
INSTITUTO TECNOLÓGICO DE MORELIA

"José María Morelos y Pavón"

Materia

Inteligencia Artificial

Unidad

III

PRESENTAN:

Víctor Daniel Guzmán Mondragón

21121529

PROFESOR:

Jesús Eduardo Alcaraz Chávez

MORELIA, MICHOACÁN, A MARTES 6 DE ENERO DE 2025.

Proyecto de Análisis de Datos mediante RAG

Introducción

El presente proyecto esta enfocado en realizar un análisis de datos mediante técnicas de RAG (Retrieval-Augmented-Generation) para explorar 2 problemas filosóficos que se nos brindó a continuación:

1. La posible crisis de sentido en la Generación Z debido a la hiperconectividad.
2. La pérdida gradual de autonomía humana frente al avance de los algoritmos y la inteligencia artificial.

Este proyecto integra análisis cualitativo y cuantitativo, minería de texto, recuperación semántica y modelos generativos.

1. Planteamiento General del Proyecto

Se utilizará un sistema RAG para:

Recuperar información desde bases de datos.

Analizar discursos y tendencias en redes sociales.

Procesar texto mediante embeddings.

Generar interpretaciones fundamentadas mediante IA.

El propósito central es comprender cómo los algoritmos moldean el sentido de vida, identidad y autonomía en la Generación Z.

2. Marco teórico

Las herramientas que se utilizaron durante este proyecto fueron las siguientes:

Lenguaje de Scripting: Python (para la extracción de datos).

Plataforma RAG: AnythingLLM (versión Desktop).

Motor de Inferencia: Ollama.

Modelo de Lenguaje (LLM): Qwen 2.5 (excelente en razonamiento y codificación.)

Modelo de Embeddings: Nomic Embed Text. (Transforma el texto en vectores numéricos para que la IA pueda buscar similitudes).

3. Metodología de implementación

Fase A: Ingeniería de Datos y Pipeline ETL

1. Descripción del Módulo

Para la construcción de la base de conocimiento del modelo RAG, se desarrolló un pipeline de extracción automatizada de datos (Web Scraping) diseñado para recopilar opiniones y discursos sociales no estructurados provenientes de YouTube. El objetivo fue generar un dataset (CSV) que complementara la teoría académica de los libros con "sentimiento social real" sobre la Generación Z.

2. Stack Tecnológico

Para esta fase se utilizaron las siguientes librerías de Python:

- **youtube_comment_downloader:** Seleccionada por su capacidad de iteración rápida sin las restricciones de cuota de la API oficial de Google, permitiendo una extracción masiva y eficiente.
- **pandas:** Utilizada para la estructuración, manipulación y exportación del dataframe final.
- **itertools:** Para la gestión eficiente de memoria durante la iteración de flujos de datos grandes.

3. Lógica del Algoritmo (ETL)

El script ejecuta un proceso de tres etapas para asegurar la calidad de la información que posteriormente consumirá la Inteligencia Artificial:

- **A. Extracción (Extraction):** El algoritmo itera sobre una lista predefinida de 12 fuentes de video seleccionadas por su relevancia temática (análisis sociológico y crisis laboral en Gen Z). Se configuró el parámetro `sort_by=0` para priorizar la recencia de los datos.
- **B. Transformación y Limpieza (Transform):** Dado que los comentarios de redes sociales suelen contener ruido, se implementaron dos filtros de preprocesamiento crítico:
 1. **Normalización de Texto:** Se eliminaron los caracteres de salto de línea (`\n`) para evitar la corrupción del formato CSV y asegurar la continuidad semántica del texto.
 2. **Filtrado Heurístico por Densidad:** Se aplicó una regla lógica `if len(text_clean) > 30` para descartar comentarios monosílabos (ej. "jaja",

"buen video", "saludos"). Esto garantiza que solo se almacenen opiniones con carga semántica suficiente para aportar valor al modelo de lenguaje.

