

Weighted Finite-State Transducers

Alexander Seward
Centre for Speech Technology
KTH
Stockholm, Sweden

Why Weighted Finite-State Transducers?

1. Efficiency and Generality of Classical Automata Algorithms

Efficient algorithms for a variety of problems (e.g. string-matching, compilers, parsing, pattern matching, process control, design of controllability systems in aircrafts).

General algorithms: rational operations, optimizations.

2. Weights

Handling uncertainty: text, handwritten text, speech, image, biological sequences.

Increased generality: finite-state transducers, multiplicity/indeterminism.

3. Applications

Text: pattern-matching, indexation, compression.

Speech: Large-vocabulary speech recognition, speech synthesis.

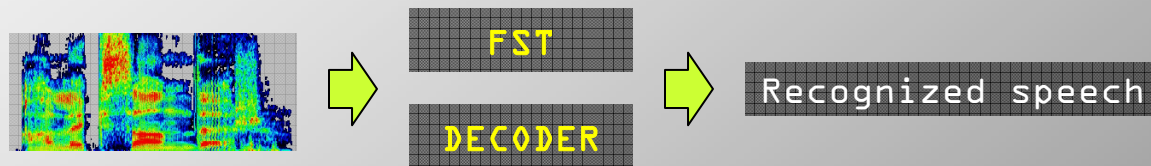
Image: image compression, filters.

* credits to M.Mohri

TRANSDUCERS

IN AUTOMATIC SPEECH RECOGNITION

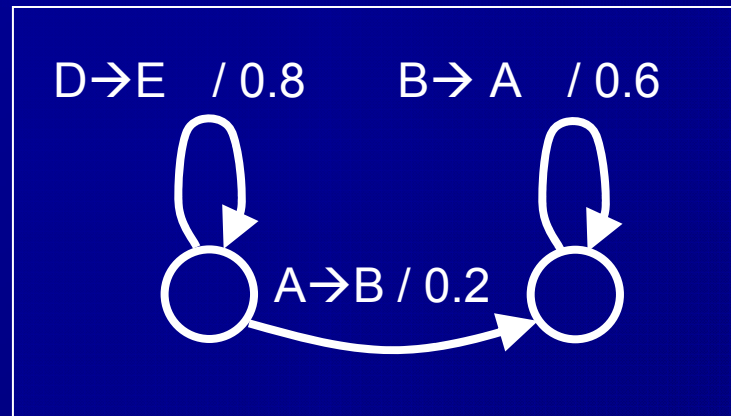
In ASR: Mathematical models for speech-to-text translation



A uniform composition of different information sources:
HMM-data, lexica, language models, etc..

Flexible: reduces decoder dependencies,
multiple layers,
generic optimization methods.

What is a Weighted Finite-State Transducer (WFST) ?



A finite-state machine where each arc is a weighted transduction consisting of an input, an output, and a probability/weight

Simply put: *A translation device*

A WFSA is a transducer without output

WFSTs in recognition

- I want a ticke..#noise" ..Boston from New York

FST trained on acoustics and language corpus:

#noise" must be "t...to" !

WFSTs in recognition

The **bare** was **bear** naked?

- FST trained on language corpus:
- The bear was bare naked

Recognition Cascade (simplified)

