

Informe Técnico: SurfSense – Plataforma IA de Investigación Open-Source

1. Génesis y Filosofía del Proyecto

Origen e inspiración: SurfSense nació en 2024 como iniciativa de un desarrollador independiente (alias "Uiqueblhats" en Reddit) para crear un asistente de investigación personal al navegar la web ¹. La motivación inicial fue resolver el "brain freeze" de recordar gran cantidad de contenido guardado durante la navegación ¹. SurfSense se concibió como un "asistente IA personal para cualquier cosa que veas en la web" ², permitiendo capturar páginas visitadas mediante una extensión de navegador y construir una **base de conocimiento personal** con ellas ¹. Su creador se inspiró en herramientas como Google NotebookLM y Perplexity, buscando ofrecer capacidades similares pero en un entorno **privado y personalizable** ³. De hecho, el lema del proyecto es tener *"tu propio NotebookLM/Perplexity privado, integrado con tus fuentes de conocimiento"* ⁴. Inicialmente, SurfSense puso énfasis en no olvidar nada visto en Internet, integrando una **gráfica de conocimiento personal** ("Knowledge Graph Brain") para relacionar información ⁵ ⁶.

Comparativa con NotebookLM y Perplexity: SurfSense toma como referencia las funciones de Google NotebookLM (IA sobre documentos personales) y Perplexity AI (búsqueda con IA y citas), pero busca *"elevar esas capacidades"* añadiendo integración con la **base de conocimiento personal del usuario** ³. Al igual que NotebookLM, permite **cargar y consultar documentos propios**, y al estilo de Perplexity genera respuestas con **referencias citadas** ⁷. Sin embargo, SurfSense amplía el alcance al conectarse con fuentes externas diversas (por ejemplo Slack, Notion, motores de búsqueda, GitHub, YouTube, etc.) que no están presentes en NotebookLM ⁸ ⁹. A diferencia de servicios cerrados, SurfSense es **auto-hospedable** y soporta LLMs locales, proporcionando privacidad total de los datos ¹⁰. También soporta más formatos de archivos (27+ extensiones de documentos, imágenes, audio, etc. vs. los PDF limitados de NotebookLM) ¹¹ ¹². Otro punto diferenciador es su extensión de navegador para capturar páginas web (incluidas páginas detrás de autenticación) sin recurrir a scraping tradicional – lee directamente el DOM para mayor precisión ¹³. En resumen, SurfSense *emula* la experiencia de NotebookLM/Perplexity (IA conversacional con citas) **mejorando** la integración con datos privados, la personalización de modelos y la posibilidad de funcionar offline o a bajo costo (por ejemplo con modelos locales en lugar de depender de GPT-4 en la nube) ¹⁴.

Visión a largo plazo: El equipo de SurfSense (principalmente su creador y una creciente comunidad) manifiesta en Discord, GitHub y foros que la visión es convertirlo en un *"asistente de investigación personal universal"*. En su servidor de Discord invitan a la comunidad a *"dar forma al futuro de SurfSense"* activamente ¹⁵. Entre los planes a futuro destaca la incorporación de **más conectores** (integraciones con otras herramientas y plataformas) ¹⁶, mejorar características pendientes (por ejemplo, rehacer la UI de Podcasts) y optimizar la experiencia de chat con documentos ¹⁷. El desarrollador ha sido receptivo a sugerencias – por ejemplo, ante pedidos de la comunidad incorporó rápidamente soporte Docker para facilitar la instalación ¹⁸ ¹⁹. En discusiones recientes se mencionó la idea de crear un "MCP" (*Master Control Program*), es decir, un *hub* central para orquestar múltiples herramientas IA desde SurfSense ²⁰. También planean funciones avanzadas como análisis de contenido de video (diapositivas en YouTube) a largo plazo, aunque conscientes de los costos, posiblemente habilitándolas como opciones bajo demanda ²¹ ²². En esencia, la visión de SurfSense es ser una **plataforma abierta y extensible** que integre IA de lenguaje con **razonamiento simbólico** y múltiples fuentes de

conocimiento, sirviendo como base para experimentación y aplicaciones personalizadas en el ámbito de la investigación asistida por IA.

2. Arquitectura Técnica

Backend (FastAPI) y endpoints críticos: El núcleo de SurfSense es un backend Python construido sobre FastAPI, que expone una API REST para todas las funciones ²³. Este backend maneja la autenticación de usuarios (mediante JWT/OAuth2, integrado con FastAPI Users ²⁴), la ingestión de datos y la gestión del flujo de preguntas y respuestas. Entre los endpoints críticos se incluyen: endpoints de **ingesta** (para subir archivos o recibir páginas web desde la extensión, convirtiéndolos en fragmentos de texto e insertando en la base de conocimiento), endpoints de **embeddings** (para generar y almacenar vectores de nuevos documentos), endpoints de **consulta** (que reciben la pregunta del usuario, desencadenan el pipeline de búsqueda y generación de respuesta) y endpoints para gestionar **fuentes externas** (por ejemplo, sincronizar datos de Slack, Notion u otras integraciones). La lógica de aplicación está organizada en servicios y routers de FastAPI; por ejemplo, existe un módulo principal (e.g. `main.py`) que inicia la aplicación Uvicorn ²⁵ y diferentes sub-rutas para funcionalidades (autenticación, documentos, búsqueda, etc.). El backend también expone una documentación Swagger (en `/docs`) que facilita probar operaciones como registrar un usuario o realizar consultas ²⁶. Gracias a FastAPI, el diseño es modular, permitiendo añadir nuevos endpoints de forma relativamente sencilla, y soporta respuestas **streaming** para el chat en tiempo real (integrado con el protocolo de streaming de Vercel en el front) ²⁷.

Frontend (Next.js 15) – modularidad y extensibilidad: SurfSense provee una interfaz web construida con Next.js 15 (App Router) y React 19, escrita en TypeScript ²⁸. La arquitectura del front end está dividida en **módulos** o páginas para las distintas secciones: por ejemplo, una sección de “Spaces” o colecciones de búsqueda, una para administrar documentos cargados, otra para el chat de investigación, etc. (imágenes divulgadas muestran pantallas de *Search Spaces*, *Manage Documents*, *Research Agent*, *Agent Chat* ²⁹ ³⁰). Esta separación de lógica en componentes React facilita la extensión de la UI: es posible añadir nuevas vistas o paneles (por ejemplo, integrar una nueva fuente de datos) sin afectar al resto. El uso de **Tailwind CSS** y librerías UI headless (Shadcn) agiliza la creación de componentes de interfaz consistentes ³¹ ³². Además, el front end utiliza el *Vercel AI SDK* (Kit UI Stream) para manejar eficientemente las respuestas en streaming del backend ²⁷, lo que indica que la lógica de recepción de mensajes del modelo y su renderizado en forma de chat está bien abstraída. Esto permitiría, por ejemplo, ajustar fácilmente el formato de visualización de las respuestas o insertar nuevas funcionalidades de interacción en la conversación. En términos de extensibilidad, dado que Next.js soporta incrementally adding API routes, uno podría crear endpoints front-end adicionales para nuevas integraciones y consumir nuevos endpoints del backend. En resumen, la elección de Next.js 15 + React moderna proporciona una base sólida, con **código dividido automáticamente**, componentes de estado (React Hook Form para formularios, etc. ³³) y un entorno familiar para desarrolladores, haciendo más sencillo **personalizar o ampliar** la interfaz de SurfSense según las necesidades.

Pipeline de RAG (Retrieval-Augmented Generation): SurfSense implementa un robusto *pipeline* de recuperación y generación que combina técnicas avanzadas de **RAG**. Emplea tanto **LangChain** como **LangGraph** para la orquestación de agentes y flujos ³⁴. En esencia, cuando el usuario realiza una pregunta, el backend sigue estos pasos (simplificados):

1. **Búsqueda híbrida:** La pregunta se utiliza para realizar una búsqueda tanto semántica como de texto completo en la base de conocimiento. Para ello, SurfSense almacena embeddings vectoriales de cada fragmento de documento en PostgreSQL (vía pgvector) y también indexa texto para búsqueda lexical. Al consultar, ejecuta una búsqueda de **similitud vectorial**

(embeddings) y una búsqueda tradicional (palabras clave) en paralelo ³⁵ ³⁶ . Luego combina los resultados usando **Reciprocal Rank Fusion (RRF)**, técnica que fusiona ambos rankings para obtener los documentos más relevantes en conjunto ³⁵ ³⁷ . Esto asegura que se contemplen tanto la relevancia semántica como la coincidencia exacta de términos.

2. **Índices jerárquicos (2-tier RAG):** SurfSense aplica un enfoque jerárquico en la recuperación ³⁸ . Primero identifica qué **documentos o fuentes** son más relevantes, luego dentro de esos documentos recupera los **fragmentos** específicos. Esta “doble capa” aumenta la precisión al filtrar contexto: por ejemplo, si la base tiene libros enteros, el sistema encuentra los libros relevantes y luego pasajes dentro de ellos en lugar de buscar a nivel de párrafo globalmente. Este concepto de *Hierarchical Index* está incorporado en el pipeline ³⁸ .