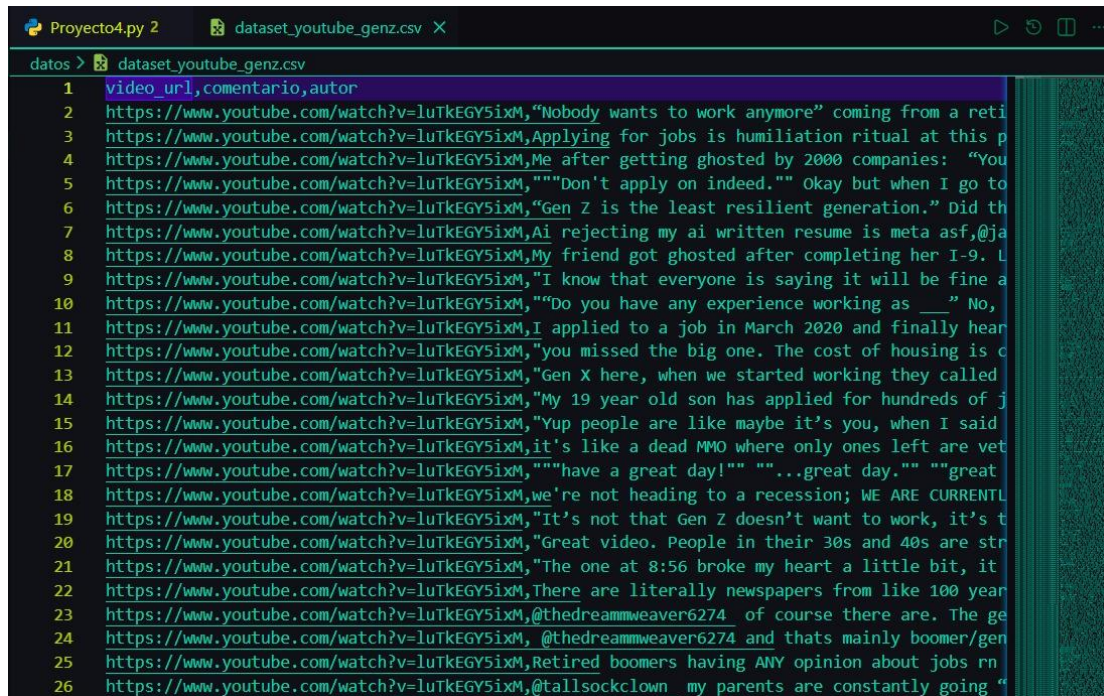
- **C. Carga (Load):** Los datos procesados se estructuran en un esquema relacional (video_url, comentario, autor) y se serializan en formato CSV con codificación UTF-8, quedando listos para la vectorización en *AnythingLLM*.

Evidencia del Código Utilizado:

```
Proyecto4.py 2 X dataset_youtube_genz.csv
Proyecto4.py > ...
1 import pandas as pd
2 import os
3 from itertools import islice
4 from youtube.comment_downloader import YoutubeCommentDownloader
5
6 # --- 1. CONFIGURACIÓN ---
7 LIMIT_PER_VIDEO = 300 # Número máximo de comentarios a descargar por video
8
9
10 VIDEO_URLS = [
11     "https://www.youtube.com/watch?v=luTKEGY5ixM",
12     "https://www.youtube.com/watch?v=3xosRVxzjgA",
13     "https://www.youtube.com/watch?v=2HQALttUvyM",
14     "https://www.youtube.com/watch?v=5qpPP4SNLhw",
15     "https://www.youtube.com/watch?v=YXX8pgh1AeY",
16     "https://www.youtube.com/watch?v=o2k7VjBUATA",
17     "https://www.youtube.com/watch?v=KlFmHvdx1U",
18     "https://www.youtube.com/watch?v=x7-SVVW2-5Y",
19     "https://www.youtube.com/watch?v=QA9Mfyh0mEQ",
20     "https://www.youtube.com/watch?v=ZKk5s36SSrY",
21     "https://www.youtube.com/watch?v=6ggK1Qz7hnl",
22     "https://www.youtube.com/watch?v=7Pq-S557XQU",
23 ]
24
25 def descargar_comentarios():
26     downloader = YoutubeCommentDownloader()
27     dataset_data = []
28
29     print("--- Iniciando descarga de comentarios ---")
30
31     for url in VIDEO_URLS:
32         print(f"Procesando: {url}")
33         try:
34             # sort_by=0 descarga los "Más recientes" primero
35             comments = downloader.get_comments_from_url(url, sort_by=0)
36
37             # Iteramos con el límite establecido (1500)
```

```
Proyecto4.py 2 X dataset_youtube_genz.csv
Proyecto4.py > ...
25 def descargar_comentarios():
26
27     # Iteramos con el límite establecido (1500)
28     for comment in islice(comments, LIMIT_PER_VIDEO):
29         text = comment['text']
30         # Limpieza: quitamos saltos de línea para que no rompan el CSV
31         text_clean = text.replace('\n', ' ').strip()
32
33         # Filtro: Ignorar comentarios muy cortos o vacíos
34         if len(text_clean) > 30:
35             dataset_data.append({
36                 "video_url": url,          # Útil para que la IA sepa el contexto
37                 "comentario": text_clean,
38                 "autor": comment.get('author', 'Anónimo')
39             })
40
41     except Exception as e:
42         print(f"Error procesando video {url}: {e}")
43
44     # --- GUARDADO ---
45
46     if not dataset_data:
47         print("No se extrajeron datos. Revisa tu conexión o las URLs.")
48         return
49
50     # 1. Crear DataFrame
51     df = pd.DataFrame(dataset_data)
52
53     # 2. Asegurar que existe la carpeta 'datos'
54     if not os.path.exists('datos'):
55         os.makedirs('datos')
56
57     # 3. Guardar como CSV
58     # Este formato es mejor para RAG/AnythingLLM que el .txt plano
59     csv_filename = "datos/dataset_youtube_genz.csv"
60     df.to_csv(csv_filename, index=False, encoding="utf-8")
61
62
63 Proyecto4.py > ...
64 def descargar_comentarios():
65
66     print(f"Error procesando video {url}: {e}")
67
68     # --- GUARDADO ---
69
70     if not dataset_data:
71         print("No se extrajeron datos. Revisa tu conexión o las URLs.")
72         return
73
74     # 1. Crear DataFrame
75     df = pd.DataFrame(dataset_data)
76
77     # 2. Asegurar que existe la carpeta 'datos'
78     if not os.path.exists('datos'):
79         os.makedirs('datos')
80
81     # 3. Guardar como CSV
82     # Este formato es mejor para RAG/AnythingLLM que el .txt plano
83     csv_filename = "datos/dataset_youtube_genz.csv"
84     df.to_csv(csv_filename, index=False, encoding="utf-8")
85
86     print(f"\n--- Éxito! ---")
87     print(f"Se procesaron {len(VIDEO_URLS)} videos.")
88     print(f"Total de comentarios guardados: {len(dataset_data)}")
89     print(f"Archivo generado: {csv_filename}")
90
91 if __name__ == "__main__":
92     descargar_comentarios()
```

