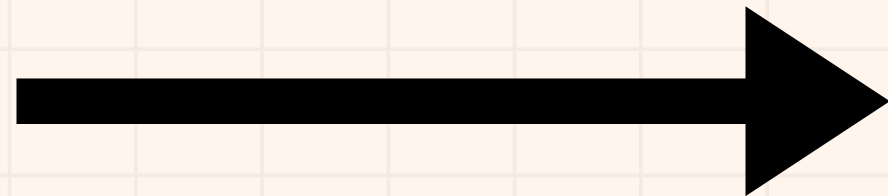
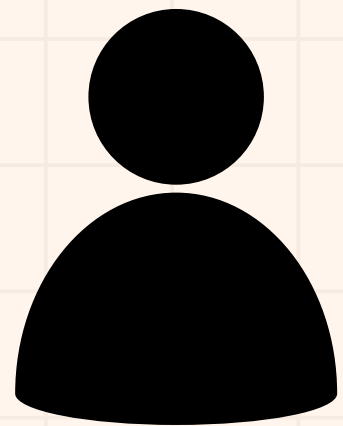


#

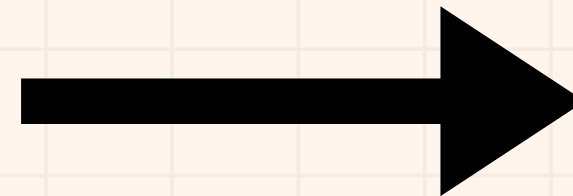
2007

Chris Messina introduce por primera vez el Hashtag (#) en Twitter

¿CÓMO FUNCIONA?



#cats



MAINSTREAM HASHTAGS

Hashtags normales, pero que además tienen:

- **Alto volumen** de uso global
- **Amplia visibilidad** en la plataforma

#KOBEDDEAD

v.s.

#KOBEDDEATH

PROBLEMAS A RESOLVER

Comprender el contexto a tiempo real

- Los Tweets nuevos reaccionan a **eventos emergentes** (ej. #NewPope)
- Los métodos estáticos carecen de **contexto actualizado**

Selección de etiquetas populares

- Muchos hashtags pueden describir el mismo tema, pero solo unos pocos son **mainstream**.
- Mantener manualmente una lista de hashtags populares es costoso y poco escalable

¿CÓMO SE HA ABARCADO ESTE PROBLEMA?

ESTADO DEL ARTE

MÉTODOS DE RECUPERACIÓN

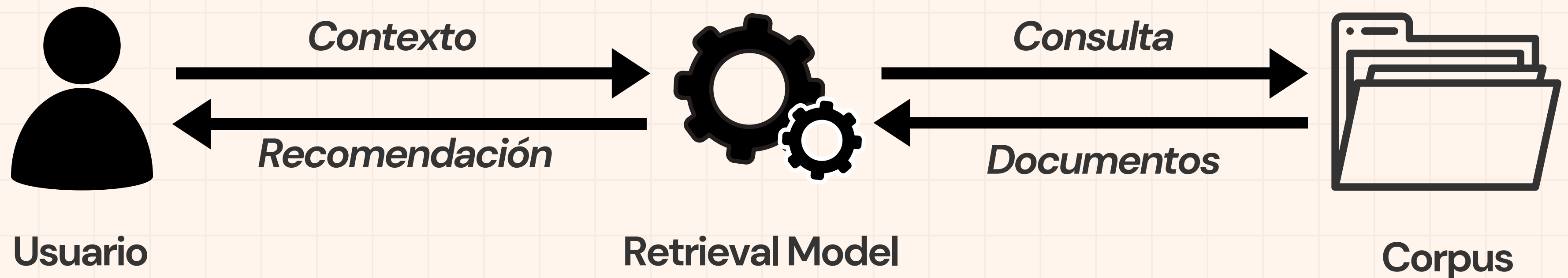
- BM25, SimCSE sobre un corpus histórico

Lo bueno:

- Asegura calidad de etiquetas preexistentes

Lo malo:

- No captan temas emergentes y requiere listas fijas



ESTADO DEL ARTE

MÉTODOS DE GENERACIÓN

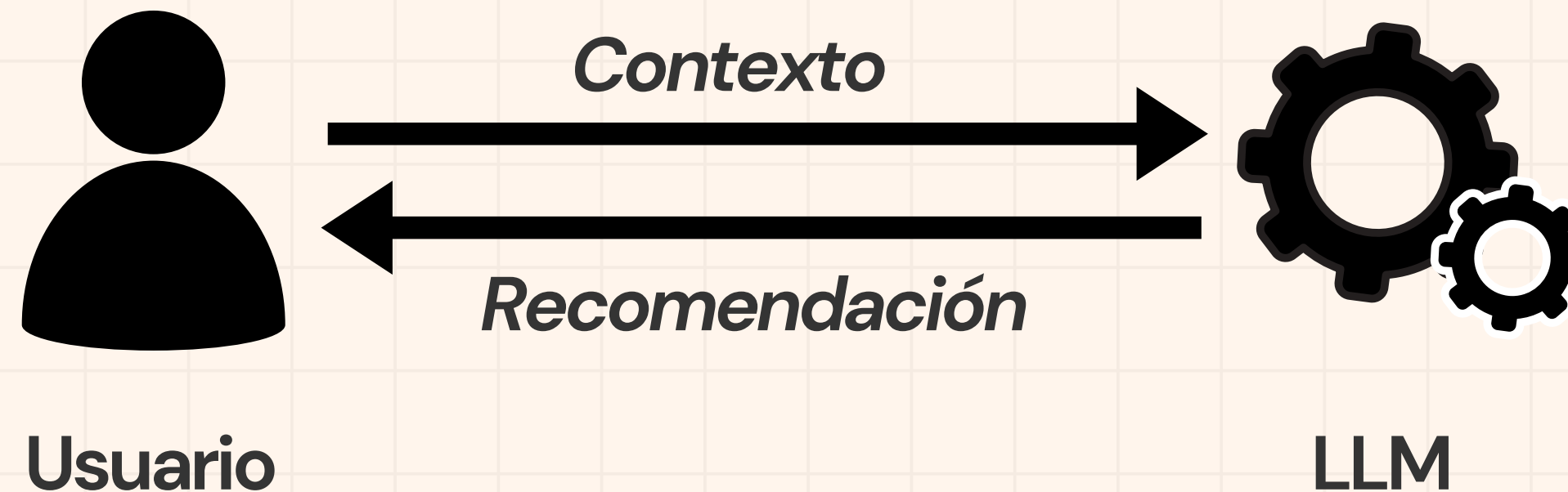
- seq2seq (T5, mT5), LLMs (ChatGPT).

Lo bueno:

- Buena comprensión de Tweets, generan etiquetas semánticamente válidas.

Lo malo:

- No garantiza popularidad.
- Puede **alucinar** hashtags



ESTADO DEL ARTE

MÉTODOS HÍBRIDOS

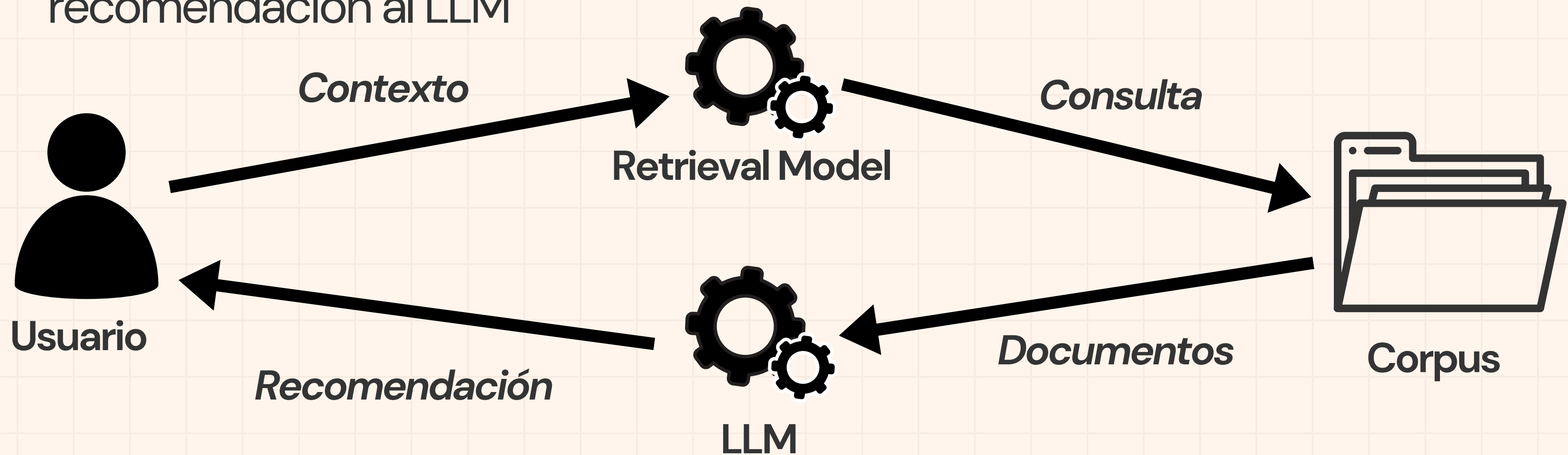
- Utilizan Recuperación temprana (p.e. Noticias) para dar una mejor recomendación al LLM