- I: Input feature vectors

Feature vectors

- H: HMM

Feature vectors

\xrightarrow{w}

CD-HMMs

- C: Context-Dependency Model

CD-HMMs

\xrightarrow{w}

Transcr. syms

- L: Lexicon

Transcr. syms

\xrightarrow{w}

Words

- G: Grammars

Words

\xrightarrow{w}

Words

Use Weighted FST Composition
to compose the parts into one

Weighted FST operations

Best-path

Difference

Weight pushing

Closure

Equivalence

Label pushing

Compaction

Hadamard product

Reversal

Composition

Inversion

Epsilon removal

Concatenation

Minimization

Topological sort

Connection

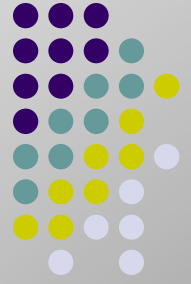
Projection

Union

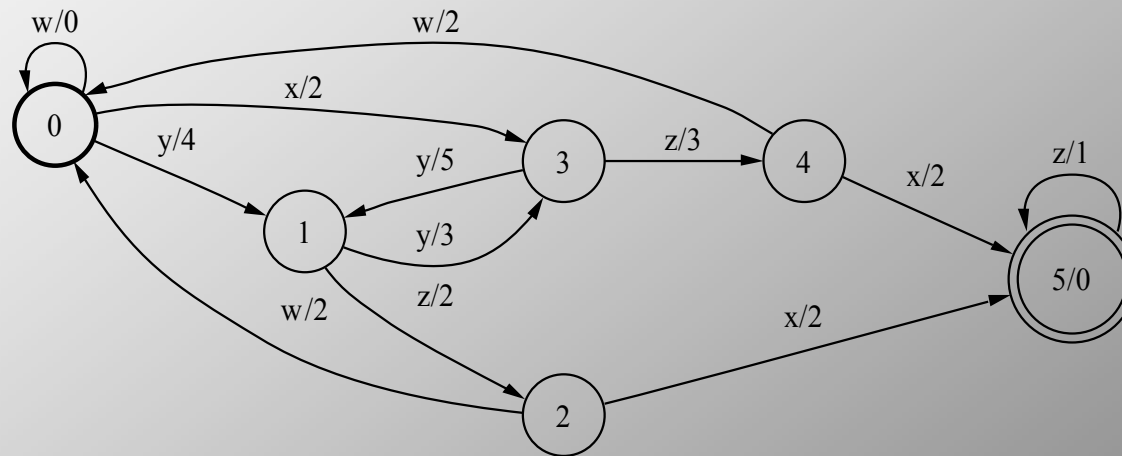
Determinization

Pruning

Language model WFSA

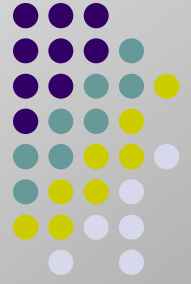


Model a priori weights for different word sequences (n-grams)

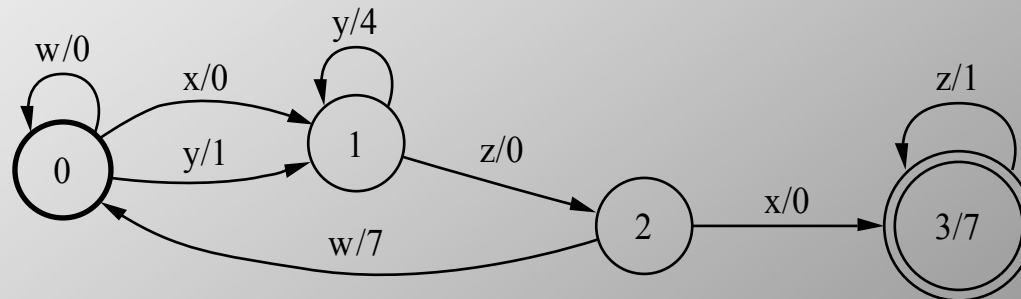


4 fictitious words: w , x , y , z

LM WFSA - Minimized



Model a priori weights to different word sequences (n-grams)



