

Paquete 5 – Procesamiento de Lenguaje Natural y Grandes Modelos de Lenguaje

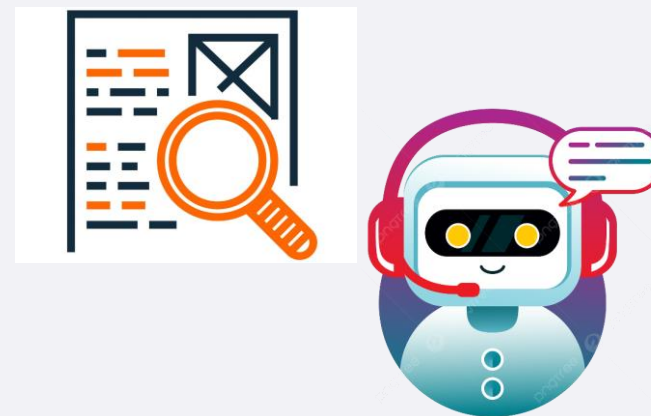
Introducción

Grupo de Investigación INTELIA: <https://intelia.uah.es/>

- **Antonio García Cabot**, Eva García López,
José Manuel Lanza Gutiérrez, Ana Castillo Martínez, Sergio Caro Álvaro
- 2 Investigadores Contratados (Estudiantes de doctorado)

El Paquete de Trabajo PT5 – Procesamiento de Lenguaje Natural (NLP por sus siglas en inglés de Natural Language Processing) en Sistemas Aeronáuticos y Aeroespaciales tiene como objetivo desarrollar diversas actividades y acciones en materia de I+D relacionadas a la aplicación de NLP en los diferentes procesos relacionados con el ámbito aeronáutico y aeroespacial:

- Creación de chatbots
- Análisis de sentimientos
- Extracción de información



Actividades de formación interna y externa

Dentro de las actividades de formación de este paquete de trabajo, se están llevando a cabo dos tesis doctorales:

- Título "Advancing Evaluation and Reliability Methods for Large Language Models". Doctorando: Pablo Trull Báguena
- Título "Estudio y optimización de los métodos de generación aumentada por recuperación basada en grandes modelos del lenguaje". Doctorando: Antonio Moreno Cediel

Realización de 6 Trabajos Fin de Máster (TFM):

- Detección automática de técnicas de persuasión o manipulación
- Análisis de sentimientos dirigido (Targeted sentiment analysis)
- Análisis de sentimiento de los comentarios de aerolíneas
- Comparación de técnicas de optimización de modelos aplicadas al análisis de sentimientos
- Análisis de opiniones en Reddit usando Topic Modeling y Análisis de sentimientos
- Detección de movimientos relacionados con estrategias de desinformación en redes sociales

Distintas propuestas de Trabajos Fin de Grado (TFG)

Estudio del estado del arte y datasets

Durante los primeros meses de trabajo se ha llevado a cabo un estudio exhaustivo del estado del arte en relación con la Generación Aumentada por Recuperación (RAG) y el Aprendizaje por refuerzo (RL).

