

A SPECULATIVE PARALLEL DFA MEMBERSHIP TEST FOR MULTICORE, SIMD AND CLOUD COMPUTING ENVIRONMENT

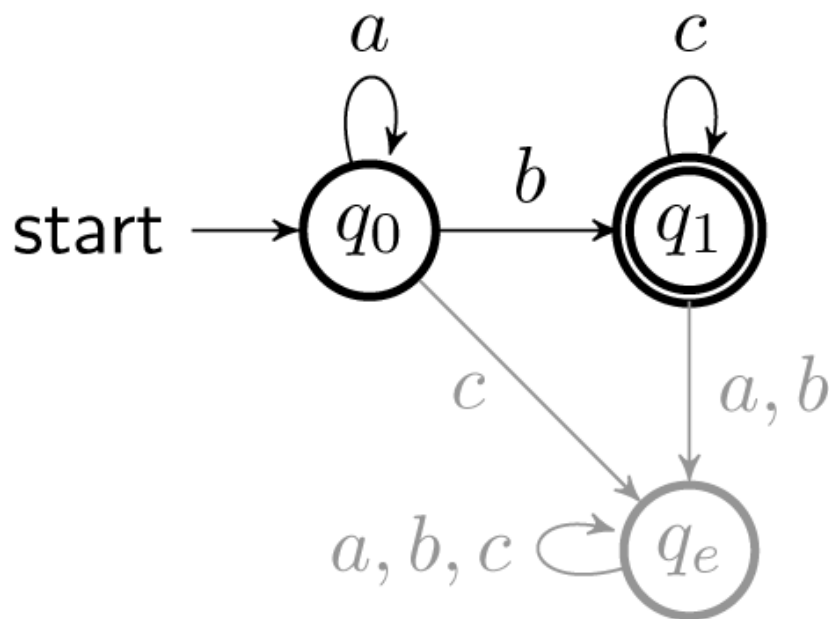
Yousun Ko

Yonsei University, The University of Sydney

Membership Test

2

- DFA: Deterministic Finite Automaton
 - ▣ Mathematical model of computation to test membership of the regular languages



$$L = a^*bc^*$$

$$M = (Q, \Sigma, \delta, q_0, \{q_1\})$$

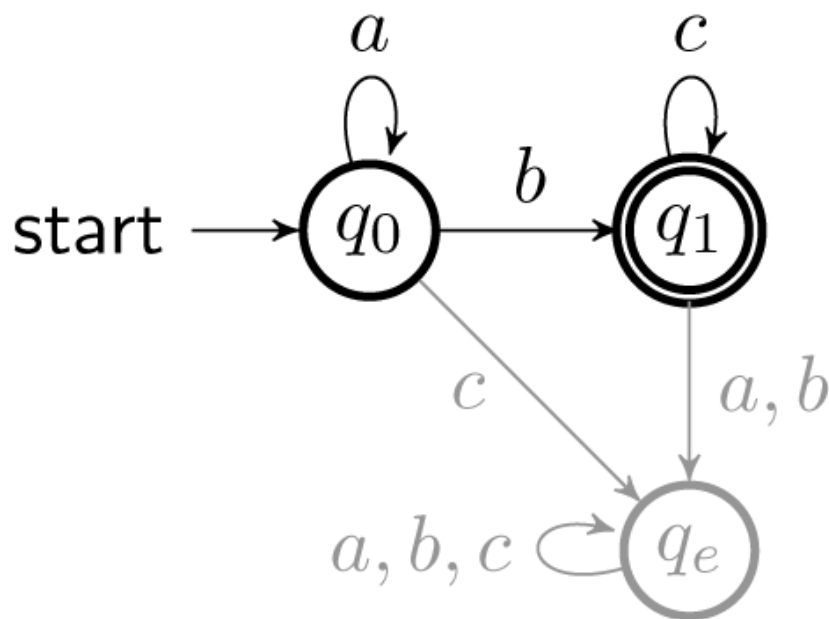
$$Q = \{q_0, q_1\}$$

$$\Sigma = \{a, b, c\}$$

Membership Test

3

- DFA: Deterministic Finite Automaton
 - ▣ Mathematical model of computation to test membership of the regular languages



$$L = a^*bc^*$$

$$M = (Q, \Sigma, \delta, q_0, \{q_1\})$$

$$Q = \{q_0, q_1\}$$

$$\Sigma = \{a, b, c\}$$

$Str_0 = aaaaaaabccc$

$Str_1 = bbaaaaaabccc$

$\in L?$

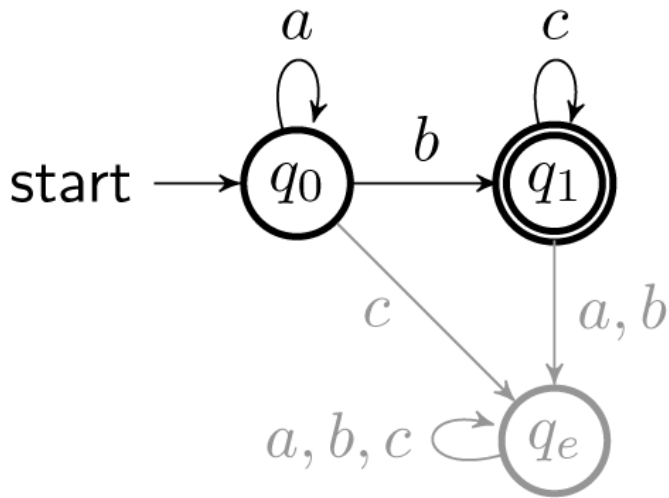
Applications of DFA

4

- Text editing, compiler front-ends, web browsers, scripting languages, file-search, command-processors, databases, internet search engines, computer security, DNA sequence analysis aso...
- Why parallelization of DFA membership test is needed?
 - Execution time is originated by transition of states
 - Execution time is proportional to size of input
 - Long input in general
 - 30,000,000,000 long input for DNA sequences analysis