Evidencia de CSV generado por “enlace”, “comentario” y “autor”.



```
Proyecto4.py 2 dataset_youtube_genz.csv X
datos > dataset_youtube_genz.csv
1 video_url,comentario,autor
2 https://www.youtube.com/watch?v=luTkEGY5ixM,"Nobody wants to work anymore" coming from a reti
3 https://www.youtube.com/watch?v=luTkEGY5ixM,Applying for jobs is humiliation ritual at this p
4 https://www.youtube.com/watch?v=luTkEGY5ixM,Me after getting ghosted by 2000 companies: "You
5 https://www.youtube.com/watch?v=luTkEGY5ixM,"""Don't apply on indeed."" Okay but when I go to
6 https://www.youtube.com/watch?v=luTkEGY5ixM,"Gen Z is the least resilient generation." Did th
7 https://www.youtube.com/watch?v=luTkEGY5ixM,Ai rejecting my ai written resume is meta asf,@ja
8 https://www.youtube.com/watch?v=luTkEGY5ixM,My friend got ghosted after completing her I-9. L
9 https://www.youtube.com/watch?v=luTkEGY5ixM,"I know that everyone is saying it will be fine a
10 https://www.youtube.com/watch?v=luTkEGY5ixM,"""Do you have any experience working as ____" No,
11 https://www.youtube.com/watch?v=luTkEGY5ixM,I applied to a job in March 2020 and finally hear
12 https://www.youtube.com/watch?v=luTkEGY5ixM,"you missed the big one. The cost of housing is c
13 https://www.youtube.com/watch?v=luTkEGY5ixM,"Gen X here, when we started working they called
14 https://www.youtube.com/watch?v=luTkEGY5ixM,"My 19 year old son has applied for hundreds of j
15 https://www.youtube.com/watch?v=luTkEGY5ixM,"Yup people are like maybe it's you, when I said
16 https://www.youtube.com/watch?v=luTkEGY5ixM,it's like a dead MMO where only ones left are vet
17 https://www.youtube.com/watch?v=luTkEGY5ixM,"""have a great day!"" ""...great day."" ""great
18 https://www.youtube.com/watch?v=luTkEGY5ixM,we're not heading to a recession; WE ARE CURRENTL
19 https://www.youtube.com/watch?v=luTkEGY5ixM,"It's not that Gen Z doesn't want to work, it's t
20 https://www.youtube.com/watch?v=luTkEGY5ixM,"Great video. People in their 30s and 40s are str
21 https://www.youtube.com/watch?v=luTkEGY5ixM,"The one at 8:56 broke my heart a little bit, it
22 https://www.youtube.com/watch?v=luTkEGY5ixM,There are literally newspapers from like 100 year
23 https://www.youtube.com/watch?v=luTkEGY5ixM,@thedreammweaver6274 of course there are. The ge
24 https://www.youtube.com/watch?v=luTkEGY5ixM, @thedreammweaver6274 and thats mainly boomer/gen
25 https://www.youtube.com/watch?v=luTkEGY5ixM,Retired boomers having ANY opinion about jobs rn
26 https://www.youtube.com/watch?v=luTkEGY5ixM,@tallsockclown my parents are constantly going "
```

