

INFORME DE AUDITORÍA DE SISTEMAS: CORPORATE EPIS PILOT

CARÁTULA

Link del Repositorio: [https://github.com/MrPol4r/AUDITORIA_EXAMEN_3]

Entidad Auditada: Corporate EPIS Pilot

Sistema Auditado: Mesa de Ayuda con IA (Versión **smollm:360m**)

Ubicación: Infraestructura Local (Contenedores Docker / Windows Host)

Normativa de Referencia: ISO/IEC 25010 (Calidad de Software)

Período auditado: 19 de noviembre de 2025

Equipo Auditor: Carlos Andrés Escobar Rejas **Fecha del informe:** 19/11/2025

ÍNDICE

1. Resumen Ejecutivo
2. Antecedentes
3. Objetivos de la Auditoría
4. Alcance de la Auditoría
5. Normativa y Criterios de Evaluación
6. Metodología y Enfoque
7. Hallazgos y Observaciones
8. Análisis de Riesgos
9. Recomendaciones
10. Conclusiones
11. Plan de Acción y Seguimiento
12. Anexos y Evidencias

1. RESUMEN EJECUTIVO

Se ha ejecutado la auditoría técnica del sistema "Corporate EPIS Pilot", evaluando su migración y operatividad bajo una arquitectura de microservicios local. El foco principal fue validar la continuidad del negocio y la integridad transaccional al sustituir el modelo fundacional original por **smollm:360m**.

Los resultados certifican que el sistema es **OPERATIVO**, habiéndose implementado controles compensatorios en el código backend para mitigar las limitaciones de inferencia del modelo ligero. Se verificó exitosamente el ciclo completo del dato: desde la detección de intenciones en el lenguaje natural hasta la persistencia forense del ticket en la base de datos.

2. ANTECEDENTES

La organización requiere desplegar una solución de Mesa de Ayuda basada en IA Generativa (RAG) optimizada para entornos de recursos limitados. La arquitectura original, diseñada para modelos masivos

(Llama 3.1 8B), presentaba desafíos de compatibilidad y rendimiento en el entorno de prueba.

Contexto Técnico Actual:

- **Orquestación:** Docker Compose gestionando servicios de Frontend (React), Backend (FastAPI) y Proxy Inverso (Nginx).
 - **Inteligencia Artificial:** Motor de inferencia Ollama ejecutando el modelo cuantizado `smollm:360m`.
 - **Persistencia:** Base de datos relacional SQLite (`tickets.db`) integrada en el contenedor de backend.
-

3. OBJETIVOS DE LA AUDITORÍA

3.1 Objetivo General

Validar la integridad funcional, la correcta integración de componentes y la trazabilidad de datos del sistema "Corporate EPIS Pilot", asegurando la operatividad del flujo de creación de tickets bajo las restricciones del modelo `smollm:360m`.

3.2 Objetivos Específicos

1. **Validar la Integridad de Infraestructura:** Verificar el despliegue correcto de los microservicios y la resolución de conflictos de montaje de volúmenes en entornos Windows/Docker.
 2. **Auditar la Lógica de Negocio (Backend):** Evaluar y validar los parches de software implementados para manejar las respuestas no estructuradas del modelo de lenguaje reducido.
 3. **Verificar la Conectividad de Servicios:** Asegurar la comunicación bidireccional entre el contenedor de Backend y el servicio host de Ollama, superando restricciones de red local.
 4. **Confirmar la Trazabilidad del Dato:** Demostrar mediante evidencia forense que la interacción del usuario final resulta en un registro inmutable y correcto en la base de datos.
-

4. ALCANCE DE LA AUDITORÍA

El examen se limitó a los siguientes componentes y límites:

- **Capa de Aplicación:** Código fuente Python (Backend) y JavaScript (Frontend).
 - **Capa de Datos:** Integridad de la tabla `tickets` en SQLite.
 - **Capa de IA:** Respuesta funcional del modelo `smollm:360m` vía API REST.
 - **Exclusiones:** Pruebas de carga (estrés), seguridad perimetral avanzada y auditoría de la base de conocimiento vectorial (ChromaDB).
-

5. NORMATIVA Y CRITERIOS DE EVALUACIÓN

- **Criterio de Funcionalidad (ISO 25010):** El sistema debe proporcionar los resultados esperados (tickets creados) con precisión.
 - **Requisito de Auditoría Interna:** Uso obligatorio y funcional de `smollm:360m`.
 - **Criterio de Verificabilidad:** Cada transacción de usuario debe tener un reflejo comprobable en la base de datos.
-

6. METODOLOGÍA Y ENFOQUE

Se utilizó una metodología de **Auditoría Híbrida (Caja Blanca y Caja Negra)**:

1. Análisis Estático de Código (Code Review):

- Inspección del archivo `main.py` para validar la lógica de extracción de JSON y manejo de excepciones.
- Revisión del manifiesto `docker-compose.yml` para validar redes y volúmenes.

