

# Efficient Algorithms for Regular Expression Constrained Sequence Alignment

Yun-Sheng Chung<sup>†</sup>, Chin Lung Lu<sup>¶</sup>, and Chuan Yi Tang<sup>†</sup>

<sup>†</sup>Department of Computer Science

National Tsing Hua University, Taiwan

<sup>¶</sup>Department of Biological Science and Technology

National Chiao Tung University, Taiwan

17th Annual Symposium on Combinatorial Pattern Matching

# Regular Expression Constraints

- PROSITE contains biologically important sites and motifs.
- PROSITE motifs can be represented by regular expressions.

## Technical section:

PROSITE method (with tools and information) covered by this documentation:

ATP\_GTP\_A, [PS00017](#); ATP/GTP-binding site motif A (P-loop) (PATTERN *with a high probability of occurrence*)

Consensus pattern: [AG] - x(4) - G - K - [ST]

P-loop motif: [AG]-x(4)-G-K-[ST]

Regular expression: (A + G)ΣΣΣΣGK(S + T)

# Regular Expression Constraints

- In an alignment, it is reasonable to expect functional sites to be aligned together.
- Introduced by Arslan (CPM 2005).

```

T  G  F  P  S  V  G  K  T  K  D  D  -  -  -  -  A
|      |      |  |      |  |  |      |
T  -  F  -  S  V  A  -  -  K  D  D  D  G  K  S  A
    
```

```

T  -  -  -  G  F  P  S  V  G  K  T  K  D  D  A
|                        |  |      |
T  F  S  V  A  K  D  D  D  G  K  S  -  -  -  A
          *  *  *  *  *  *  *  *
    
```

$(G + A)\Sigma\Sigma\Sigma G K(S + T)$ : the P-loop motif.

# Outline

- Related works
- Weighted automata
- Arslan's algorithm
- Our algorithms
- Conclusion and future works

# Related Works

- Tang et al. (CSB 2003, JBCB 2003) introduced the Constrained Multiple Sequence Alignment (CMSA) problem
  - Constraint: a sequence of characters
  - Example: Sequences from RNase family share the conserved sequence of residues H, K, H. In this case the constraint is H, K, H.

W A Q H — — K P — H C  
W — — H A Q K P Y H C

- Chin et al. (CSB 2004, JBCB 2005) then proposed a more efficient algorithm for pairwise CSA, along with a 2-approximation algorithm for CMSA.

# Related Works

- Tsai et al. (Bioinformatics 2004) generalized the definition of constraints.
  - Constraint: Sequence of strings allowing mismatches
  - $P_1 = AGCC, P_2 = CG$   
 $\epsilon = 0.25$   
 $\dots AGCC \dots CG \dots$   
 $\dots AGUC \dots CG \dots$
  - MuSiC: a web tool
- Lu and Huang (Bioinformatics 2005) gave a more space efficient algorithm
  - MuSiC-ME

# RECSA

- Arslan (CPM 2005) introduced the regular expression constrained sequence alignment (RECSA) problem.
  - Constraint: Regular expression  $R$ .
  - $R = (G + A)\Sigma\Sigma\Sigma\Sigma GK(S + T)$

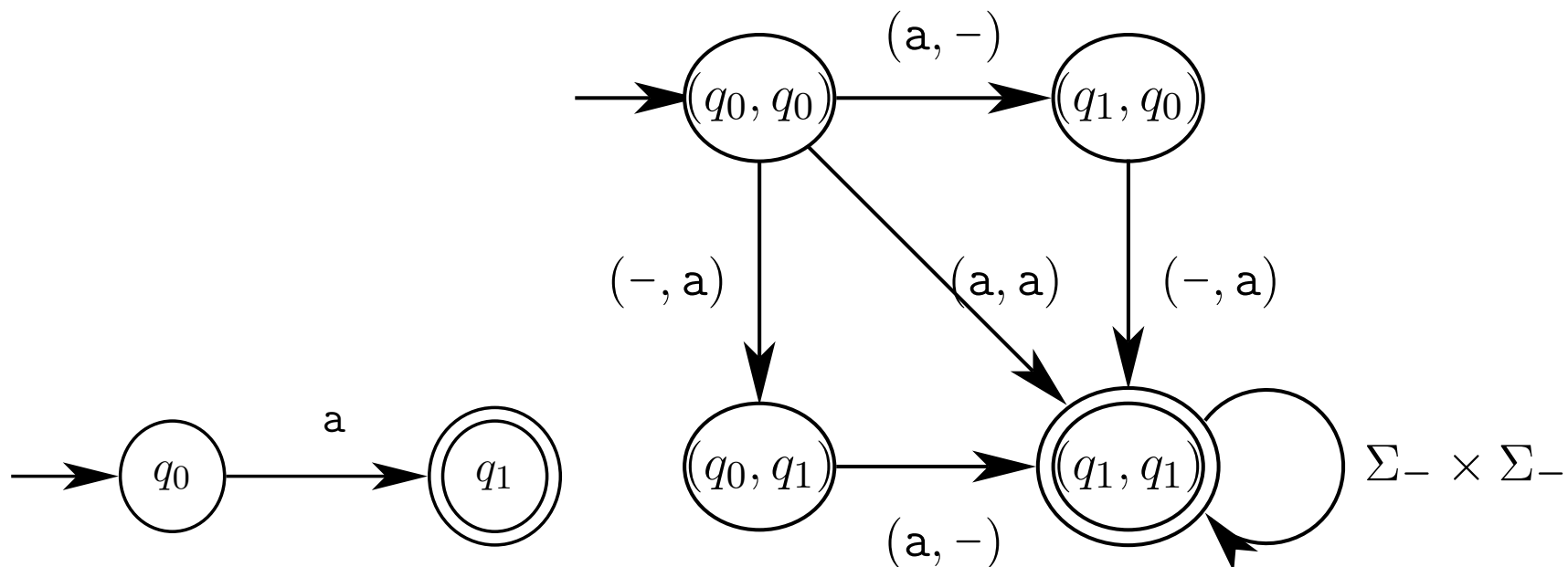
T - - - G F P S V G K T K D D A  
T F S V A K D D D G K S - - - A

- In this current work, we give more time and space efficient algorithms.

# Weighted Automata: Topology

- Regular expression constraint  $R$  is converted into an  $\epsilon$ -free NFA  $A$ .
- The topology of a weighted automaton  $M$  is essentially the same as the topology of the product machine of  $A$ .

$R = a; \Sigma_- = \Sigma \cup \{-\}$ :

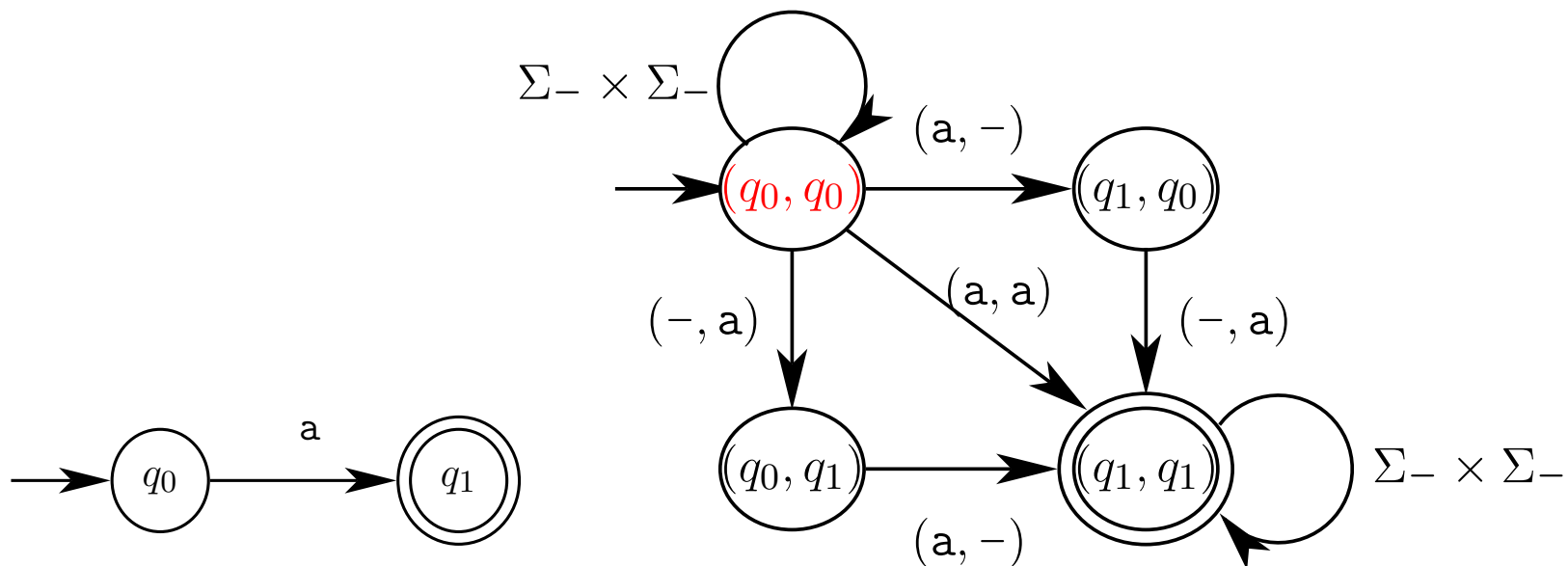




# Weighted Automata

Let  $\delta$  be the transition function of  $A$ . Then  $\delta^M$ , the transition function of  $M$ , is defined as

