OpenFst: An Open-Source, Weighted Finite-State Transducer Library and its Applications to Speech and Language

Part I. *Theory and Algorithms*

# Overview

1. Preliminaries

   - Semirings

   - Weighted Automata and Transducers

2. Algorithms

   - Rational Operations

   - Elementary Unary Operations

   - Fundamental Binary Operations

   - Optimization Algorithms

   - Search Operations

   - Fundamental String Algorithms
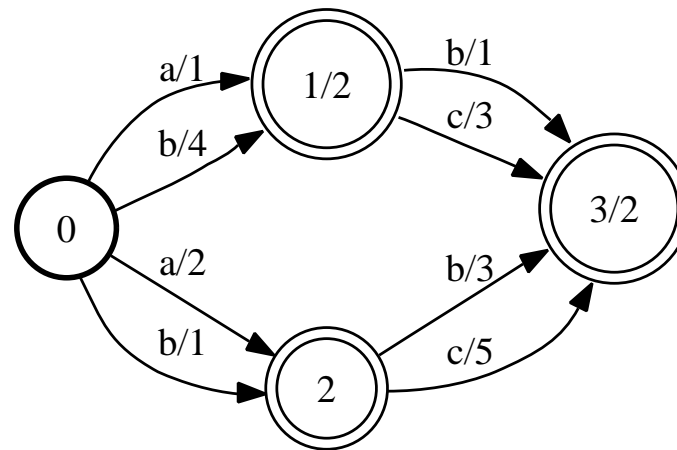
# Weight Sets: Semirings

A *semiring* $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$ = a ring that may lack negation.

- Sum: to compute the weight of a sequence (sum of the weights of the paths labeled with that sequence).

- Product: to compute the weight of a path (product of the weights of constituent transitions).

| SEMIRING | SET | $\oplus$ | $\otimes$ | $\bar{0}$ | $\bar{1}$ |
|----------|-----|----------|-----------|-----------|-----------|
| Boolean | $\{0, 1\}$ | $\vee$ | $\wedge$ | $0$ | $1$ |
| Probability | $\mathbb{R}_+$ | $+$ | $\times$ | $0$ | $1$ |
| Log | $\mathbb{R} \cup \{-\infty, +\infty\}$ | $\oplus_{\log}$ | $+$ | $+\infty$ | $0$ |
| Tropical | $\mathbb{R} \cup \{-\infty, +\infty\}$ | $\min$ | $+$ | $+\infty$ | $0$ |
| String | $\Sigma^* \cup \{\infty\}$ | $\wedge$ | $\cdot$ | $\infty$ | $\epsilon$ |

$\oplus_{\log}$ is defined by: $x \oplus_{\log} y = -\log(e^{-x} + e^{-y})$ and $\wedge$ is longest common prefix. The string semiring is a *left semiring*.

# Weighted Automaton/Acceptor



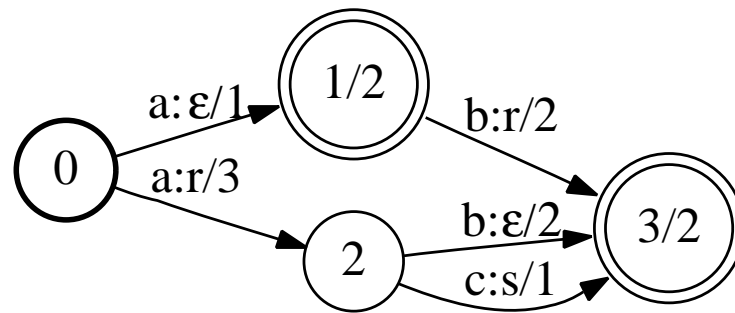| Probability semiring $(\mathbb{R}_+, +, \times, 0, 1)$ | Tropical semiring $(\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$ |
| :---: | :---: |
| $[\![A]\!](ab) = 14$ | $[\![A]\!](ab) = 4$ |
| $(1 \times 1 \times 2 + 2 \times 3 \times 2 = 14)$ | $(\min(1+1+2, 3+2+2) = 4)$ |

# Weighted Transducer



| Probability semiring $(\mathbb{R}_+, +, \times, 0, 1)$ | Tropical semiring $(\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$ |
| :---: | :---: |
| $\llbracket T \rrbracket(ab, r) = 16$ | $\llbracket T \rrbracket(ab, r) = 5$ |
| $(1 \times 2 \times 2 + 3 \times 2 \times 2 = 16)$ | $(\min(1 + 2 + 2, 3 + 2 + 2) = 5)$ |

# Transducers as Weighted Automata

A transducer $T$ is *functional* iff for each $x$ there exists at most one $y$ such that $[\![T]\!](x,y) \neq \overline{0}$

- An unweighted functional transducer can be seen as as:
  $\rightarrow$ a weighted automata over the string semiring $(\Sigma^* \cup \{\infty\}, \wedge, \cdot, \infty, \epsilon)$

- A weighted functional transducer over the semiring $\mathbb{K}$ can be seen as:
  $\rightarrow$ a weighted automata over the cartesian product of the string semiring and $\mathbb{K}$



$$[\![T]\!](ab, r) = 5 \qquad\qquad [\![A]\!](ab) = (r, 5)$$

[Tropical semiring $(\mathbb{R}_+ \cup \{\infty\}, \min, +, \infty, 0)$]

- Path $\pi$

    - Origin or previous state: $p[\pi]$.

    - Destination or next state: $n[\pi]$.

    - Input label: $i[\pi]$.

    - Output label: $o[\pi]$.



- Sets of paths

    - $P(R_1, R_2)$: set of all paths from $R_1 \subseteq Q$ to $R_2 \subseteq Q$.

    - $P(R_1, x, R_2)$: paths in $P(R_1, R_2)$ with input label $x$.

    - $P(R_1, x, y, R_2)$: paths in $P(R_1, x, R_2)$ with output label $y$.

# Definitions and Notation – Automata and Transducers

1. **General Definitions**

    - Alphabets: input $\Sigma$, output $\Delta$.
    - States: $Q$, initial states $I$, final states $F$.
    - Transitions: $E \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Delta \cup \{\epsilon\}) \times \mathbb{K} \times Q$.
    - Weight functions:

        initial weight function $\lambda : I \rightarrow \mathbb{K}$

        final weight function $\rho : F \rightarrow \mathbb{K}$.

2. **Machines**

    - Automaton $A = (\Sigma, Q, I, F, E, \lambda, \rho)$ with for all $x \in \Sigma^*$:

    $$[\![A]\!](x) = \bigoplus_{\pi \in P(I, x, F)} \lambda(p[\pi]) \otimes w[\pi] \otimes \rho(n[\pi])$$

    - Transducer $T = (\Sigma, \Delta, Q, I, F, E, \lambda, \rho)$ with for all $x \in \Sigma^*, y \in \Delta^*$:

    $$[\![T]\!](x, y) = \bigoplus_{\pi \in P(I, x, y, F)} \lambda(p[\pi]) \otimes w[\pi] \otimes \rho(n[\pi])$$

# Rational Operations – Algorithms

- Definitions

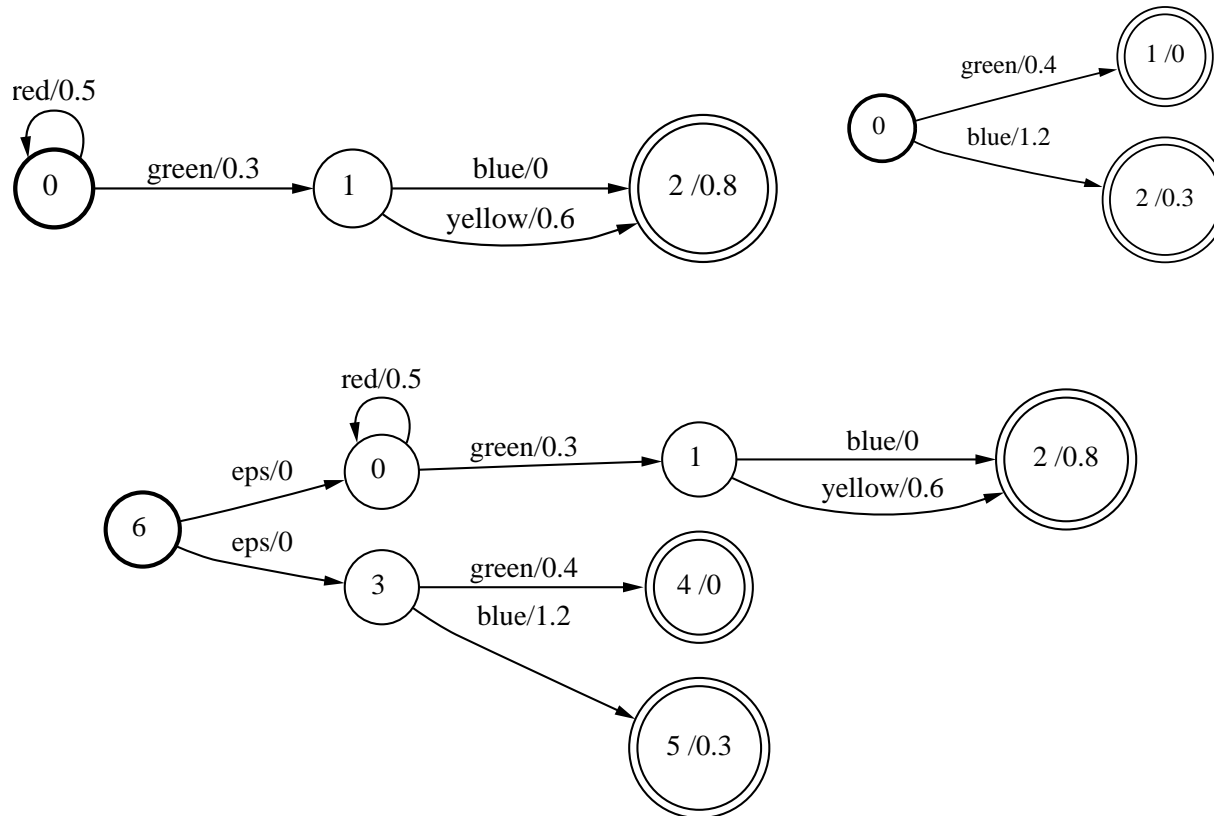| OPERATION | DEFINITION AND NOTATION |
|-----------|-------------------------|
| Sum | $[\![T_1 \oplus T_2]\!](x, y) = [\![T_1]\!](x, y) \oplus [\![T_2]\!](x, y)$ |
| Product | $[\![T_1 \otimes T_2]\!](x, y) = \displaystyle\bigoplus_{x=x_1 x_2, y=y_1 y_2} [\![T_1]\!](x_1, y_1) \otimes [\![T_2]\!](x_2, y_2)$ |
| Closure | $[\![T^*]\!](x, y) = \displaystyle\bigoplus_{n=0}^{\infty} [\![T^n]\!](x, y)$ |

- Conditions on the closure operation: condition on $T$: e.g. weight of $\epsilon$-cycles $= \overline{0}$ (*regulated transducers*), or semiring condition: e.g. $\overline{1} \oplus x = \overline{1}$ as with the tropical semiring (*locally closed semirings*).