4 fictitious words: w, x, y, z

Pronunciation knowledge



because			about		
IPA	ARPAbet	%	IPA	ARPAbet	%
[bɪkʌz]	[b iy k ah z]	27%	[əbau]	[ax b aw]	32%
[bɪkʌz]	[b ix k ah z]	14%	[əbaut]	[ax b aw t]	16%
[kʌz]	[k ah z]	7%	[bau]	[b aw]	9%
[kəz]	[k ax z]	5%	[ʌbau]	[ix b aw]	8%
[bɪkəz]	[b ix k ax z]	4%	[ɪbaut]	[ix b aw t]	5%
[bɪkʌz]	[b ih k ah z]	3%	[ɪbæ]	[ix b ae]	4%
[bəkʌz]	[b ax k ah z]	3%	[əbær]	[ax b ae dx]	3%
[kʊz]	[k uh z]	2%	[baʊr]	[b aw dx]	3%
[ks]	[k s]	2%	[bæ]	[b ae]	3%
[kɪz]	[k ix z]	2%	[baut]	[b aw t]	3%
[kɪz]	[k ih z]	2%	[əbaʊr]	[ax b aw dx]	3%
[bɪkʌʒ]	[b iy k ah zh]	2%	[əbæ]	[ax b ae]	3%
[bɪkʌs]	[b iy k ah s]	2%	[bɑ]	[b aa]	3%
[bɪkʌ]	[b iy k ah]	2%	[bær]	[b ae dx]	3%
[bɪkʌz]	[b iy k aa z]	2%	[ɪbaʊr]	[ix b aw dx]	2%
[əz]	[ax z]	2%	[ɪbat]	[ix b aa t]	2%

Figure 5.7 The 16 most common pronunciations of *because* and *about* from the hand-transcribed Switchboard corpus of American English conversational telephone speech (Godfrey et al., 1992; Greenberg et al., 1996).

respe
the t
not c
like
merr
spea
diffe
worc
glisht
of E

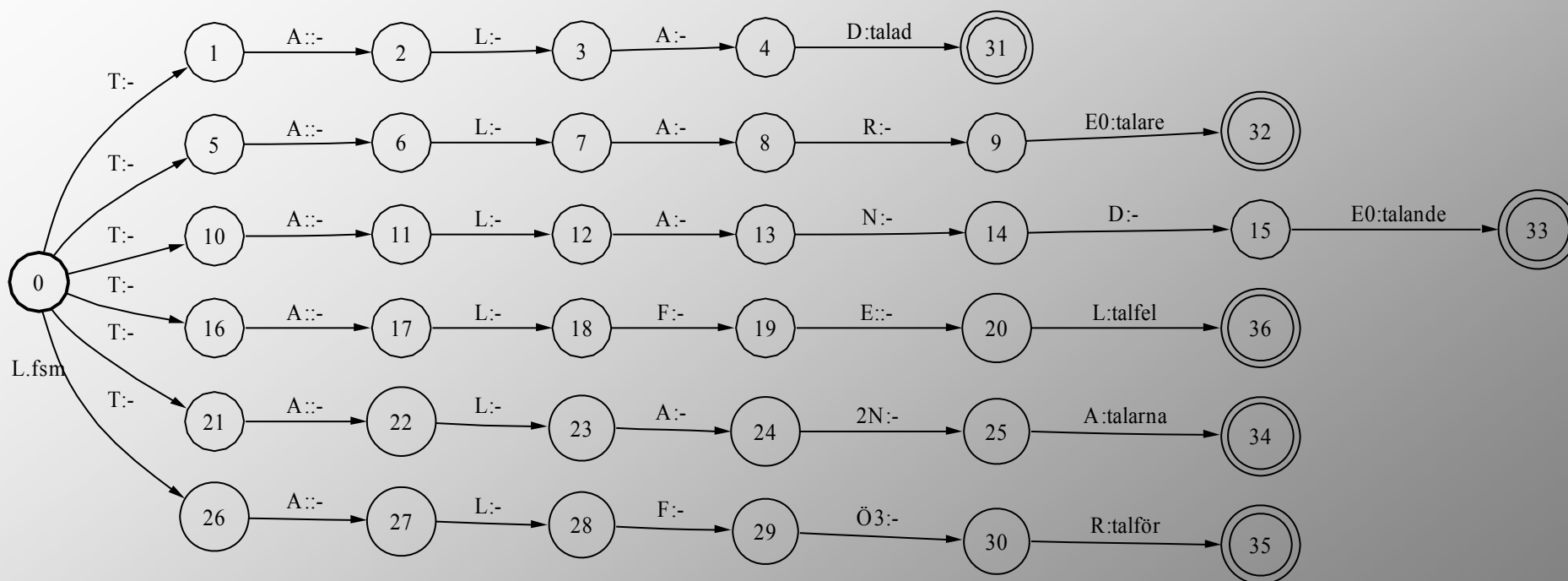
diale
mig
talki
choo
exar
pron
both
forr
[m]
gen
Shc
[m]

Use different weights to model likelihood of pronuciations!

