

# Filtrado Basado en Contenido

IIC 3633 - Sistemas Recomendadores

Denis Parra  
Profesor Asistente, DCC, PUC Chile

# TOC

## En esta clase

1. Contenido en lugar de ratings
2. Representación de Espacio Vectorial
3. TF-IDF
4. Buscando Items Similares
5. Representación en Espacio Latente

# Por Qué un Recomendador Basado en Contenido

- El filtrado colaborativo tiene algunas desventajas: cold-start, sparsity, transparency.

## PROS

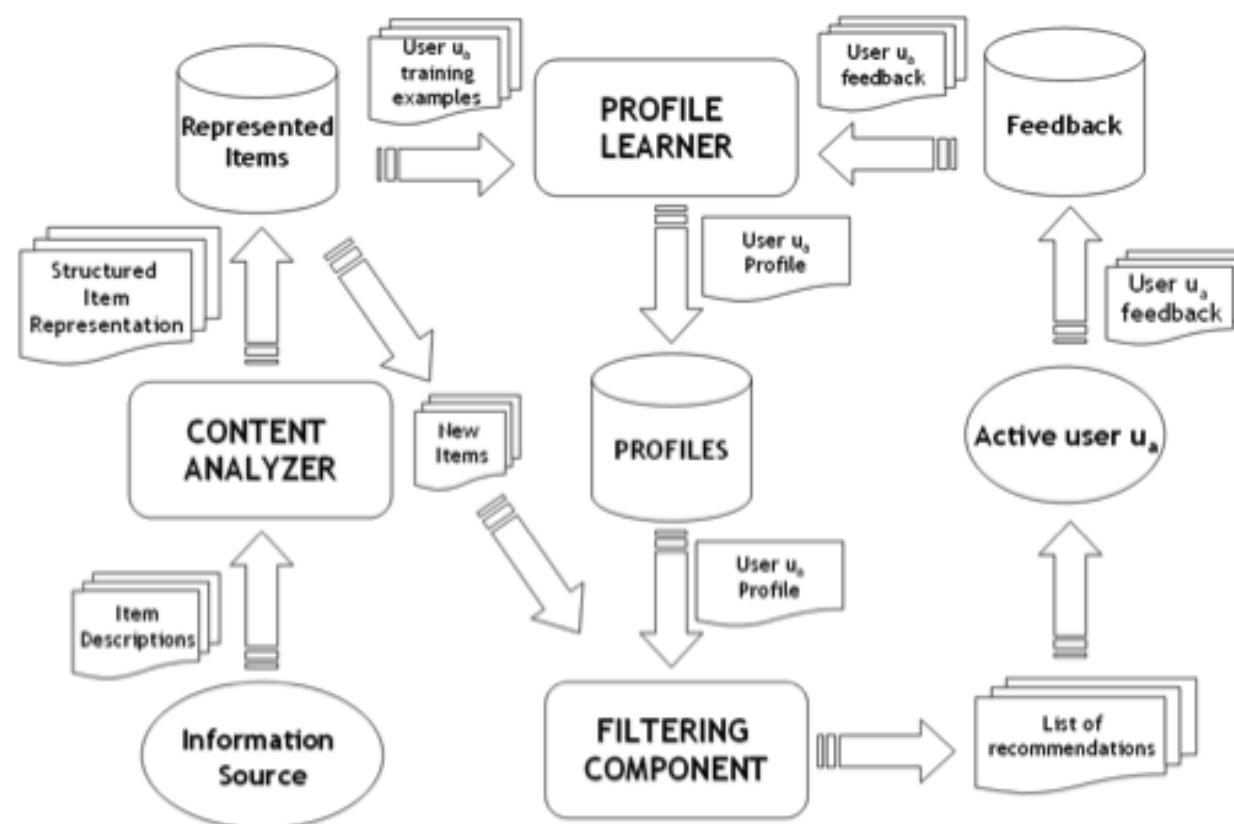
- A diferencia del Filtrado Colaborativo, si los items tienen descripciones suficientes, nos evitamos el "new-item problem"
- Las representaciones del contenido son variadas y permiten utilizar diversas técnicas de procesamiento del texto, uso de información semántica, inferencias, etc.
- Es sencillo hacer un sistema más transparente: usamos el mismo contenido para explicar las recomendaciones.

## CONS

- Tienden a la sobre-especialización: va a recomendar items similares a los ya consumidos, creando una tendencia al "filter bubble".
- Los métodos basados en filtrado colaborativo han mostrado ser, empíricamente, más precisos al momento de generar recomendaciones.

# Arquitectura de un Sistema de Recomendación CB

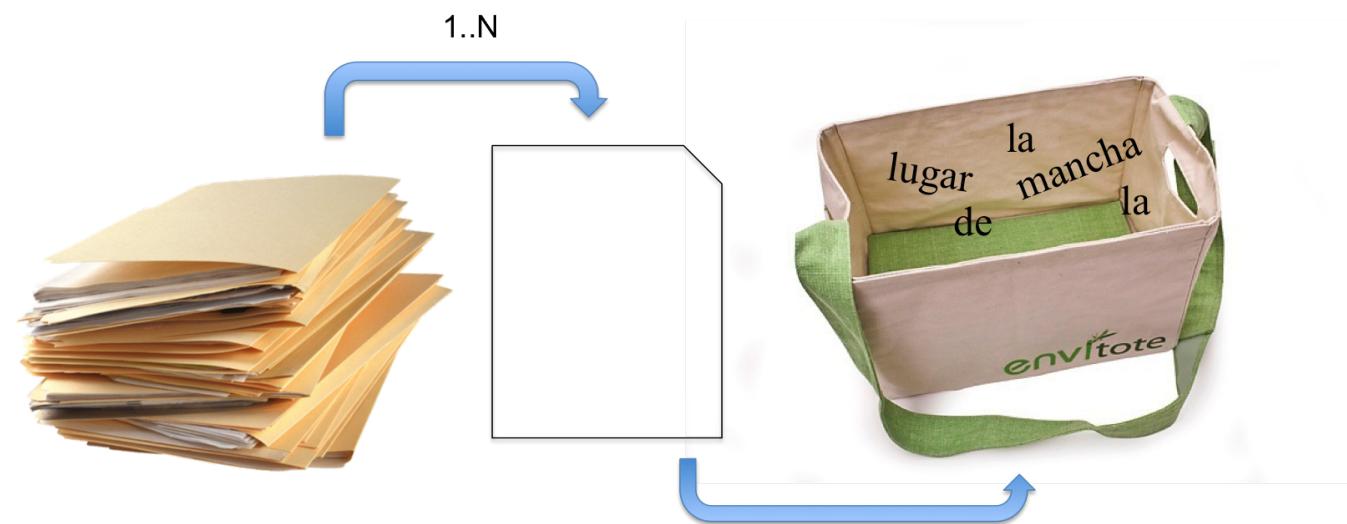
- Los componentes principales son: (1) Analizador del Contenido, (2) Aprendizaje del Perfil de Usuario, (3) Filtrado de Contenido



**Fig. 3.1:** High level architecture of a Content-based Recommender

# Representación del Contenido: Bolsa de Palabras

- Se suele representar a los documentos como "bolsas de palabras"; de esta forma es fácil pasar a representar cada documento como un vector (Vector Space Model)



# Representación del Contenido: VSM

- El corpus completo puede entonces representarse como una matriz donde las filas son términos y las columnas son documentos.

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
ANTHONY	1	1	0	0	0	1	
BRUTUS	1	1	0	1	0	0	
CAESAR	1	1	0	1	1	1	
CALPURNIA	0	1	0	0	0	0	
CLEOPATRA	1	0	0	0	0	0	
MERCY	1	0	1	1	1	1	
WORSER	1	0	1	1	1	0	
...							

- Luego, ¿Cuál es la mejor forma de representar los pesos de los términos?

# Representación del Contenido: VSM II

## Frecuencia de los términos

Cada documento se representa como un vector, el "peso" de cada palabra para ese documento

$$\text{TF}(t_k, d_j) = \frac{f_{k,j}}{\max_z f_{z,j}}$$

puede darse en base a la frecuencia del término en el documento.

