

RAG de la Coppermind, una Wikipedia del Cosmere

Ferrán Plana Caminero

Introducción

Objetivos

Desarrollo

I. Creación de la base de datos

En primer lugar, he elegido la lista de entradas de la *wikipedia* que quiero utilizar como base de datos. Me he decantado por centrarme no en una saga en concreto, sino en los conceptos que son comunes a toda saga, a todo el universo del *Cosmere*.

Adjunto la *lista* con las urls de las entradas elegidas. He hecho un script que descarga el html a partir de la lista de urls. La descarga de cada entrada, al ser de una web, ha sido en formato html.

(Script en Creación de base de datos-TFB.ipynb en local)

I.II. Tratamiento de los htmls

Para optimizar el RAG he convertido los htmls en markdown. Además, he hecho cierto procesamiento para eliminar varias cosas:

- Eliminación de vínculos a otras entradas de la *wikipedia*.
- Eliminación de citas.
- Eliminación de imágenes.
- Eliminación de la sección de notas en adelante (sección donde se ponen todas las referencias, no aporta valor).

De momento no voy a hacer más tratamiento al texto. Más adelante, con la ayuda del Test Automático, estudiaré si sería conveniente.

(Script en html_to_markdown.ipynb en local, luego subir markdowns al drive)

I.III. Embeddings

Una vez procesado el contenido de cada entrada, he dividido el texto completo de la base de datos en *chunks* con un *overlap* del 20% (el tamaño del *chunk* está por ver, a estudiar con el Test Automático). He hecho los *embeddings* con el modelo *paraphrase-MiniLM-L6-v2*, un modelo de transformers que es relativamente compacto y eficiente cuya principal aplicación es la búsqueda semántica. He obtenido el modelo de *Hugging Face*.

El código divide el texto en *chunks*, hace los *embeddings* de cada *chunk* y guarda en *Google Drive* cada *chunk* y *embedding* con varios metadatos asociados: el número de *chunk* (entre todos los *chunks* del documento), el número de documento o *DOCID* (número de entrada de las totales) y su nombre (la *url* asociada, o algo similar).

(Script en Creación de la base de datos.ipynb, Google Colab)

II. Funcionamiento del *RAG*

Un modelo *RAG* (*Retrieval-Augmented Generation*) consiste en dos partes principales: búsqueda semántica y generación de la respuesta. La búsqueda semántica consiste en hallar un cierto número de *chunks* de la base de datos que tienen mayor similitud con la pregunta del usuario. Una vez hallados, se le pide a un *LLM* componer una respuesta en base a esa pregunta y los chunks encontrados.

II.I. Búsqueda semántica

Una vez hechos y guardados los *embeddings*, podemos acceder a ellos mediante *Google Drive*. Al hacer una pregunta, se hace el *embedding* de esa pregunta y se calcula la similitud entre esa pregunta y todos los *embeddings* de la base de datos. Otro de los hiperparámetros, junto con el tamaño de *chunk* y el *overlap*, es *top-n*, el número de chunks con más similitud que nos quedamos. El valor óptimo de este hiperparámetro también lo decidiremos con la ayuda del Test Automático.

Al hacer una pregunta, obtenemos los *chunks* con mayor similitud, el número de *chunk*, el *DOCID*, su nombre y la similitud entre ese *chunk* y la pregunta. Quedaría esto:

Question: ¿Se puede leer algún libro del Cosmere sin haber leído otras sagas? The 3 chunks most similar to your question are:

1. DOCID: 12 | Document Name: <https://es.coppermind.net/wiki/Cosmere> | Chunk number: 3 | Similarity: 0.6855
subyacentes, apareciendo algunos personajes en otros mundos ajenos al suyo. A pesar de las conexiones, Brandon ha dejado claro que uno no necesita ningún conocimiento del Cosmere en general para leer, entender, o disfrutar de los libros que tienen lugar en él. Las secuencia principal del Cosmere consistirá en la saga *Dragonsteel*”), la trilogía de *Elantris*, al menos cuatro eras de la saga *Nacidos de la bruma*”) y *El archivo de las tormentas*. La historia del Cosmere no incluye ningún libro que haga referencia a la Tierra, puesto que la tierra no está en el Cosmere. Para una lista completa
2. DOCID: 19 | Document Name: <https://es.coppermind.net/wiki/Esquirla del Amanecer> | Chunk number: 49 | Similarity: 0.6744
no se refiera a Sigzil. En *Viento y verdad*, se confirmó que Hoid tuvo en su poder la Esquirla del Amanecer Existe durante los sucesos de los libros 1-5 de *El archivo de las tormentas*.
Notes 1. ↑ a b c d e f g h i j k l m n o p q r s t u *Esquirla del Amanecer (novella)* capítulo 19”)Summary: Esquirla del Amanecer (novella)/Chapter_19/Chapter 19 (la página no existe)“)#/Chapter 19”) 2. ↑ a b General Reddit 2022 — Arcanum - 2022-12-02Cite: Arcanum-15961# 3. ↑ a b c Dawnshard Annotations Reddit Q&A — Arcanum -
3. DOCID: 24 | Document Name: <https://es.coppermind.net/wiki/Hoid> | Chunk number: 184 | Similarity: 0.6562
qué libro fue eso respondió *Brazales de Duelo*”). * La misión de Hoid quizás sea: «hacer aquello que una vez fue». * Hoid no está impresionado por los Sangre Espectral. * Hoid detesta al Grupo, y los miembros de este último que le conocen también le detestan. * Aunque una vez le dijo a Kaladin (bastante acertado) que su vida comenzó como palabras en una página, este hecho no tenía la intención de romper la cuarta pared. * Hoid se ha travestido en el pasado, «muchas veces». * Antes de los eventos de *Palabras Radiantes*, a Hoid no le habían

Además, se puede obtener una gráfica interesante: el histograma de la similitud entre cada *chunk* de la base de datos y la pregunta.