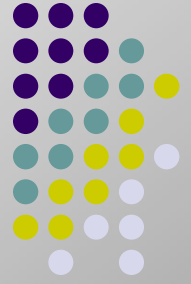


Lexicon transducer

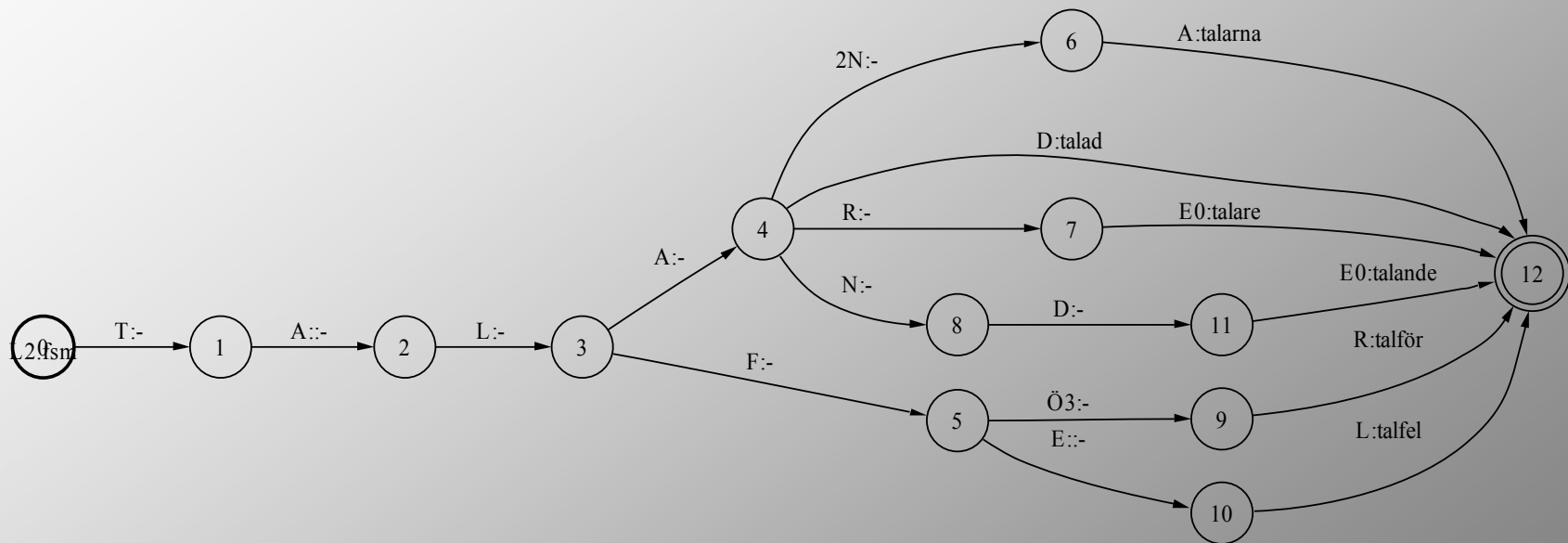
Some phonetically similar words



36 states

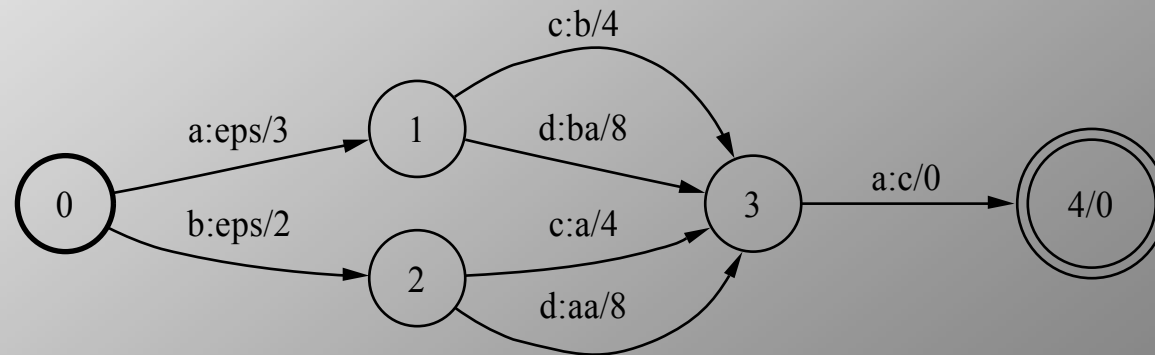
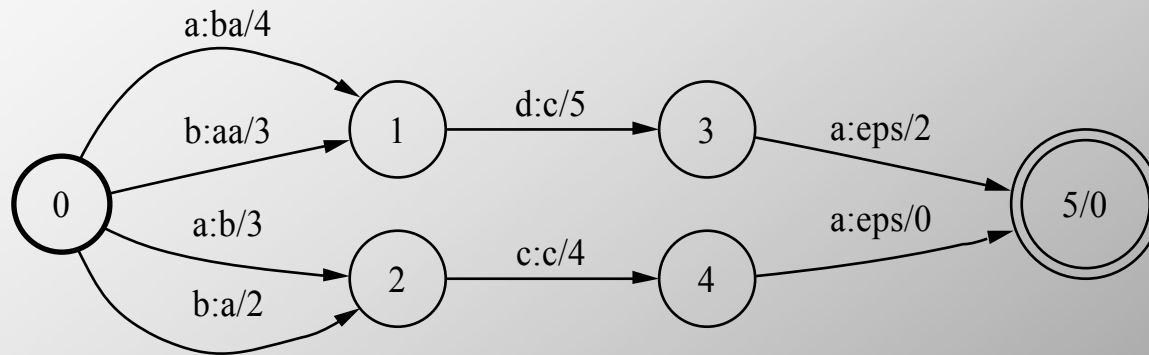
