

# Analyzing NIH Funding Patterns with Statistical Text Analysis

---

**Margaret Blume-Kohout**

New Mexico Consortium

**Jihyun Park**

**Eric Nalisnick**

**Padhraic Smyth**

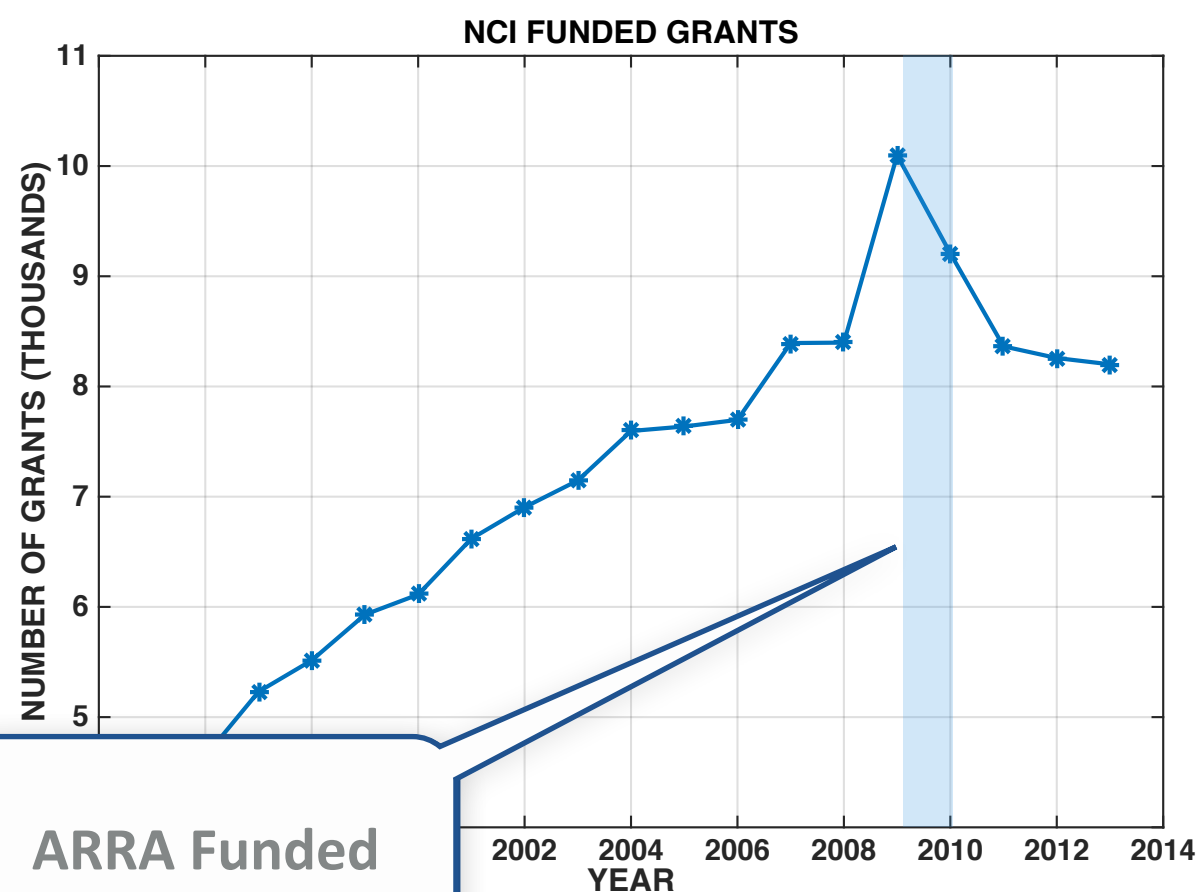
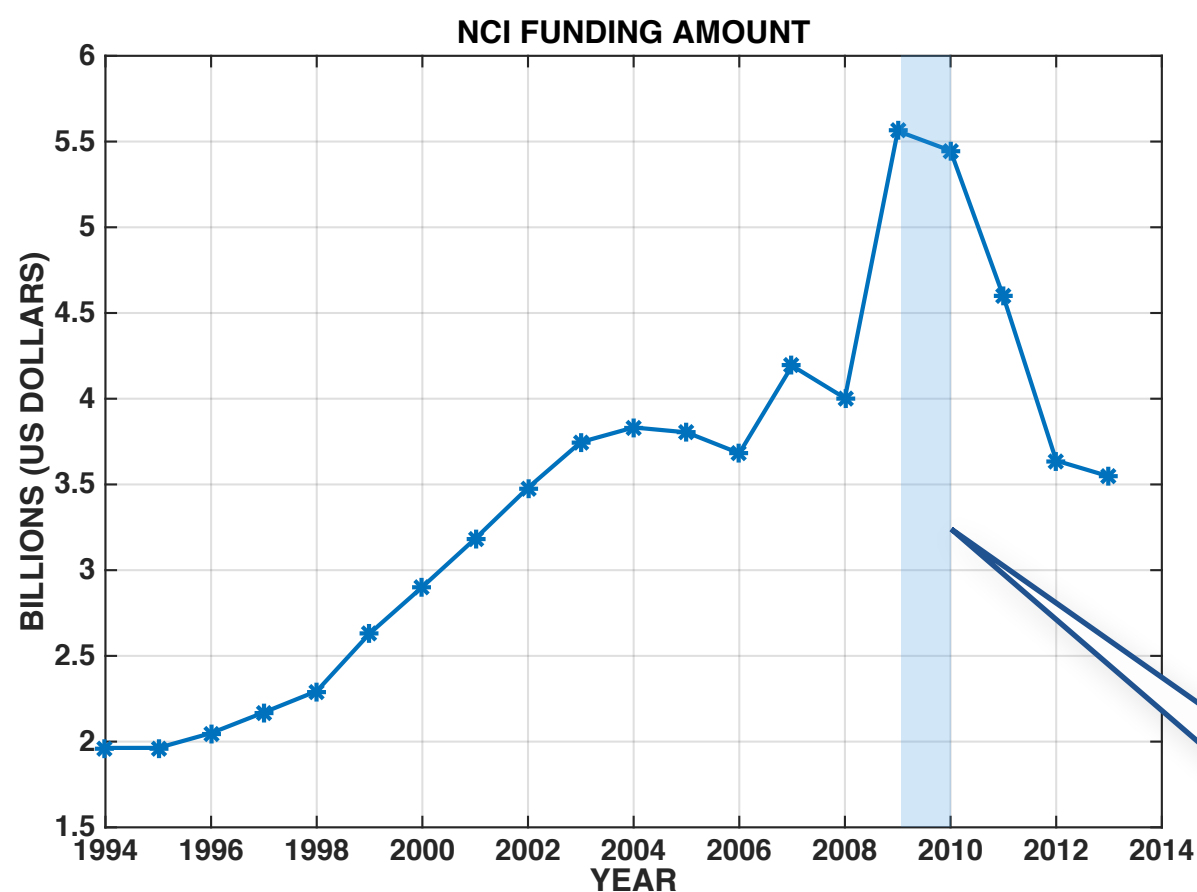
Dept. Of Computer Science  
University Of California, Irvine