Parallelization of DFA Membership Test

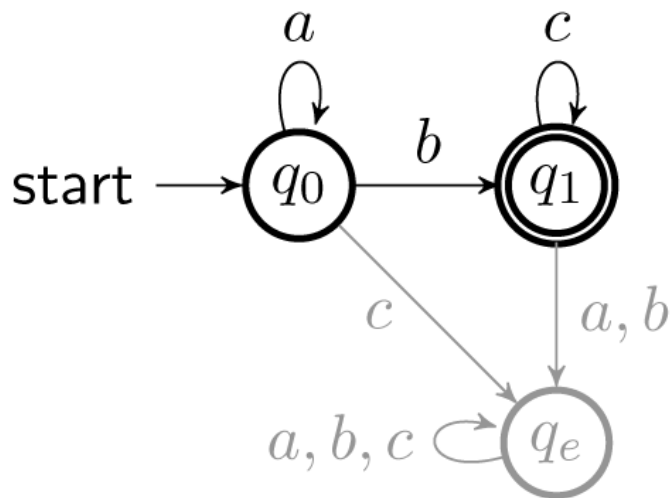
5



$Str = \boxed{a a a a a a a b b c c c}$
 $p_0 : q_0$

Parallelization of DFA Membership Test

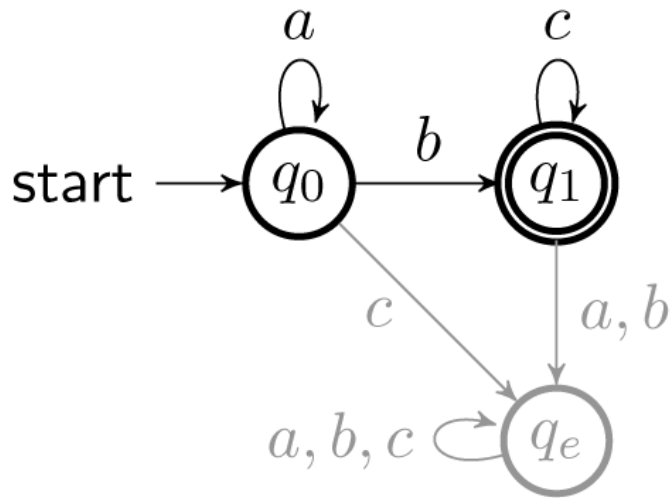
6



$Str = \overset{q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_1 \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ a \ a \ a \ a \ a \ b \ b \ c \ c \ c}} : 12$
 $p_0 : q_0$

Parallelization of DFA Membership Test

7



$Str = \overset{q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_1 \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ a \ a \ a \ a \ b \ b \ c \ c \ c}} : 12$
 $p_0 : q_0$

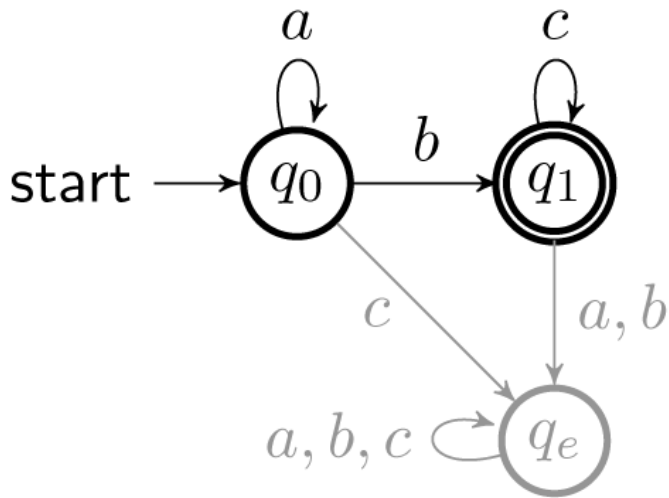
$Str = \overset{chk_0}{\boxed{a \ a \ a \ a}}$
 $p_0 : q_0$

$\overset{chk_1}{\boxed{a \ a \ a \ b}}$
 $p_1 : ?$

$\overset{chk_2}{\boxed{b \ c \ c \ c}}$
 $p_2 : ?$

Parallelization of DFA Membership Test

8



$$Str = \overset{q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_1 \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ a \ a \ a \ a \ b \ b \ c \ c \ c}} : 12$$

$p_0 : q_0$

$$Str = \overset{chk_0}{\boxed{a \ a \ a \ a}}$$

$p_0 : q_0$

$$\overset{chk_1}{\boxed{a \ a \ a \ b}}$$

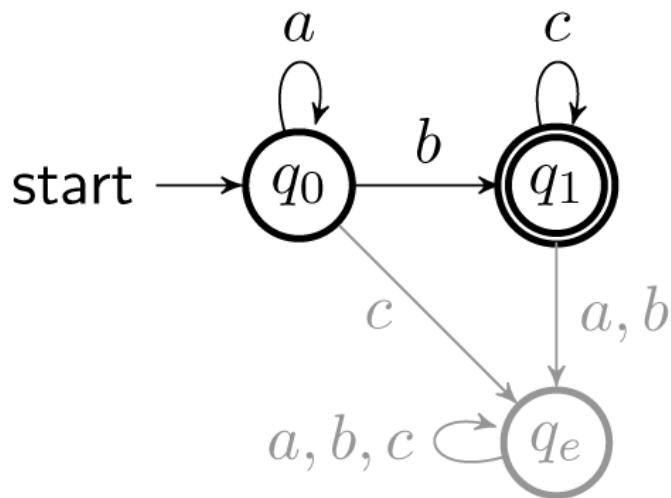
$p_1 : q_0, q_1$

$$\overset{chk_2}{\boxed{b \ c \ c \ c}}$$

$p_2 : q_0, q_1$

Parallelization of DFA Membership Test

9



$Str = \boxed{a a a a a a a b b c c c} : 12$
 $p_0 : q_0$

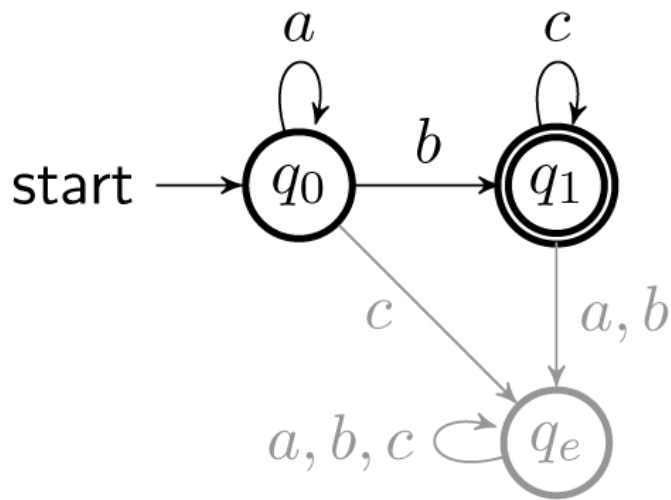
$q_0 q_0 q_0 q_0 q_0 : 4$
 $Str = \boxed{a a a a}$
 $p_0 : q_0$