3. **Re-ranking avanzado:** Adicional a RRF, se menciona soporte para *rerankers* externos (p. ej. modelos de Cohere Rerank o **FlashText/Flashrank**) que pueden refinar el orden de resultados ³⁹ . Esto sugiere que tras obtener candidatos, SurfSense puede aplicar un modelo de reranqueo (por ejemplo, un modelo cross-encoder) para ordenar las piezas de contexto por relevancia antes de pasarlas al LLM. De hecho, se configuran variables de entorno para el modelo de reranqueo a usar (nombre y tipo) ⁴⁰ , e.g. `RERANKERS_MODEL_NAME` = `"ms-marco-MiniLM-L-12-v2"` y tipo `flashrank` para usar un modelo MiniLM de Microsoft ⁴⁰ .

4. **Llamada final al LLM:** Una vez obtenido el contexto más relevante (p. ej. los 3-5 fragmentos mejor puntuados), el agente IA elabora un *prompt* que incluye la pregunta del usuario y dichos fragmentos (citando sus fuentes). Luego realiza la consulta a un **LLM** para generar la respuesta final con citas. SurfSense ha abstraído esta capa mediante **LiteLLM**, una librería que permite dirigir las llamadas a múltiples modelos de lenguaje de forma intercambiable ⁴¹ . En la práctica, el “investigador” de SurfSense es un agente LangGraph que utiliza herramientas (búsqueda, lectura de docs) y finalmente invoca al LLM a través de LiteLLM ⁴² . Para el desarrollador, esto significa que el *punto de intervención* donde puede personalizar el comportamiento es precisamente antes de la llamada al LLM – por ejemplo, modificando el *prompt* final o incluso sustituyendo la función de generación. Dado que LiteLLM soporta “cualquier LLM que puedas imaginar” (OpenAI, locales vía Ollama, Azure, etc.) ⁴² , **reemplazar el modelo** es tan sencillo como cambiar una cadena de configuración (como se detalla en la sección de customización). Incluso es posible definir *roles de LLM* en la configuración para distintos propósitos: SurfSense permite especificar un LLM rápido para respuestas breves, otro más potente para razonamiento complejo y otro especializado en contexto largo ⁴³ . Esta modularidad hace viable interceptar la llamada final: por ejemplo, un desarrollador podría apuntar la llamada a un servidor propio que aplique razonamiento simbólico extra antes de devolver la respuesta, *inyectándose* en el pipeline estándar. En resumen, el pipeline RAG de SurfSense es altamente configurable y ofrece varios puntos de extensión, garantizando respuestas **contextualizadas con citas** como en Perplexity, pero con la posibilidad de adaptar cada etapa (búsqueda, re-rank, selección de LLM) a necesidades particulares ⁴⁴ .

Base de datos PostgreSQL + pgvector: Todo el almacenamiento persiste en una base de datos PostgreSQL que combina datos tradicionales y vectores. SurfSense emplea el módulo **pgvector** para manejar embeddings, permitiendo realizar consultas vectoriales (similaridad) directamente en SQL ⁴⁵ ³⁷ . El esquema de la base incluye tablas para **usuarios** (manejado por FastAPI Users), **documentos** o elementos de conocimiento, **fragmentos o chunks** de texto con sus embeddings, y posiblemente tablas auxiliares para registrar fuentes o metadatos (como a qué “Space” o colección pertenece cada documento). El ORM SQLAlchemy facilita la definición de estos modelos y la portabilidad de la base de datos ⁴⁶ , mientras que Alembic se usa para migraciones del esquema ⁴⁷ . Un flujo típico de datos sería: cuando se sube un archivo o se guarda una página web, el backend crea entradas en la DB para el

documento (con su título, fuente, etc.), lo trocea en párrafos o porciones (utilizando la librería Chonkie, ver abajo) y calcula un **vector de embedding** para cada chunk, almacenándolo en un campo vectorial provisto por pgvector ⁴⁵ ⁴⁸ . Estos vectores luego se indexan (pgvector permite índices aproximados tipo IVF) para una recuperación rápida. PostgreSQL también almacena el texto completo, permitiendo las búsquedas lexicográficas. SurfSense incluye incluso **pgAdmin** en su entorno Docker para administrar la base fácilmente ⁴⁹ . En cuanto a rendimiento, se eligió pgvector para evitar dependencias adicionales (como Pinecone o milvus) manteniendo todo el stack en una base unificada; el trade-off es un rendimiento menor en colecciones masivas, pero en contextos personales suele ser suficiente. En definitiva, la DB actúa como **almacén vectorial + relacional** central: las consultas de embeddings se traducen a SELECT con operadores de distancia, combinadas con consultas de texto completo, todo gestionado dentro de Postgres ⁴⁵ ⁵⁰ . Esta integración simplifica el despliegue (un solo servicio de datos) a costa de requerir configurar la extensión pgvector, lo cual es sencillamente un paso de instalación adicional ⁵¹ .

Agente de Podcasts: Una de las características distintivas de SurfSense es un agente capaz de convertir chats o contenidos en **podcasts de audio** de forma casi instantánea. Arquitectónicamente, este “Podcast Agent” toma el historial de la conversación o un resumen de contenido guardado y lo pasa por un proceso de dos fases: (1) **Generación de guion** y (2) **Síntesis de voz (TTS)**. En la fase de guion, se utiliza un LLM (posiblemente una variante optimizada, incluso se menciona un agente GPT-researcher específico ⁵²) para reescribir la conversación en un formato narrativo adecuado para audio (por ejemplo, convirtiendo preguntas y respuestas en un monólogo o diálogo estructurado). Luego, ese guion se envía a un servicio de Text-to-Speech para obtener el audio final. SurfSense soporta múltiples proveedores de TTS: OpenAI (vos voces de *Whisper* o voces nuevas de GPT-4), Azure Cognitive Services (voz neural) y Google Vertex AI, configurables mediante variables de entorno ⁵³ ⁵⁴ . Gracias a la eficiencia del pipeline, pueden generar ~3 minutos de audio en menos de 20 segundos ⁵⁵ , un rendimiento notable que sugiere que usan modelos rápidos o resumidos para el guion. El agente maneja también Speech-to-Text (STT) si fuese necesario, pero principalmente se enfoca en producir audio (*de texto a voz*). La arquitectura aprovecha LiteLLM para unificar las llamadas a TTS y STT (también vía proveedores) ⁵⁶ . Actualmente, la funcionalidad de podcast está **marcada como temporalmente en revisión** en la hoja de ruta ⁵⁷ – están mejorando la interfaz y estabilidad antes de relanzarla. No obstante, el código ya implementa la lógica básica. En suma, el agente de podcasts ejemplifica la extensibilidad de SurfSense: usando la misma conversación guardada, aplica un agente IA para formatear un nuevo producto (audio) de manera **automática y veloz** ⁵⁵ . Es una manera innovadora de reutilizar el conocimiento acumulado, convirtiendo las respuestas del asistente en un medio más consumible. Cuando se reactive, permitirá a los usuarios, por ejemplo, obtener una síntesis hablada de sus sesiones de investigación en segundos, con potencial para elegir diferentes voces o estilos (dado el soporte multi-proveedor de TTS) ⁵⁴ .

(Nota: A noviembre de 2024, el módulo de podcast estaba deshabilitado para mejoras de UI ⁵⁸ ; se espera su reintroducción pronto con mejor experiencia de usuario.)

3. Customización y Extensibilidad

Uso con modelos locales (Ollama u otros): SurfSense fue diseñado para ser agnóstico al modelo de lenguaje subyacente, lo que facilita enormemente usar LLMs personalizados o locales. En particular, soporta *de fábrica* la integración con **Ollama**, un servidor de modelos locales, y en general cualquier

modelo accesible vía API OpenAI-like mediante LiteLLM ⁴². Para configurar SurfSense con un LLM local, los pasos son los siguientes:

- 1. Instalación sin Docker:** Dado que la integración con Ollama no funciona dentro de Docker (limitación conocida) ⁵⁹, se debe desplegar SurfSense en modo manual o bare-metal. Es decir, instalar los componentes directamente en el host.
- 2. Instalar y ejecutar Ollama:** Ollama es un servidor que ejecuta modelos como LLaMA, Qwen, etc., localmente. José deberá instalar Ollama en su máquina y descargar el modelo deseado (por ejemplo, Qwen-3 si estuviera disponible, o cualquier LLM soporte de su elección). Luego iniciar el servicio (por defecto escucha en `localhost:11434`). Alternativamente, podría usar otro runtime local (ej: LocalAI, text-generation-webui) mientras exponga una API compatible con OpenAI o un endpoint para LiteLLM.
- 3. Configurar LiteLLM en SurfSense:** En el archivo de configuración del backend (`.env` en `surfsense_backend`), se definen las rutas de modelos para LiteLLM. SurfSense permite especificar hasta tres LLM diferentes para distintos propósitos: `FAST_LLM` (modelo rápido para respuestas breves), `STRATEGIC_LLM` (modelo principal para razonamiento complejo) y `LONG_CONTEXT_LLM` (modelo con ventana de contexto amplia) ⁴³. Para usar un modelo local, se debe apuntar uno de estos a la instancia de Ollama. Por ejemplo, se puede setear `STRATEGIC_LLM = "ollama/Qwen-3"` (suponiendo que Ollama tenga un modelo llamado "Qwen-3" registrado). De hecho, la sintaxis suele ser `"ollama/<nombre_modelo>"` para indicar a LiteLLM que use la API de Ollama ⁴³. Igualmente, si se quisiera usar un modelo de HuggingFace local vía otra herramienta, LiteLLM documenta cómo integrar proveedores (vía URL de base, etc.) ⁶⁰.
- 4. Claves y endpoints:** Si el modelo local no requiere API key (p.ej. Ollama no necesita clave), no hace falta configurar token. Pero si se usan modelos de nube en paralelo (ej: GPT-4 de OpenAI), se deben incluir sus API keys en el `.env` correspondiente (`OPENAI_API_KEY`, `GEMINI_API_KEY`, etc., según el proveedor) ⁶⁰. En el caso puramente local, probablemente José omitirá las variables de OpenAI y solo definirá las de Ollama. LiteLLM detectará la presencia de un modelo "ollama" y enrutará las peticiones ahí. (Nota: En Docker, como se dijo, Ollama no está soportado; por tanto, la ejecución debe ser con `uvicorn` localmente.)
- 5. Ejecución y prueba:** Con todo configurado, se lanza el backend (`uvicorn main:app ...`) y el frontend (`pnpm run dev` en `surfsense_web`) en modo local. Al hacer una consulta en la UI, el backend llamará al modelo local a través de LiteLLM. Por ejemplo, si se configuró `FAST_LLM = openai/gpt-4o-mini` y `STRATEGIC_LLM = ollama/Qwen3`, es esperable que SurfSense use primero el modelo mini (vía OpenAI API) para una respuesta rápida/resumen, y luego use Qwen3 en Ollama para la respuesta final elaborada ⁴³ ⁶⁰. José puede verificar en los logs que las peticiones salen hacia `localhost:11434` (endpoint de Ollama) en lugar de `api.openai`.

Estos pasos permiten que **LLMs personalizados** (ej. Qwen-3 fine-tuneado o Llama 3 entrenado por José) sean el motor de SurfSense sin modificar el código, únicamente mediante configuración. La flexibilidad de LiteLLM hace que prácticamente cualquier modelo con API (local o remoto) pueda integrarse – "cualquier LLM que se te ocurra" según el autor ⁴². Por último, recordar que en la instalación manual conviene tener hardware adecuado si el modelo es grande, y ajustar parámetros como *LateChunker* (tamaño de fragmento) si el modelo admite contextos más largos.

Añadir conectores de datos externos: SurfSense ya soporta múltiples conectores (Slack, Notion, Linear, motores de búsqueda personalizados, YouTube, GitHub, etc.) ⁸ ⁹, y su arquitectura permite agregar nuevos con relativa facilidad. Los conectores existentes sirven de ejemplo: por lo general consisten en un **módulo backend** que sabe interactuar con la API externa para obtener datos, y luego formatea esos datos como documentos ingestables en SurfSense. Por ejemplo, el conector de Slack probablemente use la Web API de Slack (con un token de bot) para leer mensajes de ciertos canales, luego los convierte en texto y los envía al pipeline de ingesta (quizá asignándolos a un “Space” dedicado a Slack). De igual manera, el de Notion utilizará la API de Notion para extraer páginas o bases de datos en formato markdown.

Para agregar un nuevo conector (digamos, para una herramienta propietaria de “mapas mentales” que José maneja), habría que seguir estos pasos a nivel de código:

- **Backend:** Crear un **endpoint** o tarea programada que se conecte a la API de la nueva fuente. Esto implicará autenticarse (por ejemplo, con una API key u OAuth de ese servicio) y luego **extraer los datos** (p. ej., obtener la lista de elementos del mapa mental y sus contenidos). Este código residiría en `surfsense_backend`, posiblemente en un submódulo tipo `external_sources` o dentro del router de ingesta. Podría basarse en FastAPI background tasks si es algo programado, o en un endpoint REST que se activa bajo demanda.
- **Procesamiento e ingesta:** Los datos obtenidos deben transformarse a un formato aceptado – típicamente texto plano, posiblemente estructurado en secciones. Si el contenido está en HTML o formato propio, habría que parsearlo (aquí se podría reutilizar las utilidades de *Unstructured.io* o *LlamaIndex* si aplica, o escribir un parser). Luego se crearían entradas en la DB para cada fragmento, generando embeddings igual que con un documento subido. Es decir, reutilizar las funciones de ingesta ya existentes: probablemente exista una función `ingest_document(text, metadata)` que toma texto y lo indexa. Se podría llamar a esa función para cada elemento extraído del nuevo conector.
- **Frontend (opcional):** Si se desea exponer control al usuario, habría que añadir en la interfaz opciones para configurar y activar el nuevo conector. Por ejemplo, un formulario para ingresar la API key o credenciales de la nueva fuente, botones para “Sincronizar ahora” o seleccionar qué datos traer. La UI de SurfSense es modular, así que se podría añadir una pestaña en la sección de “Sources” o en configuración.

En cuanto a *rutas de código exactas*, habría que identificar dónde están implementados Slack y demás. Aunque no tenemos la lista de archivos, es probable que en el backend existan archivos como `slack_service.py` o similares. Extender SurfSense con otro conector implicará **imitar la estructura** de uno de los existentes. Por ejemplo, si Slack está implementado como un **worker que guarda mensajes** periódicamente, replicar esa lógica para la nueva fuente. Los *APIs externos* a modificar son simplemente los de la nueva herramienta: se usará su SDK o API HTTP para obtener datos. SurfSense no requiere modificar su API interna para soportar nuevos datos; más bien se agrega código que **consume la API externa** y vierte resultados en SurfSense. En síntesis, la arquitectura abierta permite añadir conectores escribiendo unas cuantas funciones de extracción y llamando al pipeline de almacenamiento. Además, en *Future Work* ya indican “Add More Connectors” como prioridad ¹⁶, por lo que el diseño seguramente anticipa esta extensibilidad.

Punto de inyección o reemplazo del LLM: Si bien cambiar el modelo de lenguaje vía configuración es sencillo, puede haber casos en que José quiera interceptar la llamada al LLM para introducir lógica adicional – por ejemplo, pasar la pregunta y contexto por un *motor de inferencia basado en grafo de*

conocimiento propio antes de obtener la respuesta final. Afortunadamente, gracias a LangGraph/LangChain, el flujo de agentes es transparente y modificable. El *entry point* para la llamada estándar al LLM es la configuración de LiteLLM mencionada (FAST_LLM, etc.) ⁶⁰. Una forma “limpia” de interceptar es aprovechar que LiteLLM permite definir un **base URL personalizado** para cualquier proveedor ⁶⁰. Por ejemplo, José podría tener un servidor local cuyo endpoint actúe como un pseudo-LLM: cuando recibe el prompt de SurfSense, en vez de simplemente completar texto, podría realizar razonamiento simbólico consultando sus grafos y luego generar una respuesta. Configurando `STRATEGIC_LLM` con un identificador apuntando a ese endpoint (vía docs de LiteLLM) ⁶⁰, SurfSense enviará allí las peticiones de completado. En otras palabras, se “engaña” al sistema haciéndole creer que es un modelo más. Esta es una vía sin tocar el código.

Otra opción es **modificar el agente LangGraph** directamente: LangGraph permite componer pasos, por lo que José podría insertar un paso extra en el agente de investigación. Por ejemplo, después de la búsqueda de documentos y antes de la pregunta final al LLM, agregar un paso que llame a su motor de grafos con la información recopilada, obteniendo conclusiones adicionales, y luego incluir esas conclusiones en el prompt final. Esto requeriría editar el código del agente (posiblemente definido en algún JSON o clase Python dentro de `surfsense_backend`). Dado que LangGraph es un framework orientado a flujos, debería ser posible añadir una herramienta custom (una función de Python anotada) que realice la consulta a la base de conocimiento grafo de José, y registrarla en el agente. Sin entrar en demasiado detalle, la capacidad está ahí gracias a la naturaleza modular de LangChain/LangGraph ³⁴.

En resumen, el *hook* para interceptar la llamada al LLM estándar es flexible: bien sea configurando LiteLLM para redirigir hacia un **modelo/middleware personalizado**, o alterando la secuencia del agente. La primera vía (LiteLLM) es sencilla: por ejemplo, apuntar `OPENAI_API_BASE` a la URL de un servidor proxy propio ⁶⁰. La segunda requiere forkar ligeramente el repo. Cualquiera de las dos le permitirá a José **reemplazar la lógica de generación** por la suya, integrando su LLM afinado o incluso un pipeline neuro-simbólico completo, manteniendo el resto de funcionalidades de SurfSense (ingesta, búsqueda, UI) intactas. Esto convierte a SurfSense en una excelente plataforma base sobre la cual montar sistemas experimentales de IA integrada con grafos.