RAG

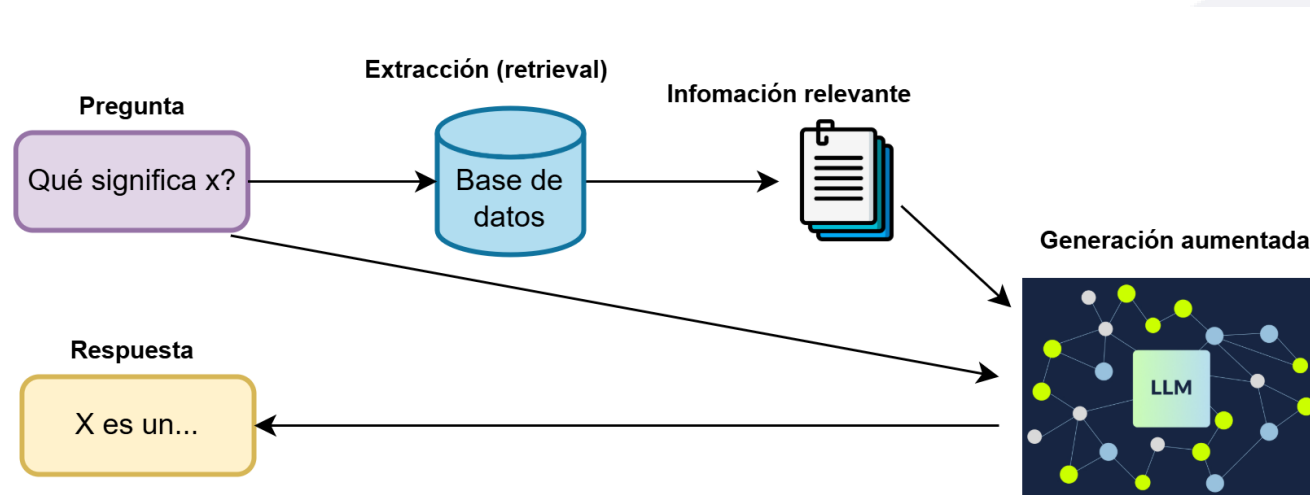


Figura 1. Descripción de la arquitectura de RAG

Estudio del estado del arte y datasets

RAG

Identificando las principales fases del proceso de RAG:

- Indexación: estrategias de segmentación.
- Recuperación: tipos, diferentes fuentes de datos, optimización.
- Generación: diferentes arquitecturas.

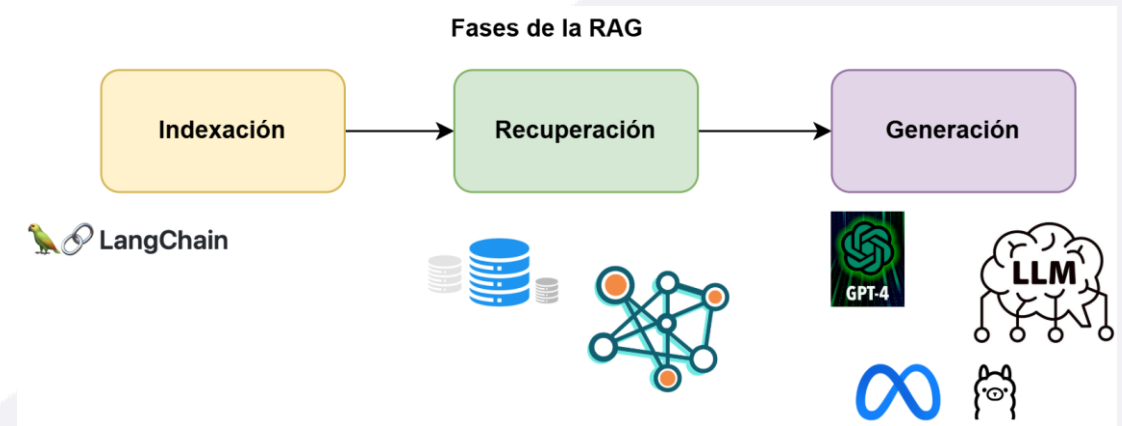


Figura 2. Fases de la RAG

Estudio del estado del arte y datasets

RL (Reinforcement Learning)

Las principales formas de utilizar de forma conjunta los LLMs y el RL:

- RL4LLM: RL para mejorar LLM en tareas de NLP
- LLM4RL: LLM para mejorar RL
- LLM+RL: LLM+RL