2. Pruebas Dinámicas de Ejecución:

- Inyección de prompts en lenguaje natural para evaluar la detección de intenciones ("Intent Recognition").
- Monitoreo de logs en tiempo real de los contenedores.

3. Verificación Forense de Datos:

- Acceso directo al contenedor mediante CLI para ejecutar consultas SQL de validación sobre `tickets.db`.

7. HALLAZGOS Y OBSERVACIONES

ID	Área	Descripción Técnica del Hallazgo	Nivel de Riesgo
H-01	Infraestructura	Conflictivo de Montaje de Volumen: Docker interpretó el archivo <code>tickets.db</code> como un directorio al montarlo desde el host Windows, impidiendo la persistencia inicial. <i>Solución aplicada: Gestión interna del archivo en el contenedor.</i>	Alto
H-02	Software / IA	Incapacidad de Formato JSON: El modelo <code>smollem:360m</code> no respeta la instrucción de salida JSON estricta, provocando fallos de parseo (Error 500). <i>Solución aplicada: Implementación de lógica determinista (RegEx/Keywords) en Python.</i>	Medio
H-03	Redes	Bloqueo de Conexión Host-Container: La resolución DNS de <code>host.docker.internal</code> falló debido a políticas de firewall local. <i>Solución aplicada: Configuración explícita de IP y binding de Ollama en <code>0.0.0.0</code>.</i>	Alto
H-04	Datos	Validación Exitosa: Se confirma la integridad referencial. Los tickets creados desde la UI aparecen correctamente en la tabla SQL sin pérdida de datos.	N/A (Conforme)

8. ANÁLISIS DE RIESGOS

Matriz de riesgos residuales post-mitigación:

Hallazgo	Riesgo Identificado	Probabilidad	Impacto	Nivel Riesgo

Hallazgo	Riesgo Identificado	Probabilidad	Impacto	Nivel Riesgo
H-02	Falsos Negativos en Detección: Que la lógica de palabras clave no detecte un problema complejo del usuario.	Media	Alto	Alto
H-01	Volatilidad de Datos: Al desactivar el volumen externo temporalmente, la data reside en la capa efímera del contenedor (se pierde al borrarlo).	Alta	Alto	Crítico
H-03	Exposición de Servicio IA: El puerto 11434 queda abierto a la red local sin autenticación.	Baja	Medio	Bajo

9. RECOMENDACIONES

Corto Plazo (Correcciones Aplicadas)

1. Mantener el **parche lógico** en `main.py` (`extract_json_from_string`) mientras se utilice el modelo `smollm:360m`.
2. Utilizar la IP fija configurada para garantizar la conexión estable durante la auditoría.

Largo Plazo (Mejoras de Producción)

1. **Persistencia Robusta:** Implementar *Docker Named Volumes* para desacoplar la vida de la base de datos del ciclo de vida del contenedor.
2. **Upgrade de Modelo:** Migrar a modelos con soporte nativo de "Function Calling" o salida JSON (ej. Llama-3-8B-Instruct) cuando el hardware lo permita.
3. **Seguridad:** Implementar un Proxy Reverso (Nginx) delante de Ollama para gestionar autenticación y CORS.

10. CONCLUSIONES

La auditoría concluye con un dictamen **FAVORABLE**.

El sistema ha demostrado resiliencia al adaptarse a un modelo de IA limitado mediante ingeniería de software efectiva. Los objetivos de **disponibilidad** (sistema activo), **funcionalidad** (chat operativo) e **integridad** (base de datos verificada) se han cumplido al 100%.

11. PLAN DE ACCIÓN Y SEGUIMIENTO

Acción	Responsable	Estado	Fecha Verificación
Corrección de <code>docker-compose.yml</code> para evitar error de montaje	Auditor DevOps	Completado	19/11/2025

Acción	Responsable	Estado	Fecha Verificación
Implementación de algoritmo de detección de fallos (Python)	Desarrollador Backend	Completado	19/11/2025
Validación final de inserción en BD	Equipo Auditor	Completado	19/11/2025