4. Ecosistema y Comunidad

Actividad en GitHub: El repositorio de SurfSense (ModSetter/SurfSense) ha mostrado una actividad constante y creciente. A octubre de 2024 contaba con cientos de commits (>300) ⁶¹ y se volvió popular en la comunidad open-source de LLMs, llegando a obtener *miles de estrellas* en GitHub (reportes indican ~4.7k estrellas hacia finales de 2024) ⁶². Esto refleja un interés significativo y aportes de múltiples desarrolladores. Aunque el desarrollo principal proviene del autor original, hay forks activos (por ejemplo, *xpertdev/SurfSense* fue un fork temprano notable) y usuarios que han contribuido con *issues* y *pull requests*. Al revisar los issues, se observan solicitudes de funcionalidades y reportes de bugs que generalmente reciben respuesta rápida. Por ejemplo, existe un *feature request* para soporte de *podcasts de larga duración* (#87) abierto por el mismo ModSetter ⁶³, indicando que el autor lleva un registro público de futuras mejoras. También se han visto PRs de la comunidad corrigiendo pequeños errores o añadiendo documentación. Un fork destacado es el de la organización Decentralised-AI, que clonó el proyecto difundiendo su descripción, lo que ayudó a visibilizarlo ⁶⁴. En general, la **salud del repositorio** es buena: actualizaciones frecuentes, integración continua configurada (GitHub Actions) ⁶⁵, aunque aún no han publicado releases formales (0 releases publicados hasta la fecha en GitHub ⁶⁶, el proyecto está en fase activa de desarrollo).

Soporte en Discord: SurfSense mantiene una comunidad en Discord donde el desarrollador y usuarios interactúan. La invitación a Discord está presente en la web oficial ⁶⁷ y en Devpost se anima a unirse

para colaborar ¹⁵. Según testimonios y la actividad en Reddit, el creador es **muy receptivo** al feedback. Vimos un ejemplo en Reddit donde tras sugerencias de usuarios sobre Docker, respondió el mismo día comprometiéndose a agregar soporte y efectivamente lo hizo poco después ¹⁸ ¹⁹. Es de esperar que en Discord ese dinamismo sea aún mayor, con respuestas directas a dudas de instalación, etc. La calidad del soporte parece alta dada la dedicación del desarrollador (quien incluso reescribió gran parte del backend/frontend cuando el proyecto creció más rápido de lo esperado ⁴², mostrando compromiso con la mantenibilidad). En Discord probablemente se discuten nuevas ideas (como la integración con Zapier que un usuario mencionó y el dev consideró ⁶⁸ ⁶⁹) y se brinda ayuda mutua. Aunque no tenemos citas directas del Discord, la invitación abierta y la alusión a “ayúdanos a acelerar el desarrollo” ¹⁵ sugieren que los colaboradores son bienvenidos y que hay transparencia en la hoja de ruta dentro de la comunidad.

Contribuidores externos destacados: Al ser un proyecto relativamente joven pero muy popular, han emergido algunos contribuidores notables. Uno de ellos es posiblemente **xpertdev**, cuyo fork se menciona como origen en GitHub ⁷⁰ – podría tratarse de un co-desarrollador temprano o simplemente del dueño de un fork que luego el autor fusionó. También la comunidad *AI Automators* ha mostrado interés (difundiendo InsightsLM en paralelo, ver sección 5). Sin embargo, el principal **mantenedor** es ModSetter (alias del autor en GitHub), quien impulsa prácticamente todas las funcionalidades clave. Entre los externos, se han visto aportes en issues por usuarios como `cionut` (sugiriendo mejoras en ingestión de video) ²¹, o `Even_End2275` (discutiendo integraciones adicionales) ⁶⁸. Estos nombres indican la participación de expertos/profesionales de distintos ámbitos aportando ideas. También cabe destacar que SurfSense participó en hackathons (tiene una página en Devpost ⁷¹), lo cual pudo atraer contribuidores o al menos testers de la comunidad hacker. En general, aunque **la comunidad está creciendo**, aún no hay “contribuidores estrella” claramente identificados aparte del autor. Posiblemente en el futuro forks significativos surjan (como adaptar SurfSense para otras bases de datos, etc.), pero hasta 2024 el enfoque parece centrado en consolidar el proyecto base con la guía del creador original.

En resumen, el ecosistema alrededor de SurfSense es **vibrante**: fuerte respaldo de la comunidad de IA local (ver hilos entusiastas en r/LocalLLM y r/selfhosted), soporte directo vía Discord, un GitHub activo, y un proyecto que pese a ser joven ya se considera uno de los más prometedores en el espacio de asistentes de investigación open-source. Esta comunidad activa reduce riesgos de abandono y asegura que nuevas ideas (p.ej. “¿puedo integrarlo con Lyzr?” ⁶⁸) se discutan y eventualmente se implementen, manteniendo a SurfSense en la frontera de las herramientas de investigación asistida por IA.

5. Comparativa Estratégica

SurfSense vs. InsightsLM: InsightsLM es otro proyecto open-source inspirado en NotebookLM, con el objetivo declarado de ser una alternativa auto-hospedada y personalizable ⁷². La principal diferencia radica en la **arquitectura y público objetivo**. Mientras SurfSense es un proyecto de código orientado a desarrolladores (requiere configurar entornos Python/Node, BD, etc.), InsightsLM fue construido con una filosofía *no-code*: utiliza Supabase (base de datos y autenticación en la nube) y n8n (plataforma de automatización visual) para la lógica backend ⁷³. Esto hace que InsightsLM sea más fácil de desplegar para usuarios sin conocimientos de programación (sus creadores enfatizan que se puede personalizar “sin escribir ni una línea de código” ⁷⁴), a costa de depender de servicios externos y posiblemente ser menos flexible en el largo plazo en cuanto a añadir funcionalidades complejas.

En **capacidades**, hay solapamiento pero también diferencias: Ambos permiten *chatear con documentos y obtener citas verificables* ⁷⁵. InsightsLM menciona también generación de **audio resúmenes** estilo podcast ⁷⁶, similar a SurfSense, por lo que busca paridad en características principales. Sin embargo,

SurfSense ofrece integraciones más amplias (múltiples fuentes externas como Slack, GitHub, buscadores, etc.) que InsightsLM no menciona explícitamente. InsightsLM parece centrarse más en documentos de empresa y casos de uso de base de conocimiento corporativa ⁷⁷, mientras SurfSense apunta a una gestión de conocimiento personal más diversa (favoritos web, notas, videos, código). Además, en soporte de modelos: SurfSense soporta 150+ LLMs y 6000+ modelos de embeddings mediante LiteLLM y AutoEmbeddings ⁷⁸ ³⁹, incluyendo locales; InsightsLM, al menos en su anuncio, no detalla tantos modelos – presumiblemente usa OpenAI por defecto pero es extensible a modelos locales si se configura (dice “usa modelos locales si lo deseas” ⁷⁹).

Otra distinción es que SurfSense implementa técnicas avanzadas como **RRF, índices jerárquicos y rerankers** ³⁵ ³⁹, lo que le da posiblemente mayor precisión en recuperación. InsightsLM, al estar construido sobre n8n, probablemente use un pipeline más sencillo (quizá una búsqueda vectorial en Supabase y luego llamada al LLM). Esto podría impactar la calidad de respuestas: SurfSense puede ofrecer resultados más relevantes gracias a ese pipeline optimizado, a costa de mayor complejidad de infraestructura. En cuanto a **extensibilidad estratégica**, ambos son open-source, pero la aproximación difiere: SurfSense es un **framework** para construir sobre él (requiere editar código Python/TypeScript para cambios profundos); InsightsLM se presenta casi como un **producto listo** que uno puede usar y comercializar directamente ⁸⁰, personalizando principalmente el contenido y flujo a nivel superficial.

En síntesis, **fortalezas de SurfSense** frente a InsightsLM: - Integración con más fuentes de datos (multi-conector) - Pipeline de búsqueda/QA más sofisticado (híbrido con RRF) - Mayor control sobre modelos (soporta cualquier LLM fácilmente, mientras InsightsLM viene prediseñado con ciertos componentes) - Comunidad más orientada a desarrolladores experimentando con nuevas funciones (lo cual acelera innovación).

Fortalezas de InsightsLM: - Despliegue más sencillo (Supabase y n8n, con tutoriales para no-coders ⁸¹ ⁸²) - Posiblemente interfaz más enfocada a uso empresarial out-of-the-box (presentación pulida, integrable sin programar con flujos corporativos) - Tiene detrás una comunidad (The AI Automators) que promueve su uso comercial, lo que puede darle estabilidad en entornos de negocio.

En cuanto a **debilidades**, SurfSense puede ser percibido como más complejo de instalar (muchos requisitos previos: Postgres+pgvector, credenciales de varias APIs, compilación manual) mientras InsightsLM aprovecha servicios gestionados (pero ello también puede ser una debilidad: depende de Supabase y de la cuenta del usuario allí, etc.). SurfSense aún no es “production-ready” según admiten ¹⁵, le falta pulir algunas aristas y documentación completa; InsightsLM, aunque joven, intenta dar la impresión de un producto más terminado.

SurfSense vs. Local-NotebookLM: El proyecto *Local-NotebookLM* (por Goekdeniz Guelmez) es otra alternativa abierta, con un enfoque un tanto distinto: está orientado a convertir **documentos PDF en audio estilo podcast usando solo recursos locales** ⁸³. Es decir, es más una herramienta de “lectura inteligente” de documentos que un asistente general de investigación. De hecho, su funcionalidad central es: el usuario provee PDFs, el sistema genera resúmenes y audio (usando LLMs locales y modelos TTS locales) ⁸³.