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
ANTHONY	157	73	0	0	0	1	
BRUTUS	4	157	0	2	0	0	
CAESAR	232	227	0	2	1	0	
CALPURNIA	0	10	0	0	0	0	
CLEOPATRA	57	0	0	0	0	0	
MERCY	2	0	3	8	5	8	
WORSER	2	0	1	1	1	5	
...							

Podemos normalizar el valor en función de la frecuencia máxima de cualquier término en el documento.

# Representación del Contenido: VSM III

## Log de Frecuencia de los términos

Pero el hecho que un término  $x$  aparece 100 veces y otro término  $y$  sólo 10 veces, no hace a  $x$  10 veces más relevantes; por lo tanto podemos usar un logaritmo.

- La log-frecuencia del término  $t$  en  $d$  se define como

$$w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d} & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- $tf_{t,d} \rightarrow w_{t,d}$ :

$$0 \rightarrow 0, \quad 1 \rightarrow 1, \quad 2 \rightarrow 1.3, \quad 10 \rightarrow 2, \quad 1000 \rightarrow 4, \quad \text{etc.}$$

# Representación del Contenido: VSM IV

## TF-IDF

Bajo la intuición de que un término que aparece en sólo unos pocos documentos podría ser descriptivo, podemos considerar la "Inverse Document Frequency" y combinarla con la "Term Frequency":

$$\text{TF-IDF}(t_k, d_j) = \underbrace{\text{TF}(t_k, d_j)}_{\text{TF}} \cdot \underbrace{\log \frac{N}{n_k}}_{\text{IDF}}$$

Donde  $t_k$  es el término  $k$ ,  $d_j$  es el documento  $j$ .

# Resumen de Componentes del TF-IDF

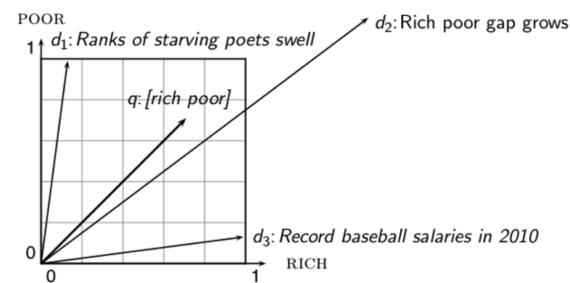
Term frequency		Document frequency		Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf)	$\max\{0, \log \frac{N-df_t}{df_t}\}$	u (pivoted unique)	$1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/CharLength^\alpha$ , $\alpha < 1$
L (log ave)	$\frac{1+\log(tf_{t,d})}{1+\log(\text{ave}_{t \in d}(tf_{t,d}))}$				

# Representación Semántica de Contenido

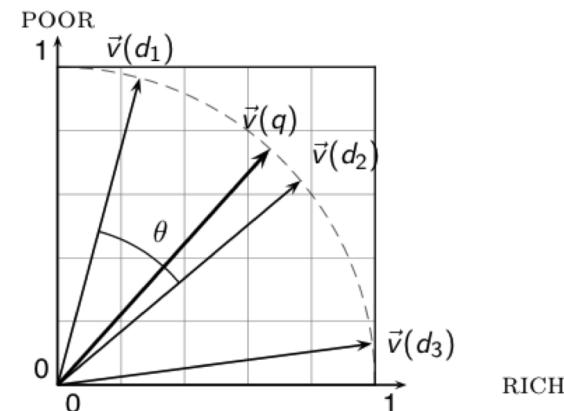
- No todo el contenido del documento corresponde a la misma categoría.
- Autor, palabras clave, fechas, tópicos pueden dar una noción adicional de filtrado.
- Opción 1: Representación semántica explícita (No lo veremos en detalle en esta clase)
  - Ontologías
  - WordNet
  - ConceptNet
- Opción 2: Inferir representación semántica (LSI, LDA)
- Opción 3: Word Vectors (Word2Vec, Glove)

# Buscando Items Similares

## Distancia Euclidiana

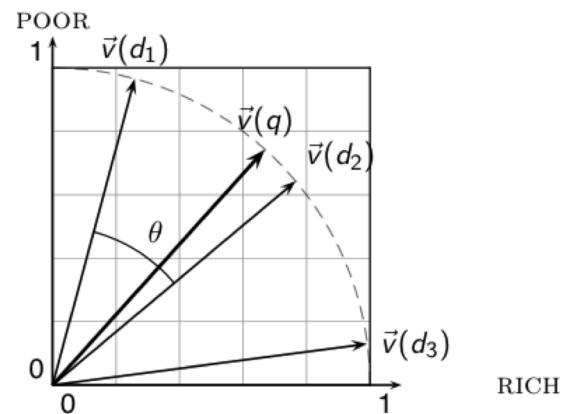


## Distancia Coseno



# Buscando Items Similares

Distancia Coseno



Fórmula

$$\text{sim}(d_i, d_j) = \frac{\sum_k w_{ki} \cdot w_{kj}}{\sqrt{\sum_k w_{ki}^2} \cdot \sqrt{\sum_k w_{kj}^2}}$$

# Buscando Items Similares II

## Okapi BM25

$$RSV_d = \sum_{t \in q} IDF \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d / L_{ave})) + tf_{td}} \cdot \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}}$$

Ref: Denis Parra and Peter Brusilovsky. 2009. Collaborative filtering for social tagging systems: an experiment with CiteULike. In Proceedings of the third ACM conference on Recommender systems (RecSys '09) <http://doi.acm.org/10.1145/1639714.1639757>

# Buscando Items Similares III

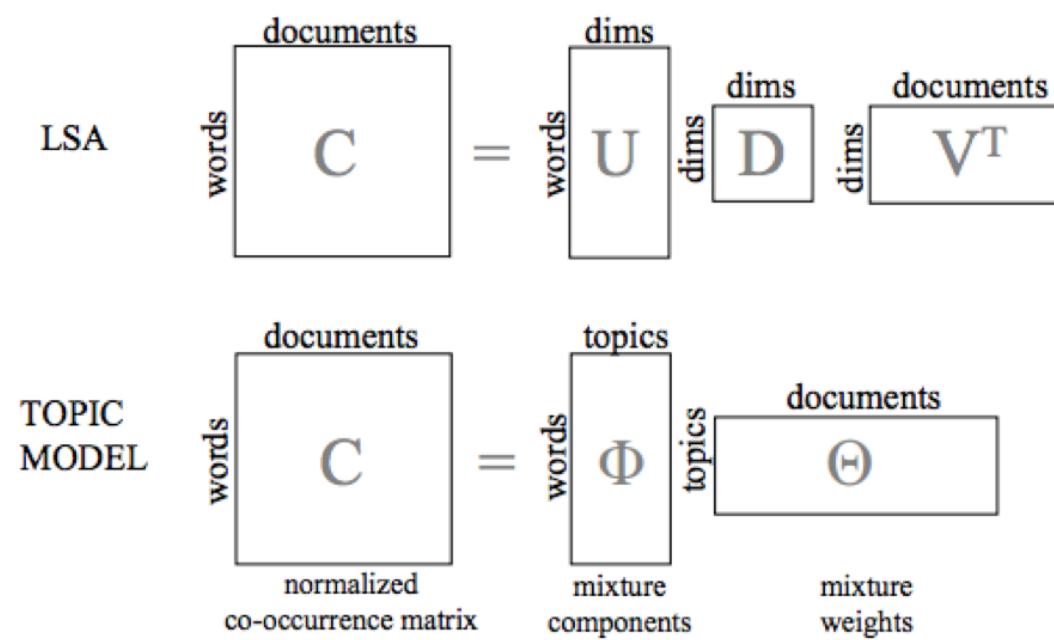
## Técnicas de Procesamiento adicionales

- Pasar a mayúsculas/minúsculas
- Tokenization
- Stemming (Porter, Krovetz)
- Lemmatization