Lo bueno:

- Mejor comprensión multidominio.
- Reduce alucinaciones

Lo malo:

- No optimizado a tiempo real.
- Alta **complejidad** y costo.



RIGHT:

Retrieval-augmented Generation for Mainstream
Hashtag Recommendation

SOLUCIÓN PROPUESTA

SOLUCIÓN PROPUESTA

3 PRINCIPALES COMPONENTES

SOLUCIÓN PROPUESTA

3 PRINCIPALES COMPONENTES

1) BUSQUEDA RELEVANTE

SOLUCIÓN PROPUESTA

3 PRINCIPALES COMPONENTES

1) BUSQUEDA RELEVANTE

2) IDENTIFICAR PRINCIPALES TRENDS

SOLUCIÓN PROPUESTA

3 PRINCIPALES COMPONENTES

1) BUSQUEDA RELEVANTE

2) IDENTIFICAR PRINCIPALES TRENDS

3) INCORPORAR SELECCIONADOS Y GENERAR

SOLUCIÓN PROPUESTA

3 PRINCIPALES COMPONENTES

1) BUSQUEDA RELEVANTE

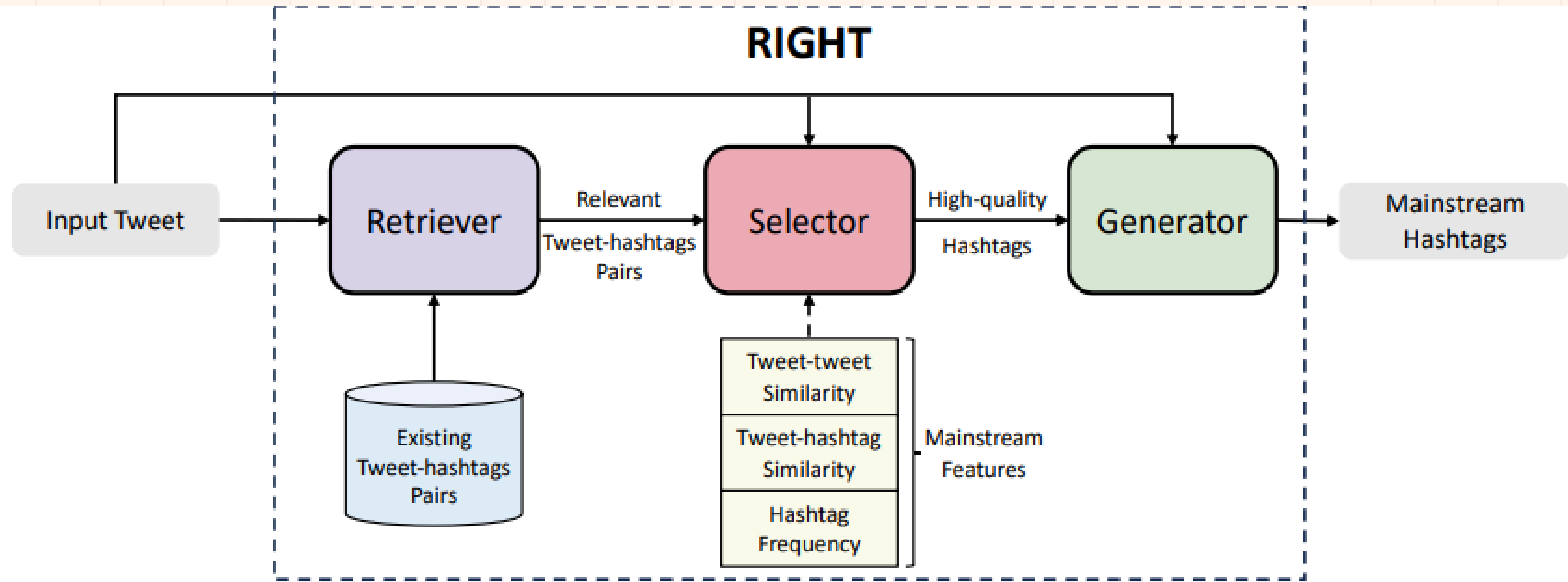
RETRIEVER

2) IDENTIFICAR PRINCIPALES TRENDS

SELECTOR

3) INCORPORAR SELECCIONADOS Y GENERAR

GENERATOR



EXTRAIDO DE [HTTPS://ARXIV.ORG/PDF/2312.10466](https://arxiv.org/pdf/2312.10466)

RETRIEVER

CONJUNTO DE DATOS (CORPUS): $C = \{(\tilde{t}_i, \tilde{H}_i)\}$ $\tilde{t}_i = \text{tweet}_i$
 $\tilde{h}_i = \text{hashtag}_i$

RETRIEVER

CONJUNTO DE DATOS (CORPUS): $C = \{(\tilde{t}_i, \tilde{H}_i)\}$ $\tilde{t}_i = \text{tweet}_i$
 $\tilde{h}_i = \text{hashtag}_i$

