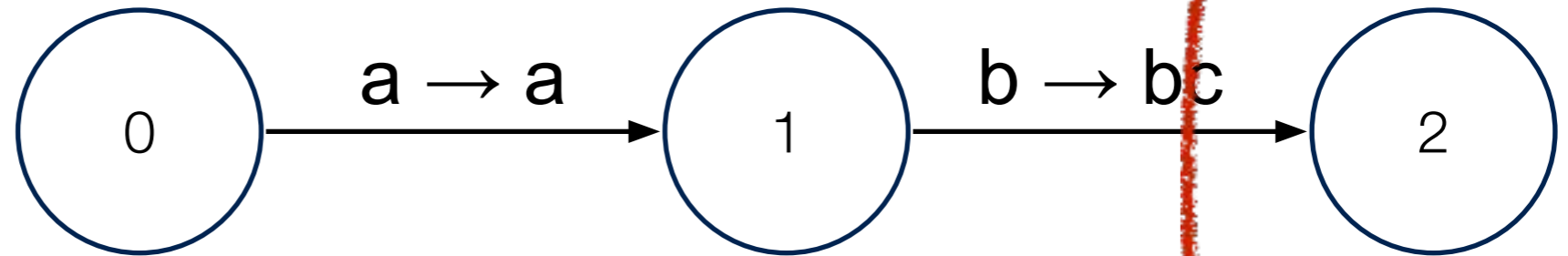


$$\text{sim}(q^I_1, q^T_1, q^O_1)$$

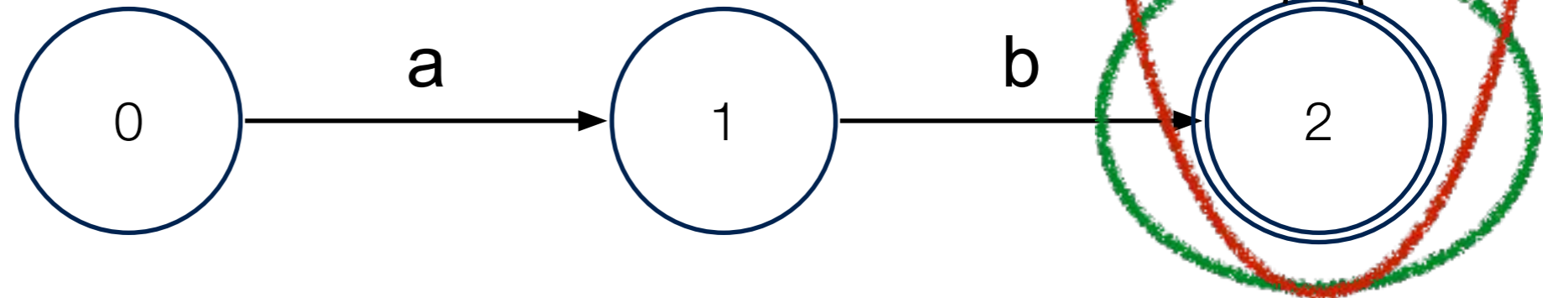
{ / }



T



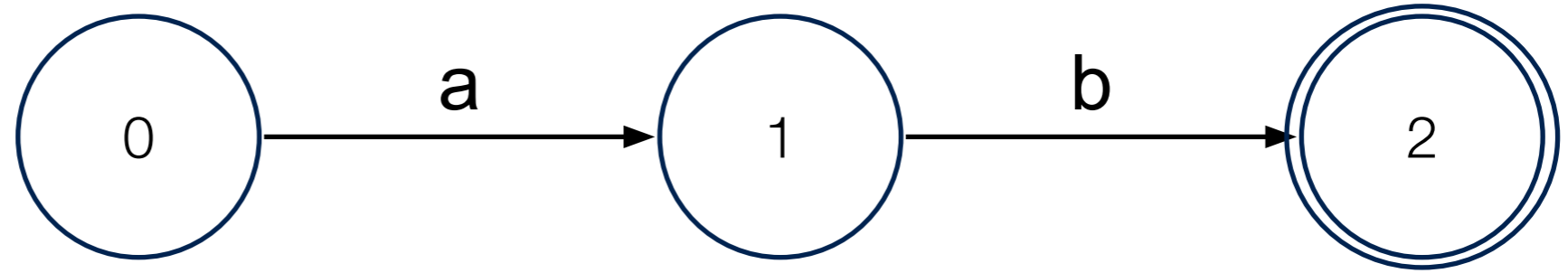
{ 0 }



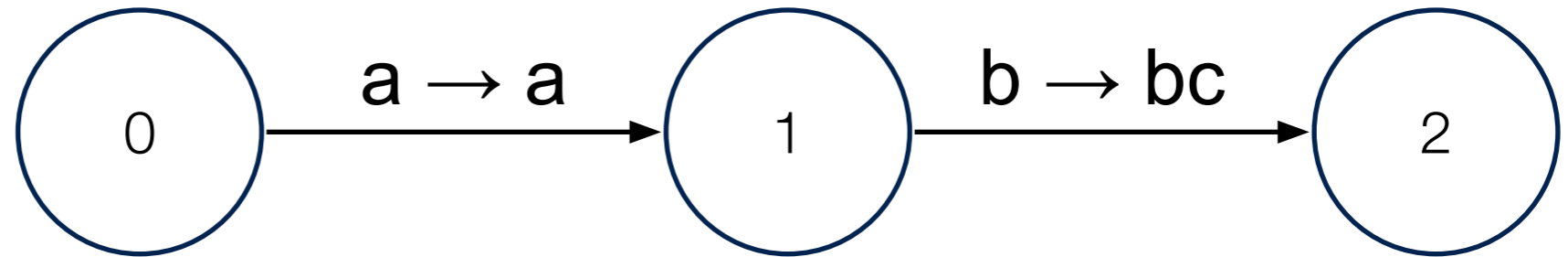