# Buscando Items Similares

Representación en espacio latente

- Latent Semantic Indexing
- Latent Dirichlet Allocation



LSI |

$$\mathbf{X} = \mathbf{U} \times \Sigma \times \mathbf{V}^T$$

$\overset{N \times d}{=}$ 
  $\times$ 
  $\times$ 
  $\overset{r \times d}{}$

---

$X$	$U$	$\Sigma$	$V^T$
$(\mathbf{d}_j)$			$(\hat{\mathbf{d}}_j)$
$\downarrow$			$\downarrow$
$(\mathbf{t}_i^T) \rightarrow$	$\begin{bmatrix} x_{1,1} & \dots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{m,1} & \dots & x_{m,n} \end{bmatrix}$	$= (\hat{\mathbf{t}}_i^T) \rightarrow$	$\left[ \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_l \end{bmatrix} \dots \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_l \end{bmatrix} \right] \cdot \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_l \end{bmatrix} \cdot \begin{bmatrix} [\mathbf{v}_1 \\ \vdots \\ \mathbf{v}_l] \end{bmatrix}$

## LSI II

		$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	
	<b>ship</b>	1	0	1	0	0	0	
	<b>boat</b>	0	1	0	0	0	0	
	<b>ocean</b>	1	1	0	0	0	0	
	<b>voyage</b>	1	0	0	1	1	0	
	<b>trip</b>	0	0	0	1	0	1	

## LSI III

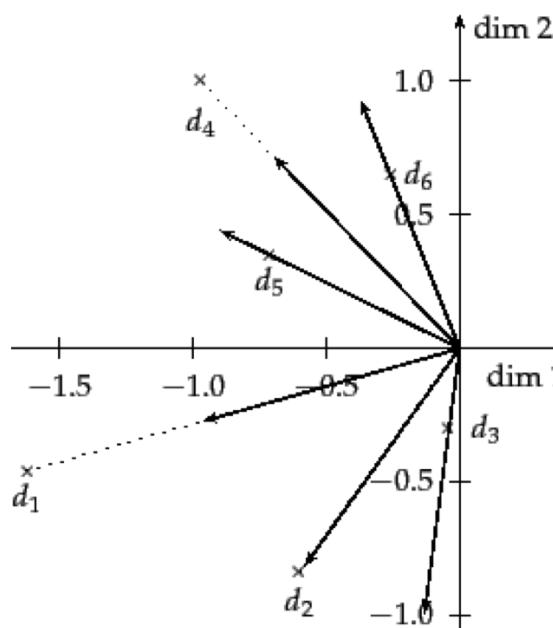
	1	2	3	4	5
ship	-0.44	-0.30	0.57	0.58	0.25
boat	-0.13	-0.33	-0.59	0.00	0.73
ocean	-0.48	-0.51	-0.37	0.00	-0.61
voyage	-0.70	0.35	0.15	-0.58	0.16
trip	-0.26	0.65	-0.41	0.58	-0.09

2.16	0.00	0.00	0.00	0.00
0.00	1.59	0.00	0.00	0.00
0.00	0.00	1.28	0.00	0.00
0.00	0.00	0.00	1.00	0.00
0.00	0.00	0.00	0.00	0.39

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
1	-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2	-0.29	-0.53	-0.19	0.63	0.22	0.41
3	0.28	-0.75	0.45	-0.20	0.12	-0.33
4	0.00	0.00	0.58	0.00	-0.58	0.58
5	-0.53	0.29	0.63	0.19	0.41	-0.22

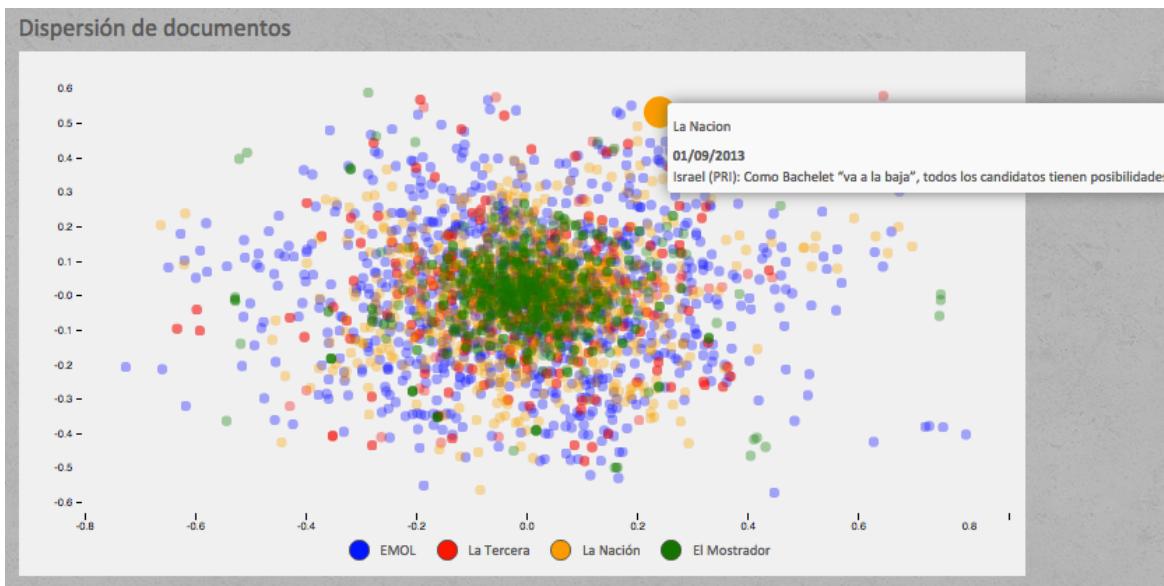
## LSI IV

		$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
1		-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2		-0.29	-0.53	-0.19	0.63	0.22	0.41
3		0.28	-0.75	0.45	-0.20	0.12	-0.33
4		0.00	0.00	0.58	0.00	-0.58	0.58
5		-0.53	0.29	0.63	0.19	0.41	-0.22



