

# A Decision Procedure for Subset Constraints over Regular Languages

Pieter Hooimeijer  
Westley Weimer

University of Virginia

# Motivation

“String values have lost  
their innocence...”  
[Thiemann05]

# An Example Program

```
// v1 and v2 are user inputs

if (!ereg('o(pp)+', v1)) {exit;}
if (!ereg('p*q', v2)) {exit;}

v3 = v1 . v2; // concat
if (v3 != 'oppqq') {exit;}
```

# An Example Program

```
// v1 and v2 are user inputs
```

```
1 if (!ereg('o(pp)+', v1)) {exit;}
```

```
2 if (!ereg('p*q', v2)) {exit;}
```

```
v3 = v1 . v2; // concat
```

```
3 if (v3 != 'oppqq') {exit;}
```

# An Example Program

$$\begin{array}{lll} c_1 = L(o(pp)+) & \textcircled{1} & v_1 \subseteq c_1 \\ c_2 = L(p*q) & \textcircled{2} & v_2 \subseteq c_2 \\ c_3 = L(oppqq) & \textcircled{3} & v_1 \circ v_2 \subseteq c_3 \end{array}$$

# Some Applications

- Automated randomized **testing**
- Abstraction refinement
- **Fault localization**
- Optimization & refactoring

# Some Questions

- Is the theory **expressive** enough?
- Does it apply to **interesting** problems?
- Is the decision procedure **fast** enough?

# Outline

- Problem definition
- Algorithm
  - example execution
  - correctness
- Experiment



# Subset Constraints

$S ::= E \subseteq C$

$E ::= E \circ E$

$C$

constants

$V$

variables

# Subset Constraints

$S ::= E \subseteq C$

$E ::= E \circ E$

|  
|

$C$

$V$

concatenation

# Regular Matching Assignments

Definition: Let  $P = \{s_1, \dots, s_p\}$  be a system of **subset constraints** over a set of shared variables  $V$ . Find all assignments that simultaneously satisfy each constraint  $s_1, \dots, s_p$

# Algorithm

Input

$$\begin{array}{l} \textcircled{1} \quad v_1 \subseteq c_1 \\ \textcircled{2} \quad v_2 \subseteq c_2 \\ \textcircled{3} \quad v_1 \circ v_2 \subseteq c_3 \end{array}$$

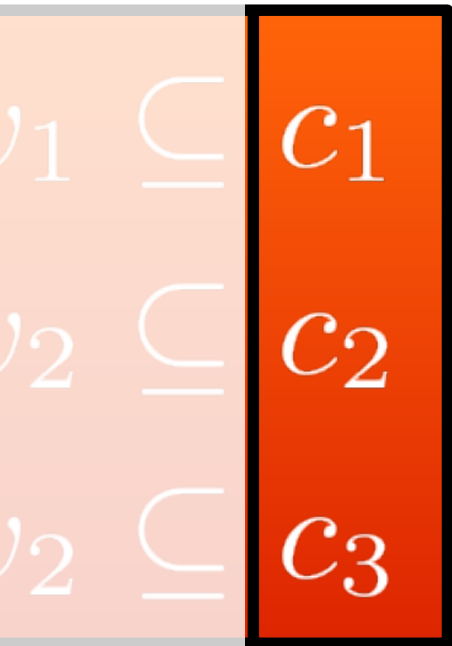
# Algorithm

Input

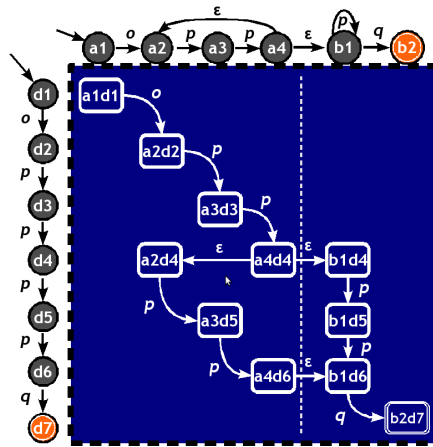
①	$v_1$	$\subseteq$	$C_1$
②	$v_2$	$\subseteq$	$C_2$
③	$v_1 \circ v_2$	$\subseteq$	$C_3$

# Algorithm

Input



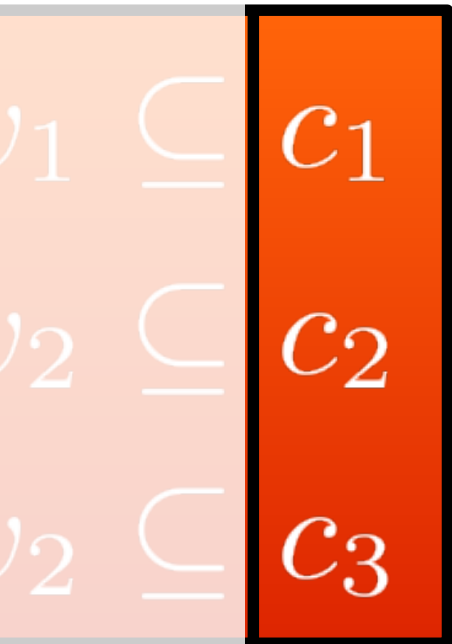
Cross Product



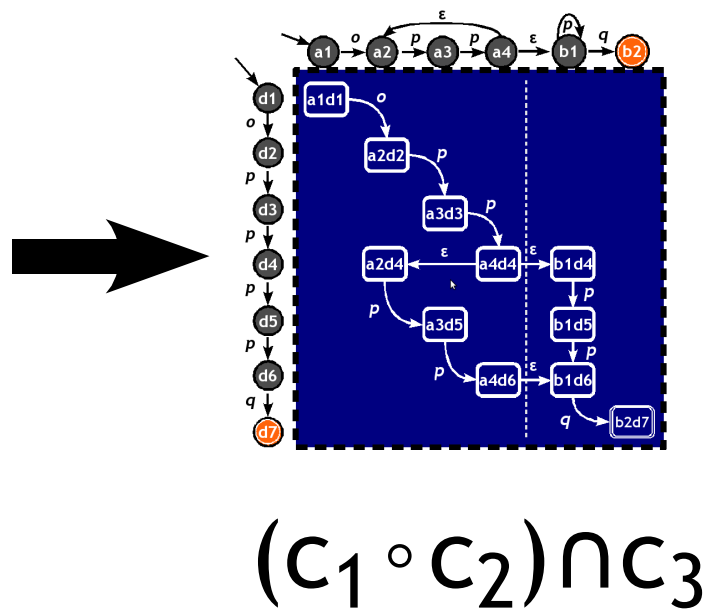
$$(C_1 \circ C_2) \cap C_3$$

# Algorithm

Input



Cross Product



✓ Sat.

