# Weighted Finite State Transducers in Automatic Speec Recognition
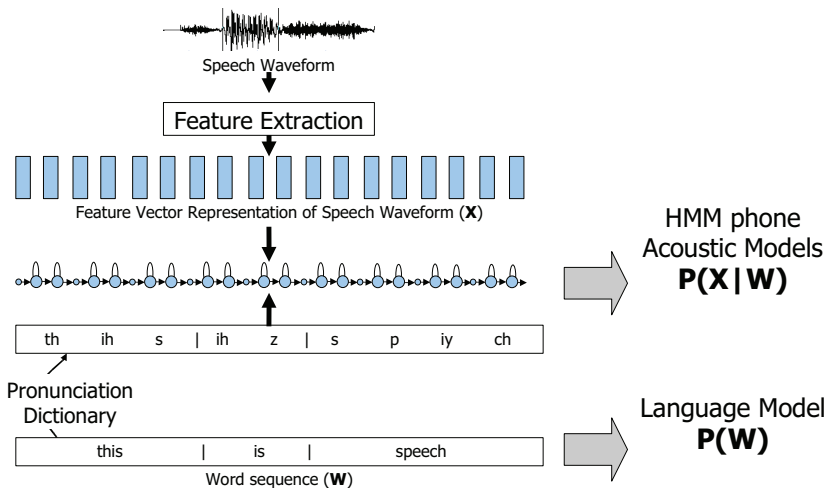
## ZRE lecture 24.04.2012

Mirko Hannemann

Slides provided with permission, Daniel Povey
some slides from T. Schultz, M. Mohri and M. Riley

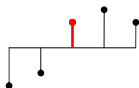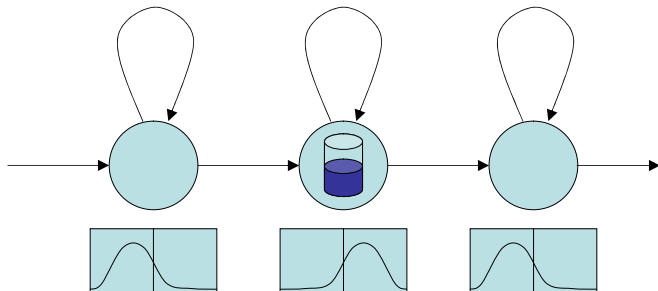24.04.2012

## Automatic speech recognition



Speech Waveform

Feature Extraction

Feature Vector Representation of Speech Waveform (**X**)

| th | ih | s | | ih | z | | s | p | iy | ch |

Pronunciation Dictionary

| this | | is | | speech |

Word sequence (**W**)

HMM phone Acoustic Models
**P(X|W)**

Language Model
**P(W)**

## Automatic speech recognition

- Determine the most probable word sequence $\tilde{W}$ given the observed acoustic signal Y

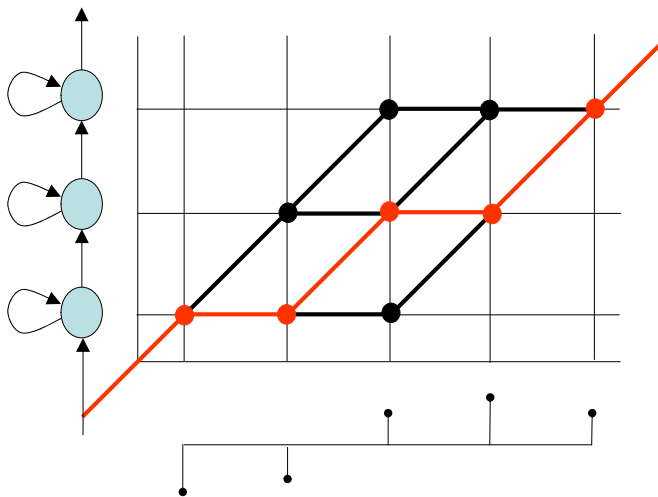$$\tilde{W} = argmax\, P(W|Y) = \frac{argmax\, P(W)\, P(Y|W)}{P(Y)}$$