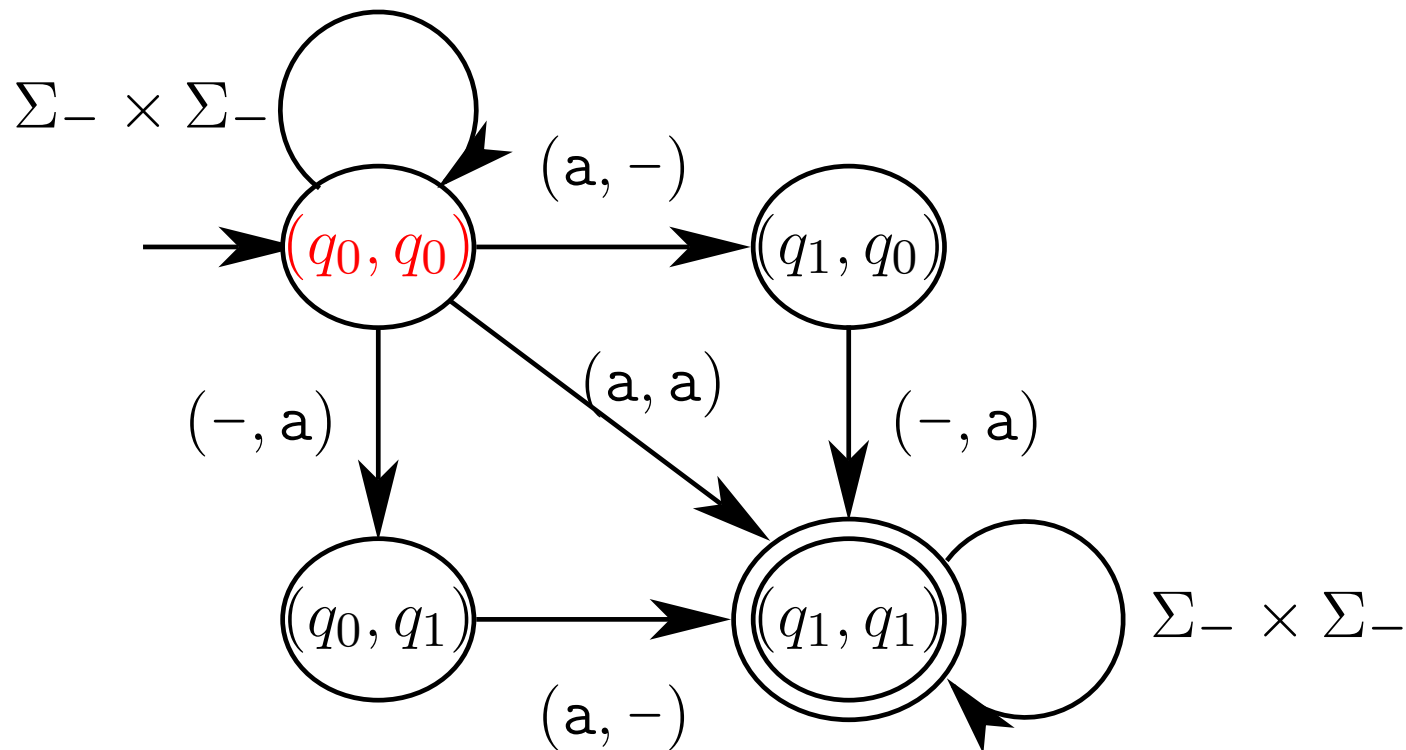
$$\delta^M((p, q), (a, b)) = \begin{cases} \delta(p, a) \times \delta(q, b) & \text{if } (p, q) \notin (F \times F) \cup \{(q_0, q_0)\} \\ (\delta(p, a) \times \delta(q, b)) \cup \{(p, q)\} & \text{otherwise} \end{cases}$$



# Weighted Automata and Alignments

$S_1[1..3] = tca, S_2[1..3] = atc, R = a$

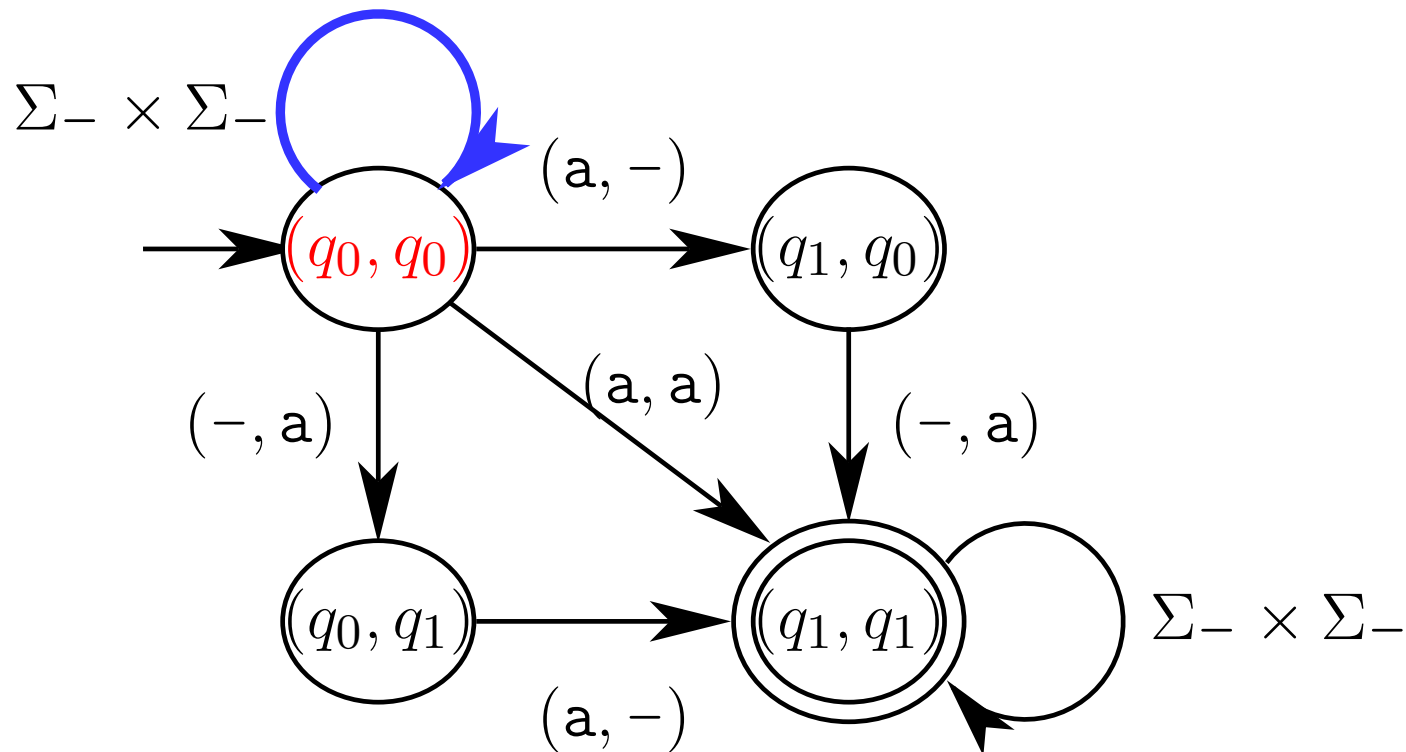
$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$



# Weighted Automata and Alignments

$S_1[1..3] = tca, S_2[1..3] = atc, R = a$

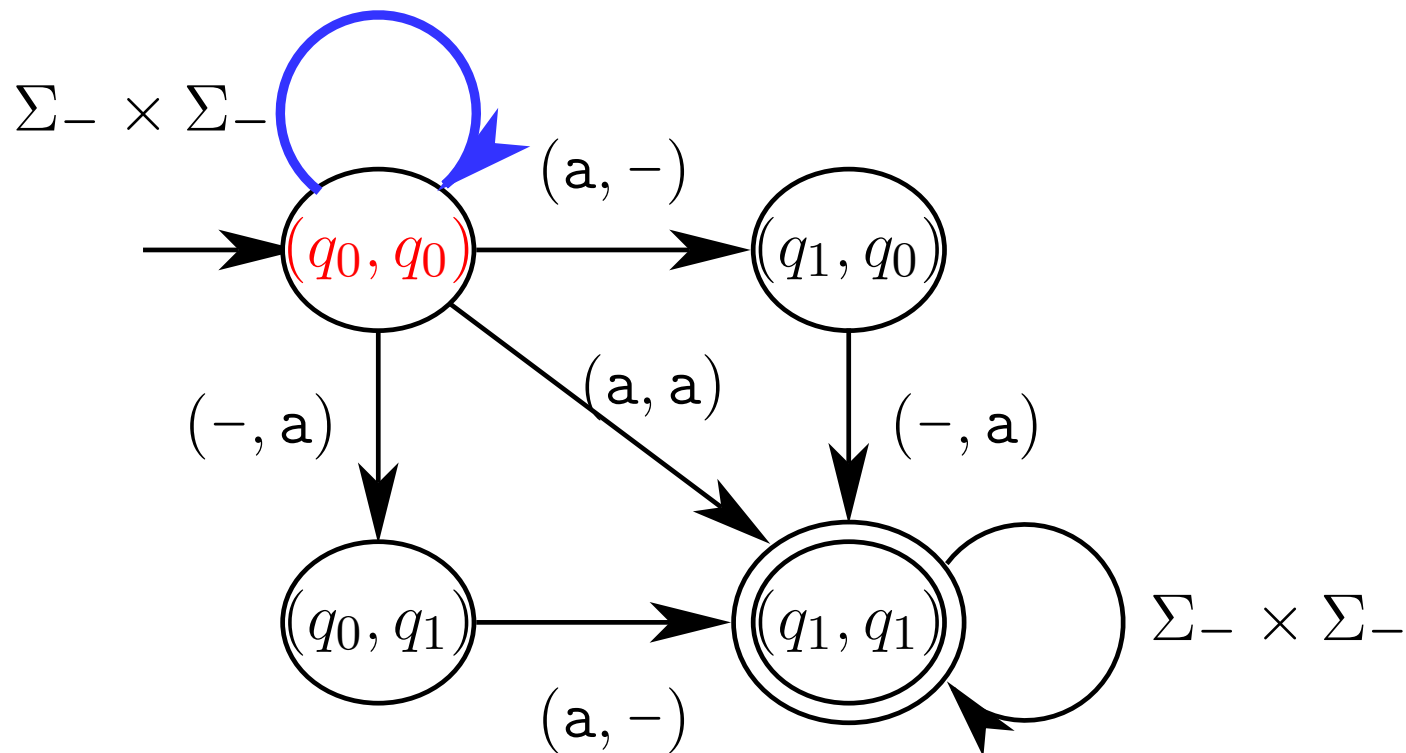
$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$



