

# Entity Matching across Heterogeneous Sources

Yang Yang\*, Yizhou Sun<sup>†</sup>, Jie Tang\*, Bo Ma<sup>#</sup>, and Juanzi Li\*

\*Tsinghua University



清华大学

Tsinghua University

<sup>†</sup>Northeastern University



<sup>#</sup>Carnegie Mellon University



# Apple Inc. VS Samsung Co.

- A patent infringement suit starts from 2012.
  - Lasts 2 years, involves \$158+ million and 10 countries.
  - 7 / 35546 patents are involved.

Apple's patent

How to find **patents** relevant  
to a specific **product**?


Galaxy S II Skyrocket	X	X	X	X	X
Galaxy S III		X	X*		X
Galaxy Tab 2 (10.1)		X	X*		X
Stratosphere	X*	X	X*	X*	X

**SAMSUNG devices accused by APPLE.**

# Cross-Source Entity Matching

- Given an entity in a **source** domain, we aim to find its matched entities from **target** domain.
  - Product-patent matching;
  - Cross-lingual matching;
  - Drug-disease matching.

**Siri** →



**USPTO PATENT FULL-TEXT AND IMAGE DATABASE**

United States Patent  
Arrouye, et al.  
8,180,204  
December 27, 2011

**Abstract**

Universal interface for retrieval of information in a computer system

The present invention provides convenient access to items of information that are related to various descriptors input by a user, by means of a unitary interface which is capable of accessing information in a variety of locations, through a number of different techniques. Using a plurality of heuristic algorithms to operate upon information descriptors input by the user, the present invention locates and displays candidate items of information for selection and/or retrieval. Thus, the advantages of a search engine can be exploited, while listing only relevant object candidate items of information.

**Claim**

What is claimed is:

1. A method for locating information in a network using a computer, comprising: receiving by the computer an inputted information descriptor from a user-input device; providing said information descriptor received from the user-input device to a plurality of heuristic modules, wherein each heuristic module corresponds to a respective area of search and employs a different, predetermined heuristic algorithm corresponding to said respective area to search the area for information that corresponds to the received information descriptor, and the search areas include storage media accessible by the computer; searching by the heuristic modules, based on the received information descriptor, the respective areas of search using the predetermined heuristic algorithms corresponding to each respective area of search; providing at least one candidate item of information located by the heuristic modules as a result of said searching; and displaying by the computer a representation of said candidate item of information on a display device.
2. The method of claim 1, wherein the step of providing the at least one item of located information comprises: ranking each candidate item when a plurality of candidate items are located; and providing the plurality of candidate items for display based on the ranking of each candidate item.
3. The method of claim 2, wherein the step of ranking each candidate item comprises: ranking each candidate item according to a number of the heuristic modules locating the same item.
4. The method of claim 2, wherein the step of ranking each candidate item comprises: ranking each of the heuristic modules; and ranking the candidate items located by each heuristic module according to the ranking of the corresponding heuristic module.
5. The method of claim 2, wherein the step of ranking each candidate item comprises: ranking each candidate item according to a confidence factor associated with the located candidate item.

## Product-Patent matching

# Problem

$C_1$

$C_2$

Source 1: Siri's Wiki page

**Siri (Software)**

iOS iPhone  
iPod iPad  
intelligent personal assistant  
Cydia knowledge navigator  
voice control Apple server  
natural language user interface



Source 2: Patents

**Method for improving voice recognition**

heuristic algorithms  
distribution system  
speech recognition  
data source text-to-speech



**Universal interface for retrieval of information in a computer system**

search engine  
descriptors  
object relevant area  
ranking module  
rank candidate



**Voice menu system**

synchronize database  
host device media  
customized processor  
graphical user interface



...

Input 1: Dual source corpus

$\{C_1, C_2\}$ , where  $C_t = \{d_1, d_2, \dots, d_n\}$  is a collection of entities

Input 2: Matching relation matrix

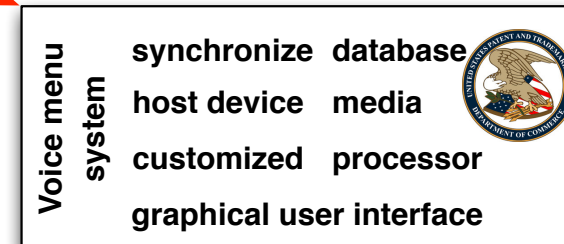
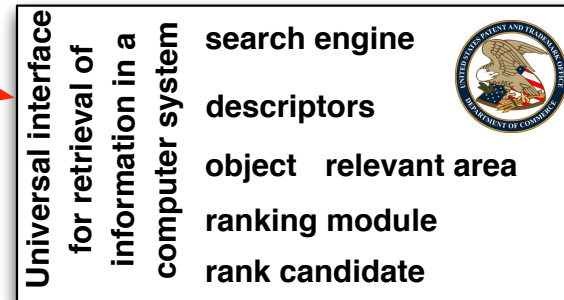
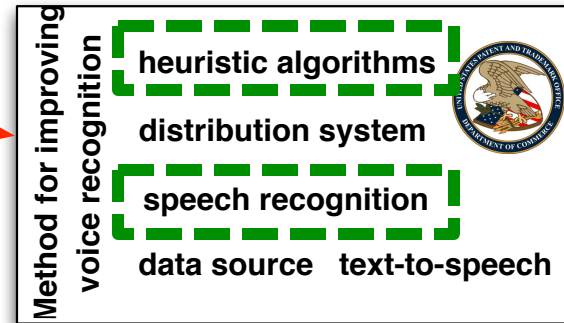
$$L_{ij} = \begin{cases} 1, & d_i \text{ and } d_j \text{ are matched} \\ 0, & \text{not matched} \\ ?, & \text{unknown} \end{cases}$$

# Challenges

Source 1: Siri's Wiki page



Source 2: Patents



...

