Finite state transducers, String-to-Weight Transducers and their application in speech processing

Miloš Janda, ijanda@fit.vutbr.cz

Department of Computer Graphics and Multimedia Faculty of Information Technology Brno University of Technology

December 19, 2010

S Q P

-∢∃≯

Speech processing

2 Finite State Automatons

- Nondeterministic Finite State Automaton (NFSA)
- Deterministic Finite State Automaton (DFSA)

3 Finite State Transducers

- Sequential String-to-String Transducer
- Subsequential Transducer
- String-to-Weight Transducer
- String-to-Weight Subsequential Transducer
- Operations under WSFTs
 - Determinization
 - Minimization

5 Speech recognizer composition, Lattices

 $\mathcal{A} \mathcal{A} \mathcal{A}$

Speech processing

2 Finite State Automatons

- Nondeterministic Finite State Automaton (NFSA)
- Deterministic Finite State Automaton (DFSA)

3 Finite State Transducers

- Sequential String-to-String Transducer
- Subsequential Transducer
- String-to-Weight Transducer
- String-to-Weight Subsequential Transducer
- Operations under WSFTs
 - Determinization
 - Minimization

5 Speech recognizer composition, Lattices

 $\mathcal{A} \mathcal{A} \mathcal{A}$

I = ▶ < E ▶</p>

Speech processing

2 Finite State Automatons

- Nondeterministic Finite State Automaton (NFSA)
- Deterministic Finite State Automaton (DFSA)

3 Finite State Transducers

- Sequential String-to-String Transducer
- Subsequential Transducer
- String-to-Weight Transducer
- String-to-Weight Subsequential Transducer
- 4 Operations under WSFTs
 - Determinization
 - Minimization

5 Speech recognizer composition, Lattices

S Q P

< 글 ▶ < 글 ▶

Speech processing

2 Finite State Automatons

- Nondeterministic Finite State Automaton (NFSA)
- Deterministic Finite State Automaton (DFSA)

3 Finite State Transducers

- Sequential String-to-String Transducer
- Subsequential Transducer
- String-to-Weight Transducer
- String-to-Weight Subsequential Transducer
- 4 Operations under WSFTs
 - Determinization
 - Minimization

5 Speech recognizer composition, Lattices

S Q A

 $\Xi \rightarrow$

-∢ ≣ ▶

Speech processing

2 Finite State Automatons

- Nondeterministic Finite State Automaton (NFSA)
- Deterministic Finite State Automaton (DFSA)

3 Finite State Transducers

- Sequential String-to-String Transducer
- Subsequential Transducer
- String-to-Weight Transducer
- String-to-Weight Subsequential Transducer
- 4 Operations under WSFTs
 - Determinization
 - Minimization

5 Speech recognizer composition, Lattices

S Q A

 $\Xi \rightarrow$

-∢ ≣ ▶

Speech processing is the study of speech signals and the processing methods of these signals

Categories of speech processing

- ★ Speech recognition (LVCSR)
- ★ Keyword spotting (KWS)
- ★ Languge (LID) and speaker (SID) identification
- ★ Speaker verification
- ★ Phoneme detection

S Q A

- < E ▶ < E ▶

Spech is main form of human communication

Aplications

- ★ Military (LID, SID, KWS)
- \star Telephony and other domains
- ★ Remote controlling
- ★ Teaching (Dictionaries, Pronunciation dicts)

* ...

~ Q (?

·< 토 ▶ < 토 ▶

Simple schema of ASR



Figure 1: Recognition system schema.

590

Overall schema of ASR



Figure 2: Overall schema of speech recognizer.