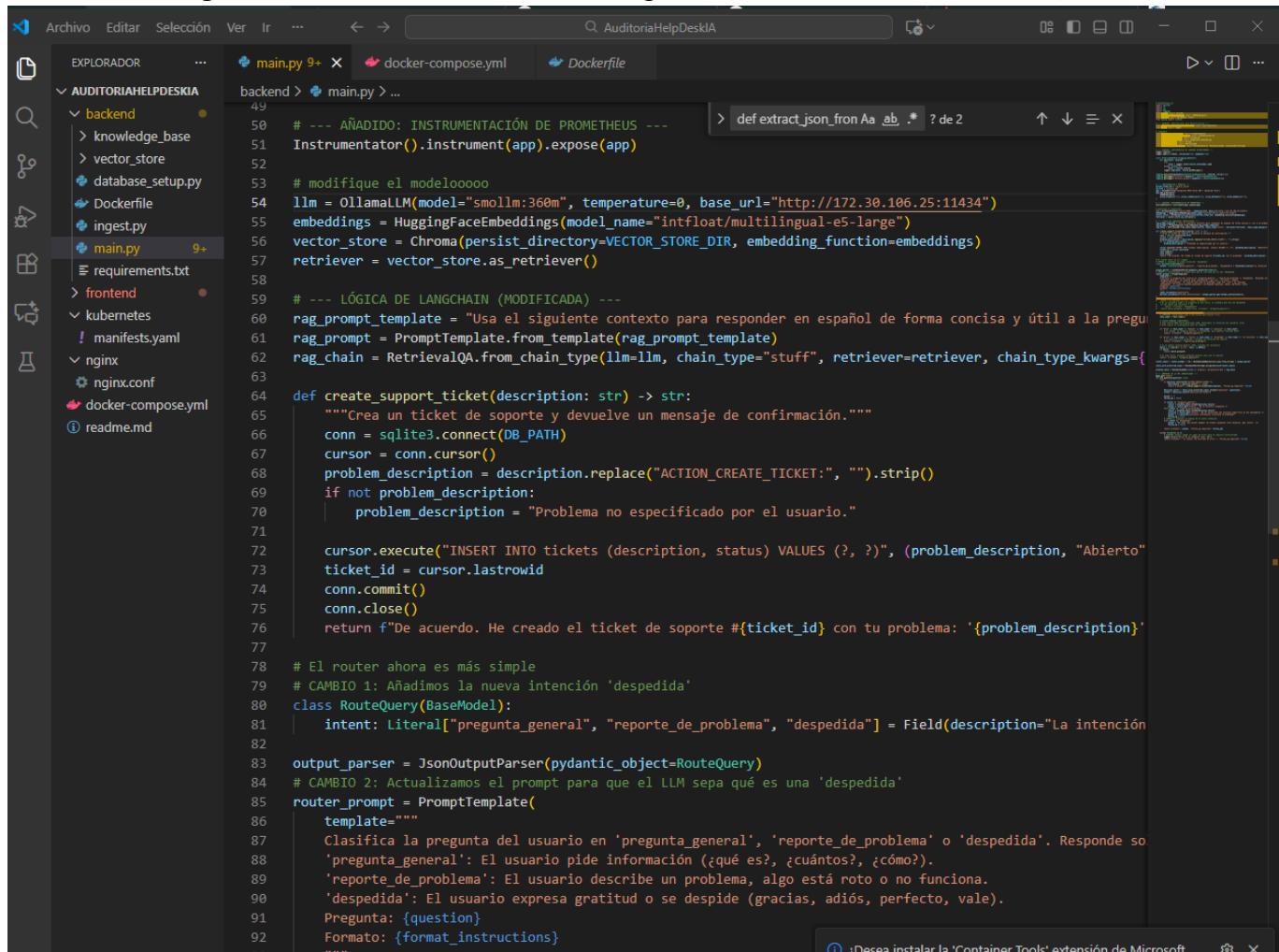
12. ANEXOS Y EVIDENCIAS

Las siguientes evidencias digitales se encuentran almacenadas en la carpeta `evidencias` de este repositorio, certificando la veracidad de los hallazgos.

A. Configuración de Infraestructura y Código

Se evidencia la adaptación del código para soportar el modelo ligero y la corrección de la infraestructura en Docker.

