# Latent Topic Models for Hypertext

Amit Gruber[1]
Michal Rosen-Zvi[2]
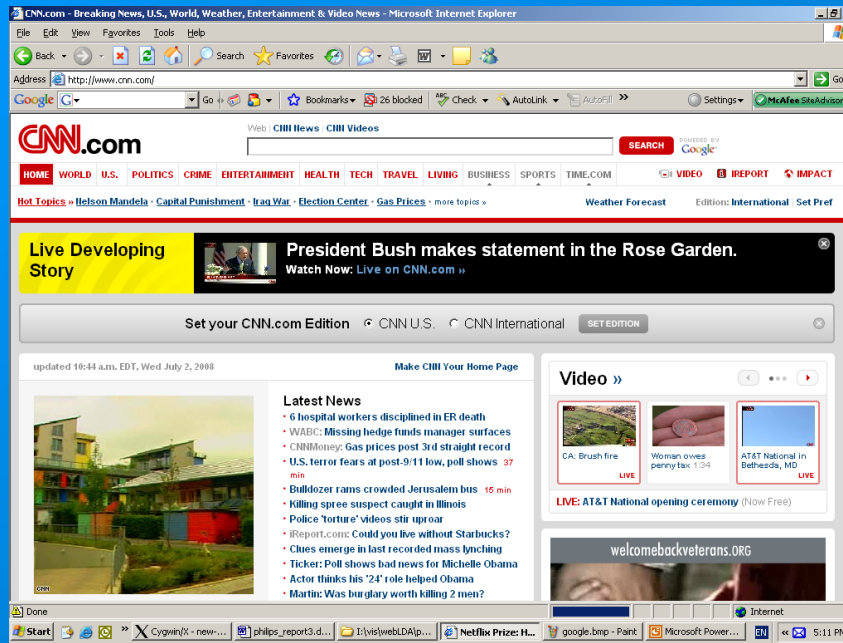Yair Weiss[1]