**Ralf Krestel**

Web Science Research Group  
Hasso-Plattner-Institut

# Measuring the Impact of NIH(National Institute of Health) Funding

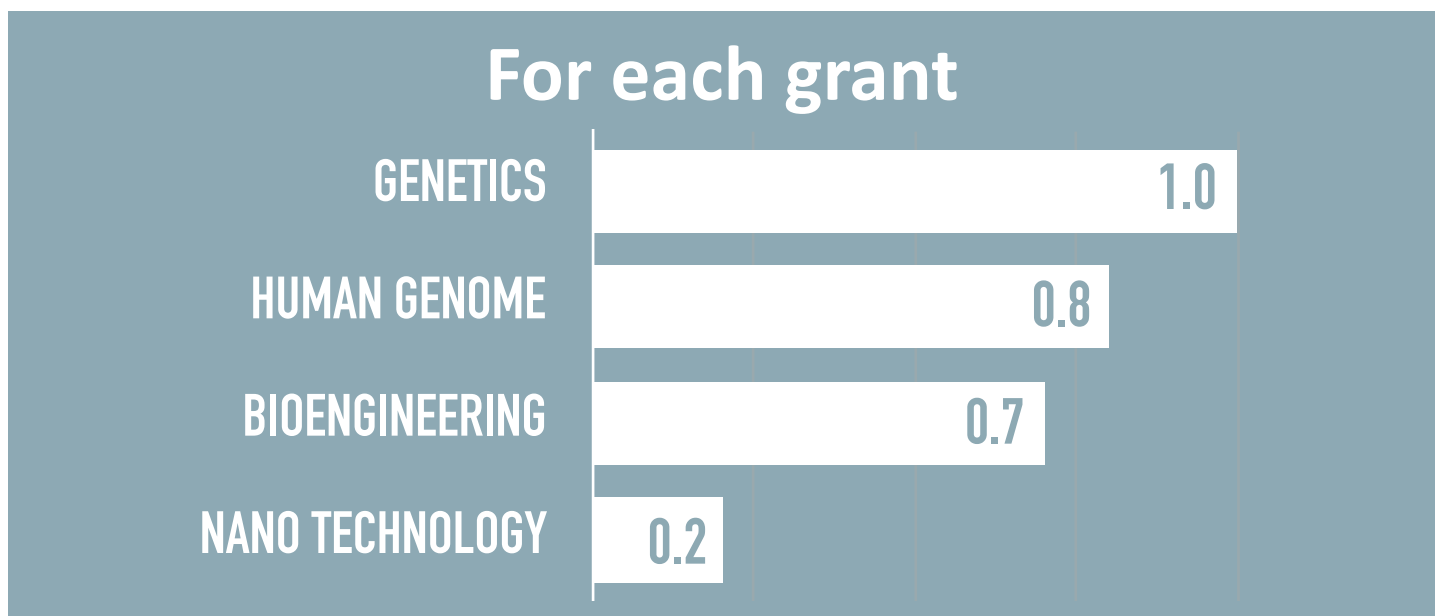
- ▶ NIH invests over \$30 billion each year
- ▶ Can we gain insight into this process using text and metadata?
- ▶ Our approach is to use statistical topic modeling
- ▶ We used grants data from NCI (National Cancer Institute)



ARRA Funded

# Overview

NCI Data					
PROJECT ID	GRANT ABSTRACT	RCDC Labels	FUNDING	YEAR	...
	...				

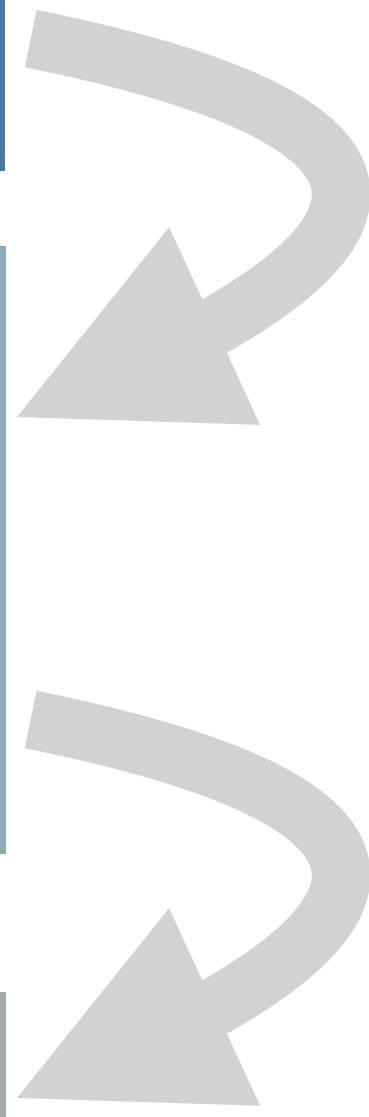


Probability of each label being associated with the grant

Funding patterns over time for each area

Text Classification Techniques

Funding + Year information



# NCI (National Cancer Institute) Data

---

- ▶ **Grant abstracts from 1994 through 2013**
- ▶ **Text Processing**
  - ▶ BOW representation
  - ▶ Removed 500 common stopwords
  - ▶ Extracted noun-phrase terms using a NLP parser
- ▶ **BOW Data**
  - ▶ Total 149,901 documents
  - ▶ Number of documents with labels (training data) : 31,628 (2008~2011)
  - ▶ Number of documents without labels : 118,273
  - ▶ Size of vocabulary (W) : 29,713

# LDA: Topics are Represented as Distributions over Words

Figures from Mark Steyvers

## Terrorism

SEPT\_11  
WAR  
SECURITY  
IRAQ  
TERRORISM  
NATION  
KILLED  
AFGHANISTAN  
ATTACKS  
OSAMA\_BIN\_LADEN

## Wall Street Firms

WALL\_STREET  
ANALYSTS  
INVESTORS  
FIRM  
GOLDMAN\_SACHS  
FIRMS  
INVESTMENT  
MERRILL\_LYNCH  
COMPANIES  
SECURITIES

## Stock Market

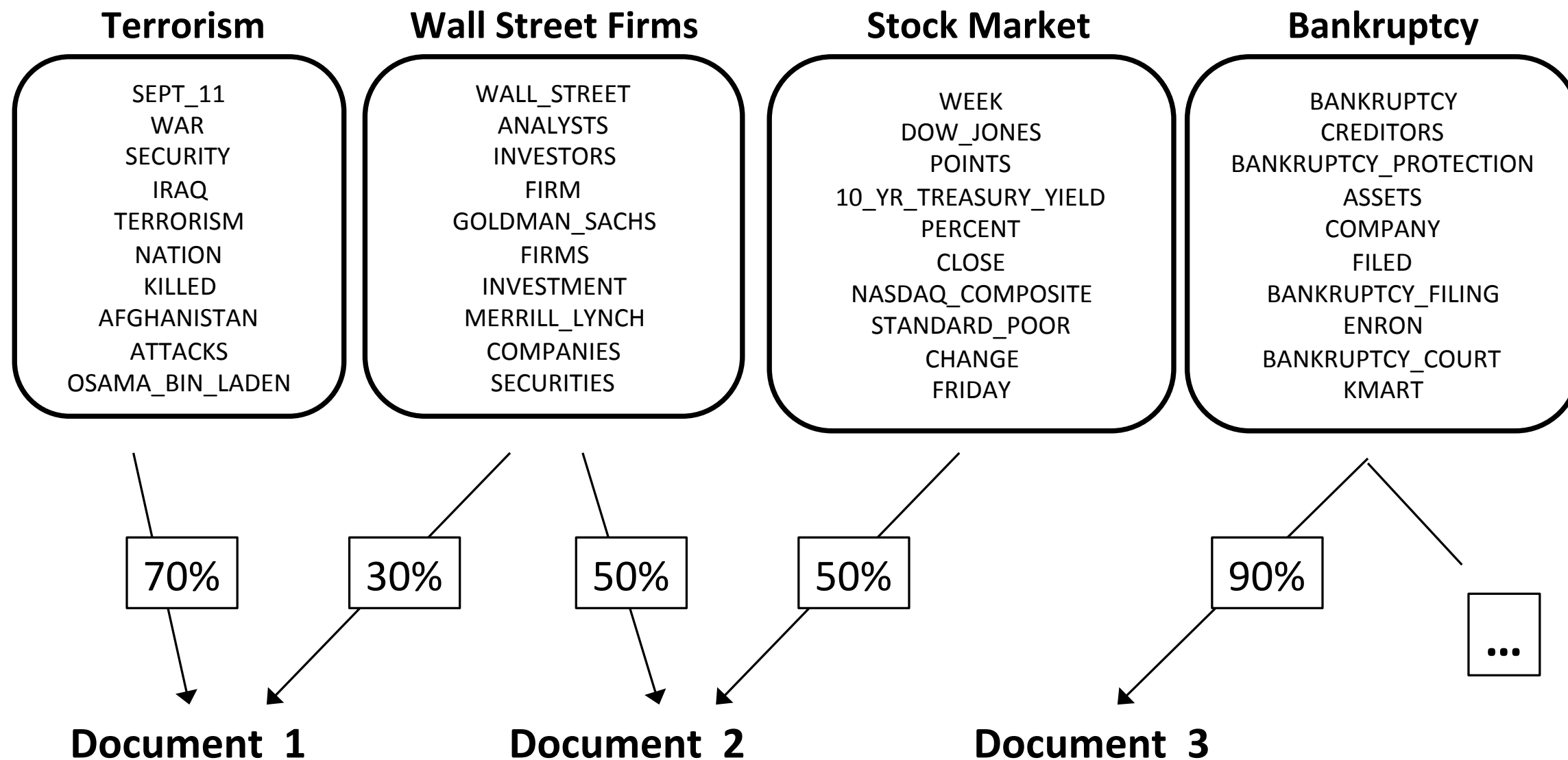
WEEK  
DOW\_JONES  
POINTS  
10\_YR\_TREASURY\_YIELD  
PERCENT  
CLOSE  
NASDAQ\_COMPOSITE  
STANDARD\_POOR  
CHANGE  
FRIDAY