## LSI IV

Demo: <http://dfaο-uc.github.io/>



## Proyección de documentos o términos nuevos

- Folding in: Using Linear Algebra for Intelligent Information Retrieval

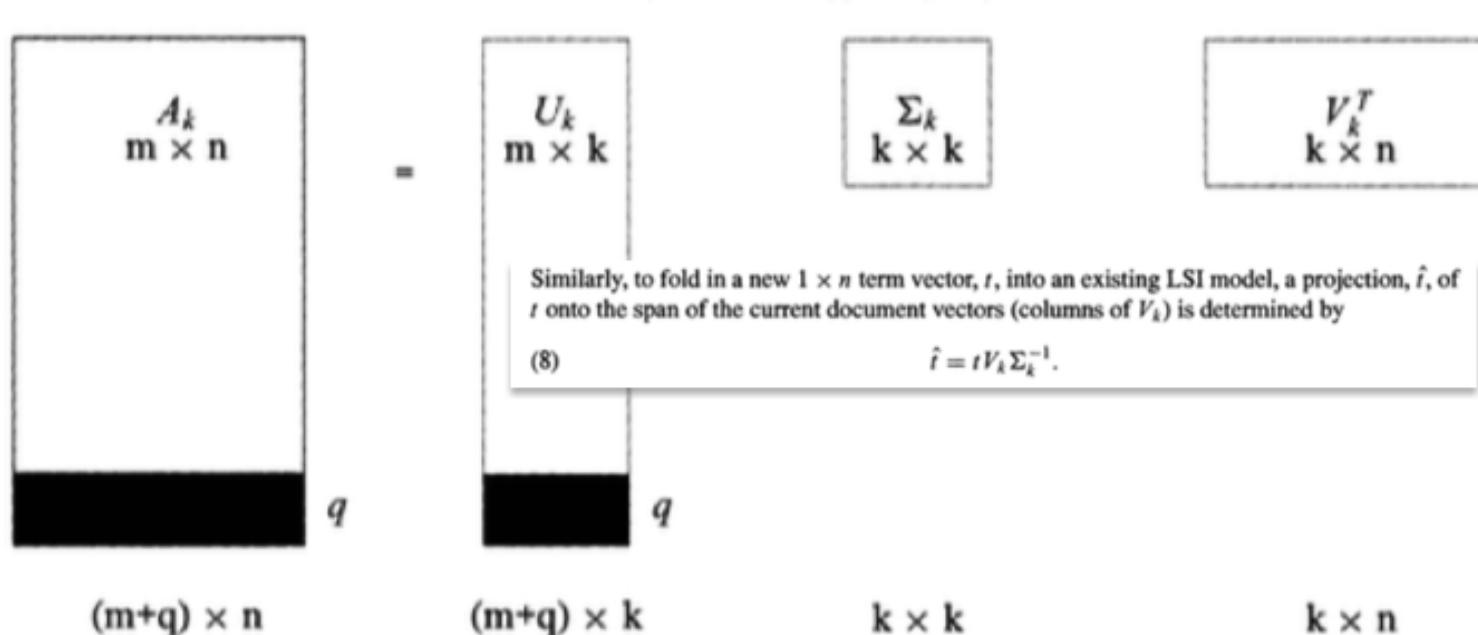
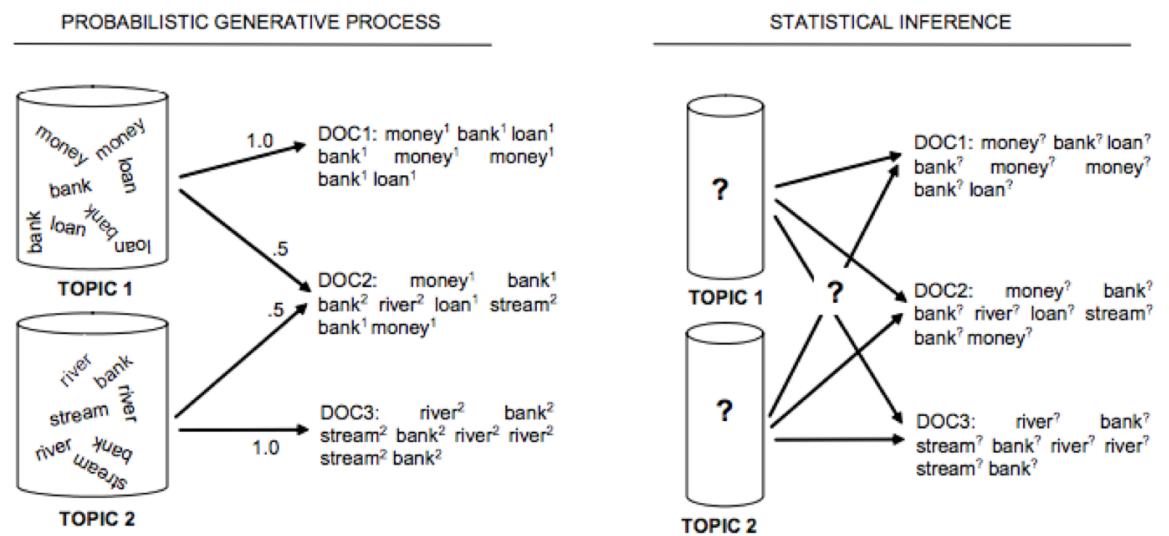


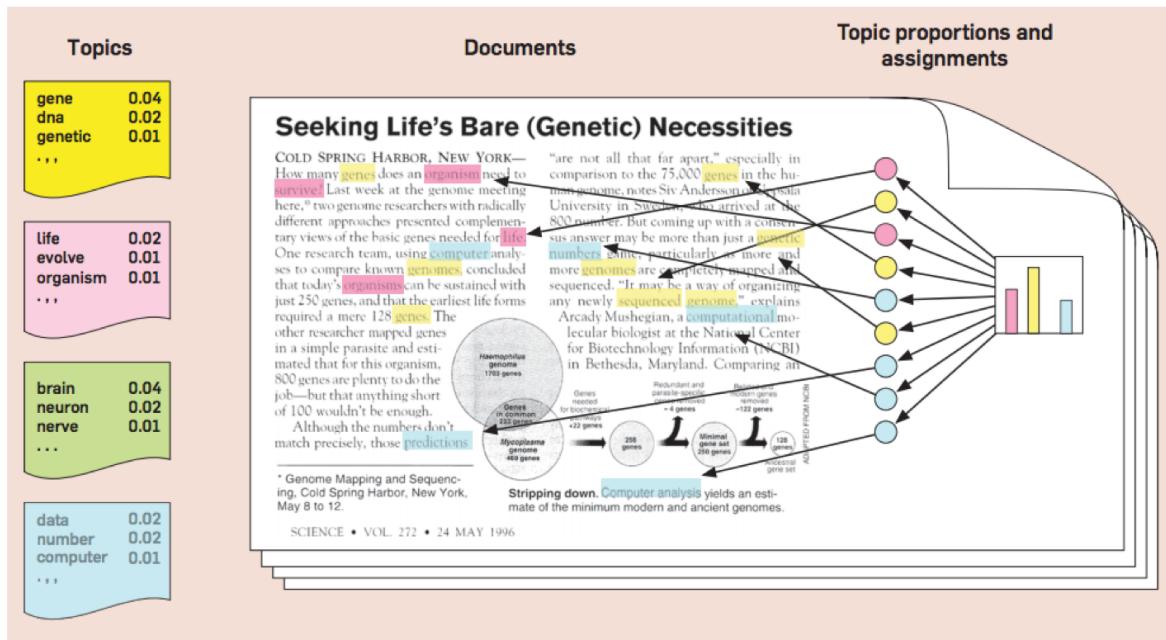
FIG. 3. Mathematical representation of folding-in  $q$  terms.

## LDA |



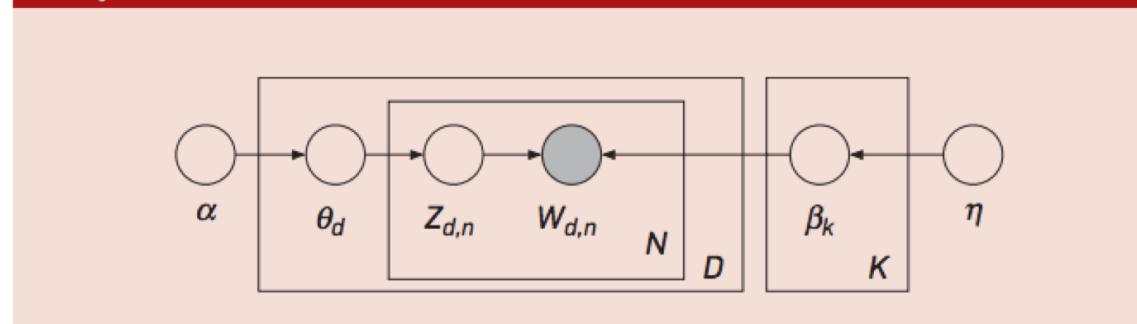
**Figure 2.** Illustration of the generative process and the problem of statistical inference underlying topic models

## LDA II

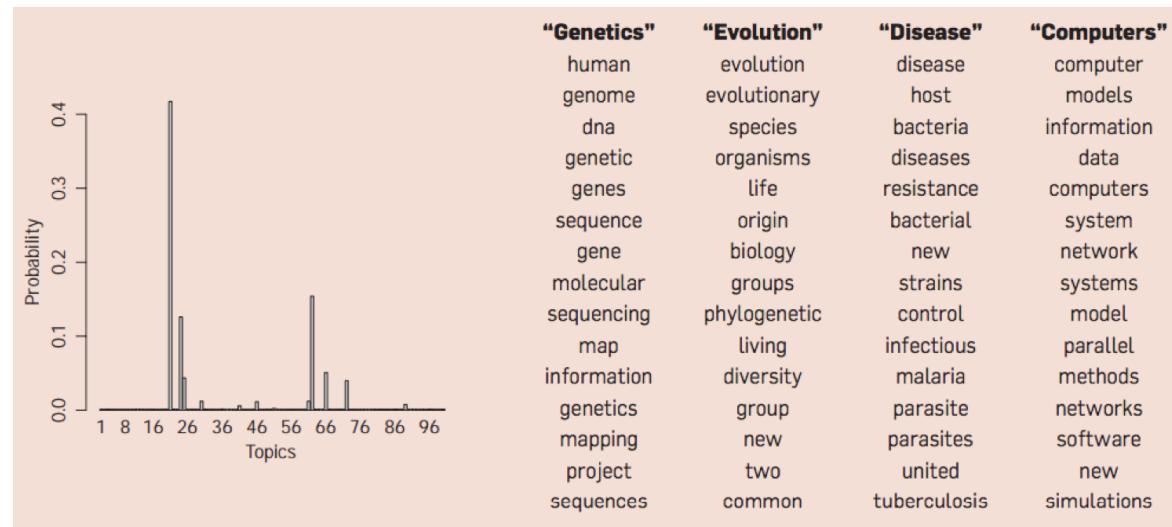


## LDA III

**Figure 4. The graphical model for latent Dirichlet allocation. Each node is a random variable and is labeled according to its role in the generative process (see Figure 1). The hidden nodes—the topic proportions, assignments, and topics—are unshaded. The observed nodes—the words of the documents—are shaded. The rectangles are “plate” notation, which denotes replication. The  $N$  plate denotes the collection words within documents; the  $D$  plate denotes the collection of documents within the collection.**