- Search for word sequence $\tilde{W}$ that maximizes P(W) and P(Y|W)
  - P(W) = language model (likelihood of word sequence)
  - P(Y|W) = acoustic model
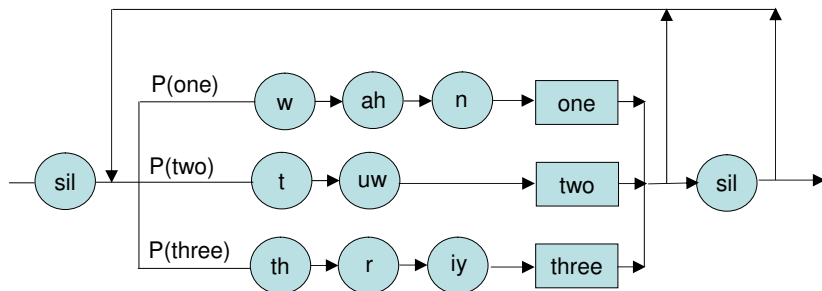    (likelihood of observed acoustic signal given word sequence)

# Hidden Markov Model, Token passing

# Viterbi algorithm, trellis
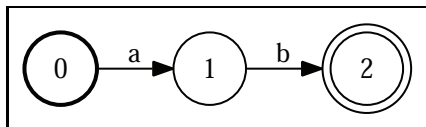
# Decoding graph/recognition network
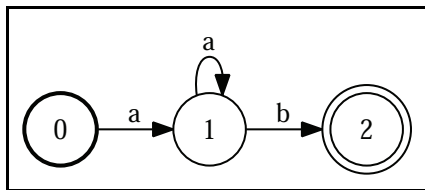
# Why Finite State Transducers?

Motivation:

- most components (LM, lexicon, lattice) are finite-state
- unified framework for describing models
- integrate different models into a single model via composition operations
- improve search efficiency via optimization algorithms
- flexibility to extend (add new models)
- $\rightarrow$ *speed*: pre-compiled search space, near realtime performance on embedded systems
- $\rightarrow$ *flexibility*: same decoder used for hand-held devices and LVCSR
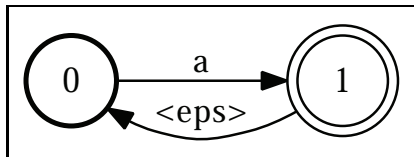
# Finite State Acceptor (FSA)



- ▶ An FSA "accepts" a set of strings
- ▶ (a string is a sequence of symbols).
- ▶ View FSA as a representation of a possibly infinite set of strings.
- ▶ This FSA accepts just the string *ab*, i.e. the set {*ab*}
- ▶ Numbers in circles are state labels (not really important).
- ▶ Labels are on arcs are the symbols.
- ▶ Start state(s) bold; final/accepting states have extra circle.
  - ▶ Note: it is sometimes assumed there is just one start state.
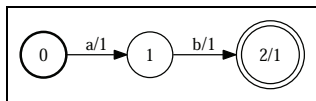
# A less trivial FSA



- ▶ The previous example doesn't show the power of FSAs because we could represent the set of strings finitely.
- ▶ This example represents the infinite set $\{ab, aab, aaab, \ldots\}$
- ▶ Note: a string is "accepted" (included in the set) if:
    - ▶ There is a path with that sequence of symbols on it.
    - ▶ That path is "successful' (starts at an initial state, ends at a final state).

## The epsilon symbol



- The symbol $\epsilon$ has a special meaning in FSAs (and FSTs)
- It means "no symbol is there".
- This example represents the set of strings $\{a, aa, aaa, \ldots\}$
- If $\epsilon$ were treated as a normal symbol, this would be $\{a, a\epsilon a, a\epsilon a\epsilon a, \ldots\}$
- In text form, $\epsilon$ is sometimes written as <eps>
- Toolkits implementing FSAs/FSTs generally assume <eps> is the symbol numbered zero