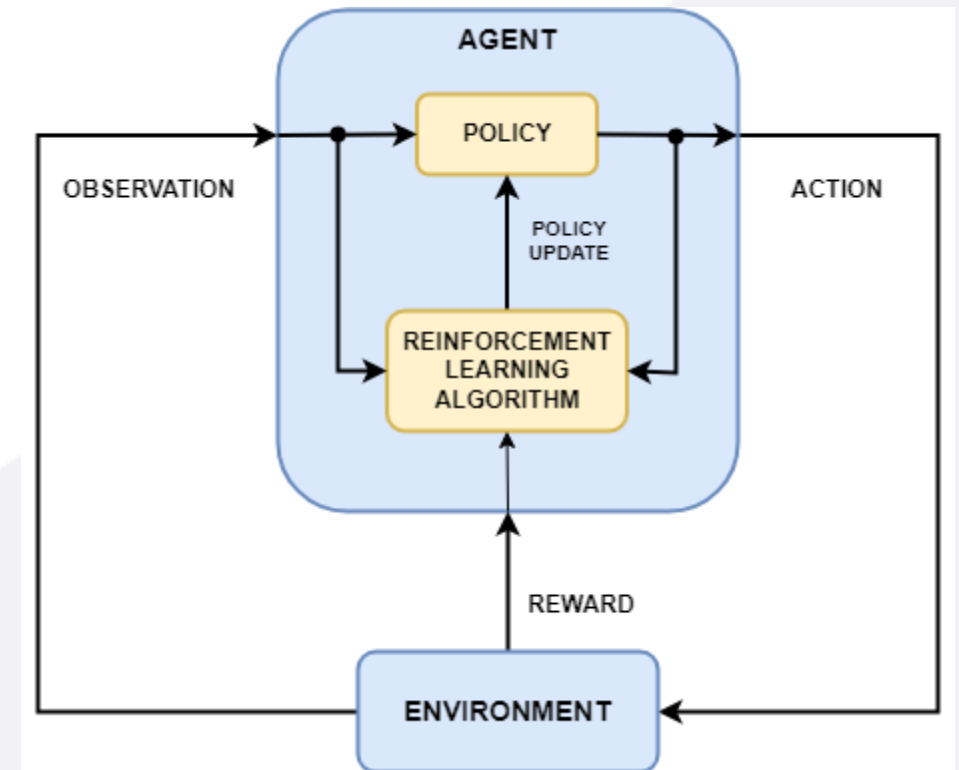


Figura 3. Diagrama de la interacción entre el agente y el entorno

Creación de los modelos de IA y ajuste

Topic Modeling

- Optimización de topic models basados en clustering a través de *content augmentation*, realizado mediante *zero-shot prompting LLMs*.

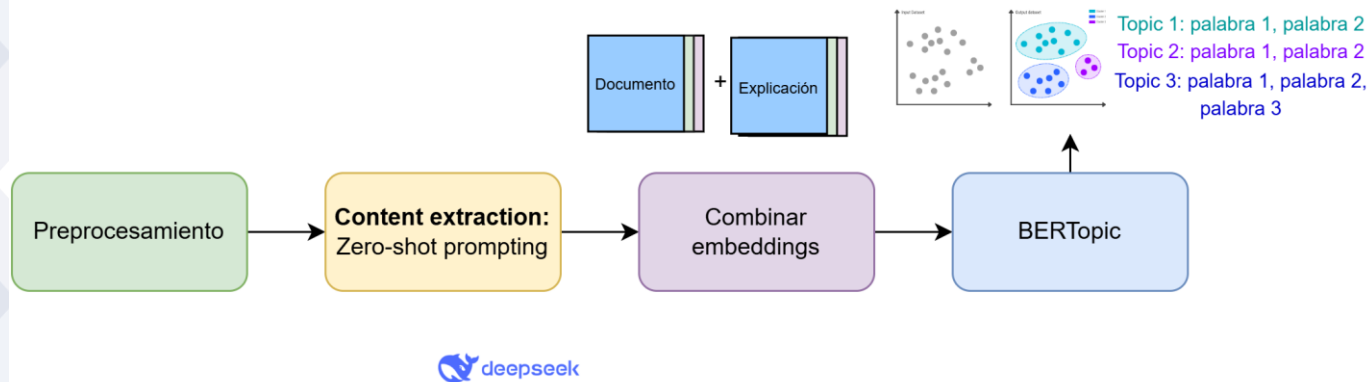


Figura 4. Metodología propuesta

- 1.- Preprocesamiento:** se eliminan los enlaces, puntuación excesiva, etc.
- 2.- Content extraction:** Se extraen explicaciones de cada texto del dataset con LLMs (Deepseek y Llama3)
- 3.- Combinar embeddings:** Se computan embeddings de los textos originales y de las explicaciones y se combinan
- 4.- BERTtopic:** los embeddings de la fase anterior se pasan a BERTopic y este los agrupa en clusters temáticos y extrae las palabras más representativas de cada grupo

Creación de los modelos de IA y ajuste

Generación de queries de *SPARQL*

- Se propone una nueva metodología, basada en *few-shot prompting*, para generar *queries* de *SPARQL* dada una consulta en lenguaje natural.