 $\mathcal{O} \mathcal{Q} \mathcal{O}$

Finite State Automaton (nondeterministic)

```
Let M is a 5-tuple= (Q, \Sigma, \delta, i, F) where
```

- *Q* is a finite set of states
- Σ is an input alphabet
- δ is a transition function that maps state-symbol pairs to set of states:

$$\delta \colon \boldsymbol{Q} imes \Sigma o \mathbf{2}^{\boldsymbol{Q}}$$

 $(\delta \colon \boldsymbol{Q} imes \Sigma \cup \{\epsilon\} o \mathbf{2}^{\boldsymbol{Q}})$

- *i* is initial state: $i \in Q$
- *F* is set of accepting (final) states: $F \subseteq Q$

S Q A

Finite State Automaton (deterministic)

```
Let M is a 5-tuple= (Q, \Sigma, \delta, i, F)
where
```

- *Q* is a finite set of states
- Σ is an input alphabet
- δ is a transition function that maps state-symbol pairs to state:

$$\delta\colon \boldsymbol{Q}\times\boldsymbol{\Sigma}\to\boldsymbol{Q}$$

- *i* is initial state: $i \in Q$
- *F* is set of accepting (final) states: $F \subseteq Q$

S Q P

글▶ ∢ 글▶

Sequential String-to-String Transducer

```
Let T is a 7-tuple= (Q, \Sigma, \Delta, i, F, \delta, \sigma)
where
```

- *Q* is a finite set of states
- Σ and Δ , finite sets corresponding to input and output alphabets
- $i \in Q$ is the initial state
- $F \subseteq Q$, the set of final states
- δ , the state transition function which maps $Q \times \Sigma \rightarrow Q$
- σ , the output function which maps $Q \times \Sigma \to \Delta$

500

Example of Sequential String-to-String Transducer



Figure 3: Sequential String-to-String Transducer.

Miloš Janda (FIT BUT)

Speech processing

2010 10/33

590

王

Subsequential Transducer

Generalized sequential transducer with possibility of generating an additional output string at final states.

p-Subsequential Transducer

Same as Subsequential transducer, but limited to generate at most *p* final output strings at each final state.

500

Example of 2-Subsequential Transducer



Figure 4: 2-Subsequential Transducer.

590

王

<日 ▶ < □ ▶ < □ ▶

< □ ▶

String-to-Weight Transducer (Weighted Acceptor - WFSA)

String-to-Weight Transducer

```
Let T is a 7-tuple= (Q, \Sigma, i, F, E, \lambda, \rho) with
```

- *Q* a finite set of states
- Σ the input alphabet
- $i \in Q$ is the initial state
- $F \subseteq Q$, the set of final states
- $E \subseteq Q \times \Sigma \times \mathbb{R}_+ \times Q$ a finite set of transitions
- λ the initial weight function mapping *i* to \mathbb{R}_+
- ρ the final weight function mapping F to \mathbb{R}_+

 $\mathcal{A} \mathcal{A} \mathcal{A}$

String-to-Weight Transducer (Weighted Acceptor - WFSA)

Transition function

Transition function δ mapping $Q \times \Sigma$ to 2^Q :

$$orall (\boldsymbol{q}, \boldsymbol{a}) \in \boldsymbol{Q} imes \boldsymbol{\Sigma}, \ \delta(\boldsymbol{q}, \boldsymbol{a}) = \{ \boldsymbol{g}^{'} | \exists \boldsymbol{x} \in \mathbb{R}_{+} \colon (\boldsymbol{q}, \boldsymbol{a}, \boldsymbol{x}, \boldsymbol{q}^{'}) \in \boldsymbol{E} \}$$

Output function

Output function σ assigning to each transition its weight:

$$\forall t = (p, a, x, q) \in E, \sigma(t) = x$$

 $\mathcal{A} \mathcal{A} \mathcal{A}$

3

- 4 回 ト 4 三 ト 4 三 ト

String-to-Weight Transducer (Weighted Acceptor - WFSA)

Path

A path π is a sequence of successive transitions from one state to another:

$$\pi = ((q_0, a_0, x_0, q_1), ..., (q_{m-1}, a_{m-1}, x_{m-1}, q_m))$$

with $(q_i, a_i, x_i, q_{i+1}) \in E$

The label of a path π is $a_0...a_{m-1}$. The set of paths form q to q' labeled with input string w is abbreviated by:

$$q \rightsquigarrow^w q'$$

500

'로▶ ◀ 로▶

Example of Weighted Acceptor (WFSA)



Figure 5: Weighted Acceptors (WFSA).