RELEVANCIA ENTRE TWEETS BASANDOSE EN SEMANTICA/LEXICA

RETRIEVER

CONJUNTO DE DATOS (CORPUS): $C = \{(\tilde{t}_i, \tilde{H}_i)\}$ $\tilde{t}_i = \text{tweet}_i$
 $\tilde{h}_i = \text{hashtag}_i$

RELEVANCIA ENTRE TWEETS BASANDOSE EN SEMANTICA/LEXICA

RETORNAR TOP-N PARES MAS SIMILARES: $\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_i, \tilde{H}_i, \tilde{s}_i)\} = R(t|C)$
 $s_i = \text{similitud entre } t \text{ y } t_i$

BM25 (RECUPERACIÓN BASADA EN TÉRMINOS) Y **SIMCSE** (RECUPERACIÓN DENSA BASADA EN EMBEDDINGS)

SELECTOR

FILTRAR LOS HASHTAG DE BAJA CALIDAD Y/O NO MAINSTREAM ENCONTRADOS

MAINSTREAM FEATURES:

SELECTOR

FILTRAR LOS HASHTAG DE BAJA CALIDAD Y/O NO MAINSTREAM ENCONTRADOS

MAINSTREAM FEATURES:

SIMILITUD TUIT INPUT Y TUIT RECUPERADO

SIMILITUD TUIT INPUT Y HASHTAG RECUPERADO

FRECUENCIA DE USO DE HASHTAG

SELECTOR

TRAINING

SELECTOR

TRAINING

EMBEDDINGS Y SU IMPORTANCIA

SELECTOR

TRAINING

EMBEDDINGS Y SU IMPORTANCIA

AJUSTES Y PROPUESTA

t_i = muestra original (tuit)

t_i^+ = hashtag positivo (relevantes)

t_i^- = hashtag positivo perturbada (no relevantes)

$$\{(t_i, t_i^+, t_i^-) | i = 1, \dots, N\}$$

SELECTOR

TRAINING

t_i

"KOBE BRYANT DIES IN HELICOPTER CRASH".

