

Multiple Word Alignment with Profile Hidden Markov Models

Aditya Bhargava and Grzegorz Kondrak

Department of Computing Science, University of Alberta

{abhargava, kondrak}@cs.ualberta.ca

MULTIPLE WORD ALIGNMENT

kwatro	→	kwa-tro
kwattro		kwattro
katr		k-a-tr-
den		d--e--n-
deny	→	d--e--ny
dzen		dz-e--n-
dzien		dzie--n-
giorno		g--iorno
corteza		-c-o-rtez--a-
cortex		-c-o-rt---ex
cortica	→	-c-o-rtic--a-
corteccia		-c-o-rteccia-
scorza		sc-o-r--z--a-

USES & APPROACHES

Alignment of two words useful for:

- String similarity (Mackay and Kondrak, 2005)
- Dialect distances (Nerbonne and Heeringa, 1997)
- Cognate identification (Mackay and Kondrak, 2005)
- Comparative reconstruction (Covington, 1996)

Multiple alignment gets us:

- String similarity vs. multiple words
- Better-informed cognate identification
- Better-informed comparative reconstruction (Covington, 1998)
- Sentence-level paraphrasing (Barzilay and Lee, 2003)

How to do it?

- One way: hand-crafted scales of similarity phoneme classes (Covington, 1998)
- Iterative pairwise
- Copy the computational biologists! (Durbin et al., 1998)

MMIIM

AG...C

A-AG.C

AG.AA-

--AAAC

AG...C

EXPERIMENTS

Data:

- Comparative Indo-European Data Corpus (Dyen et al., 1992)
- cognation data for words in 95 languages corresponding to 200 languages
- English orthography

Multiple alignment:

- Initialize a model (e.g. sample parameters from Dirichlet distributions)
- Train model to words using Baum-Welch
- Align words to model using Viterbi

Cognate set matching:

- Build model from candidate sets
- Score word to sets using forward algorithm
- Choose set with highest score

Smoothing:

