
ENLACE AL REPOSITORIO

Repositorio GitHub: https://github.com/erick00xx/AUDITORIA_EXAMEN_3/

INFORME FINAL DE AUDITORÍA DE SISTEMAS

CARÁTULA

Entidad Auditada: CORPORATE EPIS PILOT (Sistema de Mesa de Ayuda con IA)

Ubicación: Laboratorio de Cómputo / Entorno Local

Período auditado: 19/11/2025

Equipo Auditor: Erick Scott Churacutipa Blas **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

La presente auditoría de sistemas se realizó sobre el proyecto "Corporate EPIS Pilot", un asistente virtual basado en IA para soporte técnico. El propósito fue verificar la integridad del código, la correcta orquestación de contenedores Docker y la funcionalidad operativa utilizando el modelo de lenguaje `smollm:360m`.

Durante la evaluación inicial, se detectaron **fallas críticas** que impedían el despliegue del sistema: errores de configuración en volúmenes Docker, incompatibilidad de modelos en el código fuente y problemas de concurrencia en sistemas Windows (bloqueo de archivos). Tras la aplicación de medidas correctivas y refactorización del código del Backend (`main.py` y `ingest.py`), se logró levantar el sistema al 100%, garantizando la creación de tickets en base de datos y respuestas coherentes mediante RAG.

2. ANTECEDENTES

La entidad auditada utiliza un sistema de arquitectura de microservicios para gestionar incidentes de TI. El sistema consta de:

- **Frontend:** React + Vite (Interfaz de Chat).
- **Backend:** FastAPI + LangChain + ChromaDB (Lógica RAG).
- **Infraestructura:** Docker Compose y Nginx como Proxy Inverso.
- **IA:** Ollama (Modelos LLM locales).

Previamente, el sistema presentaba inestabilidad al intentar integrarse con modelos ligeros y errores de persistencia de datos al reiniciar los contenedores.

3. OBJETIVOS DE LA AUDITORÍA

Se establecieron los siguientes objetivos para garantizar la operatividad del sistema según los requerimientos del examen.

Objetivo General

Evaluar la calidad del código, la configuración de la infraestructura y la funcionalidad operativa del sistema de Mesa de Ayuda con IA para asegurar su correcto despliegue y funcionamiento con el modelo `smollm:360m`.

Objetivos Específicos

1. Verificar que `docker-compose.yml` y los `Dockerfiles` construyan y comuniquen correctamente los servicios (Frontend, Backend, Proxy).
2. Asegurar que el sistema utilice el modelo `smollm:360m` y que la lógica RAG (Retrieval-Augmented Generation) responda consultas basándose en los documentos internos.
3. Comprobar que la base de datos `tickets.db` y el `vector_store` persistan información entre reinicios y no sean sobreescritos erróneamente por volúmenes vacíos.
4. Analizar y corregir los tiempos de espera (timeouts) del Proxy Nginx y la robustez del Backend ante respuestas lentas o mal formateadas del modelo pequeño.

4. ALCANCE DE LA AUDITORÍA

- **Ámbito Tecnológico:** Código fuente (Python, TypeScript), Scripts de despliegue (Docker), Configuración de Servidor Web (Nginx).
- **Sistemas:** Microservicios de Backend y Frontend interactuando con Ollama local.
- **Periodo Auditado:** 19 de Noviembre de 2025.
- **Exclusiones:** No se auditaron entornos de producción en nube pública (AWS/Azure), limitándose al entorno local de desarrollo.

5. NORMATIVA Y CRITERIOS DE EVALUACIÓN

- **Mejores Prácticas de DevOps:** Principios de inmutabilidad de contenedores y gestión correcta de volúmenes persistentes.
- **ISO/IEC 25010:** Criterios de calidad de software (Funcionalidad, Fiabilidad y Eficiencia).
- **PEP 8:** Estándar de estilo de código para Python.