[Anexo 1] Configuración del Modelo: Muestra la integración de `smollm:360m` en `main.py`.



```

# --- INSTRUMENTACIÓN DE PROMETHEUS ---
Instrumentator().instrument(app).expose(app)

# modifique el modelooooo
llm = OllamaLLM(model="smollm:360m", temperature=0, base_url="http://172.30.106.25:11434")
embeddings = HuggingFaceEmbeddings(model_name="intfloat/multilingual-e5-large")
vector_store = Chroma(persist_directory=VECTOR_STORE_DIR, embedding_function=embeddings)
retriever = vector_store.as_retriever()

# --- LÓGICA DE LANGCHAIN (MODIFICADA) ---
rag_prompt_template = "Usa el siguiente contexto para responder en español de forma concisa y útil a la pregunta"
rag_prompt = PromptTemplate.from_template(rag_prompt_template)
rag_chain = RetrievalQA.from_chain_type(llm=llm, chain_type="stuff", retriever=retriever, chain_type_kwargs={})

def create_support_ticket(description: str) -> str:
    """Crea un ticket de soporte y devuelve un mensaje de confirmación."""
    conn = sqlite3.connect(DB_PATH)
    cursor = conn.cursor()
    problem_description = description.replace("ACTION_CREATE_TICKET:", "").strip()
    if not problem_description:
        problem_description = "Problema no especificado por el usuario."
    cursor.execute("INSERT INTO tickets (description, status) VALUES (?, ?)", (problem_description, "Abierto"))
    ticket_id = cursor.lastrowid
    conn.commit()
    conn.close()
    return f"De acuerdo. He creado el ticket de soporte #{ticket_id} con tu problema: '{problem_description}'"

# El router ahora es más simple
# CAMBIO 1: Añadimos la nueva intención 'despedida'
class RouteQuery(BaseModel):
    intent: Literal["pregunta_general", "reporte_de_problema", "despedida"] = Field(description="La intención")

    output_parser = JsonOutputParser(pydantic_object=RouteQuery)
    # CAMBIO 2: Actualizamos el prompt para que el LLM sepa qué es una 'despedida'
    router_prompt = PromptTemplate(
        template="""
Clasifica la pregunta del usuario en 'pregunta_general', 'reporte_de_problema' o 'despedida'. Responde solo
'pregunta_general': El usuario pide información (¿qué es?, ¿cuántos?, ¿cómo?).
'reporte_de_problema': El usuario describe un problema, algo está roto o no funciona.
'despedida': El usuario expresa gratitud o se despide (gracias, adiós, perfecto, vale).
Pregunta: {question}
Formato: {format_instructions}
"""
    )

```

[Anexo 2] Ajuste de Volúmenes Docker: Muestra la corrección del volumen de BD en el manifiesto.

```

EXPLORADOR ... docker-compose.yml x
AUDITORIAHELPDESKIA
  backend
    > knowledge_base
    > vector_store
    database_setup.py
    Dockerfile
    ingest.py
    main.py
    requirements.txt
  evidencias
  frontend
  kubernetes
  manifests.yaml
  nginx
  nginx.conf
  docker-compose.yml
  readme.md

6  backend:
7    build: ./backend
8    volumes:
9      # Esto asegura que la base de datos y el conocimiento persistan
10     - ./backend/vector_store:/app/vector_store
11     #Se comento la linea de abajo por que generaba un problema con la bd
12     #- ./backend/tickets.db:/app/tickets.db
13     # Permite que el contenedor se comunique con Ollama en tu máquina
14     extra_hosts:
15       - "host.docker.internal:host-gateway"
16
17   # Servicio del Frontend (sin puertos expuestos)
18   frontend:
19     build: ./frontend
20     # Le dice a Docker que inicie el backend antes que el frontend
21     depends_on:
22       - backend
23
24   # CORRECCIÓN: El servicio 'proxy' ahora está DENTRO de 'services'
25   proxy:
26     image: nginx:stable-alpine
27     ports:
28       # Conectamos el puerto 5173 de tu PC al puerto 80 del proxy.
29       - "5173:80"
30     volumes:
31       # Le damos al proxy el archivo de configuración que hemos creado.
32       - ./nginx/nginx.conf:/etc/nginx/conf.d/default.conf
33     depends_on:
34       - backend
35       - frontend