[1]The Hebrew University of Jerusalem
[2]IBM Research Labs, Haifa

# Introduction



0.3 sports

0.4 crime

0.3 politics
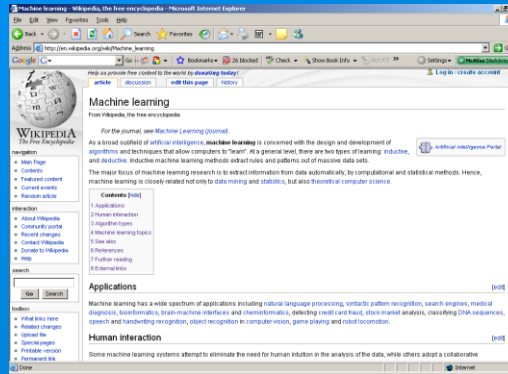
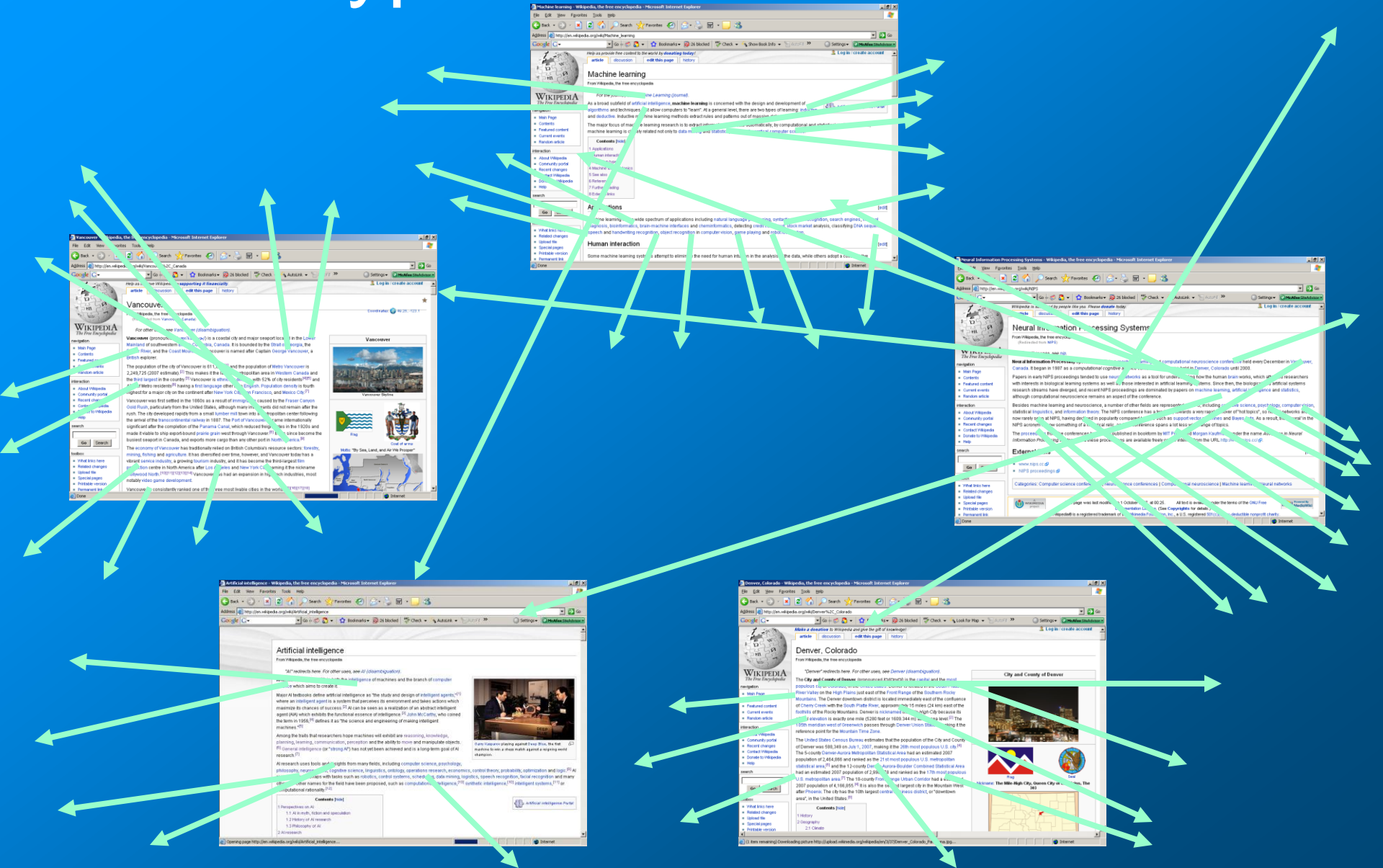- In this work we focus on hypertext documents, i.e. documents with links
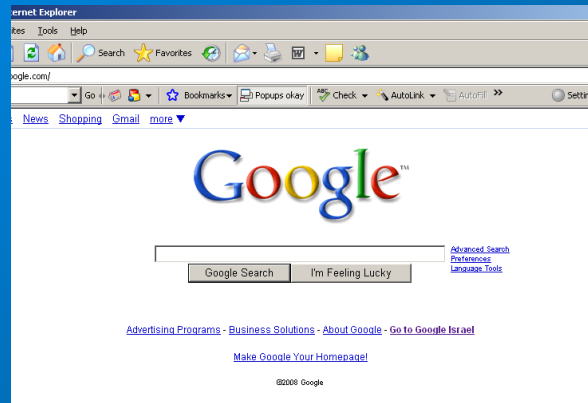
# Hypertext Documents

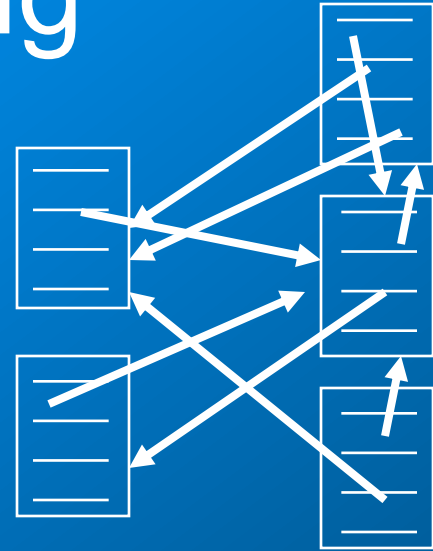# Hypertext Documents

# Introduction

- Hypertext is everywhere !
  - web pages, refs. in scientific publications
- Connectivity is important
  - PageRank