(Script en Creación de la base de datos.ipynb, Google Colab)

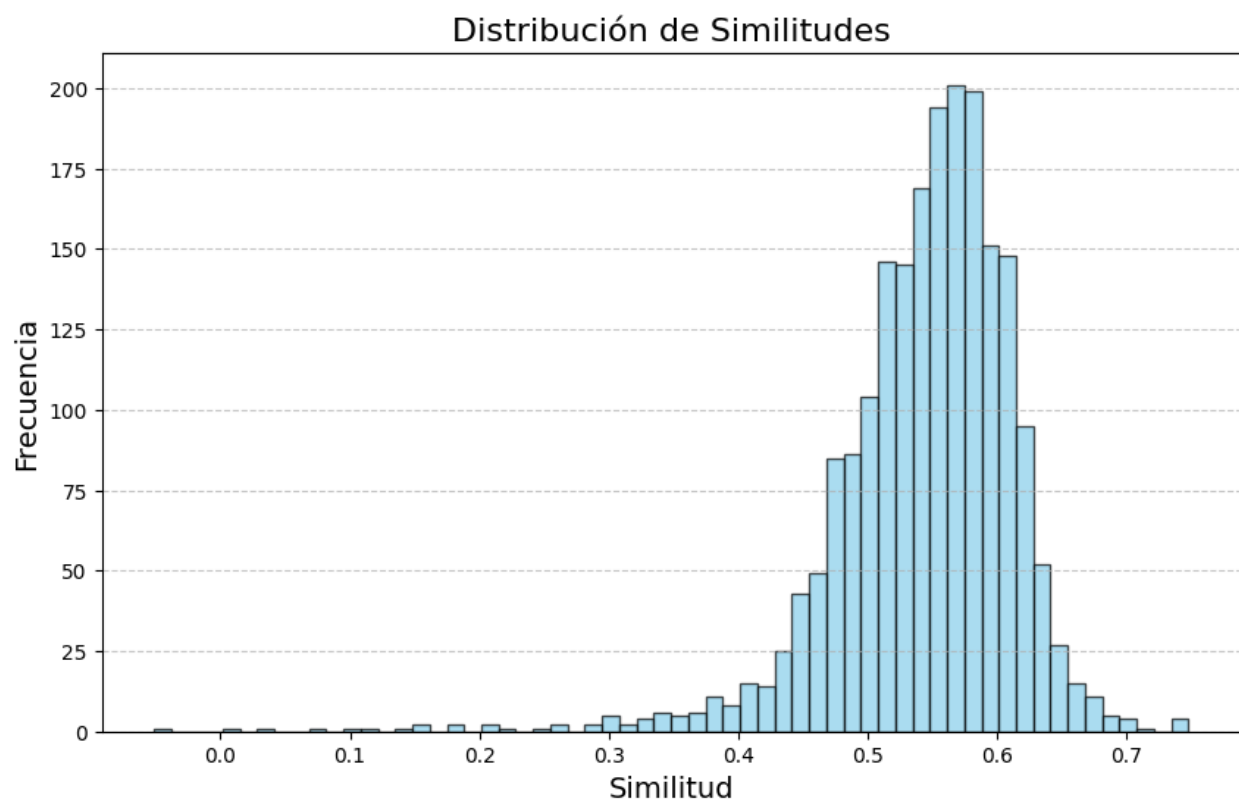


Figure 1: Distribución de *chunks* con mayor similitud respecto a una pregunta dada.

II.II. Generación de la respuesta

Esos *chunks* con mayor similitud con la pregunta se le pasan, junto con la pregunta del usuario, a un *LLM*, pidiéndole que responda a la pregunta del usuario en base a los *chunks* encontrados.

El modelo seleccionado se decidirá más adelante a raíz de una serie de tests. Para esta y posteriores llamadas a *LLMs* se utilizarán varios proveedores que permiten hacer llamadas gratis vía API (ver *LLMs_free_API_keys.ipynb* para más detalles).

El *system prompt* es el siguiente:

Eres un asistente virtual experto en responder preguntas. A continuación vas a recibir la pregunta de un usuario y el conocimiento que debes utilizar para responderla en el siguiente formato:

Pregunta: Pregunta del usuario
Conocimiento: Nombre del documento 1: contenido del documento 1
Nombre del documento 2: contenido del documento 2 ... Nombre del documento N: contenido del documento N

Responde como si fueras un chatbot de una wikipedia. Después de dar tu respuesta completa, di el nombre de los documentos en los que te has basado para elaborarla. Si te has basado dos veces en el mismo documento, no lo repitas al referenciarlo. Hazlo en este formato:

Pon aquí tu respuesta [[Pon aquí únicamente el nombre del primer documento en el que te has basado]] [[Pon aquí únicamente el nombre del segundo documento en el que te has basado, siempre y cuando no hayas puesto el mismo nombre antes]] ...

No utilices conocimiento propio de tu entrenamiento, utiliza solo el que se te proporciona. Si parte del conocimiento que se te proporciona no te sirve para responder a la pregunta, no lo utilices. Cita únicamente los documentos en los que te has basado para elaborar la respuesta. Si con el conocimiento que se te proporciona no puedes responder a la pregunta, responde únicamente “Lo siento, no puedo responderte a esa pregunta” y no cites ningún documento.