✗ Unsat.

# An Example

Solution I:

$$v_1 = \{ opp \}$$

$$v_2 = \{ ppq \}$$



Solution II:

$$v_1 = \{ opppp \}$$

$$v_2 = \{ q \}$$

$$c_1 = L(o(pp)+)$$

$$c_2 = L(p*q)$$

$$c_3 = L(oppqp)$$

$$v_1 \subseteq c_1$$

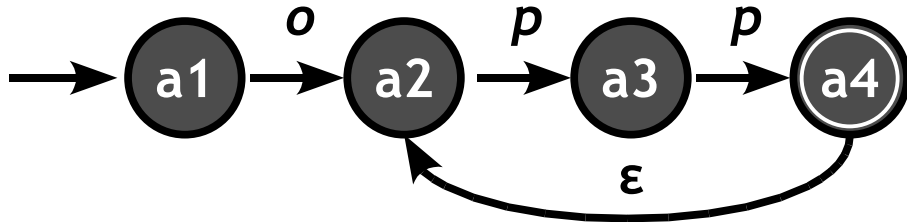
$$v_2 \subseteq c_2$$

$$v_1 \circ v_2 \subseteq c_3$$

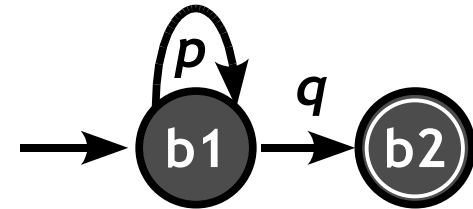


# An Example

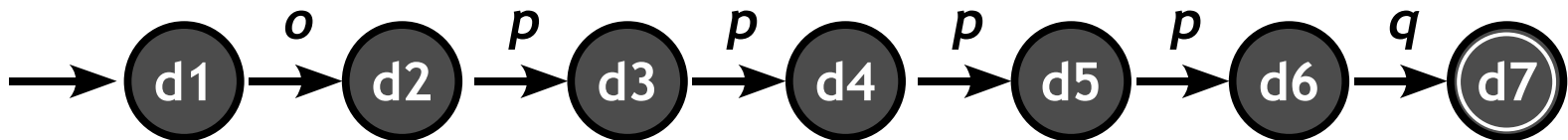
$$c_1 = L(o(pp)^+)$$



$$c_2 = L(p^*q)$$

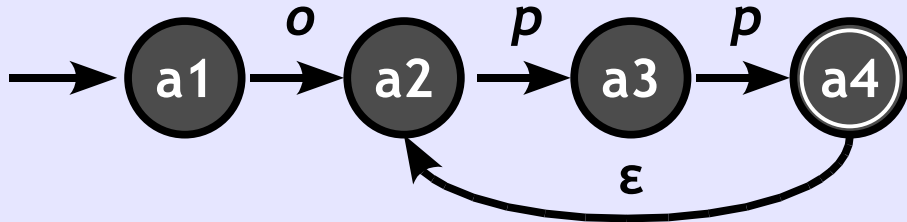


$$c_3 = L(op^4q)$$

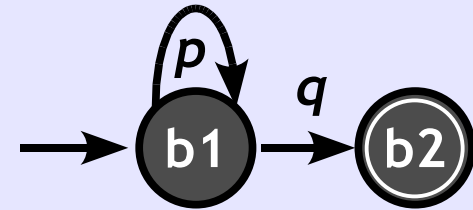


# An Example

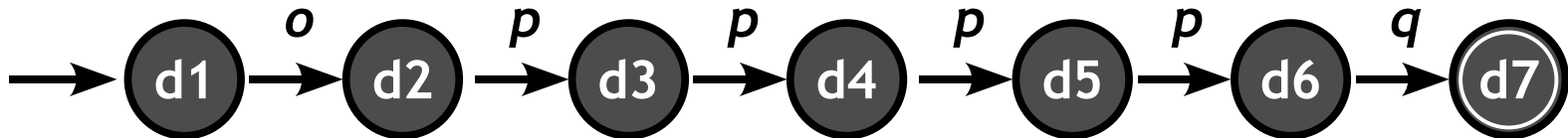