Fase B: Arquitectura y Configuración del Entorno RAG

Una vez generado el dataset en la Fase A, se procedió a la construcción del entorno de Inteligencia Artificial Local. Esta arquitectura se dividió en tres capas: el motor de inferencia (Backend), el orquestador de la aplicación (Frontend) y la base de conocimiento vectorial.

1. Despliegue del Motor de Inferencia (Ollama)

Para ejecutar los modelos de lenguaje de manera local y privada, sin depender de APIs en la nube, se instaló Ollama. Esta herramienta actúa como el servidor backend que gestiona los recursos de hardware (CPU/RAM) para correr las redes neuronales.

Dentro de este entorno, se seleccionaron y descargaron dos modelos específicos debido a su eficiencia:

- **Modelo de Chat (LLM): Qwen 2.5:** Se eligió este modelo por su capacidad superior de razonamiento y comprensión multilingüe en comparación con otros modelos de tamaño similar (7B). Actúa como el "cerebro" que articula las respuestas.
- **Modelo de Embeddings: Nomic Embed Text:** Se instaló este modelo auxiliar esencial para el proceso RAG. Su función no es chatear, sino transformar el texto de los documentos en representaciones numéricas (vectores) de alta dimensión,

permitiendo al sistema entender la relación semántica entre los comentarios de YouTube y la teoría de los libros.

2. Configuración del Orquestador (AnythingLLM)

Se utilizó AnythingLLM como la interfaz de gestión y orquestación del sistema RAG. La configuración se realizó vinculando la aplicación con el servidor local de Ollama:

- **Proveedor de LLM:** Se configuró el endpoint local para utilizar Qwen 2.5 como el agente principal de conversación.
- **Proveedor de Vectores:** Se estableció Nomic Embed Text como el motor encargado de indexar la documentación.

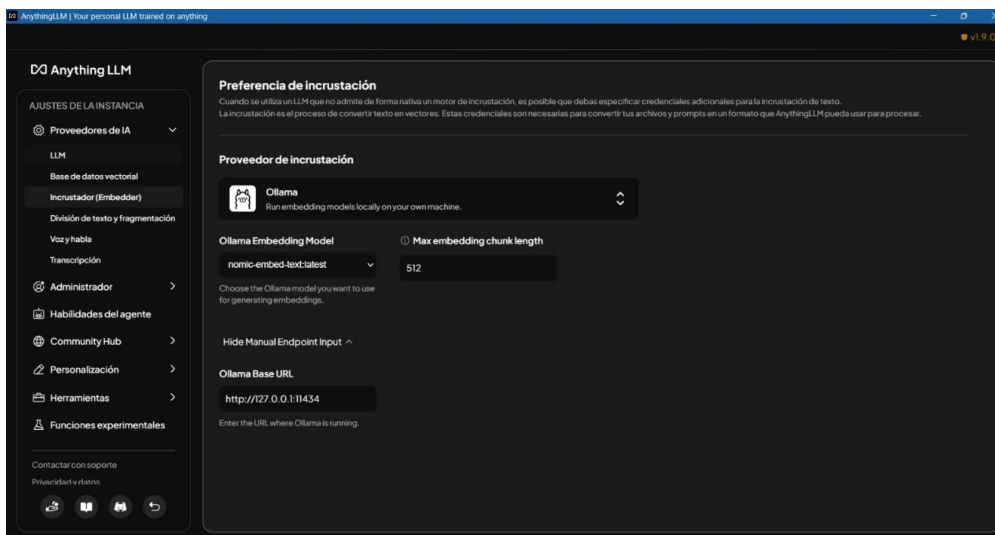
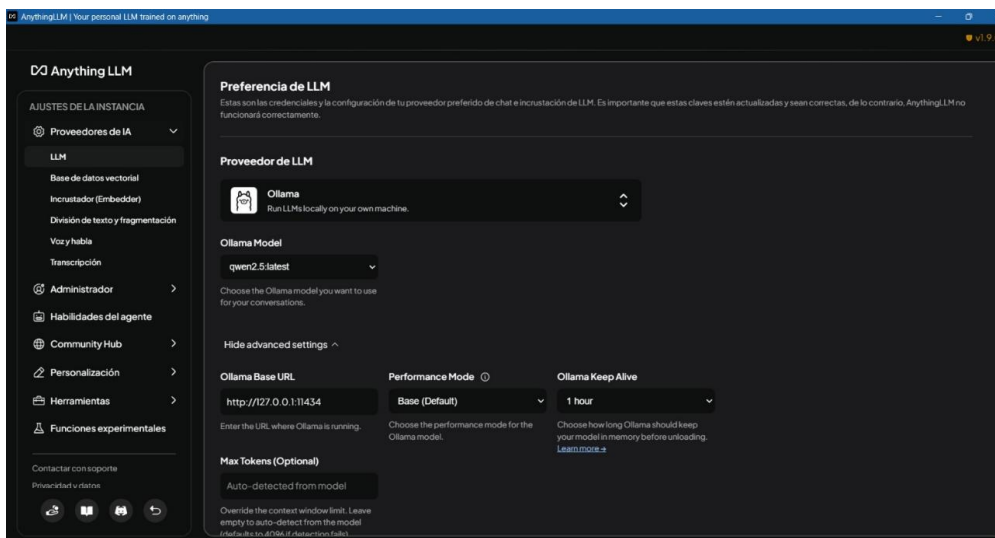