- Substitution matrix
- Added during Baum-Welch

RESULTS

```

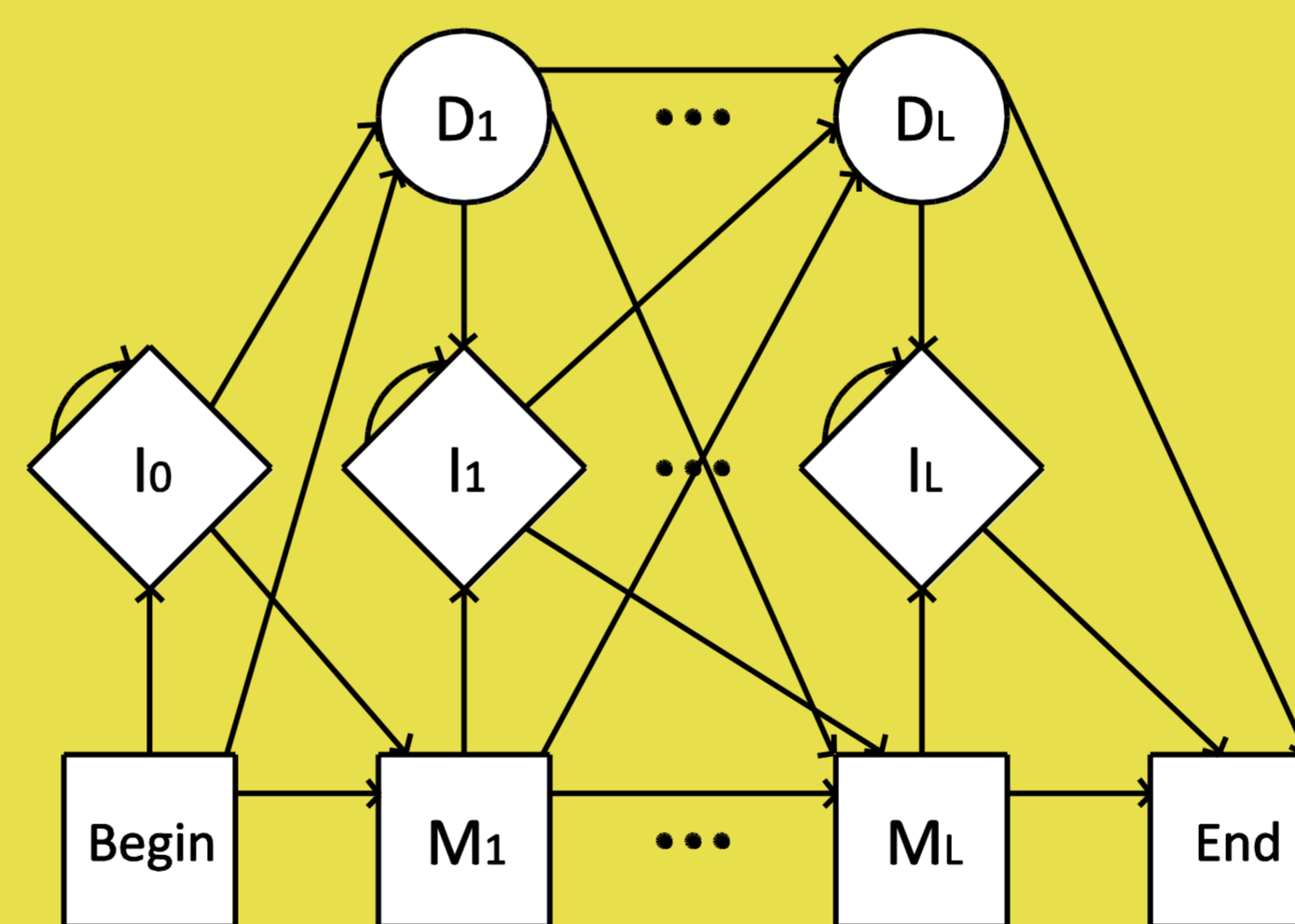
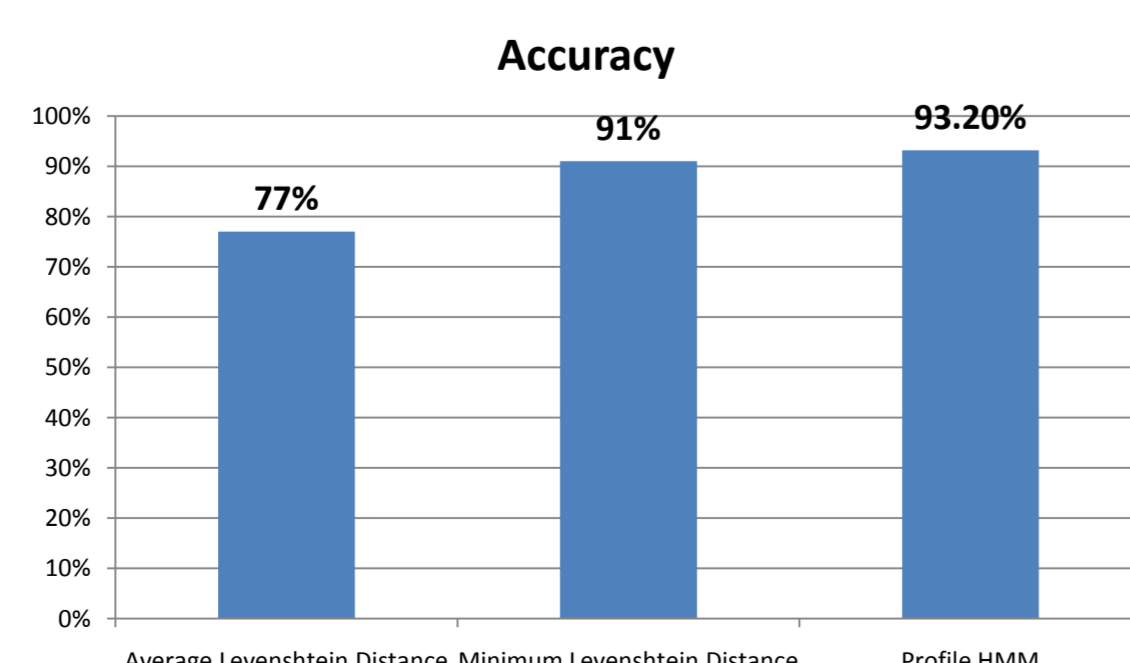
MMIIMMIIMMI
--M--R--N--U-
U--M--R-----I-T
--M--R-----I-
--M--R-----A-V
--M--R-----T-I-
--M--R-----T-I-
--M--R-----S-T
--M--R-----A-
--M--R-----C--I-
ZEM--R--I-T-I-
--M--R-----I-T-I-
--M--R-----E-S
--M--R-----E-C
U--M--R-----A-M
U--M--R-----A-C
U--M--R-----Z--E-C
U--M--R-----A-T
--M--R-----E-T
U--M--R-----E-T-I-
U--M--R-----E-T-I-
U--M--R-----A-T
--M--R-----E-T
U--M--R-----T-
--M--R-----E-T-I-
--M--R-----T-Y-
--M--R-----O-
--M--R-----O-
--M--R-----E-R-E-
--M--R-----I-R
--M--R-----I-
--M--R-----I-
--M--R-----I-R
--M--R-----I-R-E-
--M--R-----I-
--M--R-----E-R-E-
S--M--R-----I-R
--M--R-----I-
--M--R-----E-L
--M--R-----T-H-A-
--M--R-----I-
--M--R-----Y-N
--M--R-----E-N
--M--R-----N-U-
--M--R-----U-N
--M--R-----A-