590

王

String-to-Weight Subsequential Transducer (Weighted Transducer - WFST)

String-to-Weight Subsequential Transducer

Let T is a 8-tuple= $(Q, \Sigma, \Omega, i, F, E, \lambda, \rho)$ with

- *Q* a finite set of states
- Σ the input alphabet
- Ω the output alphabet
- $i \in Q$ is the initial state
- $F \subseteq Q$, the set of final states
- $E \subseteq Q \times \Sigma \times \Omega \times \mathbb{R}_+ \times Q$ a finite set of transitions
- λ the initial weight function mapping *i* to \mathbb{R}_+
- ρ the final weight function mapping F to \mathbb{R}_+

 $\mathcal{A} \mathcal{A} \mathcal{A}$

32

▲□▶ ▲□▶ ▲ □▶ ▲ 亘▶

Example of Weighted Transducer (WFST)



Figure 6: Weighted Transducer (WFST).

5900

E

<ロ> < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Determinization

The determinization algorithm takes a non-subsequential string-to-weight transducer

$$T_1 = (Q_1, \Sigma, I_1, F_1, E_1, \lambda_1, \rho_1)$$

and transform it into an equivalent, subsequential string-to-weight transducer

$$T_2 = (Q_2, \Sigma, i_2, F_2, \delta_2, \sigma_2, \lambda_2, \rho_2).$$

Multiple states reachable from one state on the same input symbol are merged.

The algorithm begins with the initial states of T_1 and build new states for T_2 by iterating over non-marked states.

 $\mathcal{A} \mathcal{A} \mathcal{A}$

· < 급 > < 트 > < 트 > · 트

 $\lambda_2 := \min_{i \in I_1} \lambda_1(i)$ $i_2 := \bigcup \{(i, \lambda_1(i) - \lambda_2)\}$ initially, i_2 is the only state in Q_2 and it is unmarked 3 while there is an unmarked state q_2 in Q_2 do begin 4 mark q_2 5if there exists (q, x) in q_2 such that $q \in F_1$ then 6 add q_2 to F_2 7 $\rho_2(q_2) := \min_{q \in F_1, (q,x) \in q_2} (x + \rho_1(q))$ 8 for each input symbol a and transition t such that $out(q_2, a, t) \neq \emptyset$ do begin 9 $\sigma_2(q_2, a) := \min_{(q, x, t) \in out(q_2, a)} x + \sigma_1(t)$ 10 $\{(q', \min_{(q,x,t)\in out(q_2,a), dest(t)=q'}(x+\sigma_1(t)-\sigma_2(q_2,a))\}$ 11 $\delta_2(q_2, a) :=$ $q' \in reachable(q_2, a)$ 12if $\delta_2(q_2, a)$ is not in Q_2 then 13add $\delta_2(q_2, a)$ as an unmarked state to Q_2 14 end 15end

Figure 7: Derminization algorithm for String-to-Weight Transducer.

Miloš Janda (FIT BUT)

 $\mathcal{A} \mathcal{A} \mathcal{A}$

<ロト < 団ト < 団ト < 団ト = 三目

Example of determinization



Figure 8: Non-deterministic weighted acceptor.



Figure 9: Equivalent WFSA obtained by determinization.

Miloš Janda (FIT BUT)

Speech processing

2010 21/33

590

E

Minimization

Process of minimization of String-to-Weight Transducer consist of two steps:

- Pushing, that adapts the transducer such that its output weights are spreading differently.
- Usual minimization (same as for FSMs)

 $\mathcal{A} \mathcal{A} \mathcal{A}$

Pushing

We define a function *d* giving the minimal output necessary to reach a final state.