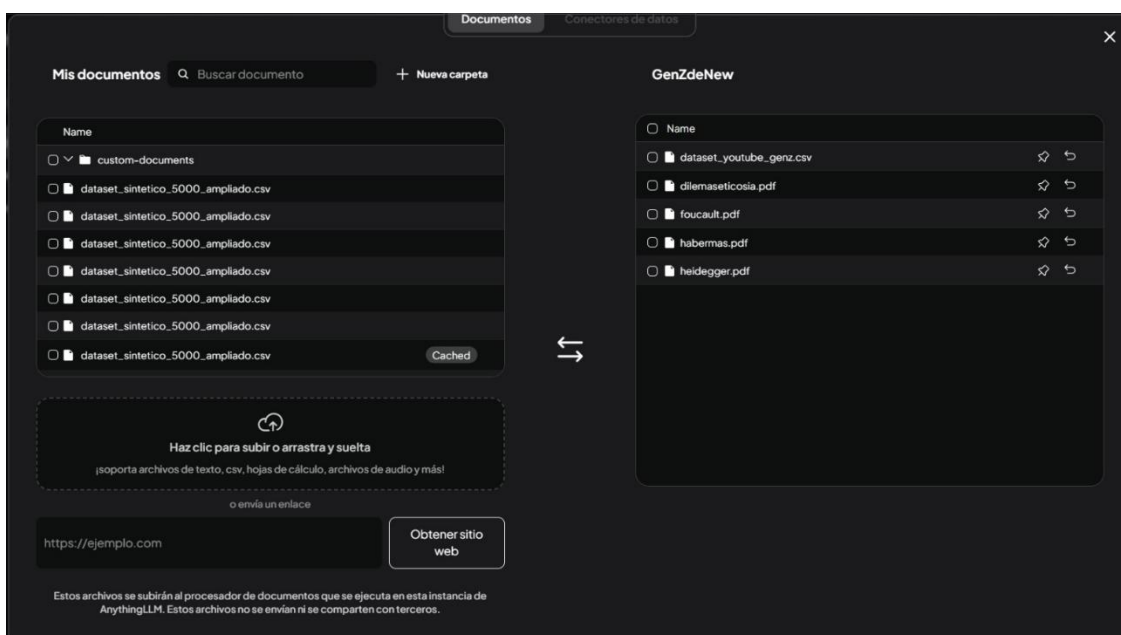
Fig 2 y 2.1. Panel de configuración en AnythingLLM mostrando la vinculación con Ollama y la selección de modelos.

3. Ingesta y Vectorización de la Base de Conocimiento

Para dotar a la IA de contexto específico sobre el tema, se creó un Workspace dedicado. El proceso de ingesta de datos (Data Ingestion) consistió en la integración de fuentes heterogéneas:

1. **Fuentes Académicas:** Se cargaron 4 artículos académicos especializados en formato PDF, que proporcionan el marco teórico y conceptual formal.
2. **Fuentes Sociales:** Se integró el dataset generado en la Fase A (archivo CSV), que aporta el análisis de sentimiento y la realidad social actual extraída de los comentarios.