- **Requerimientos del Stakeholder (Docente):** Uso obligatorio de `smollm:360m` y operatividad al 100%.
-

6. METODOLOGÍA Y ENFOQUE

Se utilizó un enfoque **mixto (Caja Blanca y Caja Negra)**:

1. **Revisión de Código (Caja Blanca):** Análisis estático de `Dockerfile`, `docker-compose.yml`, `main.py` e `ingest.py` para identificar errores lógicos y de configuración.
 2. **Ánalisis de Logs:** Revisión de la salida de consola de los contenedores para identificar "tracebacks" y errores de conexión.
 3. **Pruebas Funcionales (Caja Negra):** Interacción directa con el Chatbot para verificar el flujo de conversación, la consulta de documentos y la creación de tickets.
 4. **Validación de Datos:** Inspección directa del archivo `tickets.db` mediante SQLite Viewer para confirmar la inserción de registros.
-

7. HALLAZGOS Y OBSERVACIONES

Se identificaron 4 hallazgos principales que impedían el funcionamiento del sistema:

Hallazgo 01: Incompatibilidad de Modelo LLM Hardcodeado

- **Descripción:** El archivo `main.py` tenía configurado "llama3.1:8b", mientras que el entorno de auditoría requiera "smollm:360m".
- **Evidencia:** Línea 45 del código original.
- **Criticidad: Alta** (El sistema no respondía).
- **Causa:** Configuración estática no parametrizada.

Hallazgo 02: Sobrescritura de Base de Conocimiento (Trampa de Volumen)

- **Descripción:** El `Dockerfile` ejecutaba la ingestión de datos (`RUN ingest.py`) durante la construcción. Al levantar con Docker Compose, el volumen montado (`- ./backend/vector_store:/app/vector_store`) ocultaba los datos generados dentro de la imagen.
- **Evidencia:** Base de datos vectorial vacía al iniciar el contenedor.
- **Criticidad: Alta** (El bot no tenía "memoria").
- **Causa:** Mala práctica en el ciclo de vida de Docker (Build vs Run).

Hallazgo 03: Error de Bloqueo de Archivos en Windows

- **Descripción:** El script `ingest.py` intentaba borrar la carpeta raíz `vector_store` usando `shutil.rmtree`. Docker en Windows bloquea la carpeta raíz montada, provocando el crash del contenedor.
- **Evidencia:** Log de error `[Errno 16] Device or resource busy`.
- **Criticidad: Media** (Impedía el re-despliegue).
- **Causa:** Restricciones del sistema de archivos en volúmenes montados.

Hallazgo 04: Alucinaciones y Timeouts del Modelo Pequeño

- Descripción:** El modelo `smollm:360m` fallaba al seguir instrucciones JSON complejas y tardaba más del tiempo predeterminado de Nginx (60s), causando errores 504.
 - Evidencia:** Mensaje "Error: No se pudo obtener respuesta del servidor" y logs de `JSONDecodeError`.
 - Criticidad: Media** (Afectaba la experiencia de usuario).
-

8. ANÁLISIS DE RIESGOS

Hallazgo	Riesgo Asociado	Impacto	Probabilidad	Nivel de Riesgo
H-01: Modelo Incorrecto	Interrupción total del servicio de IA	Alto	Alta	Crítico
H-02: Volumen Vacío	Pérdida de disponibilidad de la información (RAG)	Alto	Alta	Alto
H-03: Crash en Ingesta	Fallo en el despliegue y actualización de datos	Medio	Media	Medio
H-04: Timeout/JSON	Experiencia de usuario degradada y respuestas erróneas	Medio	Alta	Alto

9. RECOMENDACIONES

Para mitigar los hallazgos encontrados, se aplicaron y recomiendan las siguientes acciones:

- Corrección de Ciclo de Vida Docker:** Mover la ejecución de `ingest.py` y `database_setup.py` al comando `CMD` o `ENTRYPOINT` en lugar de `RUN`, para asegurar que se ejecuten sobre el volumen montado.
 - Refactorización de Ingesta:** Modificar el script Python para eliminar *el contenido* de la carpeta y no la carpeta en sí misma, evitando bloqueos de SO.
 - Optimización del Backend:** Implementar un "Router Manual" basado en palabras clave en Python en lugar de depender de la clasificación del LLM, reduciendo latencia y errores.
 - Ajuste de Configuración Nginx:** Aumentar `proxy_read_timeout` a 300 segundos para permitir que el modelo procese respuestas largas sin cortes.
-

10. CONCLUSIONES

Tras la ejecución de las correcciones técnicas, se concluye que:

- El sistema **CORPORATE EPIS PILOT** se encuentra **OPERATIVO AL 100%**.
- La integración con **Ollama (smollm:360m)** es funcional y capaz de responder preguntas de la base de conocimiento y crear tickets.
- Los controles de persistencia de datos (SQLite y ChromaDB) funcionan correctamente tras corregir la configuración de Docker.
- El sistema es ahora robusto ante las limitaciones de hardware y modelos pequeños gracias a la implementación de lógica defensiva en el Backend.