- Complexity and implementation

  - Complexity (linear): $O((|E_1| + |Q_1|) + (|E_2| + |Q_2|))$ or $O(|Q| + |E|)$.
  - Lazy implementation.

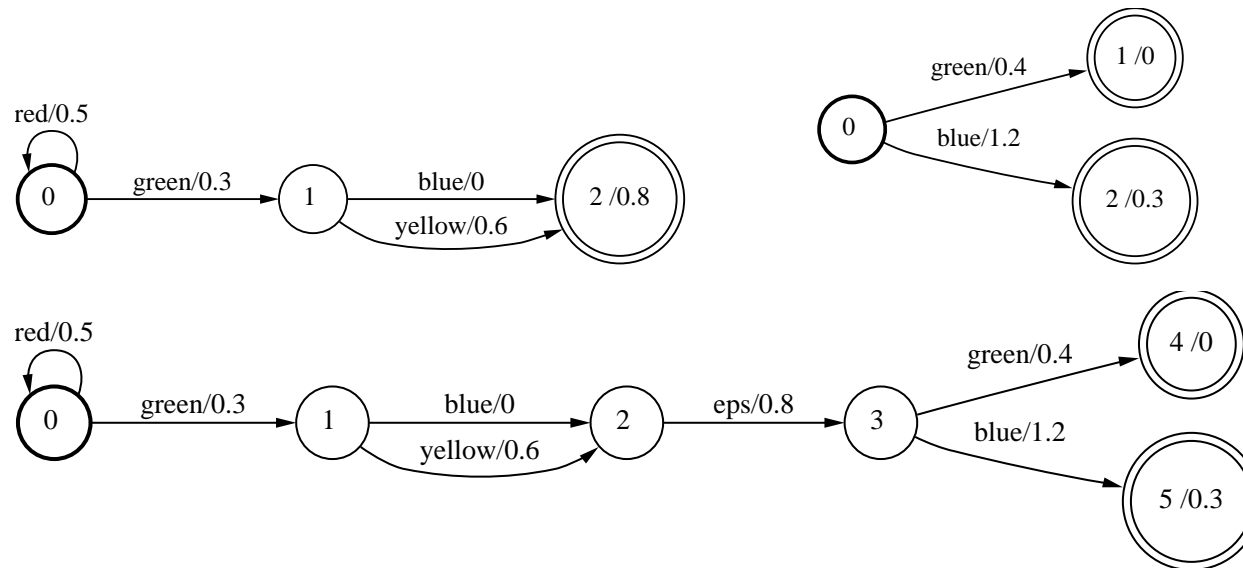# Sum (Union) – Illustration

- **Definition:** $\llbracket T_1 \oplus T_2 \rrbracket(x, y) = \llbracket T_1 \rrbracket(x, y) \oplus \llbracket T_2 \rrbracket(x, y)$

- **Example:**
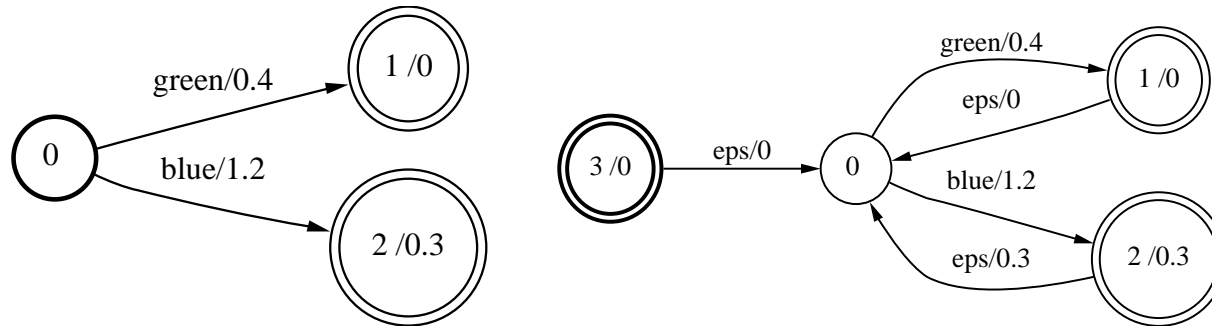
# Product (Concatenation) – Illustration

- Definition: $[\![T_1 \otimes T_2]\!](x, y) = \bigoplus\limits_{x=x_1 x_2, y=y_1 y_2} [\![T_1]\!](x_1, y_1) \otimes [\![T_2]\!](x_2, y_2)$

- Example:

# Closure – Illustration

- Definition: $[\![T^*]\!](x, y) = \bigoplus_{n=0}^{\infty} [\![T^n]\!](x, y)$

- Example:

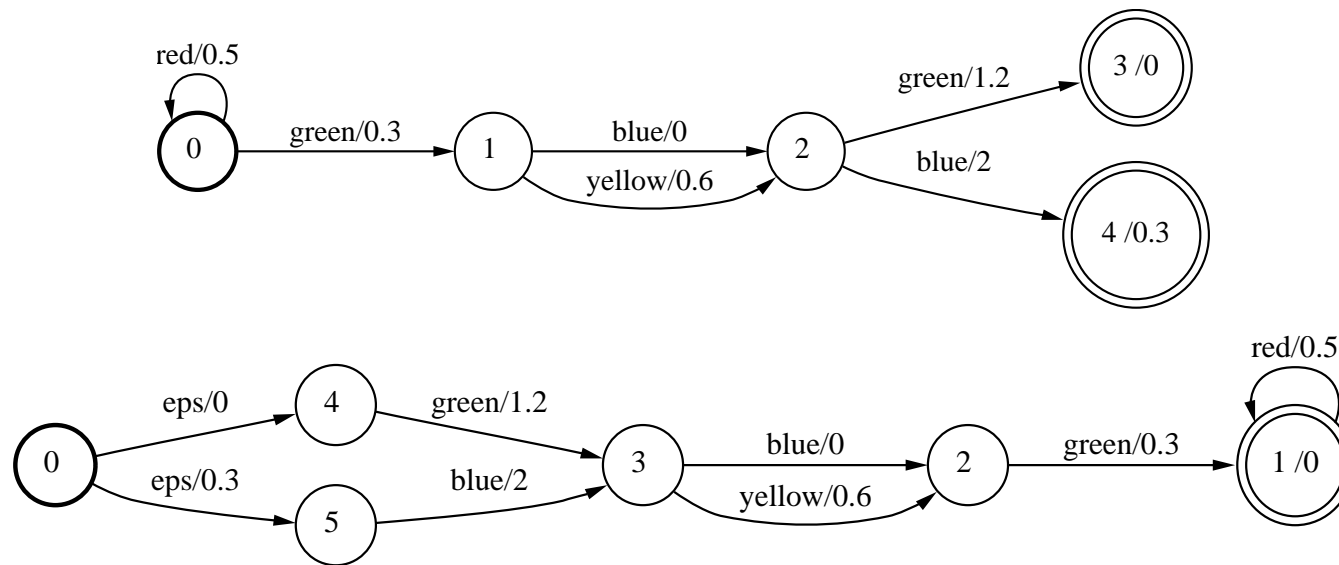# Some Elementary Unary Operations – Algorithms

- Definitions

| Operation | Definition and Notation | Lazy Implementation |
|-----------|-------------------------|---------------------|
| Reversal | $[\![\widetilde{T}]\!](x,y) = [\![T]\!](\widetilde{x},\widetilde{y})$ | No |
| Inversion | $[\![T^{-1}]\!](x,y) = [\![T]\!](y,x)$ | Yes |
| Projection | $[\![\Pi_1(T)]\!](x) = \bigoplus_{y}[\![T]\!](x,y)$ | Yes |
|  | $[\![\Pi_2(T)]\!](x) = \bigoplus_{y}[\![T]\!](y,x)$ |  |

- Complexity and implementation

  - Complexity (linear): $O(|Q| + |E|)$.
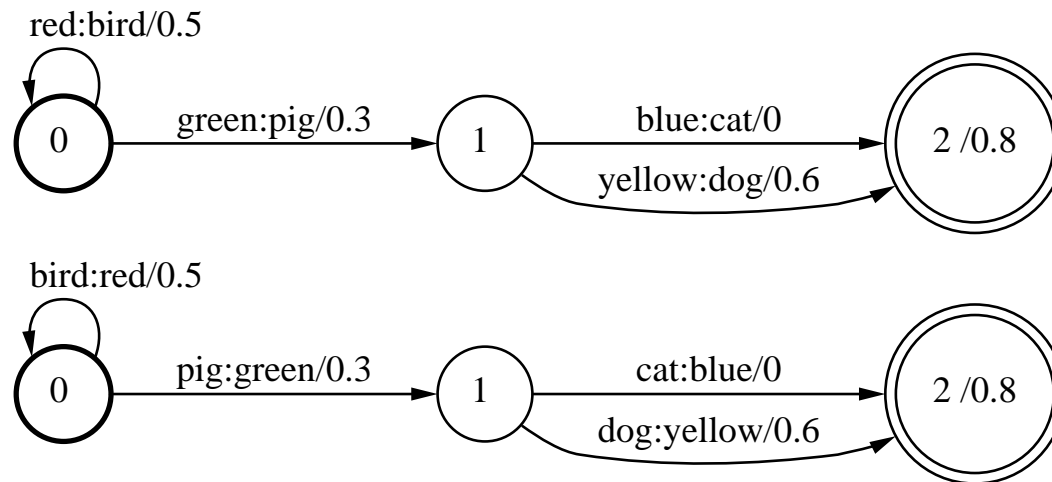
  - Lazy implementation (see table).