## Bankruptcy

BANKRUPTCY  
CREDITORS  
BANKRUPTCY\_PROTECTION  
ASSETS  
COMPANY  
FILED  
BANKRUPTCY\_FILING  
ENRON  
BANKRUPTCY\_COURT  
KMART

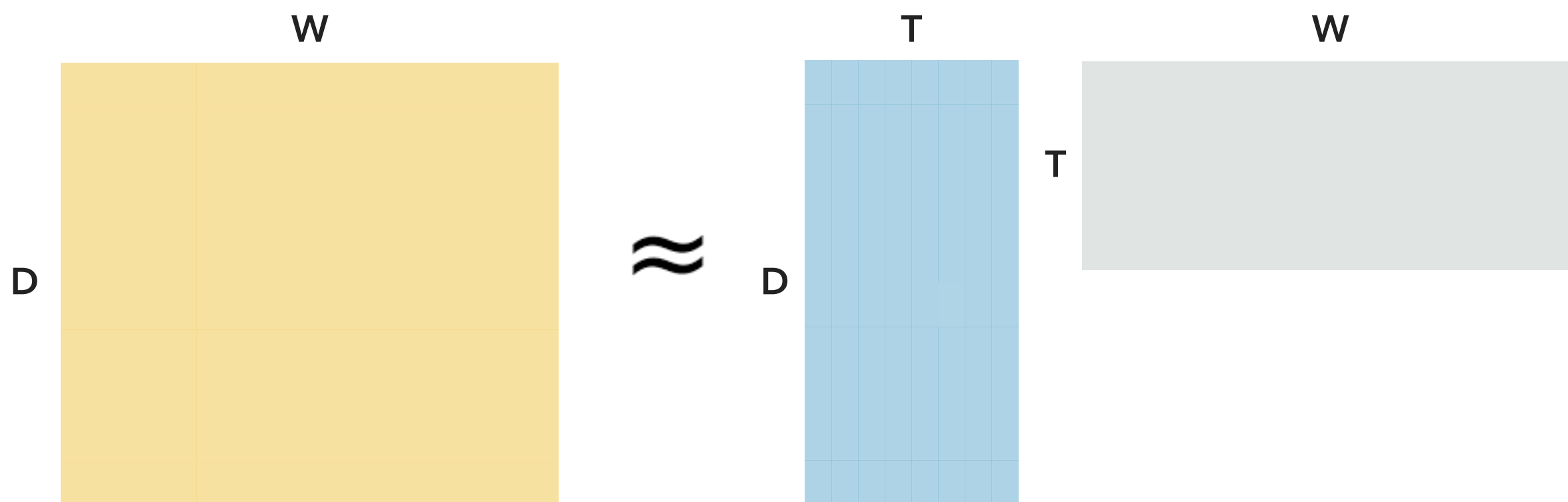
# LDA: Documents are Represented as Combinations of Topics

Figures from Mark Steyvers



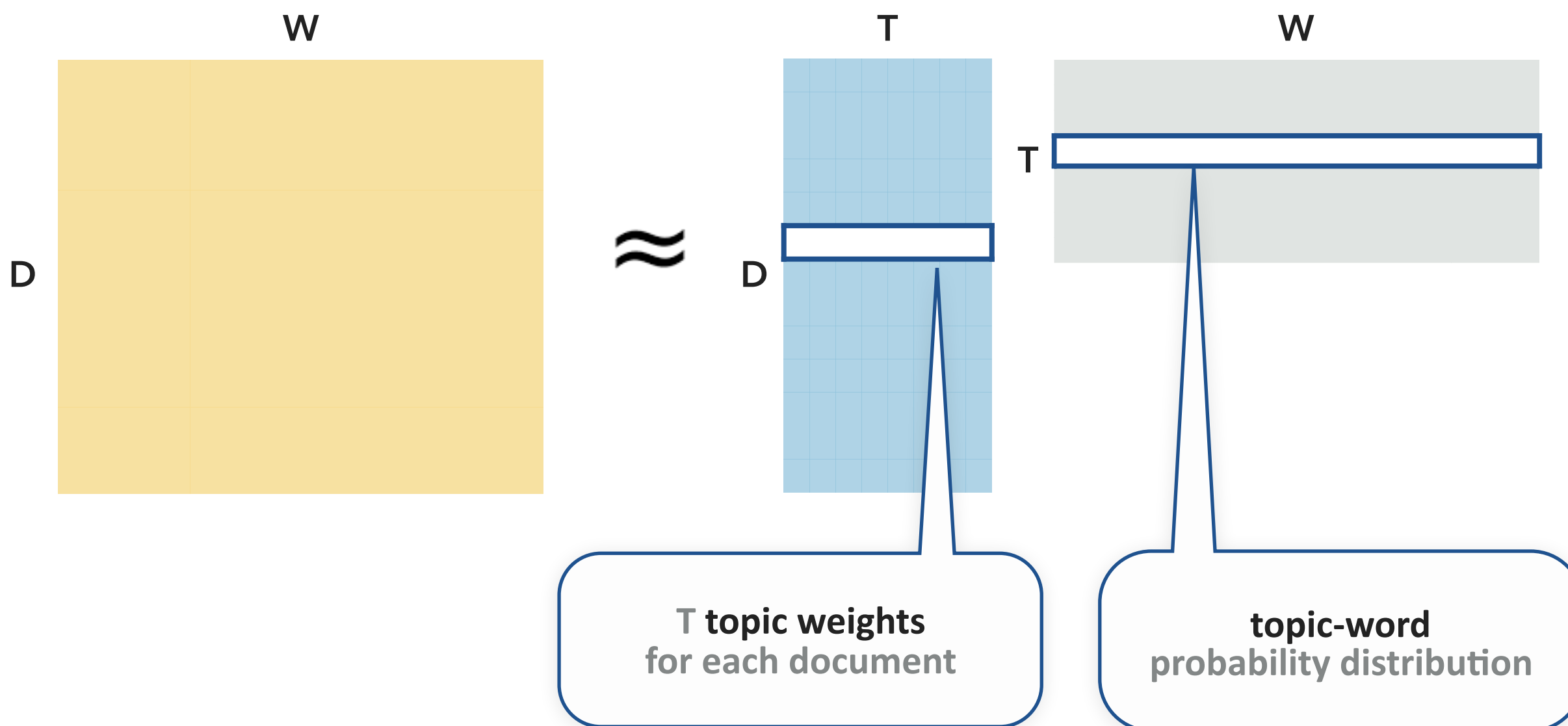
# LDA (Latent Dirichlet Allocation)