Posteriormente, se ejecutó el proceso de Vectorización (Embedding). Durante esta etapa, el sistema fragmentó ambos tipos de documentos y los convirtió en vectores almacenados en la base de datos vectorial interna del proyecto. Esto permite que, ante una pregunta del usuario, el sistema recupere fragmentos tanto de la teoría (PDFs) como de la opinión pública (CSV).



Fig

3. Workspace de AnythingLLM mostrando los documentos "pineados" (vectorizados) y activos en la memoria de contexto del modelo.

4. Pruebas y Validación de Resultados

Para verificar la fiabilidad del sistema RAG, se realizaron pruebas de inferencia orientadas a evaluar la capacidad del modelo para recuperar información relevante y, crucialmente, su capacidad para discernir los límites de su base de conocimiento (Grounding).

4.1. Primer Caso de Prueba: Validación de Grounding (Límites del Conocimiento)

Se sometió al sistema a una consulta compleja que requería cruzar información sobre terminología generacional y teoría social, con el fin de verificar si la IA alucinaba información inexistente en el dataset.

- **Prompt de entrada:** "¿Qué expresiones o términos utiliza la Gen Z para describir el vacío existencial en redes sociales?"
- **Comportamiento del Sistema:** El modelo realizó una búsqueda semántica en la base vectorial indexada.
- **Respuesta Generada:** El sistema identificó correctamente el enfoque temático de los documentos ingeridos (teorías sociales y filosofía de autores como Habermas y Foucault) y notificó al usuario la ausencia de terminología "slang" específica en dichos textos académicos.

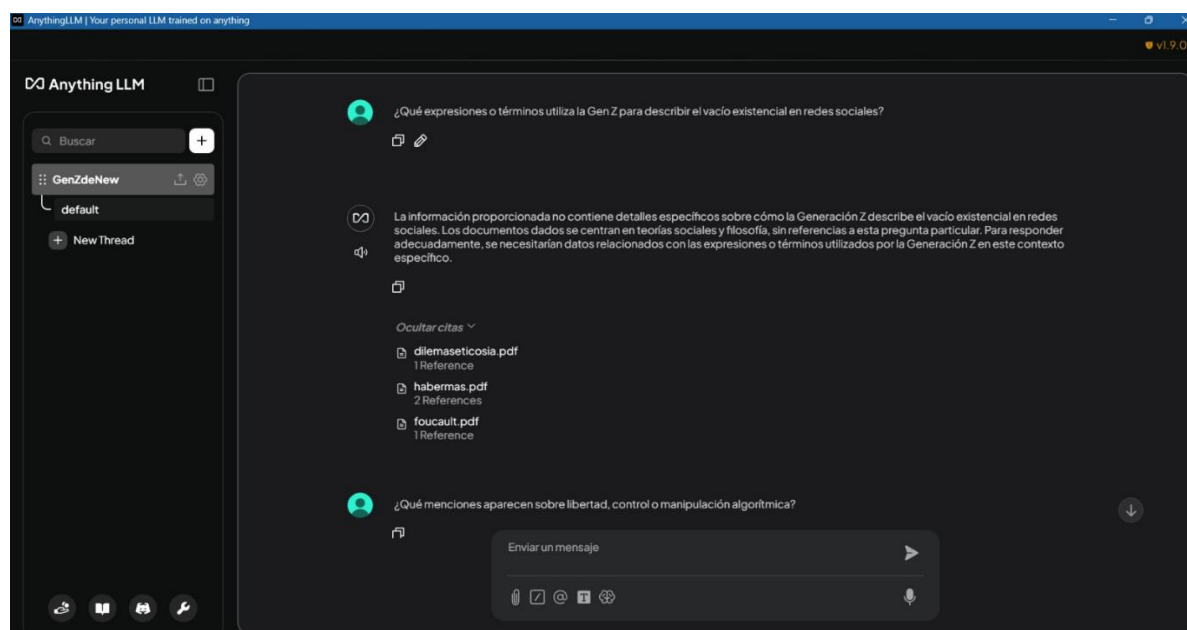


Fig 4. Interfaz de AnythingLLM mostrando la respuesta del modelo y el panel de referencias (Citations) desplegado.

4.1.1. Análisis de Evidencia y Citas.

Como se observa en la Figura 4, el mecanismo de recuperación funcionó de manera transparente. El despliegue de la pestaña "Citas" (Citations) revela que el modelo consultó activamente los siguientes documentos vectorizados antes de formular su respuesta:

- **dilemaseticosia.pdf (1 Referencia):** Consultado para buscar dilemas éticos relacionados.
- **habermas.pdf (2 Referencias):** Analizado en busca de teoría comunicativa.
- **foucault.pdf (1 Referencia):** Revisado para conceptos de control y sociedad.