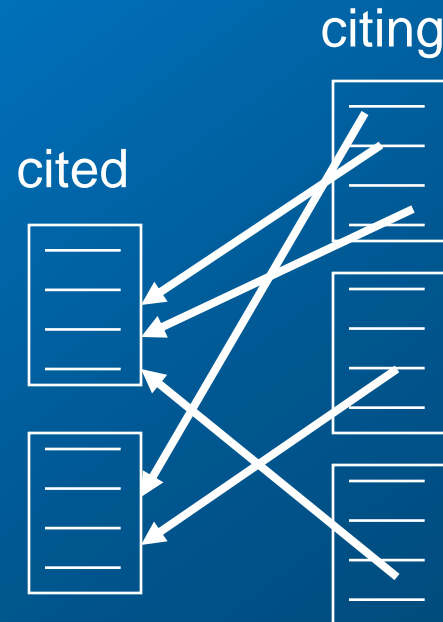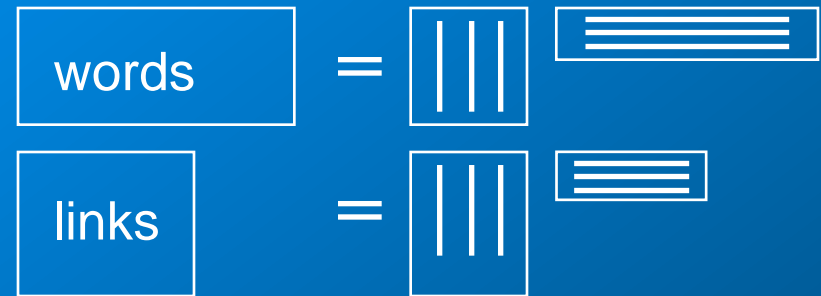- Topology of the WWW is complicated

# Problem Setting

- Input: documents and links



- Estimate:
  - Document topic mixture
  - Pr(word | topic)
  - Document importance

- Unsupervised

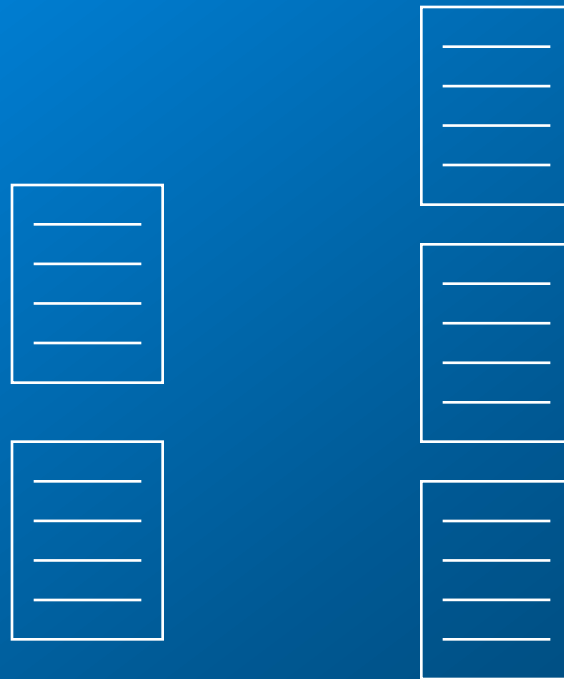# Previous Work: Topic Models for Hypertext

- Cohn and Hofmann, '01.
- Erosheva et al. '04.
  - Links are modeled similar to words
  - Links are not associated with words

- Dietz et al. '07.
- Nallapati and Cohen, '08.
  - Distinguish between citing and cited docs

words =

links =

citing

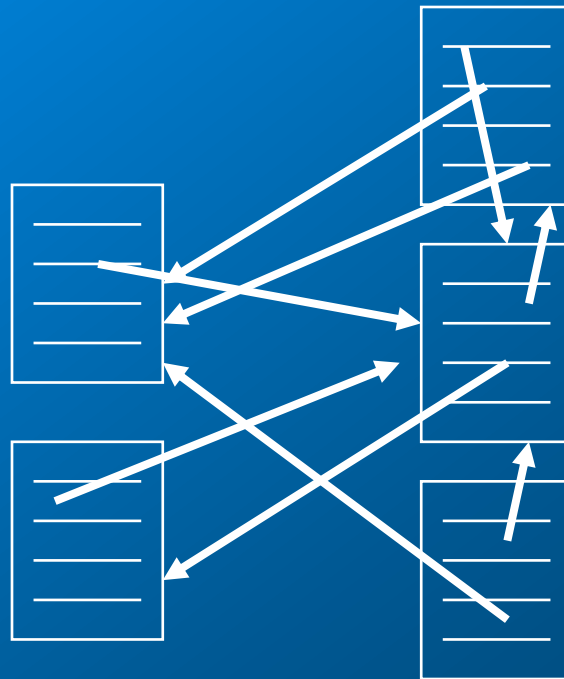cited

# The Latent Topic Hypertext Model

# (LTHM)

# LTHM: Generative Model

1. Words are created (by LDA)

# LTHM: Generative Model

1. Words are created (by LDA)

2. Links are created (our contribution)

# LTHM: Modeling Links

- Allows for arbitrary topology of the citation graph (including self links)