# Reversal – Illustration

- Definition: $[\![\widetilde{T}]\!](x, y) = [\![T]\!](\widetilde{x}, \widetilde{y})$

- Example:

# Inversion – Illustration

- Definition: $[\![T^{-1}]\!](x, y) = [\![T]\!](y, x)$

- Example:

# Projection – Illustration

- Definition: $[\![\Pi_1(T)]\!](x) = \bigoplus_y [\![T]\!](x,y)$

- Example:

# Some Fundamental Binary Operations – Algorithms

- Definitions

| OPERATION | DEFINITION AND NOTATION | CONDITION |
|---|---|---|
| Composition | $[\![T_1 \circ T_2]\!](x, y) = \bigoplus_z [\![T_1]\!](x, z) \otimes [\![T_2]\!](z, y)$ | $\mathbb{K}$ commutative |
| Intersection | $[\![A_1 \cap A_2]\!](x) = [\![A_1]\!](x) \otimes [\![A_2]\!](x)$ | $\mathbb{K}$ commutative |
| Difference | $[\![A_1 - A_2]\!](x) = [\![A_1 \cap \overline{A_2}]\!](x)$ | $A_2$ unweighted & deterministic |

- Complexity and implementation
  - Complexity (quadratic): $O((|E_1| + |Q_1|)\,(|E_2| + |Q_2|))$.
  - Path multiplicity in presence of $\epsilon$-transitions: $\epsilon$-filter.
  - Lazy implementation.

# Composition – Illustration

- Definition: $\llbracket T_1 \circ T_2 \rrbracket(x, y) = \bigoplus_z \llbracket T_1 \rrbracket(x, z) \otimes \llbracket T_2 \rrbracket(z, y)$
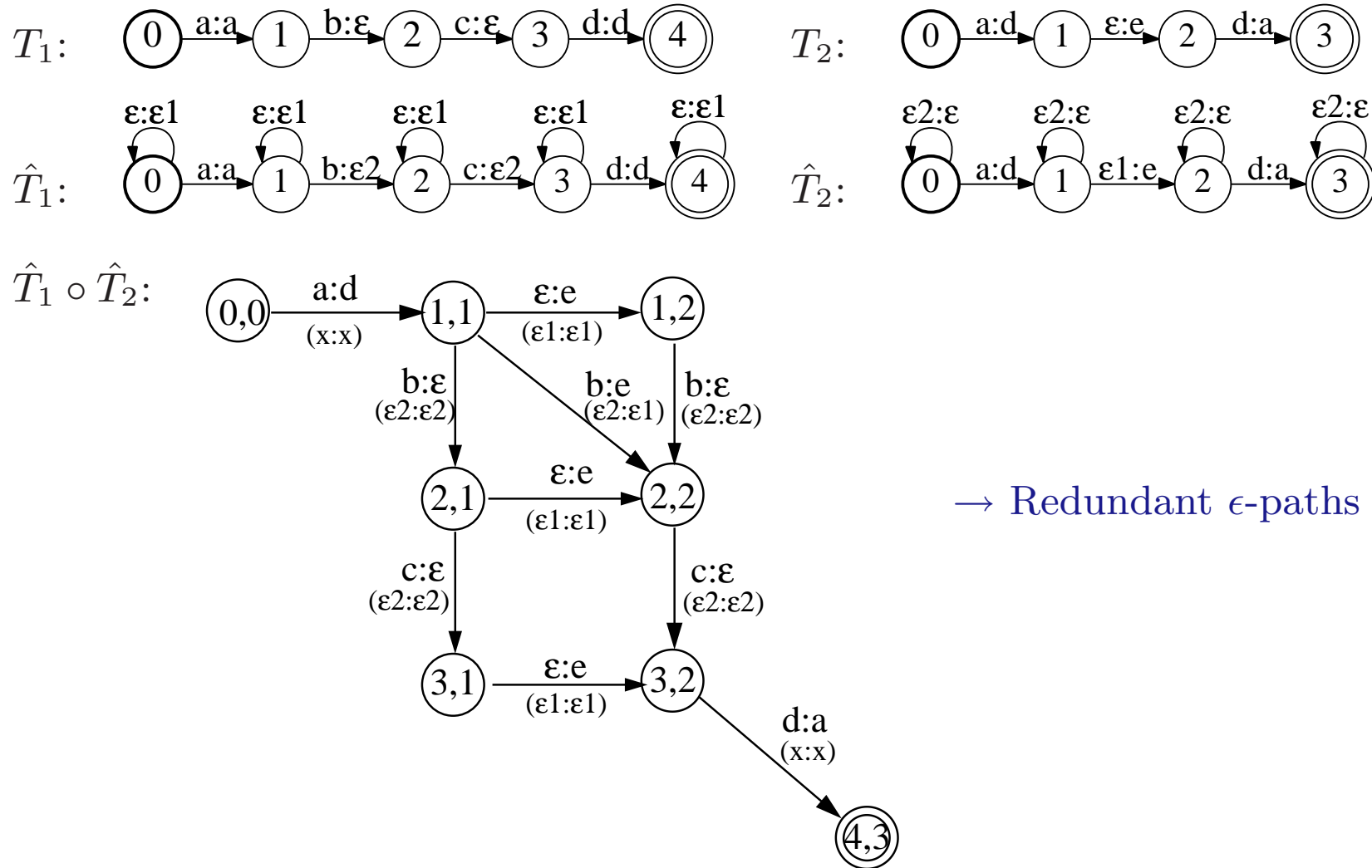
- Example:

# Composition – Pseudocode

COMPOSITION$(T_1, T_2)$
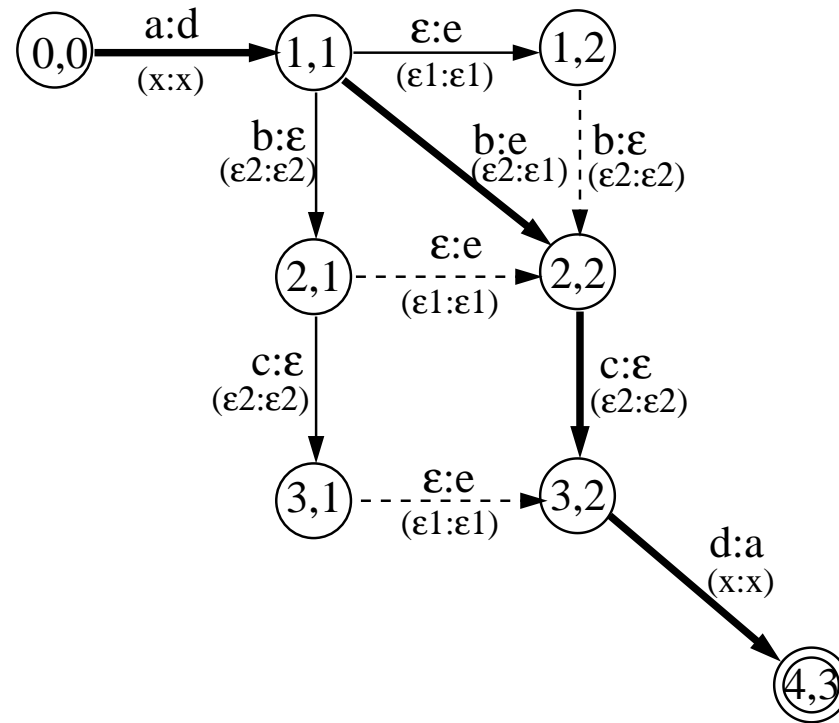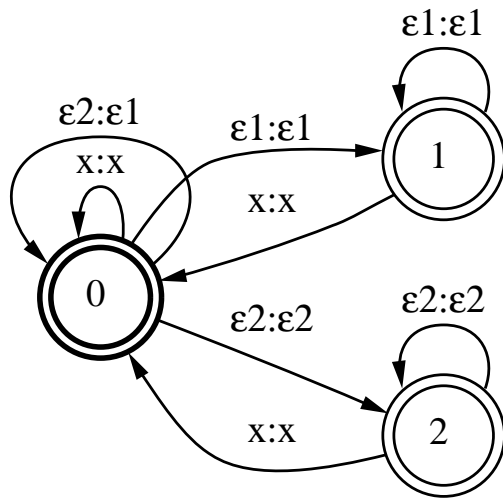
1    $S \leftarrow Q \leftarrow I_1 \times I_2$

2    **while** $S \neq \emptyset$ **do**

3        $(q_1, q_2) \leftarrow$ HEAD$(S)$

4        DEQUEUE$(S)$

5        **if** $(q_1, q_2) \in I_1 \times I_2$ **then**

6            $I \leftarrow I_1 \times I_2$

7            $\lambda(q_1, q_2) \leftarrow \lambda_1(q_1) \otimes \lambda_2(q_2)$

8        **if** $(q_1, q_2) \in F_1 \times F_2$ **then**

9            $F \leftarrow F \cup \{(q_1, q_2)\}$

10          $\rho(q_1, q_2) \leftarrow \rho_1(q_1) \otimes \rho_2(q_2)$

11      **for** each$(e_1, e_2)$ such that $o[e_1] = i[e_2]$ **do**

12         **if** $(n[e_1], n[e_2]) \notin Q$ **then**

13            $Q \leftarrow Q \cup \{(n[e_1], n[e_2])\}$

14            ENQUEUE$(S, (n[e_1], n[e_2]))$

15        $E \leftarrow E \cup \{((q_1, q_2), i[e_1], o[e_2], w[e_1] \otimes w[e_2], (n[e_1], n[e_2]))\}$

16    **return** $T = (\Sigma, \Delta, Q, I, F, E, \lambda, \rho)$

# Multiplicity & $\epsilon$-Transitions – Problem



$\rightarrow$ Redundant $\epsilon$-paths
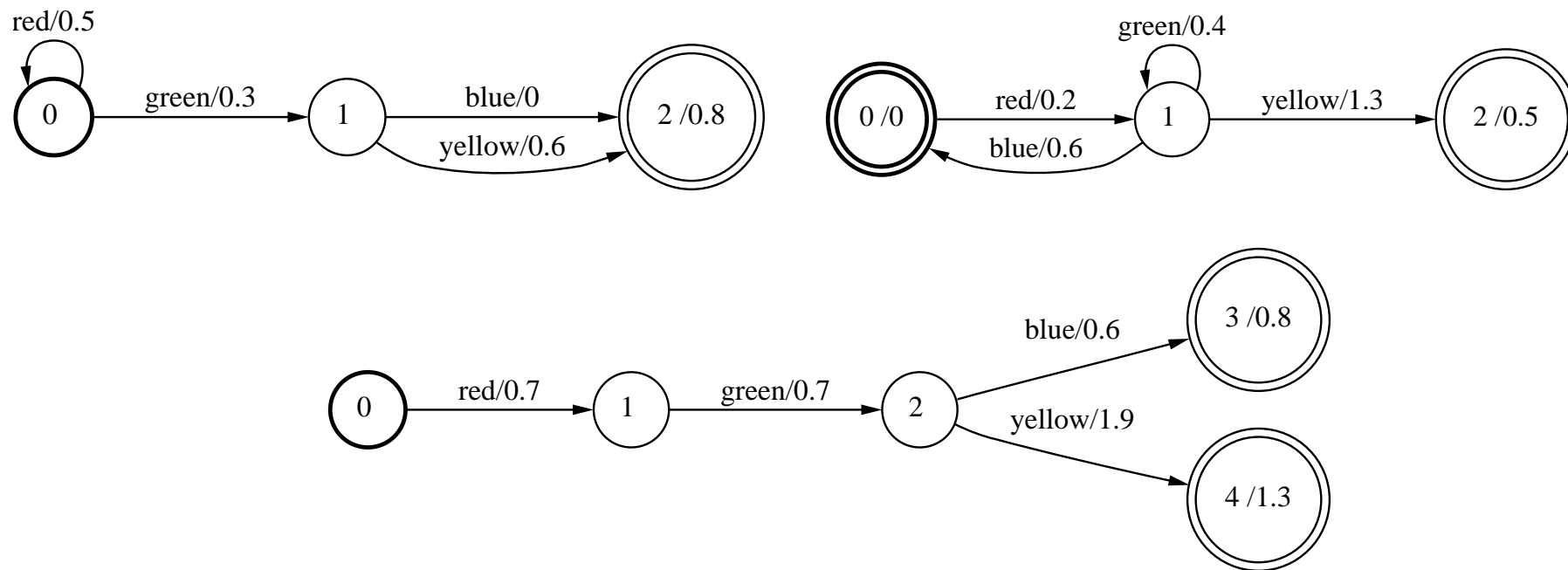
# Solution – Filter $F$ for Composition