## Weighted finite state acceptors



- ▶ Like a normal FSA but with costs on the arcs and final-states
- ▶ Note: cost comes after "/". For final-state, "2/1" means final-cost 1 on state 2.
- ▶ View WFSA as a function from a string to a cost.
- ▶ In this view, unweighted FSA is f : string $\rightarrow \{0, \infty\}$.
- ▶ If multiple paths have the same string, take the one with the lowest cost.
- ▶ This example maps *ab* to $(3 = 1 + 1 + 1)$, all else to $\infty$.
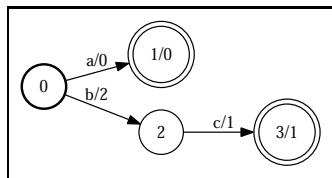
## Semirings

- ▶ The semiring concept makes WFSAs more general.
- ▶ A semiring is
    - ▶ A set of elements (e.g. $\mathbb{R}$)
    - ▶ Two special elements $\bar{1}$ and $\bar{0}$ (the identity element and zero)
    - ▶ Two operations, $\oplus$ (plus) and $\times$ (times) satsifying certain axioms.

*Semiring examples.* $\oplus_{\log}$ *is defined by:* $x \oplus_{\log} y = -\log(e^{-x} + e^{-y})$.

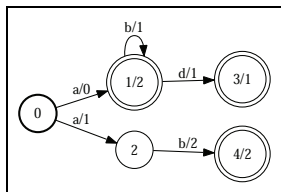| SEMIRING | SET | $\oplus$ | $\otimes$ | $\bar{0}$ | $\bar{1}$ |
|----------|-----|----------|-----------|-----------|-----------|
| Boolean | $\{0, 1\}$ | $\vee$ | $\wedge$ | 0 | 1 |
| Probability | $\mathbb{R}_+$ | $+$ | $\times$ | 0 | 1 |
| Log | $\mathbb{R} \cup \{-\infty, +\infty\}$ | $\oplus_{\log}$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

- In WFSAs, weights are multiplied along paths
- summed over paths with identical symbol-sequences
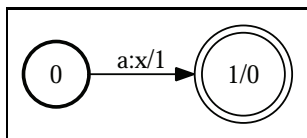
## Weights vs. costs



- ▶ Personally I use "cost" to refer to the numeric value, and "weight" when speaking abstractly, e.g.:
  - ▶ The acceptor above accepts *a* with unit weight.
  - ▶ It accepts *a* with zero cost.
  - ▶ It accepts *bc* with cost $4 = 2 + 1 + 1$
  - ▶ State 1 is final with unit weight.
  - ▶ The acceptor assigns zero weight to *xyz*.
  - ▶ It assigns infinite cost to *xyz*.

## Weights and costs



- ▶ Consider the WFSA above, with the tropical ("Viterbi-like") semiring. Take the string *ab*.
- ▶ We "multiply" ($\otimes$) the weights along paths; this means adding the costs.
- ▶ Two paths for *ab*:
    - ▶ One goes through states $(0, 1, 1)$; cost is $(0 + 1 + 2) = 3$
    - ▶ One goes through states $(0, 2, 3)$; cost is $(1 + 2 + 2) = 5$
- ▶ We add weights across different paths; tropical $\oplus$ is "take min cost" $\rightarrow$ this WFSA maps *ab* to 3

# Weighted finite state transducers (WFST)



- ▶ Like a WFSA except with two labels on each arc.
- ▶ View it as a function from a (pair of strings) to a weight
- ▶ This one maps $(a, x)$ to 1 and all else to $\infty$
- ▶ Note: view 1 and $\infty$ as costs. $\infty$ is $\bar{0}$ in semiring.
- ▶ Symbols on the left and right are termed "input" and "output" symbols.
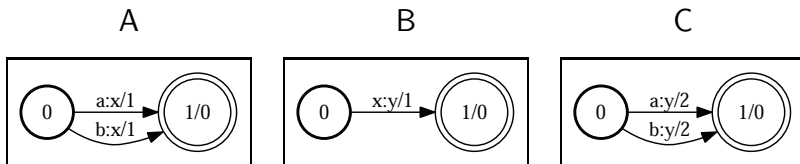
# Construction of decoding network

- WFST approach [Mohri et al.]
- exploit several knowledge sources (lexicon, grammar, phonetics) to find most likely spoken word sequence

$$HCLG = H \circ C \circ L \circ G \tag{1}$$

G probabilistic grammar or language model acceptor (word)