1

Two domains have *less or no overlapping* in content

Daily expression

VS

Professional expression


# Challenges

Source 1: Siri's Wiki page


<b>Siri (Software)</b>	<b>iOS</b>	<b>iPhone</b>	
	<b>iPod</b>	<b>iPad</b>	
	<b>intelligent personal assistant</b>		
	<b>Cydia</b>	<b>knowledge navigator</b>	
	<b>voice control</b>	<b>Apple server</b>	
	<b>natural language user interface</b>		

Topic:  
voice control  
0.83


Source 2: Patents

Method for improving voice recognition	heuristic algorithms	
	distribution system	
	speech recognition	
	data source text-to-speech	

Topic: ranking  
0.54

Universal interface for retrieval of information in a computer system	search engine	
	descriptors	
	object relevant area	
	ranking module	
	rank candidate	

Topic: ???

Voice menu system	synchronize database	
	host device media	
	customized processor	
	graphical user interface	

...

1

Two domains have *less or no overlapping* in content

2

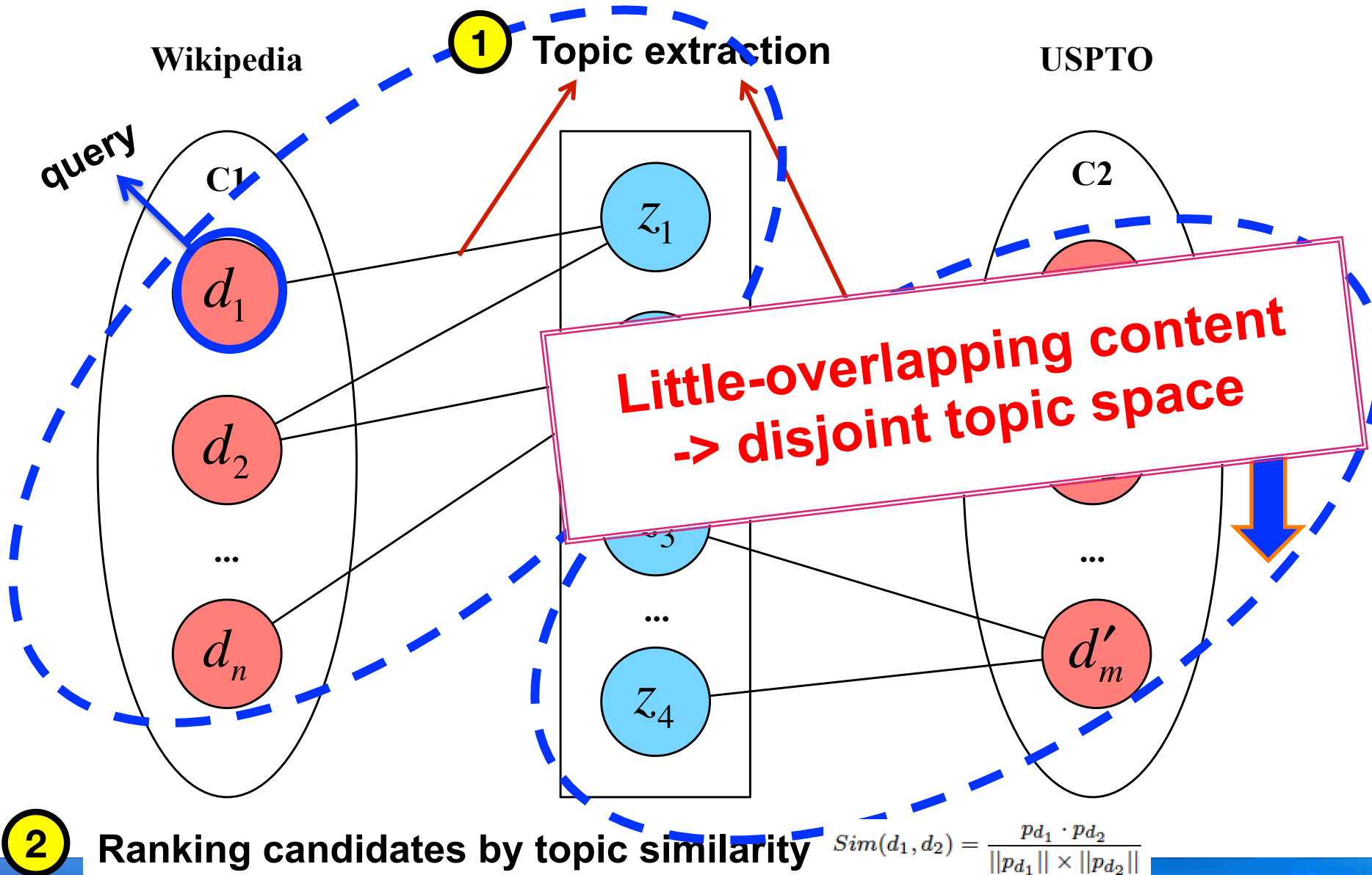
How to model the topic-level relevance probability