$q_0 q_0 q_0 q_0 q_0 \parallel q_1 q_e q_e q_e q_e : 8$
 $\boxed{a a a b}$
 $p_1 : q_0, q_1$

$q_0 q_1 q_1 q_1 q_1 \parallel q_1 q_e q_e q_e q_e : 8$
 $\boxed{b c c c}$
 $p_2 : q_0, q_1$

Parallelization of DFA Membership Test

10



$$Str = \boxed{a a a a a a a b b c c c} : 12$$

$p_0 : q_0$

$$Str = \overset{Chk_0}{\boxed{a a a a a a}}$$

$p_0 : q_0$

$$\overset{Chk_1}{\boxed{a b b}}$$

$p_1 : q_0, q_1$

$$\overset{Chk_2}{\boxed{c c c}}$$

$p_2 : q_0, q_1$

Satisfying:

$$1) l_0 = \mathcal{I} \quad l_1 = \mathcal{I} \quad l_2 \quad (\mathcal{I} = |Q| = 2)$$

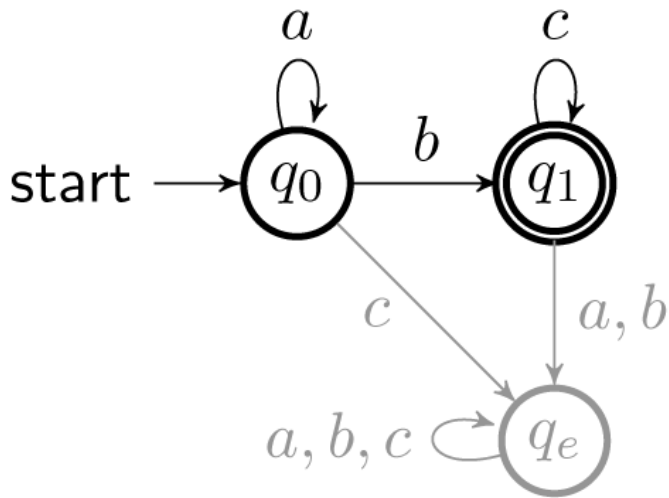
$$2) l_0 + l_1 + l_2 = 12$$

$$3) l_1 = l_2$$

$$l_0 = 6$$

$$l_1 = l_2 = 3$$

Parallelization of DFA Membership Test



$$Str = \boxed{a a a a a a b b c c c} : 12$$

$p_0 : q_0$

$$Str = \overset{q_0 q_0 q_0 q_0 q_0 q_0 q_0}{\boxed{a a a a a a}} : 6$$

$p_0 : q_0$

$$\overset{q_0 q_0 q_1 q_1}{\boxed{a b b}} \parallel \overset{q_1 q_e q_e q_e}{: 6}$$

$p_1 : q_0, q_1$

$$\overset{q_0 q_e q_e q_e}{\boxed{c c c}} \parallel \overset{q_1 q_1 q_1 q_1}{: 6}$$

$p_2 : q_0, q_1$

Satisfying:

$$1) l_0 = \mathcal{I} \quad l_1 = \mathcal{I} \quad l_2 \quad (\mathcal{I} = |Q| = 2)$$

$$2) l_0 + l_1 + l_2 = 12$$

$$3) l_1 = l_2$$

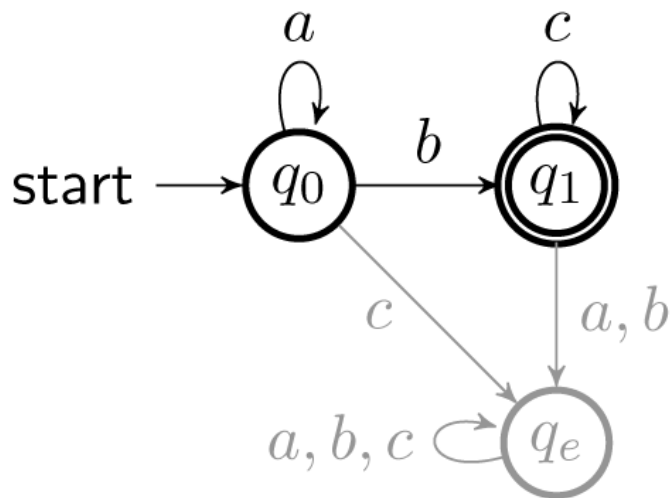
$$l_0 = 6$$

$$l_1 = l_2 = 3$$

Optimization: Reverse Lookahead

12

- Not all the states are candidates of initial possible states.



$Str = \overline{a a a a a a}$ ^{Chk₀}

$p_0 : q_0$

$\overline{a b b}$ ^{Chk₁}

$p_1 : q_0, q_1$

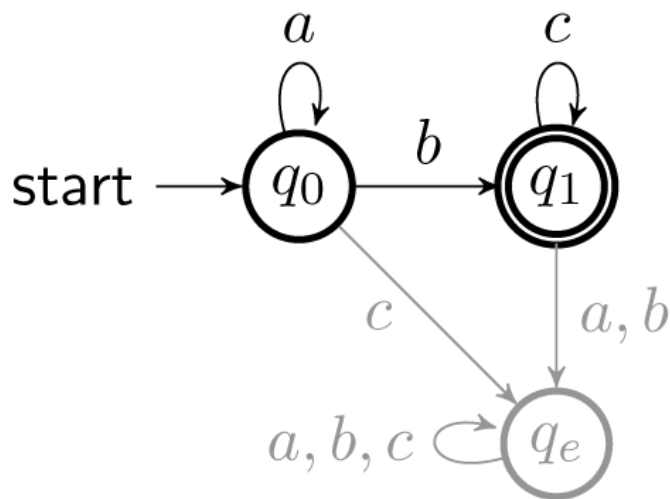
$\overline{c c c}$ ^{Chk₂}

$p_2 : q_0, q_1$

Optimization: Reverse Lookahead

13

- Not all the states are candidates of initial possible states.