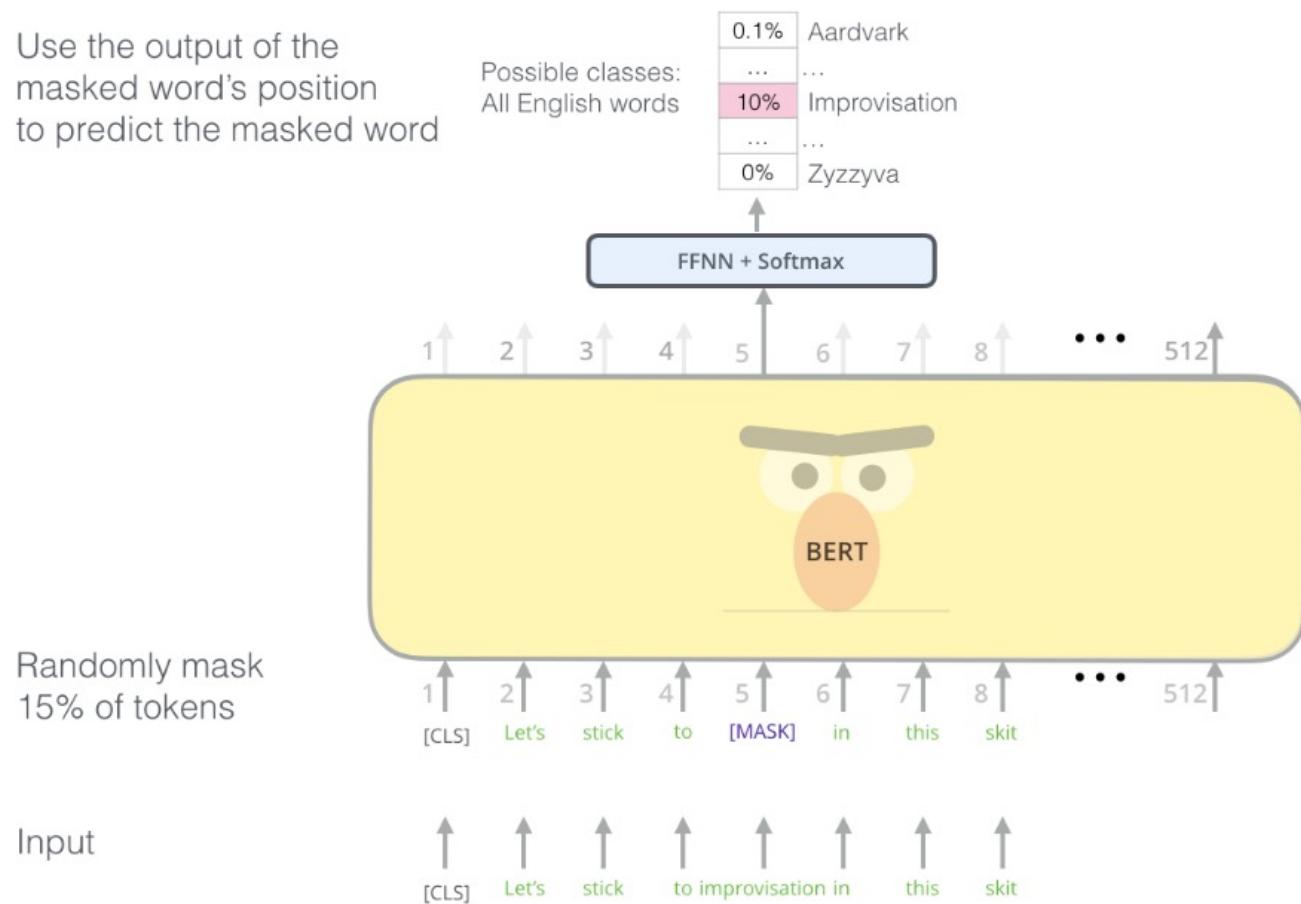


## LDA IV



# BERT

Use the output of the masked word's position to predict the masked word



BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word.

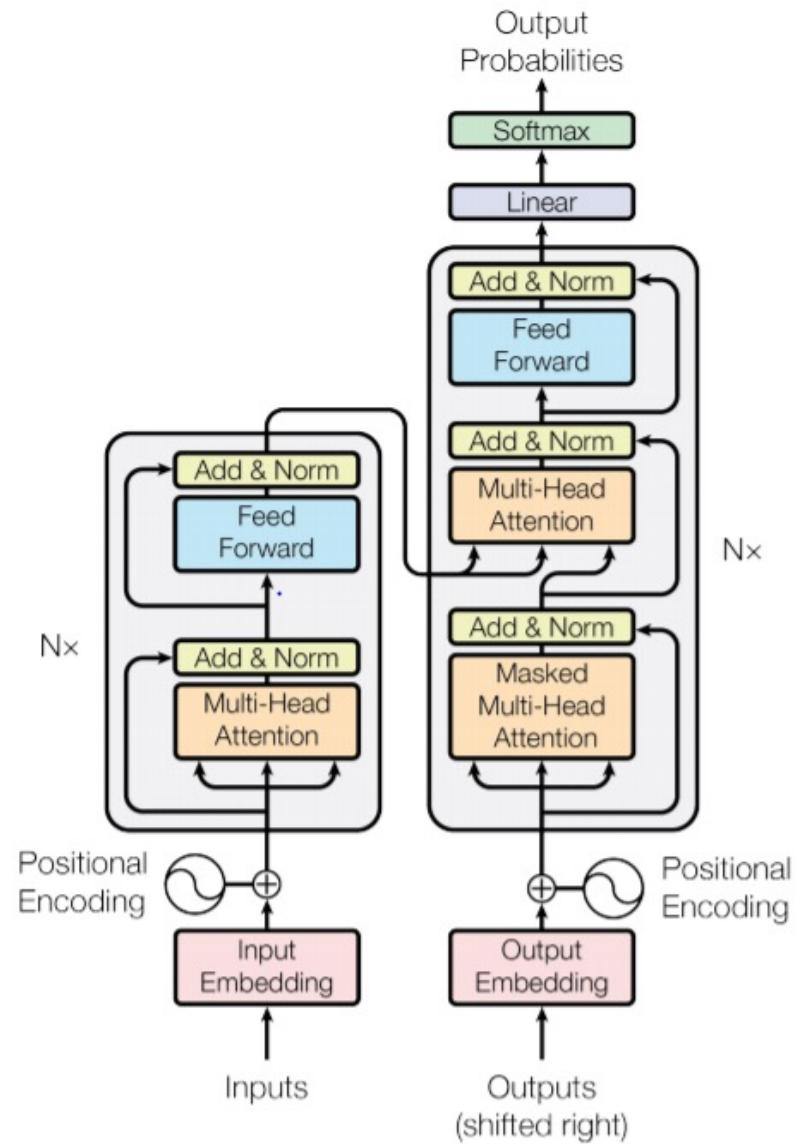
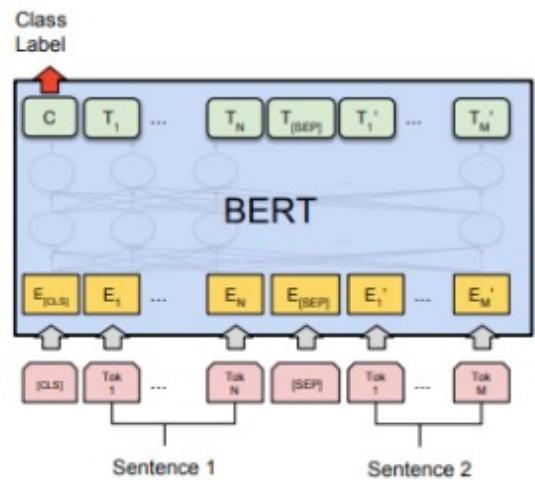
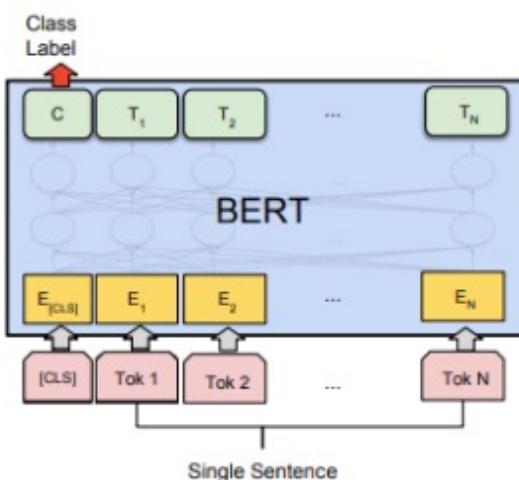