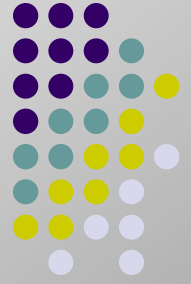


Equivalent lexicon transducer – Deterministic

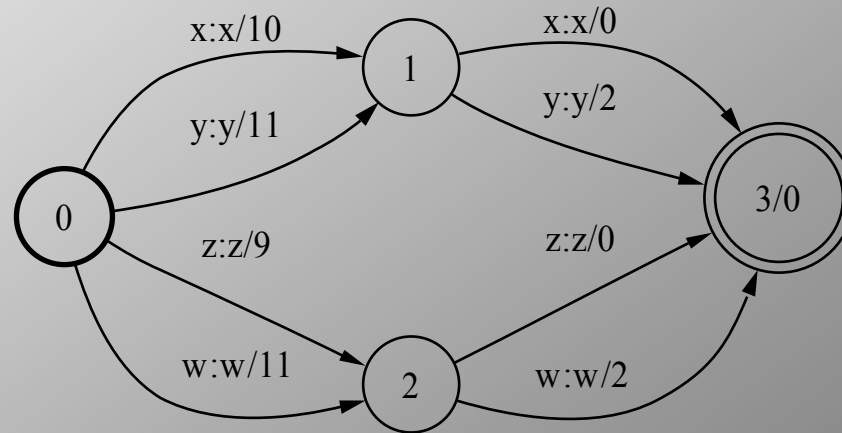
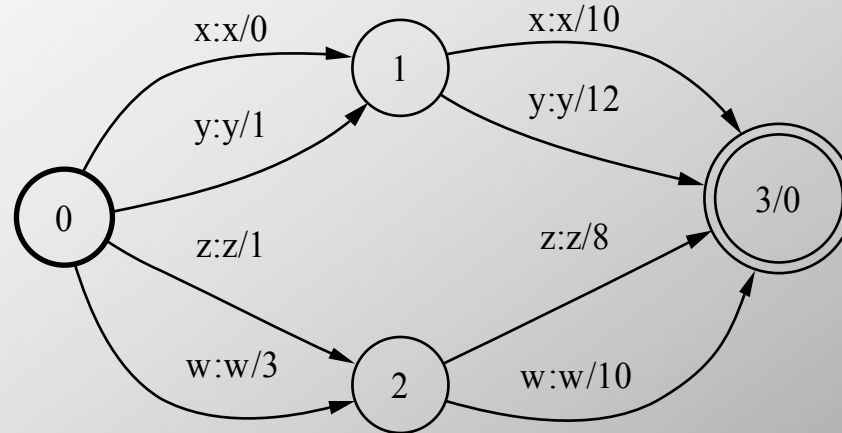
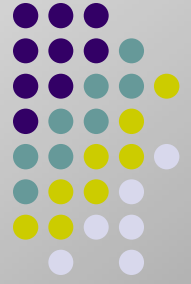