$Str =$ $\overline{a a a a a a}$

$p_0 : q_0$

$\overline{a b b}$

$p_1 : q_0, \cancel{q_1}$

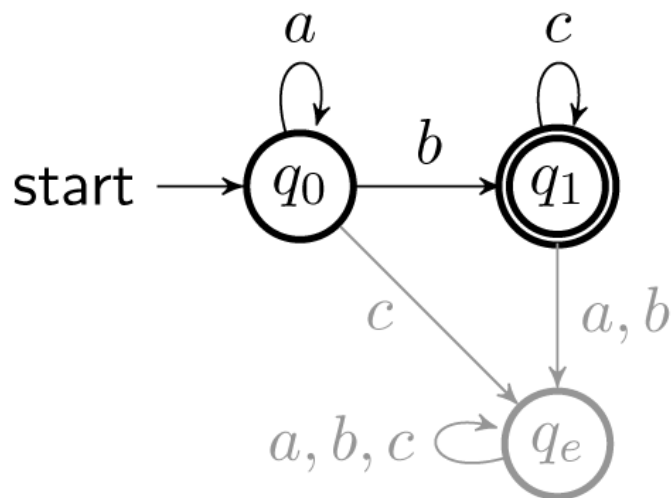
$\overline{c c c}$

$p_2 : \cancel{q_0}, q_1$

Optimization: Reverse Lookahead

14

- Not all the states are candidates of initial possible states.
- More reverse lookahead symbols for smaller \mathcal{I}



$$Str = \begin{array}{l} q_0 \ q_0 \ q_0 \ q_0 \ q_0 : 4 \\ \boxed{a \ a \ a \ a} \\ p_0 : q_0 \end{array}$$

$$\begin{array}{l} q_0 \ q_0 \ q_0 \ q_0 \ q_1 : 4 \\ \boxed{a \ a \ a \ b} \\ p_1 : q_0 \end{array}$$

$$\begin{array}{l} q_1 \ q_e \ q_e \ q_e \ q_e : 4 \\ \boxed{b \ c \ c \ c} \\ p_2 : q_1 \end{array}$$

Satisfying:

- 1) $l_0 = \mathcal{I} \ l_1 = \mathcal{I} \ l_2$ ($\mathcal{I} = 1$)
- 2) $l_0 + l_1 + l_2 = 12$
- 3) $l_1 = l_2$

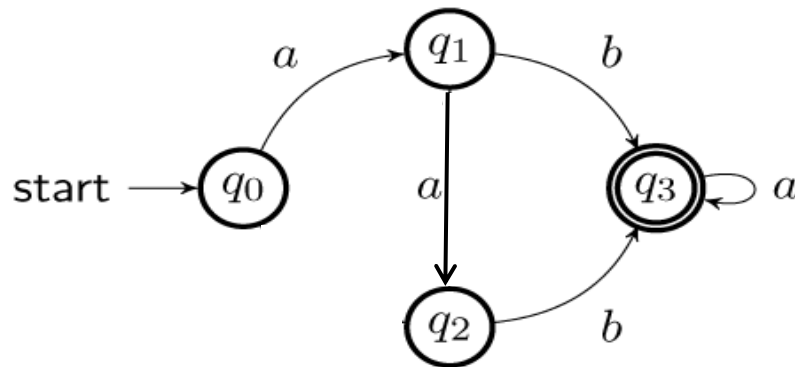
$$l_0 = 4$$

$$l_1 = l_2 = 4$$

Speculation

15

- What is the **expected** number of possible initial states?



$$\mathcal{S}_a = \{q_1, q_2, q_3\}$$

$$\mathcal{S}_b = \{q_3\}$$

$$\mathcal{I}_{min} \leq \mathcal{I} \leq \mathcal{I}_{max}$$

$$\mathcal{I}_{min} = \min(|\mathcal{S}_a|, |\mathcal{S}_b|)$$

$$\mathcal{I}_{max} = \max(|\mathcal{S}_a|, |\mathcal{S}_b|)$$

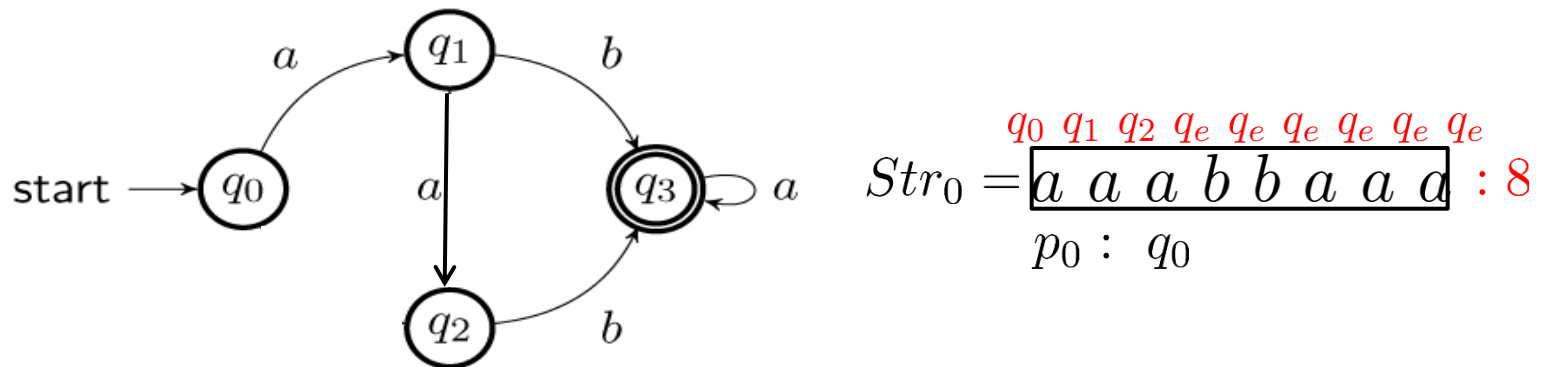
Optimistic Speculation

Safe Speculation

Speculation

16

- What is the **expected** number of possible initial states?



$$\mathcal{I} = |\{q_1, q_2, q_3\}|$$

Safe Speculation

$$Str_0 = \overset{q_0 \ q_1 \ q_2 \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ b \ b \ a}} : 6$$

$p_0 : q_0$

$$\mathcal{I} = |\{q_3\}|$$

Optimistic Speculation

$$Str_0 = \overset{q_0 \ q_1 \ q_2 \ q_e \ q_e}{\boxed{a \ a \ a \ b}} : 4$$

$p_0 : q_0$

$$\overset{q_1 \ q_2 \ q_e}{\boxed{a \ a}} \parallel \overset{q_2 \ q_e \ q_e}{\boxed{a \ a}} \parallel \overset{q_3 \ q_3 \ q_3}{\boxed{a \ a \ a}} : 6$$