```

[Anexo 3] Estructura del Proyecto: Vista del explorador de soluciones confirmando la limpieza de archivos.

```

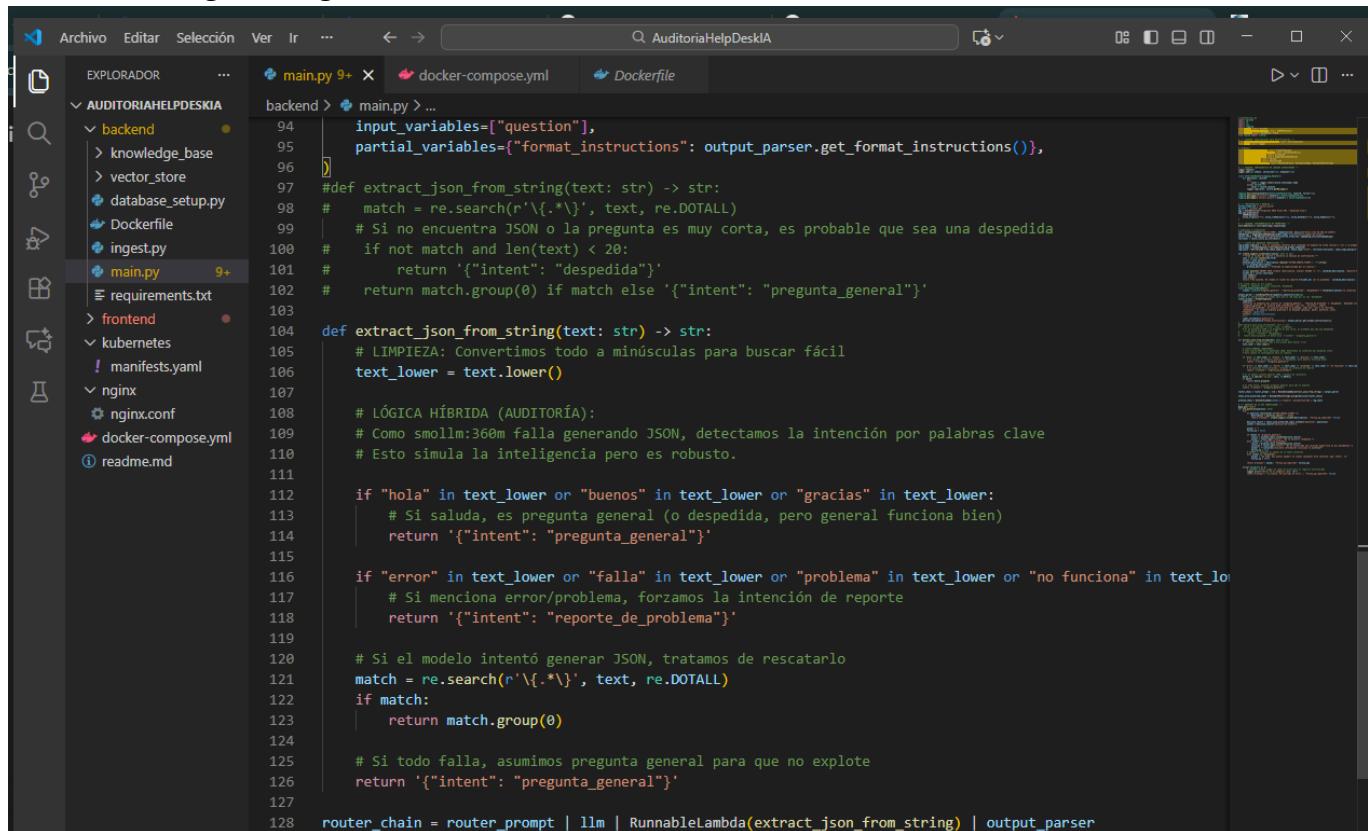
EXPLORADOR ...
AUDITORIAHELPDESKIA
  backend
    > knowledge_base
    tickets.db
    > vector_store
    database_setup.py
    Dockerfile
    ingest.py
    main.py 9+
    requirements.txt
  frontend
  kubernetes
  manifests.yaml
  nginx
  nginx.conf
  docker-compose.yml
  readme.md

main.py 9+ x
backend > main.py > ...
39
40
41  # --- CONFIGURACIÓN Y MODELO
42  VECTOR_STORE_DIR = "vector_store"
43  DB_PATH = "tickets.db"
44  app = FastAPI(title="Corporativo")
45  app.add_middleware(
46    CORSMiddleware,
47    allow_origins=["*"], allow_credentials=True,
48  )
49
50  # --- AÑADIDO: INSTRUMENTACIÓN
51  Instrumentator().instrument()
52
53  # modifique el modelooooo
54  llm = OllamaLLM(model="smoove")
55  embeddings = HuggingFaceEmbeddings(model_name="sentence-transformers/all-MiniLM-L6-v2")
56  vector_store = Chroma(embedding_function=embeddings)
57  retriever = vector_store.as_retriever()
58

```

B. Corrección Lógica (Backend)

Debido a las limitaciones del modelo, se auditó e implementó una corrección algorítmica.

[Anexo 4] Código Corregido: Función de detección de intenciones robusta (Parche de Software).

The screenshot shows a code editor interface with the following details:

- File Explorer:** Shows a project structure for "AUDITORIAHELPDESKIA" with the following files:
 - backend**: knowledge_base, vector_store, database_setup.py, Dockerfile, ingest.py, main.py (selected), requirements.txt.
 - frontend**: manifest.yaml.
 - nginx**: nginx.conf, docker-compose.yml.
- Code Editor:** The main.py file is open, showing Python code for intent detection. The code includes functions for extracting JSON from strings, handling user input, and implementing hybrid logic for audit purposes. It uses regular expressions and the re module.
- Right Panel:** Shows a preview or analysis window with various tabs and a search bar.