# Weighted Automata and Alignments

$S_1[1..3] = tca, S_2[1..3] = atc, R = a$

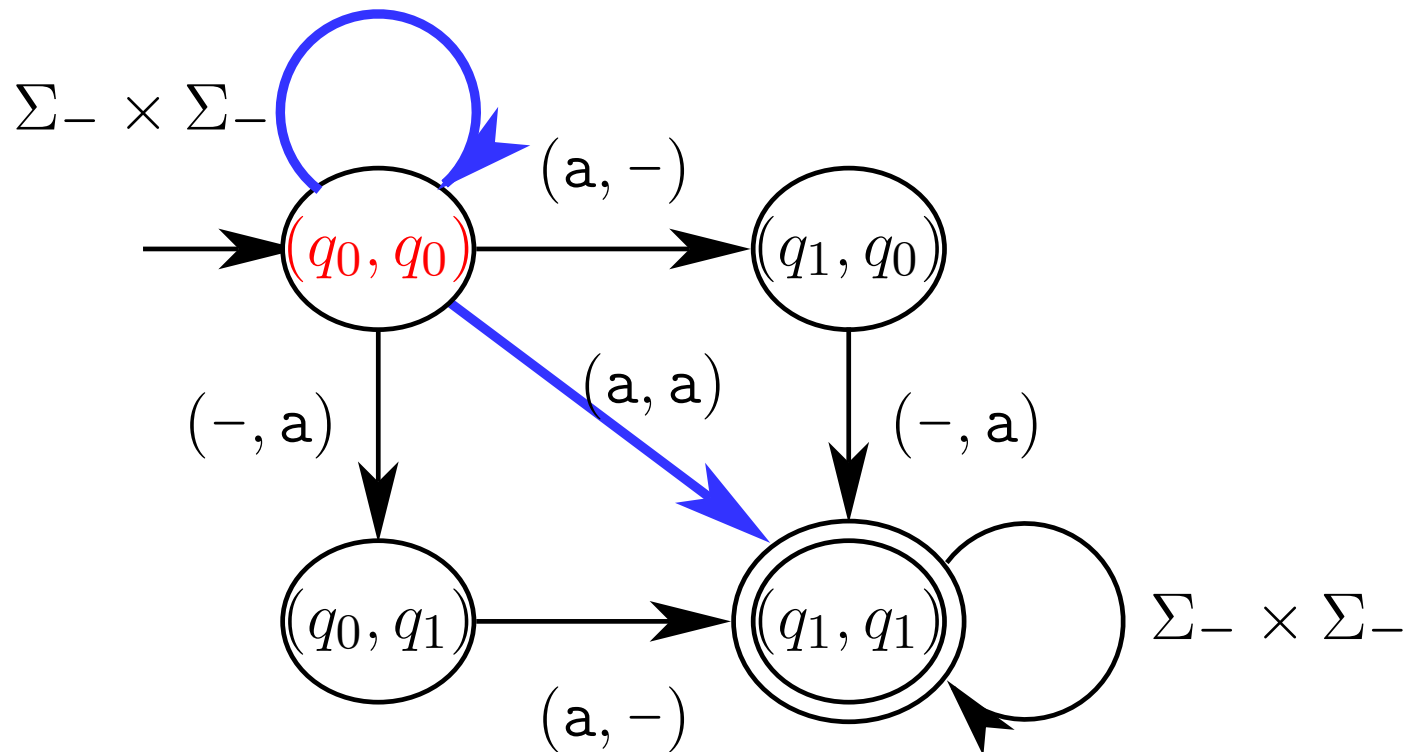
$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$



# Weighted Automata and Alignments

$S_1[1..3] = tca, S_2[1..3] = atc, R = a$

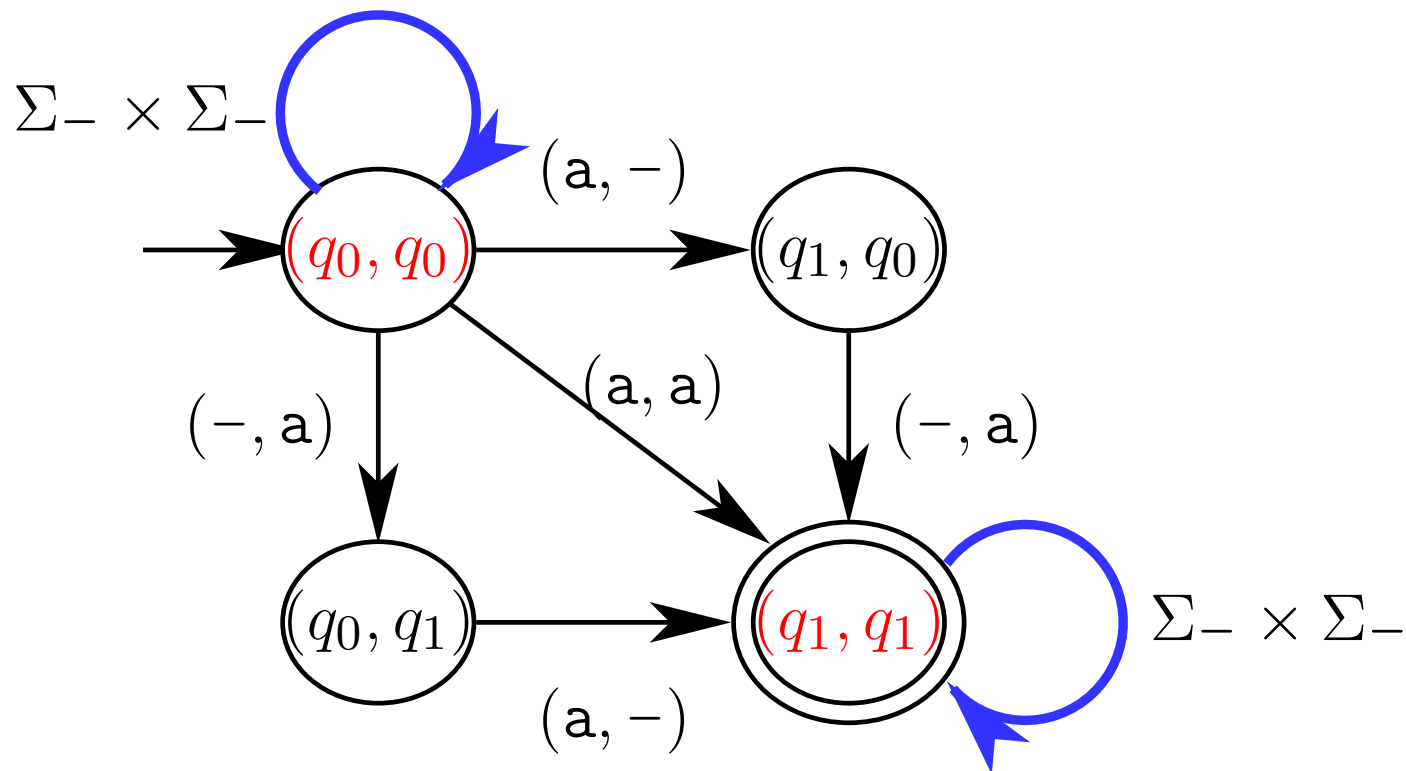
$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$



# Weighted Automata and Alignments

$S_1[1..3] = tca$ ,  $S_2[1..3] = atc$ ,  $R = a$

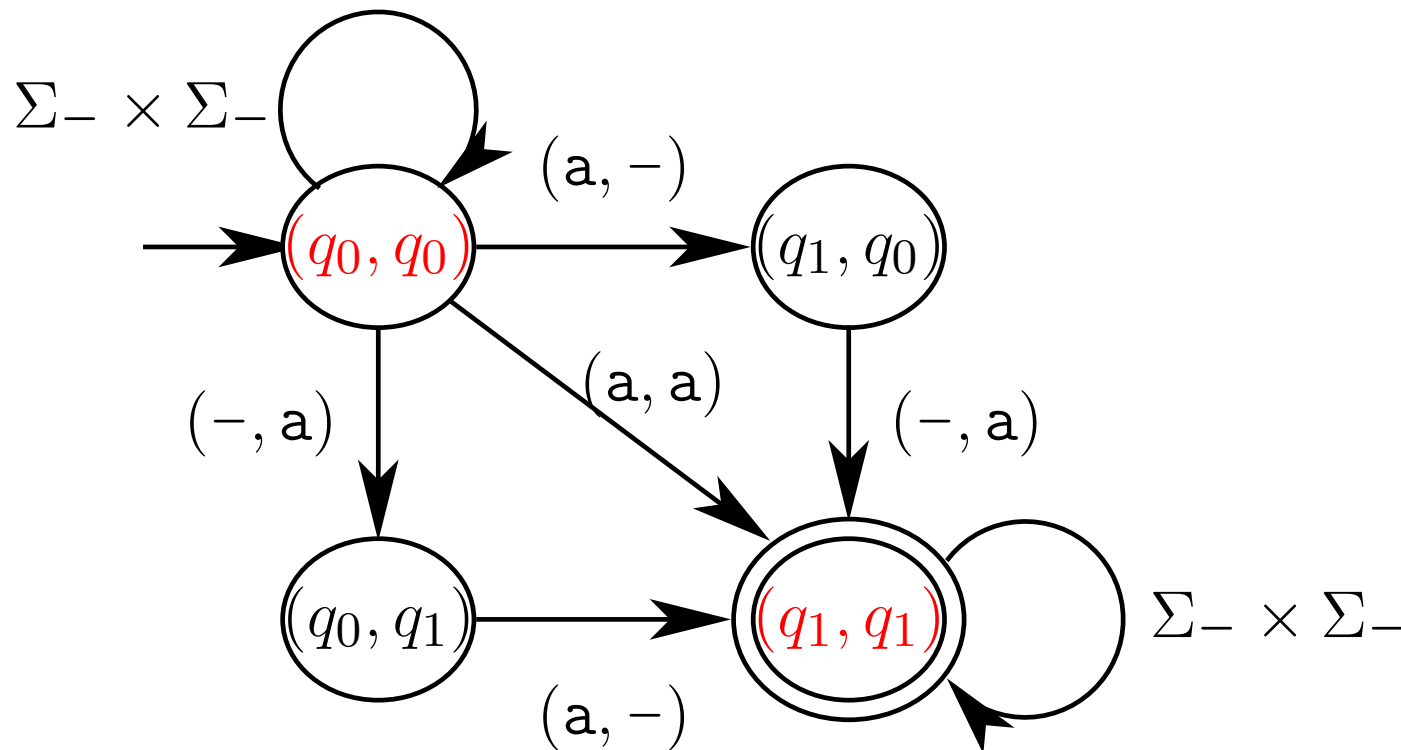
$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$