$$\text{sim}(q^I, q^T, q^O)$$

If q^I is final, q^O should also be final

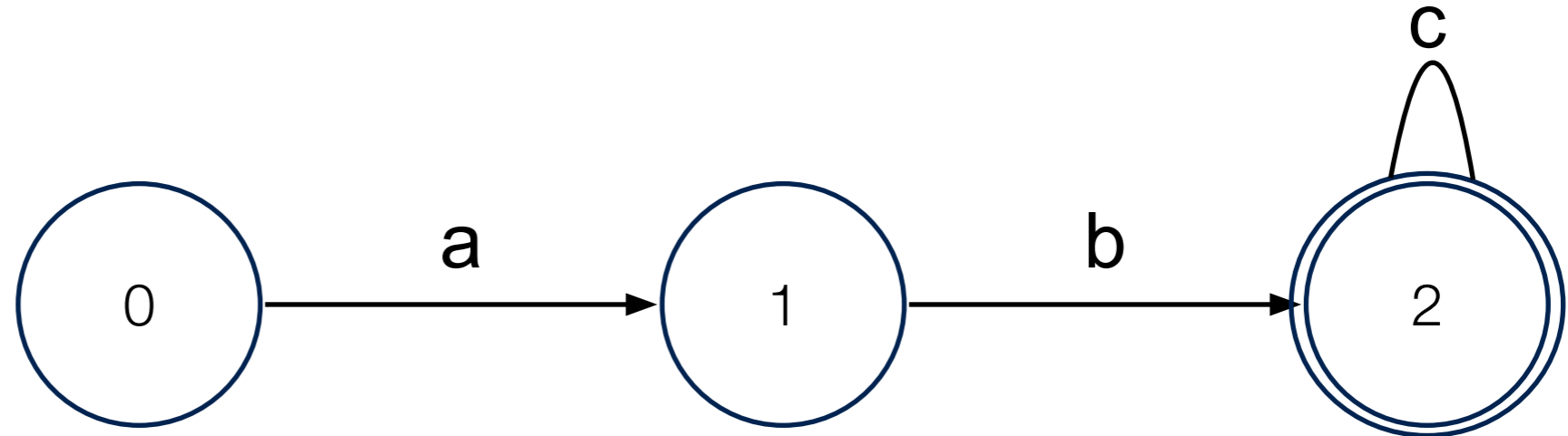
{ / }



T



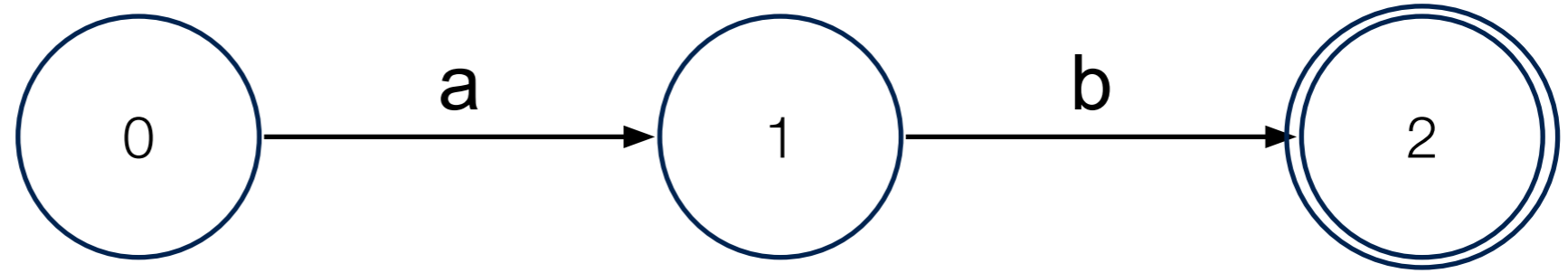
{ 0 }



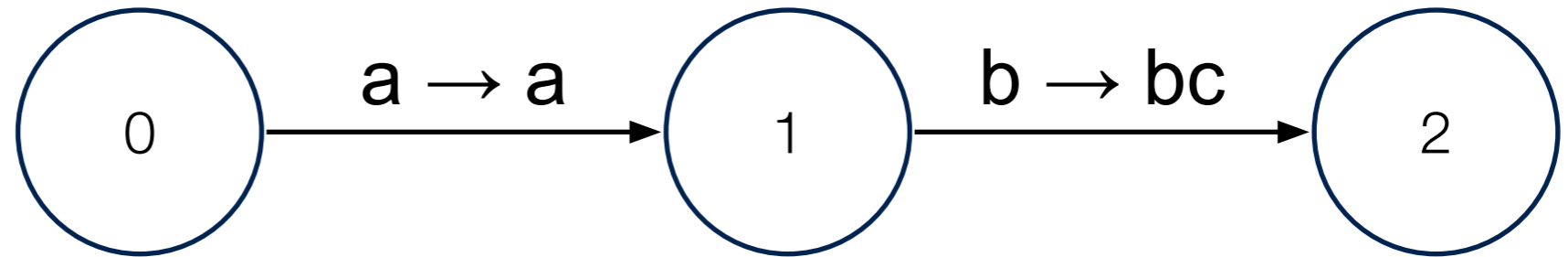
Now what about distance?