$$c_1 = L(o(pp)^+)$$



$$c_2 = L(p^*q)$$

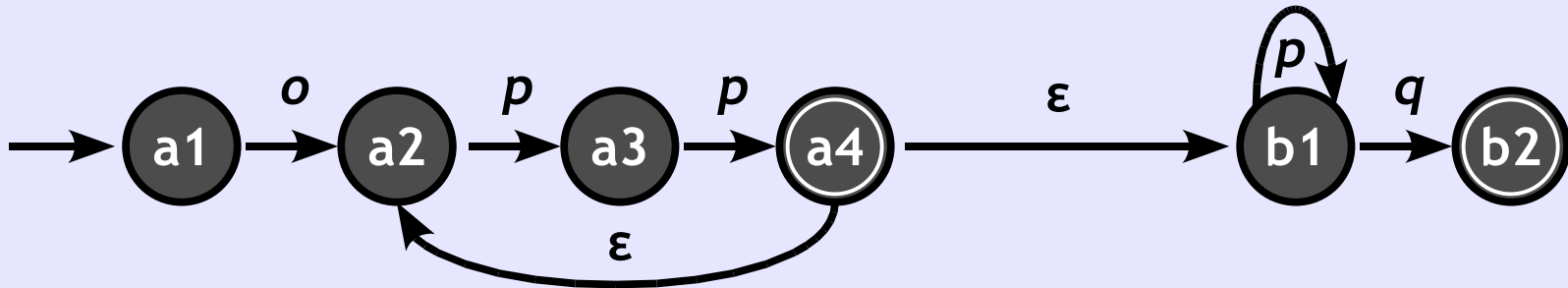


$$c_3 = L(op^4q)$$



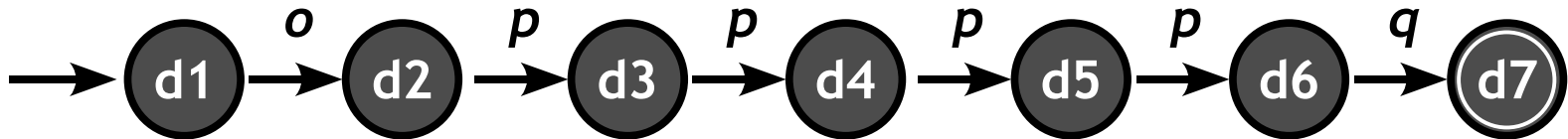
# An Example

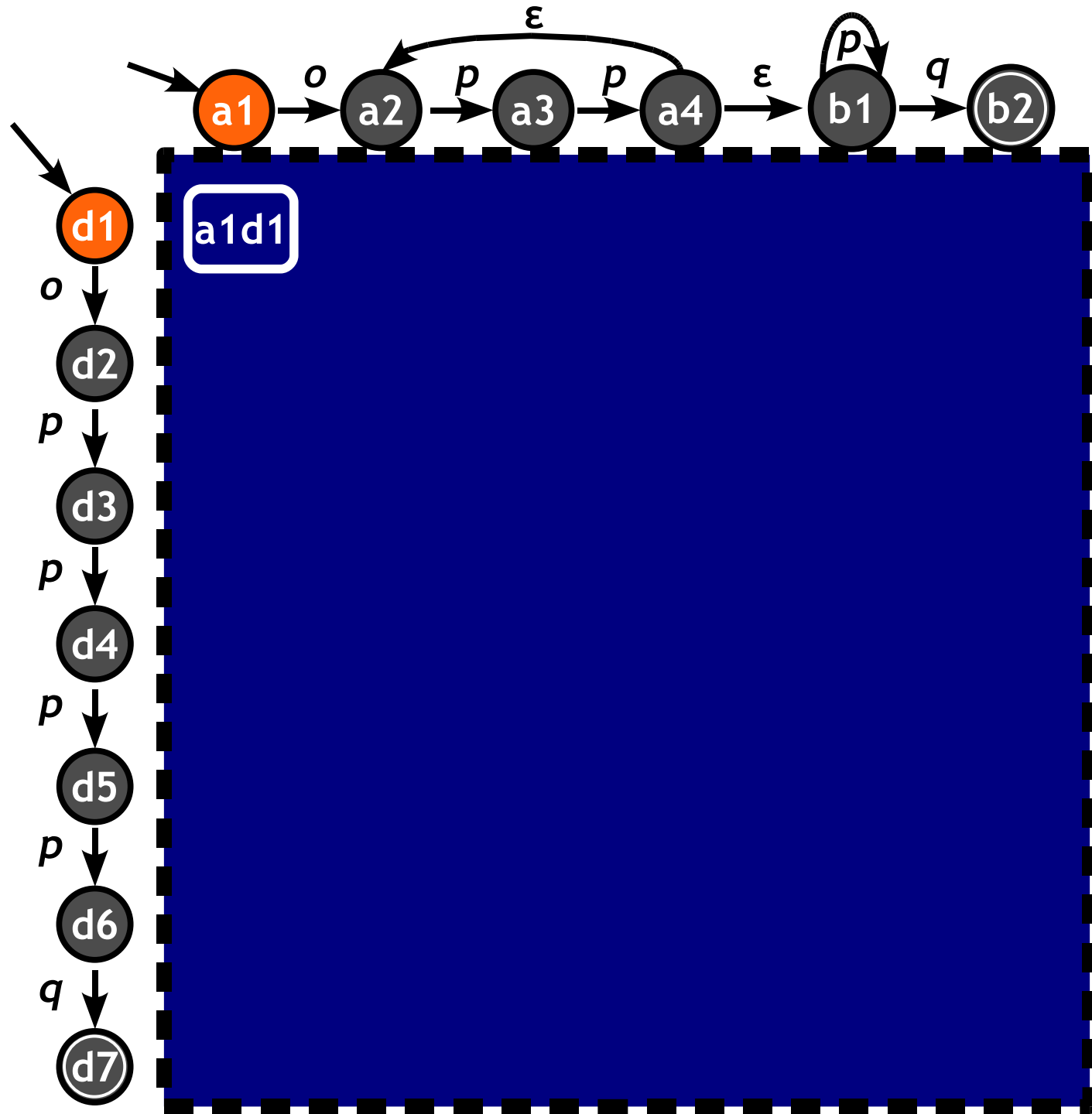
$$c_1 = L(o(pp)^+)$$

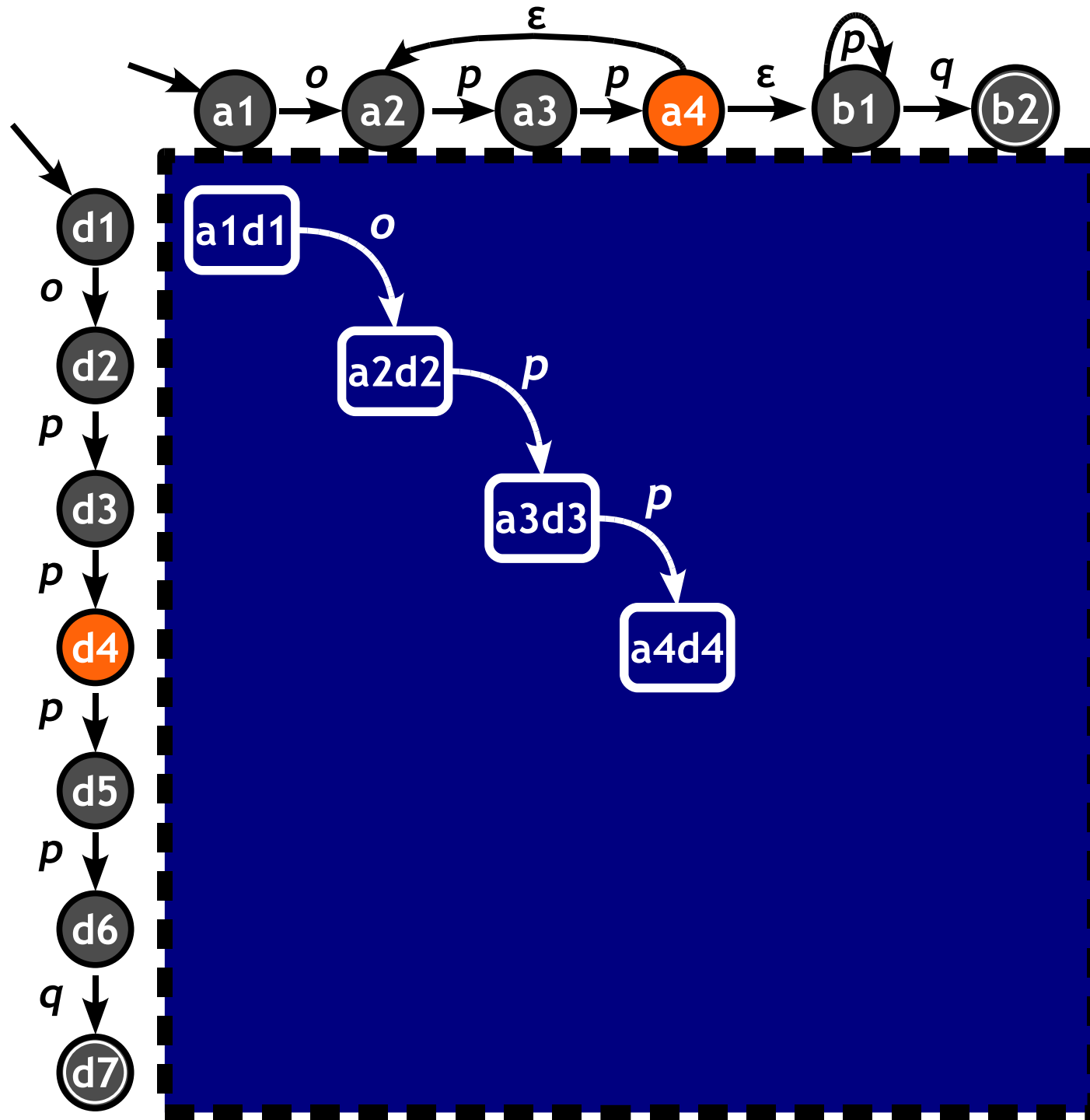


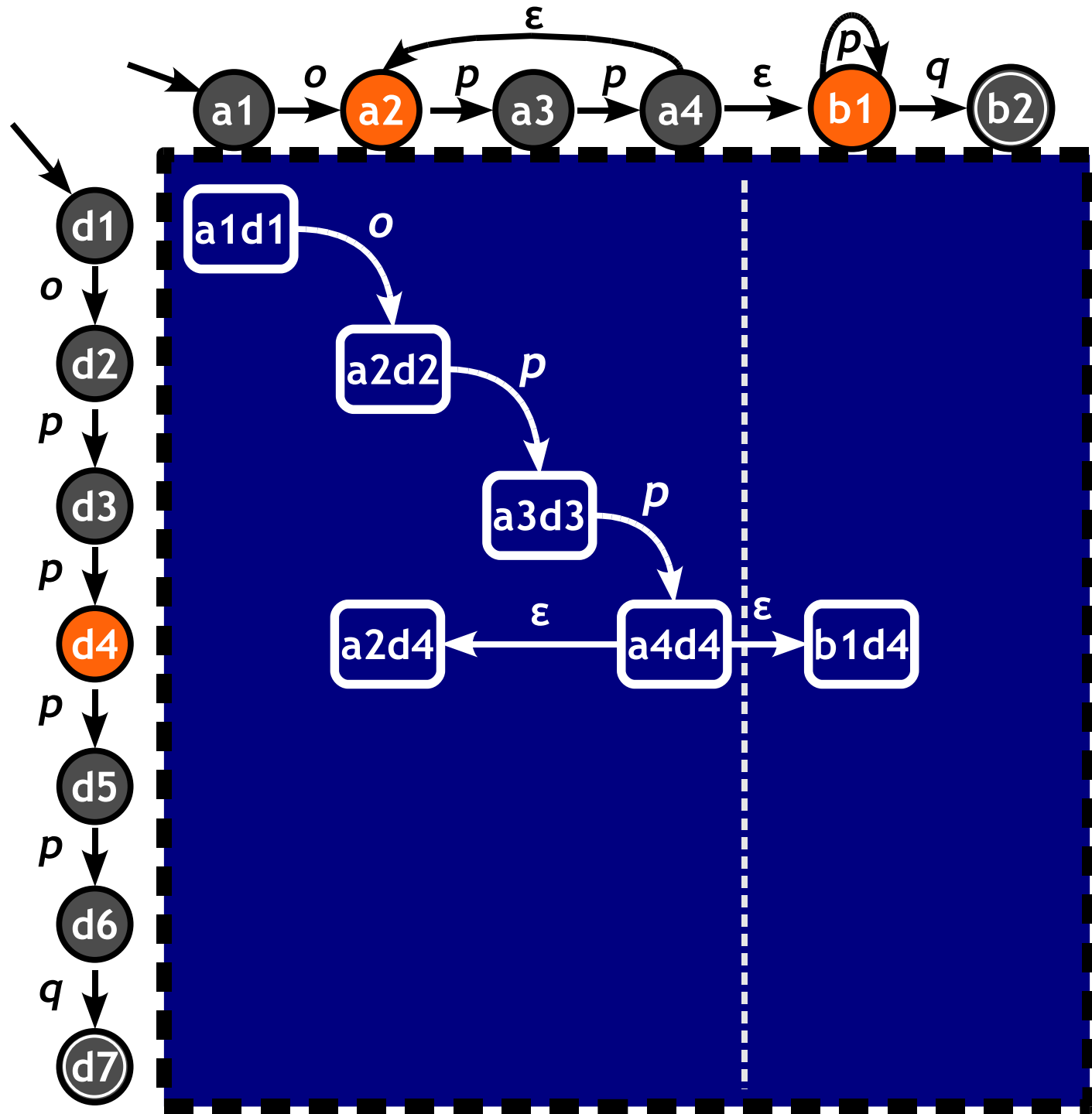
$$c_2 = L(p^*q)$$

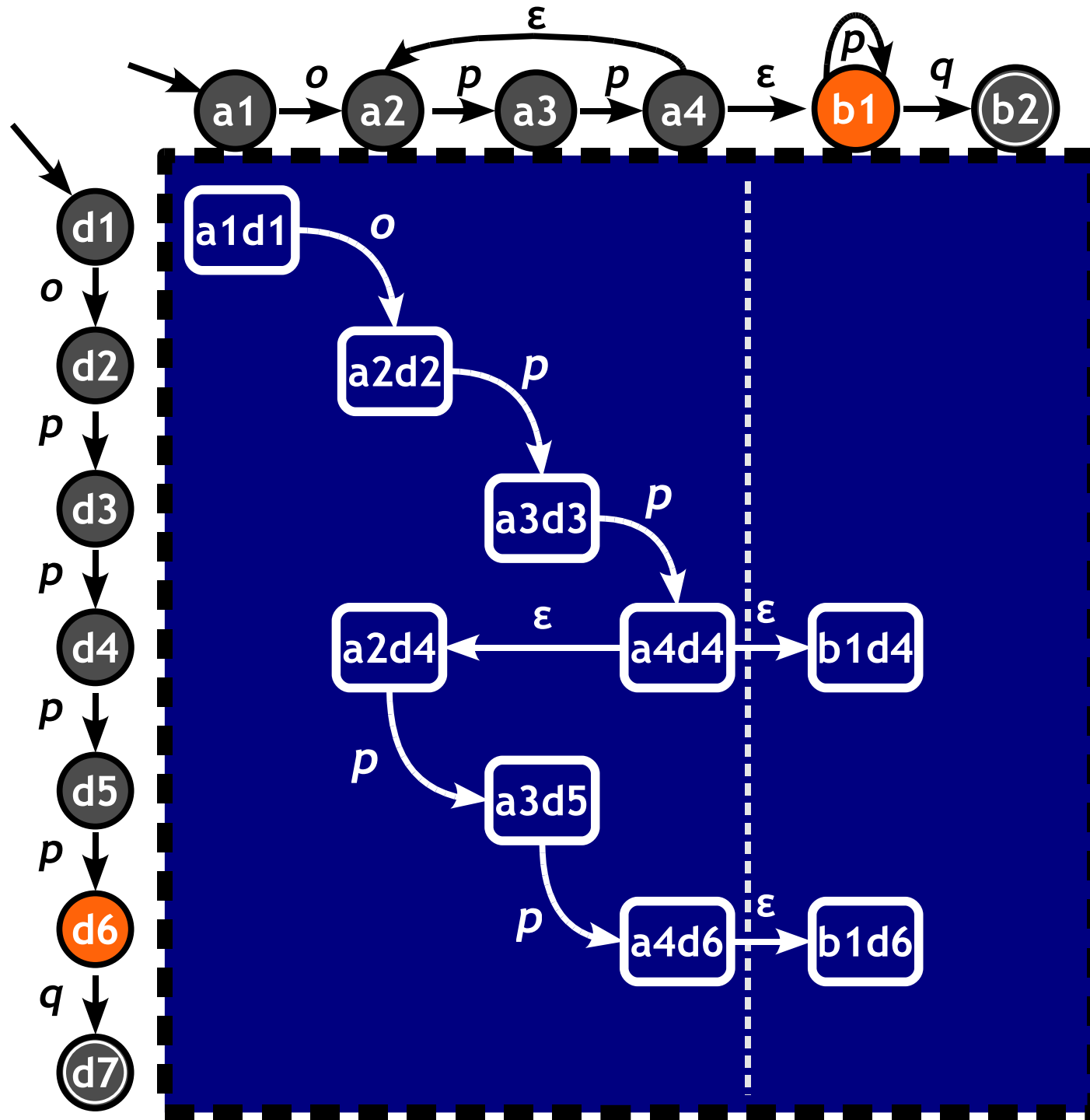
$$c_3 = L(op^4q)$$

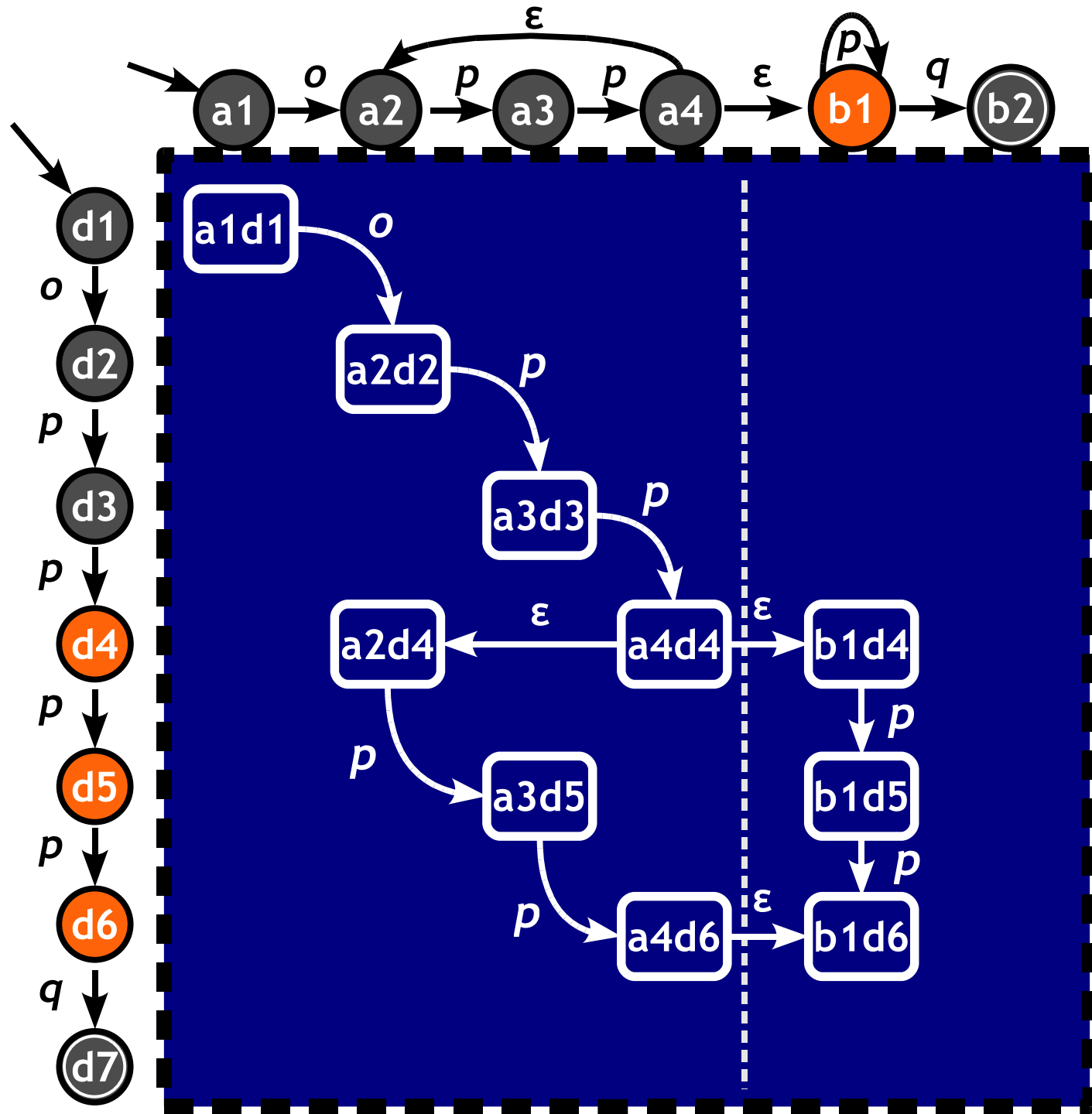




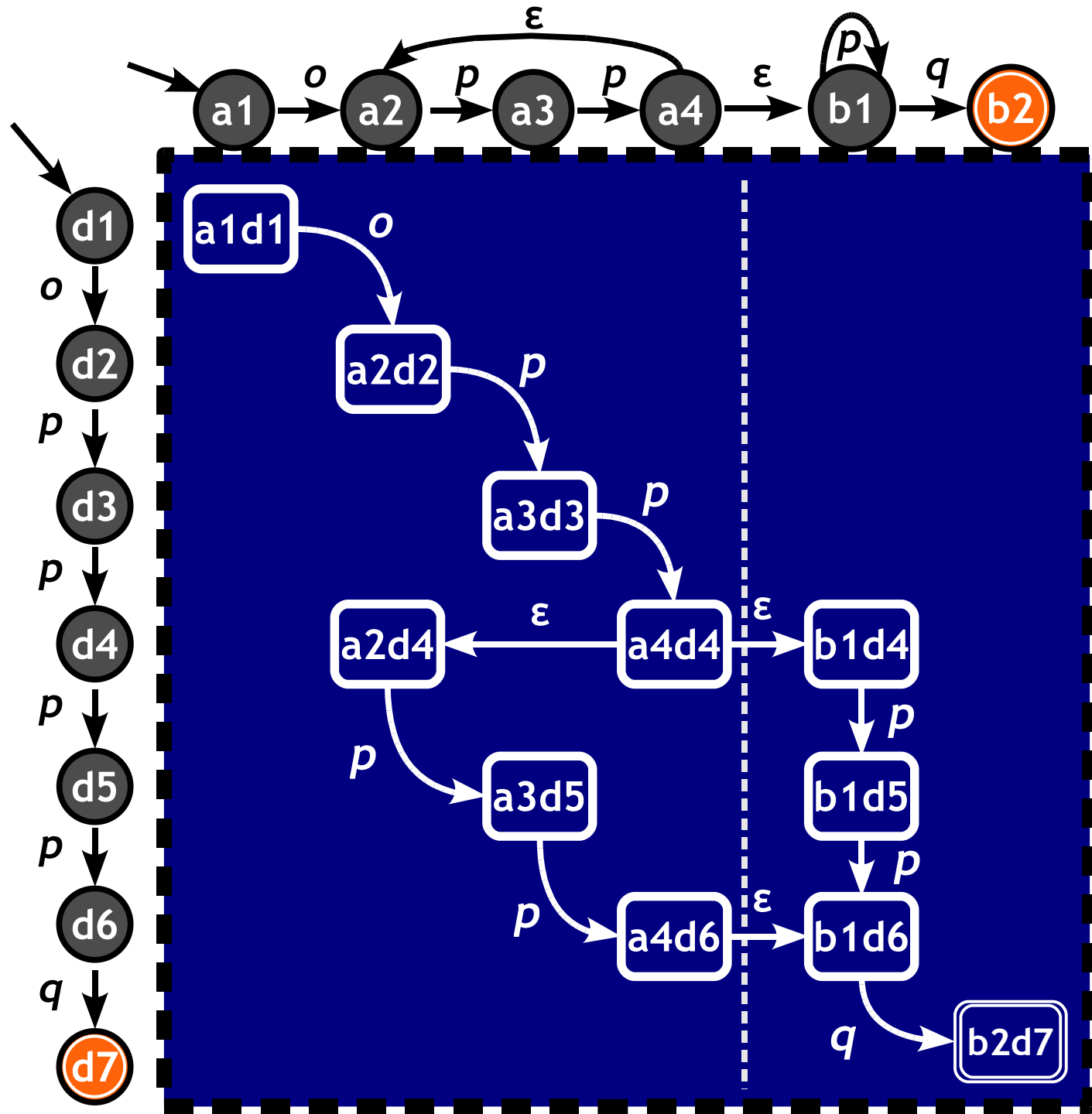


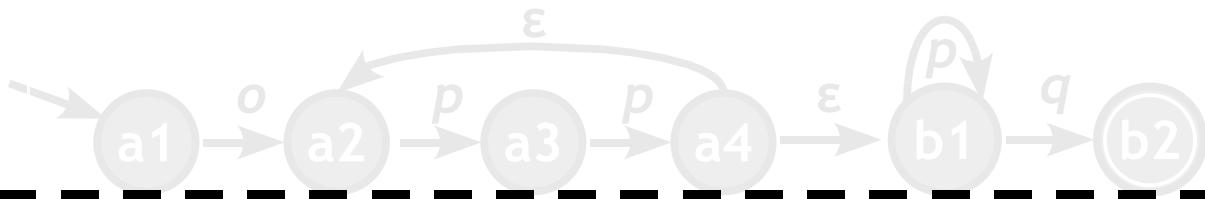




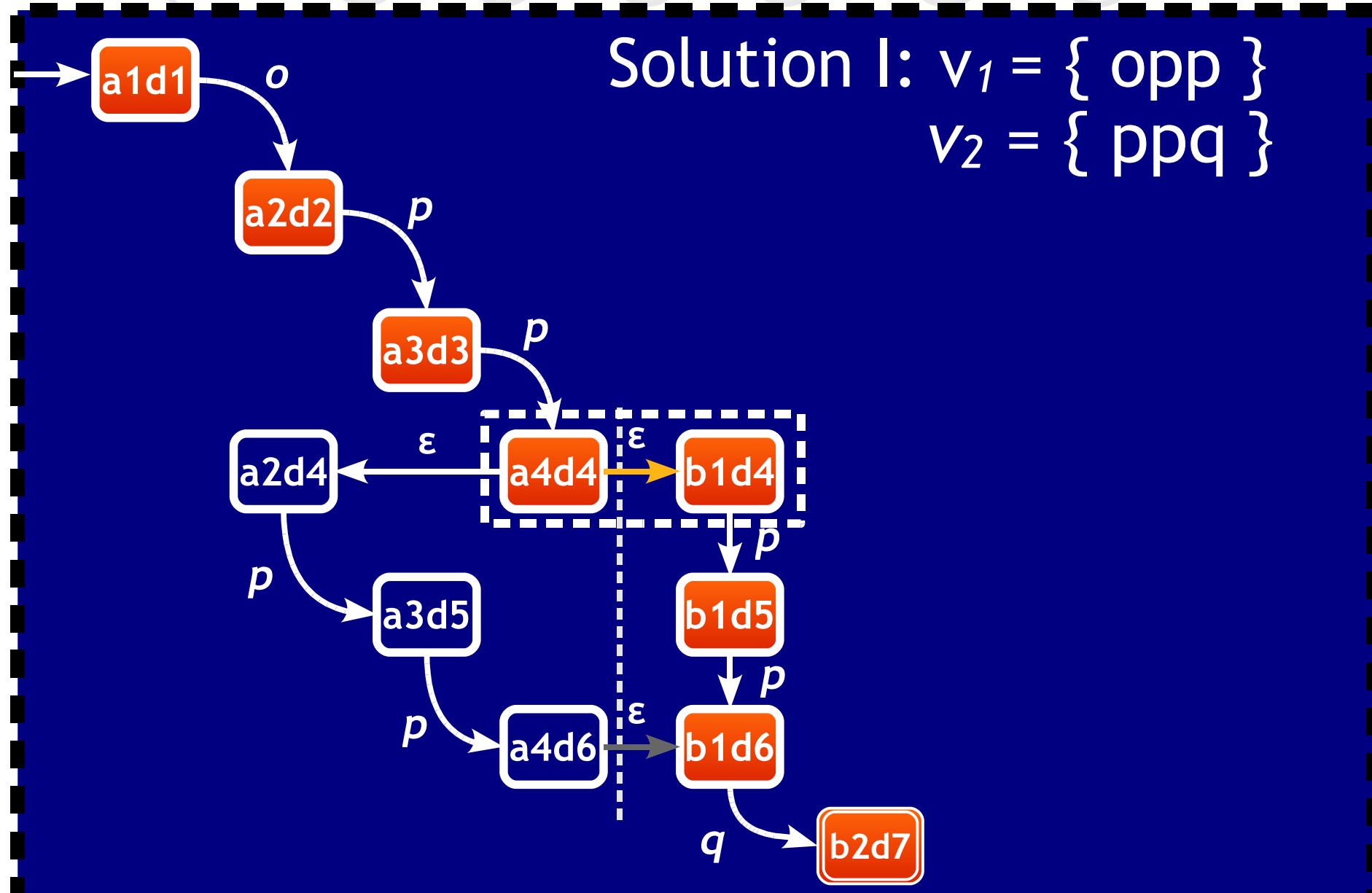


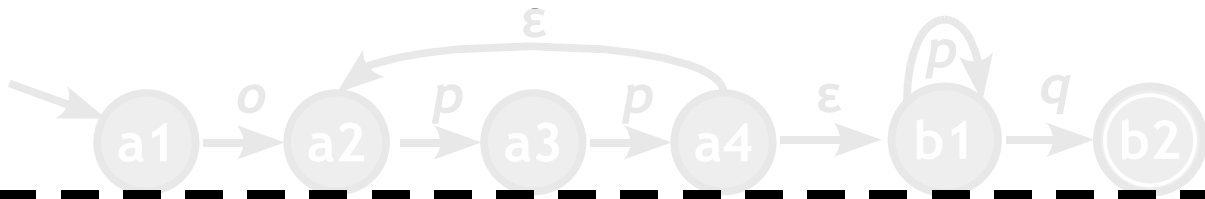