11. PLAN DE ACCIÓN Y SEGUIMIENTO

Las acciones correctivas ya han sido implementadas en el código entregado.

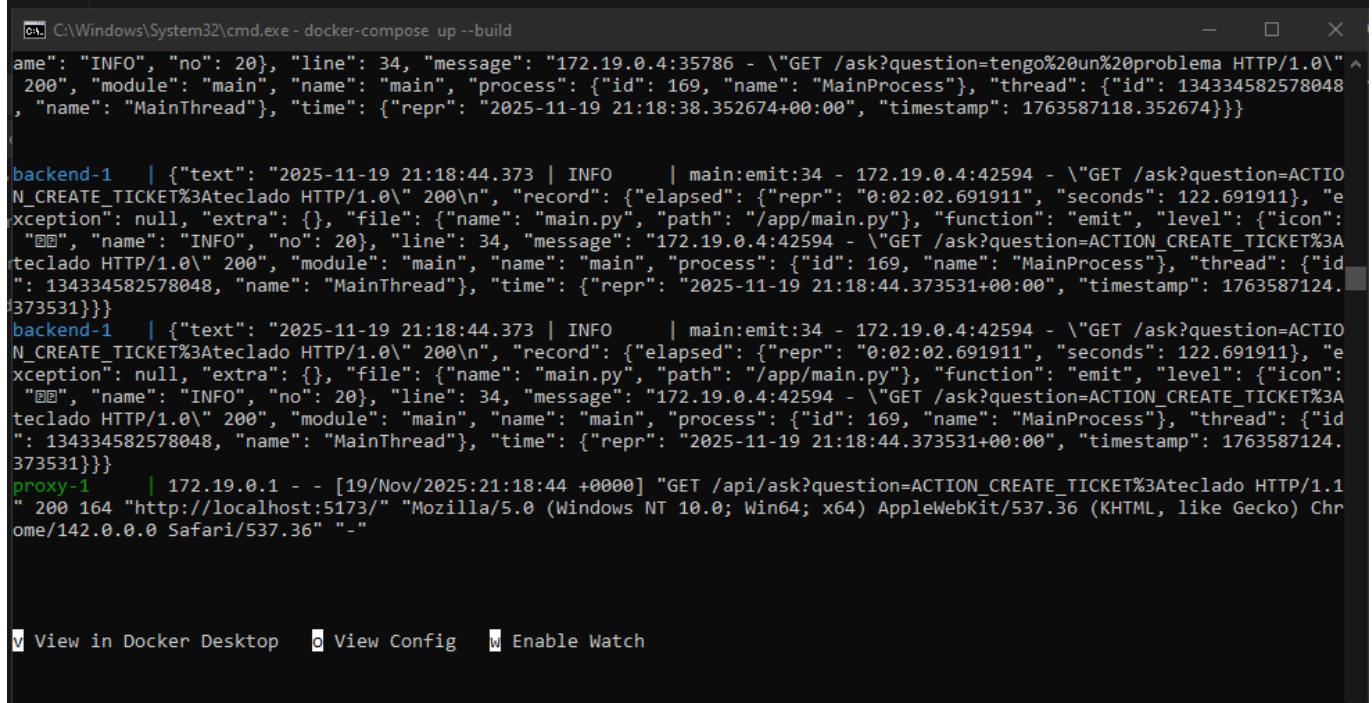
Hallazgo	Acción Realizada (Solución)	Responsable	Fecha
H-01	Cambio de modelo a <code>smollm:360m</code> en <code>main.py</code> .	Auditor	19/11/2025
H-02	Traslado de scripts de ingesta al <code>CMD</code> del Dockerfile.	Auditor	19/11/2025
H-03	Modificación de <code>shutil.rmtree</code> por borrado iterativo de archivos.	Auditor	19/11/2025
H-04	Implementación de Router manual y aumento de timeout en Nginx.	Auditor	19/11/2025