- A link points from a word to a document

# LTHM: Link Generation

- Depends on:
  - The topic of the anchor word
  - The topic mixture of the target document
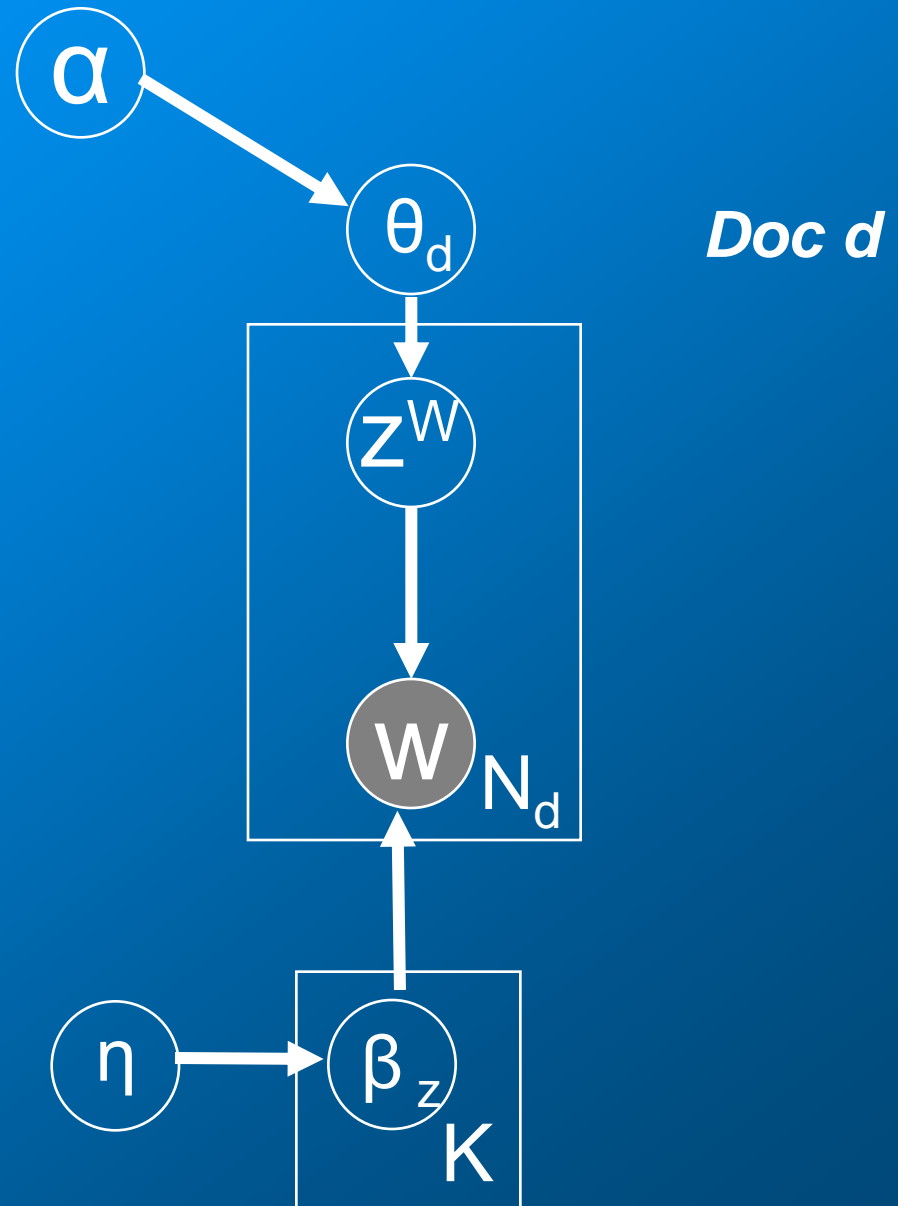  - The importance of the target document