# Weighted Automata and Alignments

$S_1[1..3] = \text{tca}$ ,  $S_2[1..3] = \text{atc}$ ,  $R = \text{a}$

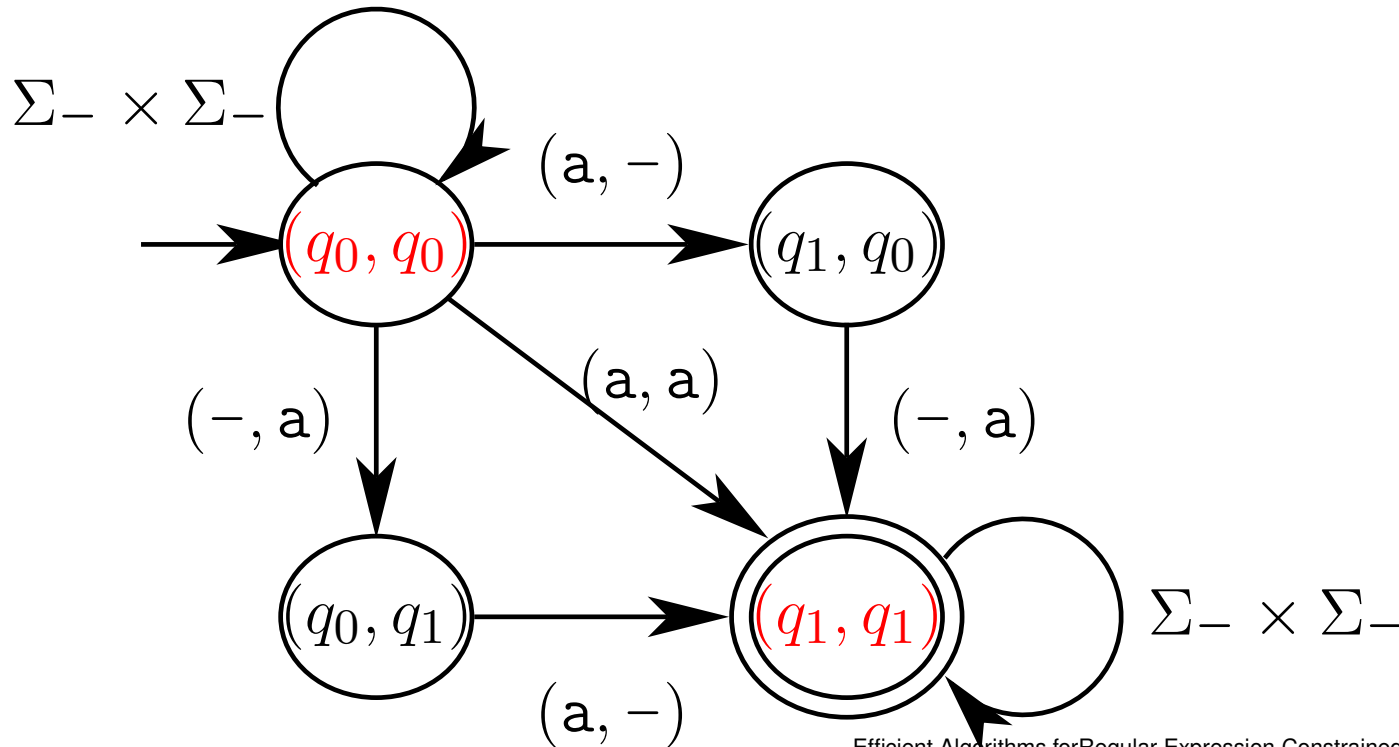
$$\mathcal{A} = \begin{bmatrix} \text{t} & \text{c} & \text{a} & - & - \\ - & - & \text{a} & \text{t} & \text{c} \end{bmatrix} \text{ score: } 1$$



# Weighted Automata and Alignments

$\mathcal{A}$  reaches the final state if and only if  
 $\mathcal{A}$  is a feasible constrained alignment

$$\mathcal{A} = \begin{bmatrix} t & c & a & - & - & - \\ - & - & - & a & t & c \end{bmatrix} \text{ score: } 0$$



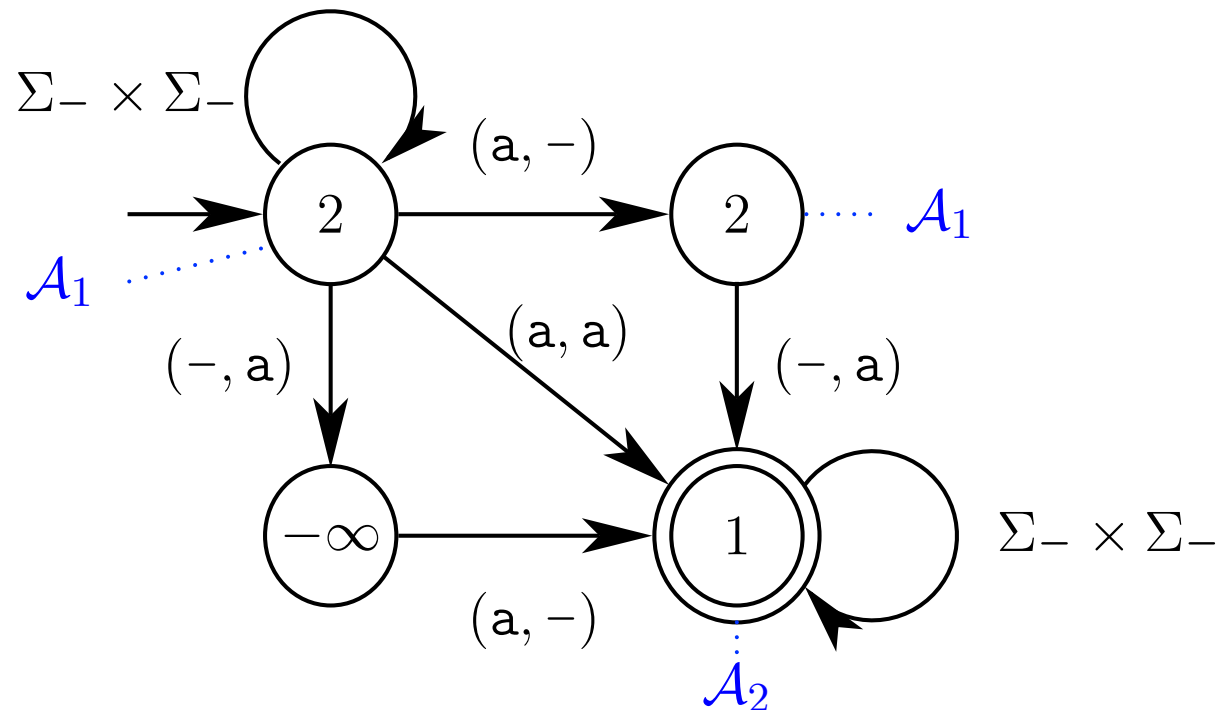


# Weighted Automata: Scores

$S_1[1..3] = tca$ ,  $S_2[1..3] = atc$ ,  $R = a$

$$\mathcal{A}_1 = \begin{bmatrix} - & t & c & a \\ a & t & c & - \end{bmatrix} \quad \mathcal{A}_2 = \begin{bmatrix} t & c & a & - & - \\ - & - & a & t & c \end{bmatrix}$$

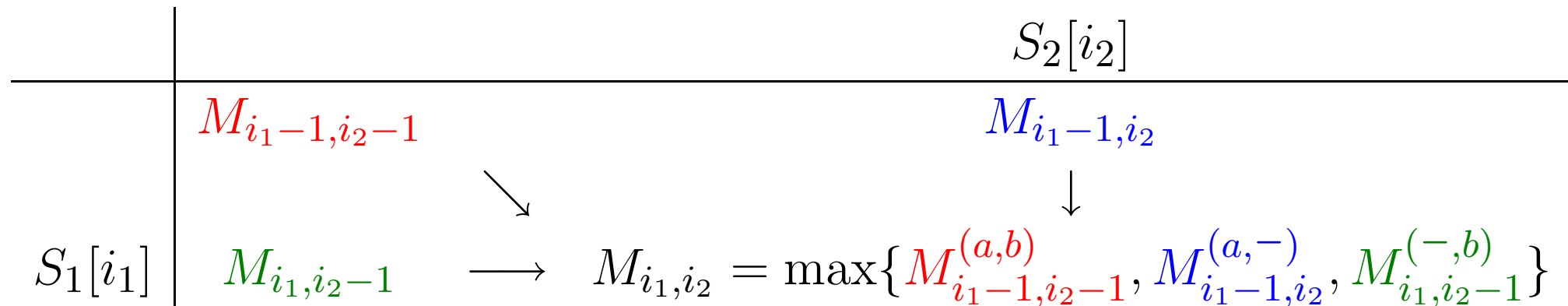
$M_{3,3}$ : (match: scored 1; other: scored 0)



# Weighted Automata: Computation

	$\epsilon$	$a$	$t$	$\dots$
$\epsilon$	$M_{0,0}$	$M_{0,1}$	$M_{0,2}$	$\dots$
$t$	$M_{1,0}$	$M_{1,1}$	$M_{1,2}$	$\dots$
$c$	$M_{2,0}$	$M_{2,1}$	$M_{2,2}$	$\dots$

Let  $a = S_1[i_1]$  and  $b = S_2[i_2]$ :

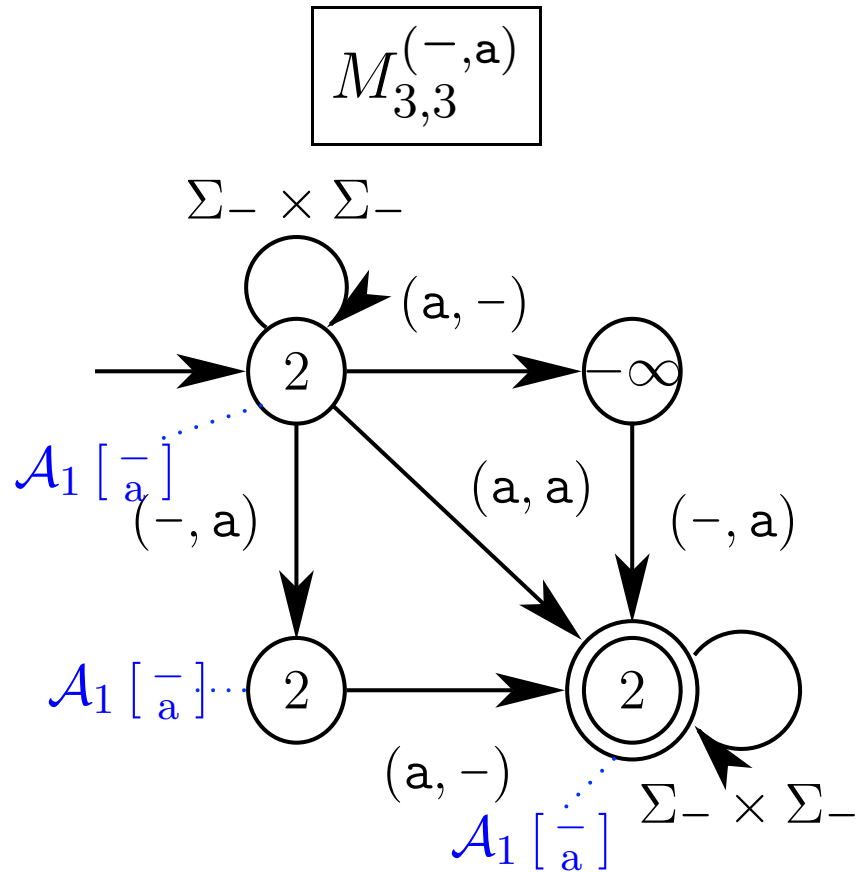
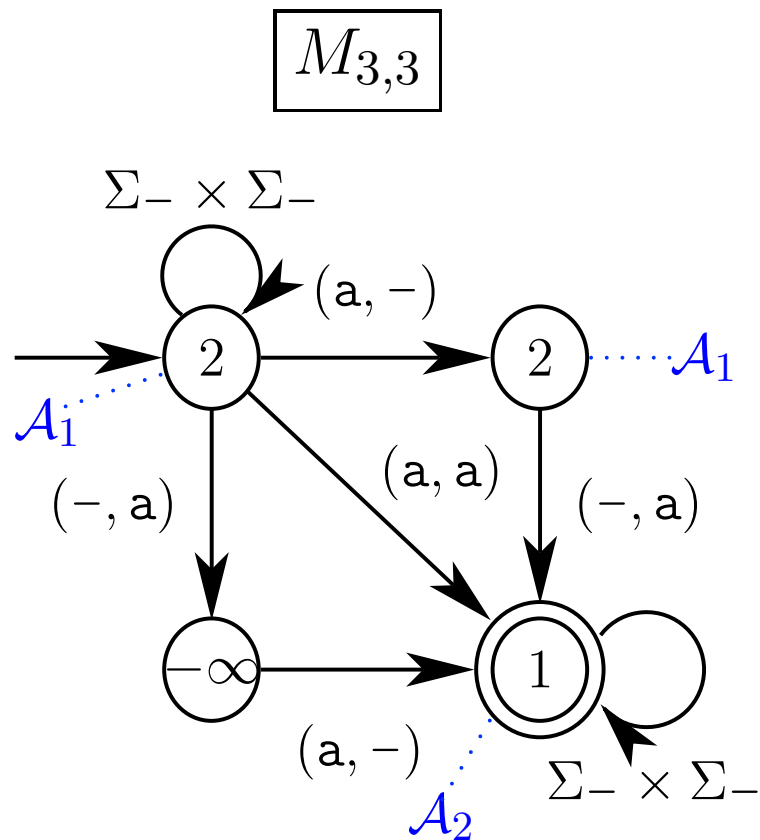