36 → 13 states

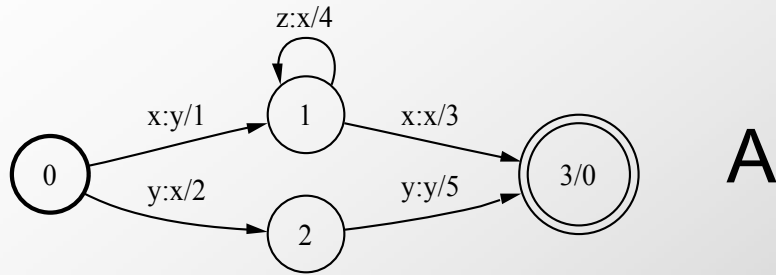
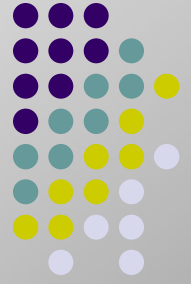
Weighted determinization



Weight pushing



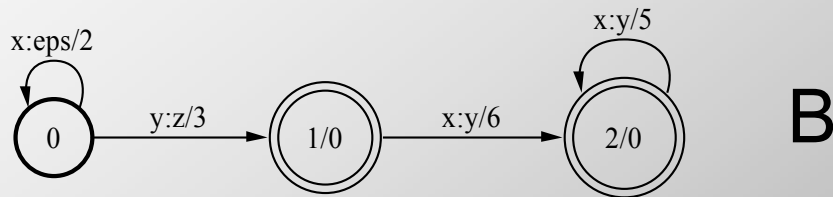
Weighted composition



A: $(x, z, z, x) \rightarrow (y, x, x, x)$
 weight: $1 + 4 + 4 + 3 = 12$.

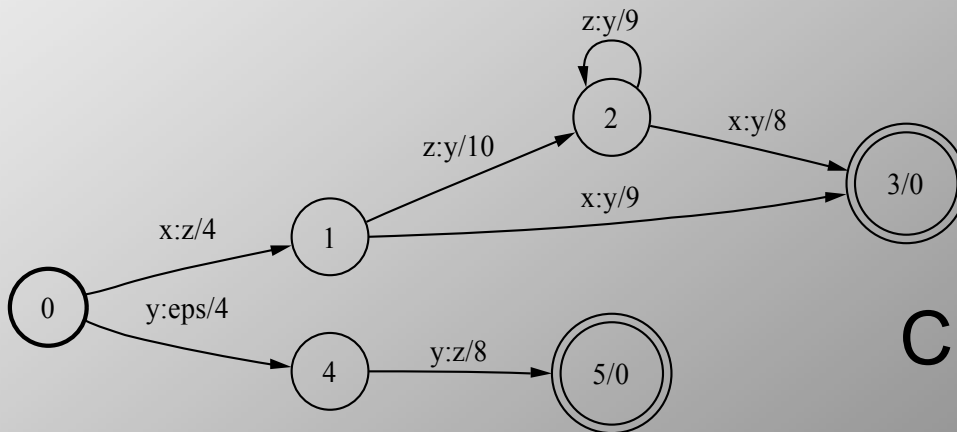
B: $(y, x, x, x) \rightarrow (z, y, y, y)$
 weight: $3 + 6 + 5 + 5 = 19$.