# Our Approach

## Cross-Source Topic Model

# Baseline





# Cross-Sampling

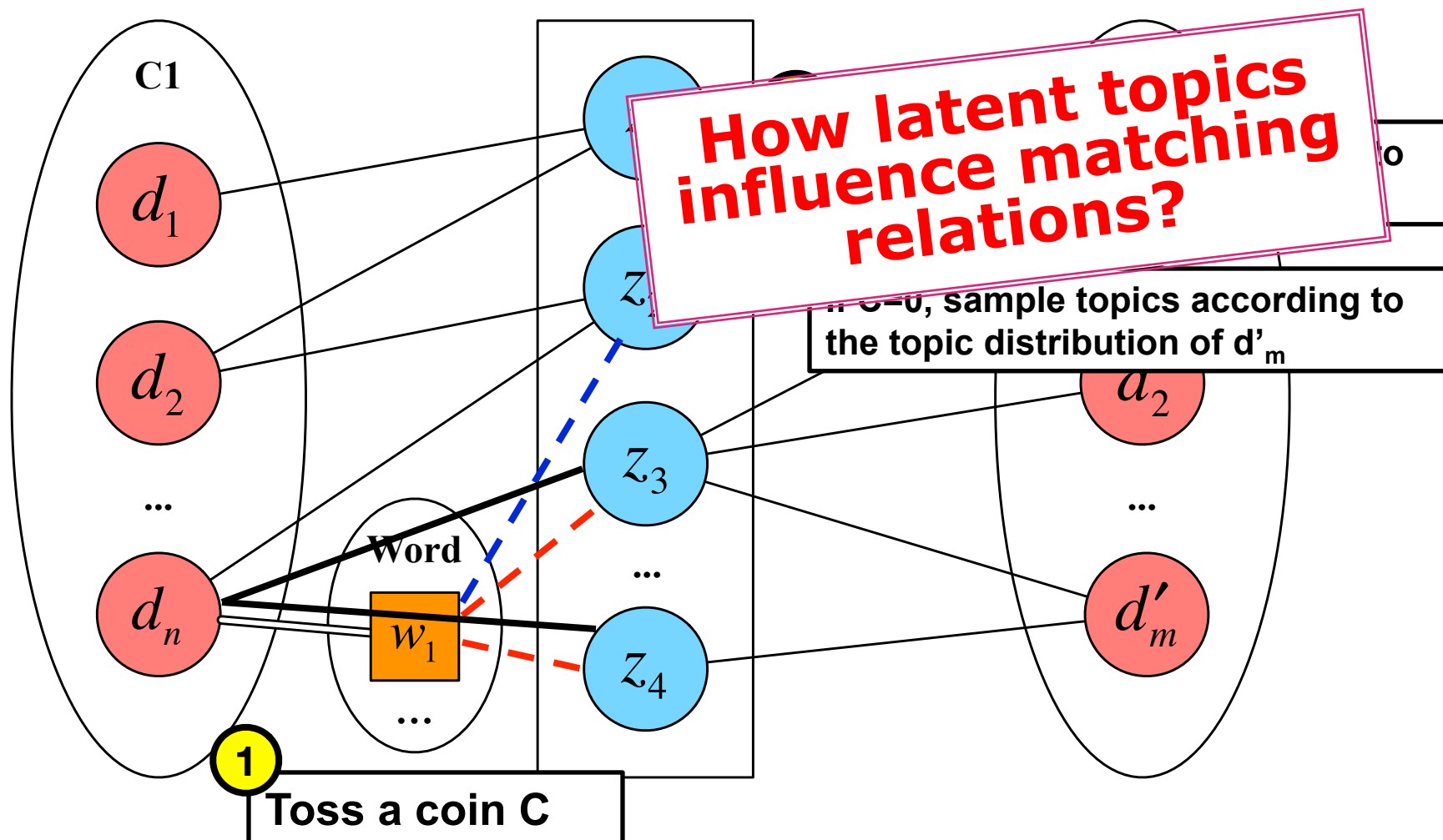
$d_n$  is matched with  $d'_m$

Bridge topic space by leveraging known matching relations.