$$W_{i_1, i_2}(p, q) = \max\{W_{i_1-1, i_2-1}^{(a,b)}(p, q), W_{i_1-1, i_2}^{(a,-)}(p, q), W_{i_1, i_2-1}^{(-,b)}(p, q)\}$$

# Weighted Automata: Insertion

$S_1[1..3] = \text{tca}$ ,  $S_2[1..4] = \text{atca}$

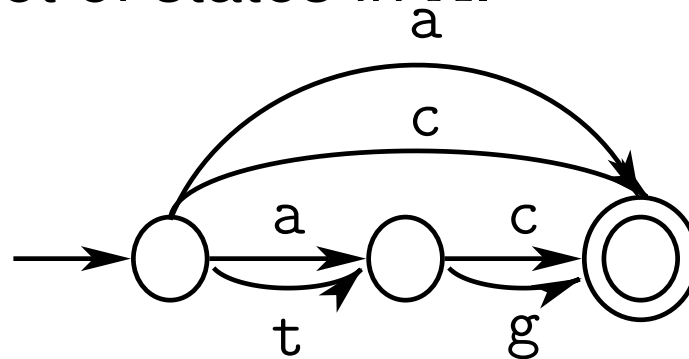
Additional edit operation:  $(-, a)$ , with score 0



The computation can be done by finding in  $M_{3,3}$  all the arcs that are labeled with  $(-, a)$  (or  $\Sigma_- \times \Sigma_-$ ).

# Arslan's Algorithm

- Let  $r$  be the size of the transition function of  $M$ .
- In the algorithm by Arslan, the whole automaton is constructed and searched for each  $(i_1, i_2)$ , hence taking  $O(r)$  time and space.
- Total time and space:  $O(rn^2)$  and  $O(rn)$ , respectively.
- Let  $V$  be the set of states in  $A$ .



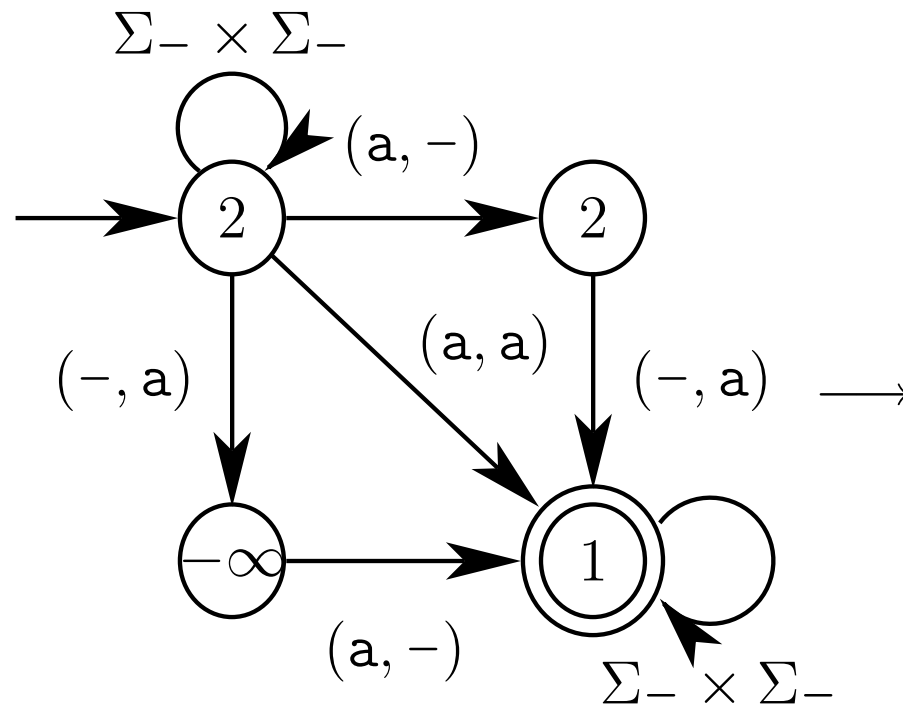
- #arcs in  $A = \Theta(|\Sigma| |V|^2)$  in the worst case.
- $r = \text{\#arcs in } M$ , which is  $\Theta(|\Sigma|^2 |V|^4)$  in the worst case.

# Time and Space Complexity

- Worst case complexity for Arslan's algorithm:  $O(|\Sigma|^2|V|^4n^2)$  time and  $O(|\Sigma|^2|V|^4n)$  space.
- The space complexity is for the computation of the optimal score only. It is not mentioned how to reconstruct the optimal alignment without affecting the space complexity.
- In this work we propose two algorithms.
  - Worst case time:  $O(|V|^3n^2)$  and  $O(|V|^2 \log |V|n^2)$ , respectively.
  - Space:  $O(|V|^2n)$ , with the optimal alignment reconstructed.

# The First Algorithm

- It is more efficient to separate the transitions from the scores.

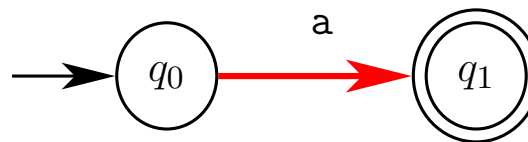
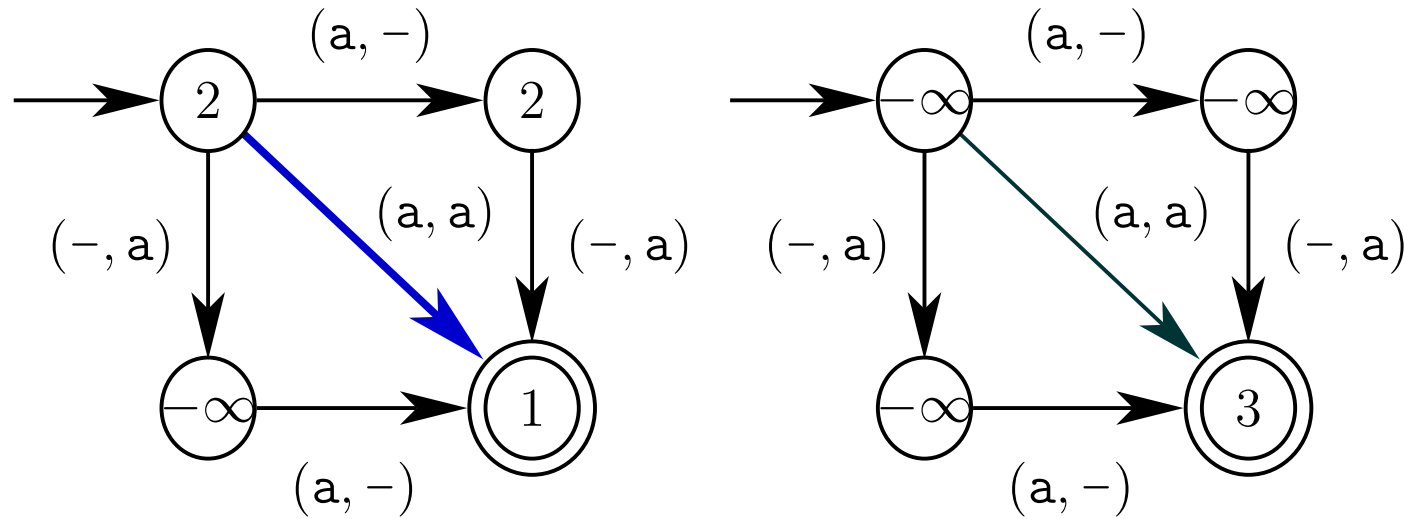


$L_{3,3}$	$q_0$	$q_1$
$q_0$	2	2
$q_1$	$-\infty$	1

- Worst case space complexity:  
 $O(|\Sigma|^2 |V|^4 n) \longrightarrow O(|V|^2 n)$

# Search in $A$ , not $M \rightarrow |V|^4 \rightarrow |V|^3$

- Consider the computation of  $M_{3,3}^{(a,a)}$  from  $M_{3,3}$ .



$L_{3,3}$	$q_0$	$q_1$
$q_0 \downarrow$	2	2
$q_1$	$-\infty$	1

temp	$q_0 \rightarrow$	$q_1$
$q_0$	$-\infty$	$-\infty$
$q_1$	2	2

$L_{3,3}^{(a,a)}$	$q_0$	$q_1$
$q_0$	$-\infty$	$-\infty$
$q_1$	$-\infty$	$2 + 1$

# The First Algorithm

Let  $L$  be a temporary  $|V| \times |V|$  table initialized to be all  $-\infty$ .