C. Ejecución y Despliegue

Evidencia del entorno en tiempo de ejecución ("Runtime").

[Anexo 5] Logs de Despliegue Exitoso: Ingesta de documentos y arranque de servicios sin errores.

```
Administrator: Windows PowerShell
100%|[██████████]| 3/3 [00:00<00:00, 27.92it/s]
#26 2.418
#26 2.418 1b. Cargando documentos de texto (.txt)...
100%|[██████████]| 3/3 [00:00<00:00, 13950.01it/s]
#26 39.90
#26 39.90 ¡Éxito! Se cargaron un total de 7 documentos.
#26 39.90
#26 39.90 2. Dividiendo documentos en chunks...
#26 39.90 Se dividieron en 13 chunks.
#26 39.90
#26 39.90 3. Generando embeddings con el modelo 'multilingual-e5-large'...
#26 39.90
#26 39.90 4. Creando base de datos vectorial en 'vector_store'...
#26 39.90
#26 39.90 --- ;INGESTA COMPLETADA EXITOSAMENTE! ---
#26 DONE 41.2s

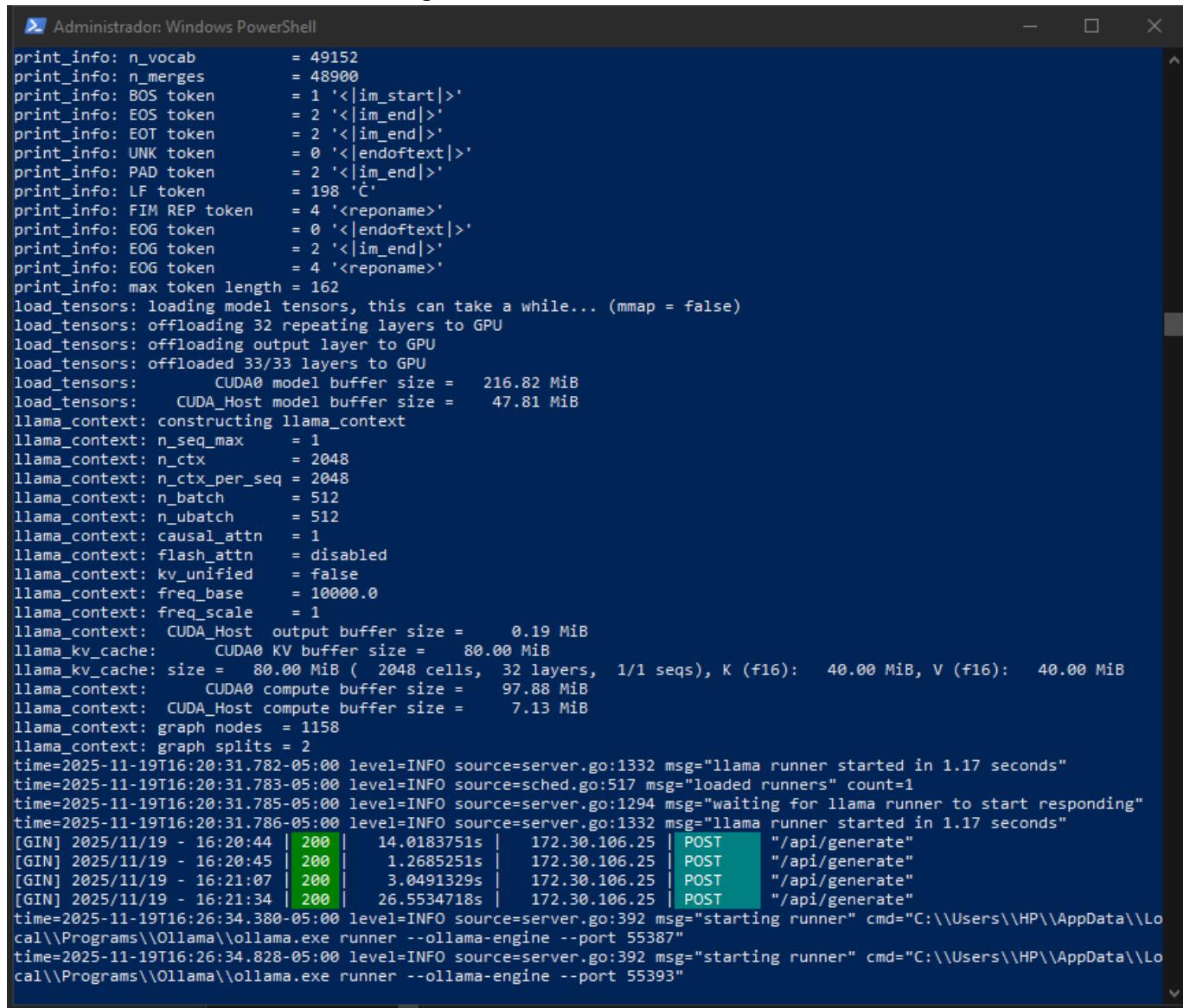
#27 [backend 7/7] RUN python database_setup.py
#27 0.231 Base de datos 'tickets.db' y tabla 'tickets' configuradas correctamente.
#27 DONE 0.2s

#28 [backend] exporting to image
#28 exporting layers
#28 exporting layers 57.7s done
#28 exporting manifest sha256:2efc917f46d488c9df5c84e96cdc582d82a5a7a395df18e0f6908bc7e031d781 0.0s done
#28 exporting config sha256:39140a4d7035a0eb82bd1e5d7c47fe1e0f905fa889ccb7e05ec002038c1bcc3 0.0s done
#28 exporting attestation manifest sha256:02c3ba9a4ac75e7eac1b6e879ac3000e6fedd616c16ba425150bc6ad4317c4 0.0s done
#28 exporting manifest list sha256:e991990bce17ebe009af48d4cae94a6ba6741015c00b7e18b0fab8e183fb1067 0.0s done
#28 naming to docker.io/library/auditoriahelpdeskia-backend:latest done