o



A o B: $(x, z, z, x) \rightarrow (z, y, y, y)$
 total weight: $12 + 19 = 31$.

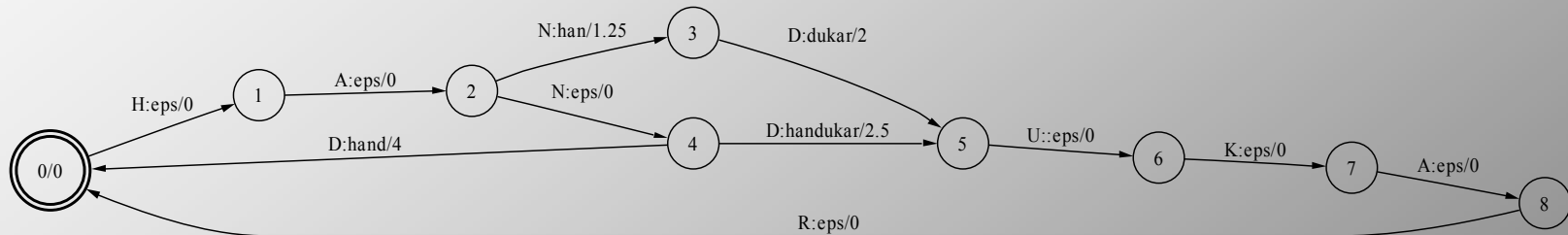
=



C = A o B: $(x, z, z, x) \rightarrow (z, y, y, y)$
 same total weight: $4 + 10 + 9 + 8 = 31$.