**for**  $p \in V$  **do**

**for**  $p' \in V$  with  $T[p', S_1[i_1]; p] = 1$  **do**

**for**  $q' \in V$  **do**

$L[p, q'] \leftarrow$

$\max\{L[p, q'], L_{i_1-1, i_2-1}[p', q'] + \gamma(S_1[i_1], S_2[i_2])\};$

**for**  $q \in V$  **do**

**for**  $q' \in V$  with  $T[q', S_2[i_2]; q] = 1$  **do**

**for**  $p \in V$  **do**

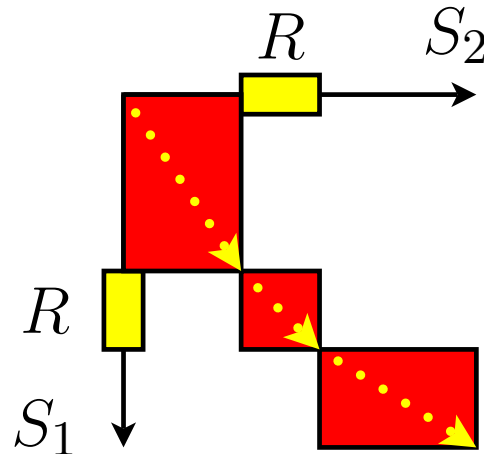
$L_{i_1-1, i_2-1}^{(S_1[i_1], S_2[i_2])}[p, q] \leftarrow \max\{L_{i_1-1, i_2-1}^{(S_1[i_1], S_2[i_2])}[p, q], L[p, q']\};$

The insertion and deletion cases are similar.



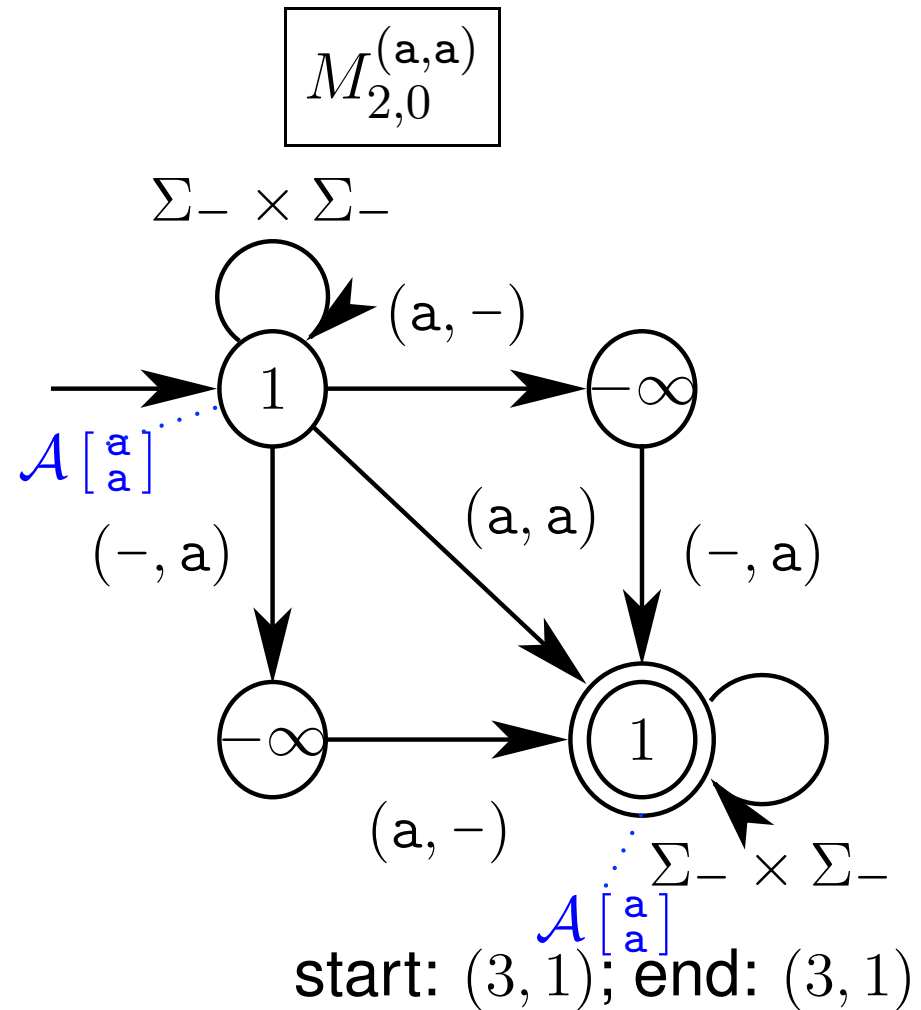
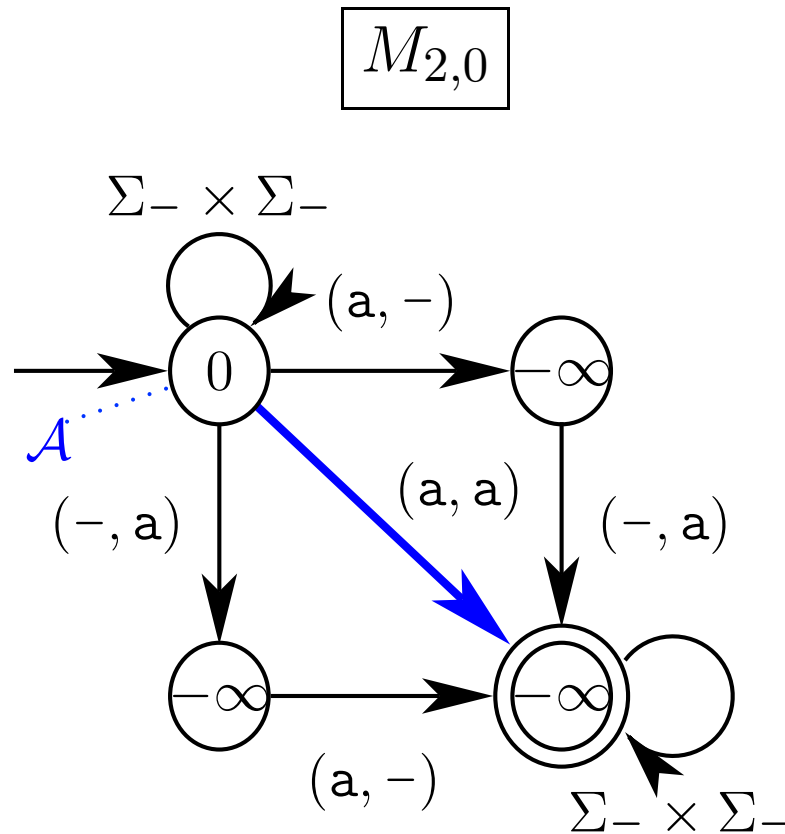
# Reconstruct the Alignment

If we know which substrings of  $S_1$  and  $S_2$  are aligned to satisfy the constraint in an optimal constrained alignment, then the optimal alignment can be reconstructed easily in  $O(n)$  space.



The indices of such substrings can be found by doing some additional bookkeeping during the computation without affecting the space complexity mentioned previously.

# Reconstruct the Alignment



$$\mathcal{A} = \begin{bmatrix} t & c \\ - & - \end{bmatrix}, \quad \mathcal{A} \begin{bmatrix} a \\ a \end{bmatrix} = \begin{bmatrix} t & c & a \\ - & - & a \end{bmatrix}$$

# The Second Algorithm

- The second algorithm requires  $|V| = O(\log n)$ .
- Let  $T[p', a]$  be the bit vector  $(T[p', a; p_1], \dots, T[p', a; p_{|V|}])$ .
- Under the assumption,  $T[p', a]$  can be stored in a single word.
- We take  $L_{i_1-1, i_2}^{(S_1[i_1], -)}[p, q]$  as an example. The other two cases are similar.

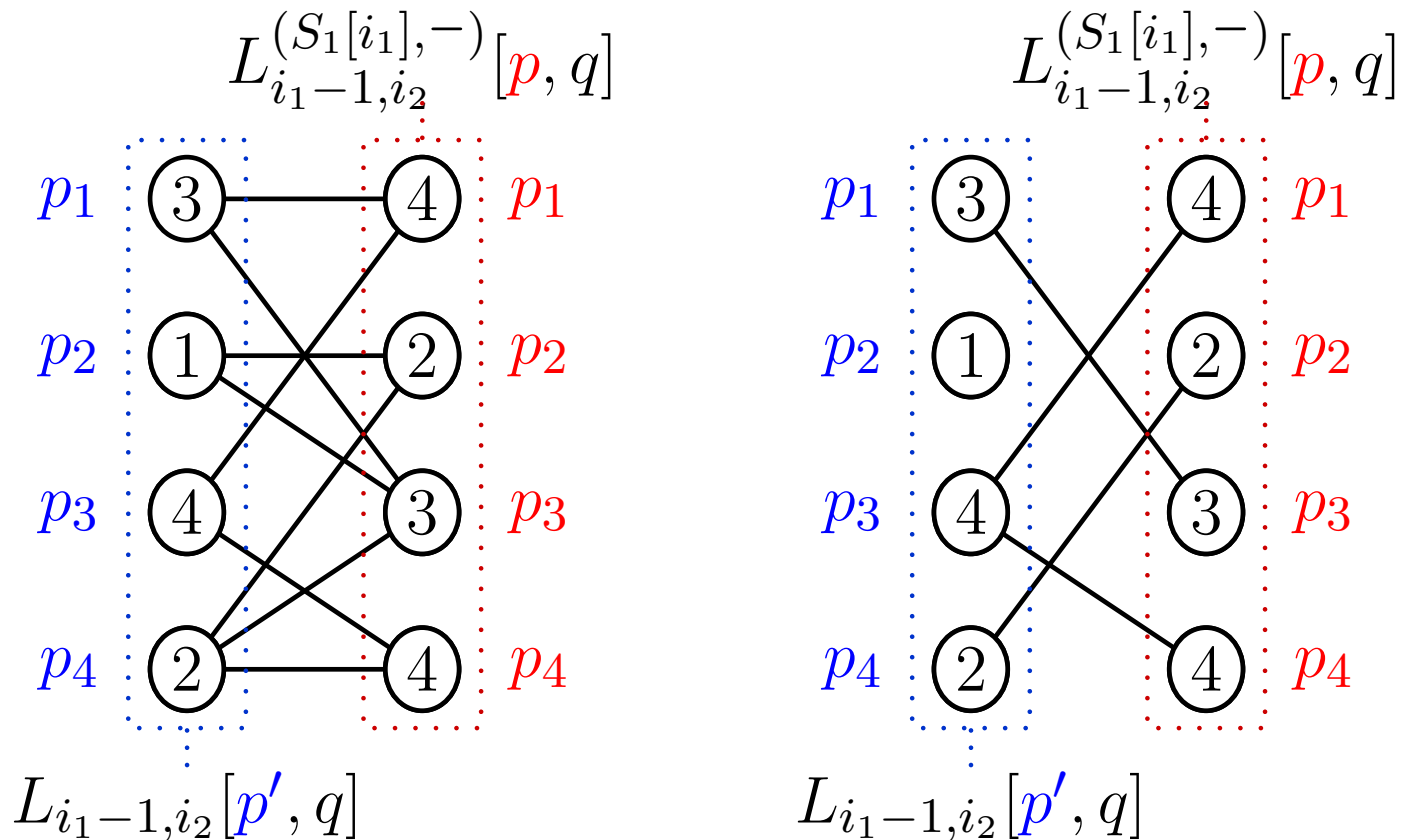
# The Second Algorithm

The arcs shown are all transitions in  $\delta$  labeled with  $S_1[i_1]$ .