Associate a budget value with each point in the simulation

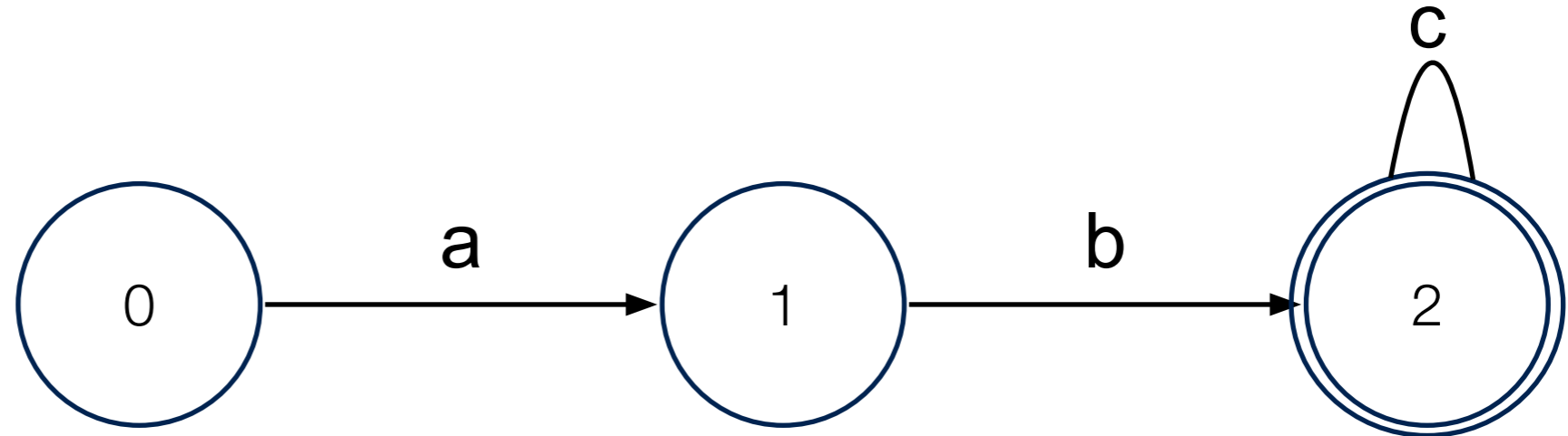
{ / }



T



{ 0 }

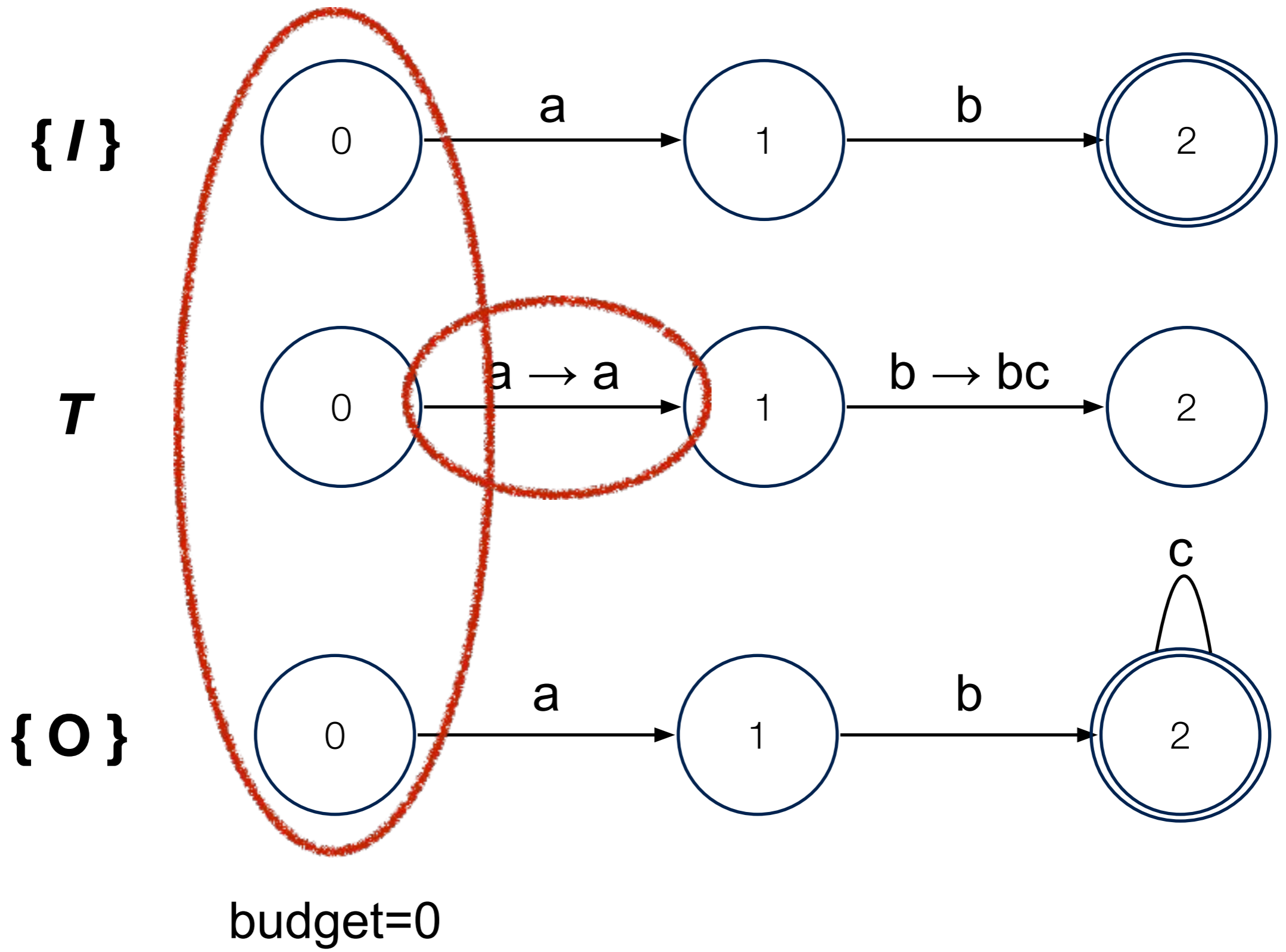


$$d = 1/2$$

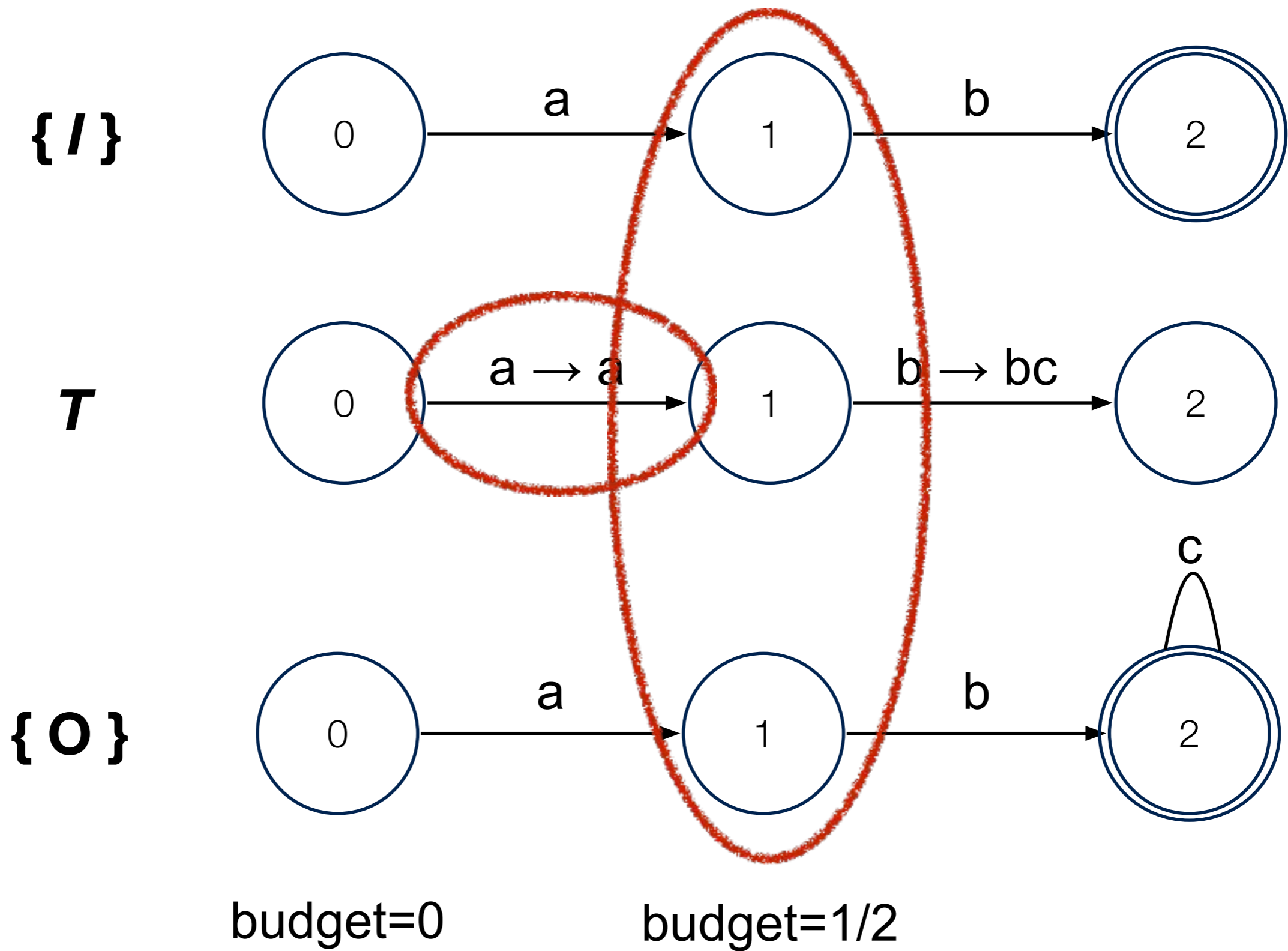
Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$