Wikipedia

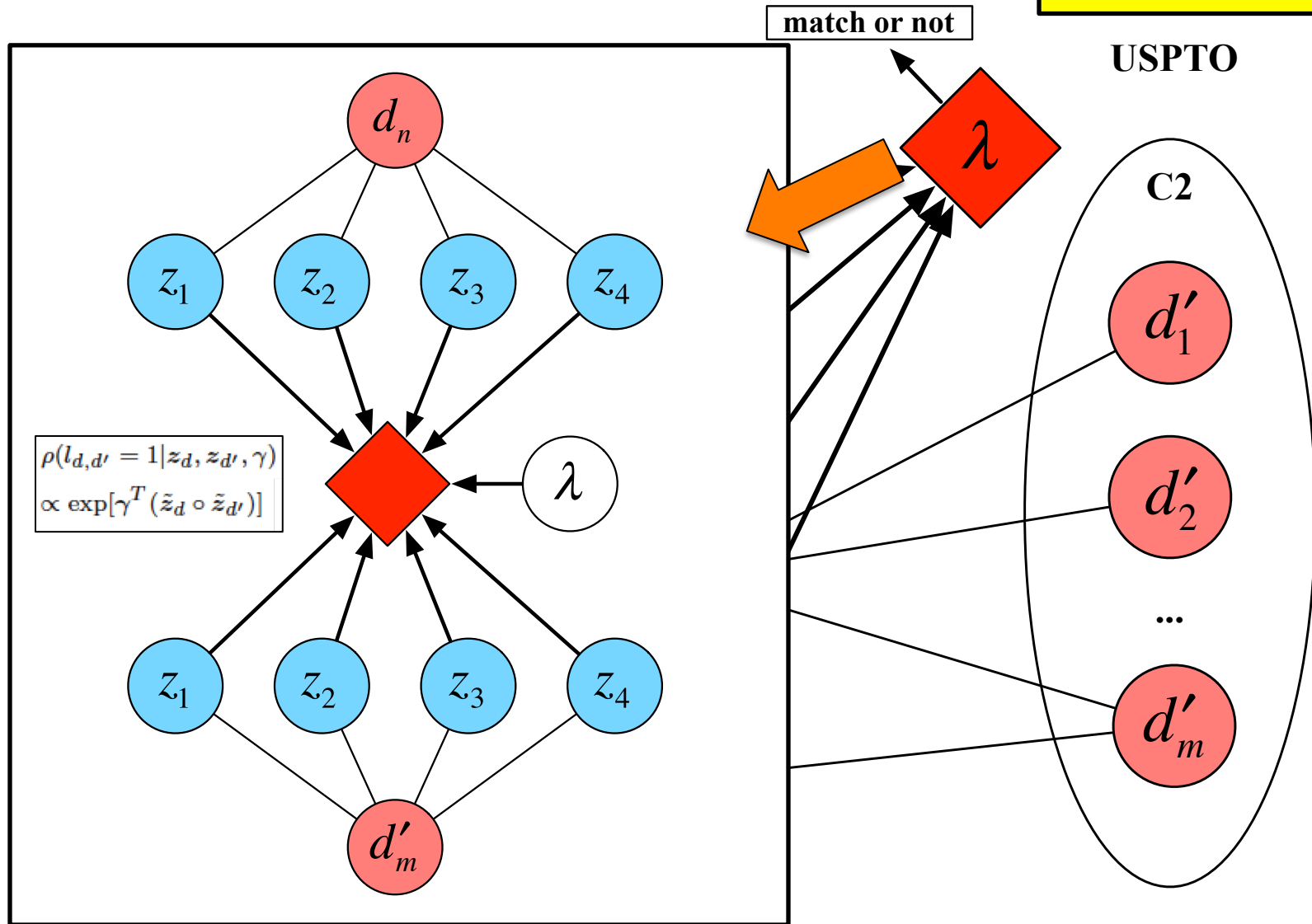
Topics

USPTO

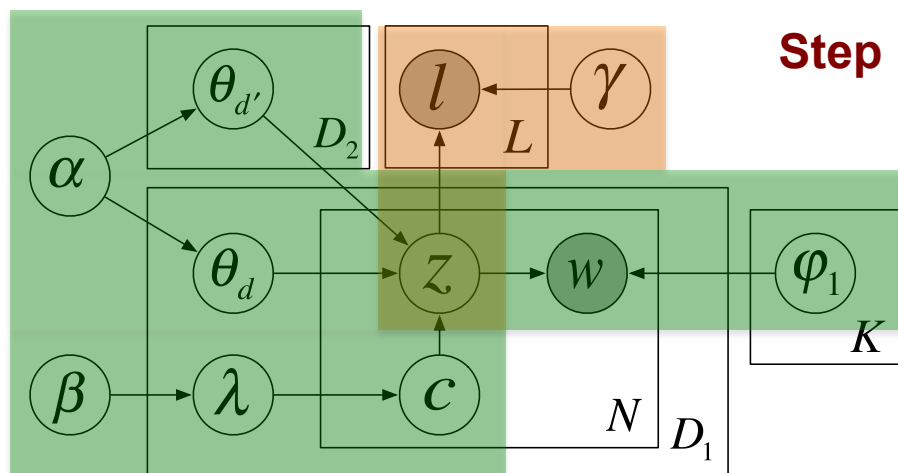


# Inferring Matching Relation

Infer matching relations by leveraging extracted topics.



# Cross-Source Topic Model



**Step 1:**

```

Input: a dual source corpus  $C$ , a matching relation matrix  $L$ ,
and hyper-parameters  $\alpha$  and  $\beta$ 
foreach entity  $d$  do
  | Generate  $\theta_d \sim \text{Dir}(\alpha)$ ;
end
% cross-sampling-based entity generation
foreach  $d$  in each source  $t$  do
  Set  $\beta$  according to  $L_d$ ;
  Generate  $\lambda_d \sim \text{Dir}(\beta)$ ;
  for  $n = 1$  to  $N_d$  do
    Generate  $c_{d,n} \sim \text{Mult}(\lambda_d)$ ,  $c_{d,n}$  can be  $d$  or the index
of matched entities with  $d$ ;
    Draw a topic  $z_{d,n} \sim \text{Mult}(\theta_{c_{d,n}})$  from the topic
distribution of the entity  $c$ ;
    Draw a word  $w_{d,n} \sim \text{Mult}(\varphi_{t,z_{d,n}})$  from  $z_{d,n}$ -specific
word distribution;
  end
end
% matching relation generation
foreach  $(d, d')$  with possible links do
  | Generate  $l_{d,d'} \sim \rho(\cdot | z_d, z_{d'}, \gamma)$ ;
end
    
```

**Step 2:**

Latent topics  $\longleftrightarrow$  Matching relations

# Model Learning

- Variational EM

- Model parameters:  $\{\varphi, \gamma\}$
- Variational parameters:  $\{\vartheta, \tau, \eta, \epsilon\}$
- E-step:

$$\begin{aligned}\eta_{d,c} &= \beta_{d,c} + N_d \times \epsilon_{d,c} \\ \tau_{d,k} &= \alpha_k + \sum_{n=1}^{N_d} \vartheta_{d,n,k} \\ \epsilon_{d,n,c} &\propto \exp\{\Psi(\eta_{d,c}) - \Psi(\sum_{i \in R(d)} \eta_{d,i})\} \\ \vartheta_{d,n,k} &\propto \sum_{d' \in \{R(d), d\}} (\exp\{\sum_{d'' \neq d'} \frac{\gamma_k \sum_{i=1}^{N_{d''}} \vartheta_{d'',i,k}}{N_{d'} N_{d''}} \\ &\quad + \Psi(\tau_{d',k}) - \Psi(\sum_{j=1}^K \tau_{d',j})\}) \epsilon_{d,n,d'} \times \varphi_{t,k,v}\end{aligned}$$

- M-step:

$$\begin{aligned}\varphi_{t,k,v} &\propto \sum_{d=1}^{D_t} \sum_{n=1}^{N_d} \vartheta_{d,n,k} \mathbf{1}(w_{d,n}^t = v) \\ \gamma_k &= \frac{\sum_{d,d'} 1}{2 \sum_{d,d'} l_{d,d'} [(\Upsilon_d - \Upsilon_{d'}) \circ (\Upsilon_d - \Upsilon_{d'})]_k}\end{aligned}$$

**Input:** a dual source corpus  $C$ , a matching relation matrix  $L$ , and hyper-parameters  $\alpha$  and  $\beta$

Initialize  $\{\vartheta, \tau, \eta, \epsilon, \varphi, \gamma\}$  randomly;

repeat

  % E-Step: optimize the ELBO;

  foreach  $d$  in each source  $t$  do

    for  $c = 0$  to 1 do

      | Update  $\eta_{d,c}$  according to Eq. 6;

    end

    for  $k = 1$  to  $K$  do

      | Update  $\tau_{d,k}$  according to Eq. 7;

    end

    for  $n = 1$  to  $N_d$  do

      for  $c = 0$  to 1 do

        | Update  $\epsilon_{d,n,c}$  according to Eq. 8;

      end

      for  $k = 1$  to  $K$  do

        | Update  $\vartheta_{d,n,k}$  according to Eq. 9;

      end

    end

  end

  % M-Step: maximize the resulting ELBO;

  foreach topic  $k$  in each source  $t$  do

    foreach term  $v$  do

      | Update  $\varphi$  according to Eq. 10;

    end

    Update  $\gamma_k$  according to Eq. 11;

  end

until Convergence;



# Experiments

Task I: Product-patent matching

Task II: Cross-lingual matching

# Task I: Product-Patent Matching

- Given a Wiki article describing a product, finding all patents relevant to the product.
- Data set:
  - 13,085 Wiki articles ;
  - 15,000 patents from USPTO;
  - 1,060 matching relations in total.