L lexicon (phones to words)

C context-dependent relabeling (ctx-dep-phone to phone)

H HMM structure (PDF labels to context-dependent phones)

Create H, C, L, G separately and compose them together

# Composition of WFSTs



A                          B                          C

- ▶ Notation: $C = A \circ B$ means, $C$ is $A$ composed with $B$.
- ▶ In special cases, composition is similar to function composition
- ▶ Composition algorithm "matches up" the "inner symbols"
  - ▶ i.e. those on the output (right) of $A$ and input (left) of $B$
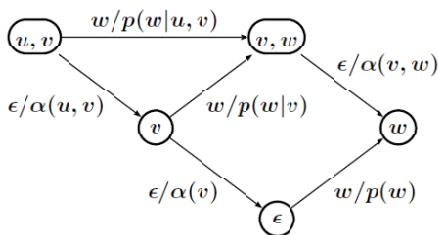
## Composition algorithm

- ▶ Ignoring $\epsilon$ symbols, algorithm is quite simple.
- ▶ States in $C$ correspond to tuples of (state in $A$, state in $B$).
  - ▶ But some of these may be inaccessible and pruned away.
- ▶ Maintain queue of pairs, initially the single pair $(0, 0)$ (start states).
- ▶ When processing a pair $(s, t)$:
  - ▶ Consider each pair of (arc $a$ from $s$), (arc $b$ from $t$).
  - ▶ If these have matching symbols (output of $a$, input of $b$):
    - ▶ Create transition to state in $C$ corresponding to (next-state of $a$, next-state of $b$)
    - ▶ If not seen before, add this pair to queue.
- ▶ With $\epsilon$ involved, need to be careful to avoid redundant paths...

## Language model acceptor G

- **G:** Grammar Transducer

Backing-off language model:

$$p(w|h) = \begin{cases} f(w|h) & : \text{if } N(w,h) > 0 \\ \alpha(h) \cdot f(w\overline{h}) & : \text{if } N(w,h) = 0 \end{cases}$$

- **Input:** word
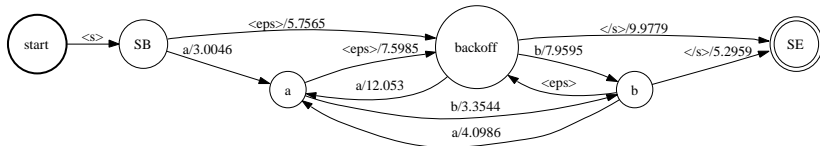- **Weight:** history dependent word probability

# Language models (ARPA back-off)

```
\1-grams:
-5.2347  a  -3.3
-3.4568  b
 0.0000  <s>  -2.5
-4.3333  </s>

\2-grams:
-1.4568  a b
-1.3049  <s> a
-1.78    b a
-2.30    b </s>
```
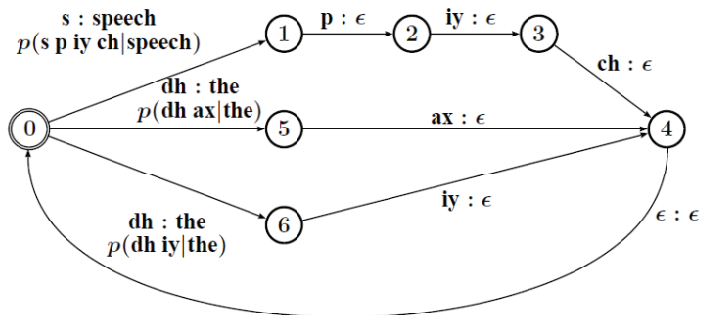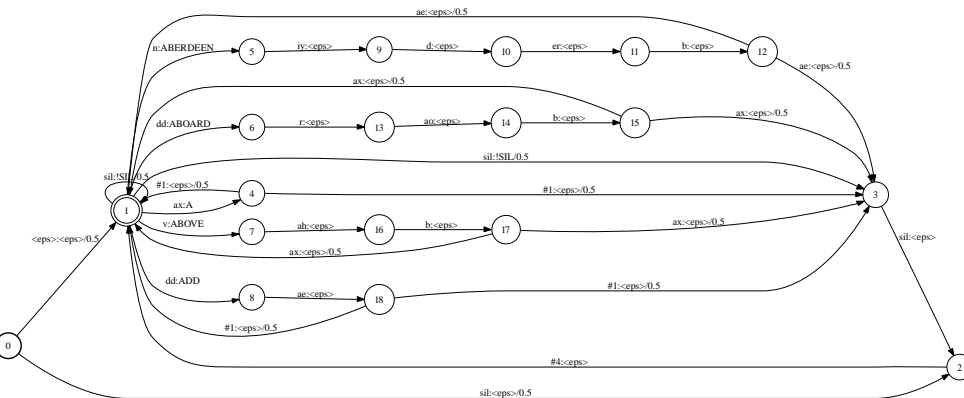
# Pronunciation lexicon L