t_i^+

#KOBEBRYANT

t_i^-

SINÓNIMOS (70%): **#KOBEBASKETBALL.**

ELIMINACIÓN (10%): **#BRYANT.**

INTERCAMBIO (10%): **#BRYANTKOBE.**

SINÓNIMOS (10%): **#KOBEELEGENDBRYANT.**

SELECTOR

TRAINING

- FUNCION DE PERDIDA →
- MINIMIZAR SIMILITUDES ENTRE INCORRECTOS
 - MAXIMIZAR SIMILITUDES CON CORRECTOS

CONSIDERANDO EMBEDDINGS

$$\mathcal{L}_S = -\log \frac{e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_i}^+)/\tau}}{\sum_{j=1}^L \left(e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_j}^+)/\tau} + e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_j}^-)/\tau} \right)}$$

ROBERTA-LARGE (INGLÉS) O **BERT-BASE-CHINESE** (CHINO)

SELECTOR

INFERENCE

INPUT : TOP-N TUIITS ANTERIOMENTE GENERADOS

$\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_i, \tilde{H}_i, \tilde{s}_i)\}$

SELECTOR

INFERENCE

INPUT : TOP-N TUIITS ANTERIOMENTE GENERADOS

$\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_i, \tilde{H}_i, \tilde{s}_i)\}$

AGRUPACIÓN DE HASHTAGS

SELECTOR

INFERENCE

INPUT : TOP-N TUIITS ANTERIOMENTE GENERADOS

$$\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_i, \tilde{H}_i, \tilde{s}_i)\}$$

AGRUPACIÓN DE HASHTAGS

$$\{\tilde{h}_1, \dots, \tilde{h}_M\}$$

CONSOLIDACION DE TODOS LOS HASHTAGS

$$\{\tilde{s}_{i,1}, \dots, \tilde{s}_{i,f_i}\}$$

SIMILITUDES ASOCIADAS

$$\{f_1, \dots, f_M\}$$

NUMERO DE TUIITS QUE LO CONTIENEN

SELECTOR

INFERENCE

CALCULO Y SELECCIÓN TOP-K

$$\ddot{s}_m = \mathcal{S}(t, \tilde{h}_m)$$

SIMILITUD ENTRE EL TUIT T Y EL MODELO ANTERIORMENTE ENTRENADO

$$s_i = \left(\left(\frac{1}{f_i} \sum_{j=1}^{f_i} \tilde{s}_{i,j} \right) + \ddot{s}_i \right) \times \left(1 + \left(\frac{f_i - 1}{10} \right) \right)$$

CALCULO SIMILITUD PONDERIZADA PARA LUEGO SELECCION FINAL

GENERATOR

GENERAR HASHTAGS DE **CALIDAD** DADO UN **TUIT** (T)
SELECCIONADO Y UN SET DE **HASHTAGS SELECCIONADOS**

$(\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_k)$

GENERATOR

GENERAR HASHTAGS DE **CALIDAD** DADO UN **TUIT** (T)
SELECCIONADO Y UN SET DE **HASHTAGS SELECCIONADOS**
POR SELECTOR

$$(\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_k)$$

CONCATENAR TUI TS Y HASHTAGS

$$I = \langle t, \text{SEP1}, \tilde{h}_1, \text{SEP1}, \tilde{h}_2, \dots, \text{SEP1}, \tilde{h}_k \rangle$$

CONCATENACION CON SEPARACIONES ENTRE LOS ELEMENTOS

GENERATOR

GENERAR HASHTAGS DE **CALIDAD** DADO UN **TUIT** (T)
SELECCIONADO Y UN SET DE **HASHTAGS SELECCIONADOS**
POR SELECTOR

$$(\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_k)$$

CONCATENAR TUIT Y HASHTAGS

$$I = \langle t, \text{SEP1}, \tilde{h}_1, \text{SEP1}, \tilde{h}_2, \dots, \text{SEP1}, \tilde{h}_k \rangle$$

CONCATENACION CON SEPARACIONES ENTRE LOS ELEMENTOS

GENERAR SECUENCIALMENTE Y CON AJUSTES

$$O = \langle h_1, \text{SEP2}, h_2, \dots, \text{SEP2}, h_{|H|} \rangle$$

OUTPUT GENERADO

GENERATOR

GENERAR HASHTAGS DE **CALIDAD** DADO UN **TUIT** (T)
SELECCIONADO Y UN SET DE **HASHTAGS SELECCIONADOS**
POR SELECTOR

$$(\tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_k)$$

CONCATENAR TUI TS Y HASHTAGS

$$I = \langle t, \text{SEP1}, \tilde{h}_1, \text{SEP1}, \tilde{h}_2, \dots, \text{SEP1}, \tilde{h}_k \rangle$$