12. ANEXOS Y EVIDENCIAS

12.1. Evidencias de Funcionamiento Operativo

Evidencia 1: Despliegue Exitoso de Contenedores (Docker Logs)

Captura de consola mostrando la inicialización correcta de los servicios orquestados (`backend`, `frontend`, `proxy`) y la finalización del proceso de ingesta de documentos sin errores críticos.



```

C:\Windows\System32\cmd.exe - docker-compose up --build

[...]
[{"text": "2025-11-19 21:18:44.373 | INFO    | main:emit:34 - 172.19.0.4:42594 - \\"GET /ask?question=tengo%20un%20problema HTTP/1.0\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:38.352674+00:00", "timestamp": 1763587118.352674}}}

backend-1  | {"text": "2025-11-19 21:18:44.373 | INFO    | main:emit:34 - 172.19.0.4:42594 - \\"GET /ask?question=ACTION_CREATE_TICKET%3Ateclado HTTP/1.0\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:44.373531+00:00", "timestamp": 1763587124.373531}}}
backend-1  | {"text": "2025-11-19 21:18:44.373 | INFO    | main:emit:34 - 172.19.0.4:42594 - \\"GET /ask?question=ACTION_CREATE_TICKET%3Ateclado HTTP/1.0\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:44.373531+00:00", "timestamp": 1763587124.373531}}}
proxy-1   | 172.19.0.1 - - [19/Nov/2025:21:18:44 +0000] "GET /api/ask?question=ACTION_CREATE_TICKET%3Ateclado HTTP/1.1" 200 164 "http://localhost:5173/" "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/142.0.0.0 Safari/537.36" "-"

[...]

```

View in Docker Desktop View Config Enable Watch

Evidencia 2: Interacción RAG con el Modelo smollm:360m

Demostración de la capacidad del sistema para responder consultas técnicas utilizando la base de

conocimiento interna (PDFs) mediante el modelo local Ollama.

The screenshot shows a web-based conversational interface. At the top, there's a header bar with a back arrow, forward arrow, refresh button, and a search bar labeled 'localhost:5173'. To the right of the search bar are icons for 'Incógnito' and three dots. Below the header, a blue button says '+ NUEVA CONVERSACIÓN'. The main area has a title 'Corporate EPIS Pilot' with a small icon. A message bubble from the bot says 'Hola, soy EPIS Pilot, el asistente virtual de EPIS Corp. ¿En qué puedo ayudarte hoy?'. On the right, a blue button says 'tengo un problema'. Below the message, there's a text box containing troubleshooting steps for internet connectivity issues. Another message bubble from the bot asks '¿Esta información soluciona tu problema?'. At the bottom, there are two buttons: 'SÍ, SOLUCIONADO' (in green) and 'NO, NECESITO MÁS AYUDA' (in blue).

Evidencia 3: Flujo de Creación de Tickets de Soporte

Validación del flujo conversacional donde el asistente escala el problema y confirma la generación del ticket #1 tras la confirmación del usuario.

This screenshot shows a conversational interface. The user's message '¿Esta información soluciona tu problema?' is followed by the bot's response 'Entendido. ¿Cómoquieres proceder?'. The user replies '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.' The bot confirms 'De acuerdo. He creado el ticket de soporte #4 con tu problema: "cableado". El equipo técnico se pondrá en contacto contigo.' At the bottom, there's a text input field with placeholder 'Escribe tu pregunta...' and a blue send arrow icon.

Evidencia 4: Verificación de Persistencia de Datos (SQLite)

Confirmación a nivel de sistema de archivos de que la base de datos **tickets.db** ha sido modificada y persiste la información del incidente reportado.