El *user prompt* es el siguiente:

Pregunta: {Pregunta del usuario}

Conocimiento:

{Nombre del documento del chunk}: {contenido del chunk}

Como se puede observar, no solo se le pide al *LLM* que componga la respuesta, sino que cite sus fuentes. De esta forma el usuario puede, con solo pinchar en el nombre de la referencia, acceder a dicha entrada de la *wikipedia* mediante la *url*.

III. Test Automático

Con el objetivo de poder valorar si los cambios tienen un impacto positivo en el modelo, he creado un test de regresión, al que llamaré Test Automático. Este test consiste en 137 preguntas de las que sé la respuesta correcta, a la que llamaré respuesta *best*, y la entrada (o entradas) de la *wikipedia* donde se responde a esa pregunta, que llamaré documento *best*. El test consistirá en 3 subtests:

- OK RAG: Porcentaje de preguntas en las que el documento/s *best* se encuentra entre los encontrados por la búsqueda semántica.
- OK FUENTES: Porcentaje de preguntas en las que el documento/s *best* se encuentra entre los elegidos y referenciados por el *LLM* para redactar la respuesta.
- OK RESPUESTA: Porcentaje de preguntas en las que un segundo *LLM* valora si la respuesta del primero se ajusta a la respuesta *best*. El *LLM* elegido para el test de *LLM as a judge* ha sido el *Llama 3.3 70B versatile*, ya que es un *LLM* grande, la última versión de los modelos *LLama* (hasta el

lanzamiento del *Llama 4*) y apto para gran variedad de tareas.

(Script en Test Automático.ipynb, Google Colab)

El *system prompt* es el siguiente:

Vas a recibir una pregunta de un usuario (Pregunta), la respuesta correcta a esta pregunta (Respuesta_Best) y una respuesta generada (Respuesta_Generada). Tus objetivos son los siguientes: 1. Determinar si la Respuesta_Generada responde a la Pregunta. 2. Determinar si la Respuesta_Generada concuerda con la Respuesta_Best y no la contradice. Tienes que hacer una valoración en detalle y razonando sobre tu valoración. Si la Respuesta_Generada no responde a la Pregunta, valorar como KO. Si la Respuesta_Generada concuerda con la Respuesta_Best y no la contradice, valorarla como OK. Si la Respuesta_Generada no concuerda con la Respuesta_Best y la contradice, valorarla como KO. No tener en cuenta posibles detalles adicionales que puedan estar incluidos en la Respuesta_Generada, siempre y cuando no contradigan la Respuesta_Best.

Una vez hayas hecho tu razonamiento, cuéntalo, y al final pon tu valoración en este formato: [[Valoración: OK/KO]]

El *user prompt* es el siguiente:

Pregunta:

{Pregunta del usuario}

Respuesta_Best:

{Respuesta marcada como correcta}

Respuesta_Generada:

{Respuesta generada por el LLM}

Las preguntas, documento/s *best* y respuestas *best* se pueden ver aquí: *Input Test Automático.xlsx*. El registro de todos los tests se puede ver aquí: PONER

III.I. *chunk size, overlap y top n*

El primer objetivo de este test es determinar el tamaño óptimo de los *chunks* de la base de datos. Para eso se ha ejecutado el test con varios *chunk size* distintos, así como para varios *top n* (el número de *chunks* pasados al LLM para redactar la respuesta). En este caso se ha fijado el LLM de elaboración de la respuesta y varios de sus parámetros. El LLM ha sido *Google Gemini 2.0 pro experimental*, aunque más adelante se comaprarán varios modelos y se elegirá el mejor. La temperatura se ha fijado a 0 para aumentar la reproducibilidad de los tests y reducir su variabilidad. Además, se ha fijado la *repetition penalty* a 0 para intentar producir respuestas concisas. Por último, el *chunk overlap* se ha fijado por defecto al 20%.

Se han hecho estas pruebas para un *chunk size* de 50, 75, 100 y 200 tokens, y *top n* de 3, 5, 10 y 15 *chunks*.

Chunk_size (tokens)	top_n (chunks)	OK RAG (%)	OK FUENTES (%)	OK RESPUESTA (%)
50	3	68,38	55,15	47,06
	5	78,68	61,76	55,15
	10	87,50	72,79	66,18
	15	90,44	74,26	69,85
75	3	63,97	48,53	42,22
	5	72,79	55,88	55,88
	10	83,09	63,97	63,97
	15	86,76	69,12	69,85
100	3	66,18	51,47	44,85
	5	75,74	60,29	52,94
	10	84,56	66,91	66,18

Chunk_size (tokens)	top_n (chunks)	OK RAG (%)	OK FUENTES (%)	OK RESPUESTA (%)
200	15	88,24	69,85	72,79
	3	61,03	47,06	44,12
	5	65,44	53,68	52,21
	10	74,26	53,68	56,62
	15	79,41	58,09	59,56

