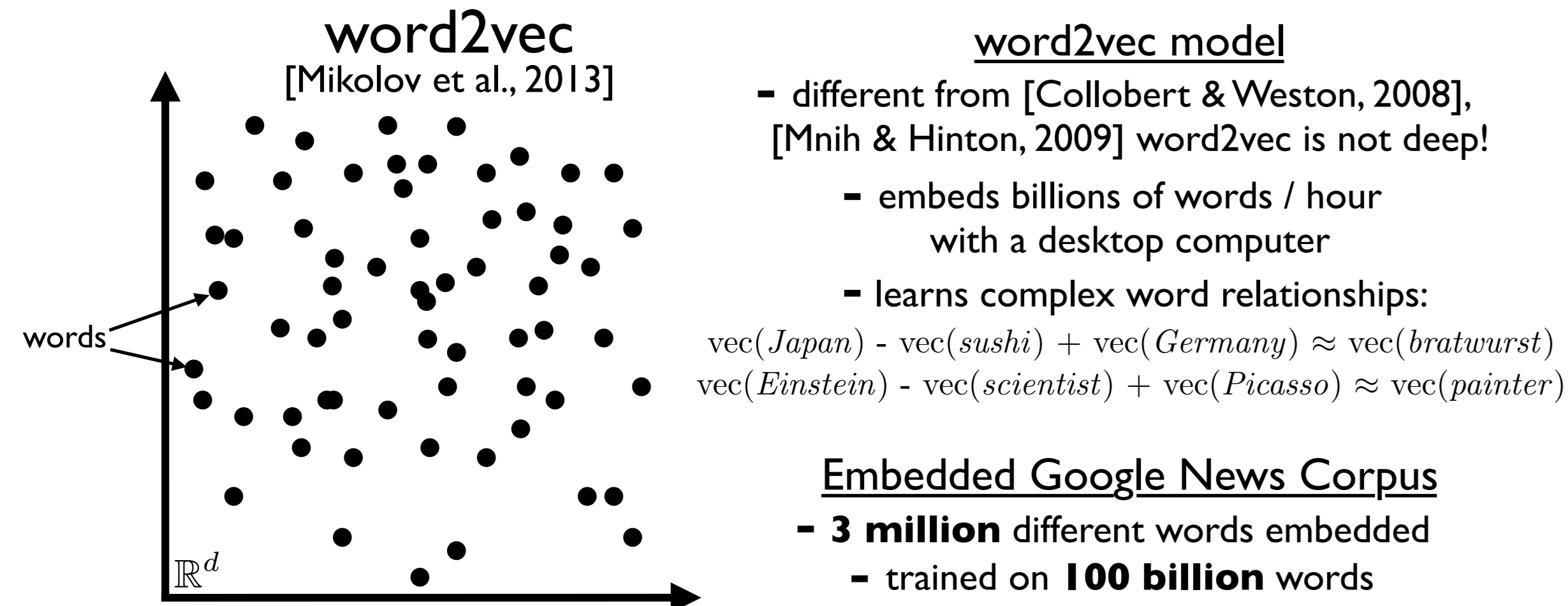