# Experimental Results

**Training:** 30% of the matching relations randomly chosen.

Method	P@3	P@20	MAP	R@3	R#20	MRR
CS+LDA	0.111	0.083	0.109	0.011	0.046	0.053
RW+LDA	0.111	0.117	0.123	0.033	0.233	0.429
RTM	0.501	0.233	0.416	0.057	0.141	0.171
RW+CST	<b>0.667</b>	0.167	0.341	<b>0.200</b>	0.333	0.668
CST	<b>0.667</b>	<b>0.250</b>	<b>0.445</b>	0.171	<b>0.457</b>	<b>0.683</b>

**Content Similarity based on LDA (CS+LDA):** cosine similarity between two entities' topic distribution extracted by LDA.

**Random Walk based on LDA (RW+LDA):** random walk on a graph where edges indicate the hyperlinks between Wiki articles and citations between patents.

**Relational Topic Model (RTM):** used to model links between documents.

**Random Walk based on CST (RW+CST):** uses CST instead of LDA comparing with RW+LDA.

# Task II: Cross-lingual Matching

- Given an English Wiki article, we aim to find a Chinese article reporting the same content.
- Data set:
  - 2,000 English articles from Wikipedia;
  - 2,000 Chinese articles from Baidu Baike;
  - Each English article corresponds to one Chinese article.



# Experimental Results

**Training:** 3-fold cross validation

Method	Precision	Recall	F1-Measure	F2-Measure
Title Only	<b>1.000</b>	0.410	0.581	0.465
SVM-S	0.957	0.563	0.709	0.613
LFG	0.661	0.820	0.732	0.782
LFG+LDA	0.652	0.805	0.721	0.769
LFG+CST	0.682	<b>0.849</b>	<b>0.757</b>	<b>0.809</b>

**Title Only:** only considers the (translated) title of articles.

**SVM-S:** famous cross-lingual Wikipedia matching toolkit.

**LFG<sup>[1]</sup>:** mainly considers the structural information of Wiki articles.

**LFG+LDA:** adds content feature (topic distributions) to LFG by employing LDA.

**LFG+CST:** adds content feature to LFG by employing CST.

# Topics Relevant to Apple and Samsung

(Topic titles are hand-labeled)

Title	Top Patent Terms	Top Wiki Terms
Gravity Sensing	Rotational, gravity, interface, sharing, frame, layer	Gravity, iPhone, layer, video, version, menu
Touchscreen	Recognition, point, digital, touch, sensitivity, image	Screen, touch, iPad, os, unlock, press
Application Icons	Interface, range, drives, icon, industrial, pixel	Icon, player, software, touch, screen, application

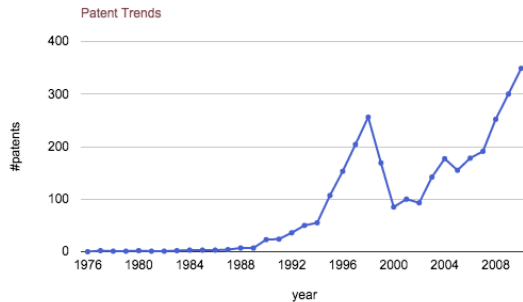
# Prototype System

competitor analysis @ <http://pminer.org>

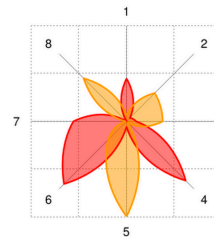
## Apple Computer, Inc.

### Trend Analysis:

Patent Trends Topic Trends Inventor Trends



### Topic Comparison



- 1.Electrical computers
- 2.Static information
- 3.Information sotrage
- 4.Data processing
- 5.Active solid-state devices
- 6.Computer graphics processing
- 7.Molecular biology and microbiology
- 8.Semiconductor device manufacturing

Radar Chart: topic comparison

### Competitors:

Global Evolution Topic Topic Evolution

International Business Machines Corporation  
# patents: 60180

home compare

Microsoft Corporation  
# patents: 14100

home compare

Hewlett-Packard Development Company, L.P.  
# patents: 21015

home compare

Sun Microsystems, Inc.  
# patents: 7291

home compare

Intel Corporation  
# patents: 18682

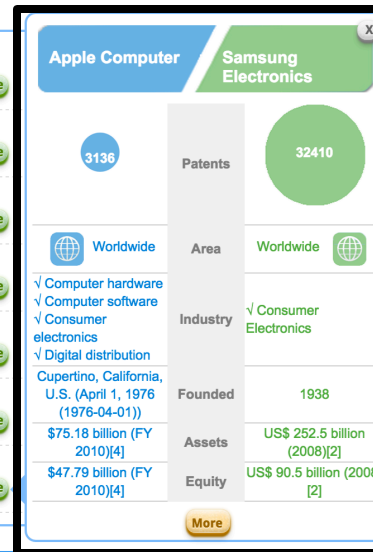
home compare

Hitachi, Ltd.  
# patents: 38863

home compare

Samsung Electronics Co., Ltd.  
# patents: 32410

home compare



Basic information comparison:

#patents, business area, industry, founded year, etc.

# Conclusion

- Study the problem of **entity matching across heterogeneous sources**.
- Propose the cross-source topic model, which integrates the **topic extraction** and **entity matching** into a unified framework.
- Conduct two experimental tasks to demonstrate the effectiveness of CST.

# Thank You!

## Entity Matching across Heterogeneous Sources

Yang Yang\*, Yizhou Sun<sup>†</sup>, Jie Tang\*, Bo Ma<sup>#</sup>, and Juanzi Li\*

\*Tsinghua University



清华大学

Tsinghua University

<sup>†</sup>Northeastern University



<sup>#</sup>Carnegie Mellon University



# Apple Inc. VS Samsung Co.

- A patent infringement lawsuit starts from 2012.
  - Nexus S, Epic 4G, Galaxy S 4G, and the Samsung Galaxy Tab, infringed on Apple's intellectual property: its patents, trademarks, user interface and style.
  - Lasts over **2 years**, involves **\$158+ million**.
- How to find **patents** relevant to a specific **product**?



# Problem

- Given an entity in a **source** domain, we aim to find its matched entities from **target** domain.
  - Given a textual description of a **product**, finding related **patents** in a patent database.
  - Given an **English Wiki page**, finding related **Chinese Wiki pages**.
  - Given a specific **disease**, finding all related **drugs**.