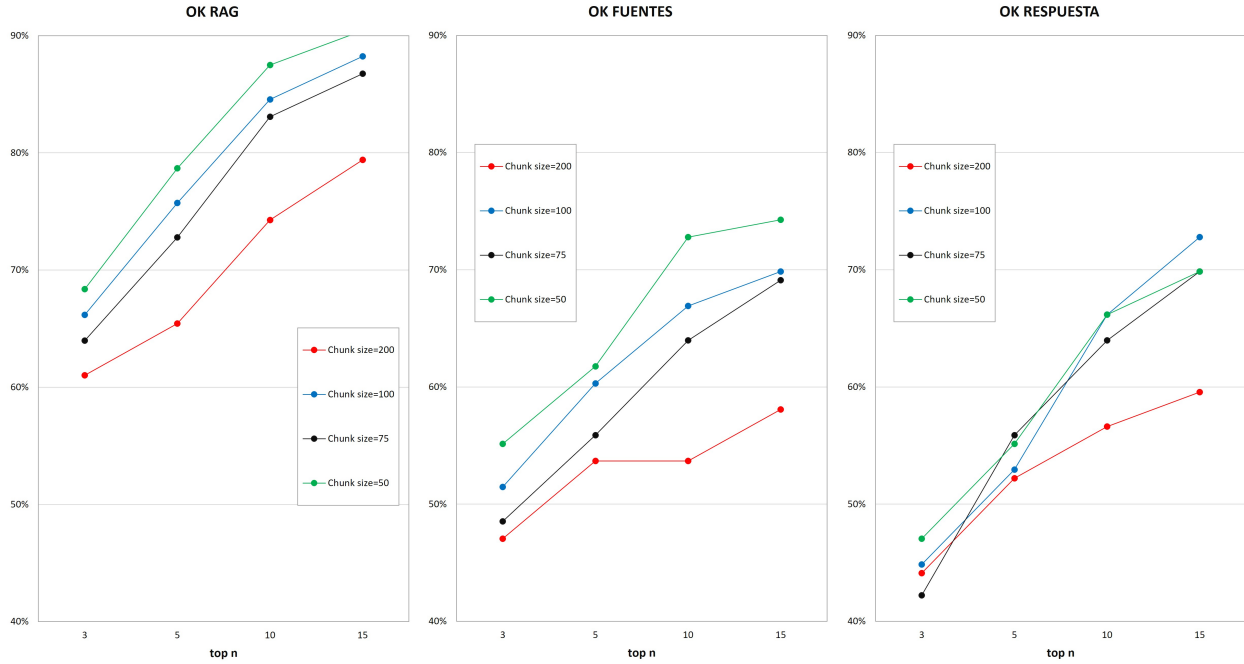


Figure 2: Valor de las métricas OR RAG, FUENTES y RESPUESTA en función del *chunk size* y *top n*, con *chunk overlap* del 20% del *chunk size*, modelo de *embeddings paraphrase-MiniLM-L6-v2* y modelo de elaboración de la respuesta *Google-Gemini 2.0 pro exp 02/05* (temperatura y *repetition penalty* a 0).

A raíz de estos resultados se utilizará un *chunk size* de 100 *tokens* (por tanto, un *chunk overlap* de 20 *tokens*) y un *top n* de 10 *chunks*. Se escoge esto por varias razones. En primer lugar, es el que mayor porcentaje de OK arroja en OK de la respuesta para *top n* de 10, junto a un *chunk size* de 50. Se elige sobre este porque, para resultados iguales, un *chunk size* de 100 ofrece más contexto. No se escoge *chunk size* de 100 y *top n* de 15 porque la mejora en esta métrica no es demasiada. Además, hay que tener en cuenta el tiempo de respuesta del modelo, una métrica que en este caso no se ha evaluado, pero que es fundamental, ya que esta aplicación de IA es un *chatbot*. En resumen, los parámetros elegidos han sido ***chunk size* = 100 *tokens*, *chunk overlap* = 20 *tokens*, *top n* = 10 *chunks***.

Por otro lado, como más adelante se va a analizar el *LLM* a utilizar para responder a las preguntas, así como su temperatura, cabría preguntarse si este análisis fijando el *LLM* es válido para otros. La realidad es que no, pero se ha hecho así para reducir el número de pruebas a hacer. Aunque los resultados del análisis del *chunk size*, *overlap* y *top n* probablemente cambien de un *LLM* a otro, se ha supuesto que no serán cambios significativos.

III.II. Modelo de *embeddings*

En el momento que se hicieron los tests III.I pensaba que el modelo de *embeddings* que estaba utilizando, *paraphrase-MiniLM-L6-v2*, era el mejor. Sin embargo, en una de las clases me hicieron saber que este

modelo no está entrenado en español, por lo que es fundamental encontrar uno que funcione mejor entrenado específicamente en español.

Lo ideal sería hacer primero este test y después el III.I, ya que es más determinante el modelo de *embeddings* utilizado. Sin embargo, como el test III.I consume mucho tiempo, asumiremos el error producido por hacerlo en este orden.

Los modelos de *embeddings* que vamos a comparar son los siguientes:

- *paraphrase-MiniLM-L6-v2* (as-is, entrenado solo en inglés)
- *paraphrase-multilingual-MiniLM-L12-v2*
- *paraphrase-multilingual-mpnet-base-v2*
- *distiluse-base-multilingual-cased-v2*
- *stsb-xlm-r-multilingual*
- *Shaharyar6/finetuned_sentence_similarity_spanish*

Por desgracia, y debido a las limitaciones de tener que usar llamadas gratis via API a *LLMs* ofrecidos por distintos proveedores, el *LLM* de elaboración de la respuesta que usé en el test III.I, *Google: Gemini Pro 2.0 Experimental (free)*, ya no está disponible. Debido a esto voy a tener que utilizar otro, *Google: Gemini 2.0 Flash Thinking Experimental 01-21 (free)*. Aún así, esto no invalida las conclusiones de este test.