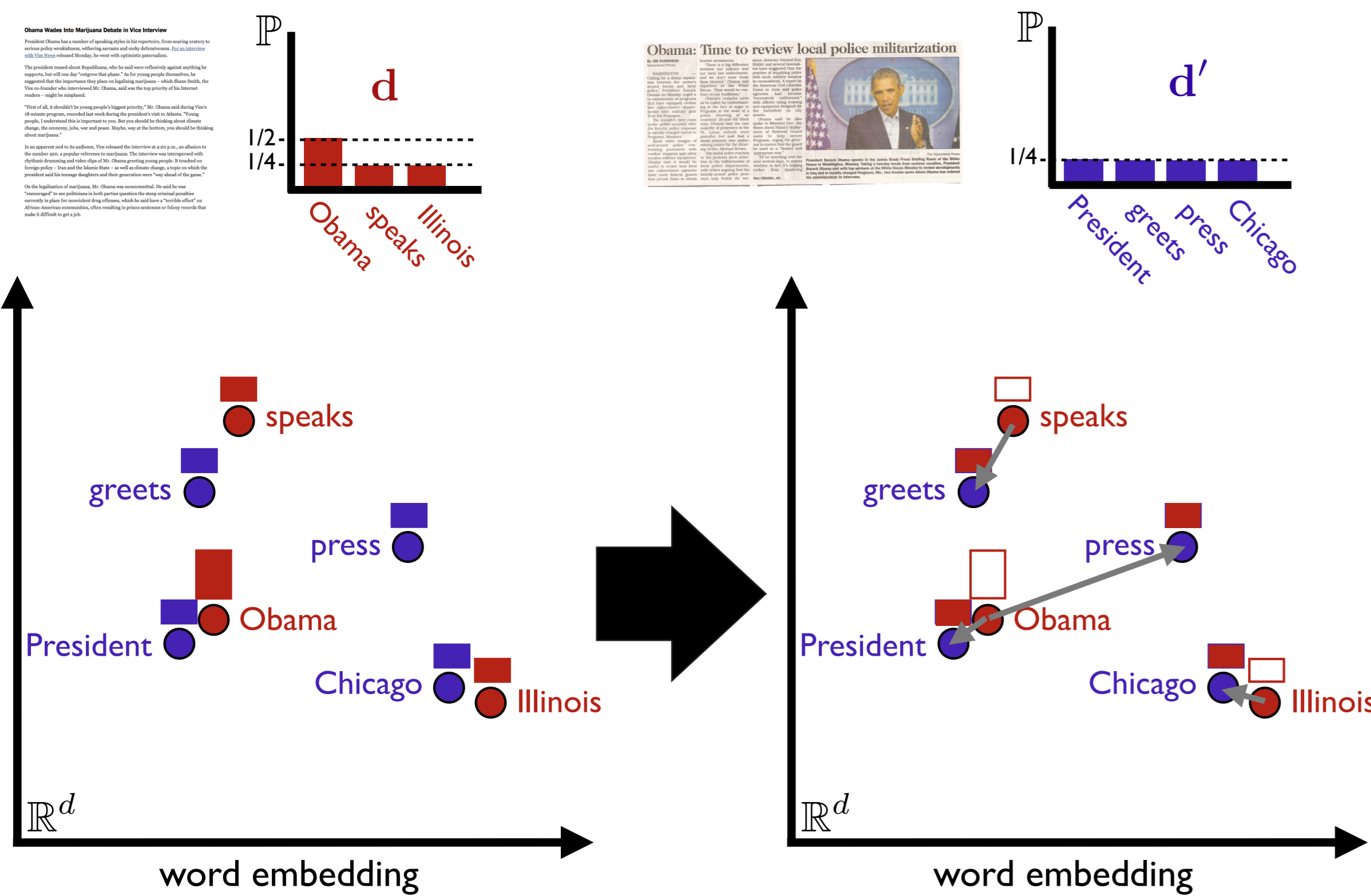