$V = \{p_1, p_2, p_3, p_4\}$ . Fix  $q \in V$ .

$T[p_1, S_1[i_1]] = (1, 0, 1, 0)$ ,  $T[p_2, S_1[i_1]] = (0, 1, 1, 0)$ ,

$T[p_3, S_1[i_1]] = (1, 0, 0, 1)$ ,  $T[p_4, S_1[i_1]] = (0, 1, 1, 1)$

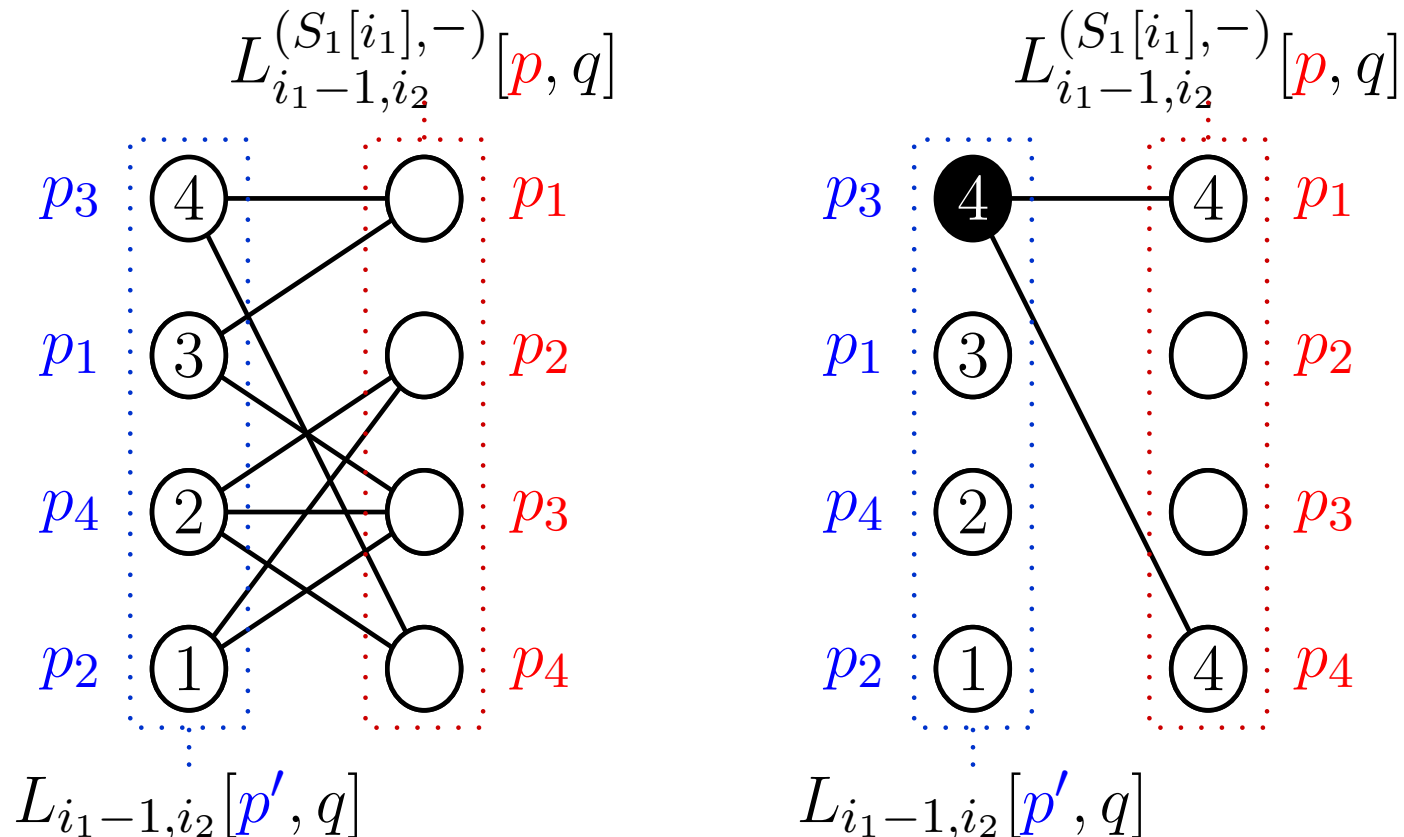


# The Second Algorithm

$Y$ : Bit vector representing the states  $p$  in  $V$  such that

$L_{i_1-1, i_2}^{(S_1[i_1], -)}[p, q]$  has not yet been set.

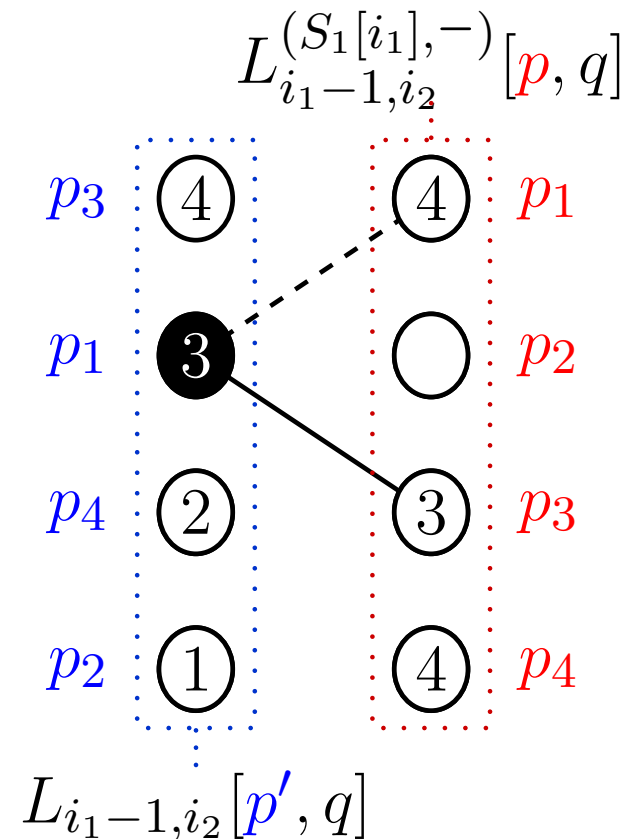
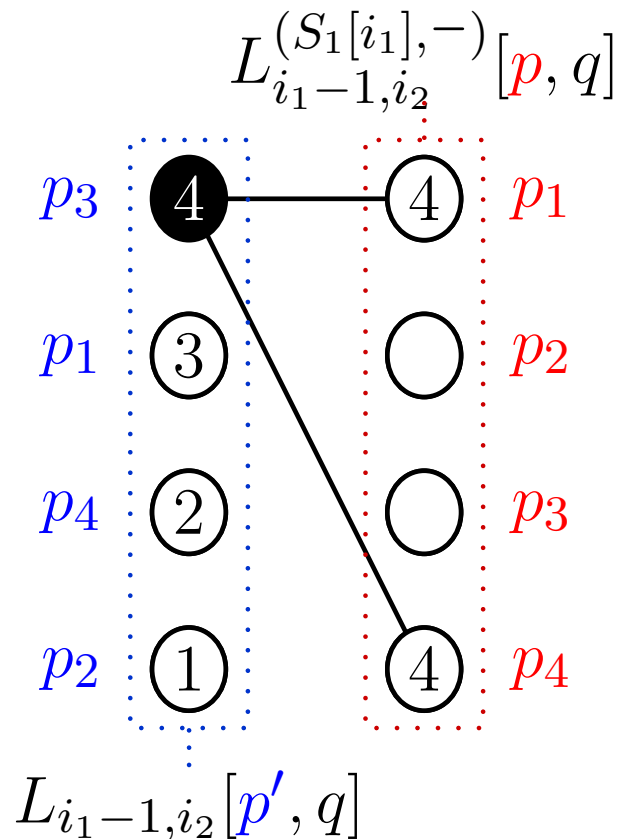
Initially:  $Y = (1, 1, 1, 1)$ .  $T[p_3, S_1[i_1]] \otimes Y = (1, 0, 0, 1)$



# The Second Algorithm

$$Y \leftarrow Y - T[p_3, S_1[i_1]] \otimes Y = (1, 1, 1, 1) - (1, 0, 0, 1) = (0, 1, 1, 0)$$