Filter $F$



$\rightarrow$ Replace $T_1 \circ T_2$ by $\hat{T}_1 \circ F \circ \hat{T}_2$.

# Intersection – Illustration

- Definition: $[\![A_1 \cap A_2]\!](x) = [\![A_1]\!](x) \otimes [\![A_2]\!](x)$

- Example:

# Difference – Illustration

- Definition: $[\![A_1 - A_2]\!](x) = [\![A_1 \cap \overline{A_2}]\!](x)$

- Example:

# Optimization Algorithms – Overview

- Definitions

| OPERATION | DESCRIPTION |
|---|---|
| Connection | Removes non-accessible/non-coaccessible states |
| $\epsilon$-Removal | Removes $\epsilon$-transitions |
| Determinization | Creates equivalent deterministic machine |
| Pushing | Creates equivalent pushed/stochastic machine |
| Minimization | Creates equivalent minimal deterministic machine |

- Conditions: There are specific semiring conditions for the use of these algorithms. Not all weighted automata or transducers can be determinized using that algorithm.

# Connection – Algorithm

- Definition

  - Input: weighted transducer $T_1$.

  - Output: weighted transducer $T_2 \equiv T_1$ with all states connected.
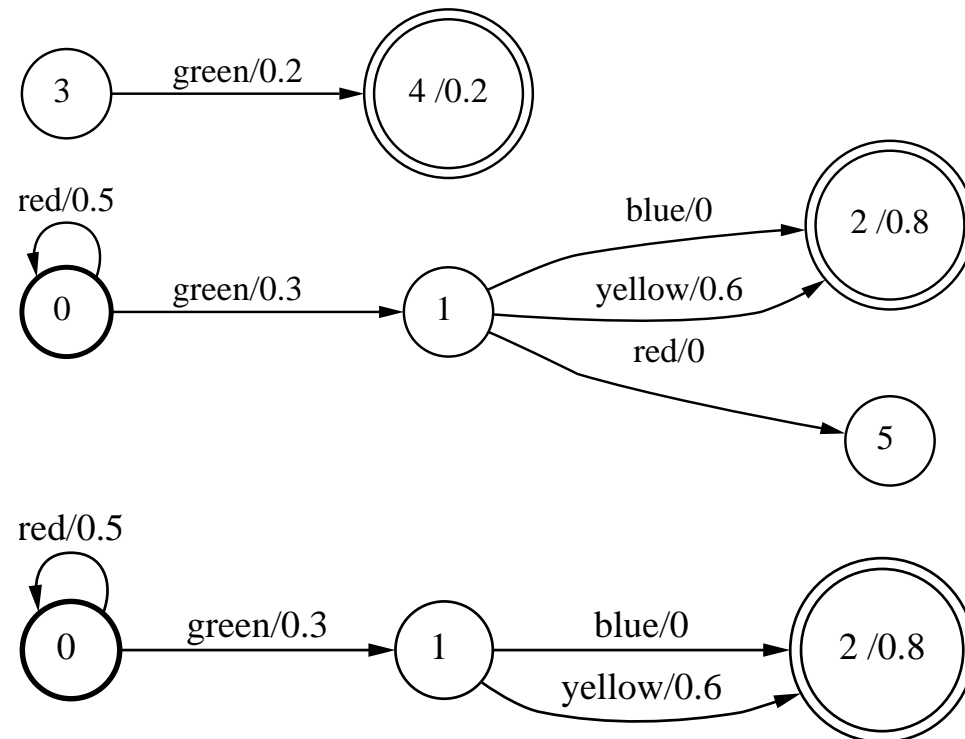
- Description

  1. Depth-first search of $T_1$ from $I_1$.

  2. Mark accessible and coaccessible states.

  3. Keep marked states and corresponding transitions.

- Complexity and implementation

  - Complexity (linear): $O(|Q_1| + |E_1|)$.

  - No natural lazy implementation.

# Connection – Illustration

- Definition: Removes non-accessible/non-coaccessible states

- Example:

# $\epsilon$-Removal – Algorithm

- **Definition**

  – Input: weighted transducer $T_1$ with $\epsilon$-transitions.

  – Output: weighted transducer $T_2 \equiv T_1$ with no $\epsilon$-transition.

- **Description** (two stages):

  1. **Computation of $\epsilon$-closures**: for any state $p$, states $q$ that can be reached from $p$ via $\epsilon$-paths and the total weight of the $\epsilon$-paths from $p$ to $q$.

  $$C[p] = \left\{ (q, w) : q \in \epsilon[p],\ d[p, q] = w \neq \overline{0} \right\}$$

  with:

  $$d[p, q] = \bigoplus_{\pi \in P(p, \epsilon, q)} w[\pi]$$

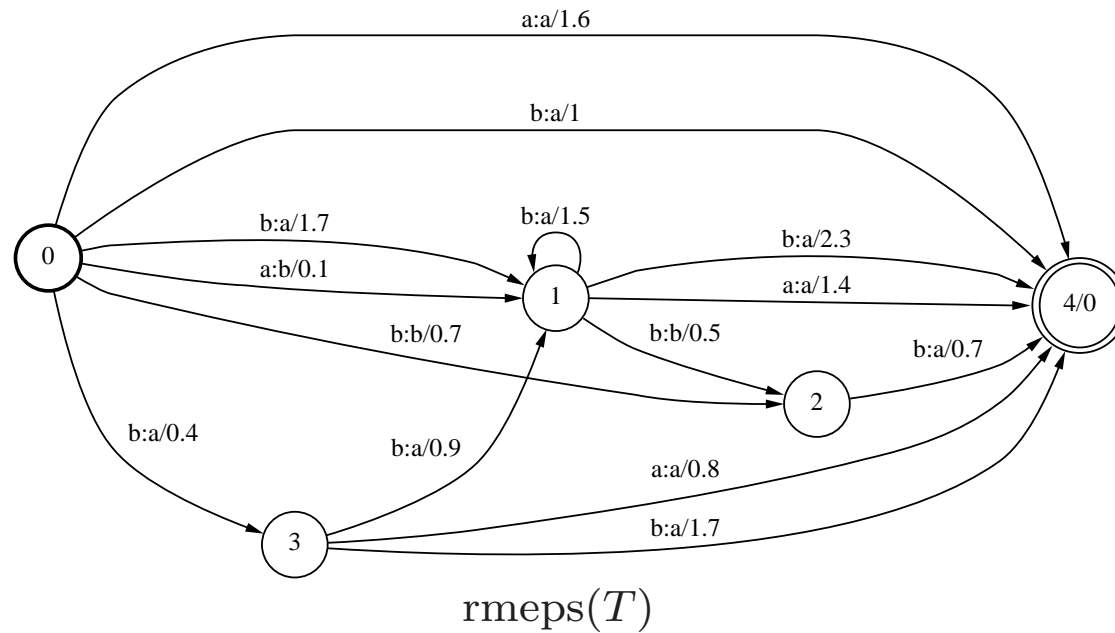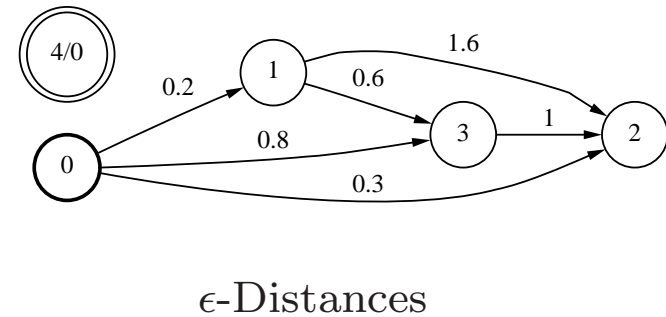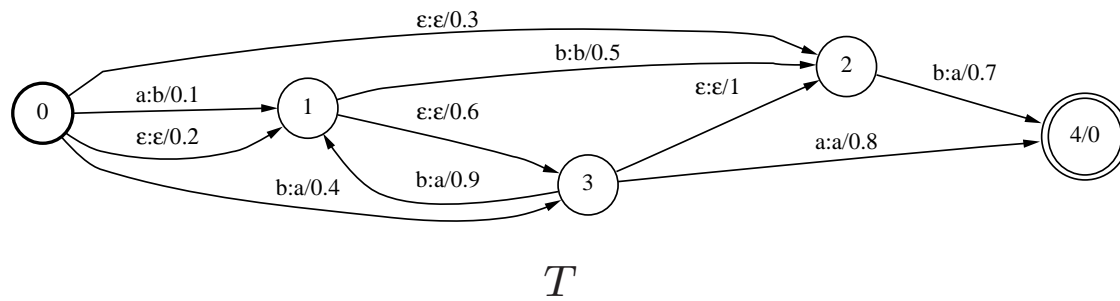  2. **Removal of $\epsilon$'s**: actual removal of $\epsilon$-transitions and addition of new transitions.

  $\implies$ All-pair $\mathbb{K}$-shortest-distance problem in $T_\epsilon$ ($T$ reduced to its $\epsilon$-transitions).