Los resultados de este test son los siguientes:

Modelo de embeddings	OK RAG (%)	OK FUENTES (%)	OK RESPUESTA (%)
<i>paraphrase-MiniLM-L6-v2</i> (inglés)	84,56	66,18	61,03
<i>paraphrase-multilingual-MiniLM-L12-v2</i>	88,97	77,21	80,74
<i>paraphrase-multilingual-mpnet-base-v2</i>	89,71	76,47	78,68
<i>distiluse-base-multilingual-cased-v2</i>	88,24	77,94	82,35
<i>stsb-xlm-r-multilingual</i>	67,65	47,06	56,62
<i>Shaharyar6/finetuned_sentence_similarity_spanish</i>	92,65	82,35	83,82

Podemos observar que, aunque OK RAG del modelo entrenado en inglés es similar al resto (entrenados en español), en las otras dos métricas es muy inferior. Es decir, los modelos de *embeddings* entrenados en español están encontrando chunks mucho más útiles para elaborar la respuesta.

Que haya más OK RESPUESTA que OK FUENTES se explica por la naturaleza de la base de datos: como es una *wikipedia* tiene mucha información redundante, así que es probable que para alguna pregunta haya más documento/s *best* de los que he puesto en el test automático.

Sin embargo, lo que más llama la atención son los pésimos resultados del modelo *stsb-xlm-r-multilingual*. EXPLICAR ESTO

Como los resultados de los modelos entrenados en español (quitando el mencionado anteriormente) son similares, vamos a elegir basándonos también en el tiempo de inferencia de cada modelo. Aunque no está implementado el cálculo del tiempo de ejecución de cada pregunta del test, podemos estimar el tiempo de otra forma. Como el número de *chunks* y *tokens* de la base de datos es fijo, podemos medir el tiempo que tarda cada modelo en crear el índice. Esto nos dará una idea de qué modelos tardan menos en hacer una inferencia.

Modelo de embeddings	Tiempo de creación del índice (s)
<i>paraphrase-MiniLM-L6-v2</i>	85
<i>paraphrase-multilingual-MiniLM-L12-v2</i>	154
<i>paraphrase-multilingual-mpnet-base-v2</i>	531
<i>distiluse-base-multilingual-cased-v2</i>	260
<i>stsb-xlm-r-multilingual</i>	365

Modelo de embeddings	Tiempo de creación del índice (s)
Shaharyar6/finetuned_sentence_similarity_spanish	447

Con estos resultados, el modelo de *embeddings* que vamos a utilizar es ***distiluse-base-multilingual-cased-v2***, el segundo que mejor OK RESPUESTA da y el tercero más rápido, además del más rápido de los modelos por encima del 70% de OK RESPUESTA.

III.III. LLM de generación de la respuesta

A continuación, se evaluará qué *LLM* se utilizará para generar la respuesta a partir del conocimiento encontrado.

Los modelos que vamos a evaluar son:

- google/gemini-2.0-flash-thinking-exp:free
- google/gemini-2.5-pro-exp-03-25:free
- meta-llama/llama-3.3-70b-instruct:free
- deepseek/deepseek-chat:free (V3)
- deepseek/deepseek-r1-zero:free
- qwen/qwen-2.5-72b-instruct:free
- microsoft/phi-3-medium-128k-instruct:free
- mistralai/mistral-7b-instruct:free
- gpt-4o-mini
- meta-llama/llama-4-maverick:free
- meta-llama/llama-4-scout:free

Los resultados han sido los siguientes (recordemos que la métrica OK RAG tiene un valor de 88,24%, que no cambiará con el *LLM* de elaboración de la respuesta, por lo que no incluiremos esta métrica aquí):

Modelo de respuesta	OK FUENTES (%)	OK RESPUESTA (%)
google/gemini-2.0-flash-thinking-exp:free	77,94	82,35
google/gemini-2.5-pro-exp-03-25:free	80,88	83,09
meta-llama/llama-3.3-70b-instruct:free	82,35	80,15
deepseek/deepseek-chat:free (V3)	79,41	82,35
deepseek/deepseek-r1-zero:free	75,00	83,09
qwen/qwen-2.5-72b-instruct:free	77,21	79,41
microsoft/phi-3-medium-128k-instruct:free	0,00	66,91
mistralai/mistral-7b-instruct:free	73,53	73,53
gpt-4o-mini	79,41	76,47
meta-llama/llama-4-maverick:free	86,76	83,82
meta-llama/llama-4-scout:free	83,09	77,94

A la vista de estos resultados se aprecia que el mejor modelo de elaboración de la respuesta es ***meta-llama/llama-4-maverick:free***, de reciente lanzamiento el 5/04/2025. Se puede apreciar que en la métrica OK RESPUESTA está cerca de otros, como los *gemini* o *deepseek*, pero parece que hay más diferencia al referenciar las fuentes.