- **L:** Context-Dependency Transducer
  - **Input:** context-independent phone (phoneme)
  - **Output:** word
  - **Weight:** pronunciation probability
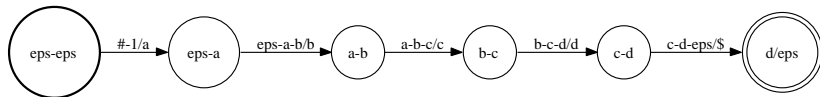
## Pronunciation lexicon L



```
A          ax #1
ABERDEEN   n iy d er b ae
ABOARD     dd r ao b ax
ADD        dd ae #1
```

# Context-dependency transducer C
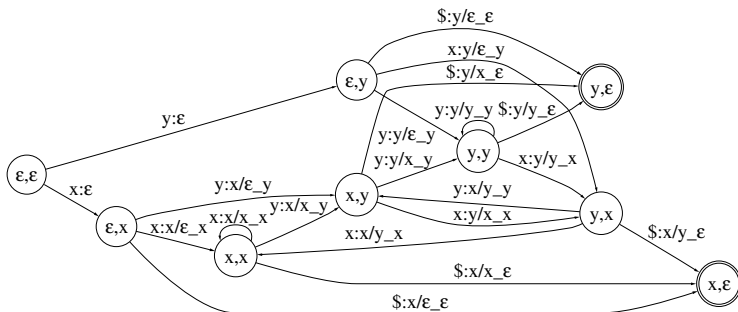
Input: context-dependent phone (triphone)

Output: context-independent phone (phone)
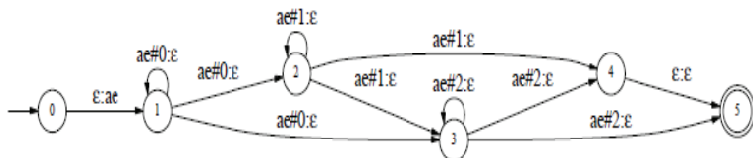
# Context-dependency transducer C



(a)

(b)

Figure 6: Context-dependent triphone transducer transition: (a) non-deterministic, (b) deterministic.
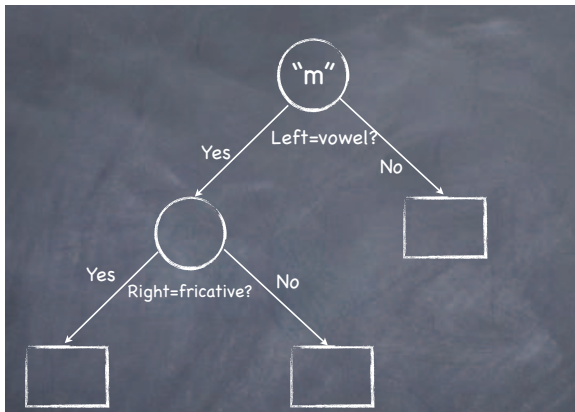
# HMM as transducer (monophone)