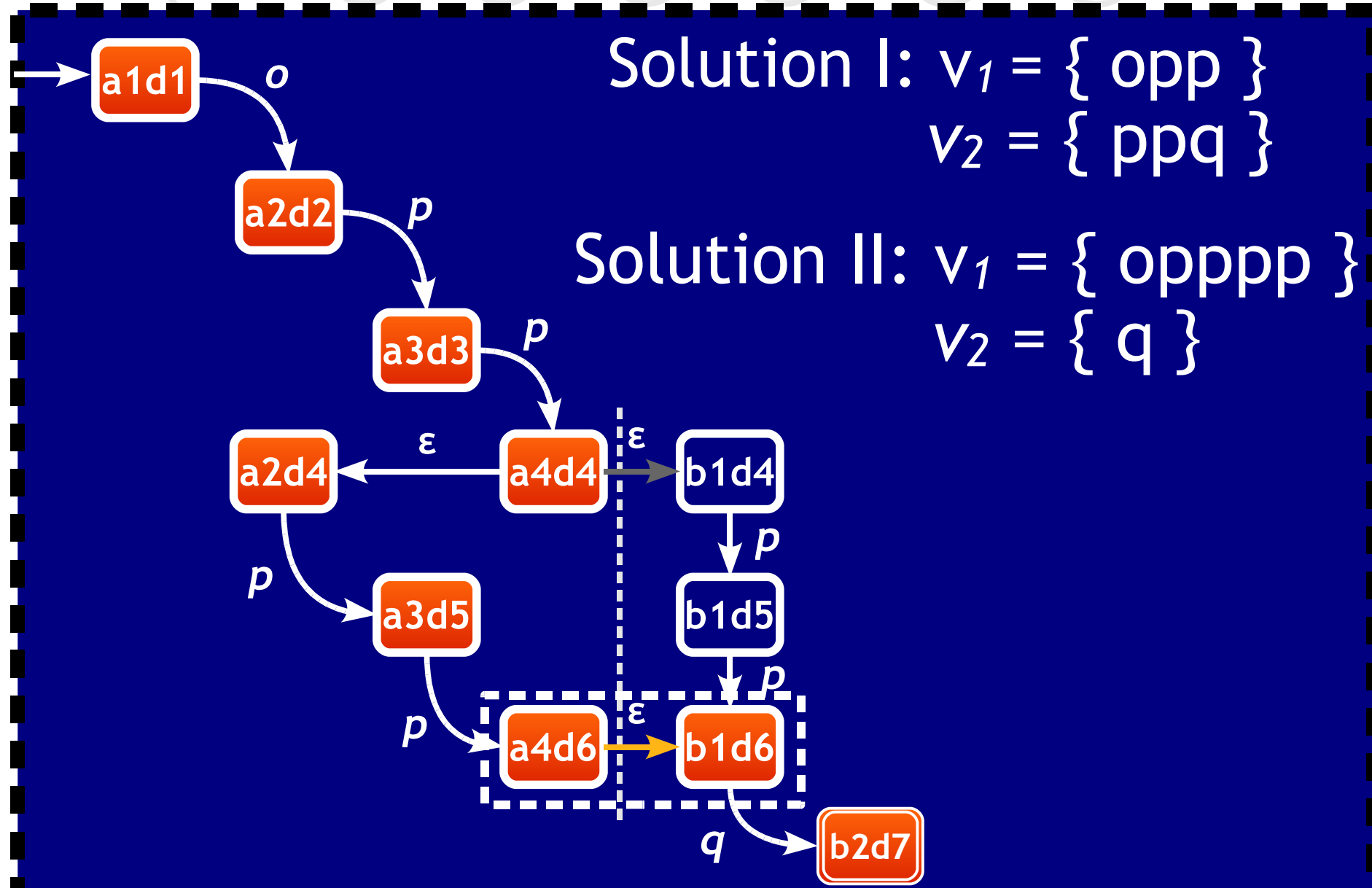
Solution I:  $v_1 = \{ opp \}$   
 $v_2 = \{ ppq \}$





Solution I:  $v_1 = \{ opp \}$   
 $v_2 = \{ ppq \}$

Solution II:  $v_1 = \{ opppp \}$   
 $v_2 = \{ q \}$



# The Algorithm

- **Concat-Intersect (CI) algorithm:**
  - two variables, three constants; fixed form
  - mechanically verified proof in Coq 8.1pl3
  - proof size is ~1300 lines
- **Regular Matching Assignments (RMA):**
  - implemented in a tool, **DPrle**
  - applies *CI* procedure inductively

# Outline

- Problem definition
- Algorithm
  - example execution
  - correctness
- **Experiment**

# Experimental Design

- Find SQL injection vulnerabilities  