A partir de una pregunta en lenguaje natural (representada en azul en la Fig. 5), nuestro método recupera las tripletas relevantes para contestar esta pregunta y combina esta información en un prompt que se pasa a Llama3-70b y este genera una query en sparql.

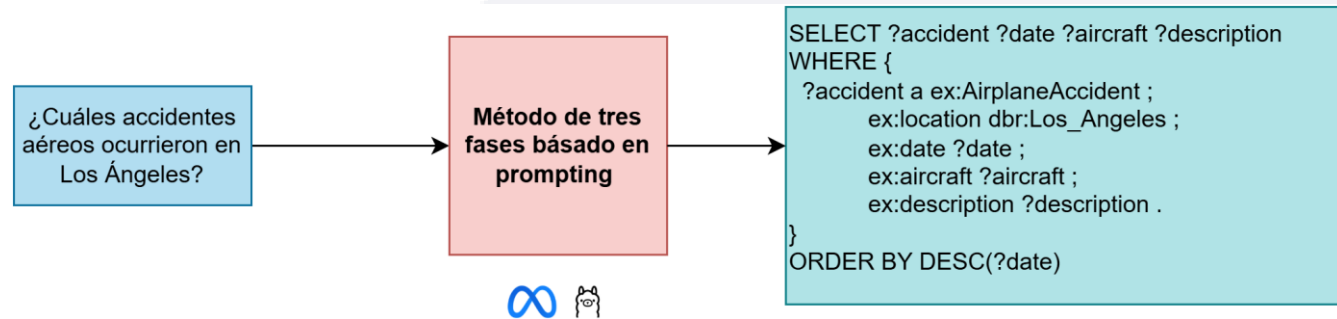


Figura 5. Ejemplo de la tarea de generación de queries

Desarrollo y ejecución de PoC

Topic Modeling

Para evaluar la técnica propuesta se emplearon 2 datasets ampliamente utilizados: 20NewsGroup y BBCNews.

	Modelo	C_v	C_{NPMI}	diversidad	ARS
20Newsgr oup	BERTopic (baseline)	0.60	0.11	0.81	0.40
	<i>Content augmentation</i> (nuestro método)	0.68	0.20	0.81	0.46
BBCNews	BERTopic (baseline)	0.53	0.07	0.83	0.89
	<i>Content augmentation</i> (nuestro método)	0.64	0.17	0.82	0.87

20NewsGroup: Se mejora la coherencia C_v un 8%, la coherencia C_{NPMI} un 9% y el ARS un 6%. La diversidad se mantiene igual en general

BBCNews: se mejora la coherencia C_v un 9% y la coherencia C_{NPMI} un 10%

Tabla 1. Resumen de resultados del estudio de Topic Modeling

Desarrollo y ejecución de PoC

Generación de queries de SPARQL

- Se empleó un grafo de conocimientos del ámbito de aviación. El trabajo donde se introdujo este grafo también introdujo un dataset de evaluación que consistía en 150 preguntas.
- Nuestro método mejoró la métrica de *Exact Match* un 6%.

Los hallazgos de este estudio y del estudio de *Topic Modeling* han sido recopilados en **dos artículos científicos** que han sido publicados en dos revistas científicas:

- <https://link.springer.com/10.1007/s10115-025-02605-0>
- <https://www.mdpi.com/2504-4990/7/2/52>

Otro artículo enviado a conferencia en 2026 que se organizará en Hakodate, Japón.