- H: HMM Topology Transducer (maps states to phonemes)
  - **Input:** state
  - **Output:** context-dependent phone (triphone)
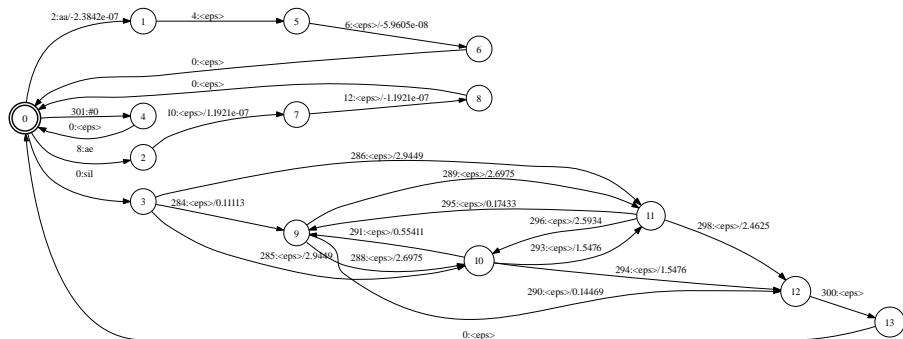  - **Weight:** HMM transition probability

## Phonetic decision tree

- too many context-dependent models → clustering
- determine model (PDF-id) based on phoneme context window and HMM-state
- sets of phones as questions

# HMM transducer $H_a$



(here shown for monophone case, without self-loops)

## Construction of decoding network

WFST approach by [Mohri et al.]

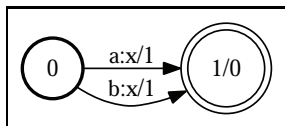$$HCLG = rds(min(det(H \circ det(C \circ det(L \circ G))))) \qquad (2)$$

rds - remove disambiguation symbols
min - minimization, includes weight pushing

det - determinization

Kaldi toolkit [Povey et al.]

$$HCLG = asl(min(rds(det(H_a \circ min(det(C \circ min(det(L \circ G)))))))) \qquad (3)$$
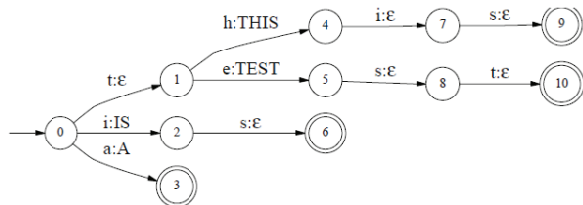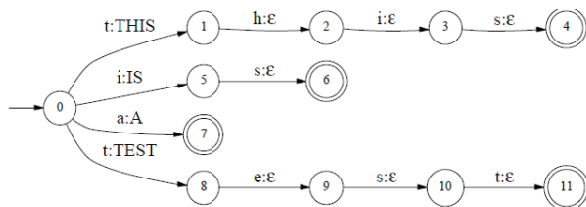
asl - add self loops

rds - remove disambiguation symbols

# Deterministic WFSTs



- ▶ Taken to mean "deterministic on the input symbol"
- ▶ I.e., no state can have $> 1$ arc out of it with the same input symbol.
- ▶ Some interpretations (e.g. Mohri/AT&T/OpenFst) allow $\epsilon$ input symbols (i.e. being $\epsilon$-free is a separate issue).
- ▶ I prefer a definition that disallows epsilons, except as necessary to encode a string of output symbols on an arc.
- ▶ Regardless of definition, not all WFSTs can be determinized.

# Determinization (like making tree-structured lexicon)

# Minimal deterministic WFSTs



- ▶ Here, the left FSA is not minimal but the right one is.
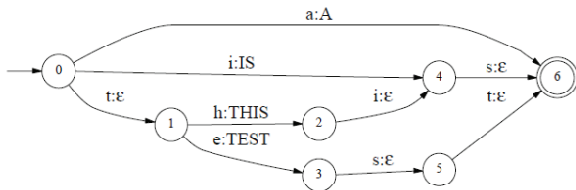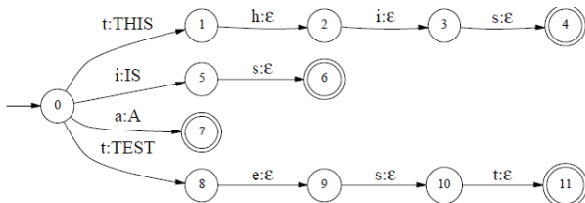- ▶ "Minimal" is normally only applied to *deterministic* FSAs.
- ▶ Think of it as suffix sharing, or combining redundant states.
- ▶ It's useful to save space (but not as crucial as determinization, for ASR).