# Basic Assumption

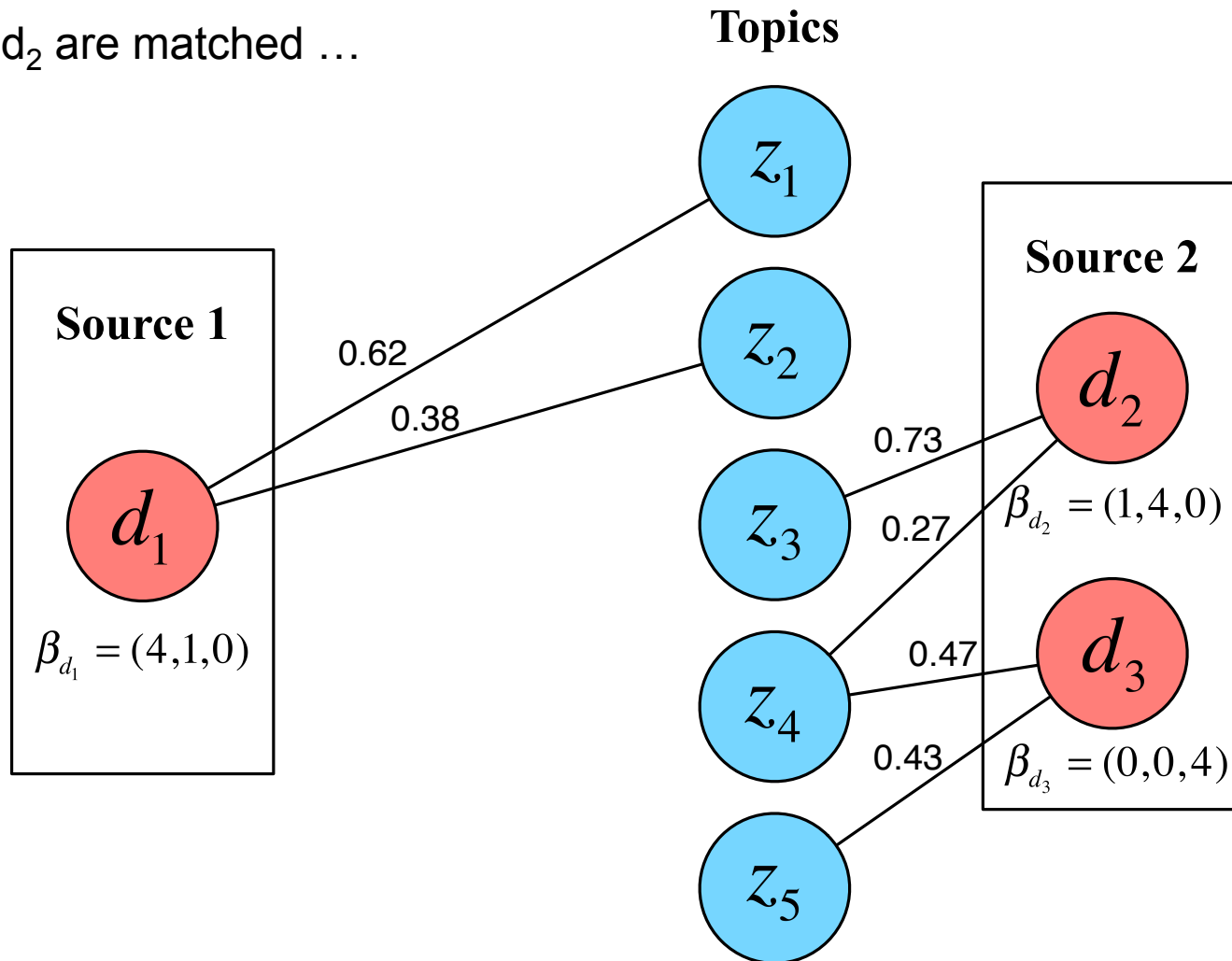
- For entities from different sources, their **matching relations** and **hidden topics** are influenced by each other.
- How to leverage the known matching relations to help **link** hidden topic spaces of two sources?



# Cross-Sampling

1

$d_1$  and  $d_2$  are matched ...