$$\Pr(\text{link} = d \mid \text{topic} = z) = \lambda_d \theta_d(z)$$
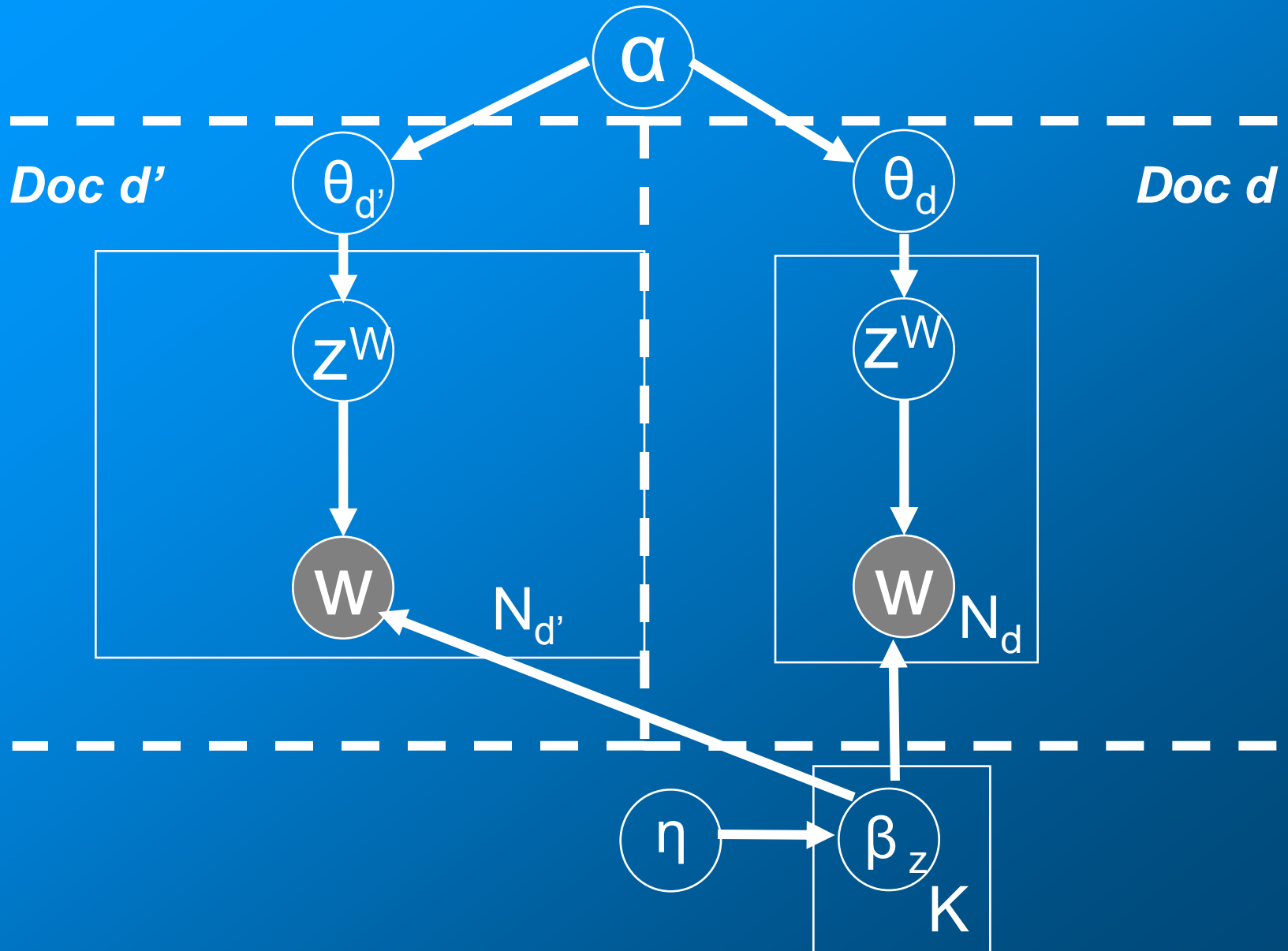