Comparado con SurfSense, Local-NotebookLM **tiene alcance más limitado**: no ofrece chat interactivo con preguntas arbitrarias, sino un flujo automático de ingesta -> resumen -> audio. También carece de integraciones con fuentes dinámicas (solo documentos, no Slack/Notion). Sin embargo, su fortaleza es ser **100% offline**: utiliza LLM local (Llama 3.1 mencionan, o similar) y posiblemente TTS open-source (MeloTTS se cita en un artículo ⁸⁴), evitando depender de cualquier API externa. SurfSense, por su parte, puede operar offline con ciertas configuraciones pero típicamente usa servicios externos para embeddings o TTS a menos que se configuren alternativas.

Otra diferencia es la **UI/UX**: Local-NotebookLM al parecer incluso ofrece una app móvil (hay menciones a Flutter app en sus topics ⁸⁵), con un objetivo de ser sencillo: arrastrar un PDF y obtener un audio. SurfSense en cambio ofrece una **suite más compleja**: gestión de espacios, chat conversacional, etc., que conlleva más elementos de interfaz. En rendimiento, SurfSense generando 3 min de audio en 20s es muy rápido ⁵⁵, probablemente porque usa servicios cloud potentes; Local-NotebookLM al usar todo local quizá sea más lento o limitado por hardware.

Fortalezas relativas: SurfSense brilla como *solución integral* para investigación interactiva (multi-fuente, con feedback en tiempo real), mientras Local-NotebookLM es una *herramienta especializada* en procesamiento de documentos a audio, útil para escuchar tus PDFs de forma inteligente. Si el objetivo de José es tener una plataforma base extensible, SurfSense claramente ofrece más en términos de personalización y crecimiento. Local-NotebookLM sería más apropiado si se buscara una solución plug-and-play para un caso de uso concreto (por ejemplo, un académico que quiera escuchar artículos largos resumidos en audio). Vale mencionar que SurfSense también tiene esa capacidad de podcasts y más, por lo que en funcionalidades superpone a Local-NotebookLM en casi todo, excepto en la independencia de recursos externos.

Dependencias críticas y riesgos: SurfSense, al ser tan completo, depende de muchos componentes. Esto conlleva **riesgos** asociados a cada dependencia: - **Dependencias de API externas:** Para funciones plenas, requiere claves de OpenAI (si se usa GPT), de Unstructured.io o LlamaIndex (para parseo de archivos ⁸⁶ ⁸⁷), de Slack, Notion, YouTube API, etc. Si cualquiera de estos servicios cambia sus políticas, introduce costos o se cae, partes de SurfSense pueden verse afectadas. Mitigación: es open-source, se podría cambiar de proveedor (ej. usar un parser local en vez de Unstructured API), pero requiere intervención. - **Frameworks de IA:** Usa LangChain y LangGraph, que son proyectos en evolución. Un cambio de versión grande podría requerir adaptar el código. LangGraph en particular es relativamente nuevo; si dejara de mantenerse o tuviera bugs, SurfSense lo sufriría (aunque siempre podría migrar a puro LangChain si fuera necesario, con esfuerzo). - **LiteLLM:** Este puente para LLMs es muy útil (soporta OpenAI, Azure, Ollama, HuggingFace, etc. en formato unificado ⁸⁸). Sin embargo, confiar en él es un riesgo: si LiteLLM tuviera un bug en producción, afectaría la interacción con modelos. Al ser open-source (BerriAI/litellm), se puede parchear, pero es otra capa más a considerar. - **Chonkie (chunking library):** Depende de Chonkie para fragmentar y embeddear ⁴⁸. Chonkie es relativamente reciente; un riesgo es que presente fallos con ciertos formatos. No obstante, su ventaja es que es local (no es un servicio externo, sino librería). - **Complejidad de stack:** SurfSense reúne Python, Node/TypeScript, Postgres, vector DB, etc. Esto significa más puntos de posible fallo e incompatibilidad. Por ejemplo, instalar pgvector en Windows requiere compilación especializada ⁸⁹ ⁹⁰ - usuarios sin experiencia podrían trabarse allí (el equipo provee guías, pero es un punto de fricción). También la coordinación de versiones (Python 3.x, Node 20+, etc.) debe manejarse. - **Bus factor:** Hasta ahora, mucho recae en el desarrollador principal. Si él reduce su involucramiento sin que la comunidad tome la posta, el proyecto podría ralentizarse. La suerte es que la comunidad es entusiasta y ya hay forks; además, la popularidad alcanzada sugiere que hay incentivos para que otros contribuyan si fuera necesario.

En cuanto a **riesgos de competencia**: NotebookLM en sí (de Google) todavía es beta cerrada, pero si Google decidiera abrirlo o integrarlo masivamente en Google Workspace, podría opacar a estas alternativas open-source. No obstante, SurfSense apuesta por nichos que los grandes pueden ignorar (soporte offline, integración con cualquier fuente/dato personalizado). Proyectos similares (InsightsLM, Local-NotebookLM, otras alternativas emergentes) también compiten en la atención de la comunidad de self-hosted AI. Si el espacio se fragmenta en muchas alternativas, existe riesgo de dispersión de contribuidores. Sin embargo, SurfSense actualmente lleva ventaja en cuanto a features y comunidad, así que mientras mantenga el paso, es probable que consolide su posición.

Resumiendo, estratégicamente SurfSense posee un fuerte **valor diferencial** (flexibilidad y potencia) pero debe navegar sus desafíos de usabilidad y depender de tantas piezas. Manejar bien las contribuciones comunitarias, mejorar documentación y quizás simplificar despliegue (ej. empaquetar todo en un instalador o imagen más simple) serán claves para mitigar riesgos y aprovechar oportunidades de convertirse en *la* plataforma estándar de asistentes de investigación IA open-source.

6. Vínculo con la Investigación Académica

SurfSense no surge en el vacío; está alineado con varias líneas de investigación académico-industrial sobre la intersección de **LLMs y grafos de conocimiento**:

- **GraphGPT**: Este término se ha usado para referirse a sistemas que convierten texto en grafos de conocimiento utilizando LLMs ⁹¹ o que integran instrucciones sobre grafos en el entrenamiento de LLMs ⁹². La idea central es aprovechar la capacidad de un GPT para **extraer entidades y relaciones** de texto y construir un grafo estructurado. SurfSense inicialmente se promocionó como un “Knowledge Graph Brain” para tus sesiones web ⁵, y de hecho en sus primeras versiones almacenaba la información en una base de datos de grafos Neo4j ⁹³. Implementaba algo denominado *GraphRAG*, es decir, Recuperación Aumentada con Grafos, para encontrar relaciones significativas entre el contenido guardado ⁶. Esto está muy en línea con enfoques tipo GraphGPT: estructurar el conocimiento en forma de grafo (nodos representando conceptos o páginas, aristas representando relaciones contextuales) para mejorar la **memoria y razonamiento** del sistema. Aunque luego SurfSense pivotó a pgvector, el concepto de GraphRAG demuestra compatibilidad con grafos: José podría reactivar o extender esa funcionalidad para que el backend mantenga un grafo de conceptos que complementen la búsqueda vectorial.
- **TEA-GLM (Token Embedding Aligned Graph LM)**: Es un marco de investigación que alinea representaciones de Graph Neural Networks (GNNs) con embeddings de LLM, permitiendo que las LLM resuelvan tareas de grafo en zero-shot ⁹⁴ ⁹⁵. En esencia, TEA-GLM entrena una proyección de los embeddings de nodos de un grafo a un espacio común con tokens de lenguaje, habilitando a un LLM pre-entrenado a “entender” grafos. ¿Cómo se conecta esto con SurfSense? Imaginemos que José tiene un grafo de conocimiento derivado de sus mapas mentales; TEA-GLM sugiere que se puede integrar ese grafo con un LLM sin re-entrenarlo desde cero, usando alineación de embeddings. SurfSense podría ser el entorno donde probar algo así: los vectores de pgvector podrían enriquecerse con información estructural del grafo (via un modelo GNN) y luego el LLM personalizado (Qwen3) ajustado para entender esos vectores alineados ⁹⁴. Si bien es avanzado, en teoría SurfSense podría incorporar una etapa donde las representaciones vectoriales no solo vengan de texto sin procesar sino de un **GNN que representa el conocimiento simbólico** de José, aplicando ideas de TEA-GLM para mejorar la pertinencia de la respuesta. Esto acercaría a SurfSense a un terreno de **IA neuro-simbólica**, combinando embeddings con interpretabilidad gráfica.
- **NOCL (Node-Oriented Conceptualization LLM)**: Esta es otra propuesta investigativa donde se reformula las tareas de grafo como problemas de comprensión textual para un LLM ⁹⁶. En lugar de usar Message Passing típico de GNNs, NOCL hace que el LLM razone sobre grafos describiéndolos en lenguaje natural. Por ejemplo, convertir la estructura de un grafo en un enunciado textual que el LLM pueda completar o analizar. SurfSense podría aprovechar una estrategia parecida: dado un grafo de conocimiento (por ejemplo, generado a partir de la información del usuario), se podría **generar “prompts enriquecidos”** que describan conexiones clave del grafo al LLM antes de la respuesta. Como SurfSense ya construye prompts con

contexto, incluir texto adicional derivado de relaciones en un grafo sería factible. Esto se alinea con NOCL en cuanto a no requerir un algoritmo de paso de mensajes, sino usar el propio LLM para inferir sobre las relaciones cuando se formulan adecuadamente ⁹⁷.