- Complexity and implementation

  - All-pair shortest-distance algorithm in $T_\epsilon$.
    * $k$-Closed semirings (for $T_\epsilon$) or approximation: generic sparse shortest-distance algorithm [See references].
    * Closed semirings: Floyd-Warshall or Gauss-Jordan elimination algorithm with decomposition of $T_\epsilon$ into strongly connected components [See references],
      space complexity (quadratic): $O(|Q|^2 + |E|)$.
      time complexity (cubic): $O(|Q|^3(T_\oplus + T_\otimes + T_*))$.

  - Complexity:
    * Acyclic $T_\epsilon$: $O(|Q|^2 + |Q||E|(T_\oplus + T_\otimes))$.
    * General case (tropical semiring): $O(|Q||E| + |Q|^2 \log |Q|)$.

  - Lazy implementation: integration with on-the-fly weighted determinization.

# $\epsilon$-Removal – Illustration

- **Definition:** Removes $\epsilon$-transitions

- **Example:**



$T$

$\epsilon$-Distances

$\mathrm{rmeps}(T)$

# Determinization – Algorithm

- Definition
  - Input: *determinizable* weighted automaton or transducer $M_1$.
  - Output: $M_2 \equiv M_1$ *subsequential* or *deterministic*: $M_2$ has a unique initial state and no two transitions leaving the same state share the same input label.

- Description
  1. Generalization of subset construction: weighted subsets $\{(q_1, w_1), \ldots, (q_n, w_n)\}$, $w_i$ remainder weight at state $q_i$.
  2. Weight of a transition in the result: $\oplus$-sum of the original transitions pre-$\otimes$-multiplied by remainders.
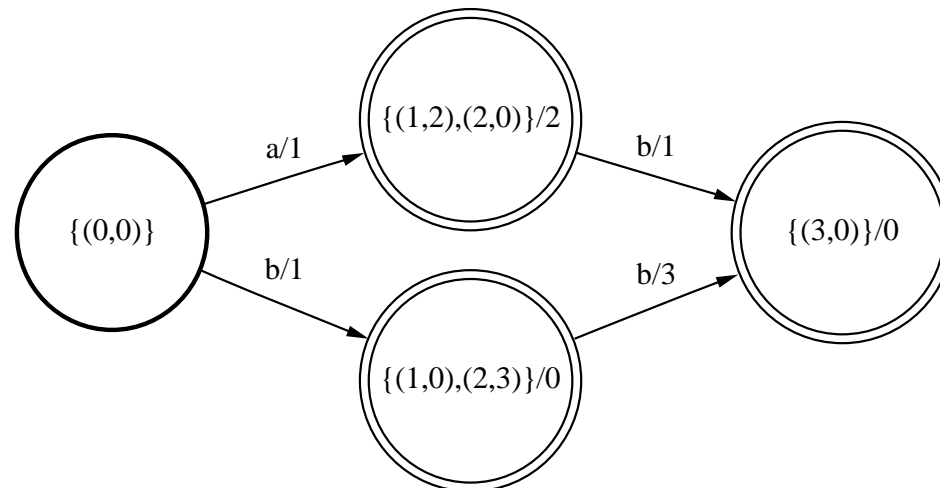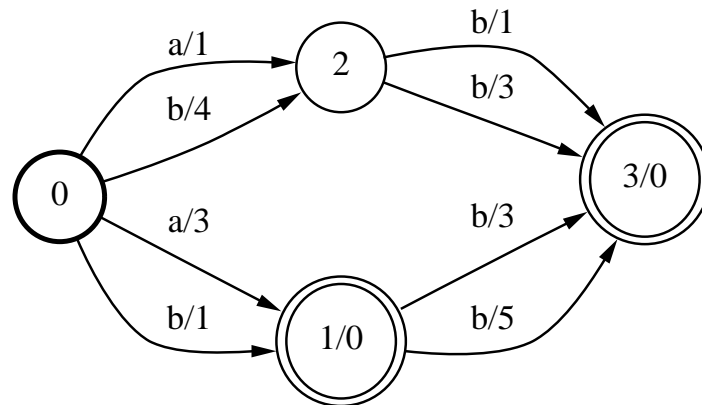
- Conditions
  - Semiring: weakly left divisible semirings.
  - $M$ is determinizable $\equiv$ the determinization algorithm applies to $M$.
  - All unweighted automata are determinizable.
  - All acyclic machines are determinizable.

– Not all weighted automata or transducers are determinizable.

– Characterization based on the *twins property*.

- Complexity and Implementation

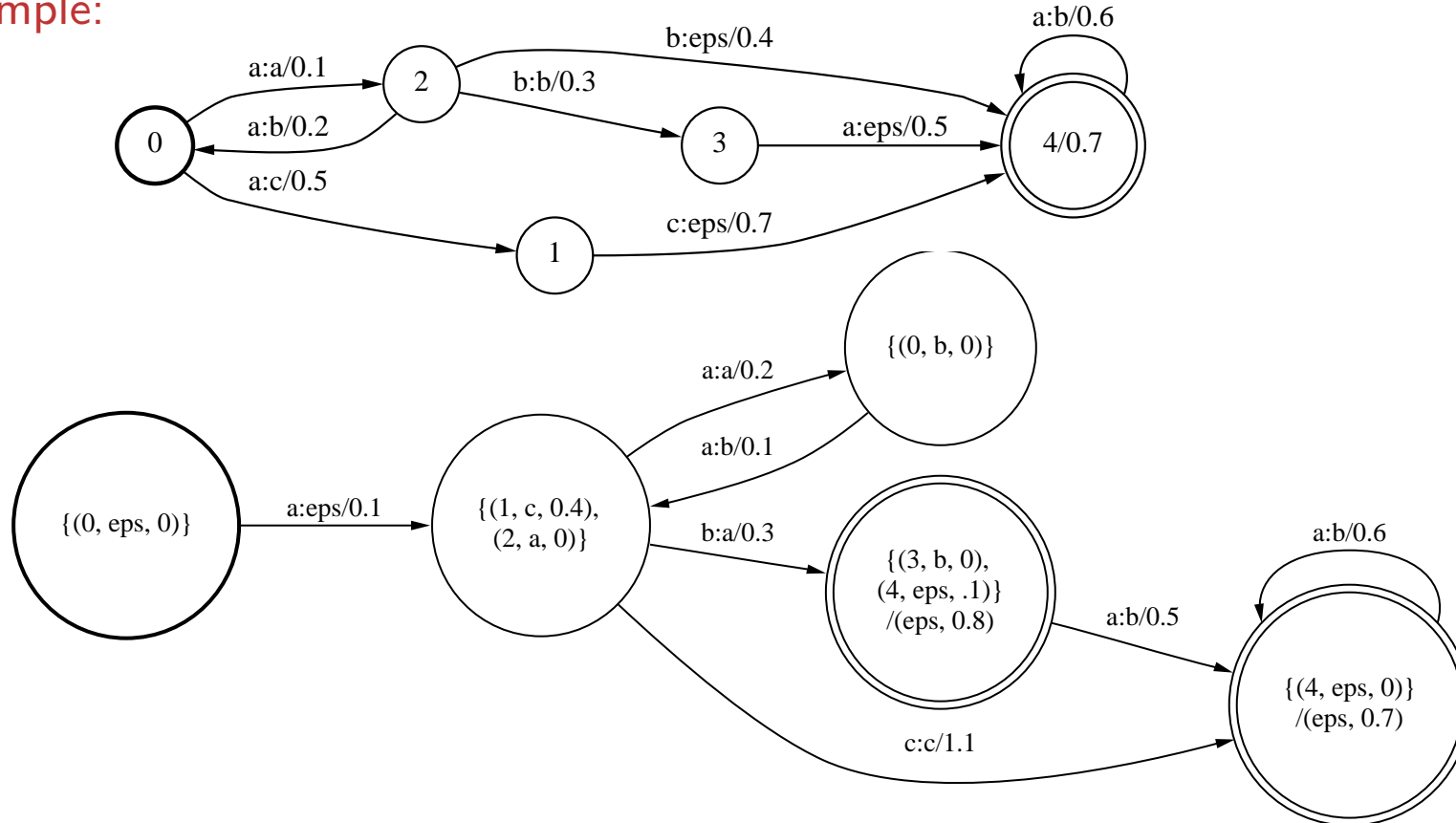– Complexity: exponential.

– Lazy implementation.

# Determinization of Weighted Automata – Illustration

- **Definition:** Creates an equivalent deterministic weighted automaton

- **Example:**

# Determinization of Weighted Transducers – Illustration

- **Definition:** Creates an equivalent deterministic weighted transducer

- **Example:**

# Determinization of Weighted Automata – Pseudocode

$\textsc{Determinization}(A)$

1    $i' \leftarrow \{(i, \lambda(i)) : i \in I\}$

2    $\lambda'(i') \leftarrow \overline{1}$

3    $S \leftarrow \{i'\}$

4    **while** $S \neq \emptyset$ **do**

5         $p' \leftarrow \textsc{Head}(S)$

6         $\textsc{Dequeue}(S)$

7        **for** each $x \in i[E[Q[p']]]$ **do**

8              $w' \leftarrow \bigoplus \{v \otimes w : (p, v) \in p', (p, x, w, q) \in E\}$

9              $q' \leftarrow \{(q, \bigoplus \{w'^{-1} \otimes (v \otimes w) : (p, v) \in p', (p, x, w, q) \in E\}) :$

                     $q = n[e], i[e] = x, e \in E[Q[p']]\}$

10        $E' \leftarrow E' \cup \{(p', x, w', q')\}$

11        **if** $q' \notin Q'$ **then**

12            $Q' \leftarrow Q' \cup \{q'\}$

13            **if** $Q[q'] \cap F \neq \emptyset$ **then**

14               $F' \leftarrow F' \cup \{q'\}$

15               $\rho'(q') \leftarrow \bigoplus \{v \otimes \rho(q) : (q, v) \in q', q \in F\}$

16            $\textsc{Enqueue}(S, q')$

17    **return** $A'$

# Pushing – Algorithm

- Definition

  - Input: weighted automaton or transducer $M_1$.

  - Output: $M_2 \equiv M_1$ such that:
    - * the longest common prefix of all outgoing paths is minimal, or
    - * the $\oplus$-sum of the weights of all outgoing transitions $= \overline{1}$ modulo the string/weight at the initial state.