Figure 1: The Transformer - model architecture.

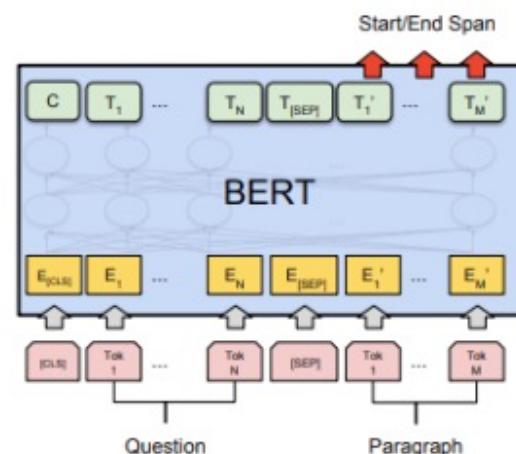
# BERT for different tasks



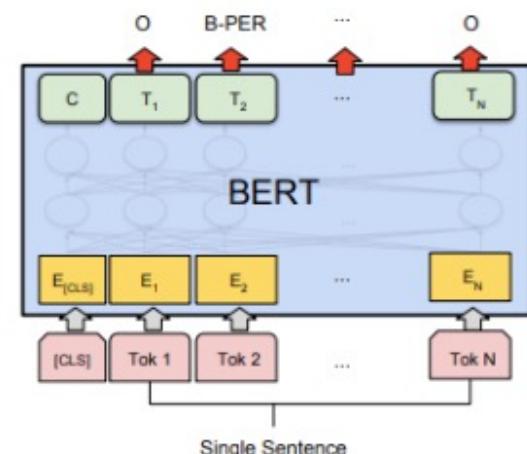
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1

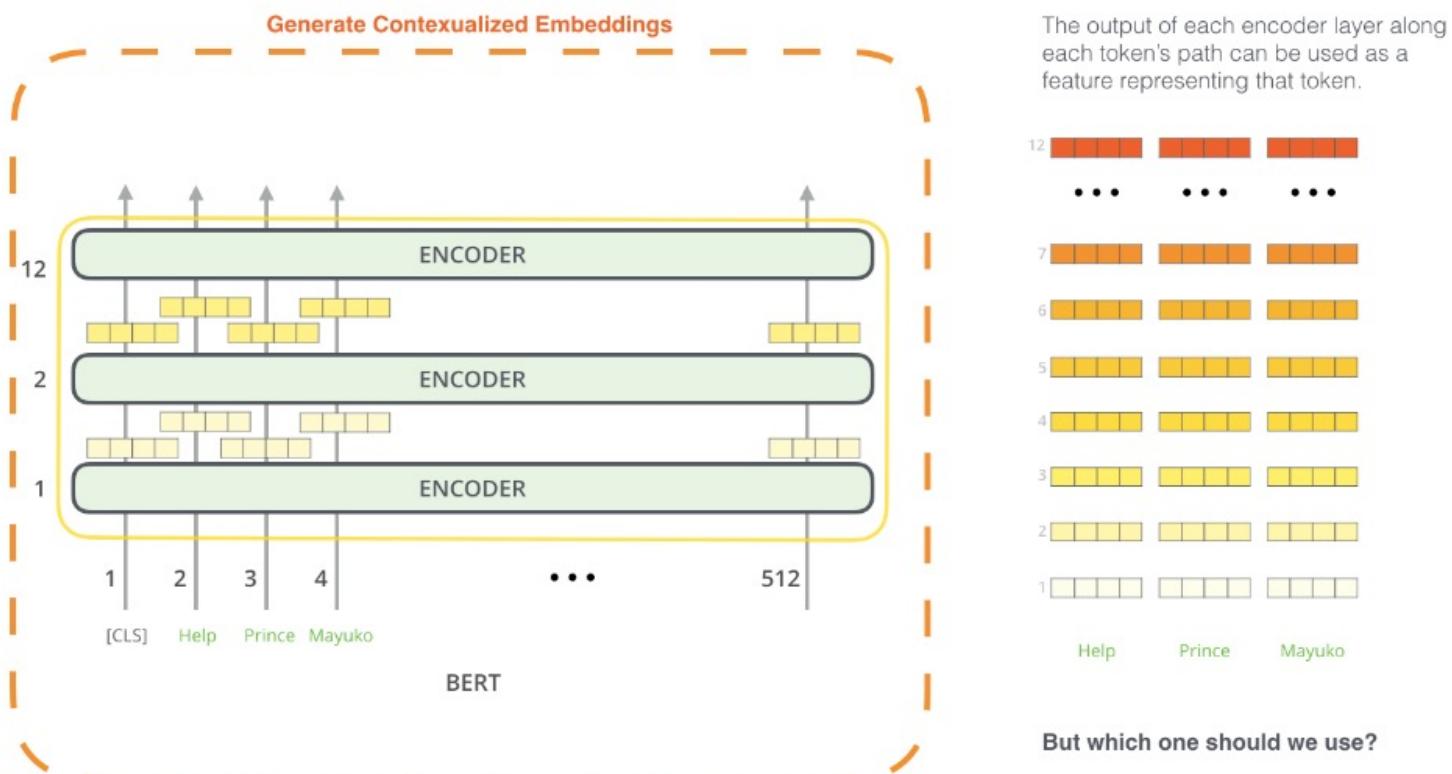


(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# BERT for feature extraction

## BERT for feature extraction

The fine-tuning approach isn't the only way to use BERT. Just like ELMo, you can use the pre-trained BERT to create contextualized word embeddings. Then you can feed these embeddings to your existing model – a process the paper shows yield results not far behind fine-tuning BERT on a task such as named-entity recognition.