- **Graph Neural Networks (GNNs) y LLMs:** Más allá de casos específicos, existe una tendencia general de combinar LLMs con conocimiento simbólico/estructurado. GNNs excel en tareas de razonamiento sobre grafos (como inferir nuevas conexiones, clasificar nodos, etc.), mientras los LLMs son buenos en conocimiento contextual y lenguaje. SurfSense puede ser un “playground” ideal para probar integración de GNNs: por ejemplo, se podría conectar un **motor de inferencia basado en grafos** (como una base de conocimiento ontológica o un GNN entrenado con los datos del usuario) para que provea resultados que luego el LLM de SurfSense verbalice. Una posibilidad concreta: si José tiene sus mapas mentales como grafos conceptuales, podría aplicar un GNN que recomiende los conceptos más relevantes a su pregunta; esos conceptos podrían añadirse como contextos adicionales en SurfSense. Esto enriquecería la respuesta con conocimiento derivado del grafo (quizá mejorando la **cobertura conceptual** y reduciendo lagunas que un embedding puro no capta).

En la práctica, ¿qué tan realizable es conectar un backend de grafos a SurfSense? Dado que ya hubo integración con Neo4j, es definitivamente posible. Podría hacerse que, además del vector store, SurfSense consultara un **Knowledge Graph**: por ejemplo, una consulta Cypher para buscar en el grafo conceptos relacionados a la pregunta, trayendo un subgrafo de contexto. Ese subgrafo podría transformarse en texto ("En tu grafo, X está relacionado con Y y Z") que el LLM reciba. O incluso, mediante LangChain, registrar el grafo como una *Tool* que el agente pueda usar (existe precedentes de agentes que consultan bases de conocimiento externas en su prompt). SurfSense, con LangGraph, puede incorporar nuevos pasos, por lo que se puede dotar al agente de la habilidad "ConsultarGrafo" que internamente llame al servicio de grafo de José. Esto permitiría un **razonamiento híbrido**: primero el agente intenta responder con sus documentos, y si necesita, invoca el grafo para deducciones adicionales, luego compone la respuesta final. Tal arquitectura combinaría lo mejor de ambos mundos: la solidez factual de la búsqueda vectorial con la **consistencia lógica** de un grafo curado.

En cuanto a **estado del arte académico**, SurfSense se encuentra en sintonía con esa ola de *augmented LLMs*. No es un experimento de laboratorio, sino una aplicación concreta, pero su estructura modular la hace apta para incorporar soluciones de papers recientes. Por ejemplo, si se desarrolla un algoritmo de explainability usando grafos (hubo investigaciones de utilizar LLMs para generar explicaciones estructuradas de respuestas, etc.), SurfSense podría incorporarlo como módulo de generación de explicaciones en grafo a posteriori.

Posibilidades reales: Concretamente para José, que dispone de mapas mentales en forma de grafo, la integración es muy factible. Podría mantener su grafo en una base Neo4j o similar y conectar SurfSense para: - Ingerir datos del grafo como textos (para al menos tenerlos en la base vectorial). Esta es la forma más simple, aunque pierde la estructura relacional. - Consultar el grafo durante las preguntas: implementando una herramienta en el agente que dado un tema devuelva los nodos relacionados de su mindmap (aprovechando relaciones *is-a*, *related-to* en ese grafo). - Incluso incorporar modelos tipo **GraphGPT** para que SurfSense pueda *expandir* el conocimiento en el grafo a partir de nuevo texto que se le alimente (aunque esto sería un paso más orientado a mantenimiento de la base de conocimiento que a la respuesta en sí).

En resumen, SurfSense proporciona la infraestructura para experimentar con **IA simbólica + conexionista** aplicada al usuario final. Sus creadores ya coquetearon con ello (GraphRAG, Neo4j) y la arquitectura actual, aunque centrada en vectores, se puede extender. Esto la hace ideal como base para

la investigación aplicada: se puede validar cómo técnicas como TEA-GLM o NOCL mejorarían las respuestas de SurfSense incorporando un grafo – por ejemplo, midiendo si la citación de hechos mejora cuando el grafo verifica consistencia, etc. Muchos proyectos académicos quedan en simulaciones; SurfSense permite implementarlos en un sistema completo y usar casos de la vida real.

En conclusión, la conexión de SurfSense con la investigación es **bidireccional**: por un lado implementa ideas avanzadas de RAG que han sido materia de papers recientes, y por otro lado ofrece un banco de pruebas para integrar de manera práctica propuestas académicas de integrar LLMs con grafos y otros conocimientos estructurados, acercando esas innovaciones al usuario final.

7. Informe SWOT

A continuación, se presenta un análisis **SWOT** (Fortalezas, Debilidades, Oportunidades, Amenazas) del proyecto SurfSense:

Fortalezas:

- **Altamente Personalizable y Privado:** Es software libre, auto-hospedable y soporta despliegue local, lo que garantiza control de datos y configuración a gusto del usuario ¹⁰. Permite usar modelos locales (Ollama) sin necesidad de servicios en la nube, preservando privacidad total.
- **Integración con Múltiples Fuentes:** SurfSense se conecta con una variedad inusual de fuentes externas – buscadores web, Slack, Notion, GitHub, YouTube, etc. ⁸ ⁹ – convirtiéndose en un hub unificado de conocimiento personal/profesional. Esto supera a muchas soluciones cerradas enfocadas en una única fuente.
- **Capacidades de Búsqueda Avanzadas:** Implementa técnicas de **RAG de punta**: búsqueda híbrida (semántica + keyword) con fusión de ranking recíproco ³⁵, índices jerárquicos de 2 niveles ³⁸, y rerankers especializados ³⁹ para optimizar la relevancia. Estas funciones suelen ser exclusivas de motores de búsqueda empresariales, dándole a SurfSense un rendimiento elevado en recuperación de información.
- **Respuestas con Citas Fuentes:** Siguiendo el modelo de Perplexity, cada respuesta viene respaldada por referencias a las fuentes originales ⁷. Esto da confianza al usuario y reduce la posibilidad de alucinaciones no detectadas.
- **Soporte masivo de Modelos (LLM y embeddings):** Gracias a LiteLLM y a la integración de *AutoEmbeddings*, puede trabajar con más de 150 modelos de lenguaje distintos y 6000+ modelos de embeddings pre-entrenados ⁷⁸ ³⁹. Es difícil encontrar otra plataforma con tal amplitud de compatibilidad, lo que lo hace futuro-proof ante nuevos LLMs (por ejemplo, es relativamente sencillo configurarlo para usar un GPT-4 local o un modelo experimental).
- **Funcionalidades Únicas (Podcast, Extensión navegador):** Ofrece características innovadoras como la generación ultrarrápida de podcasts a partir de chats ⁵⁵, convirtiendo texto en audio con mínima espera, y una extensión de navegador para capturar cualquier página web *tal cual* la ve el usuario ⁹⁸. Esto último elimina barreras de contenido dinámico o protegido (login) que otras herramientas no pueden salvar.
- **Comunidad Activa y Desarrollo Ágil:** La rapidez con que se agregaron características (ej. soporte Docker en días ⁹⁹, migración a LangGraph agente personalizado en semanas ¹⁰⁰) y la comunicación abierta en Discord/GitHub indican un proyecto vivo, con un lead developer muy dedicado y usuarios comprometidos dando feedback.

Debilidades:

- **Complejidad de Instalación y Configuración:** SurfSense requiere varios pasos previos (instalar Postgres + pgvector, obtener múltiples API keys, configurar varios archivos `.env`) ¹⁰¹ ¹⁰². Para

usuarios no técnicos, esto puede ser abrumador comparado con soluciones plug-and-play. Incluso con Docker (que simplifica parte), ciertas funcionalidades no operan (Ollama, web crawler) ⁵⁹.

- **Estado Beta/No completamente pulido:** Los propios desarrolladores indican que **no está aún listo para producción** y que varias características están en desarrollo activo ¹⁵. La función de Podcasts, por ejemplo, está temporalmente deshabilitada para mejora ⁵⁸. Pueden surgir errores o ausencias en documentación dado lo rápido del avance.
- **Elevados Requisitos de Sistema:** Para sacarle provecho, especialmente con LLMs grandes, se necesita hardware potente (RAM, quizá GPU si se usan modelos grandes en local). El stack completo (BD, backend, frontend) también consume recursos significativos. No es tan ligero como, digamos, una sola aplicación de escritorio; es una mini plataforma web.
- **Dependencia de Servicios Externos para Plenas Funcionalidades:** Aunque puede funcionar offline en modo básico, muchas características potentes usan servicios de terceros: parseo de documentos via Unstructured.io o LlamaIndex cloud ¹⁰³, TTS y STT via OpenAI/Azure ⁵⁴, búsqueda web via APIs de motores (Tavily, LinkUp). Esto introduce potenciales costos y puntos de falla externos. Sin conexión a Internet, SurfSense pierde algunas capacidades (aunque la consulta a los docs locales seguiría operativa).
- **Curva de Aprendizaje:** La interfaz, al ser completa, puede resultar algo **compleja para nuevos usuarios**. Conceptos como “Spaces”, “Late chunker”, elegir embeddings, etc., pueden abrumar. No es la experiencia simplificada de un solo chat estilo ChatGPT; es más cercana a una herramienta profesional con paneles y opciones.
- **Mantenimiento principalmente unipersonal:** Si bien hay contribuciones, gran parte recae en una persona (ModSetter). Esto puede ser un riesgo de cuello de botella en la toma de decisiones o en soporte a largo plazo. Si el autor se desvía a otros proyectos, la comunidad tendría que asumir la carga sin una transición clara.
- **Internacionalización temprana:** Hasta donde se ve, la UI y documentación están principalmente en inglés. Usuarios hispanohablantes u otros pueden encontrar barrera idiomática. (Esto podría solventarse con contribuciones de traducción en el futuro).