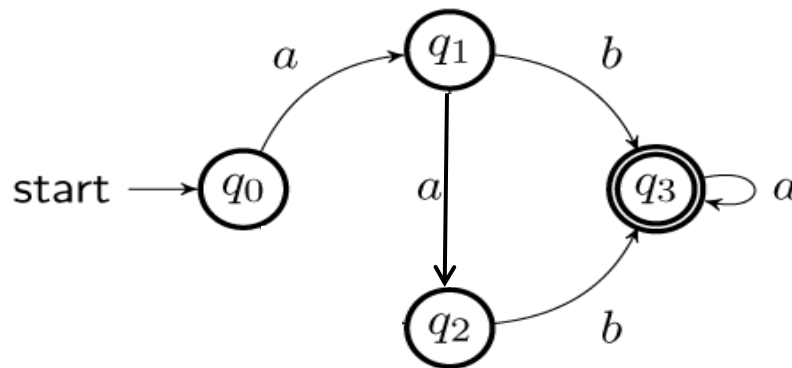
$p_1 : q_1 \ q_2, q_3$

$$\overset{q_3 \ q_e \ q_e \ q_e \ q_e}{\boxed{b \ a \ a \ a}} : 4$$

$p_1 : q_3$

Speculation

- What is the **expected** number of possible initial states?



$$Str_1 = \overset{q_0 \ q_1 \ q_2 \ q_e \ q_e \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ a \ b \ b \ b \ b}} : 8$$

$p_0 : q_0$

$$\mathcal{I} = |\{q_1, q_2, q_3\}|$$

Safe Speculation

$$Str_1 = \overset{q_0 \ q_1 \ q_2 \ q_e \ q_e \ q_e \ q_e}{\boxed{a \ a \ a \ a \ b \ b}} : 6$$

$p_0 : q_0$

$$\overset{q_3 \ q_3 \ q_3}{\boxed{b \ b}} : 2$$

$p_1 : q_3$

$$\mathcal{I} = |\{q_3\}|$$

Optimistic Speculation

$$Str_1 = \overset{q_0 \ q_1 \ q_2 \ q_e \ q_e}{\boxed{a \ a \ a \ a}} : 4$$

$p_0 : q_0$

$$\overset{q_1 \ q_3 \ q_e \ q_e \ q_e}{\boxed{b \ b \ b \ b}} \parallel \overset{q_2 \ q_3 \ q_e \ q_e \ q_e}{\boxed{b \ b \ b \ b}} \parallel \overset{q_3 \ q_e \ q_e \ q_e \ q_e}{\boxed{b \ b \ b \ b}} : 12$$

$p_1 : q_1, q_2, q_3$

Failure!



Merging Local End States

18

$$\mathcal{L}_i = [l_0, l_1, \dots, l_{|Q|-1}],$$

where $0 \leq i < |P|$ and $l_j \in Q$ for all $0 \leq j < |Q|$

$$\text{Str} = \begin{array}{c} q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \\ \boxed{a \ a \ a \ a \ a \ a} \\ p_0 : q_0 \end{array} \quad \mathcal{L}_0 = [0, \phi, \phi]$$

$$\begin{array}{c} q_0 \ q_0 \ q_1 \ q_1 \quad || \quad q_1 \ q_1 \ q_0 \ q_0 \\ \boxed{a \ b \ b} \\ p_1 : q_0, q_1 \end{array} \quad \mathcal{L}_1 = [1, 0, \phi]$$

$$\begin{array}{c} q_0 \ q_e \ q_e \ q_e \quad || \quad q_1 \ q_1 \ q_1 \ q_0 \\ \boxed{c \ c \ c} \\ p_2 : q_0, q_1 \end{array} \quad \mathcal{L}_2 = [\phi, 0, \phi]$$

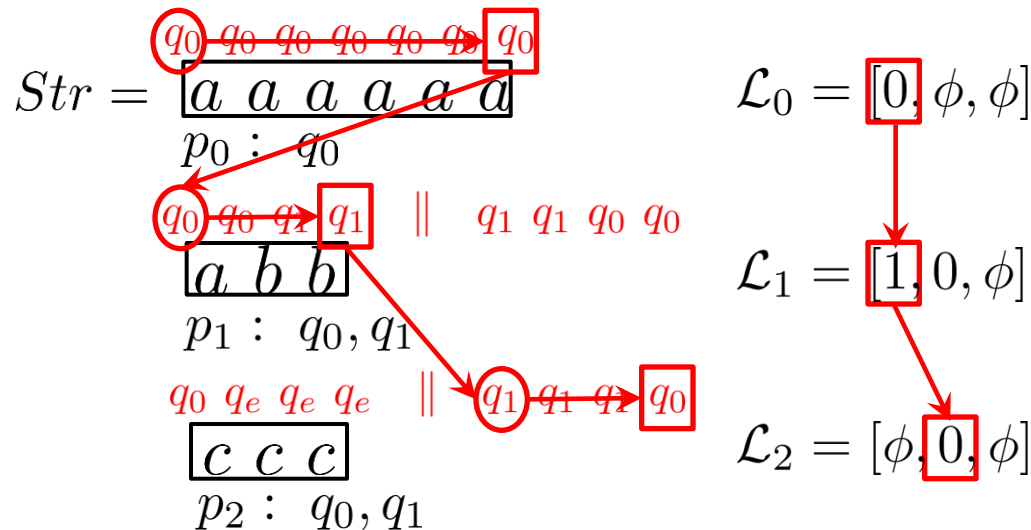
Merging Local End States

19

$$\mathcal{L}_i = [l_0, l_1, \dots, l_{|Q|-1}],$$

where $0 \leq i < |P|$ and $l_j \in Q$ for all $0 \leq j < |Q|$

□ Sequential Merging



Merging Local End States

20

$$\mathcal{L}_i = [l_0, l_1, \dots, l_{|Q|-1}],$$

where $0 \leq i < |P|$ and $l_j \in Q$ for all $0 \leq j < |Q|$

□ Parallel Merging

$$Str = \begin{array}{c} q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \ q_0 \\ \boxed{a \ a \ a \ a \ a \ a} \\ p_0 : q_0 \end{array}$$