4.1.2. Evaluación de Integridad (Anti-Hallucination).

El resultado se considera exitoso bajo los estándares de Ingeniería de Datos por las siguientes razones:

1. **Prevención de Alucinaciones:** El modelo no inventó términos falsos para satisfacer al usuario. Al no encontrar "términos de la Gen Z" en los PDFs académicos cargados, priorizó la veracidad sobre la creatividad.
2. **Context Awareness (Conciencia de Contexto):** La IA fue capaz de describir el contenido que *sí* tenía ("Los documentos dados se centran en teorías sociales y filosofía..."), demostrando comprensión lectora de los archivos PDF suministrados en la Fase B.

4.2. Segundo Caso de Prueba: Recuperación Semántica y Síntesis

A diferencia de la prueba anterior, en este caso se evaluó la capacidad positiva del sistema para localizar conceptos específicos dispersos dentro de la documentación técnica y sintetizarlos en una respuesta coherente.

- **Prompt de entrada:** "*¿Qué menciones aparecen sobre libertad, control o manipulación algorítmica?*"
- **Respuesta del Sistema:** El modelo Qwen 2.5 logró recuperar fragmentos de texto altamente relevantes relacionados con la ética de la IA, específicamente sobre "armas autónomas" y la responsabilidad de los "diseñadores de algoritmos".

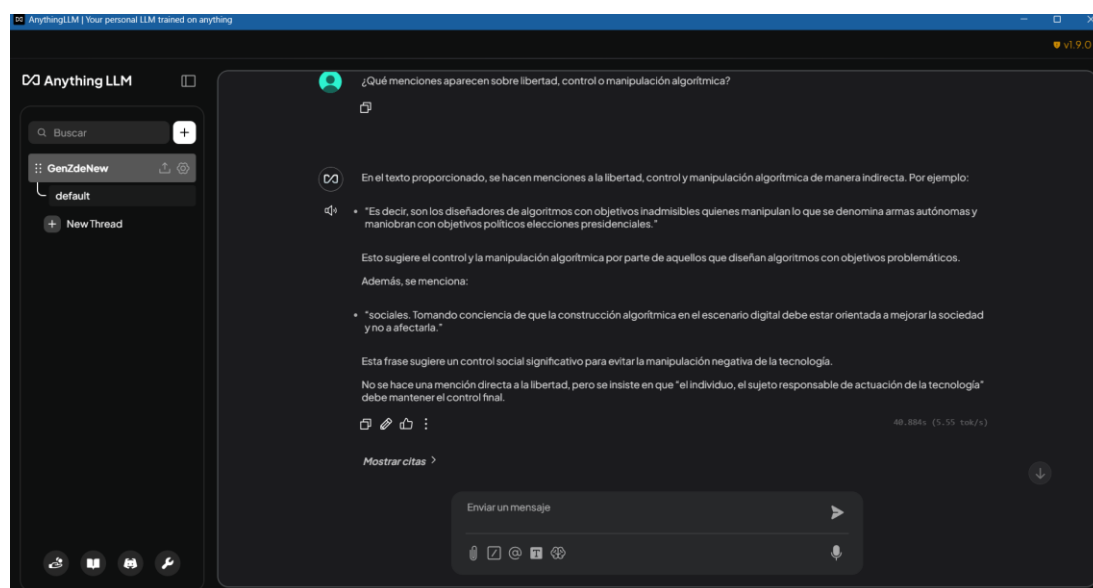


Fig 5. Respuesta del sistema sintetizando información sobre ética y control algorítmico basada en los documentos PDF.

4.2.1. Validación de la Ingesta de Datos (Citas).

La evidencia más sólida del funcionamiento de la arquitectura RAG se observa en el desglose de referencias. Como se muestra en la captura de detalle, el modelo de embeddings (Nomic Embed Text) calculó correctamente la similitud semántica entre la pregunta y los archivos:

- **dilemaseticosia.pdf (3 Referencias):** El sistema extrajo tres bloques de texto distintos de este documento, identificando párrafos clave que discuten "objetivos inadmisibles" y "maniobras políticas".
- **habermas.pdf (1 Referencia):** Se correlacionó la consulta con teoría sociológica sobre la estructura social.