---



# LDA (Latent Dirichlet Allocation)

## ► Topic Models as Factor Analysis for Count Data





# NIH Data Representation for L-LDA

Documents	Words or Terms										
	brain	lung_cancer	women	obesity	children	mice	experiment	hbv	quality	glioma	researcher
doc1	3		1				1		1		1
doc2		1		1							
doc3					2						1
doc4			1				1		1		
doc5				1				1			
doc6									2		
doc7		1				1					
doc8	2									1	
doc9					1	1	1			2	
doc10			2								

# NIH Data Representation for L-LDA

Documents

Words or Terms

	brain	lung_cancer	women	obesity	children	mice	experiment	hbv	quality	glioma	researcher
doc1	3		1				1		1		1
doc2		1		1							
doc3					2						1
doc4			1				1		1		
doc5				1				1			
doc6									2		
doc7		1				1					
doc8	2									1	
doc9					1	1	1			2	
doc10			2								

Codes or Labels

	brain cancer	breast cancer	kidney disease	lung cancer	mind and body
doc1	1				
doc2			1		
doc3					1
doc4		1		1	
doc5			1		
doc6				1	
doc7			1		
doc8	1				
doc9	1				1
doc10		1			

# NIH Data Representation for L-LDA

Documents	Words or Terms											Codes or Labels						
	brain	lung_cancer	women	obesity	children	mice	experiment	hbv	quality	glioma	researcher	brain cancer	breast cancer	kidney disease	lung cancer	mind and body	Background 1	Background 2
doc1	3		1				1		1		1	1					1	1
doc2		1		1										1			1	1
doc3					2						1				1	1	1	
doc4			1				1		1			1			1		1	1
doc5				1				1						1			1	1
doc6									2						1		1	1
doc7		1						1						1			1	1
doc8	2										1						1	1
doc9					1	1	1					1			1	1	1	
doc10			2														1	1

# Examples of Topics from NCI Abstracts (5 out of 98)

## Brain Cancer

glioma  
brain tumor  
gbm  
malignant glioma  
glioblastoma  
brain

## Breast Cancer

breast cancer  
women  
breast cancer cell  
breast  
breast cancer patient  
brca1