- Description (two stages):

  1. Single-source shortest distance computation: for each state $q$,
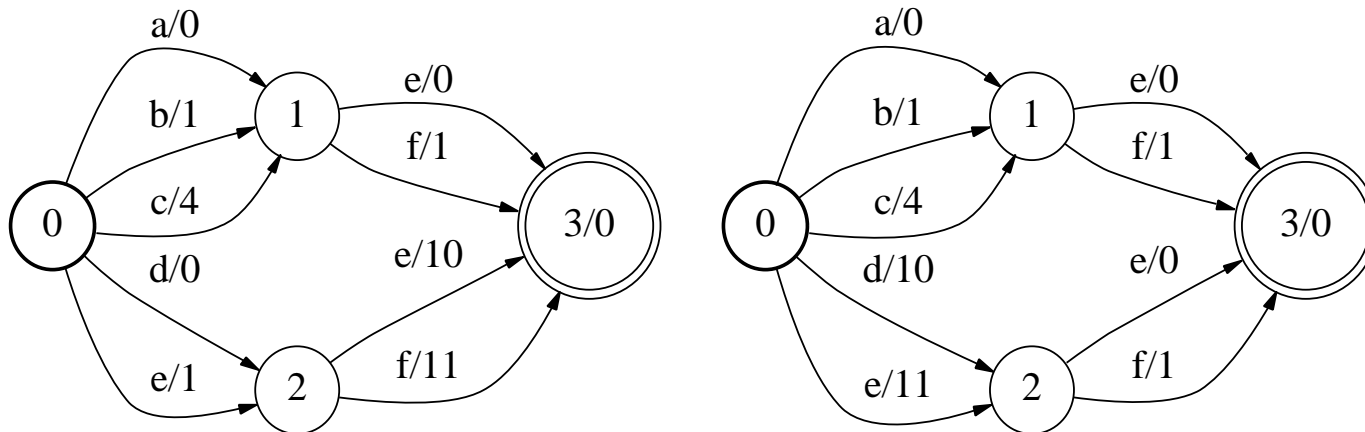
  $$d[q] = \bigoplus_{\pi \in P(q,F)} w[\pi]$$

  2. Reweighting: for each transition $e$ such that $d[p[e]] \neq \overline{0}$,

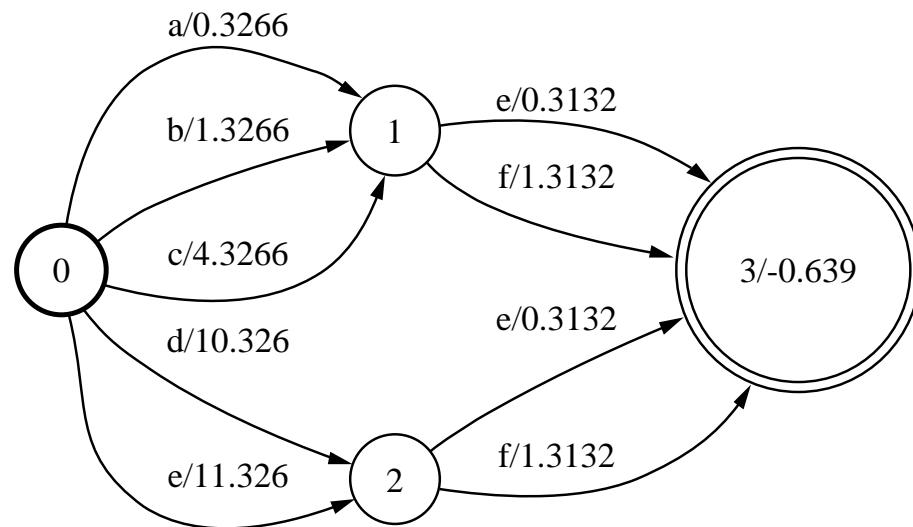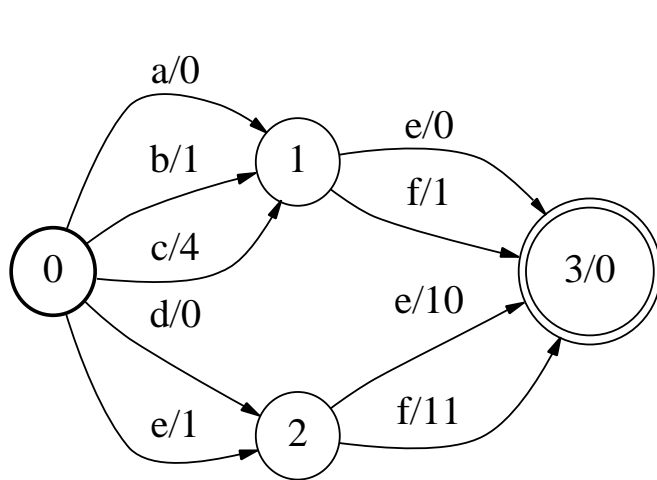  $$w[e] \leftarrow (d[p[e]])^{-1}(w[e] \otimes d[n[e]])$$

- **Conditions** (automata case)

  - Weakly divisible semiring.

  - Zero-sum free semiring or zero-sum free machine.

- **Complexity**

  - Automata case
    * Acyclic case (linear): $O(|Q| + |E|(T_\oplus + T_\otimes))$.
    * General case (tropical semiring): $O(|Q| \log |Q| + |E|)$.
  - Transducer case: $O((|P_{max}| + 1) |E|)$.

# Weight Pushing – Illustration

- **Definition:** Creates an equivalent pushed/stochastic machine

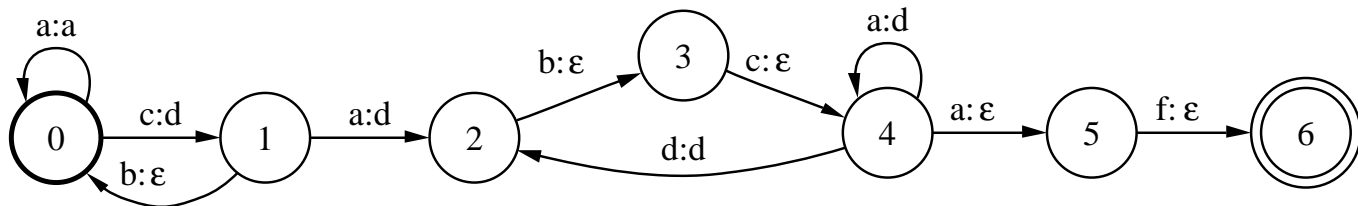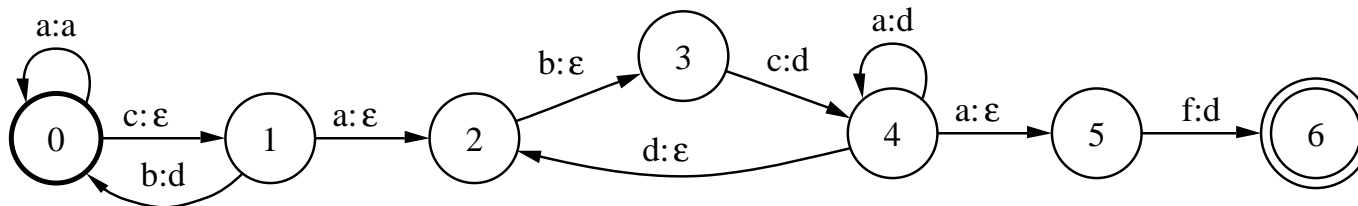- **Example:**

  - Tropical semiring

# Label Pushing – Illustration

- **Definition:** Minimizes at each state the length of the common prefix of all outgoing paths at that state.
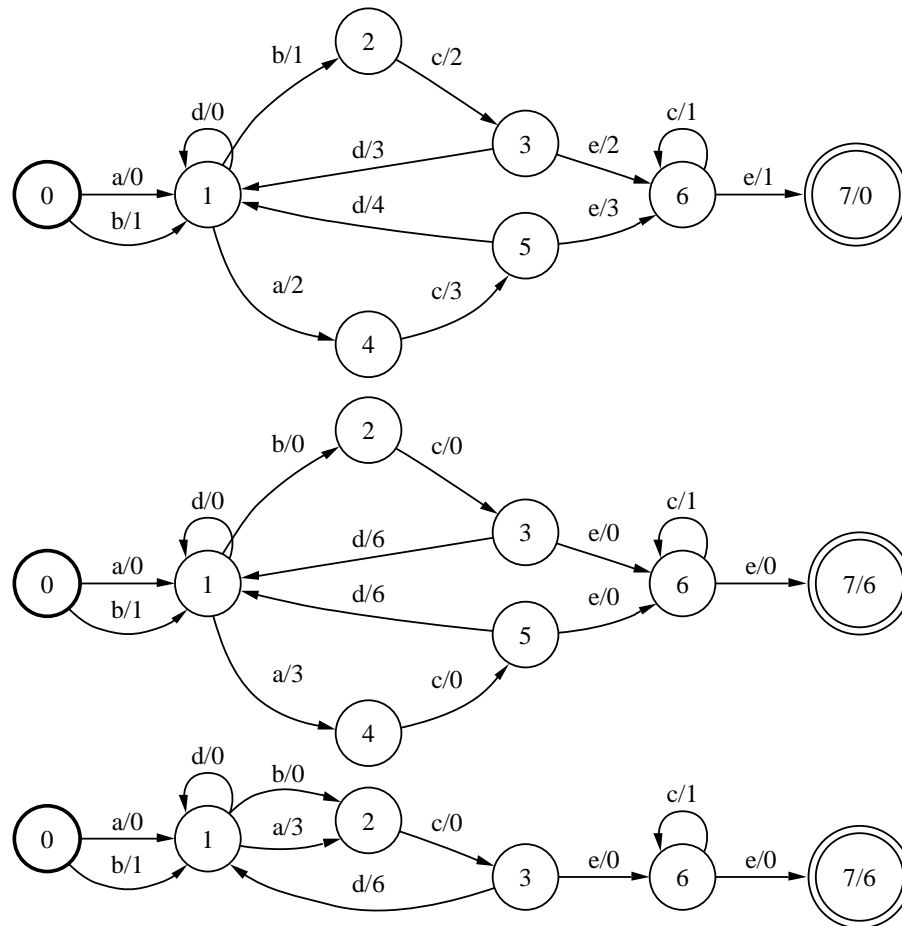
- **Example:**

# Minimization – Algorithm

- Definition

  - Input: deterministic weighted automaton or transducer $M_1$.

  - Output: deterministic $M_2 \equiv M_1$ with minimal number of states and transitions.

- Description: two stages

  1. Canonical representation: use pushing or other algorithm to standardize input automata.

  2. Automata minimization: encode pairs (label, weight) as labels and use classical unweighted minimization algorithm.

- Complexity

  - Automata case
    * Acyclic case (linear): $O(|Q| + |E|(T_\oplus + T_\otimes))$.
    * General case (tropical semiring): $O(|E| \log |Q|)$.
  - Transducer case
    * Acyclic case: $O(S + |Q| + |E|(|P_{max}| + 1))$.
    * General case: $O(S + |Q| + |E|(\log |Q| + |P_{max}|))$.