Actuaciones recientes

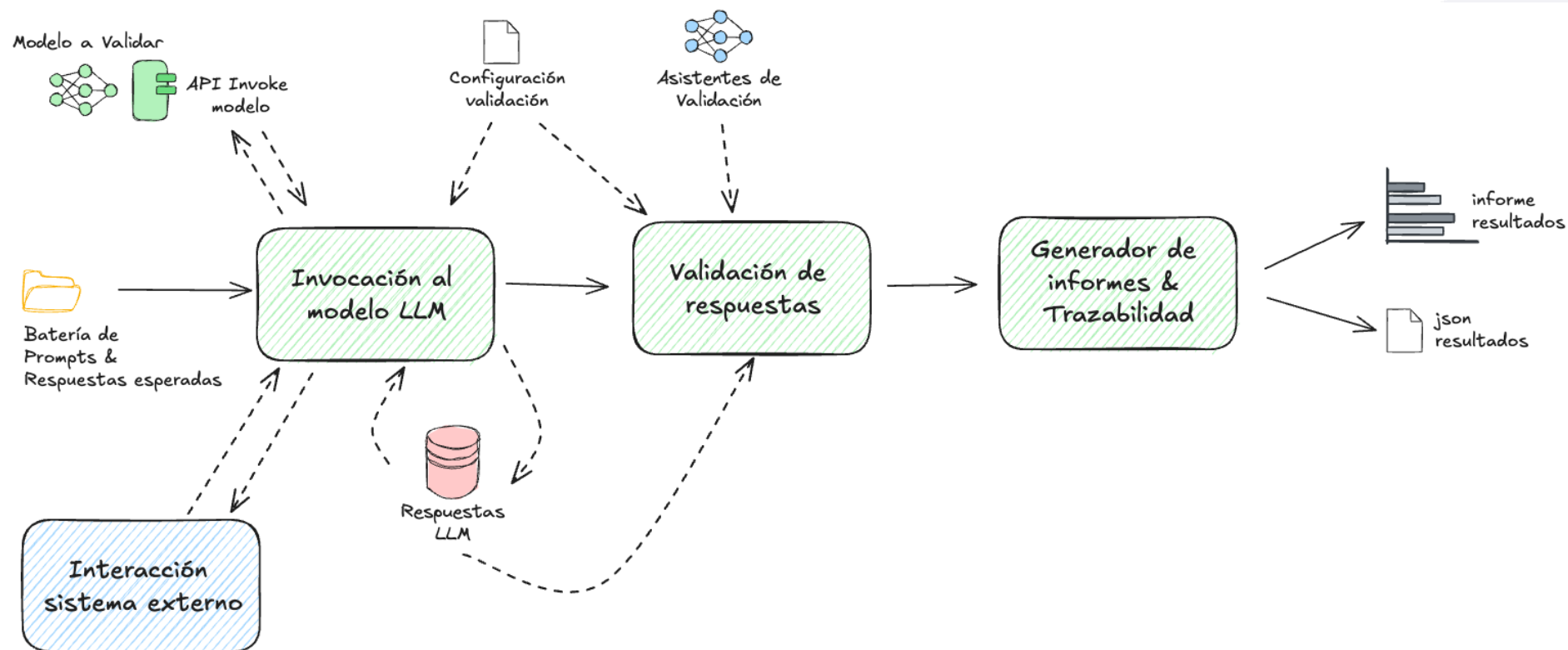
Tres líneas de actuación junto con **Indra**:

- Validación y experimentación de modelos LLM
- Generación de grafos de conocimiento
- Language Alignment by Bootstrapping (LAB)

Validación y experimentación de modelos LLM

- Validación mediante una herramienta que permita contrastar de forma sistemática las respuestas generadas por un modelo grande de lenguaje (LLM).
- Estado del arte
 - LLM-as-a-Judge, Incertidumbre, Robustez, ...
 - Métricas de evaluación (exactitud, medidas de precisión/recall, similitud semántica, métricas tradicionales en NLP como ROUGE, BLEU o METEOR)...
- Experimentación y prototipos