For a given sebsequential transducer T, for any state $q \in Q$, d is defined by:

$$d(q) = min(\sigma(q, w) + \rho(dest(q, w)))$$

Then pushing is defined by:

•
$$\lambda' = \lambda + d(i)$$

• $\forall (q, a) \in Q \times \Sigma,$
 $\sigma'(q, a) = \sigma(q, a) + d(dest(q, a)) - d(q)$
• $\forall q \in Q, \rho'(q) = 0$

Note that pushing does not change the topology of the transducer.

Example of minimization 1/2



Figure 10: Non-minimal transducer before pushing.



Figure 11: Example of pushing operation.

Miloš Janda (FIT BUT)

Speech processing

=

Ξ.

2010 24/33

S Q A

王

Example of minimization 2/2



Figure 12: Non-minimal transducer after pushing.



Figure 13: Minimal transducer.

Miloš Janda (FIT BUT)

Speech processing

< □ ▶

2010 25/33

 $\mathcal{A} \mathcal{A} \mathcal{A}$

臣

=

臣

Composition

The composition of two transducers represent their relational composition.

The composition $T = R \circ S$ of two transducers R and S

has exactly one path mapping sequence *u* to sequence *w* for each pair of paths,

the first in R mapping u to v and second in S mapping v to w.

The weight of a path in T is the \otimes -product of the weights of corresponding paths in R and S.

SQ Q

→ @ ▶ ★ 말 ▶ ★ 말 ▶ ... 말

Example of composition





(c)

Figure 14: Example of transducer composition.

Miloš Janda (FIT BUT)

590

王

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Composition

Speech recognizer consist of this parts:

- Acoustic observation sequence of an utterance O
- Acouctic modeling block A
- Context-dependency modeling block C
- Dictionary D
- Grammar G

Thus, the domain of speech recognition can be presented by transducer composition:

$G \circ D \circ C \circ A \circ O$

SQ Q

< ∃ < < </l>



Figure 15: Stages of speech recognition.

 $\mathcal{O} \mathcal{Q} \mathcal{O}$

12

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ →

Output sentense, Lattice

Sometimes, we are not interested in the best path only, but we would like to obtain a lattice of hypotheses. This is achieved by storing more than just the best token in the end-of-word state.

A word-level graph - Lattice - is created this way.



Figure 16: Lattice sample.

Miloš Janda (FIT BUT)

Speech processing

 $\mathcal{A} \mathcal{A} \mathcal{A}$

王

<ロト < 団 > < 国 > < 国 > < 国 > 、

Example of lattice



Figure 17: Word lattice W_1 , ATIS task, for the utterance Which flights leave Detroit and arrive at Saint-Petersburg around nine am ?

For determinization only we can obtain reduction factor for states \approx 3 and for transitions \approx 9.

For determ. + minim. is reduction factor \approx 5 for states and \approx 17 for transitions.

Miloš Janda (FIT BUT)

Speech processing

2010 31/33

Thank you for attention

5900

王

·< ≣ ▶ < ≣ ▶

< □ ▶ < @ ▶

- Mehryar Mohri, Fernando Pereira, Michael Riley: Weighted finite-state transducers in speech recognition. Computer Speech Language, Vol. 16, No. 1. (2002), pp. 69-88.
- Mehryar Mohri: Finite-State Transducers in Language and Speech Processing.
 Computational Linguistics, Vol. 23, No. 2. (1997), pp. 269-311.
- Gold Ben, Morgan Nelson: Speech and Audio Signal Processing. John Wiley Sons, 2000

 $\mathcal{A} \mathcal{A} \mathcal{A}$

▲ 토 ▶ ▲ 토 ▶ - 토