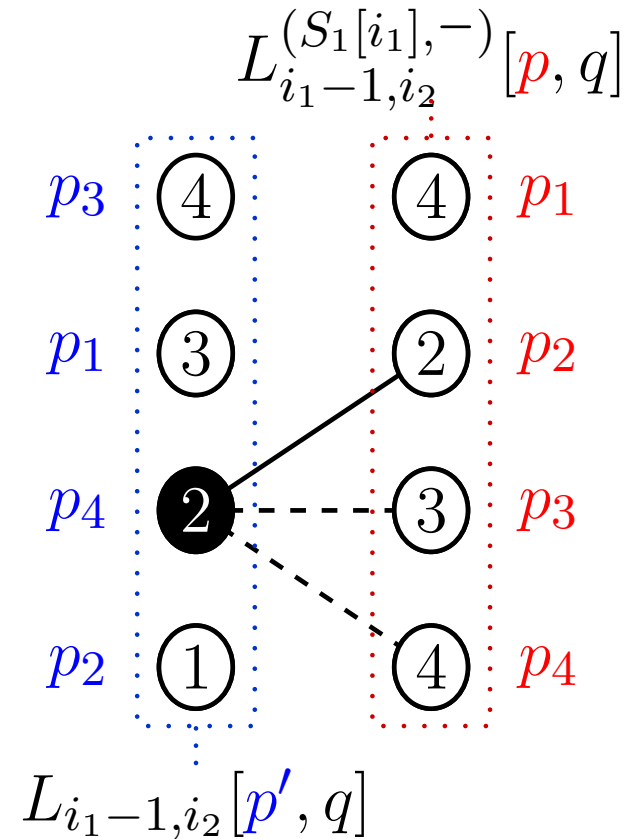
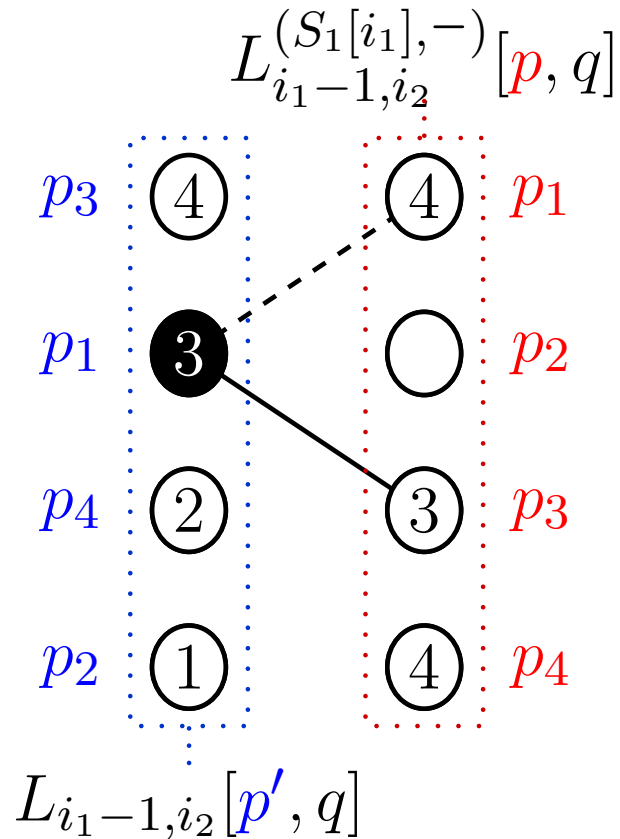
$$T[p_1, S_1[i_1]] \otimes Y = (0, 0, 1, 0)$$



# The Second Algorithm

$$Y \leftarrow Y - T[p_1, S_1[i_1]] \otimes Y = (0, 1, 1, 0) - (0, 0, 1, 0) = (0, 1, 0, 0)$$

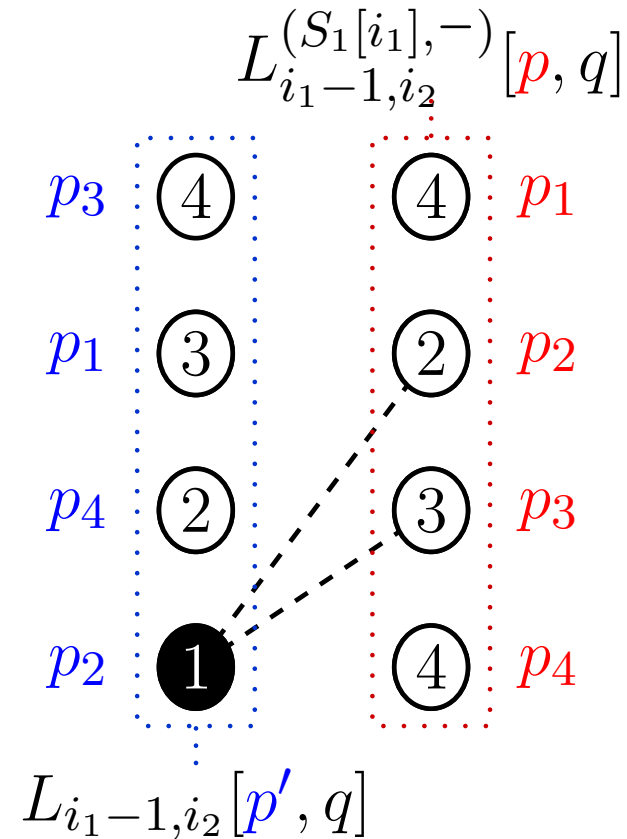
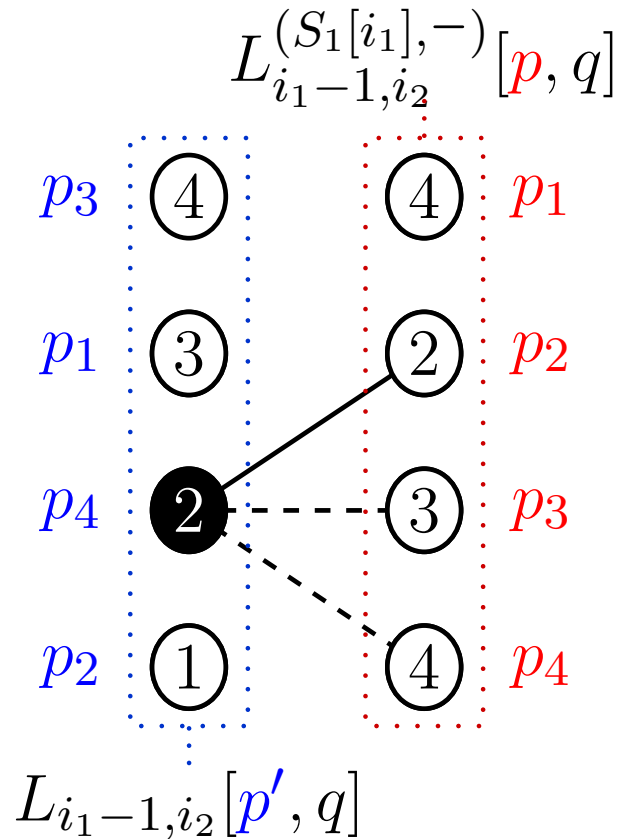
$$T[p_4, S_1[i_1]] \otimes Y = (0, 1, 0, 0)$$



# The Second Algorithm

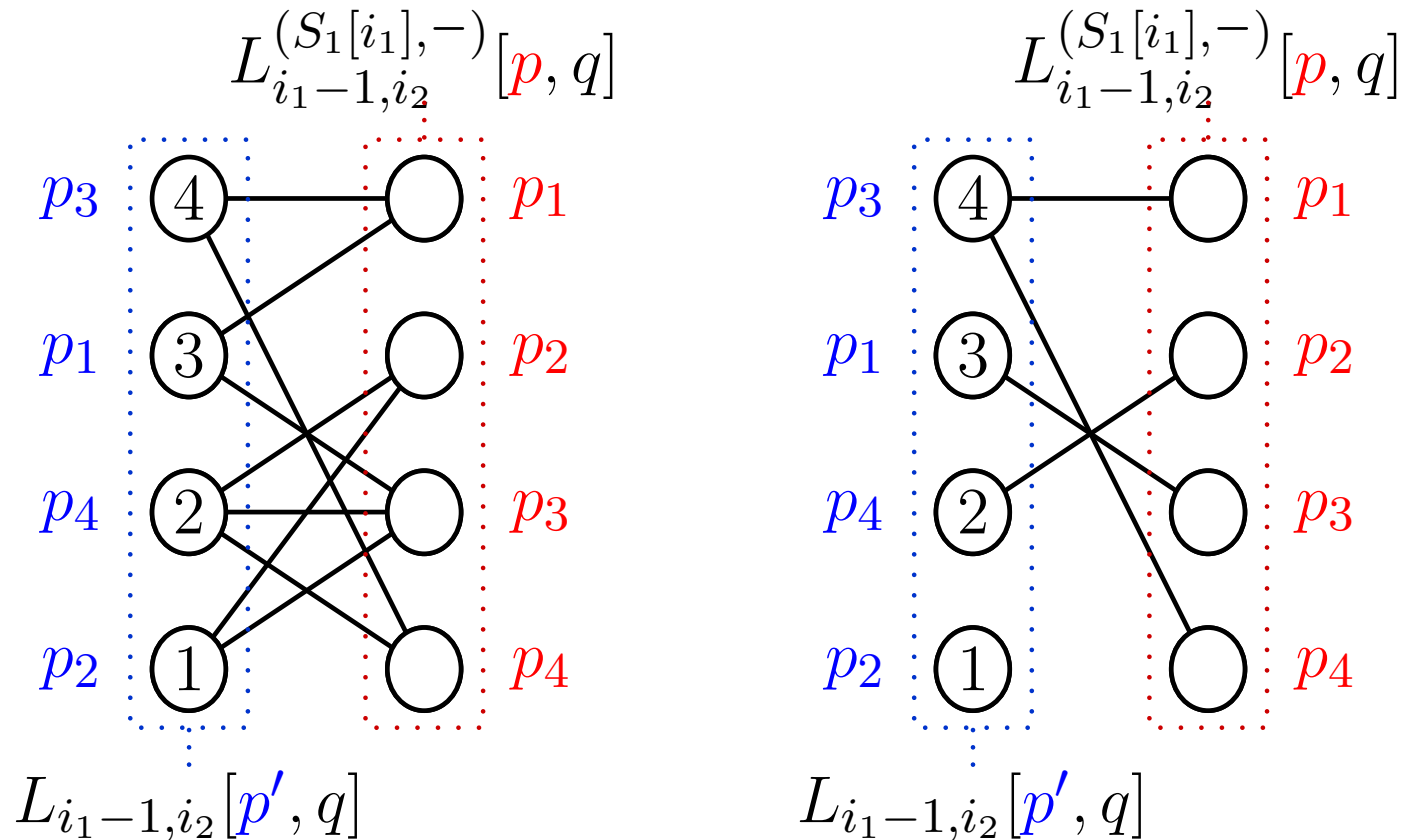
$$Y \leftarrow Y - T[p_4, S_1[i_1]] \otimes Y = (0, 1, 0, 0) - (0, 1, 0, 0) = (0, 0, 0, 0)$$

$$T[p_2, S_1[i_1]] \otimes Y = (0, 0, 0, 0)$$





# The Second Algorithm



Time:  $O(|V| \log |V|)$  for each fixed  $q$ ;  $O(|V|^2 \log |V|)$  for the entire  $L_{i_1-1, i_2}^{(S_1[i_1], -)}$

# Conclusion and Future Works

- More time and space efficient algorithms for regular expression constrained sequence alignment are proposed.
- A multiple alignment version is important.
  - Generalize the weighted automata structure to obtain optimal solutions: suggested by Arslan (CPM 2005).
  - Progressive heuristics: e.g., Lu and Huang's work (Bioinformatics 2005).
  - Approximation algorithms: in preparation.
- Local alignment version would also be useful.
  - Has a similar motivation as PHI-BLAST; may work complementarily by taking a different approach.

Thank you for your attention.