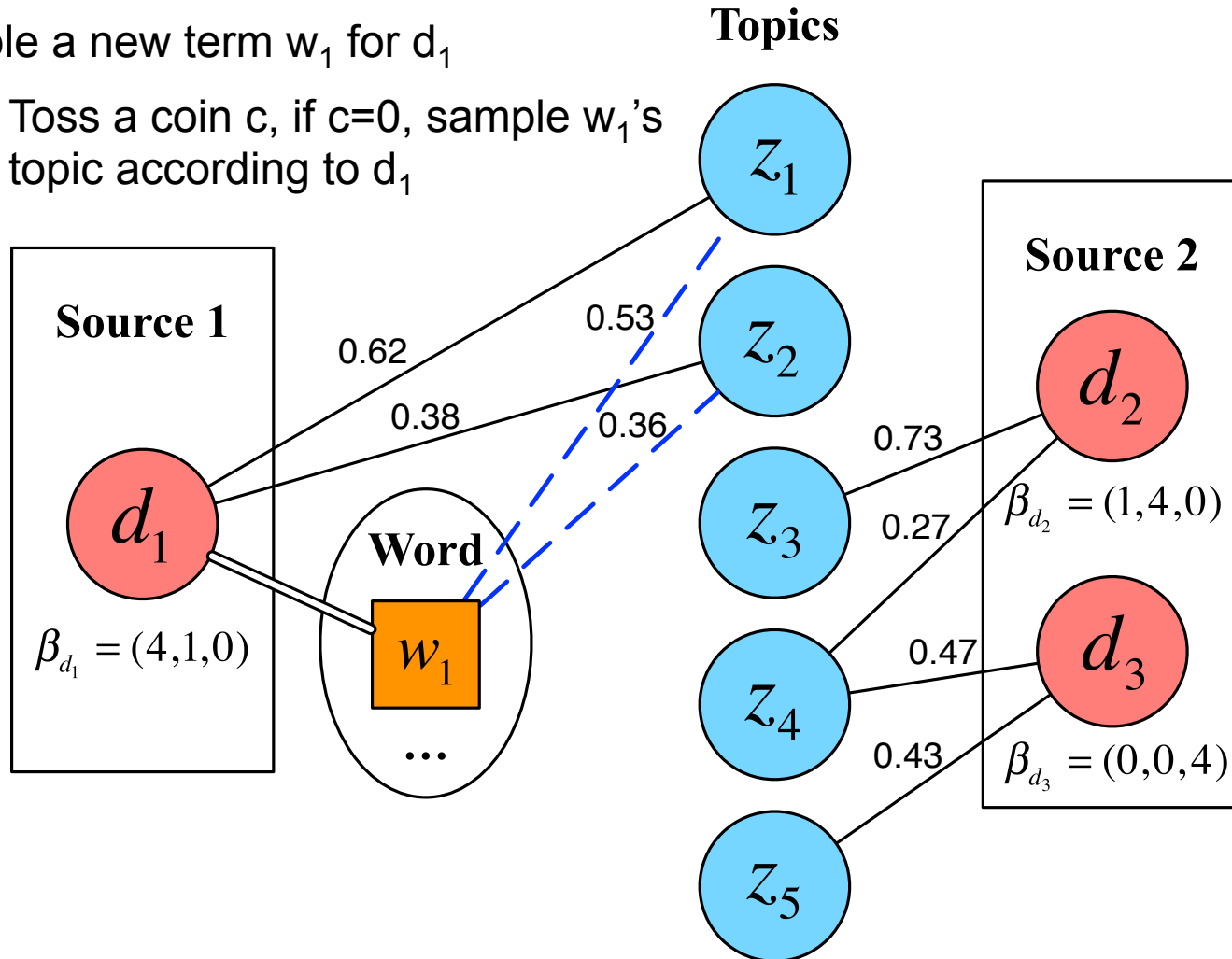


# Cross-Sampling

2

Sample a new term  $w_1$  for  $d_1$

└ Toss a coin  $c$ , if  $c=0$ , sample  $w_1$ 's topic according to  $d_1$

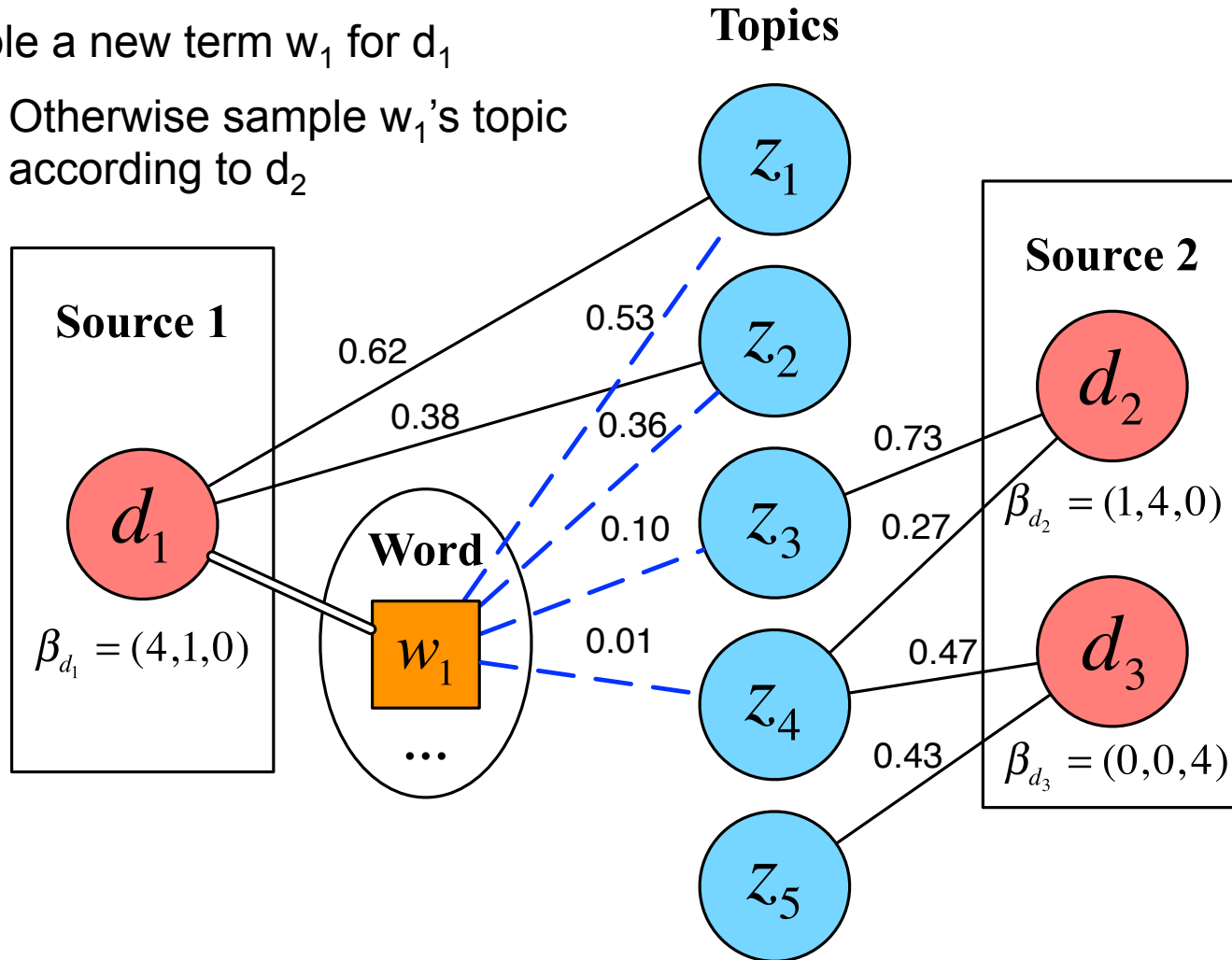


# Cross-Sampling

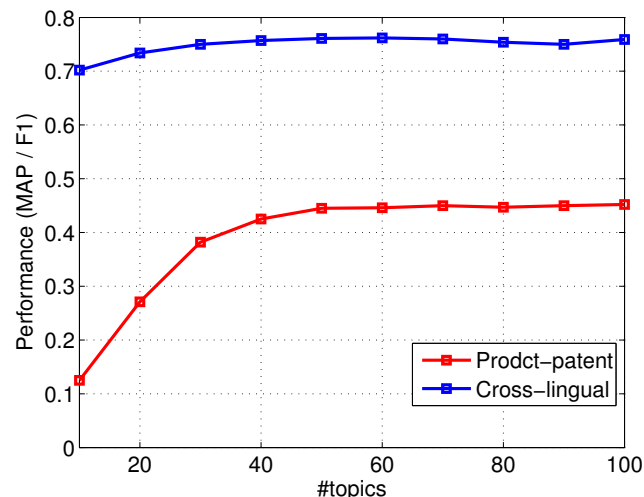
3

Sample a new term  $w_1$  for  $d_1$

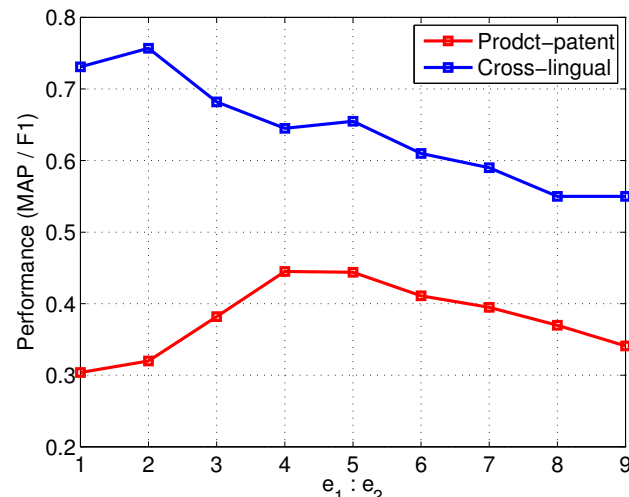
↳ Otherwise sample  $w_1$ 's topic according to  $d_2$



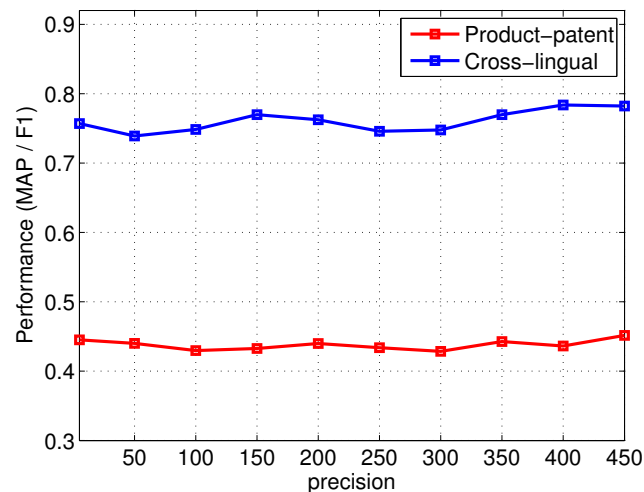
# Parameter Analysis



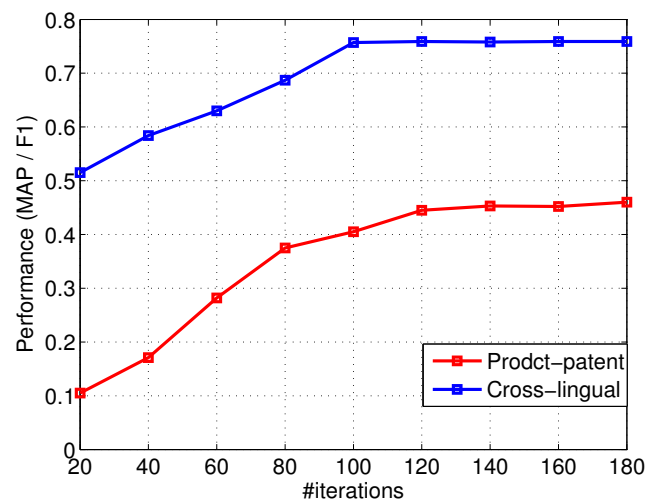
(a) Number of topics K



(b) Ratio



(c) Precision



(d) Convergence analysis

# Experimental Results

**Training:** 30% of the matching relations randomly chosen.

Method	P@3	P@20	MAP	R@3	R#20	MRR
CS+LDA	0.111	0.083	0.109	0.011	0.046	0.053
RW+LDA	0.111	0.117	0.123	0.033	0.233	0.429
RTM	0.501	0.233	0.416	0.057	0.141	0.171
RW+CST	<b>0.667</b>	0.167	0.341	<b>0.200</b>	0.333	0.668
CST	<b>0.667</b>	<b>0.250</b>	<b>0.445</b>	0.171	<b>0.457</b>	<b>0.683</b>