# Neural Networks for text-aware content-based recommendation

## Content-Based Citation Recommendation

**Chandra Bhagavatula**

Allen Institute for AI

chandrab@allenai.org

**Sergey Feldman**

Data Cowboys \*

sergey@data-cowboys.com

**Russell Power**

Independent Researcher <sup>†</sup>

russell.power@gmail.com

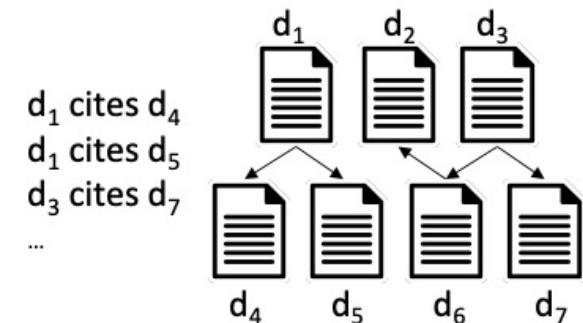
**Waleed Ammar**

Allen Institute for AI

waleeda@allenai.org

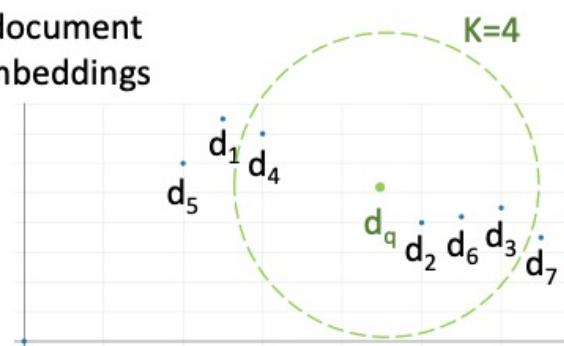
# 2-stage model: filter and re-ranking

Phase 1:  
candidate selection



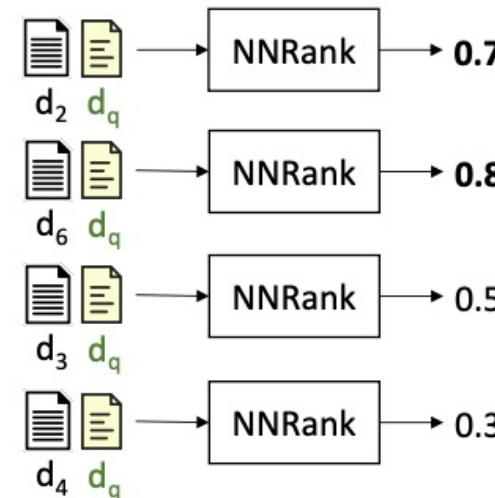
query  
document  
 $d_q$

document  
embeddings

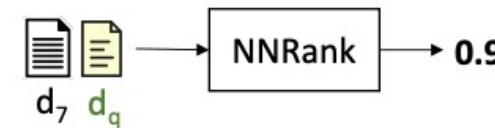


Phase 2:  
reranking

nearest neighbors of  $d_q$ :



cited in nearest neighbors:



reranked  
list

- $d_7$   
 $d_6$   
 $d_2$   
 $d_3$   
 $d_4$

top  $N=3$   
recommendations

# NNRank of stage 2

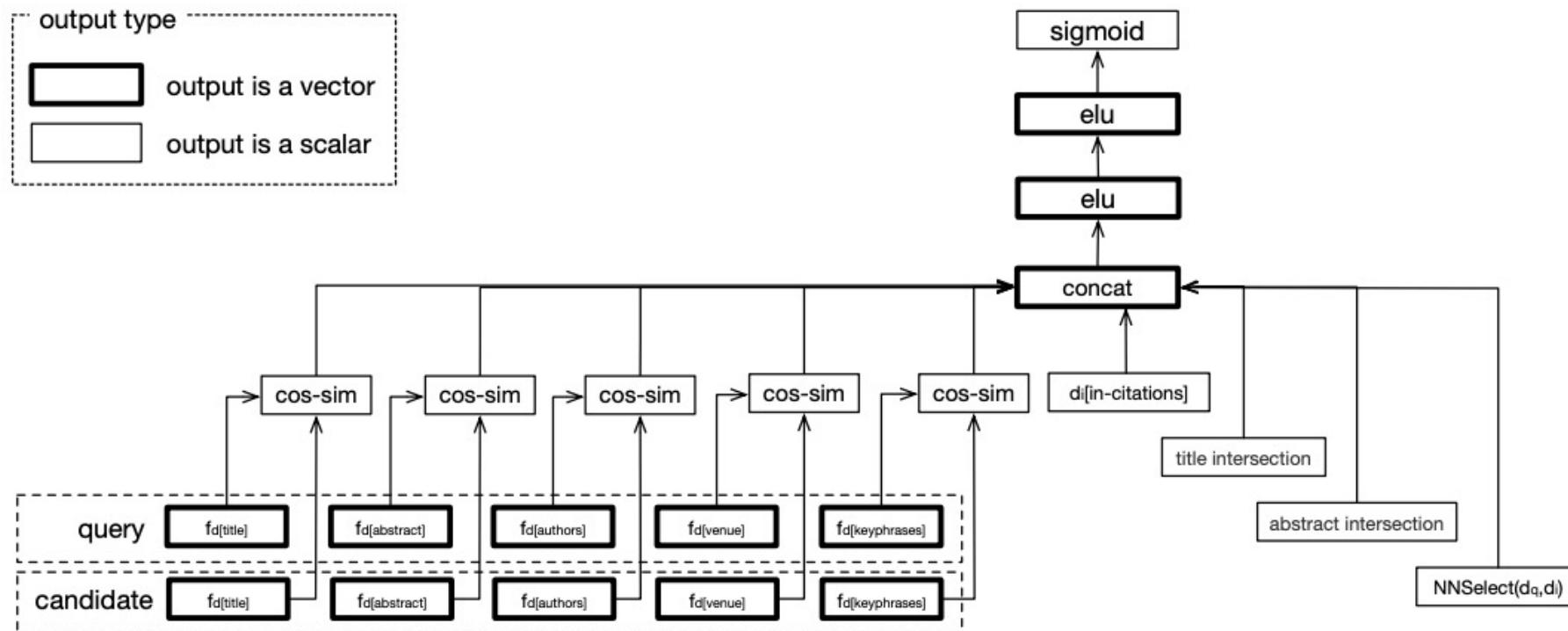


Figure 2: NNRank architecture. For each of the textual and categorical fields, we compute the cosine similarity between the embedding for  $d_q$  and the corresponding embedding for  $d_i$ . Then, we concatenate the cosine similarity scores, the numeric features and the summed weights of the intersection words, followed by two dense layers with ELU non-linearities. The output layer is a dense layer with sigmoid non-linearity, which estimates the probability that  $d_q$  cites  $d_i$ .

# Results

Method	DBLP		PubMed		OpenCorpus	
	F1@20	MRR	F1@20	MRR	F1@20	MRR
BM25	0.119	0.425	0.209	0.574	0.058	0.218
ClusCite	0.237	0.548	0.274	0.578	—	—
NNSelect	0.282±0.002	0.579±0.007	0.309±0.001	0.699±0.001	0.109	0.221
+ NNRank	0.302±0.001	0.672±0.015	0.325±0.001	0.754±0.003	<b>0.126</b>	0.330
+ metadata	<b>0.303±0.001</b>	<b>0.689±0.011</b>	<b>0.329±0.001</b>	<b>0.771±0.003</b>	0.125	<b>0.330</b>

Table 2: F1@20 and MRR results for two baselines and three variants of our method. BM25 results are based on our implementation of this baseline, while ClusCite results are based on the results reported in [Ren et al. \(2014\)](#). “NNSelect” ranks candidates using cosine similarity between the query and candidate documents in the embedding space (phase 1). “NNSelect + NNRank” uses the discriminative reranking model to rerank candidates (phase 2), without encoding any of the metadata features. “+ metadata” encodes the metadata features (i.e., keyphrases, venues and authors), achieving the best results on all datasets. Mean and standard deviations are reported based on five trials.

# News recommendation

## **UNBERT: User-News Matching BERT for News Recommendation**

**Qi Zhang , Jingjie Li\* , Qinglin Jia , Chuyuan Wang , Jieming Zhu , Zhaowei  
Wang and Xiuqiang He**

Huawei Noah's Ark Lab

{zhangqi193, lijingjie1, jiaqinglin2, wangchuyuan, jamie.zhu, wangzhaowei3,  
hexiuqiang1}@huawei.com