Lexicon + Grammar: L ◦ G

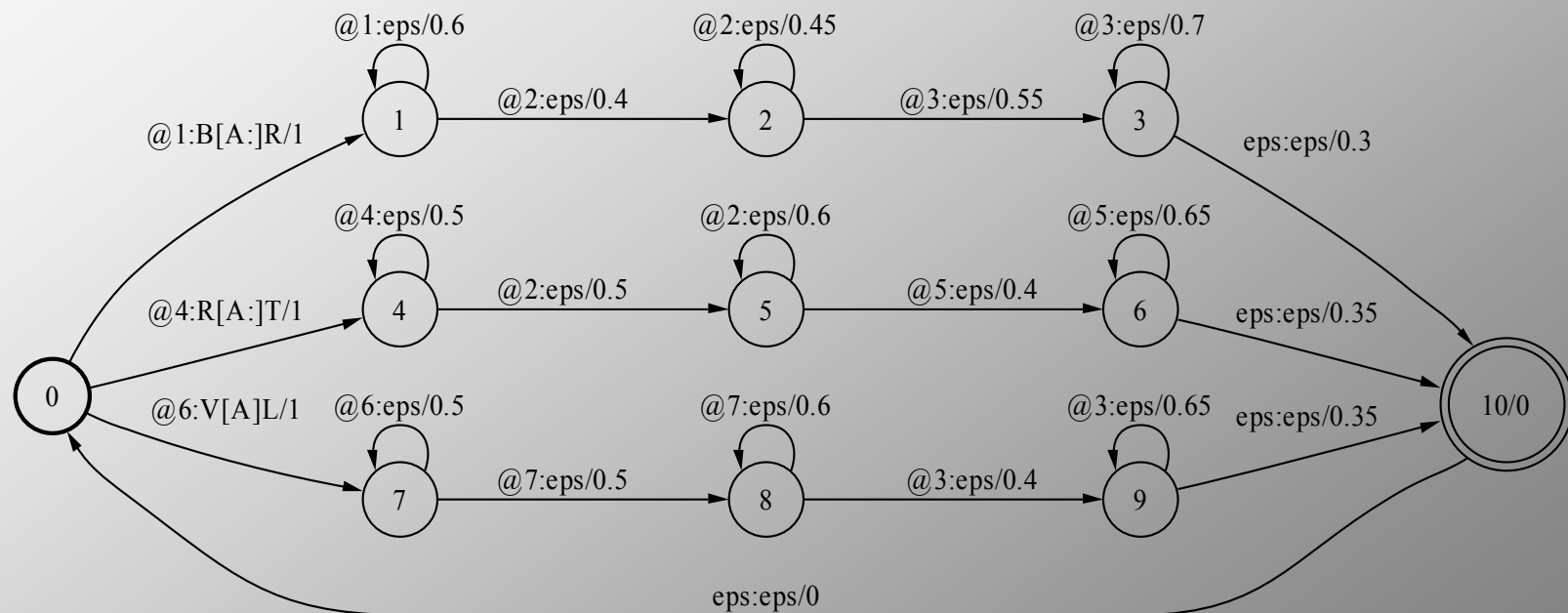


han	H A N	(Eng: he)
hand	H A N D	(Eng: hand)
handdukar	H A N D U: K A R	(Eng: towels)
dukar	D U: K A R	(Eng: set the table)

Sequence:
 han dukar H A N D U: K A R (Eng: he sets the table)



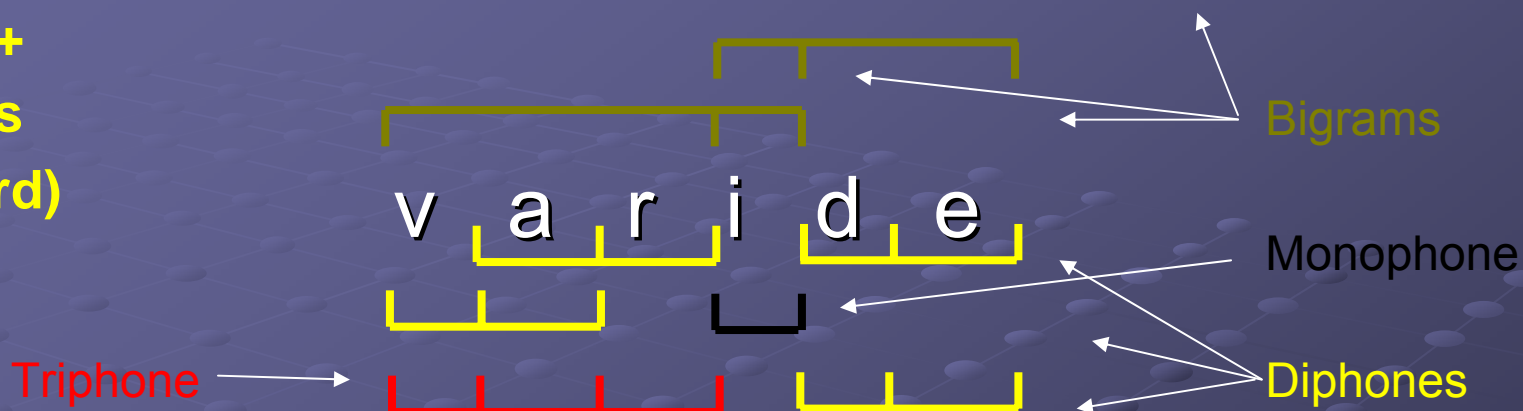
HMMs as WFSTs



Context-Dependency Modeling

“...var i det...”

**Bigrams +
Triphones
(no X-word)**



**Trigrams +
X-word Triphones**

r [i] d
a [r] i
x [v] a , i[d]e
v[a]r , d[e]x