Want (budget ≥ 0) to not be in a deficit

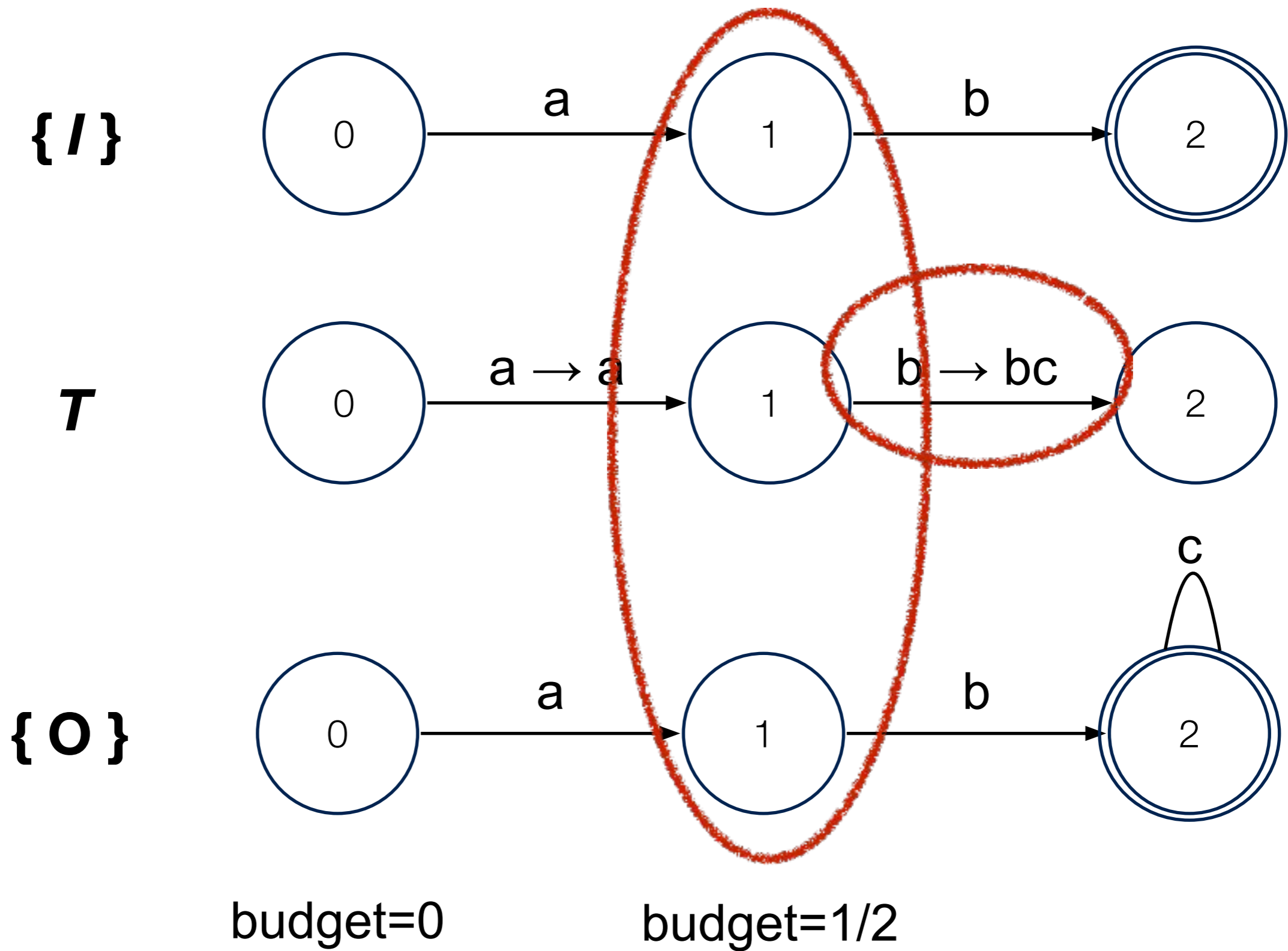
Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$