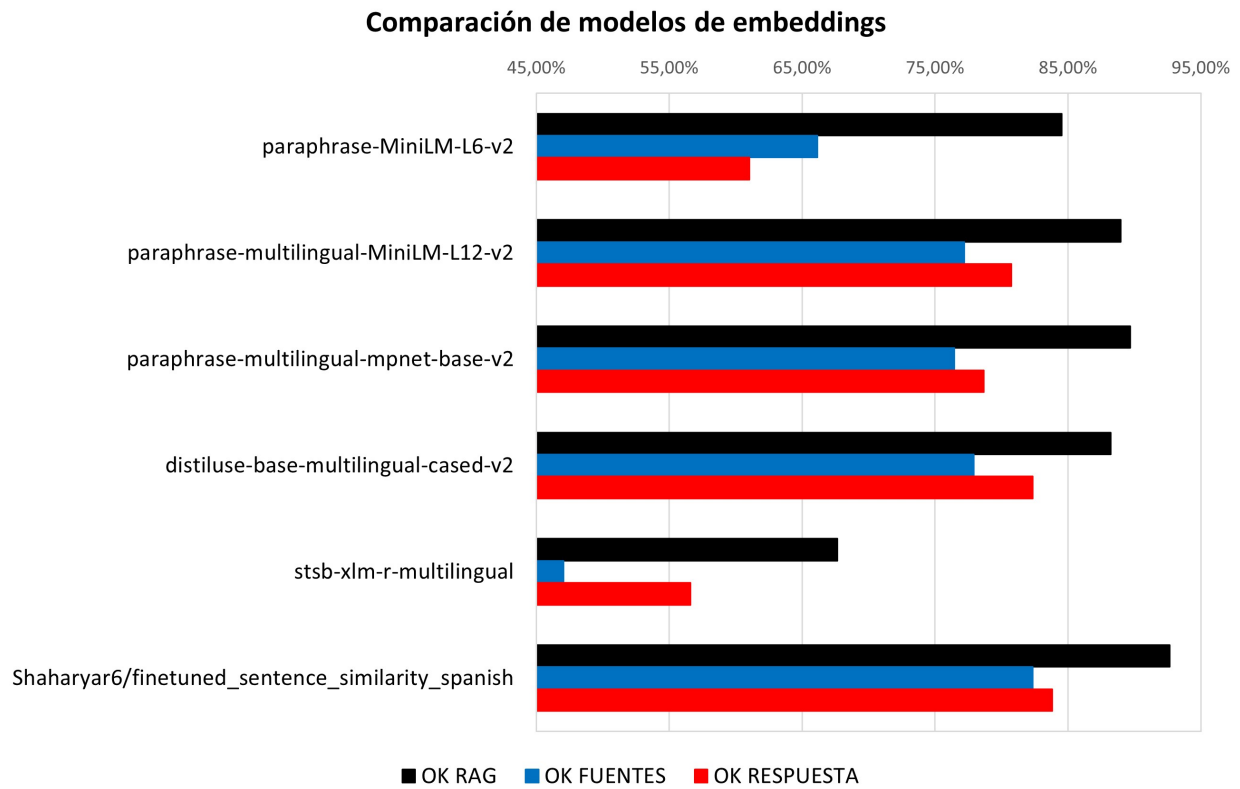


Figure 3: Valor de las métricas OR RAG, FUENTES y RESPUESTA en función del modelo de *embeddings*, con *chunk size* de 100 *tokens*, *chunk overlap* de 20 *tokens*, *top n* de 10 *chunks* y modelo de elaboración de la respuesta *Google: Gemini 2.0 Flash Thinking Experimental 01-21* (temperatura y *repetition penalty* a 0).

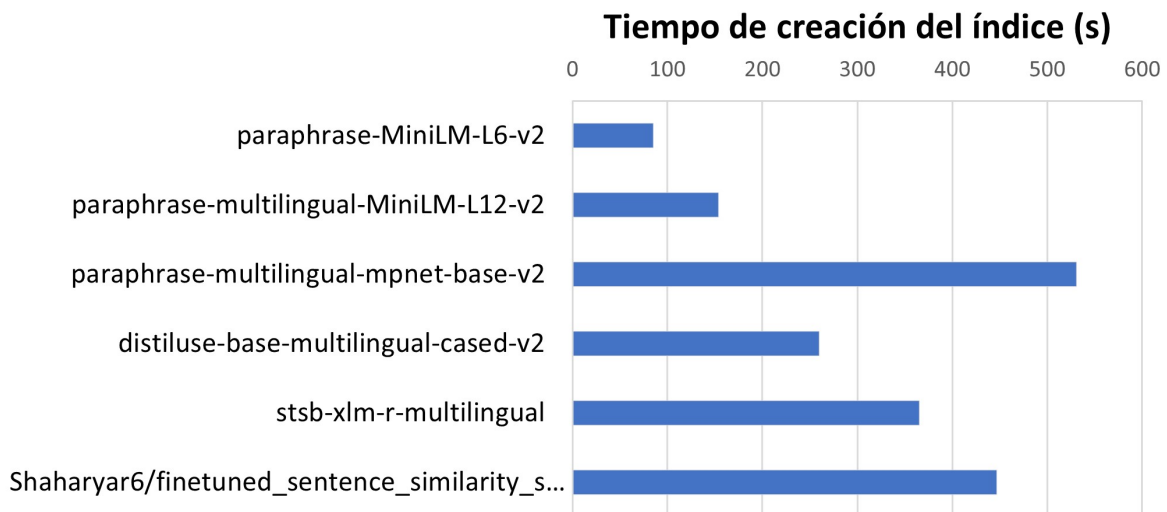


Figure 4: Tiempo de creación del índice de la base de documentos en función del modelo de *embeddings*, con *chunk size* y *overlap* de 100 y 20 *tokens* respectivamente.