**Content Similarity based on LDA (CS+LDA):** cosine similarity between two articles' topic distribution extracted by LDA.

**Random Walk based on LDA (RW+LDA):** random walk on a graph where edges indicate the hyperlinks between Wiki articles and citations between patents.

**Relational Topic Model (RTM):** used to model links between documents.

**Random Walk based on CST (RW+CST):** uses CST instead of LDA comparing with RW+LDA.

# Experimental Results

**Training:** 30% of the matching relations randomly chosen.

Method	P@3	P@20	MAP	R@3	R#20	MRR
CS+LDA	0.111	0.083	0.109	0.011	0.046	0.053
RW+LDA	0.111	0.117	0.123	0.033	0.233	0.429
RTM	0.501	0.233	0.416	0.057	0.141	0.171
RW+CST	<b>0.667</b>	0.167	0.341	<b>0.200</b>	0.333	0.668
CST	<b>0.667</b>	<b>0.250</b>	<b>0.445</b>	0.171	<b>0.457</b>	<b>0.683</b>

**Content Similarity based on LDA (CS+LDA):** cosine similarity between two articles' topic distribution extracted by LDA.

**Random Walk based on LDA (RW+LDA):** random walk on a graph where edges indicate the hyperlinks between Wiki articles and citations between patents.

**Relational Topic Model (RTM):** used to model links between documents.

**Random Walk based on CST (RW+CST):** uses CST instead of LDA comparing with RW+LDA.