$$p_1 : q_0, q_1$$



$$p_2 : q_0, q_1$$

$$\mathcal{L}_0 = [0, \phi, \phi]$$

$$\mathcal{L}_1 = [1, 0, \phi]$$

$$\mathcal{L}_2 = [\phi, 0, \phi]$$

$$\mathcal{L}_{1,2} = [0, \phi, \phi]$$

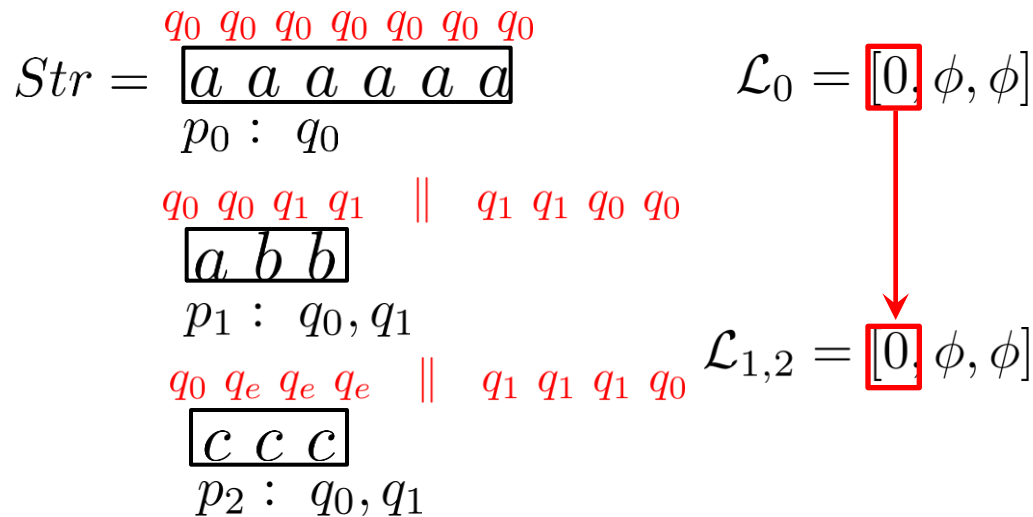
Merging Local End States

21

$$\mathcal{L}_i = [l_0, l_1, \dots, l_{|Q|-1}],$$

where $0 \leq i < |P|$ and $l_j \in Q$ for all $0 \leq j < |Q|$

□ Parallel Merging



Time Complexity

22

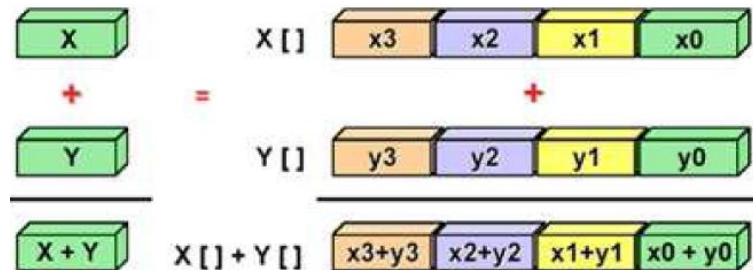
$$\mathcal{O}\left(\frac{nm}{m+p}\right),$$

where $n = |Str|$, $p = |P|$ and m is either $|Q|$ or \mathcal{I} .

Experiments: H/W Overview

23

- Shared-memory multicore
 - Intel Manycore Testing Lab (MTL): $4\text{CPUs} \times 10\text{ cores/CPU} = 40\text{ cores}$
- SIMD (Single Instruction Multiple Data)
 - AVX2 instruction set extension (8-fold vectorization)



- Cluster Computing Environment
 - Amazon Elastic Cluster Computing (EC2) Environment
 - 288 cores

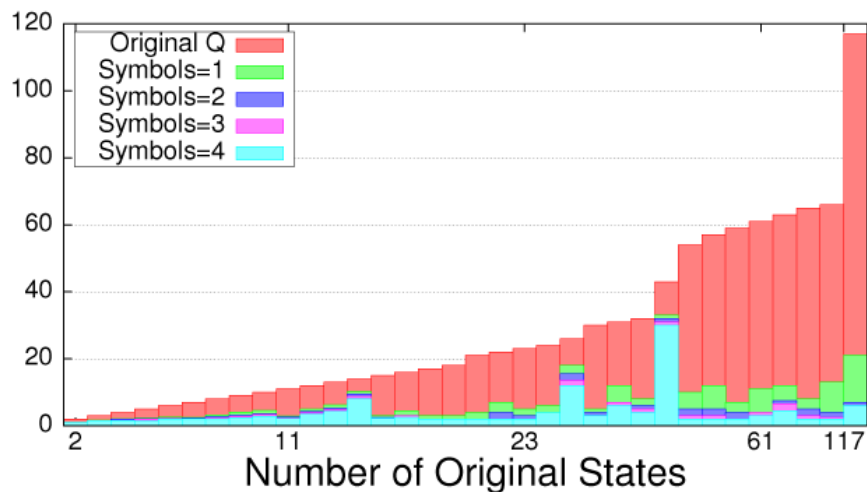
Experiments: Benchmarks

24

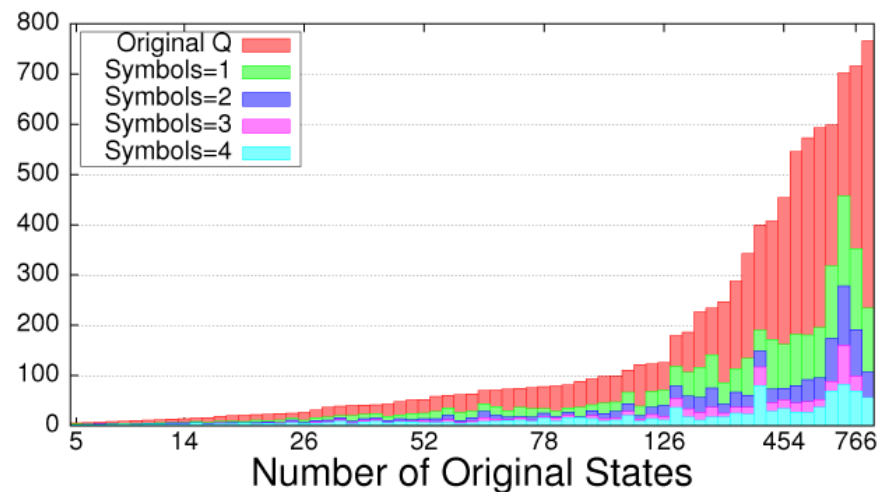
- 299 PCRE (Perl-compatible Regular Expression) RE patterns
- 110 PROSITE protein patterns
- 1MB long input
 - ▣ Longer input results better performance!