--M--R-----V-E-L
--M--R-----E-L
--M--R-----W-
--M--R-----W-
--M--R-----W-U-
U--M--R-----E-
--M--R-----D-A-N
--M--R-----A-
--M--R-----A-
--M--R-----A-
--M--R-----E-
--M--R-----D-A-N
--M--R-----E-L
U--M--R-----E-
U--M--R-----A-T-Y-
PAM--R-----A-C
--M--R-----E-R
--M--R-----E-R
--M--R-----N-I-L
--M--R-----E-L
    
```

```

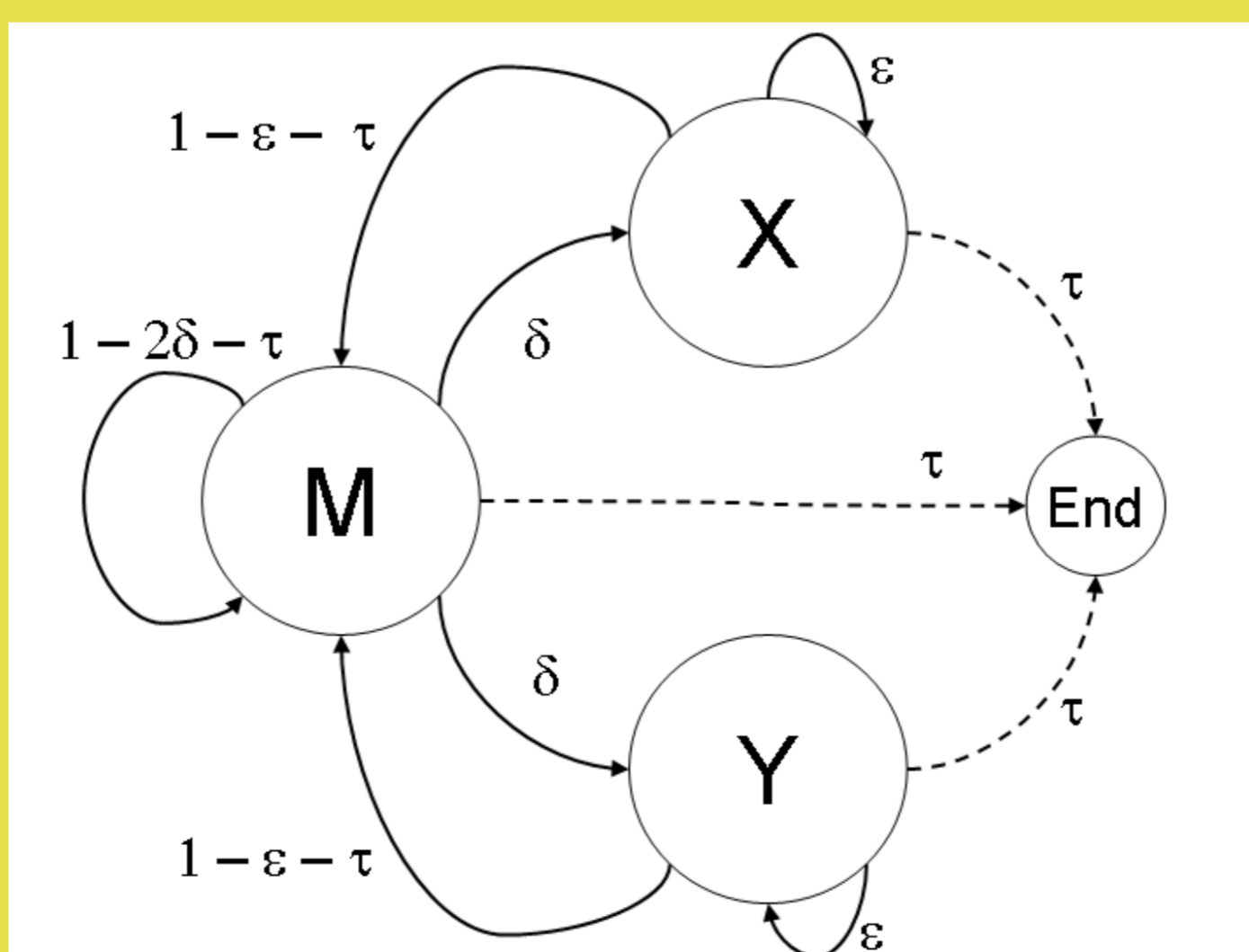
MMIIMMIIMMI
--SCHO-R-----S-
--C-O-RTEZ--A-
--C-O-RT---EX
--C-O-RTIC--A-
--C-O-RTECCIA-
--S-C-O-R--Z--A-
E-C-O-R--C--E-
--S-C-OART---A-
ISC-O-R--Z--A-
ESC-O-R--X--Z
--KRO---Z--U-
ISK-O-RT-H--A-
--SK-O-R-----A-
--K-A-R-----A-
--K-O-R-----A-
--K-U-R-----A-
    
```

```

MIIMMIIMMI
D--E--N-
D--E--NY
D--E--N-
D--E--N-
Z--E--N-
DZ--E--N-
DZIE--N-
D--E--N-
D--A--N-
D--A--N-
DI--E--NA
D--E--IZ
D--E--
D--Y--DD
D--I--A-
D--I--E-
D--I--I-
D--I--I-
Z--U--E-
Z--U--U-
J--O--UR
DJ--O--U-
J--O--UR
G--IORNO
    
```