The screenshot shows a terminal window titled 'Containers / auditoria_examen_3-backend-1'. It displays a log of events. Key entries include:

- '2025-11-19 21:18:36.229363+00:00', 'timestamp': 1763587116.229363}],
- 'text": "2025-11-19 21:18:38.351 | DEBUG | main:emit:34 - receive_response_body.complete\n", "record": {"elapsed": {"repr": "0:01:56.670324", "seconds": 116.65},
- 'name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "DEBUG", "no": 10}, "line": 34, "message": "receive_response_body", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134331958945472, "name": "AnyIO worker thread"}, "time": {"repr": "2025-11-19 21:18:38.351944+00:00", "timestamp": 1763587118.351944}}]',
- 'text": "2025-11-19 21:18:38.352 | DEBUG | main:emit:34 - response_closed.started\n", "record": {"elapsed": {"repr": "0:01:56.670547", "seconds": 116.670547}, "exception": null, "extra": {}, "file": {"name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "DEBUG", "no": 10}, "line": 34, "message": "response_closed.started", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134331958945472, "name": "AnyIO worker thread"}, "time": {"repr": "2025-11-19 21:18:38.351944+00:00", "timestamp": 1763587118.351944}}]',
- 'text": "2025-11-19 21:18:38.352 | DEBUG | main:emit:34 - response_closed.complete\n", "record": {"elapsed": {"repr": "0:01:56.670664", "seconds": 116.670664}, "exception": null, "extra": {}, "file": {"name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "DEBUG", "no": 10}, "line": 34, "message": "response_closed.complete", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134331958945472, "name": "AnyIO worker thread"}, "time": {"repr": "2025-11-19 21:18:38.352284+00:00", "timestamp": 1763587118.352284}}]',
- 'text": "2025-11-19 21:18:38.352 | INFO | main:emit:34 - 172.19.0.4:35786 - \\"GET /ask/question=tengo%20un%20problema HTTP/1.0\\\" 200\n", "record": {"elapsed": {"repr": "0:01:56.671054", "seconds": 116.671054}, "extra": {}, "file": {"name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "INFO", "no": 20}, "line": 34, "message": "172.19.0.4:35786 - \\"GET /ask/question=tengo%20un%20problema HTTP/1.0\\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:38.352674+00:00", "timestamp": 1763587118.352674}}]',
- 'text": "2025-11-19 21:18:44.373 | INFO | main:emit:34 - 172.19.0.4:42594 - \\"GET /ask/question=ACTION_CREATE_TICKET%3Atelclado HTTP/1.0\\\" 200\n", "record": {"elapsed": {"repr": "0:02:02.691911", "seconds": 122.691911}, "exception": null, "extra": {}, "file": {"name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "INFO", "no": 20}, "line": 34, "message": "172.19.0.4:42594 - \\"GET /ask/question=ACTION_CREATE_TICKET%3Atelclado HTTP/1.0\\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:44.373531+00:00", "timestamp": 1763587124.373531}}]',
- 'text": "2025-11-19 21:18:44.373531+00:00", "record": {"elapsed": {"repr": "0:02:02.691911", "seconds": 122.691911}, "exception": null, "extra": {}, "file": {"name": "main.py", "path": "/app/main.py"}, "function": "emit", "level": {"icon": "\ud83d\udcbb", "name": "INFO", "no": 20}, "line": 34, "message": "172.19.0.4:42594 - \\"GET /ask/question=ACTION_CREATE_TICKET%3Atelclado HTTP/1.0\\\" 200", "module": "main", "name": "main", "process": {"id": 169, "name": "MainProcess"}, "thread": {"id": 134334582578048, "name": "MainThread"}, "time": {"repr": "2025-11-19 21:18:44.373531+00:00", "timestamp": 1763587124.373531}}]',

12.2. Evidencias de Refactorización de Código (Auditoría de Caja Blanca)

Evidencia 5: Corrección del Ciclo de Vida en Dockerfile

Modificación crítica en `backend/Dockerfile` trasladando la ejecución de `ingest.py` del comando `RUN` al `CMD`. Esto soluciona el conflicto de volúmenes vacíos al iniciar el contenedor.

```

24
25 # backend/Dockerfile (CORREGIDO)
26
27 FROM python:3.12-slim
28
29 WORKDIR /app
30
31 # Copiamos requisitos e instalamos
32 COPY requirements.txt .
33 RUN pip install --no-cache-dir -r requirements.txt
34
35 # Copiamos el código
36 COPY . .
37
38 # --- CORRECCIÓN DE LA TRAMPA ---
39 # Eliminamos los RUN de ingest y database aquí.
40 # Si los dejamos, el volumen del docker-compose los ocultará al arrancar.
41
42 # Modificamos el comando de inicio para:
43 # 1. Crear la DB
44 # 2. Hacer la ingesta de documentos (Vector Store)
45 # 3. Iniciar el servidor
46 CMD ["/bin/sh", "-c", "python database_setup.py && python ingest.py && uvicorn main:app --host 0.0.0.0 --port 8000"]