# Minimization – Illustration

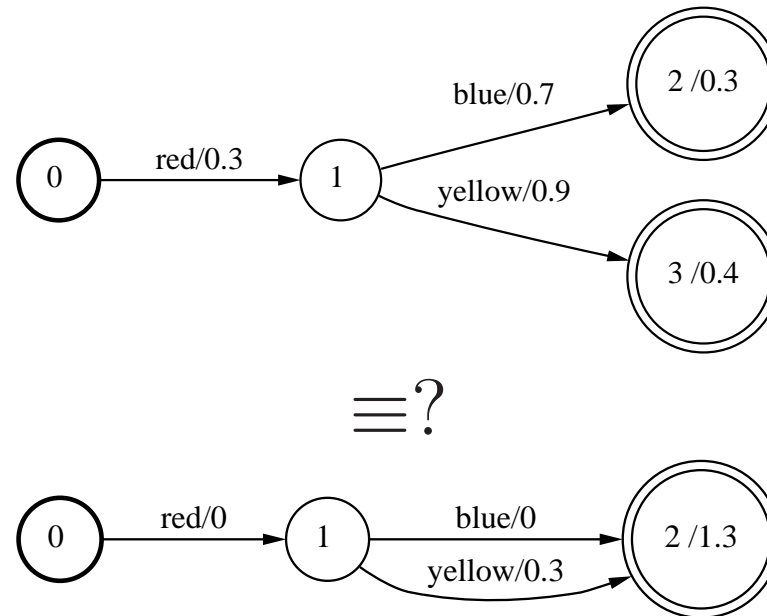- **Definition:** Computes a minimal equivalent deterministic machine

- **Example:**

# Equivalence – Algorithm

- **Definition**

  - Input: deterministic weighted automata $A_1$ and $A_2$.

  - Output: TRUE if $A_2 \equiv A_1$, FALSE otherwise.

- **Description**: two stages

  1. **Canonical representation**: use pushing or other algorithm to standardize input automata.

  2. **Test**: encode pairs (label, weight) as labels and use classical algorithm for testing the equivalence of unweighted automata.

- **Complexity**

  - First stage: $O((|E_1| + |E_2|) + (|Q_1| + |Q_2|) \log(|Q_1| + |Q_2|))$ if using pushing in the tropical semiring.

  - Second stage (quasi-linear): $O(m\, \alpha(m, n))$ where $m = |E_1| + |E_2|$ and $n = |Q_1| + |Q_2|$, and $\alpha$ is the *inverse of Ackermann's function*.

# Equivalence – Illustration

- Definition: $A_1 \equiv A_2$ iff $[\![A_1]\!](x) = [\![A_2]\!](x)$ for all $x$

- Graphical Representation:

# Single-Source Shortest-Distance Algorithms – Algorithm

- Generic single-source shortest-distance algorithm

  - Definition: for each state $q$,
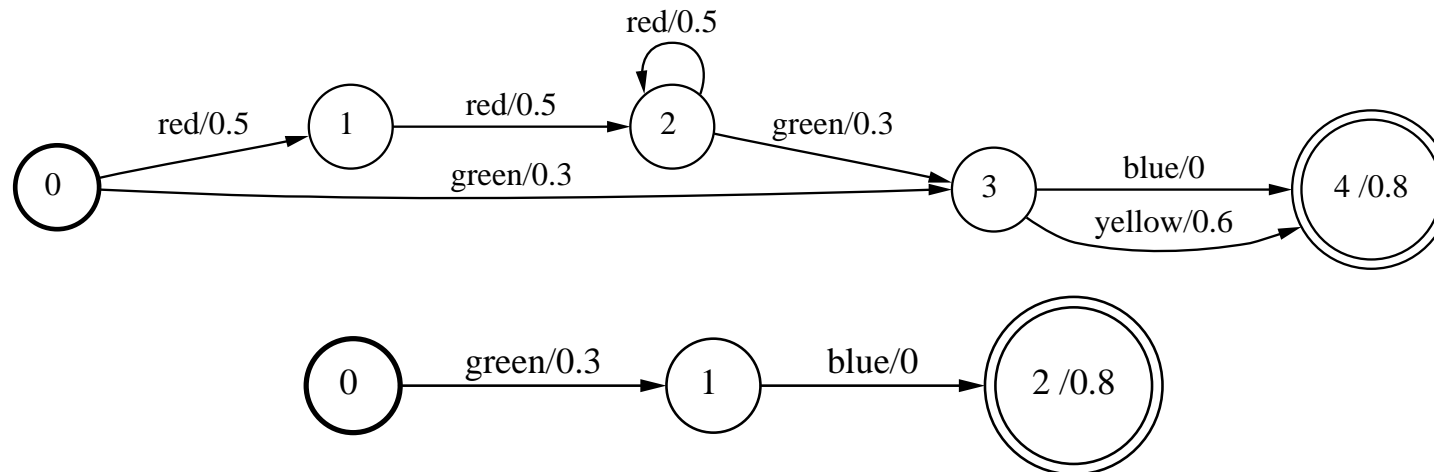
  $$d[q] = \bigoplus_{\pi \in P(q,F)} w[\pi]$$

  - Works with any queue discipline and any semiring $k$-closed for the graph.

  - Coincides with classical algorithms in the specific case of the tropical semiring and the specific queue disciplines: shortest-first (Dijkstra), FIFO (Bellman-Ford), or topological sort order (Lawler).

- $N$-shortest paths algorithm

  - General $N$-shortest paths algorithm augmented with the computation of the potentials.

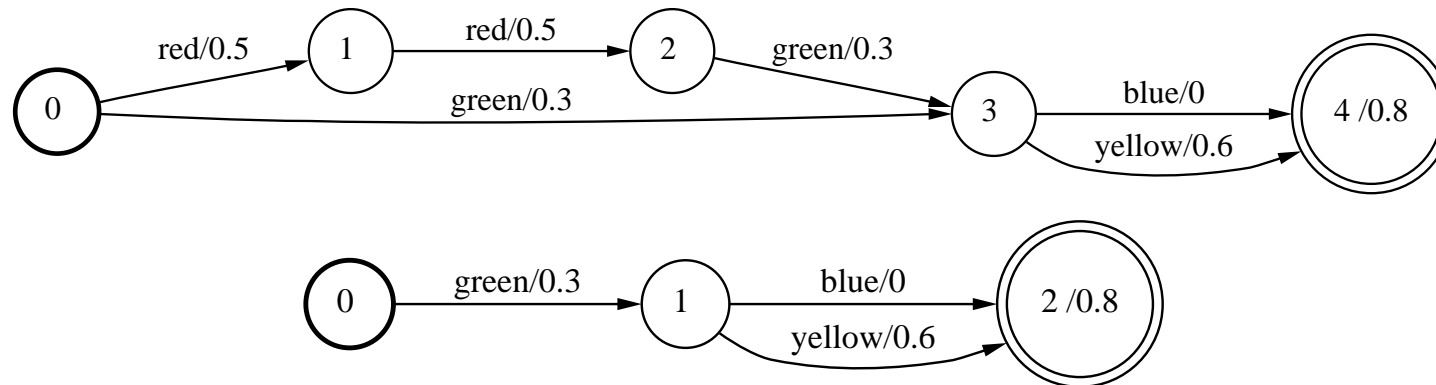  - On-the-fly weighted determinization for $n$-shortest strings.

# $N$-Shortest Paths – Illustration

- **Definition:** Computes the $N$-shortest paths in the input machine

- **Condition:** Semiring needs to have the path property: $a \oplus b \in \{a, b\}$ (e.g. tropical semiring)

- **Example:**

# Pruning – Illustration

- **Definition:** Removes any paths which weight is more than the shortest-distance $\otimes$-multiply by a specified threshold

- **Condition:** Semiring needs to be commutative and have the path property: $a \oplus b \in \{a, b\}$ (e.g. tropical semiring)

- **Example:**

# String Algorithms – Overview

- How to implement some fundamental string algorithms using the operations previously described:

  - Counting patterns (e.g. $n$-grams) in automata

  - Pattern matching using automata

  - Compiling regular expression into automata

- Benefits: generality, efficiency and flexibility

# Counting from weighted automata

- Expected count of $x$ in $A$:

  Let $A$ be a weighted automaton over the probability semiring,

  $$c(x) = \sum_{u \in \Sigma^*} |u|_x \, [\![A]\!](u)$$

  where:

  - $|u|_x$: number of occurrences of $x$ in $u$

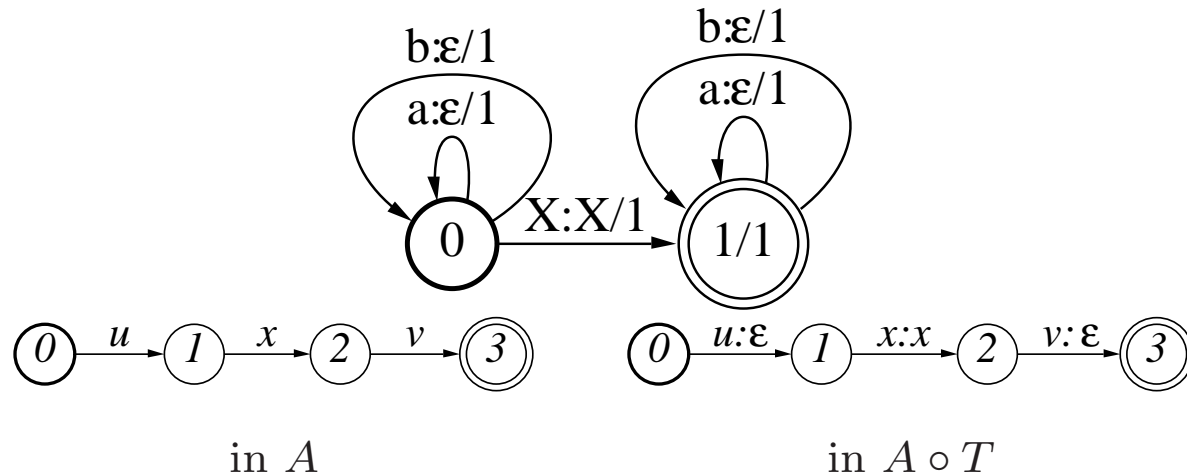  - $[\![A]\!](u)$: weight associated to $u$ by $A$
    $\rightarrow \Pr(u)$ if $A$ is pushed

- Condition:

  The weight of any cycle in $A$ should be less than 1.
  This is the case if $A$ represents a probability distribution.