Oportunidades:

- **Demanda Creciente de Alternativas NotebookLM/Perplexity Auto-gestionadas:** La tendencia de llevar LLMs a entornos privados va en aumento (como refleja la popularidad en subreddits de Local LLM). SurfSense está bien posicionado para convertirse en **la opción de referencia** en este nicho, dado que proyectos comparables son pocos y con menos funciones. Esto podría atraer colaboración de empresas o instituciones interesadas en una base sólida para asistentes de investigación.
- **Integración con Investigación Académica y Nuevas Tecnologías:** Como se analizó, SurfSense podría incorporar innovaciones de punta (p. ej. integración profunda con knowledge graphs, nuevas técnicas de memoria a largo plazo, agentes multi-modal). Esto no solo mejoraría el proyecto sino que podría convertirlo en una **plataforma de experimentación** adoptada en ámbitos académicos. Por ejemplo, grupos de investigación podrían usar SurfSense para prototipos, aportando mejoras a la base del código.
- **Colaboraciones y Extensiones Comerciales:** Existe la posibilidad de que terceros construyan sobre SurfSense soluciones especializadas – por ejemplo, una versión orientada a abogados (con conectores a bases jurídicas), o para medicina (integrando repositorios científicos). Al ser Apache 2.0, incluso puede ser usada en productos comerciales sin problema ¹⁰⁴. Esto podría generar un ecosistema de plugins o extensiones de SurfSense para distintos verticales.
- **Mejoras en Usabilidad (Low-Code):** Aprendiendo de InsightsLM, el proyecto podría desarrollar *wrappers* o instaladores que simplifiquen su adopción (por ej., un script de instalación todo-en-uno, plantillas n8n para flujos comunes, etc.). Esto bajaría la barrera de entrada y ampliaría la

base de usuarios más allá de desarrolladores. También la comunidad puede traducir la interfaz, mejorar la estética, etc., haciéndola más atractiva.

- **Expansión de la Comunidad:** Con la tracción que tiene (miles de estrellas, comentarios en Reddit), hay oportunidad de crecer la comunidad mucho más. Un foro oficial, sesiones de demostración en streaming, o hackathons dedicados podrían atraer talento que contribuya. La diversidad de interesados (desde entusiastas de IA hasta profesionales de datos) enriquece el proyecto con ideas y casos de uso.
- **Posicionarse como “RAG-as-a-Service”:** En el README mencionan ofrecer un backend de RAG como servicio ¹⁰⁵. Si se pule, SurfSense podría usarse como una API que otras apps consuman para hacer Q&A con datos propios. Dado que integra tantas fuentes, podría ser **la API definitiva de búsqueda con IA** para desarrolladores independientes. Esto abre puertas a que se incluya en listas de APIs públicas o marketplaces, aumentando su notoriedad.

Amenazas:

- **Competencia de Grandes Plataformas:** Google NotebookLM podría salir del beta y ofrecer su versión gratuita integrable con Google Drive, etc., lo que para muchos usuarios sería suficiente. Perplexity y otros asistentes cerrados siguen mejorando; aunque no ofrecen la personalización de SurfSense, podrían tratar de añadir conectividad a cuentas personales en el futuro (ej: Perplexity podría un día conectarse a tu Notion oficialmente). Si bien SurfSense tendría aún la ventaja open-source, la comodidad de las soluciones “oficiales” puede restarle usuarios potenciales.
- **Fragmentación del Espacio Open-Source:** Con varias alternativas surgiendo (InsightsLM, OpenNotebookLM, quizás proyectos como Haystack o LlamaIndex evolving con UIs propias), la comunidad de desarrolladores podría dispersarse. Si no se consolida un estándar de facto, se duplicarán esfuerzos. SurfSense necesita mantener un ritmo superior de innovación para continuar liderando; de lo contrario, otro proyecto más usable podría tomar la delantera.
- **Cambios en Políticas de APIs externas:** Un riesgo concreto es que servicios como Slack o Notion limiten el uso de sus APIs para propósitos de bots de este tipo, o que los motores de búsqueda cierren sus endpoints no oficiales (Tavily/LinkUp podrían romperse si Google cambia algo). Cualquier conector podría dejar de funcionar temporalmente por causas ajenas, lo que impacta la propuesta de valor (aunque los datos internos seguirían utilizables). La dependencia en OpenAI (si bien evitable con locales) es también un riesgo: cambios de pricing o rate limits más duros afectarían a usuarios que lo usen con GPT-4 vía API.
- **Escalabilidad y Performance:** Si SurfSense se adopta en entornos empresariales con *datasets* muy grandes, la elección de Postgres/pgvector podría quedarse corta en rendimiento frente a soluciones especializadas. Si usuarios con decenas de millones de documentos intentan usarlo, podría no escalar linealmente. Esto podría generar críticas al proyecto por no manejar big data, aunque su objetivo original sea knowledge bases personales (más pequeñas). Habría que monitorear la percepción para que no se lo utilice fuera de su scope sin los debidos ajustes (por ejemplo, usar vector DB externos si hiciera falta).
- **Seguridad y Privacidad:** Al ser un agregador de datos personales, un bug de seguridad podría exponer información sensible. Si se encontraran vulnerabilidades (en la autenticación JWT, o inyección a través de archivos maliciosos), la reputación del proyecto podría dañarse. Mantener buenas prácticas de seguridad y auditorías será clave, pero siempre existe ese riesgo en software complejo. Asimismo, cualquier brecha en las dependencias (e.g. una vulnerabilidad en FastAPI Users o en la extensión de navegador) sería una amenaza que requeriría respuesta rápida.
- **Riesgo de Estancamiento:** Dado que el campo de LLMs avanza veloz, SurfSense debe adaptarse o podría volverse menos relevante. Por ejemplo, si emergen LLMs especializados con *retrieval* interno (como modelos que consultan vectores por sí mismos o integrados en bases vectoriales),

la arquitectura actual podría requerir una revisión. O si los estándares cambian (imaginemos un nuevo formato de embedding dominante), SurfSense tendría que actualizar su soporte. La amenaza es quedarse atrás tecnológicamente; no parece inminente dado el ritmo actual, pero es algo a considerar en la planificación a futuro.

8. Hoja de Ruta para la Integración (Caso José)

A continuación, se detalla un plan paso a paso para que José integre SurfSense con su sistema de fine-tuning incremental basado en mapas mentales (grafos) y su modelo de lenguaje personalizado (por ejemplo Qwen 3 afinado). Este plan asume un nivel técnico alto por parte de José, dado su perfil de investigador en IA simbólica.

Paso 1: Clonar el repositorio y preparar el entorno.

Comenzar obteniendo el código fuente de SurfSense. Esto se hace con un simple `git clone` del repo oficial (GitHub: `MODSetter/SurfSense`). Una vez clonado, José tendrá tres componentes principales: `surfsense_backend` (FastAPI), `surfsense_web` (Next.js frontend) y `surfsense_browser_extension` (extensión). Antes de instalar nada, revisar las dependencias: necesitará Python (3.10+), Node.js (18 o 20) y PostgreSQL. También asegurarse de tener acceso a las API keys de los servicios que vaya a usar (si usa Unstructured para PDF, su clave; si usará OpenAI para TTS, su API key, etc.). En el caso de modelo local puro, muchas API keys de LLM no serán necesarias, pero otras funcionales (como Google OAuth si se desea login con Google) sí podría configurarlas.

Paso 2: Configurar PostgreSQL con pgvector.

Instalar PostgreSQL 15+ y habilitar la extensión pgvector. En Linux/Mac, esto implica compilar el módulo pgvector tal como indican los docs ¹⁰⁶. En Windows, los pasos son un poco más arduos (requiere Visual Studio C++ tools, luego `nmake`) ¹⁰⁷ ⁹⁰. Como alternativa, José podría usar Docker solo para la base de datos (ejecutar un contenedor Postgres con pgvector preinstalado) mientras corre el resto nativo; pero dado que es cómodo compilando código, puede hacerlo localmente. Crear una base de datos `surfsense` y un usuario (o usar el default `postgres`), recordando cargar la extensión con `CREATE EXTENSION vector;`. Verificar que puede conectar a la DB.

Paso 3: Instalación del Backend (FastAPI).