## Kidney Disease

rcc  
kidney cancer  
renal cell carcinoma  
vhl  
renal cancer  
pvhl

## Background 1

program  
trainee  
university  
training  
candidate  
field

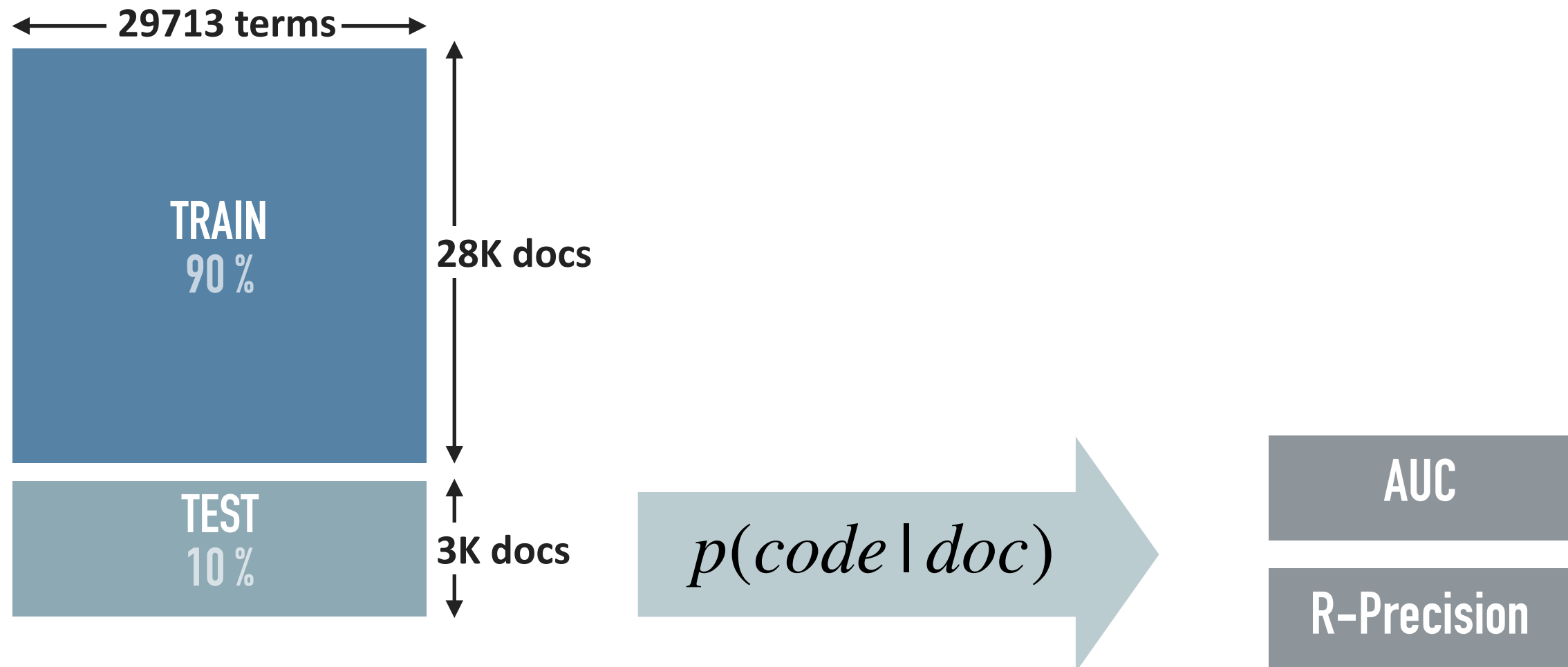
## Background 7

model  
mice  
work  
experiment  
human  
mouse model

**88 Topics from RCDC labels**  
**10 Background topics**

# Evaluation

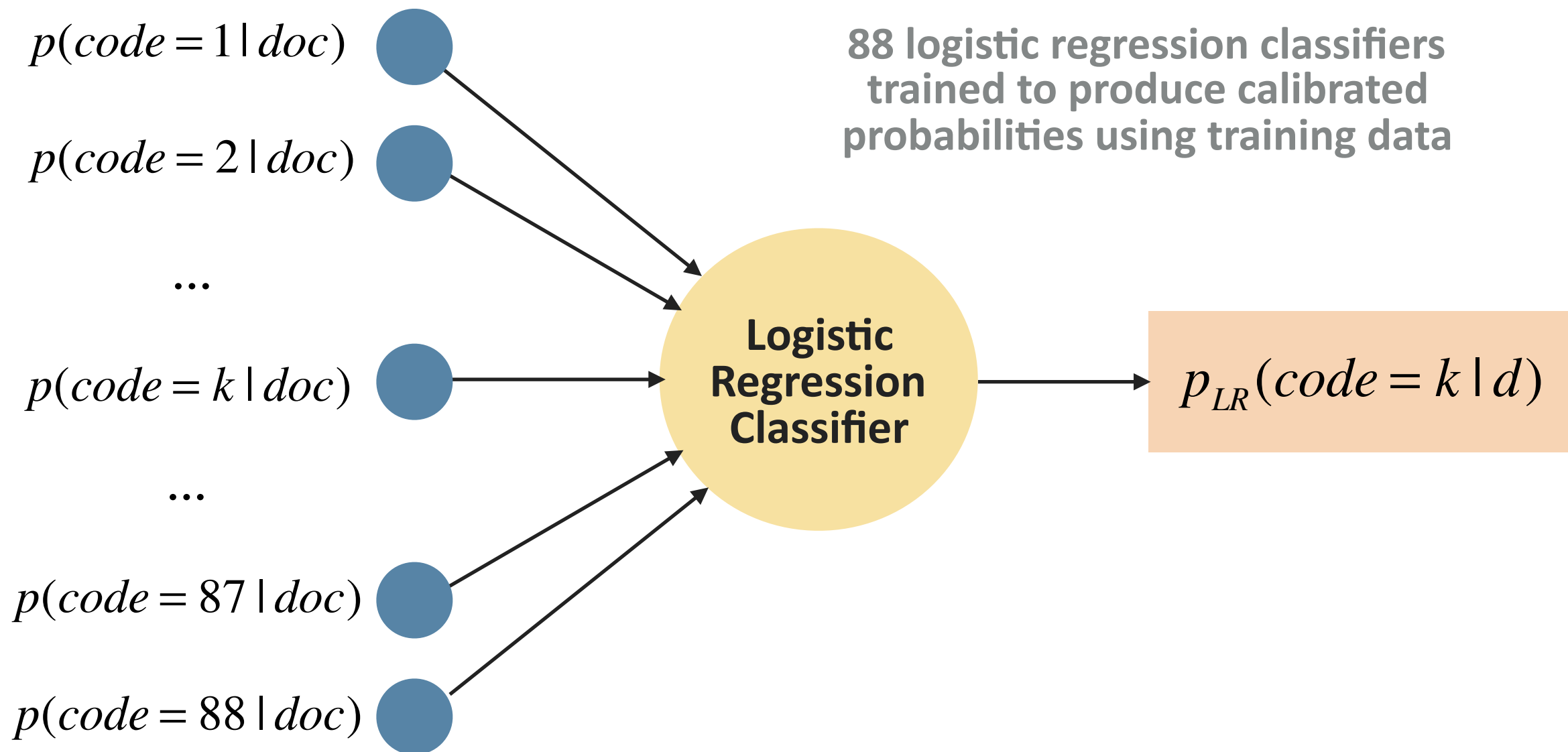
## Grants with RCDC labels (31,628 documents)



$$p(\text{code} \mid \text{doc}) \propto \sum_i p(\text{code} \mid \text{word}_i, \text{doc})$$

Sampling probabilities were averaged over the words in a document to calculate **AUC** and **R-precision** scores

# Logistic Regression Classifier



L-LDA TOPIC PROBABILITY

CALIBRATED TOPIC PROBABILITY

# Evaluation Result

---

	$p(c   d)$	$p_{LR}(c   d)$
	L-LDA	L-LDA + Logistic Regression
AUC	0.80	0.89
R-Precision	0.56	0.64

# Analyzing Funding Patterns over Time

- ▶ Fractionally assign the funds in direct proportion to the probabilities from the logistic regression classifiers  $p_{LR}(code | doc)$

$$w_{cd} = \frac{p_{LR}(c | d)}{\sum_k p_k(c = k | d)} \quad c = 1, 2, \dots, 88$$

$$F_c^y = \sum_{d: y_d=y} w_{cd} x_d$$

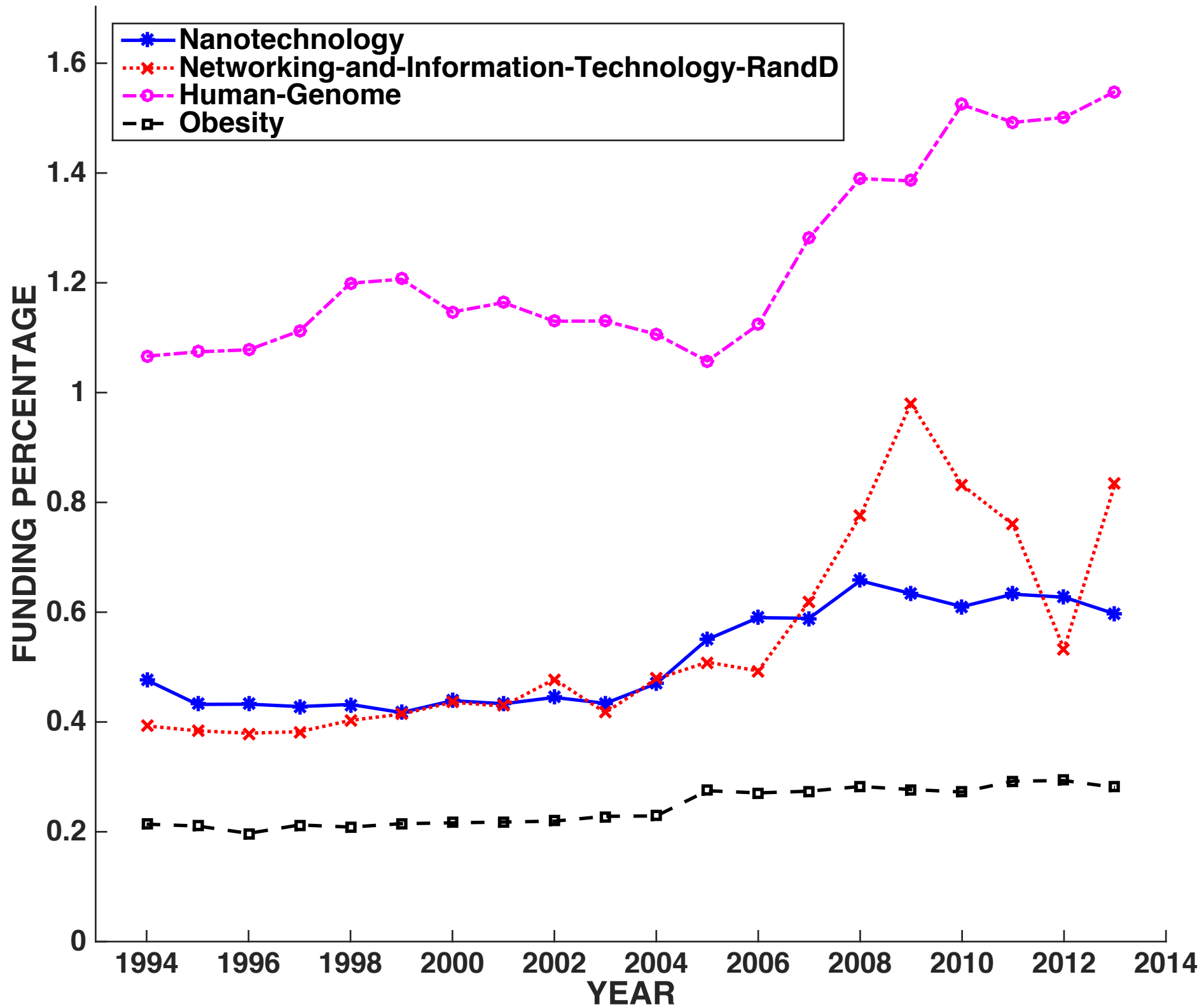
$w_{cd}$  : weight for the category  $c$  for document  $d$

$x_d$  : amount of funding for document  $d$  (considered inflation)