<https://www.ijcai.org/proceedings/2021/0462.pdf>

# Objetivo de UNBERT

- Separar la contribución a nivel de palabra y a nivel de noticia

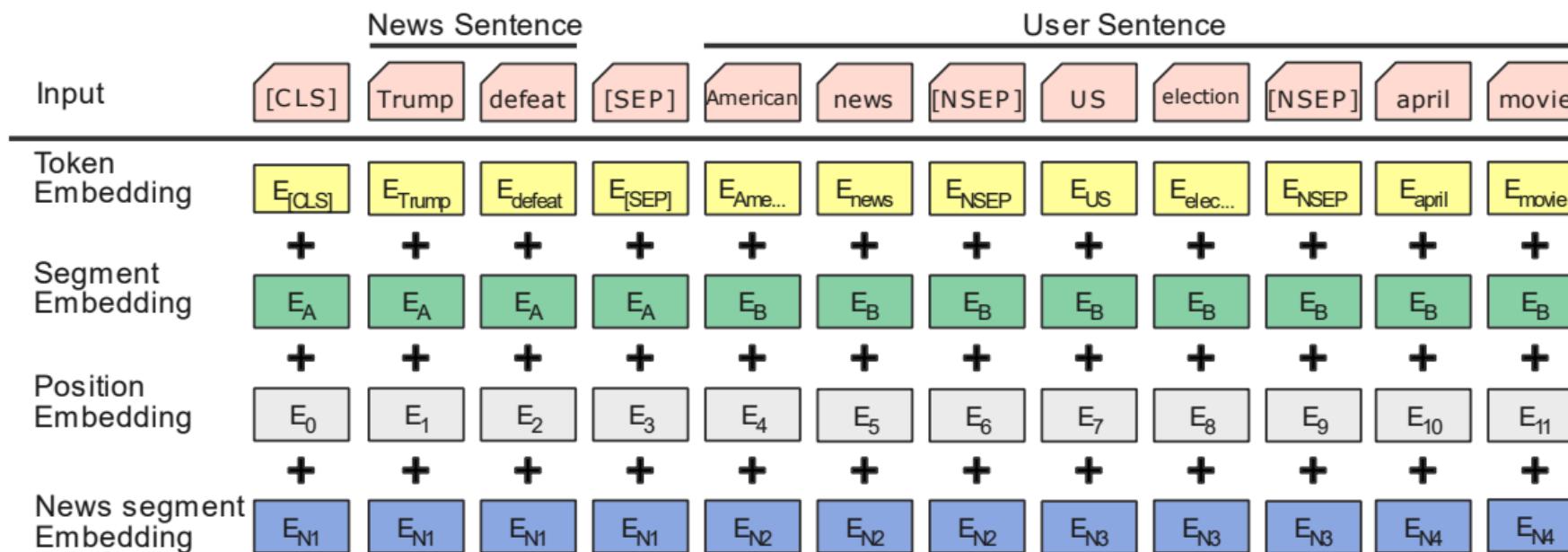


Figure 2: UNBERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings, the position embeddings, and the News segmentation embeddings.

# Arquitectura de UNBERT

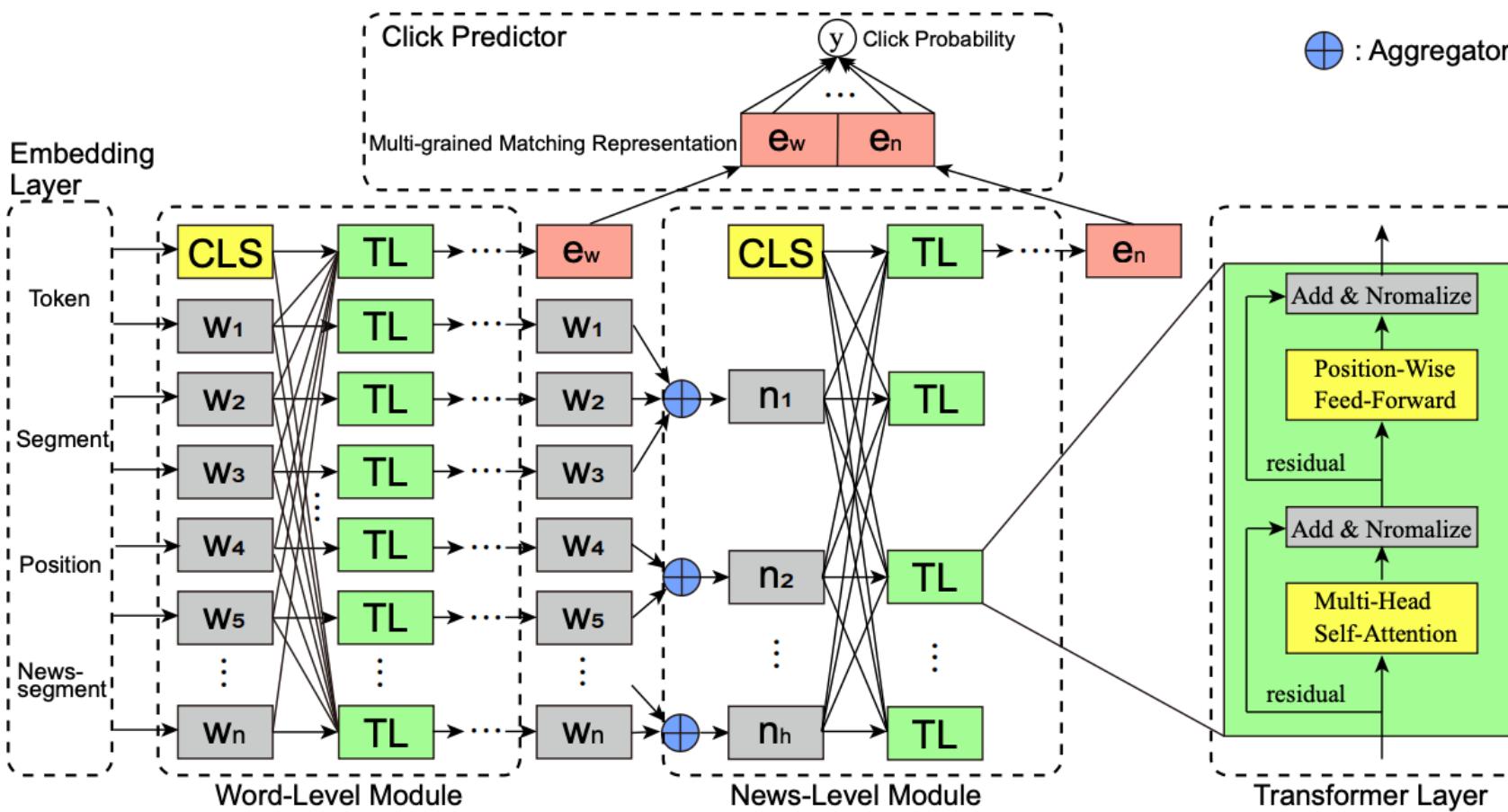


Figure 3: The overall architecture of our UNBERT approach.

# Resultados de UNBERT de UNBERT

Method	MIND-small				MIND-large			
	AUC	MRR	nDCG@5	nDCG@10	AUC	MRR	nDCG@5	nDCG@10
LibFM	0.5974	0.2633	0.2795	0.3429	0.6185	0.2945	0.3145	0.3713
DeepFM	0.5989	0.2621	0.2774	0.3406	0.6187	0.2930	0.3135	0.3705
DKN	0.6175	0.2705	0.2890	0.3538	0.6407	0.3042	0.3292	0.3866
NPA	0.6321	0.2911	0.3170	0.3781	0.6592	0.3207	0.3472	0.4037
NAML	<u>0.6550</u>	<u>0.3039</u>	<u>0.3308</u>	<u>0.3931</u>	0.6646	0.3275	0.3566	0.4140
LSTUR	0.6438	0.2946	0.3189	0.3817	0.6708	0.3236	0.3515	0.4093
NRMS	0.6483	0.3001	0.3252	0.3892	0.6766	0.3325	0.3628	0.4198
FIM	0.6502	0.3026	0.3291	0.3910	<u>0.6787</u>	<u>0.3346</u>	<u>0.3653</u>	<u>0.4221</u>
UNBERT	<b>0.6762</b>	<b>0.3172</b>	<b>0.3475</b>	<b>0.4102</b>	<b>0.7068</b>	<b>0.3568</b>	<b>0.3913</b>	<b>0.4478</b>
%Improv.	2.12	1.33	1.67	1.71	2.81	2.22	2.60	2.57
UNBERT-en <sup>△</sup>	-	-	-	-	0.7183	0.3659	0.4020	0.4581

Boldface indicates the best results (the higher, the better), while the second best is underlined. UNBERT-en<sup>△</sup> represents the ensemble score based on UNBERT which is at the top of <https://msnews.github.io/#leaderboard>.

# Referencias

- Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to information retrieval (Vol. 1, p. 6). Cambridge: Cambridge university press.
- Steyvers, M., & Griffiths, T. (2007). Probabilistic topic models. *Handbook of latent semantic analysis*, 427(7), 424-440.
- Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4), 77-84.