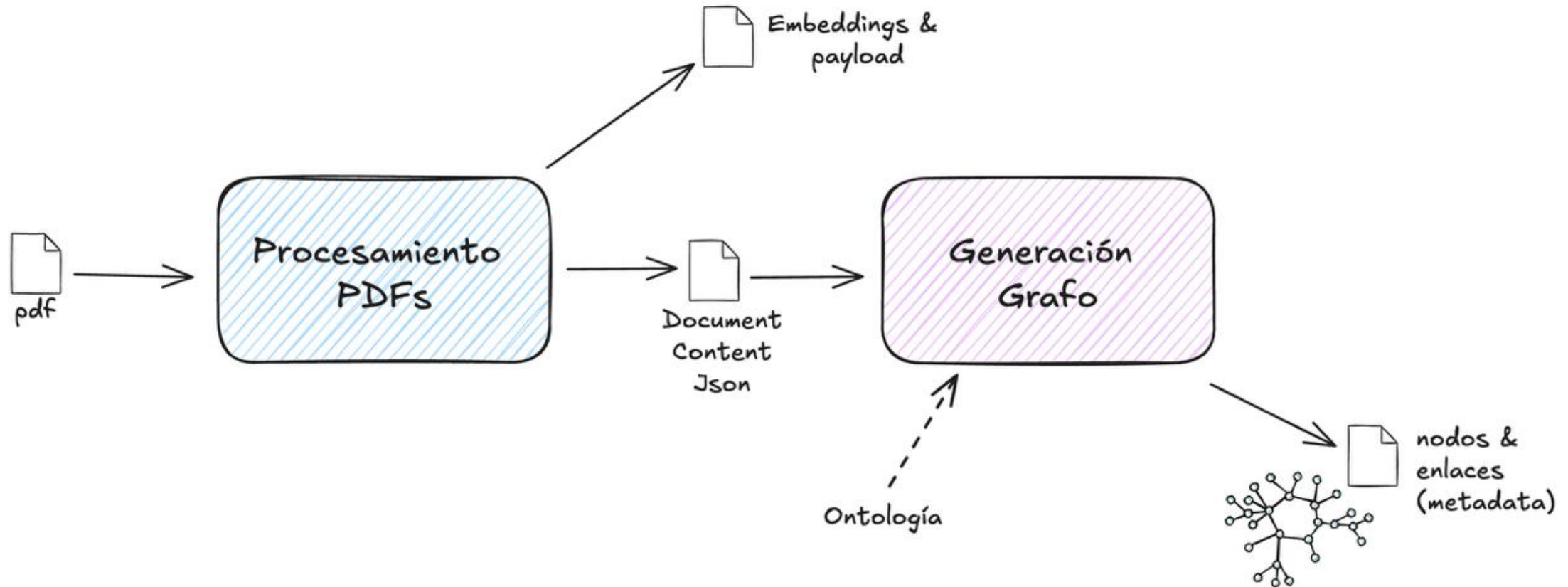
Validación y experimentación de modelos LLM



Generación de grafos de conocimiento

- El objetivo principal es poder almacenar en una estructura de tipo grafo, la información contenida en documentos expresados en lenguaje natural como PDF, MS Word, ...
- Capaz de trabajar tanto con una ontología objetivo con formato libre.
- Sistema que realiza un análisis orientado al análisis de cada frase intentando extraer de cada una de ellas los sintagmas nominales y verbales para extraer los nodos y enlaces del grafo.

Generación de grafos de conocimiento



Language Alignment by Bootstrapping (LAB)

- LAB es la estrategia para alinear modelos LLM con comportamientos contextualizados a determinados contextos (conocimiento específico), de manera más eficiente que los procesos de *fine-tuning* supervisado.
- Se propone utilizar InstructLab como parte de la colaboración con **Indra** y realizar experimentación y pruebas.

Otras iniciativas o investigaciones de interés

En el laboratorio de investigación se están llevando a cabo otras investigaciones relacionadas con el paquete de trabajo de la Cátedra IA3. Algunas de ellas son:

- Reparación automática de código fuente utilizando LLMs.
- Generación de datos sintéticos de calidad.
- Creación de chatbots y RAG para el acceso a documentación y bases de datos de conocimiento.
- Interpretabilidad y explicabilidad de LLM (Autoencoders y Transcoders).
- Optimización y reducción de modelos LLM.

¡Gracias!

- ¿Preguntas/Dudas?
- Contacto:
 - Dr. Antonio García Cabot (a.garciac@uah.es)
 - LinkedIn:



- Contacto genérico: intelia@uah.es
- Sitio web: <https://intelia.uah.es>