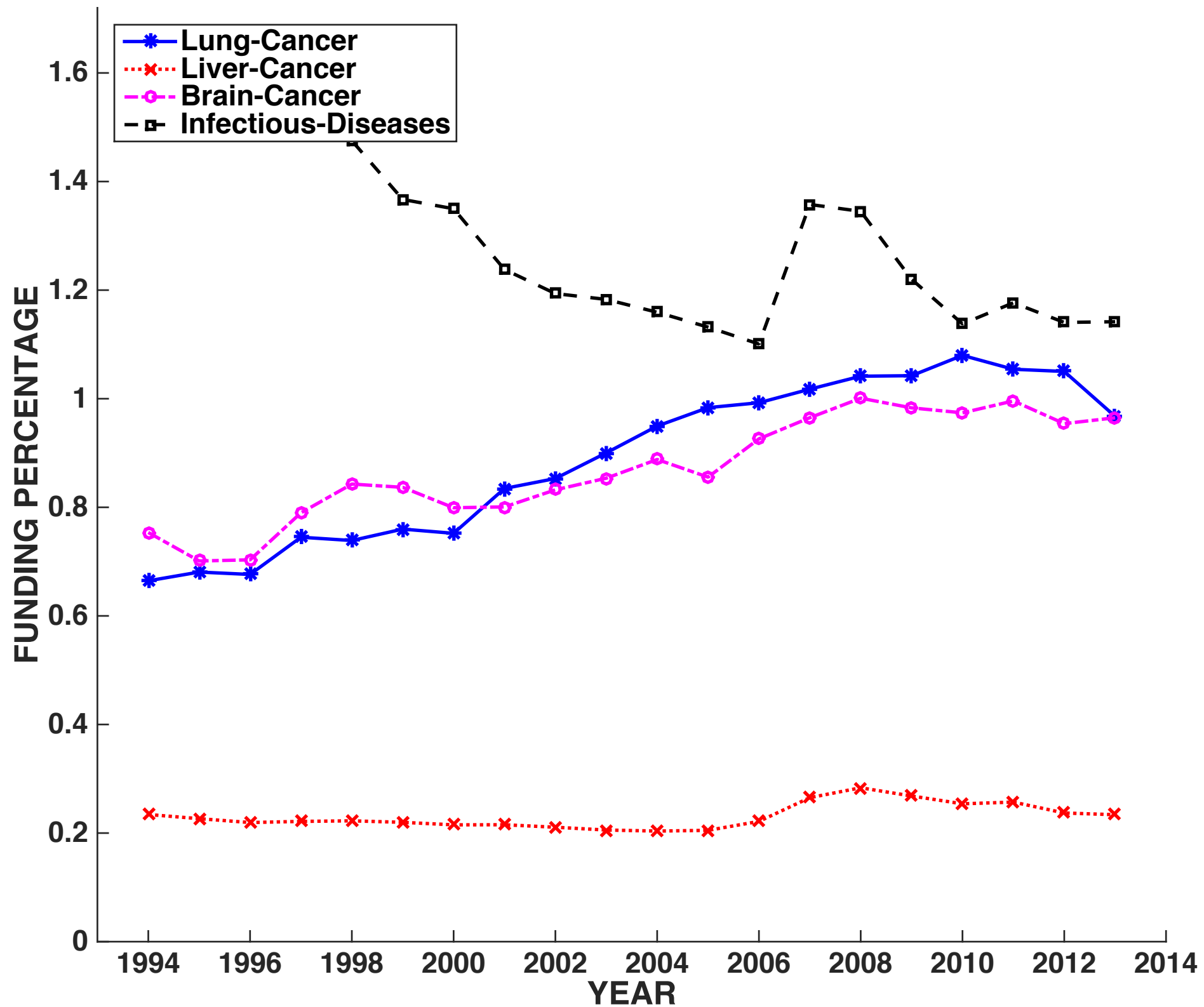
$y_d$  : year when document  $d$  was funded