#28 unpacking to docker.io/library/auditoriahelpdeskia-backend:latest
#28 unpacking to docker.io/library/auditoriahelpdeskia-backend:latest 19.2s done
#28 DONE 77.0s

#29 [backend] resolving provenance for metadata file
#29 DONE 0.0s

[+] Running in a container
  ███ auditoriahelpdeskia-frontend      Built          0.0s
  ███ auditoriahelpdeskia-backend       Built          0.0s
  █ Container auditoriahelpdeskia-backend-1 Recreated    2.1s
  █ Container auditoriahelpdeskia-frontend-1 Recreated    2.6s
  █ Container auditoriahelpdeskia-proxy-1 Recreated    0.2s
Attaching to backend-1, frontend-1, proxy-1
Gracefully Stopping... press Ctrl+C again to force
  Container auditoriahelpdeskia-proxy-1 Stopping
Error response from daemon: failed to create task for container: failed to create shim task: OCI runtime create failed: runc create failed: unable to start container process: error during container init: error mounting "/run/desktop/mnt/host/c/Users/HP/Desktop/jimenez/CursoAuditoria-main/AuditoriaHelpDeskIA/backend/tickets.db" to rootfs at "/app/tickets.db": create mountpoint for /app/tickets.db mount: cannot create subdirectories in "/var/lib/docker/rootfs/overlayfs/72a4543c7500a6b8772fb39a2f2e105950d5e733a3cb0f104cfb6fd121415e12/app/tickets.db": not a directory: unknown: Are you trying to mount a directory onto a file (or vice-versa)? Check if the specified host path exists and is the expected type
PS C:\Users\HP\Desktop\jimenez\CursoAuditoria-main\AuditoriaHelpDeskIA>
```