importance of d        prevalence of z in d
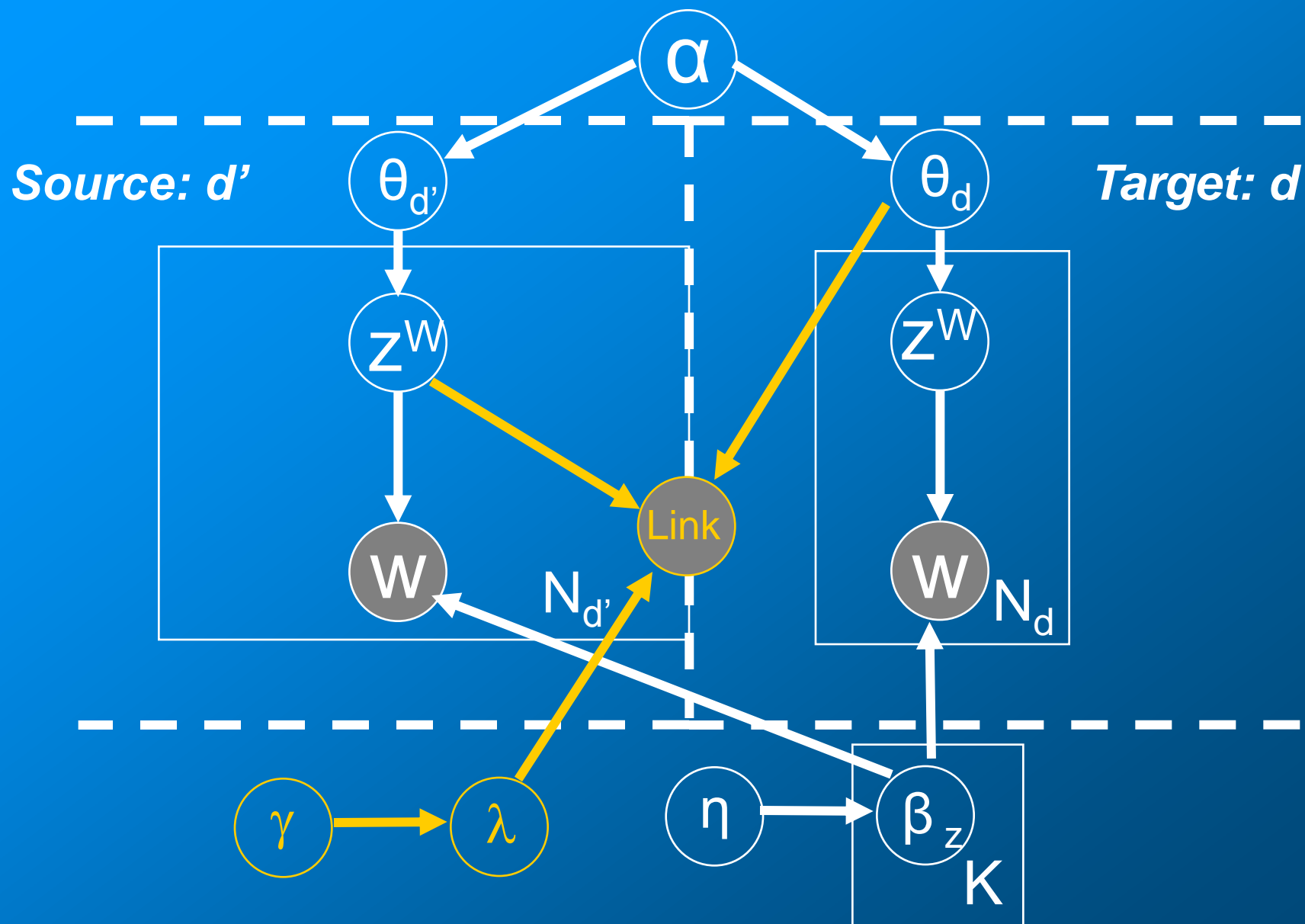
# Generating a single document d