$F_c^y$  : total estimated amount of funding for category  $c$  in year  $y$

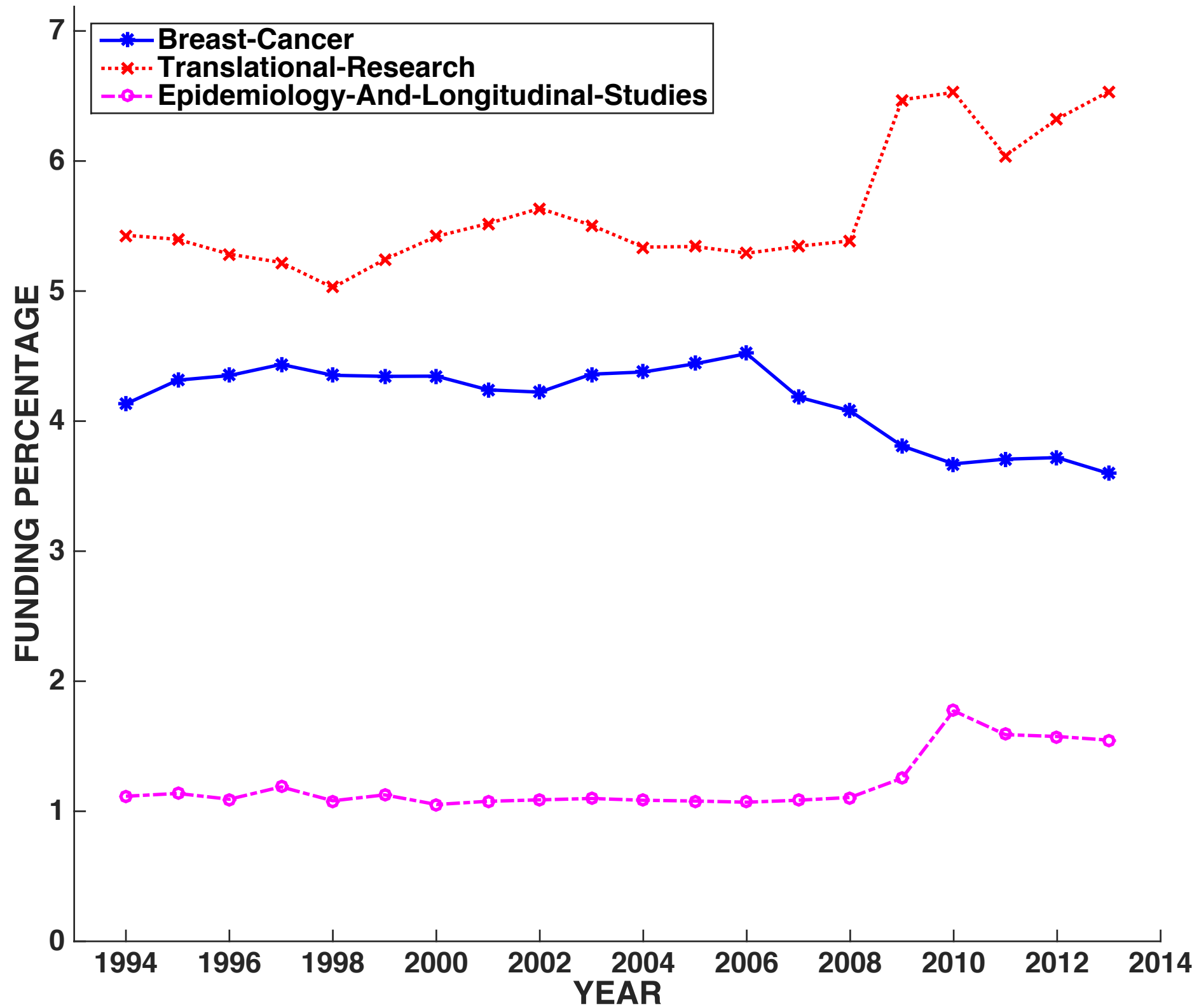


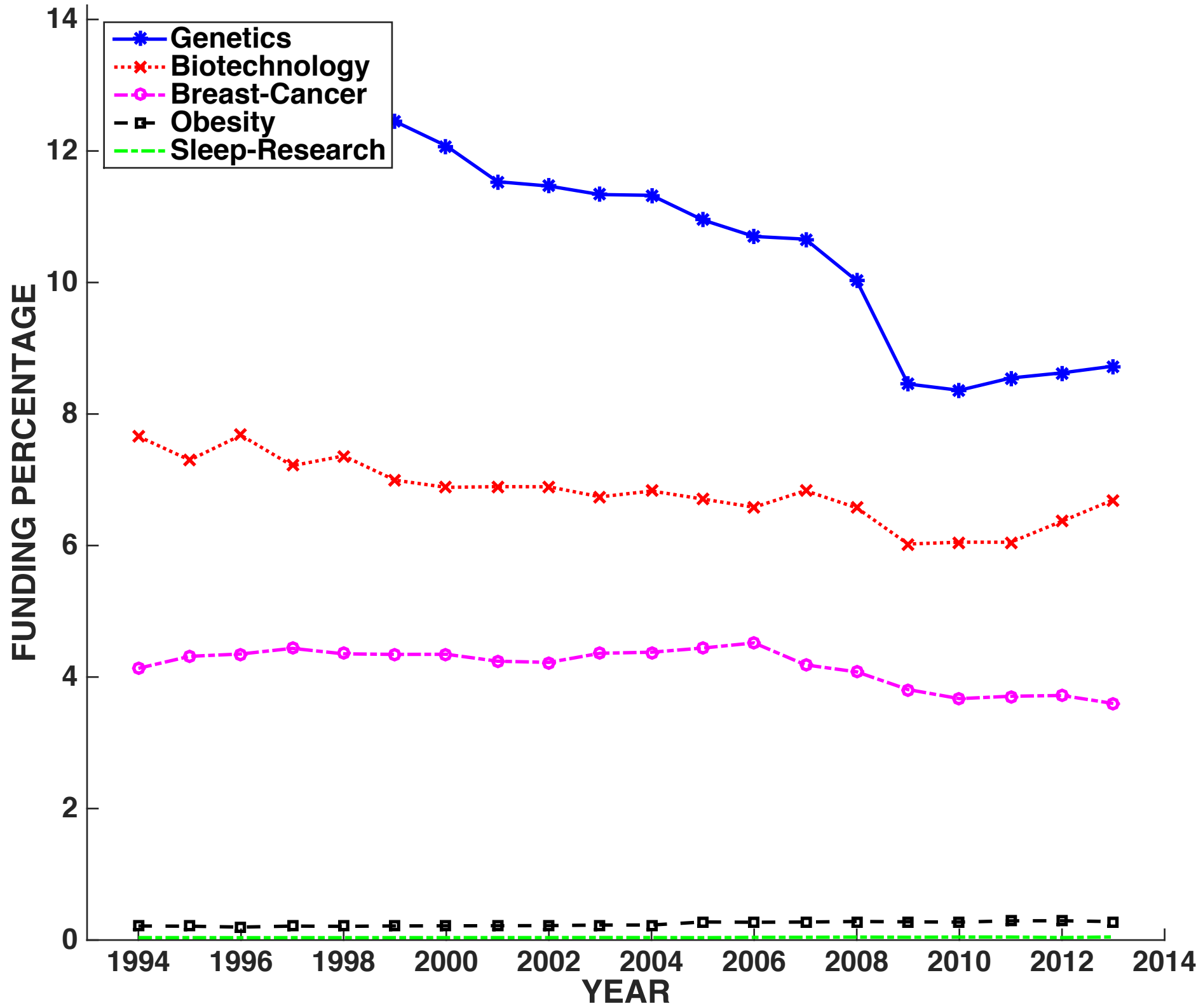


Estimated percentage of funding allocated to 4 general RCDC categories



Estimated percentage of funding allocated to 4 specific RCDC disease categories





# Conclusions

---

## ▶ Summary

- ▶ Labeled topic modeling and logistic classifiers can be combined to analyze NIH grant funding data
- ▶ Statistical topic modeling allows linking of text with metadata in a quantifiable manner

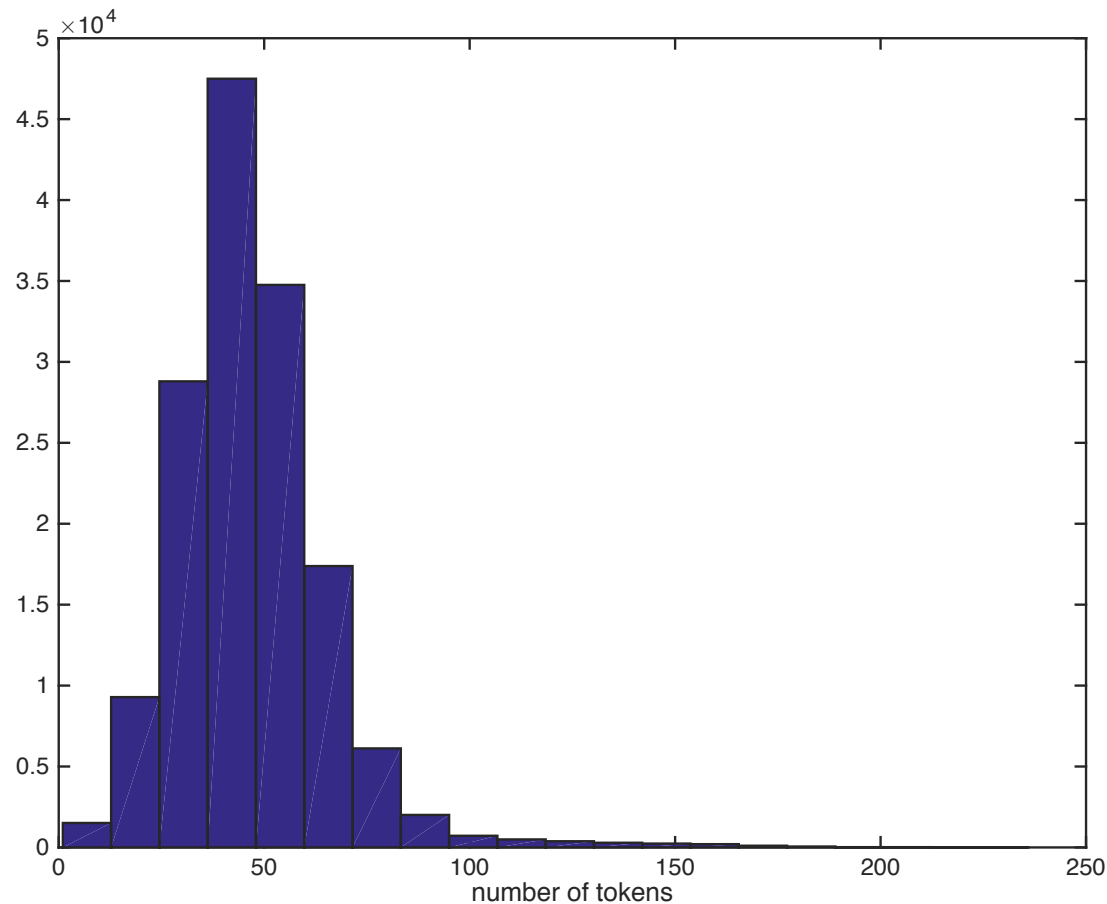
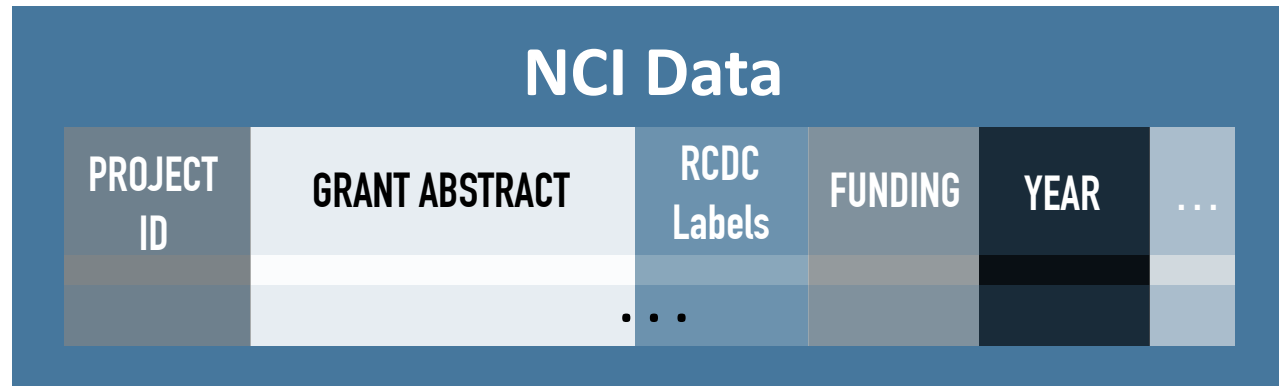
## ▶ Future Work