A prototypical Profile HMM of length L .



...contrast to Pair HMMs, which have been used for word similarity & cognate identification (Mackay and Kondrak, 2005).

CONCLUSIONS

- Profile HMMs can work for word-related tasks
- Multiple alignments are reasonable
- Cognate set matching performance exceeds that of average and minimum Levenshtein distance
- If multiple words need to be considered, Profile HMMs present a viable method

FUTURE WORK

- Model construction from aligned sequences: e.g. maximum a posteriori model construction
- Initial models for unaligned sequences: more informed, decrease guesswork
- Smoothing methods
- N-gram output symbols

REFERENCES

- Regina Barzilay and Lillian Lee. 2003. Learning to paraphrase: An unsupervised approach using multiple sequence alignment. In *Proc. of NAACL-HLT*, pages 16–23.
- Michael A. Covington. 1996. An algorithm to align words for historical comparison. *Computational Linguistics*, 22(4):481–496.
- Michael A. Covington. 1998. Alignment of multiple languages for historical comparison. In *Proc. of COLING-ACL*, pages 275–279.
- Richard Durbin, Sean R. Eddy, Anders Krogh, and Graeme Mitchison. 1998. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge University Press.
- Isidore Dyen, Joseph B. Kruskal, and Paul Black. 1992. An Indo-European classification: A lexicostatistical experiment. *Transactions of the American Philosophical Society*, 82(5).
- Wesley Mackay and Grzegorz Kondrak. 2005. Computing word similarity and identifying cognates with pair hidden Markov models. In *Proc. of CoNLL*, pages 40–47.
- John Nerbonne and Wilbert Heeringa. 1997. Measuring dialect distance phonetically. In *Proc. of the Third Meeting of ACL SIGPHON*.