[Wassermann and Su; PLDI07]
- For each vulnerability:
  - generate bad SQL + **path**
  - check path consistency (Simplify)
  - solve string constraints (*this*)

# The Data

Name	LOC	No. Vuln.
eve 1.0	905	1
utopia 1.3.0	5438	4
warp 1.2.1	24365	12

# Example

```
// $userid is untrusted
if (!eregi('[0-9]+', $userid)) {
    exit;
}

$user = query("SELECT * FROM `unp_user`".
              "WHERE userid='$userid'");
```



# Example

```
// $userid is untrusted
if (!eregi('[0-9]+', $userid)) {
    exit;
}

$user = query("SELECT * FROM `unp_user`".
              "WHERE userid='$userid'");
```

'^[0-9]+\$'

	C	T(s)
1	29	0.320
2	16	0.052
3	16	0.001
4	156	0.650
5	102	0.260
6	10	0.054
7	4	0.010
8	10	0.530
9	10	0.063

	C	T(s)
10	31	0.170
11	41	0.430
12	81	577.000
13	10	0.057
14	66	0.220
15	387	0.500
16	10	0.052
17	10	0.180

	C	T(s)
1	29	0.320
2	16	0.052
3	16	0.001
4	156	0.650
5	102	0.260
6	10	0.054
7	4	0.010
8	10	0.530
9	10	0.063

	C	T(s)
10	31	0.170
11	41	0.430
12	81	577.000
13	10	0.057
14	66	0.220
15	387	0.500
16	10	0.052
17	10	0.180

	C	T(s)
1	29	0.320
2	16	0.052
3	16	0.001
4	156	0.650
5	102	0.260
6	10	0.054
7	4	0.010
8	10	0.530
9	10	0.063

	C	T(s)
10	31	0.170
11	41	0.430
12	81	577.000
13	10	0.057
14	66	0.220