## Word Embedding & Document Vectors



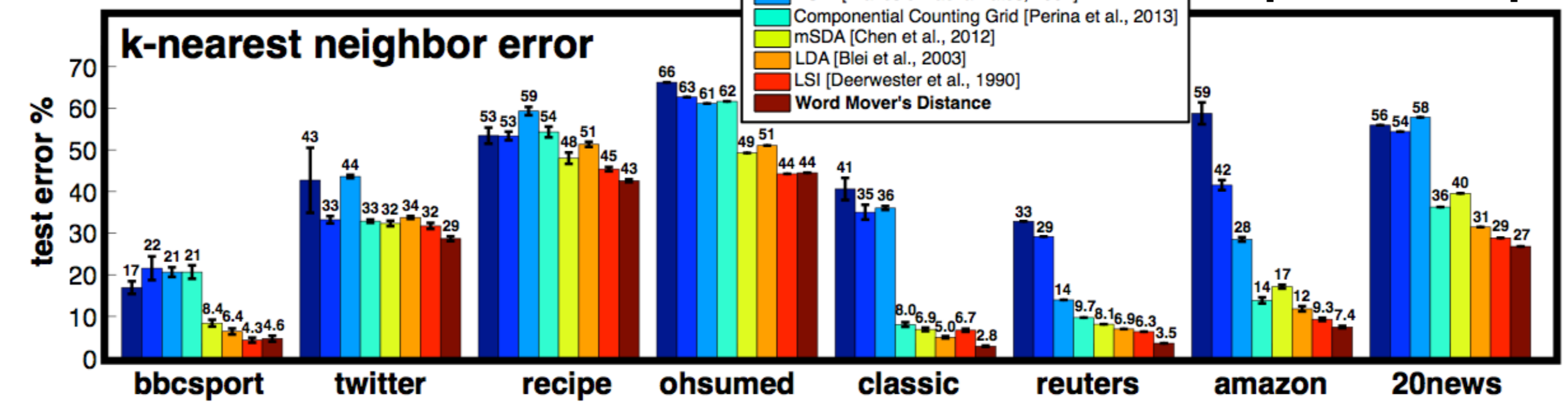
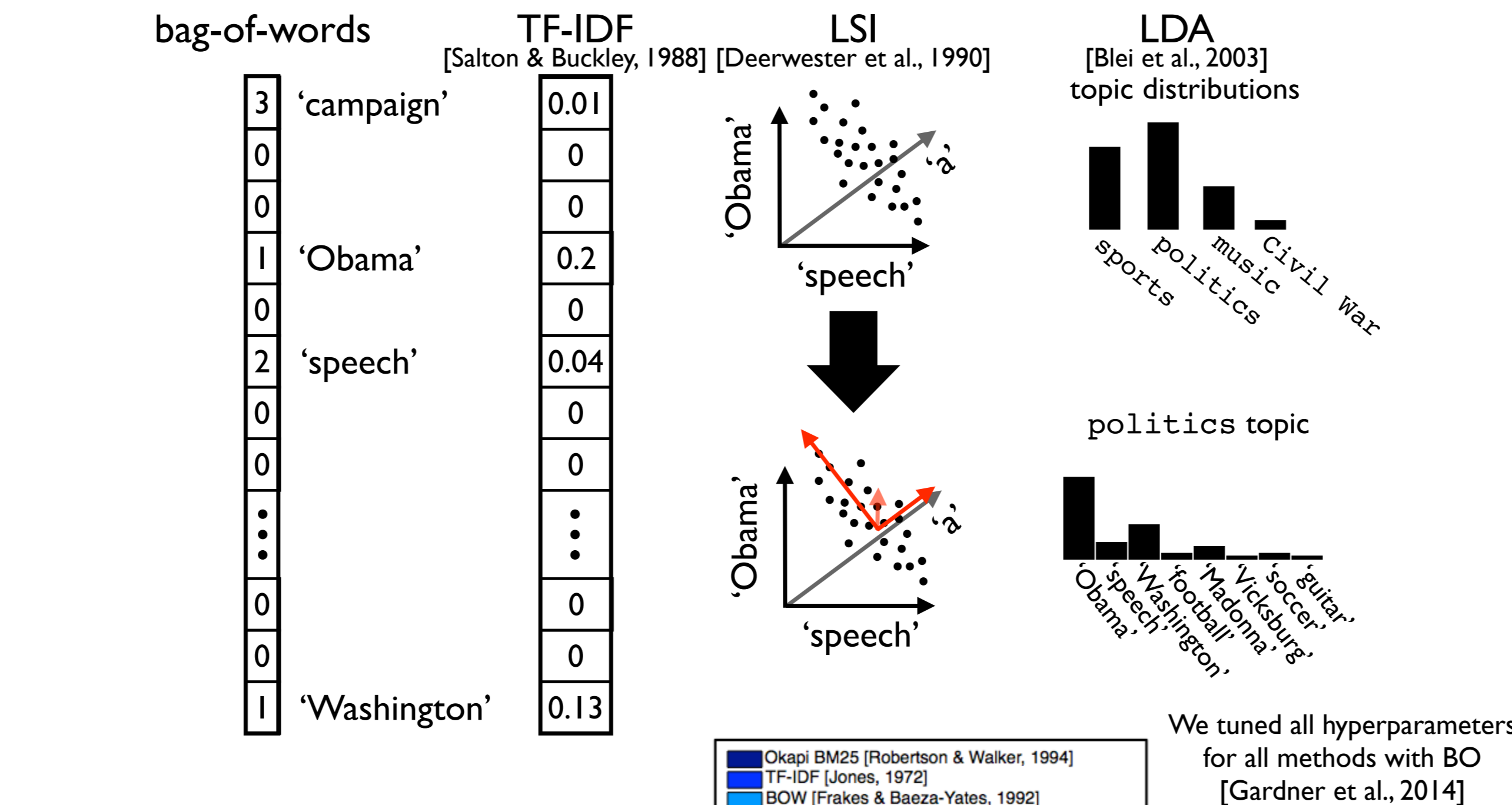
How can we leverage this high quality word embedding to compute **document distances**?