α

θ_d

*Doc d*

d is generated by
Latent Dirichlet
Allocation

$z^w$

$w_{N_d}$

η

$β_z$

K

# Generating two documents d' and d

Generating Links from d' to d

# Properties of the Model

- D additional parameters ($\lambda$) for links vs. $D_x K$ parameters in previous models

- The existence (or non-existence) of a link is an observation

- A link shares the same topic with the word

- Link affects topic estimation in both the source and target documents

# Approximate Learning

- Exact inference is intractable in hierarchical models such as LDA

- Approximate inference in LTHM is even more challenging as non-links are also observations

- Using symmetries, we derived an O(K*corpus size) EM algorithm

# Experiments

- WebKB dataset

  8282 documents

  12911 links

- Wikipedia

  A new data set, collected by crawling from the NIPS Wikipedia page

  105 documents

  790 links

# Experiments: Wikipedia

# Experiments: Wikipedia

# Experiments: Wikipedia



## Topic 3

| | | |
|---|---|---|
| vancouver | 0.051 | Denver, Colorado |
| denver | 0.043 | |
| city | 0.041 | 0.0008 |
| retrieved | 0.024 | |
| colorado | 0.011 | |
| area | 0.009 | Vancouver |
| population | 0.009 | 0.0002 |
| canada | 0.008 | |

# Experiments: Wikipedia



**Topic 4**

| | |
|---|---|
| brain | 0.047 |
| cognitive | 0.026 |
| science | 0.016 |
| press | 0.011 |
| neurons | 0.010 |
| mind | 0.010 |
| systems | 0.010 |
| human | 0.010 |

Cognitive science 0.003

Neuroscience 0.002

# Journal of Machine Learning Research

## LDA

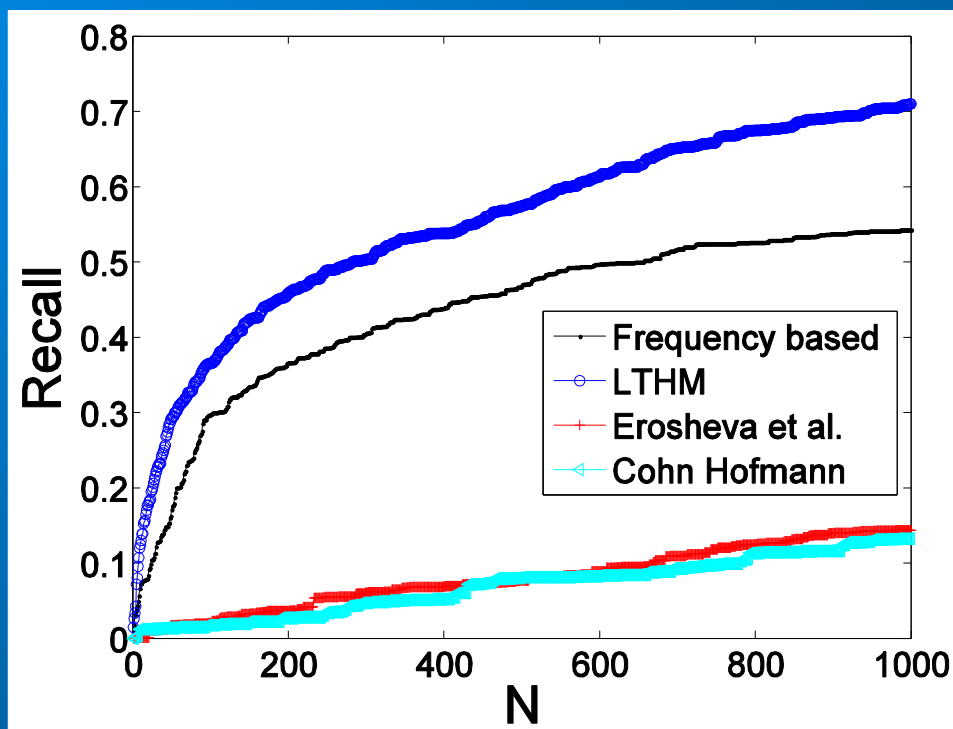| Topic prob | Top words |
| --- | --- |
| 0.1504 | search, article, navigation |
| 0.0798 | press, university, new |
| 0.0652 | learning, machine, algorithms |
| 0.0594 | fixes, skins, import |
| 0.0533 | model, regression, reasoning |