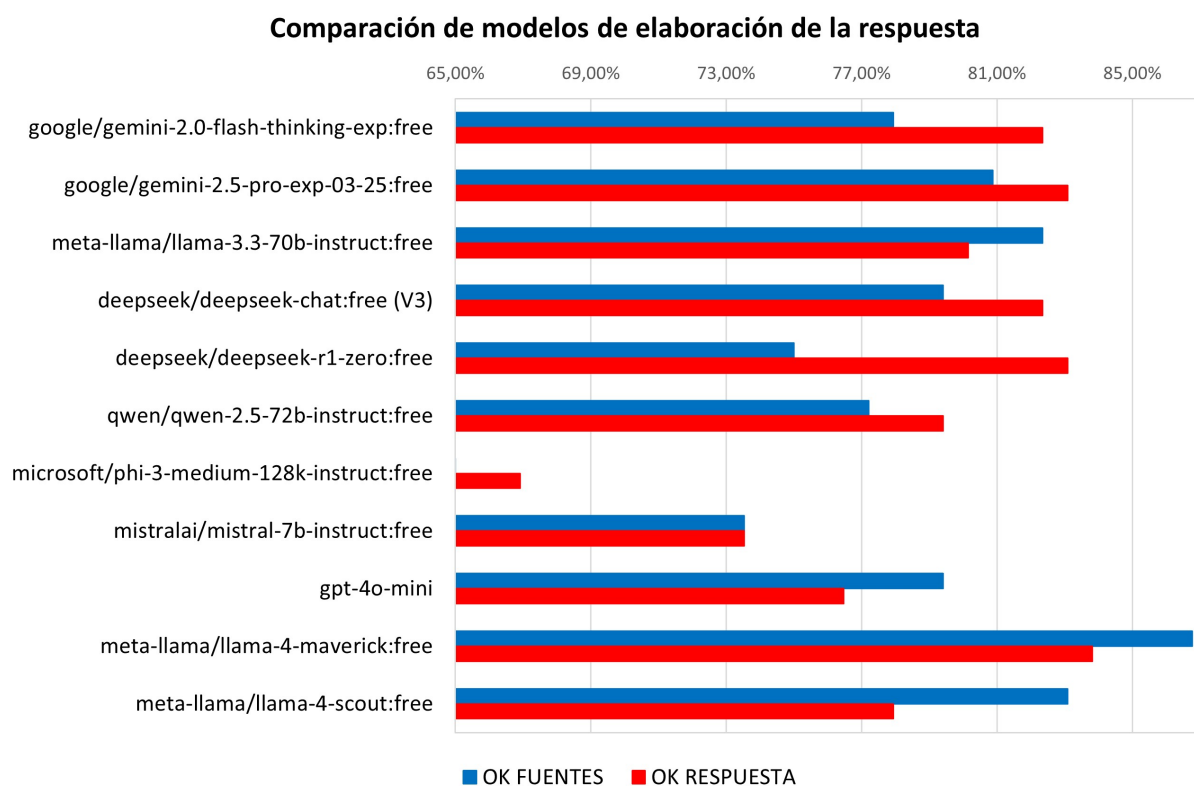


Figure 5: Valor de las métricas OK FUENTES y RESPUESTA en función del modelo de elaboración de la respuesta (temperatura y *repetition penalty* a 0), con *distiluse-base-multilingual-cased-v2* como modelo de *embeddings*, *chunk size* de 100 *tokens*, *chunk overlap* de 20 *tokens* y *top n* de 10 *chunks*.

Puesta en producción

Conclusiones

Posibles mejoras

En esta sección voy a detallar las posibles mejoras que hacer a este trabajo. Son ideas que han salido durante la realización del mismo o gracias a las clases recibidas, que no se alejan demasiado de los objetivos del trabajo, pero que o bien son demasiado ambiciosas o bien no ha dado tiempo a hacerlas.

1. *LLMs* de pago

Si ha habido algo que haya lastrado este trabajo es la limitación que teníamos de utilizar llamadas gratis via *API* a *LLMs* ofrecidos por distintos proveedores. Esto ha hecho que estemos restringidos en cuanto a los modelos que utilizar, tengamos que crear varias cuentas por proveedor para poder sortear esos *rate limit*, no podamos automatizar del todo los test automáticos, etc. Además, los mejores *LLMs* no se ofrecen gratis, por lo que contratar alguno puede aumentar también el desempeño del *RAG*, además de probablemente reducir los tiempos de inferencia.

Usar *LLMs* de pago me permitiría además aplicar *Structured Outputs* para que en la elaboración de la respuesta y el *LLM as a judge* den, respectivamente, las referencias separadas de la respuesta y la valoración del razonamiento.

Por otro lado, utilizar el modelo de *embeddings* de *OpenAI* mejoraría significativamente los resultados de la parte del *retrieval*, ya que es un modelo grande pero que no necesita ser alojado en local, por lo que además es rápido. En clases y prácticas anteriores hemos discutido y comprobado la mejora notable por utilizar este modelo.

2. *Prompts*

En cuanto a los *prompts*, el prompt que utilizo en este trabajo, tanto para la elaboración de la respuesta como para la parte del *LLM as a judge* del test automático, son los que han funcionado para el modelo *Google: Gemini Pro 2.0 Experimental (free)*. Que funcionaran bien con ese modelo no implica que funcionen también con el resto, por lo que una posible mejora podría ser encontrar el *prompt* ideal para cada *LLM* utilizado. Además, se podría haber aplicado la técnica *Few-Shot Prompting* para incluir algún ejemplo de cómo elaborar la respuesta y referenciar los documentos adecuados.

Otra cosa que me gustaría haber hecho mejor es la gestión de los prompts. En el repositorio del código están incluidas en un *script* de *python*, pero deben poder guardarse y tratar las versiones con alguna herramienta externa que sea más idónea.

3. Test Automático

Una mejora clara en esta parte es incluir el tiempo de inferencia medio de cada test. Esta es una métrica fundamental para soluciones tipo *RAG*. Esta es además una métrica que, de poner esta solución en producción, mejoraría mucho, ya que los *LLMs* y modelos de *embeddings* de pago son mucho más rápidos.