Reverse Lookahead Symbols

25



(a) PCRE

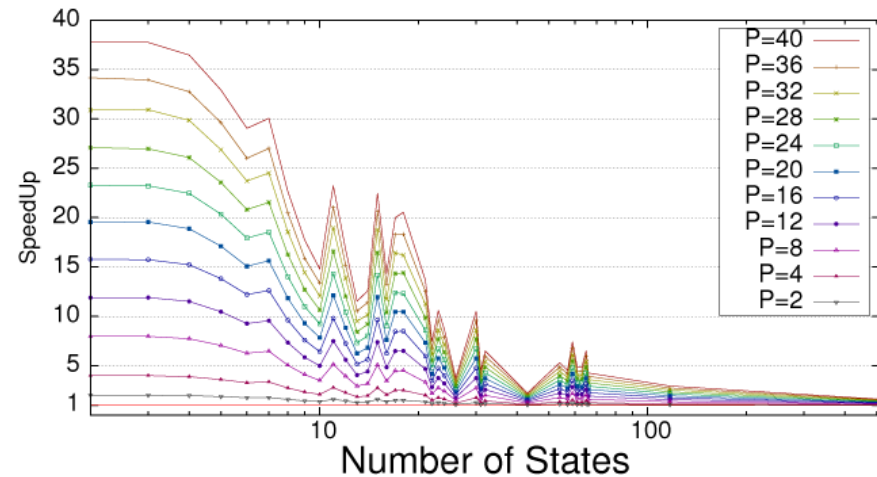
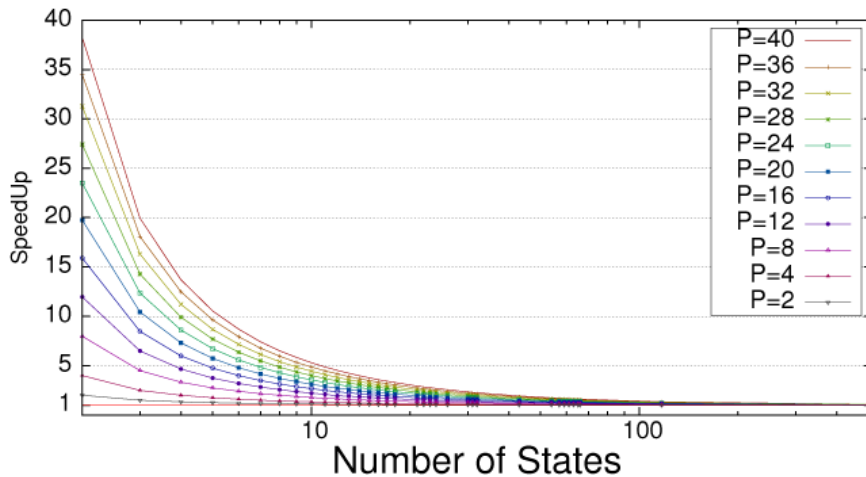


(b) PROSITE

r	0	1	2	3	4
PCRE	100%	33.7%	26.4%	23.7%	21.7%
PROSITE	100%	47.2%	29.2%	20.5%	16.0%

Shared-Memory Multicore

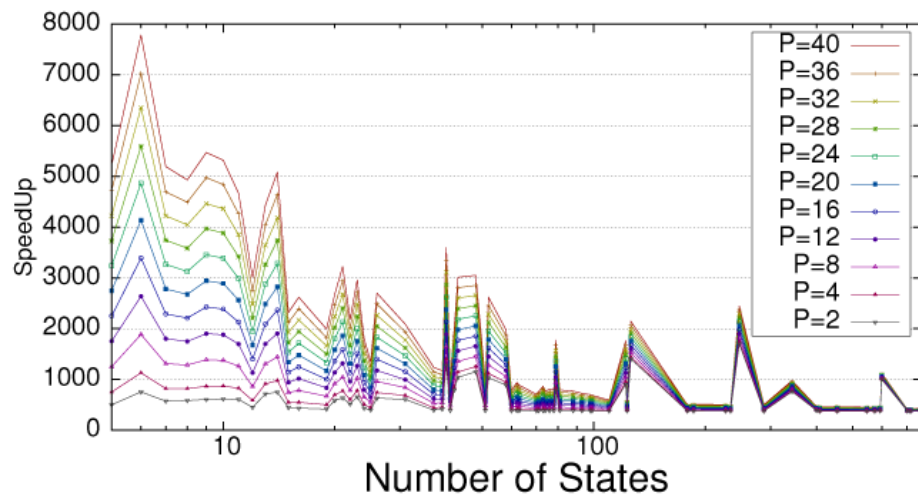
26



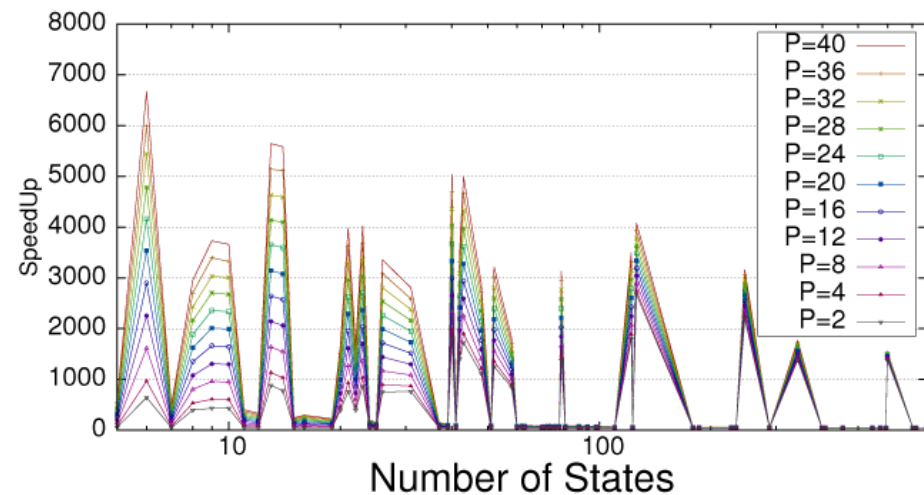
(a) Speedup without optimization for PCRE (b) Speedup with optimization for PCRE