15	387	0.500
16	10	0.052
17	10	0.180

# Conclusion

- Presented a theory and decision procedure for modeling string values
- Algorithm relies on basic NFA operations and tracking transitions
- Implementation is available now!

# Future Work

- Finding single string assignments quickly (*cf.* Hampi [ISSTA '09])
- Support for common idioms, including integer indexing
- Applications beyond testcase generation

```
int lastSlash = str.LastIndexOf('/');  
if (lastSlash < 0) { return false; }  
  
var rest = str.Substring(lastSlash + 1);  
if (! rest.Contains("EasyChair"))  
    { return false; }  
  
if (! str.StartsWith("http://"))  
    { return false; }  
  
// ...  
return true;
```

1 `int lastSlash = str.LastIndexOf('/');  
if (lastSlash < 0) { return false; }`

2 `var rest = str.Substring(lastSlash + 1);  
if (! rest.Contains("EasyChair"))  
{ return false; }`

3 `if (! str.StartsWith("http://"))  
{ return false; }`

`// ...  
return true;`



1

$$v_1 \subseteq \Sigma^* /$$

2

$$v_2 \subseteq \Sigma^* \text{EasyChair} \Sigma^*$$

3

$$v_1 \circ v_2 \subseteq \text{http} : // \Sigma^*$$

# Tool: Decision Procedure for Regular Language Equations

**DPrle**

<http://dprle.sf.net/>

<http://www.cs.virginia.edu/~ph4u/dprle/>