- ▶ Jointly analyze grants and scientific articles related to the grants (ongoing)
- ▶ Broader analysis of the economic and policy implications
- ▶ Improvements on topic model
  - ▶ How to best calibrate
  - ▶ Selecting the right Hyper parameters
  - ▶ Methods such as using seed words

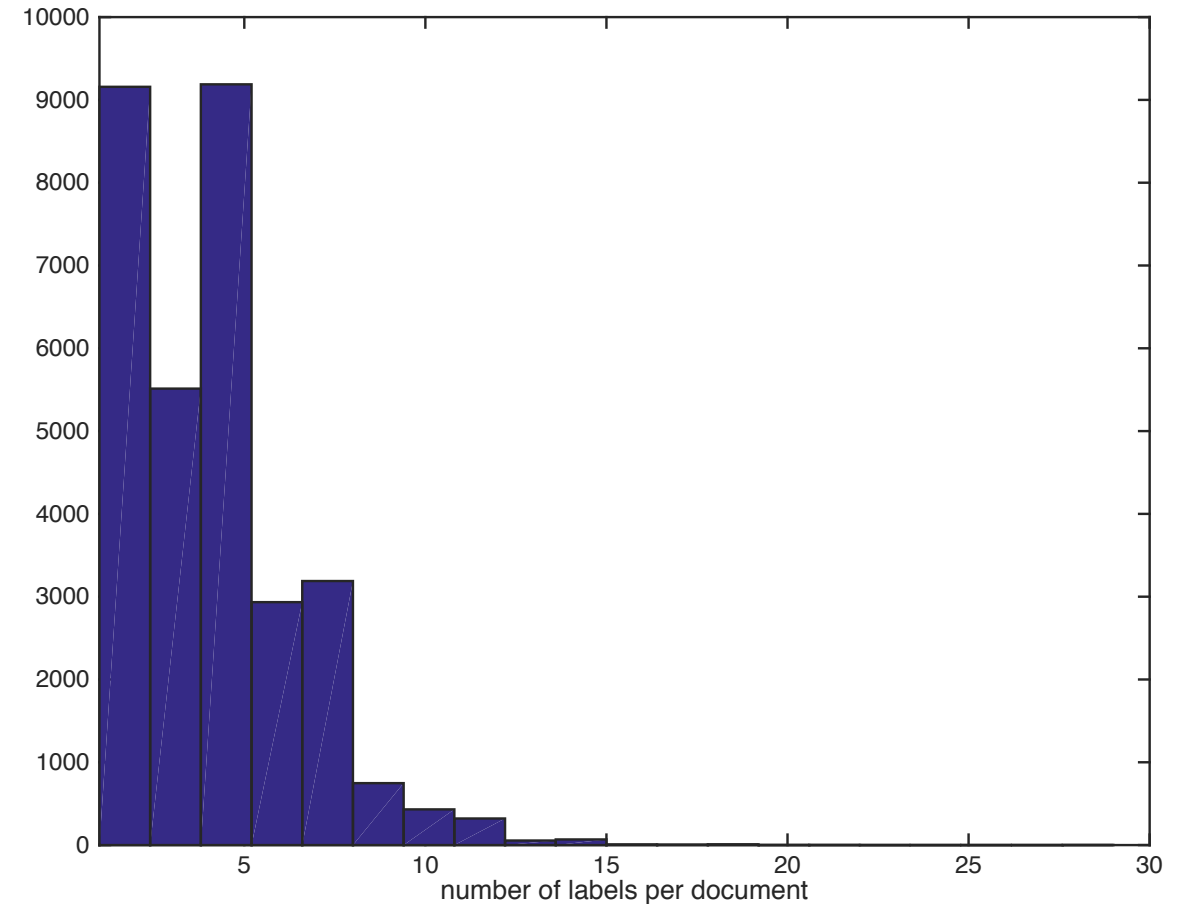
**THANK YOU**

# BACKUP SLIDES

# NCI (National Cancer Institute) Data



Number of tokens in a document



Number of labels in a document



# Examples of Topics from NCI Abstracts (5 out of 88)

## Brain Cancer

glioma  
brain tumor  
gbm  
malignant glioma  
glioblastoma  
brain

## Breast Cancer

breast cancer  
women  
breast cancer cell  
breast  
breast cancer patient  
brca1

## Kidney Disease

rcc  
kidney cancer  
renal cell carcinoma  
vhl  
renal cancer  
pvhl

## Hepatitis

hcv  
hbv  
liver cancer  
hepatitis virus  
hbv infection  
hbv replication

## Lung Cancer

gliomalung cancer  
nslc  
lung  
leading cause  
cancer death  
egfr

# Labeled-LDA for NIC Grants

---

- ▶ 88 Topics (RCDC Codes)
- ▶ 10 Background Topics
  
- ▶ Hyper parameters
- ▶ Dirichlet prior for word-topic distribution
  - ▶  $\beta_w = 0.01$
- ▶ Dirichlet prior for doc-topic distribution
  - ▶ Used proportional alphas

$$\sum_{c=1}^{88} \alpha_c = 5 \quad \sum_b^B \alpha_b = 1$$

## Improving Gibbs Sampling Predictions on Unseen Data for Latent Dirichlet Allocation

**Yannis Papanikolaou**

*Department of Informatics*

*Aristotle University of Thessaloniki*

*Thessaloniki, Greece*

YPAPANIK@CSD.AUTH.GR

**Timothy N. Rubin**

*Cognitive Computing Laboratory*

*Indiana University*

*Bloomington, IN, USA*

TIMRUBIN@INDIANA.EDU

**Grigorios Tsoumakas**

*Department of Informatics*

*Aristotle University of Thessaloniki*

*Thessaloniki, Greece*

GREG@CSD.AUTH.GR

### 3.2 Averaging the conditional probabilities

Another approach we consider is averaging over the conditional probabilities from Equation 1 for every document. As stated before,  $p$  expresses the probability of assigning to a word token a given topic, conditioned on the word, the document, and the other assignments of words to topics in both the given document (through  $n_{dl}$ ) and in the whole corpus (through  $n_{wl}$  during training and  $\phi$  during prediction). The idea is to average these word-level probabilities for every document and obtain a total probability  $p(l|d)$  such that:

$$\theta_p(d, l) = p(l|d, \cdot) = \sum_{w=1}^{N_d} p(z_{w_i} = l | w_i, d) = \sum_{w=1}^{N_d} \phi_l(w) \times \frac{n_{d,l} + \alpha}{\sum_{d=1}^D (n_{d,l} + \alpha)} \quad (6)$$

The motivation for this approach is that the conditional probabilities should provide a richer representation of the document-topics distribution. We note here, that as the probability  $p(z_{w_i} = l | w_i, d)$  is unnormalized, an additional step of normalization is required so that  $\sum_{l=1}^L p(z_{w_i} = l | w_i, d) = 1$ . We will refer to this approach as CGS<sub>p</sub>.

# Analyzing Funding Patterns over Time

---

- ▶ Fractionally assign the funds in direct proportion to the probability

$$w_{cd} = \frac{p_l(c | d)}{\sum_k p_k(c = k | d)} \quad c = 1, 2, \dots, 88$$

$$F_c^y = \sum_{d: y_d = y} w_{cd} x_d$$

# NCI (National Cancer Institute) Data

- ▶ **149,901 grants in total**
- ▶ **for FY1994 ~ FY2013**
- ▶ **Number of grants with labels : 31,628 (2008~2011)**
- ▶ **Number of grants without labels : 118,273**
- ▶ **Size of vocabulary (W) : 29,713**

		W						
D		3	1		1	1	1	
			1	1				
					2		1	
			1			1	1	
				1			1	
							2	
		1			1			
		2					1	
					1	1	1	2
		2						