CONCATENACION CON SEPARACIONES ENTRE LOS ELEMENTOS

GENERAR SECUENCIALMENTE Y CON AJUSTES

$$O = \langle h_1, \text{SEP2}, h_2, \dots, \text{SEP2}, h_{|H|} \rangle$$

OUTPUT GENERADO

GENERATOR

$$\mathcal{L}_g = \sum_{(I,O) \in \mathcal{D}} -\log p(O|I; \theta_g)$$

PERDIDA ASOCIADA AL GENERAR

GENERATOR

$$\mathcal{L}_g = \sum_{(I,O) \in \mathcal{D}} -\log p(O|I; \theta_g)$$

PERDIDA ASOCIADA AL GENERAR SOBRE LOS TOKENS DE SALIDA.

T5 , BART

GPT

RESULTADOS

TESTEO PARA CADA COMPONENTE, DATASET EMPLEADO

THG

ENGLISH TWITTER

WHG

CHINESE WEIBO

RESULTADOS

TESTEO PARA CADA COMPONENTE, DATASET EMPLEADO

THG Y WHG

METRICAS

ROUGE

SUPERPOSICIÓN EN TÉRMINOS DE PALABRAS

F1@K

PRESICION Y RECALL

RESULTADOS

TESTEO PARA CADA COMPONENTE, DATASET EMPLEADO

THG Y WHG

MÉTRICAS **ROUGE** Y **F1@K**

COMPARACIÓN CON OTROS MODELOS

MÉTODOS BASADOS EN **RECUPERACIÓN**

(BM25 Y SIMCSE)

MÉTODOS BASADOS EN **GENERACIÓN**

(CHATGPT , SEQ2SEQ , SEGTRM SOFT)

MÉTODOS **HÍBRIDOS**

(RIGHT)

RESULTADOS

Model	THG					WHG			
	RG-1	RG-2	RG-L	F1@1	F1@5	RG-1	RG-2	RG-L	F1@1
<i>Retrieval-based Methods</i>									
BM25	16.23	4.17	15.11	5.92	9.84	61.98	58.76	61.81	48.20
SimCSE	28.43	10.34	26.38	12.40	15.15	59.71	55.81	59.54	47.65
<i>Generation-based Methods</i>									
ChatGPT	44.60	27.67	39.29	9.72	26.08	32.27	24.54	31.80	7.9
SEGTRM Soft	51.18	37.15	47.05	27.17	29.02	55.51	51.28	54.30	30.72
Seq2Seq	59.90	41.39	59.15	29.75	41.71	66.64	61.71	66.39	48.60
<i>Retrieval-augmented Generative Methods (Ours)</i>									
RIGHT ^{ChatGPT}	47.54	25.63	44.47	22.39	31.09	48.17	41.51	47.75	26.15
RIGHT _{BM25}	<u>61.60</u>	<u>43.77</u>	<u>60.85</u>	<u>30.27</u>	<u>42.98</u>	70.62*	66.12*	70.35*	53.85*
RIGHT _{SimCSE}	62.11*	43.86*	61.39*	30.58*	43.23*	<u>68.84</u>	<u>64.19</u>	<u>68.56</u>	<u>51.50</u>