# Minimization (like suffix sharing)

# Decoding graph construction (complexities)

- Have to do things in a careful order or algorithms "blow up"
- Determinization for WFSTs can fail
    - Need to insert "disambiguation symbols" into the lexicon.
    - Need to "propagate these through" H and C.
- Need to make sure final HCLG is stochastic, for optimal pruning performance.
    - i.e. sums to one, like a properly normalized HMM.
    - Standard algorithm to do this (weight-pushing) can fail
    - . . . because FST representation of backoff LMs is non-stochastic
- want to recover the phone sequence from path through HCLG (also PDF-index and transition arcs)

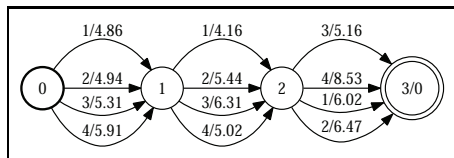## Decoding with WFSTs

- Solve

$$W' = argmax_W P(X|W)P(W)$$

- Compose recognizer as (H o C o L o G) which maps states to word sequences

- Decode by aligning the feature vectors X with HCLG ''
  i.e.,

$$W' = argmax_W X \circ (H \circ C \circ L \circ G)$$

## Decoding with WFSTs



- First– a "WFST definition" of the decoding problem.
- Let $U$ be an FST that encodes the acoustic scores of an utterance (as above).
- Let $S = U \circ HCLG$ be called the search graph for an utterance.
- Note: if $U$ has $N$ frames (3, above), then
  - #states in $S$ is $\leq (N + 1)$ times #states in $HCLG$.
  - Like $N + 1$ copies of $HCLG$.

## Decoding with WFSTs

- ▶ With beam pruning, we search a subgraph of $S$.
- ▶ The set of "active states" on all frames, with arcs linking them as appropriate, is a subgraph of $S$.
- ▶ Let this be called the *beam-pruned subgraph* of $S$; call it $B$.
- ▶ A standard speech recognition decoder finds the best path through $B$.
- ▶ In our case, the output of the decoder is a linear WFST that consists of this best path.
- ▶ This contains the following useful information:
    - ▶ The word sequence, as output symbols.
    - ▶ The state alignment, as input symbols.
    - ▶ The cost of the path, as the total weight.

## Lattices as WFSTs

The word "lattice" is used in the ASR literature as:

- Some kind of compact representation of the alternate word hypotheses for an utterance.
- Like an N-best list but with less redundancy.
- Usually has time information, sometimes state or word alignment information.
- Generally a directed acyclic graph with only one start node and only one end node.

# Lattice Generation with WFSTs [Povey12]

Basic process (not memory efficient):

- Generate beam-pruned subgraph B of search graph S
- The states in B correspond to the active states on particular frames.
- Prune B with beam $\alpha$ to get pruned version P.
- Convert P to acceptor and lattice-determinize to get A (deterministic acceptor)
- $\rightarrow$ No two paths in L have same word sequence (take best)
- Prune A with beam $\alpha$ to get final lattice L (in acceptor form).

## Finite State Transducers for ASR

Pro's:

Fast: compact/minimal search space due to combined minimization of lexicon, phonemes, HMM's

Simple: easy construction of recognizer by composition from states, HMMs, phonemes, lexicon, grammar

Flexible: whatever new knowledge sources, the compose/optimize/search remains the same

Con's:

- composition of complex models generates a huge WFST
- search space increases, and huge memory is required
- esp. how to deal with huge language models

## Ressources

OpenFST http://www.openfst.org/ Library, developed at Google Research (M. Riley, J. Schalkwyk, W. Skut) and NYU's Courant Institute (C. Allauzen, M. Mohri)

Mohri08 M. Mohri et al., *"Speech Recognition with weighted finite state transducers."*

Povey11 D. Povey et al., *"The Kaldi Speech Recognition Toolkit."*

Povey12 D. Povey et al., *"Generating exact lattices in the WFST framework."*