Shared-Memory Multicore

27



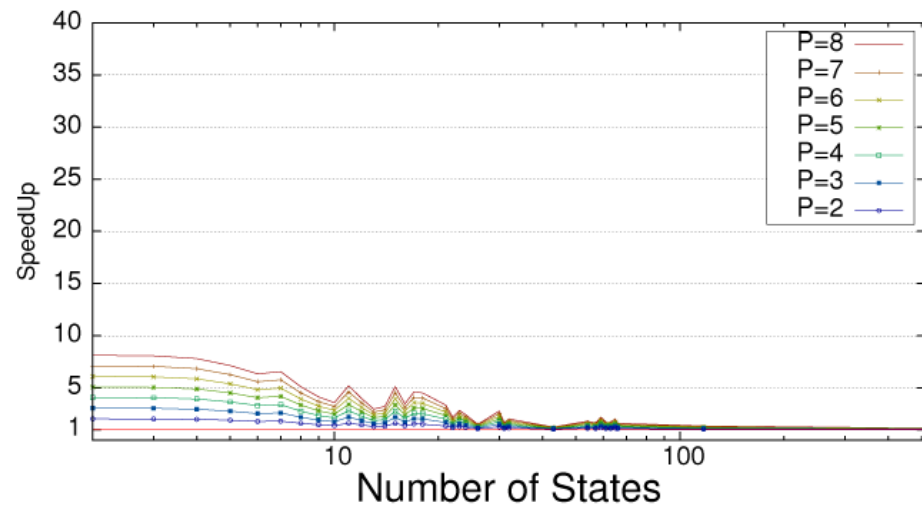
(a) Speedup over ScanProsite



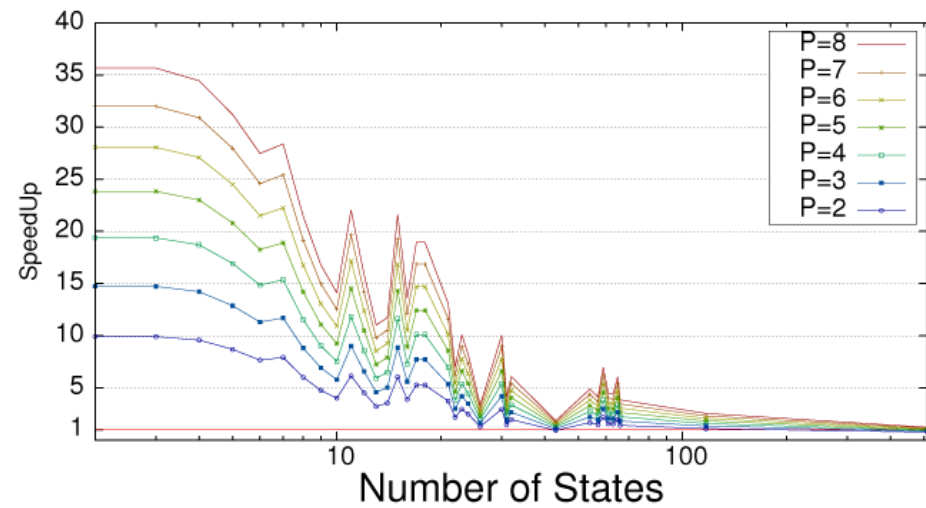
(b) Speedup over the UNIX grep utility

SIMD

28



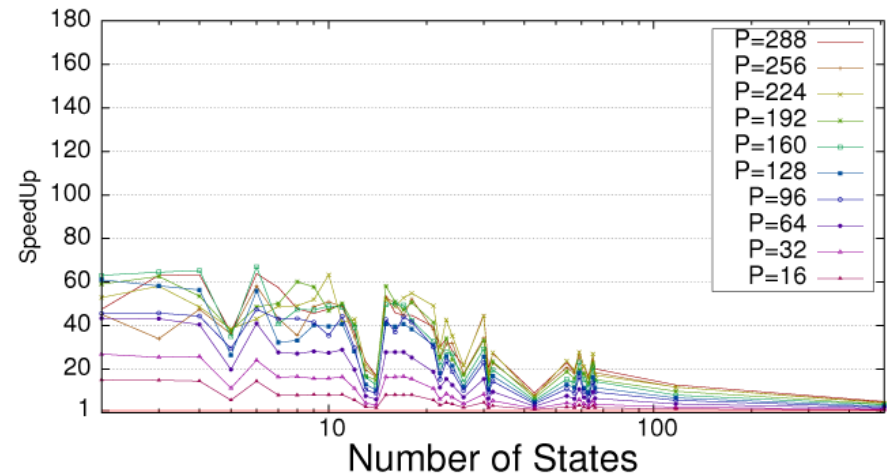
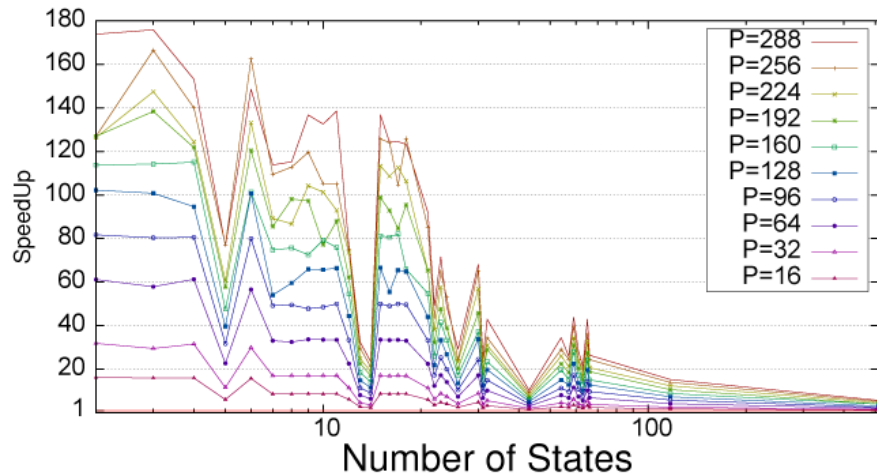
(a) Speedup without AVX2 on PCRE



(b) Speedup with AVX2 on PCRE

Cloud Computing Environment

29



(a) Speedup without MPI communication cost on PCRE

(b) Speedup with MPI communication cost on PCRE

Thanks!