## The Word Mover's Distance

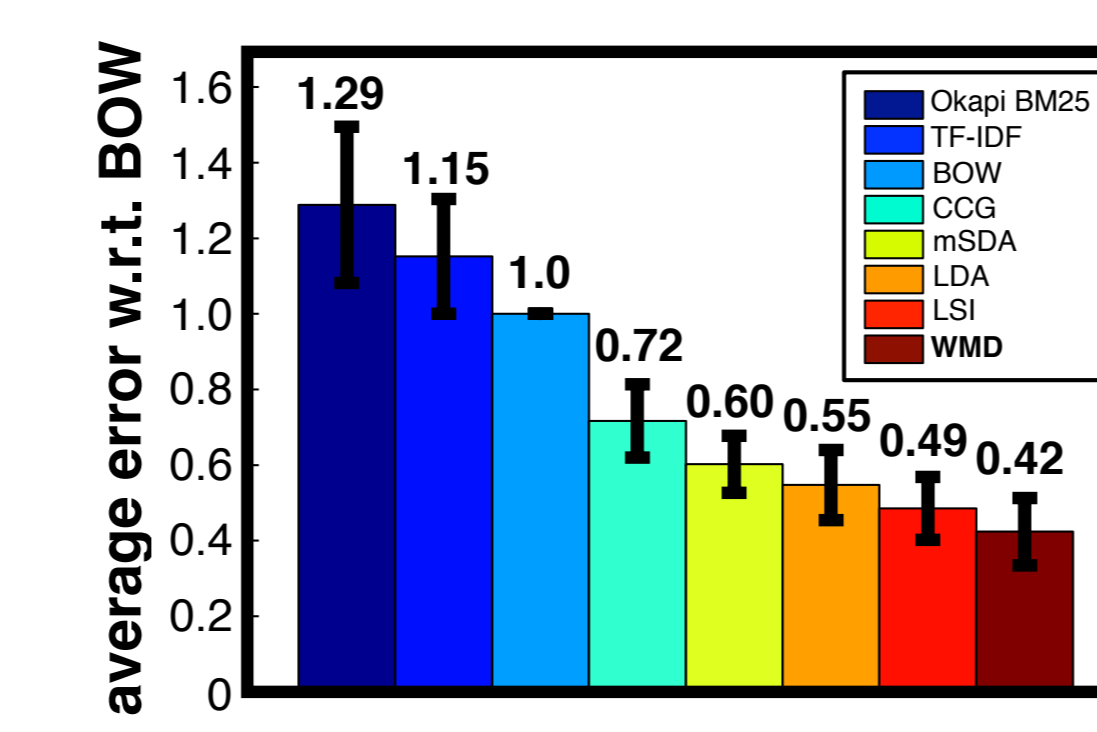


## k-Nearest Neighbor Results

### Prior Art



### average kNN error



### different word embeddings

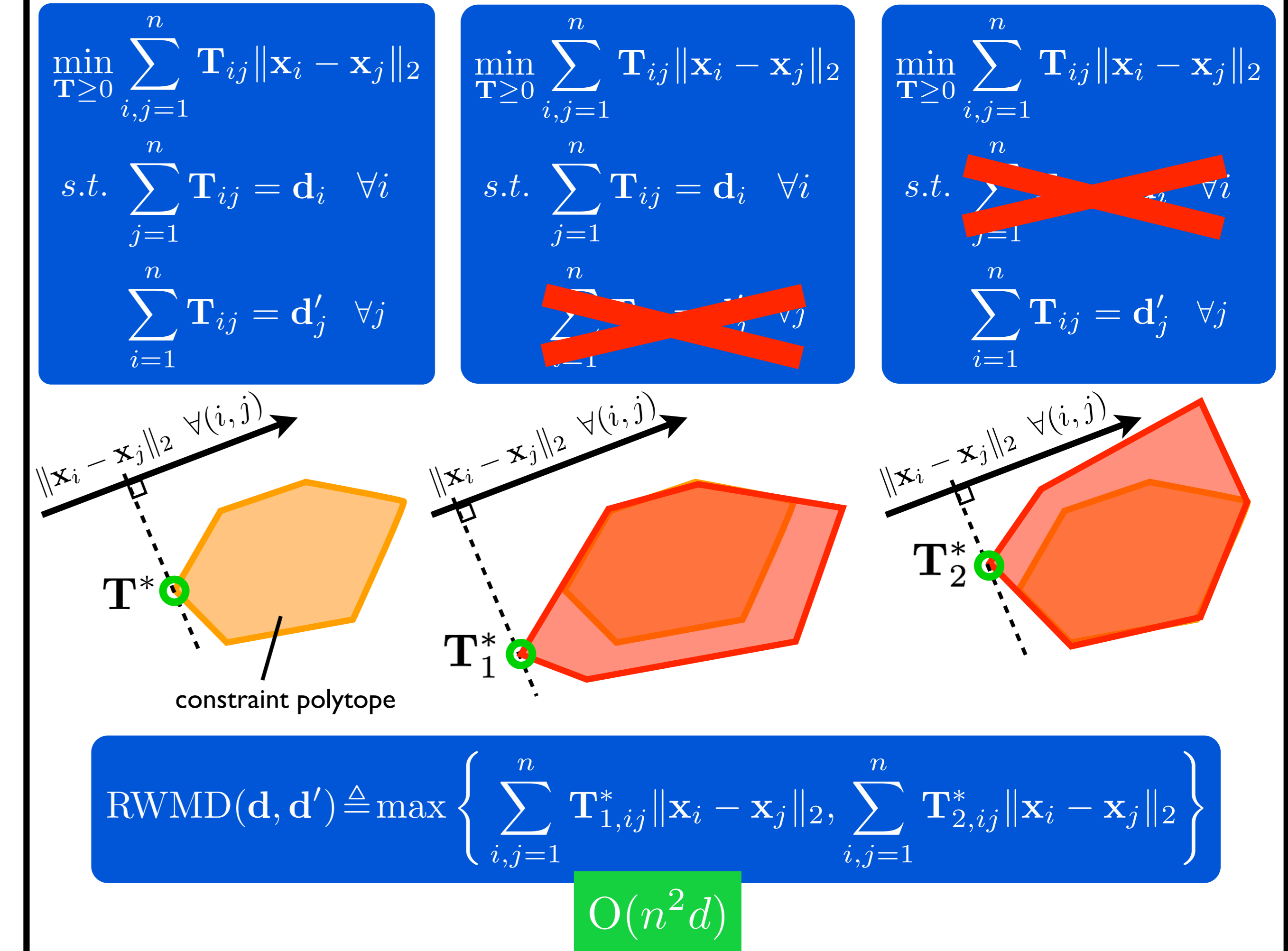
**DOCUMENT k-NEAREST NEIGHBOR RESULTS**

DATASET	HLBL (w2v)	CW (w2v)	NIPS (w2v)	AMZ (w2v)	NEWS (w2v)
BBCSPORT	4.5	8.2	9.5	4.1	5.0
TWITTER	33.3	33.7	29.3	28.1	28.3
RECIPE	47.0	51.6	52.7	47.4	45.1
OHSUMED	52.0	56.2	55.6	50.4	44.5
CLASSIC	5.3	5.5	4.0	3.8	3.0
REUTERS	4.2	4.6	7.1	9.1	3.5
AMAZON	12.3	13.3	13.9	7.8	7.2