Por otro lado, el test puede no ser todo lo representativo que pretende, ya que el número de preguntas de cada documento no lo he decidido de forma rigurosa. Podría hacerse que el porcentaje de preguntas sobre cada documento dependa de la longitud de cada documento de la base de datos. Además, por supuesto, las preguntas pueden estar hechas de forma sesgada, ya que he sido yo mismo el que las ha diseñado. Una manera de hacerlo menos sesgado es quizás pasarle fragmentos del documento a un *LLM* y pedirle que elaborara una pregunta que fuera respondida con algo de ese fragmento.

A su vez, los resultados de cada test se han enviado a un *Excel*, cosa que no es ni muy limpia ni muy escalable. Lo que podría hacerse es enviar los resultados y parámetros de cada test a *MLflow*.

4. Métricas de *Ragas* y *Groundedness* de *Microsoft*

Una mejora que sería muy buena es utilizar *Ragas*, una herramienta de código abierto diseñada para evaluar sistemas *RAG*. En la carpeta de Apoyo deo un notebook donde hago un análisis de las distintas métricas que ofrece, además de poder usarse de soporte para elaborar métricas propias.

Además de estas métricas, podría calcularse la métrica *Groundedness* de *Microsoft*. Esta métrica mide lo desviada que está la respuesta de un sistema *RAG* respecto del contexto que se le pasa. Es al fin y al cabo una forma de medir las alucinaciones, o lo que añade el *LLM* que elabora la respuesta al contexto recibido. Esta es una métrica que se calcula via *API* y en la que no se utiliza un *LLM* para calcularla, por lo que es muy rápida (unos 300 ms). Es por esto que podría incluso utilizarse para avisar al usuario de lo fiable que puede ser la respuesta, pintándola por ejemplo en una escala de color del rojo al verde.

5. Llamadas a *LLMs* con *LiteLLM*

En el código de este trabajo hago las llamadas a los *LLMs* de los distintos proveedores de forma algo sucia; cada uno necesita una estructura diferente. *LiteLLM* es una librería de código abierto que actúa como una interfaz unificada para hacer estas llamadas, de forma que lo hace mucho más escalable (es más fácil añadir otros proveedores) y limpia. En la carpeta de Apoyo deo un pequeño tutorial de cómo se haría.

Esto finalmente me ha dado tiempo a incluirlo.

Líneas a futuro

En esta sección voy a detallar las líneas a futuro que surgen de este trabajo. Son ideas que han salido durante la realización del mismo o gracias a las clases recibidas, como las de la sección de posibles mejoras, pero estas suponen cambios grandes en el funcionamiento de la solución actual y podrían formar parte de un nuevo trabajo por sí mismas.

1. Implementación de un *chatbot*

Actualmente este *RAG* funciona sin poder hacer varias iteraciones sobre las mismas preguntas; únicamente recibe una *query* de entrada a la que da respuesta. En un primer momento se pretendía que se pudiera *chatear*, pero esto supone un gran cambio en la lógica, como el manejo del historial de la conversación, una búsqueda más dinámica en el *RAG*, etc. Una de las piezas nuevas que habría que pensar en incluir es un detector de *chit-chat* que, si uno de los mensajes del usuario no tiene intención de realizar una búsqueda en el *RAG*, sepa manejarlo y no haga la búsqueda.

2. Implementación de un agente

La solución actual es muy rígida por construcción: para cada *query* obtiene un número fijo de *chunks*, con los que tiene que elaborar la respuesta. Si ahí no está la respuesta no puede hacerse nada. Todos los hiperparámetros de esta solución son fijos, aunque puedan no ser los más idóneos para alguna pregunta. Un ejemplo de pregunta para la que esta solución no es idónea es si se le pidiera hacer una lista de características de varios elementos: lo ideal sería que buscara cada elemento por separado, pero tal y como está planteada la solución actualmente hace una única búsqueda, pudiendo no encontrar toda la información.

En este sentido, sería muy interesante elaborar una solución en la que sea un agente el que gestione las búsquedas en la base de datos y la elaboración de la respuesta. De esta forma, el agente puede adaptar la *query* del usuario y hacer búsquedas con *queries* más adecuadas, no hacer solamente una, decidir si tiene que hacer más porque aún no tiene la respuesta o modificar hiperparámetros como *top_n*.

3. Mejoras en las *queries* del *RAG*

Una técnica muy habitual en *RAGs* productivos es *RAG-fusión*. *RAG-fusión* consiste en, con la ayuda de un *LLM*, generar más *queries* a partir de la inicial, hacer varias búsquedas para cada *query*, reordenar los

chunks encontrados y generar la respuesta con todos ellos. Esta solución aporta mayor cobertura, precisión y respuestas más robustas a las soluciones de *RAG*.

4. Mejoras en la creación del índice

Una de las limitaciones de este *RAG* es que en la creación del índice se hacen *embeddings* directamente sobre el contenido de la base de datos, pero este contenido no tiene por qué ser muy similar a las preguntas que hacen los usuarios. Otra opción es hacer los *embeddings* sobre las preguntas que podría responder cada *chunk* de la base de datos, generando estas preguntas con la ayuda de un *LLM*.

Otra de las posibles mejoras en la creación del índice de este tipo de soluciones es incluir un pequeño resumen de cada documento como metadato de los *chunks* que lo conforman. De esta forma, el *LLM* que elabora la respuesta tiene más contexto y puede responder mejor. Estos resúmenes pueden crearse con otro *LLM*.