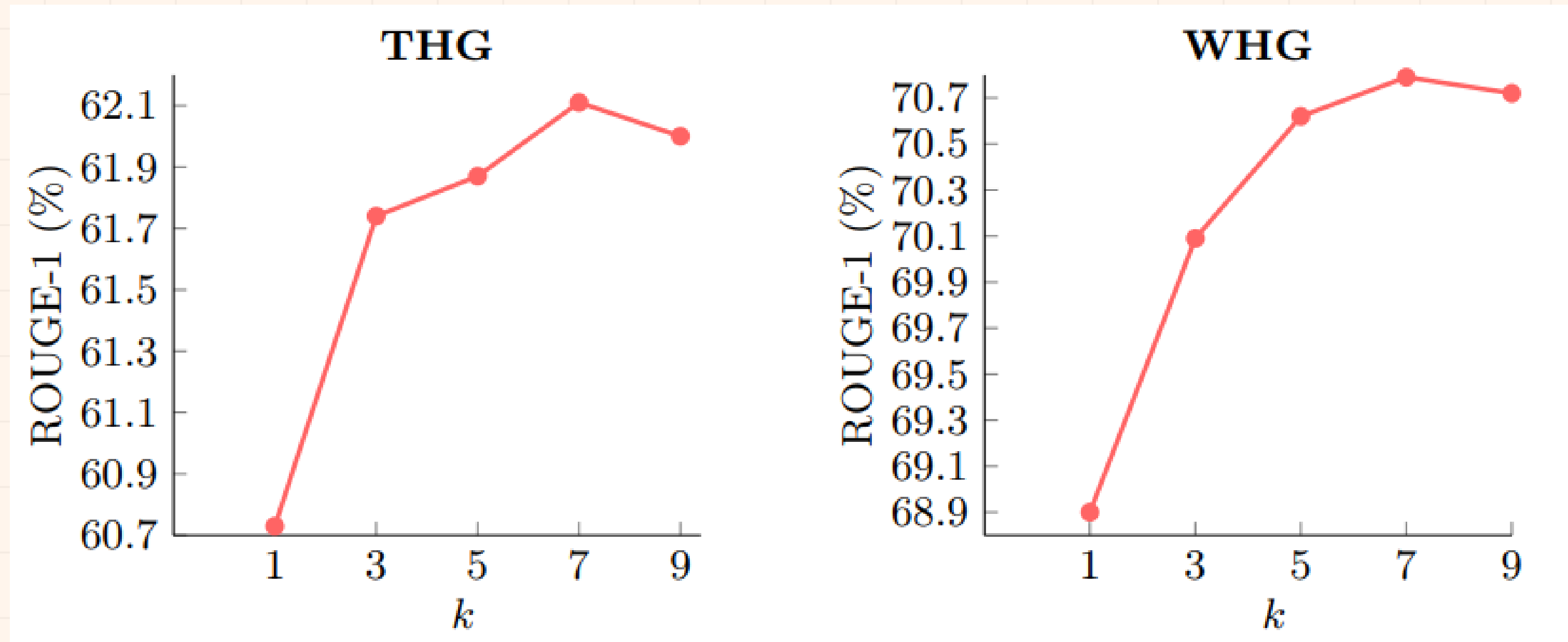
RESULTADOS

PRUEBAS DE ABLACIÓN PARA JUSTIFICAR CADA COMPONENTE DEL MODELO.

Model	ROUGE-1	ROUGE-2	ROUGE-L	F1@1	F1@5
RIGHT	62.11	43.86	61.39	30.58	43.23
w/o Retriever	59.91	41.63	59.23	29.66	41.70
w/o Selector	60.49	42.06	59.76	30.22	41.95
w/o Generator	36.24	16.02	32.86	24.61	26.73

RESULTADOS

IMPACTO DEL NÚMERO DE HASHTAG AUMENTADOS



RESUMEN Y CONTRIBUCIONES

- Arquitectura **híbrida modular**
- Selector con señales de popularidad
- Actualización automática de Corpus
- Ganancias empíricas significativas
- **Impacto en el mundo real**

MUCHAS GRACIAS

¿PREGUNTAS?

POR:
RENÉ SAAVEDRA
CARLOS OLGUÍN
GERALDINE COLI ACEVEDO

BIBLIOGRAFÍA

- Parker. A, (2011). Twitter's Secret Handshake. The New York Times
<https://www.nytimes.com/2011/06/12/fashion/hashtags-a-new-way-for-tweets-cultural-studies.html>
- Fan, R.-Z., Fan, Y., Chen, J., Guo, J., Zhang, R., & Cheng, X. (2023). RIGHT: Retrieval-augmented generation for mainstream hashtag recommendation. arXiv. <https://arxiv.org/abs/2312.10466>