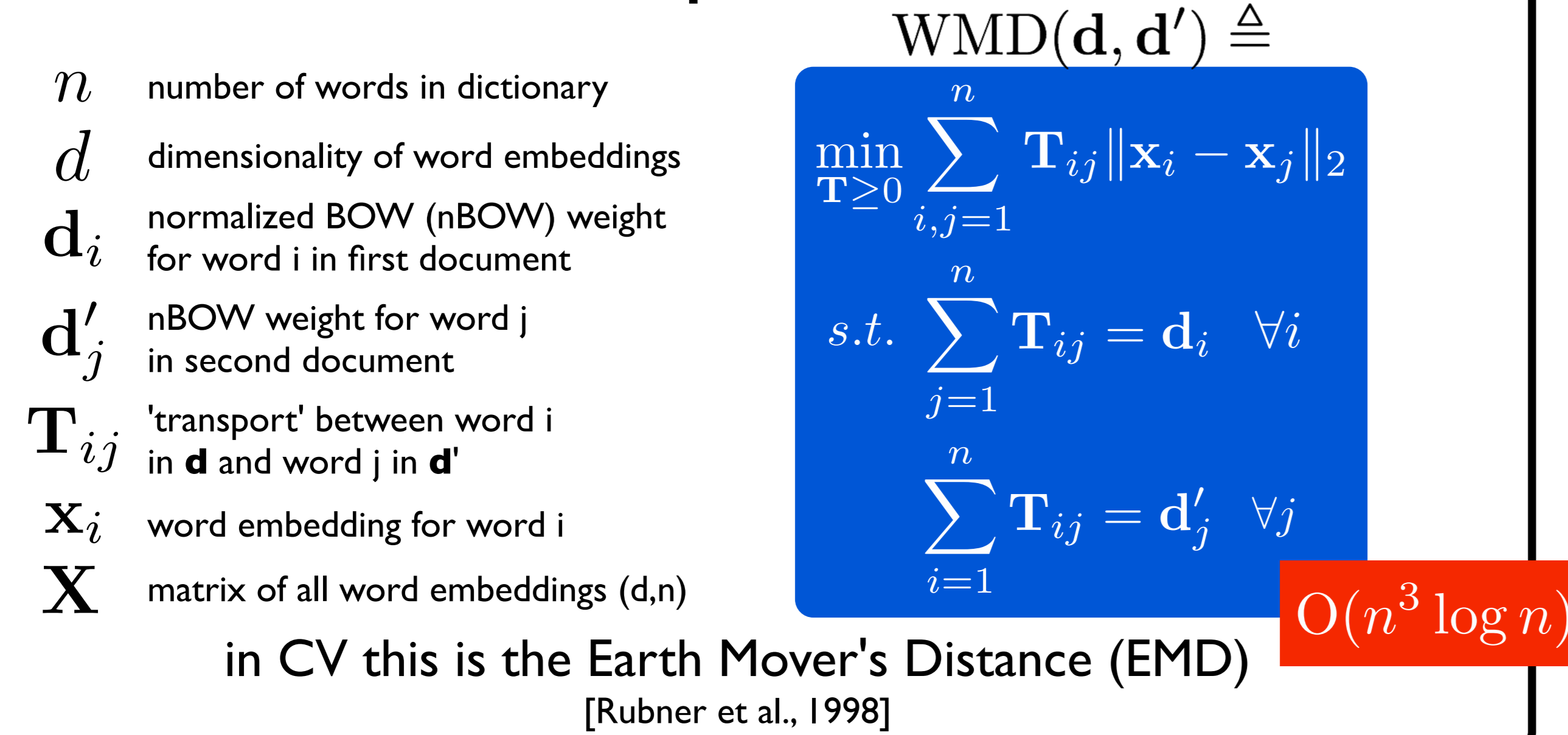
HLBL: [Mnih & Hinton, 2009]  
 CW: [Collobert & Weston, 2008]  
 W2V: [Mikolov et al., 2013]