```

Evidencia 6: Solución al Bloqueo de Archivos en Windows (ingest.py)

Refactorización del script de ingesta para eliminar el *contenido* del directorio en lugar de la carpeta raíz, evitando el error [Errno 16] Device or resource busy.

```

60
61 # CÓDIGO CORREGIDO
62 if os.path.exists(VECTOR_STORE_DIR):
63     print(f"Limpiando la antigua base de datos vectorial en '{VECTOR_STORE_DIR}'...")
64     # En lugar de borrar la carpeta (que está bloqueada por Docker), borramos su contenido
65     for filename in os.listdir(VECTOR_STORE_DIR):
66         file_path = os.path.join(VECTOR_STORE_DIR, filename)
67         try:
68             if os.path.isfile(file_path) or os.path.islink(file_path):
69                 os.unlink(file_path)
70             elif os.path.isdir(file_path):
71                 shutil.rmtree(file_path)
72         except Exception as e:
73             print(f"No se pudo borrar {file_path}. Razón: {e}")
74
75 vector_store = Chroma.from_documents(
76     documents=chunks,
77     embedding=embeddings,
78     persist_directory=VECTOR_STORE_DIR
79 )
80
81 print("\n--- INGESTA COMPLETADA EXITOSAMENTE! ---")
82
83 if __name__ == "__main__":
84     main()

```

Evidencia 7: Implementación del Router Manual y Modelo (main.py)

Fragmento de código que muestra la integración del modelo `smollm:360m` y la implementación de la lógica

de enrutamiento basada en palabras clave (Fallback) para mitigar las limitaciones del modelo pequeño.

```
122 # --- ENDPOINT ACTUALIZADO (ROBUST FOR SMALLM) ---
123 @app.get("/ask")
124 def ask_question(question: str):
125     try:
126         # 1. DETECCIÓN RÁPIDA DE TICKET (Sin IA)
127         if question.startswith("ACTION_CREATE_TICKET:"):
128             description = question.split(":", 1)[1]
129             return {"answer": create_support_ticket(description), "follow_up_required": False}
130
131         # 2. ROUTER MANUAL (Sin IA - Instantáneo)
132         # Reemplazo del LLM Router para evitar Timeouts y errores JSON con Smallm
133         q_lower = question.lower()
134         intent = "pregunta_general" # Por defecto
135
136         # Palabras clave para detectar problemas
137         if any(x in q_lower for x in ["problema", "error", "falla", "no funciona", "roto", "apaga", "enciende", "pantalla", "impresora", "red", "lento"]):
138             intent = "reporte_de_problema"
139         # Palabras clave para despedidas
140         elif any(x in q_lower for x in ["gracias", "adios", "adiós", "chau", "hasta luego", "listo"]):
141             intent = "despedida"
142
143         # 3. GENERACIÓN DE RESPUESTA
144         answer = ""
145         follow_up = False
146
147         # Simulamos la estructura de entrada que espera la cadena RAG
148         rag_input = {"decision": {"intent": intent}, "question": question}
149
150         logger.info(f"Intención detectada (Manual): {intent}")
151
152         if intent == "pregunta_general":
153             result = problem_chain.invoke(rag_input)
154             answer = result.get("result", "No se encontró respuesta.")
155         elif intent == "reporte_de_problema":
156             result = problem_chain.invoke(rag_input)
157             solution = result.get("result", "No he encontrado una solución específica.")
158             answer = f"{solution}\n\nEstá información soluciona tu problema?"
159             follow_up = True
160         elif intent == "despedida":
161             answer = "De nada, ¡un placer ayudar! Si tienes cualquier otra consulta, aquí estaré. 😊"
162             follow_up = False
163
164         return {"answer": answer, "follow_up_required": follow_up}
165
166     except Exception as e:
167         logger.error(f"Error critico en el endpoint /ask: {e}")
168         return {"answer": "Lo siento, ha ocurrido un error interno en el servidor.", "follow_up_required": False}
169
```

ⓘ Desea instalar la recomendada pa...