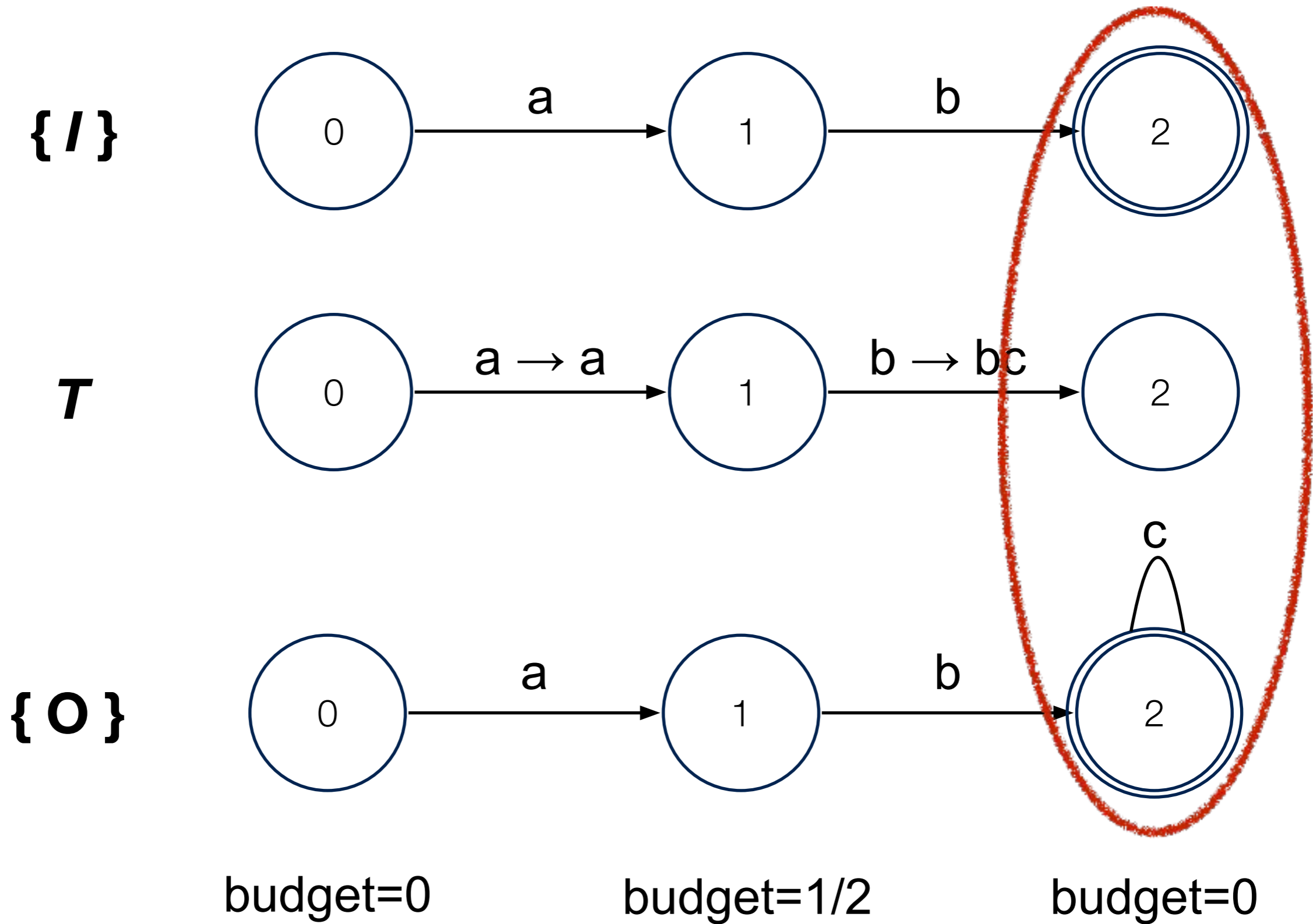
Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$



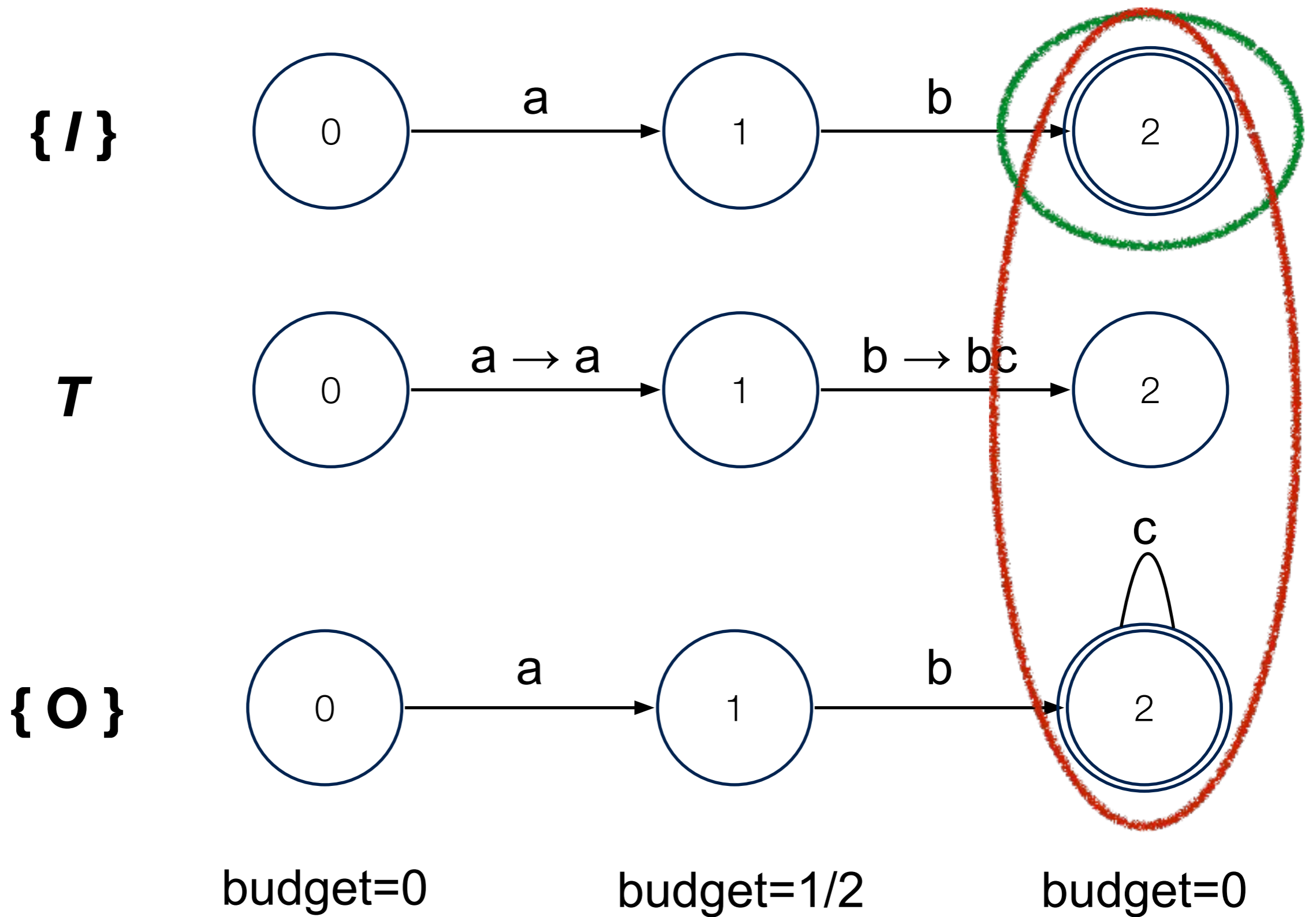
Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$



Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$

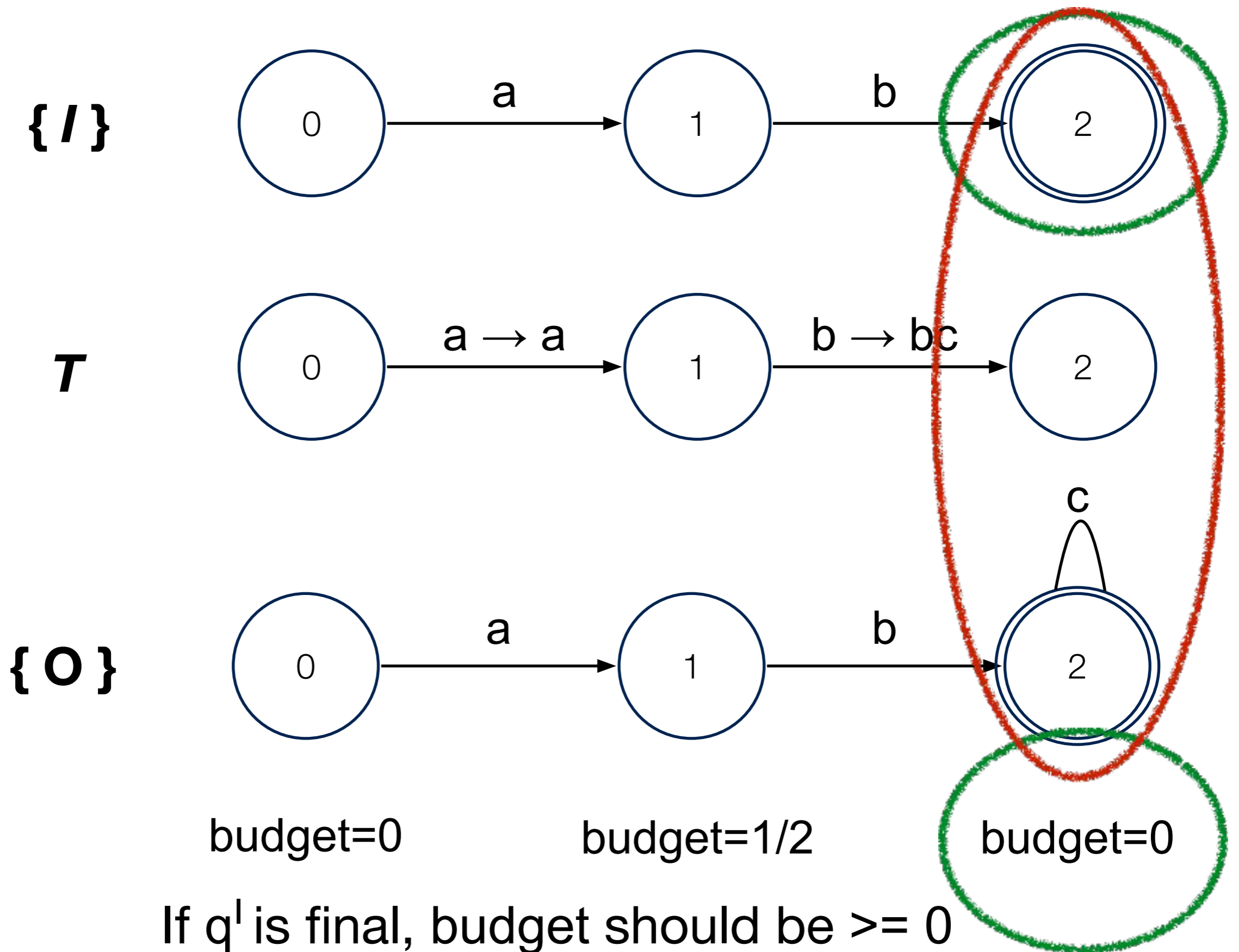


Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$



If q^i is final, budget should be ≥ 0

Goal: $\forall i \in I. \text{dist}(i, T(i)) \leq 1/2$



Constraints use uninterpreted functions:

$$\text{sim: } Q_I \times Q_T \times Q_O \rightarrow \{0, 1\}$$

$$\text{budget: } Q_I \times Q_T \times Q_O \rightarrow \mathbb{Z}$$

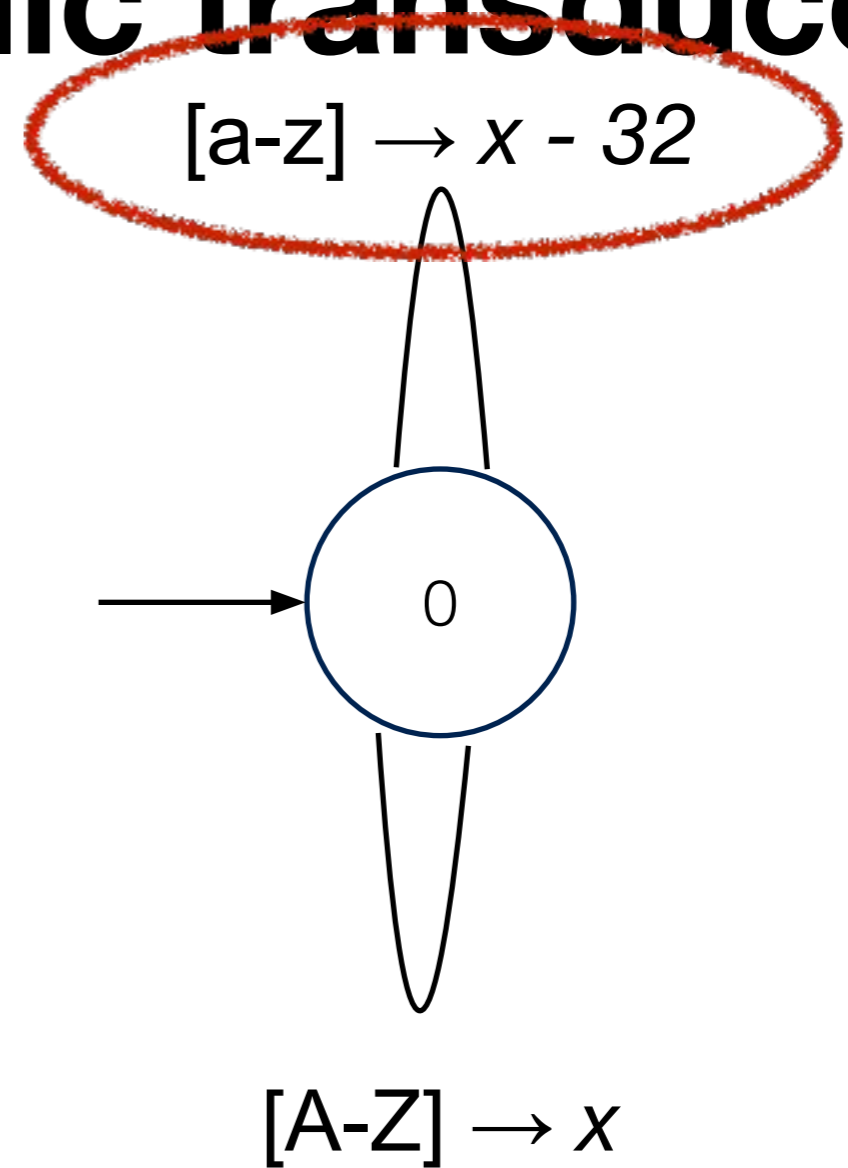
budget computation is sound but incomplete:

We operate at transition level instead of string level

Also in the paper:

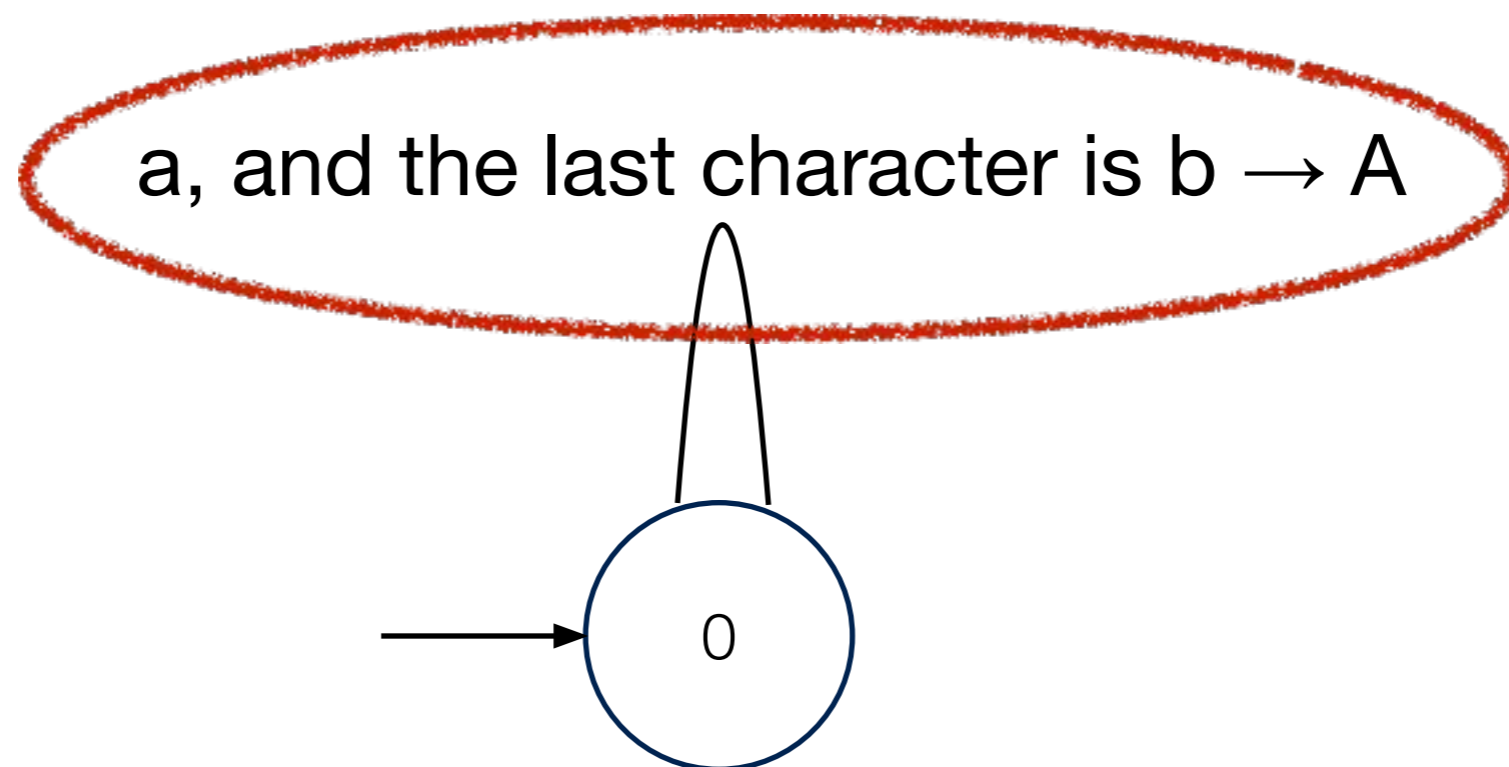
Synthesizing symbolic transducers

symbolic toUpperCase



Also in the paper:

Synthesizing transducers with lookahead (non-determinism)



Lookahead: Know what is going to come later

Evaluation

Implemented in a tool: Astra

Benchmarks

What can we evaluate on?