## Approximations

### 2. The Relaxed Word Mover's Distance

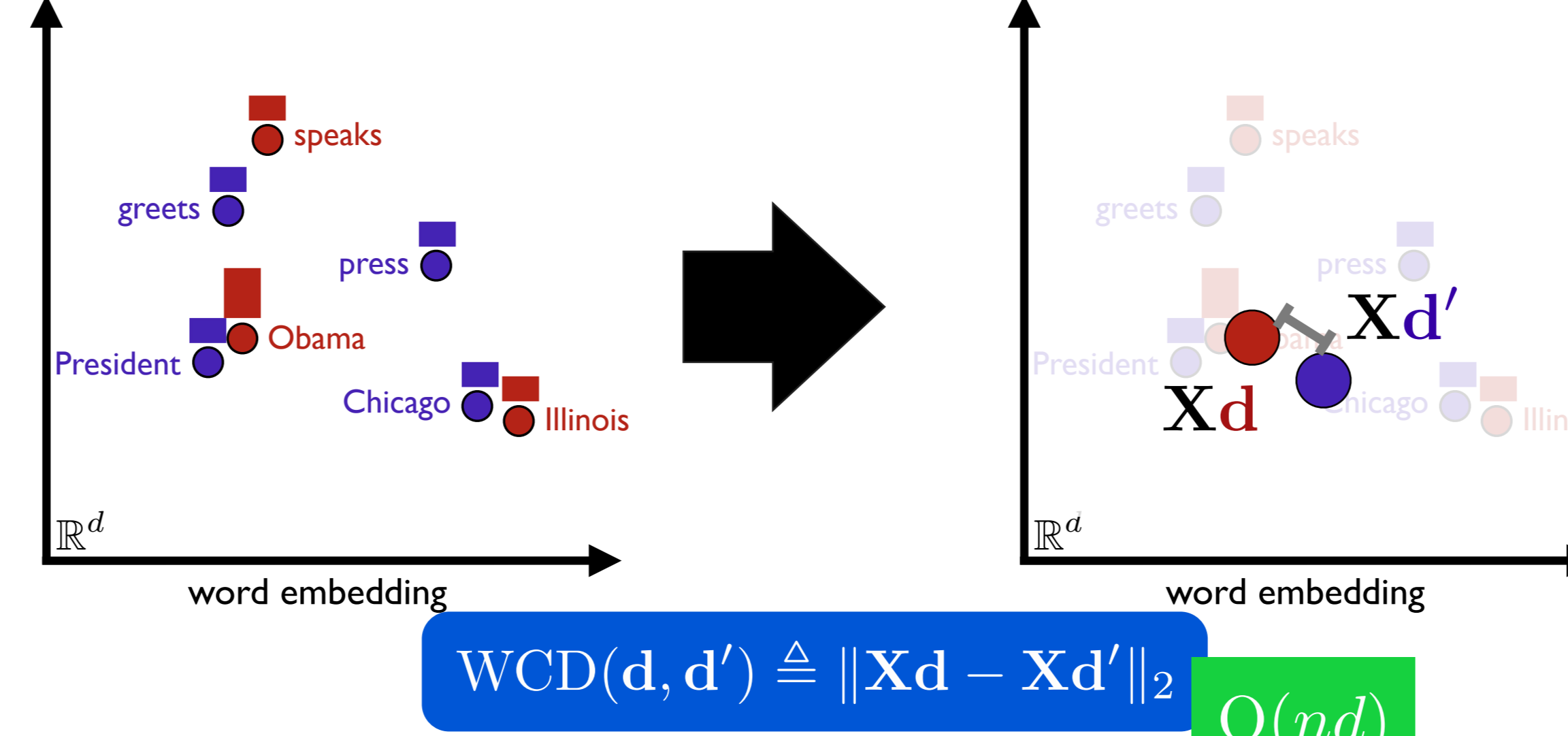


## WMD Optimization

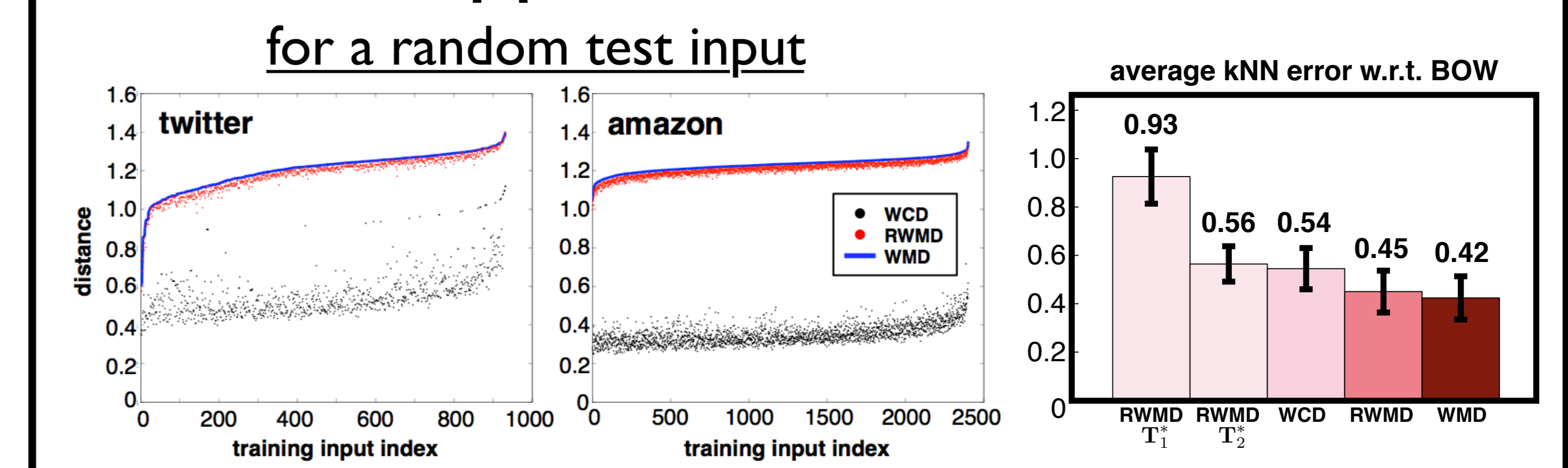


## Approximations

### 1. The Word Centroid Distance



## Approximation Results



## References

- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. In NIPS, pp. 3111–3119, 2013
- Collobert, R. and Weston, J. A unified architecture for natural language processing: Deep neural networks with multitask learning. In ICML, pp. 160–167. ACM, 2008.
- Mnih, A. and Hinton, G. E. A scalable hierarchical distributed language model. In NIPS, pp. 1081–1088, 2009.
- Rubner, Y., Tomasi, C., and Guibas, L. J. A metric for distributions with applications to image databases. In ICCV, pp. 59–66. IEEE, 1998.
- Salton, G. and Buckley, C. Term-weighting approaches in automatic text retrieval. Information processing & management, 24(5):513–523, 1988.
- Deerwester, S. C., Dumais, S. T., Landauer, T. K., Furnas, G. W., and Harshman, R. A. Indexing by latent semantic analysis. Journal of the American Society of Information Science, 41(6):391–407, 1990.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. Journal of Machine Learning Research, 3: 993–1022, 2003.
- Gardner, J., Kusner, M. J., Xu, E., Weinberger, K. Q., and Cunningham, J. Bayesian optimization with inequality constraints. In ICML, pp. 937–945, 2014.