## LTHM

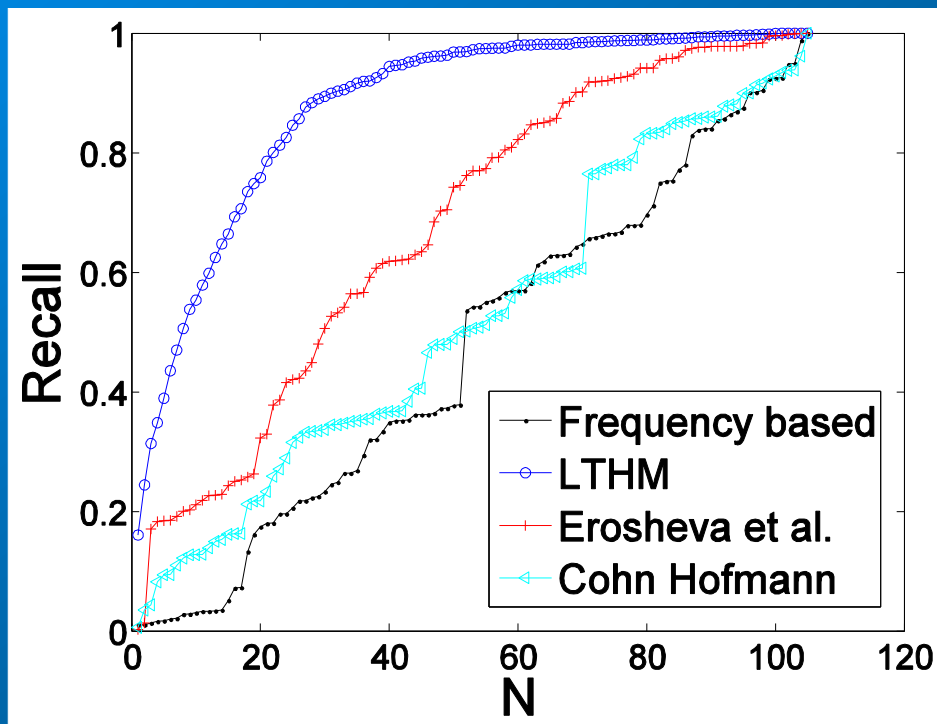| Topic prob | Top words |
| --- | --- |
| 0.4136 | learning, machine, engineering |
| 0.0943 | card, conference, credit |

# Experiments – link prediction on test set

- Wikipedia corpus: 105 documents with 790 links
- 20 hidden aspects
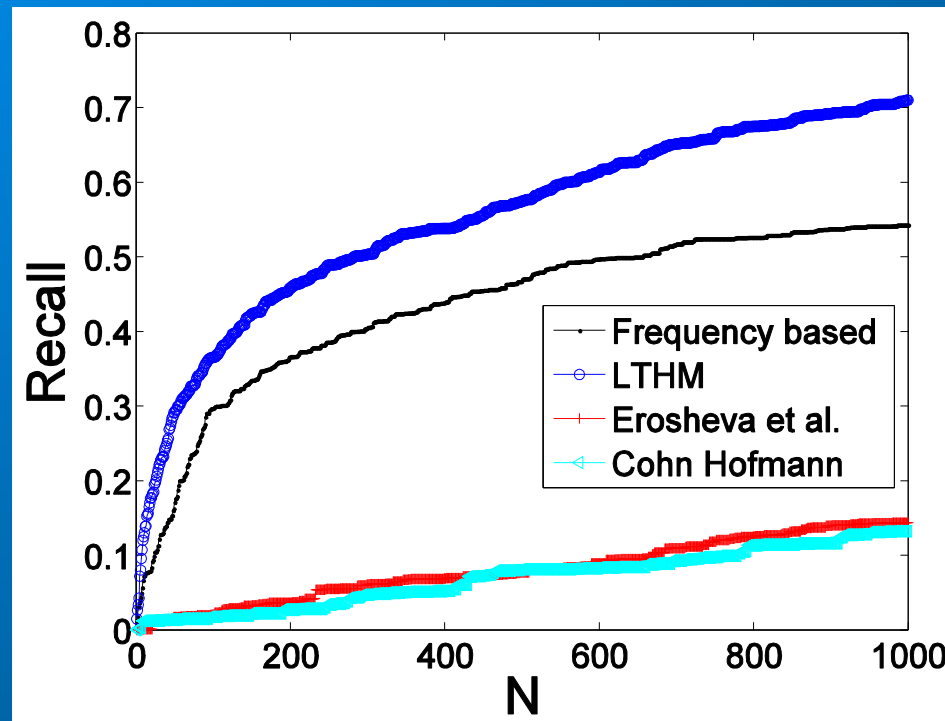- Test set: 11 documents, outgoing links are invisible

# Experiments – link prediction on train set

- Wikipedia corpus: 105 documents with 790 links
- 20 hidden aspects

# Experiments – link prediction

- Webkb corpus: 8282 documents with 12911 links

- 20 hidden aspects

- Test set: 10%

# Summary

- Explicit modeling of link generation in an LDA like model

- Efficient approximate inference algorithm

- Performs better than previous topic models in link recommendation

- Code and data available online at http://www.cs.huji.ac.il/~amitg/lthm.html