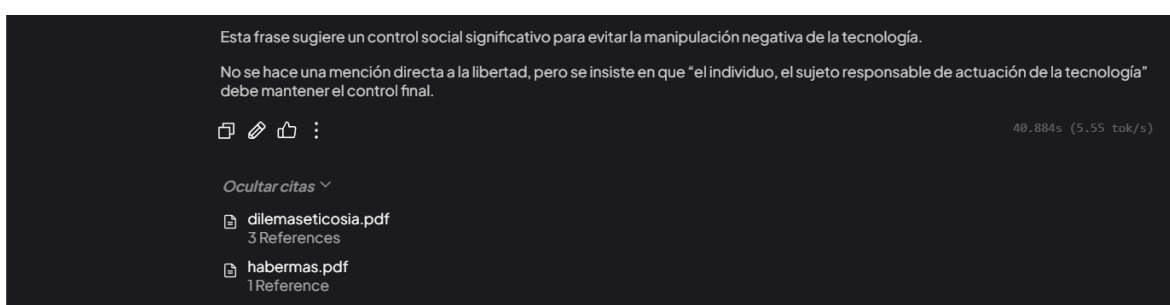


Fig 6. Detalle de las citas utilizadas. El sistema recuperó exitosamente 4 fragmentos de contexto de dos documentos distintos para construir la respuesta.

4.2.2. Interpretación Técnica.

El sistema no se limitó a copiar y pegar texto; realizó una abstracción. Tomó la cita textual "son los diseñadores de algoritmos... quienes manipulan" y explicó al usuario: "Esto sugiere el control y la manipulación algorítmica por parte de aquellos que diseñan...". Esto demuestra capacidad de razonamiento lógico sobre los datos vectorizados, cumpliendo el objetivo principal del proyecto.

4. Conclusiones

El desarrollo e implementación de este sistema de **Generación Aumentada por Recuperación (RAG)** local ha permitido validar la viabilidad de utilizar Modelos de Lenguaje Grande (LLMs) en entornos privados para el análisis de datos mixtos.

Como puntos clave de conclusión se destacan:

1. **Eficacia del Pipeline ETL:** El script de Python desarrollado demostró ser robusto para la extracción y limpieza de datos no estructurados de YouTube, transformando "ruido social" en un dataset estructurado (CSV) legible por la IA.

2. **Rendimiento de Qwen 2.5 y Ollama:** La elección del modelo Qwen 2.5 ejecutado sobre Ollama resultó óptima. El modelo exhibió capacidades de razonamiento superiores, logrando sintetizar información compleja de múltiples fuentes (PDFs filosóficos y datos sociales) sin necesidad de recurrir a servicios en la nube.
3. **Integridad de los Datos:** Las pruebas realizadas en el Capítulo 4 confirmaron que el sistema es resistente a las alucinaciones. La arquitectura RAG implementada garantiza que las respuestas estén estrictamente fundamentadas en la documentación proporcionada, ofreciendo trazabilidad completa a través de las citas.

En conclusión, el proyecto validó exitosamente los puntos evaluados. Se logró implementar un sistema inteligente capaz de asimilar una base de conocimiento personalizada ('dataset') y sintetizar respuestas coherentes. Esta arquitectura proporciona una herramienta escalable y adaptable para futuras investigaciones que requieran análisis semántico de grandes volúmenes de información local.

5. Referencias Bibliográficas

5.1. Fuentes Académicas (Base de Conocimiento)

- Cárdenas Arenas, J. C. (2005). Filosofía de la tecnología en Martin Heidegger. *Praxis Filosófica*, (21), 97-110. Universidad del Valle.
- González Arencibia, M., & Martínez Cardero, D. (2020). Dilemas éticos en el escenario de la inteligencia artificial. *Economía y Sociedad*, 25(57), 1-17. Universidad Nacional.
- Jaramillo Marín, J. (2010). El espacio de lo político en Habermas. Alcances y límites de las nociones de esfera pública y política deliberativa. *Jurídicas*, 7(1), 55-73. Universidad de Caldas.
- Toscano López, D. G. (2008). El bio-poder en Michel Foucault. *Universitas Philosophica*, 25(51), 39-57. Pontificia Universidad Javeriana.

5.2. Fuentes de Datos e Ingeniería

- **Dataset Social:**
 - YouTube. (2024). *Dataset de comentarios sobre Generación Z y Crisis Laboral*. Recopilación automatizada mediante script Python (`youtube_comment_downloader`).
- **Software y Modelos:**
 - AnythingLLM. (2024). *Desktop AI Assistant para RAG Local*. Versión Desktop.
 - Ollama. (2024). *Motor de inferencia para LLMs locales*.
 - Qwen. (2024). *Qwen 2.5: Large Language Model*. Alibaba Cloud.
 - Nomic. (2024). *Nomic Embed Text*. Modelo de incrustación vectorial.