Optician: Tool for synthesizing
bidirectional transformations

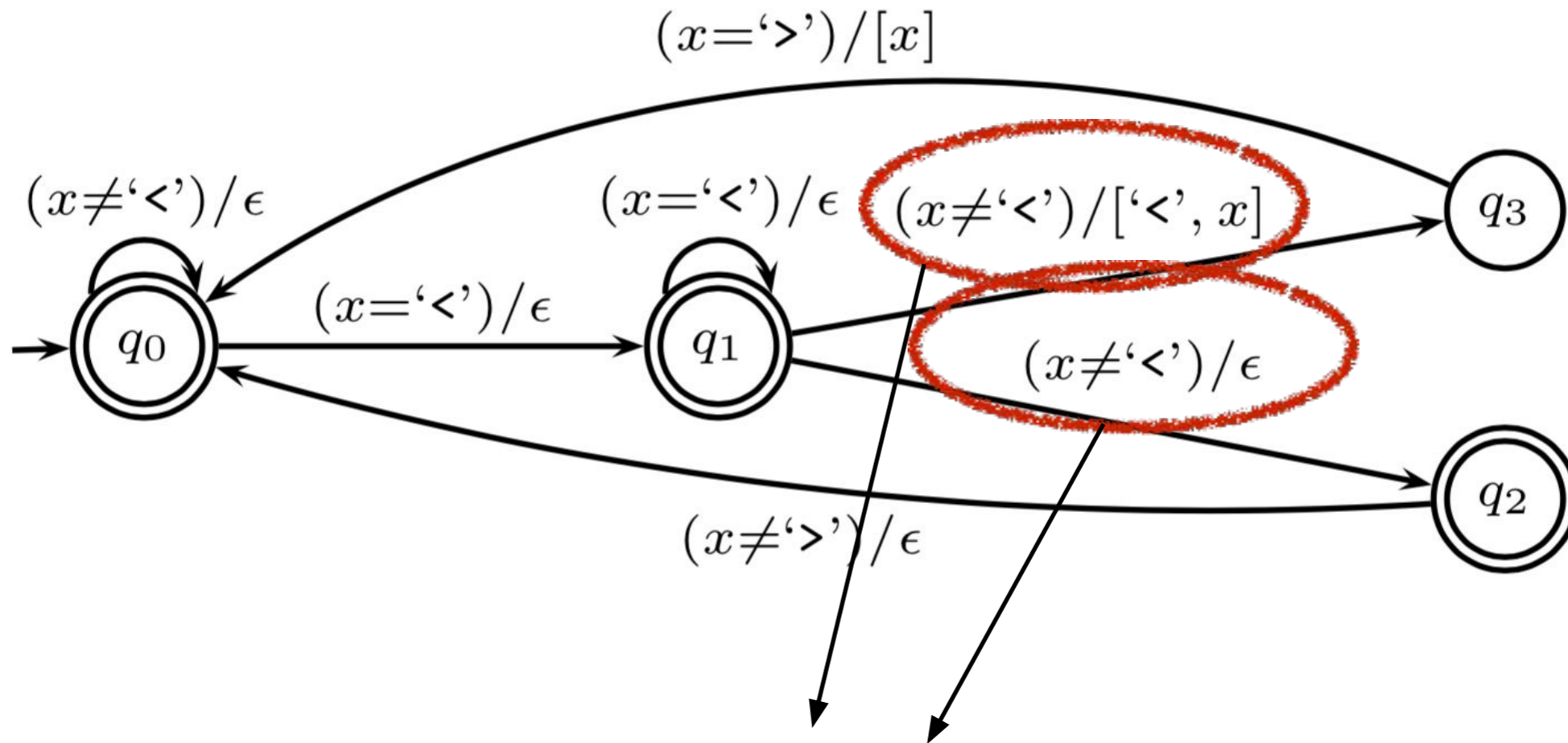
Also uses IO types+examples

Benchmarks:

1. Optician's benchmark suite
(data transformation tasks)
2. Examples from work on
symbolic transducers

The 'getTags' transducer

[Symbolic Finite State Transducers: Algorithms and Applications] POPL '12



Same input predicate (nondeterministic)

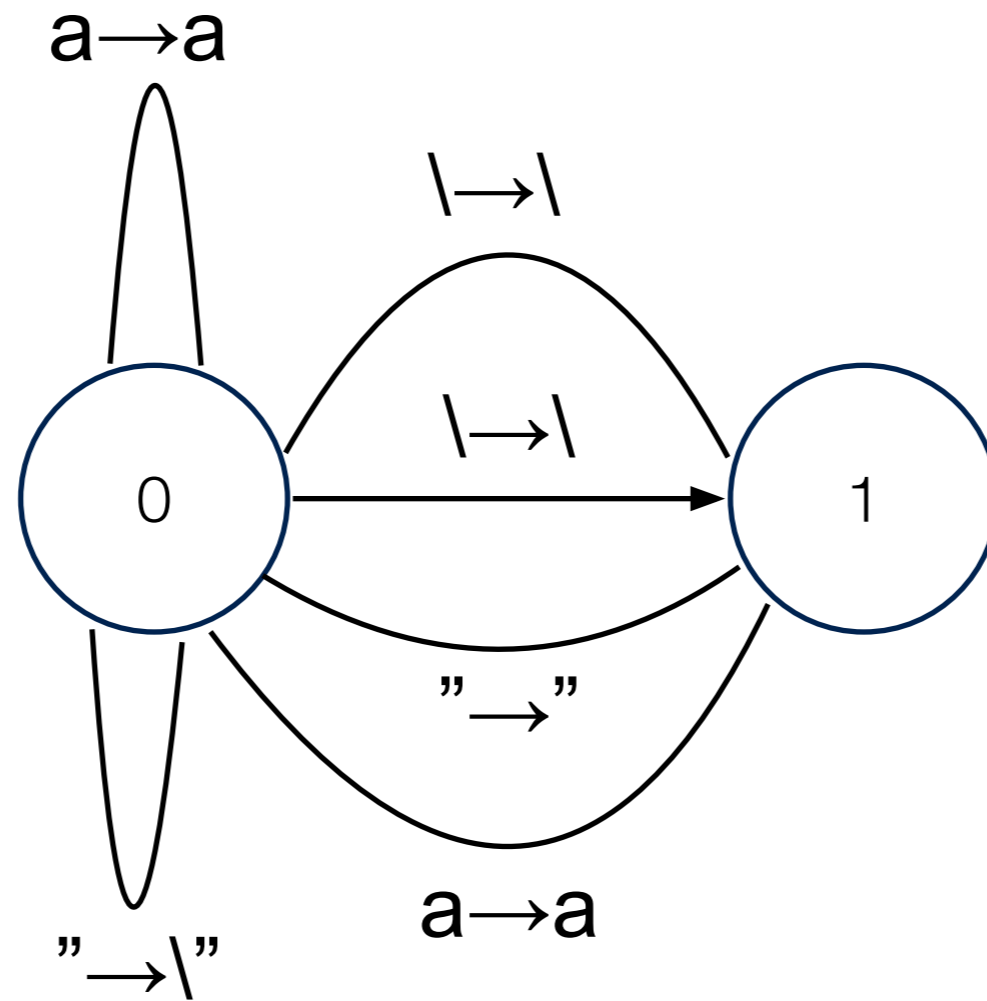
Benchmark	ASTRA (s)	Optician (s)
extrAcronym	0.11	0.05
extrAcronym2	0.42	—
extrNum	0.93	0.05
extrQuant	0.19	0.09
normalizeSpaces	0.46	16.64
extrOdds	15.87	0.12
capProb	0.05	0.05
removeLast	0.21	0.15
sourceToViews	50.92	0.06
normalizeNamePos	X	0.05
titleConverter	X	0.07
bibtexToReadable	X	0.64



Benchmark	ASTRA (s)	Optician (s)
extrAcronym	0.11	0.05
extrAcronym2	0.42	—
extrNum	0.93	0.05
extrQuant	0.19	0.09
normalizeSpaces	0.46	16.64
extrOdds	15.87	0.12
capProb	0.05	0.05
removeLast	0.21	0.15
sourceToViews	50.92	0.06
normalizeNamePos	X	0.05
titleConverter	X	0.07
bibtexToReadable	X	0.64



>10 states in the input automaton



escapeQuotes: replace every " with \"

Optician cannot synthesize this

Working with unstructured data

regexes I , O for Input-Output Types



Automata



Back to regex

$[A-Z][a-z]^*$



$((([A-Z] \cdot [a-z]) \cdot [a-z]^*) \mid [A-Z])$

Benchmark	ASTRA (s)	Optician (s)	Optician-re (s)
extrAcronym	0.11	0.05	X
extrAcronym2	0.42	—	—
extrNum	0.93	0.05	0.07
extrQuant	0.19	0.09	X
normalizeSpaces	0.46	16.64	X
extrOdds	15.87	0.12	X
capProb	0.05	0.05	X
removeLast	0.21	0.15	0.07
sourceToViews	50.92	0.06	X
normalizeNamePos	X	0.05	0.10
titleConverter	X	0.07	X
bibtexToReadable	X	0.64	0.15

We find T such

$\{I\} T \{O\}$ that:

Types

$$T(i_1) = o_1$$

...

$$T(i_n) = o_n$$

Examples

$\forall i \in I. \text{dist}(i, T(i)) \leq d$ **Distance**

Summary

Astra synthesizes:

- Finite Transducers
- Symbolic Transducers
- Transducers with Lookahead

Future work:

- Scalability
- Invariants for string manipulating programs