# Counting by composition with a transducer

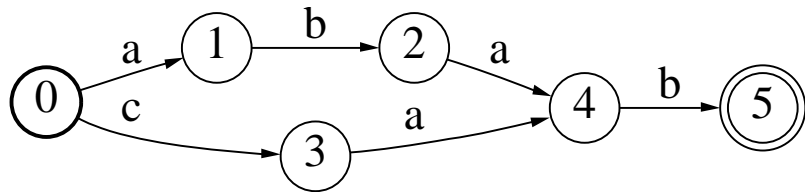- Counting transducer $T$ for set of sequences $X$ with $\Sigma = \{a, b\}$:



in $A$           in $A \circ T$

To each successful path $\pi$ in $A$ and each occurence of $x$ along $\pi$
$\rightarrow$ corresponds a successful path with output $x$ in $A \circ T$
$\rightarrow$ $c(x)$ is the sum of the weight of all the successful path
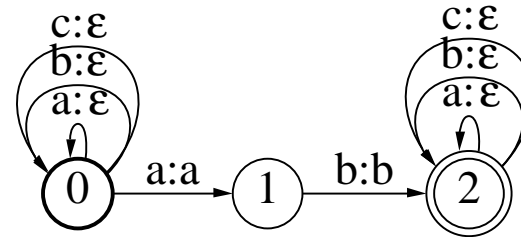     with output $x$ in $A \circ T$

- Theorem:

Let $\Pi_2$ denote projection onto output labels. For all $x \in X$,
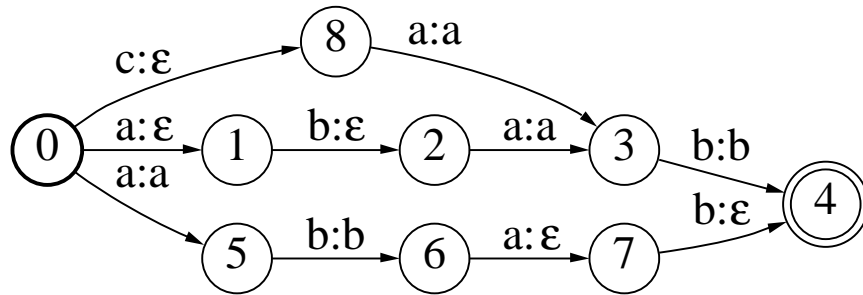
$$c(x) = [\![\Pi_2(A \circ T)]\!](x)$$

# Counting with transducers – Example



Automaton $A$
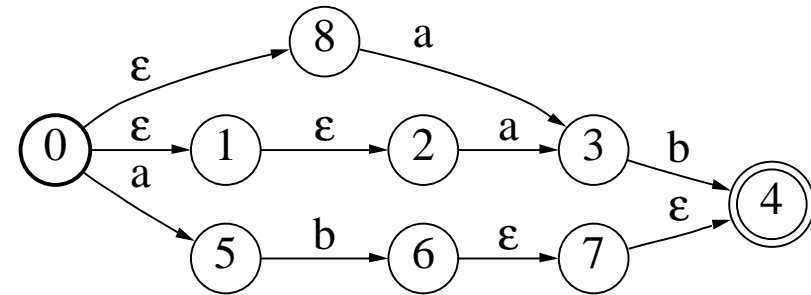
Counting transducer $T$

$A \circ T$

$\Pi_2(A \circ T)$
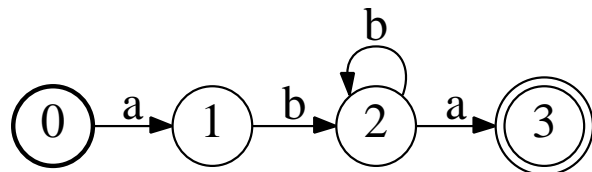
$$[\![\Pi_2(A \circ T)]\!](ab) = 3 = c(ab)$$

# Local Grammar – Algorithm

[Mohri, 94]
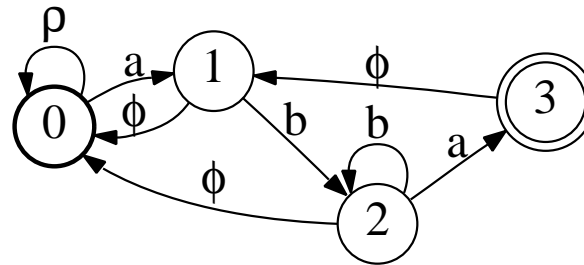
- Definition

  - Input: a deterministic finite automaton $A$

  - Output: a compact representation of $\det(\Sigma^* A)$

- Description

  - A generalization of [Aho-Corasick, 75]

  - Failure transitions: labeled by $\phi$, non-consuming, traversed when no transition with required label is present

  - Default transitions: labeled by $\rho$, consuming, traversed when no transition with required label is present, only present at the initial state
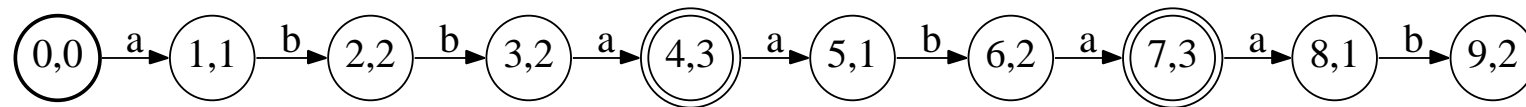
# Local Grammar – Illustration

- **Pattern matching:** find all occurences of pattern $A$ in text $T$

  $T = abbaabaab,\ A = ab^{+}a$



$$A$$

$$\det(\Sigma^{*}A)$$

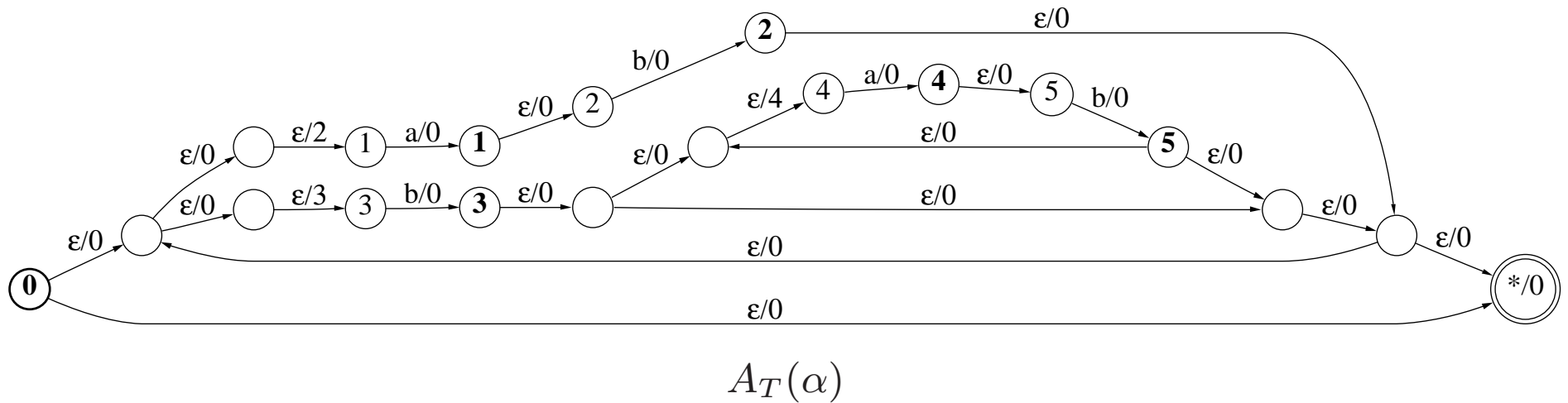$$T \cap \det(\Sigma^{*}A)$$

- **Complexity:** search time linear in $|T|$

# Regular Expression Compilation – Algorithms

- **Definition**
  - Input: a (weighted) regular expression $\alpha$
  - Output: a (weighted) automata representing $\alpha$

- **Description: Thompson construction**
  1. Build a sparse tree for $\alpha$
  2. Walk the tree bottom-up and apply the relevant rational operation at each node

- **Complexity and implementation**
  - Linear in the length of $\alpha$
  - Admits lazy implementation

- Other constructions (Glushkov, Antimirov, Follow) can be obtained from Thompson using epsilon-removal and minimization
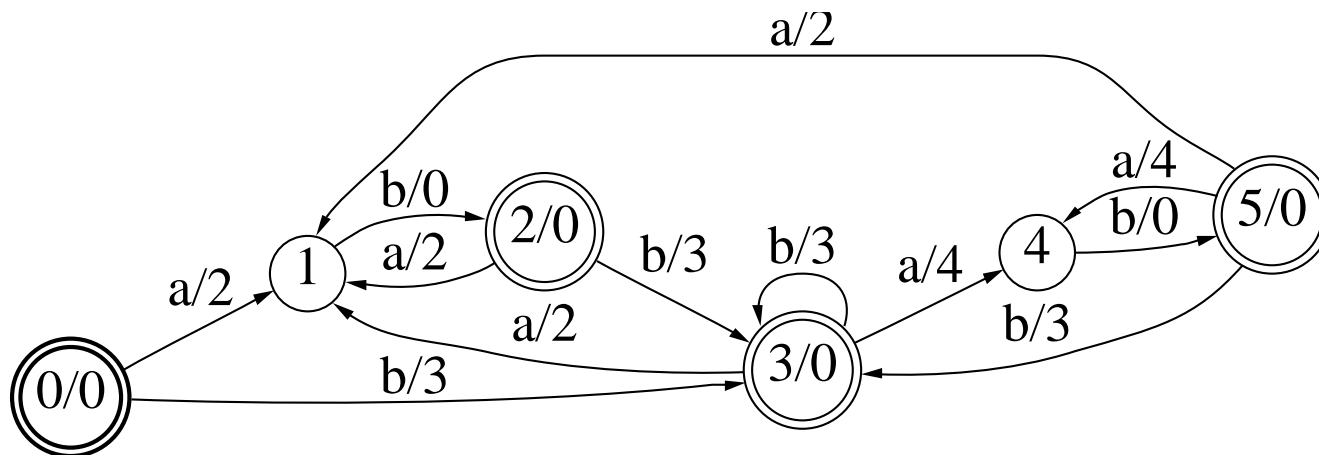
# Regular Expression Compilation – Thompson

- Regular expression: $\alpha = (2ab + 3b(4ab)^*)^*$

- Thompson automaton:



$$A_T(\alpha)$$

# Regular Expression Compilation – Glushkov

- Regular expression: $\alpha = (2ab + 3b(4ab)^*)^*$

- Glushkov automaton:



$$A_G(\alpha) = \mathrm{rmeps}(A_T(\alpha))$$