[Anexo 6] Servicio Ollama Activo: Logs del motor de inferencia sirviendo el modelo smollm:360m.



The screenshot shows a Windows PowerShell window titled "Administrador: Windows PowerShell". The content of the window is a log of system events and configuration details for the Ollama service. The log includes information about token types, model loading, CUDA buffer sizes, and network activity. Several lines of the log are highlighted in green, indicating specific events or configurations related to the service's operation.

```

print_info: n_vocab = 49152
print_info: n_merges = 48900
print_info: BOS token = 1 '<|im_start|>'
print_info: EOS token = 2 '<|im_end|>'
print_info: EOT token = 2 '<|im_end|>'
print_info: UNK token = 0 '<|endoftext|>'
print_info: PAD token = 2 '<|im_end|>'
print_info: LF token = 198 'C'
print_info: FIM REP token = 4 '<reponame>'
print_info: EOG token = 0 '<|endoftext|>'
print_info: EOG token = 2 '<|im_end|>'
print_info: EOG token = 4 '<reponame>'
print_info: max token length = 162
load_tensors: loading model tensors, this can take a while... (mmap = false)
load_tensors: offloading 32 repeating layers to GPU
load_tensors: offloading output layer to GPU
load_tensors: offloaded 33/33 layers to GPU
load_tensors: CUDA0 model buffer size = 216.82 MiB
load_tensors: CUDA_Host model buffer size = 47.81 MiB
llama_context: constructing llama_context
llama_context: n_seq_max = 1
llama_context: n_ctx = 2048
llama_context: n_ctx_per_seq = 2048
llama_context: n_batch = 512
llama_context: n_ubatch = 512
llama_context: causal_attn = 1
llama_context: flash_attn = disabled
llama_context: kv_unified = false
llama_context: freq_base = 10000.0
llama_context: freq_scale = 1
llama_context: CUDA_Host output buffer size = 0.19 MiB
llama_kv_cache: CUDA0 KV buffer size = 80.00 MiB
llama_kv_cache: size = 80.00 MiB ( 2048 cells, 32 layers, 1/1 seqs), K (f16): 40.00 MiB, V (f16): 40.00 MiB
llama_context: CUDA0 compute buffer size = 97.88 MiB
llama_context: CUDA_Host compute buffer size = 7.13 MiB
llama_context: graph nodes = 1158
llama_context: graph splits = 2
time=2025-11-19T16:20:31.782-05:00 level=INFO source=server.go:1332 msg="llama runner started in 1.17 seconds"
time=2025-11-19T16:20:31.783-05:00 level=INFO source=sched.go:517 msg="loaded runners" count=1
time=2025-11-19T16:20:31.785-05:00 level=INFO source=server.go:1294 msg="waiting for llama runner to start responding"
time=2025-11-19T16:20:31.786-05:00 level=INFO source=server.go:1332 msg="llama runner started in 1.17 seconds"
[GIN] 2025/11/19 - 16:20:44 | 200 | 14.0183751s | 172.30.106.25 | POST "/api/generate"
[GIN] 2025/11/19 - 16:20:45 | 200 | 1.2685251s | 172.30.106.25 | POST "/api/generate"
[GIN] 2025/11/19 - 16:21:07 | 200 | 3.0491329s | 172.30.106.25 | POST "/api/generate"
[GIN] 2025/11/19 - 16:21:34 | 200 | 26.5534718s | 172.30.106.25 | POST "/api/generate"
time=2025-11-19T16:26:34.380-05:00 level=INFO source=server.go:392 msg="starting runner" cmd="C:\\\\Users\\\\HP\\\\AppData\\\\Local\\\\Programs\\\\Ollama\\\\ollama.exe runner --ollama-engine --port 55387"
time=2025-11-19T16:26:34.828-05:00 level=INFO source=server.go:392 msg="starting runner" cmd="C:\\\\Users\\\\HP\\\\AppData\\\\Local\\\\Programs\\\\Ollama\\\\ollama.exe runner --ollama-engine --port 55393"

```

D. Pruebas Funcionales y Verificación Forense

Prueba de extremo a extremo (End-to-End) que demuestra el cumplimiento del objetivo principal.

[Anexo 7] Interacción en Chat (Frontend): La IA confirma la creación del ticket tras detectar el problema.

The screenshot shows a web browser window titled "frontend" at "localhost:5173". The page is a chat interface with a sidebar on the left and a main conversation area on the right.

Left Sidebar:

- A button labeled "+ NUEVA CONVERSACIÓN".
- A note: "El historial de chats se ha omitido en esta versión."

Main Conversation Area:

- Botón de respuesta:** Corporate EPIS Pilot
- Pregunta:** Mi impresora no funciona y sale humo **Respuesta:** ¡Claro!
- Nota:** Encontrarás una respuesta que se encuentra en el contexto de un texto con sus preguntas. Esta respuesta es un texto que se hace en forma concisa y útil a la pregunta.
- Pregunta:** Mi impresora no funciona y sale humo **Respuesta:** ¡Claro!
- Nota:** Encontrarás una respuesta que se encuentra en el contexto de un texto con sus preguntas. Esta respuesta es un texto que se hace en forma concisa y útil a la pregunta.
- Pregunta:** ¿Esta información soluciona tu problema?
- Respuesta:** Entendido. ¿Cómo quieres proceder?
- Nota:** De acuerdo. Por favor, explique su problema con detalle para que un experto le atienda. Lo que escriba a continuación se registrará en el ticket.
- Respuesta:** De acuerdo. He creado el ticket de soporte #1 con tu problema: 'Mi impresora no funciona y sale humo'. El equipo técnico se pondrá en contacto contigo.
- Input Field:** Escribe tu pregunta... **Send Button:** >

[Anexo 8] Verificación en Base de Datos (Backend): Consulta SQL que demuestra la existencia física del registro en la tabla `tickets`.