Navegar al directorio `surfsense_backend`. Copiar el archivo de entorno de ejemplo: `cp .env.example .env` ¹⁰⁸. Abrir `.env` en un editor y configurar los valores esenciales: - `DATABASE_URL` con la cadena de conexión a Postgres (ej: `postgresql+asyncpg://user:pass@localhost:5432/surfsense`) ¹⁰⁹. Usar `asyncpg` como driver async. - `SECRET_KEY` con una cadena aleatoria robusta (para JWT) ¹⁰⁹. - `AUTH_TYPE`: decidir si usar Google OAuth o autenticación local. Para simplificar pruebas, puede usar `LOCAL` (usuario/contraseña). - Si `AUTH_TYPE=GOOGLE`, proveer `GOOGLE_OAUTH_CLIENT_ID` y secret con credenciales de la API de Google ¹¹⁰, pero esto es opcional. - `EMBEDDING_MODEL`: aquí José debe elegir un modelo de embeddings. Por facilidad, puede usar uno como `"sentence-transformers/all-MiniLM-L6-v2"` o el sugerido `"mixedbread-ai/mxbai-embed-large-v1"` ⁴⁰. Esta elección afecta cómo se indexará el texto. - `RERANKERS_MODEL_NAME` y `RERANKERS_MODEL_TYPE`: puede inicialmente dejar los valores por defecto del `.env.example` (ej: MiniLM cross-encoder y flashrank) ¹¹¹, o deshabilitar rerankers si prefiere (posiblemente poniendo un modelo vacío desactivaría esa fase). - **Configurar LLMs via LiteLLM:** Dado que José quiere usar su modelo Qwen3, asumiremos que lo tiene ejecutándose en un servicio local. Si ha optado por Ollama, por ejemplo, y lo cargó ahí con nombre "qwen3", entonces en `.env` pondrá: - `STRATEGIC_LLM = "ollama/qwen3"` (como modelo principal para respuestas) ⁴³. - `FAST_LLM` podría ser un modelo más pequeño para pre-respuestas (podría poner un modelo open-source ligero como `ollama/llama2:7b` o incluso usar el mismo Qwen3 si no

quiere complicarse). - `LONG_CONTEXT_LLM` solo si planea usar uno distinto para documentos muy largos; se puede omitir inicialmente. - Importante: si `FAST_LLM` o otros apuntan a OpenAI (ej. `"openai/gpt-4o-mini"`), asegurarse de definir `OPENAI_API_KEY`. En este caso, supongamos que José quiere todo local: pondrá modelos de Ollama en FAST y STRATEGIC. Entonces **no necesita OpenAI key**, pero LiteLLM sí requiere saber cómo conectar a Ollama. Por defecto, LiteLLM usa `http://localhost:11434` para Ollama, lo cual funciona sin más si el demonio está en la misma máquina. - Configurar TTS/STT si usará podcasts: por ejemplo, `TTS_SERVICE="google/standard"` y su `TTS_SERVICE_API_KEY`, o `TTS_SERVICE="elevenlabs"` si tuviera una clave allí. Si no va a usar la función Podcast de inmediato, puede dejar estos campos vacíos o con valores dummy. - `ETL_SERVICE`: elegir `UNSTRUCTURED` o `LLAMACLOUD` según cómo quiera parsear documentos ¹¹². Si no planea cargar PDFs complejos, podría incluso más adelante implementar su propio parser. Inicialmente, poner `UNSTRUCTURED` y configurar `UNSTRUCTURED_API_KEY` para que la carga de archivos PDF/Word funcione ¹¹² ¹¹³.

Guardados los cambios, instalar las dependencias Python. Los docs recomiendan usar `uv` (una suerte de gestor que sincroniza deps) ¹¹⁴, pero José puede simplemente crear un virtualenv e instalar con pip: `pip install -r requirements.txt`. Esto descargará FastAPI, LangChain, LangGraph, SQLAlchemy, etc., asegurando el entorno.

Finalmente, iniciar el backend: ejecutar `uvicorn server:app --host 0.0.0.0 --port 8000` desde `surfsense_backend` ¹¹⁵. (Asume que el archivo principal se llama `server.py` con la app FastAPI; si es distinto, verificar en la doc o README). Con `--reload` para desarrollo. Comprobar en la consola que se lanzó sin errores y está conectado a la DB (debería mostrar migraciones aplicándose y luego "Running on http://0.0.0.0:8000"). También navegar a `http://localhost:8000/docs` para ver la documentación Swagger del API y probar la ruta de `/register` por ejemplo ²⁶.

Paso 4: Instalación del Frontend (Next.js).

Entrar al directorio `surfsense_web`. Copiar igualmente `.env.example` a `.env` ¹¹⁶. Editar `.env` del front: principalmente debe contener `NEXT_PUBLIC_FASTAPI_BACKEND_URL = "http://localhost:8000"` ¹¹⁷ (o la URL donde corre el backend). Esto permite que la web sepa a dónde hacer las peticiones. Ejecutar `pnpm install` (o `npm install` si pnpm no está instalado, aunque recomiendan pnpm) ¹¹⁸. Esto descargará todas las dependencias Node (React, Next, Tailwind, etc.). Luego correr `pnpm run dev` ¹¹⁸. Si todo va bien, la aplicación web estará accesible en `http://localhost:3000`. Abrir ese enlace en un navegador. Debería aparecer la pantalla de inicio de SurfSense (posiblemente pidiendo login). Si José dejó `AUTH_TYPE=LOCAL`, habrá una opción para registrarse con email/password. Registrar un usuario (o hacerlo vía Swagger como indicaban, pero vía interfaz es más sencillo si funciona). Tras login, la app probablemente pedirá ingresar algunas credenciales de servicios (por ejemplo, la clave de OpenAI si no detecta configurada, etc.) – en los screenshots se ve una sección de Settings donde poner credenciales de Neo4j y OpenAI en la versión antigua ⁹³ ¹¹⁹, pero en la nueva puede que solicite la de Unstructured, etc. José debería introducir las keys necesarias (o falsas si no usará ese servicio inmediato).

Con esto, la plataforma ya estaría en funcionamiento localmente con su modelo. Es buen momento para probar **una carga de documento simple** y una **pregunta** para validar: Por ejemplo, subir un archivo de texto o Markdown (la interfaz tiene sección de "Manage Documents" ¹²⁰) – internamente, el backend llamará al parser (unstructured) y guardará el contenido. Luego en la sección de chat, realizar una pregunta sobre ese documento para verificar que el modelo Qwen3 responde adecuadamente con citas.

Paso 5: Integrar el sistema de grafos (mapas mentales).

Dado que José cuenta con mapas mentales convertidos a grafos de conocimiento, hay dos posibles niveles de integración: - *Ingesta estática*: Si tiene esos grafos en un formato textual (ej. exportados a un JSON/CSV con relaciones, o a un texto que lista conceptos y sus vínculos), podría importarlos a SurfSense como documentos. Por ejemplo, podría convertir cada relación del grafo en una frase del estilo “A está relacionado con B en tal contexto” y cargar eso como un documento (o varios). Así, al menos la información del grafo estaría indexada en los embeddings, y el LLM podría usarla en las respuestas. Esta vía es rápida (transformar los datos y usar la función de upload). - *Integración dinámica*: Para explotar plenamente su motor de razonamiento simbólico, José puede conectar su backend de grafo como describimos en la sección 3 (punto de inyección). Una forma inicial: crear un pequeño servicio API REST en su sistema de grafos que, dado un input (pregunta o tema), devuelva resultados (por ej., “nodos vecinos relevantes” o inferencias del grafo). Luego, en SurfSense backend, implementar una **Tool de LangChain** que llame a ese servicio. Esto requiere modificar el agente de LangGraph para incluir la herramienta; por ejemplo, definir una función Python que consuma la API del grafo y la registre en LangChain. Alternativamente, usar la estrategia LiteLLM fake-model: entrenar un mini-modelo o script que combine la respuesta del grafo y LLM. Esto es más complejo, así que recomendamos comenzar con la herramienta separada.

Un plan concreto: Añadir en el backend un endpoint `/graph_query` que reciba una query y responda con datos del grafo. Luego, editar la chain del agente (posiblemente en un archivo de configuración YAML/JSON si LangGraph lo usa así, o en código Python donde se construye la chain) para que antes de la llamada final, haga una llamada a `/graph_query` y agregue esa info al contexto. En paralelo, José puede experimentar con uno de sus modelos (o el mismo Qwen3) para formatear resultados del grafo en texto legible. Por ejemplo, si el grafo devuelve “NodoX relacionado con NodoY y NodoZ”, quizá el LLM podría en la respuesta citar: “Según tu mapa mental, X se conecta con Y y Z, lo que sugiere...”. Este tipo de inserción elevaría la respuesta a un nivel simbólico.

Paso 6: Pruebas y validación incremental.

Con la integración del grafo en marcha, realizar pruebas en casos concretos: preguntas donde espera que el grafo aporte valor (p.ej., preguntas de razonamiento que requieren múltiples piezas). Observar las respuestas. Ajustar en consecuencia: - Si las respuestas del LLM ignoran la info del grafo, tal vez mejorar el prompt que inserta esos datos (hacerlo más directo: “*Datos de tu Knowledge Graph:*”). - Si el modelo Qwen3 genera textos incoherentes o alucina, quizás necesita más afinación o usar un modelo más robusto para respuestas finales (por ej., Qwen3 para rapidez pero repreguntar a GPT-4 en segundo plano si se requiere validación). - Monitorizar el desempeño: SurfSense es compatible con **LangSmith** para trazas ¹²¹, José podría habilitarlo para depurar dónde la chain podría mejorar.

Paso 7: Optimizar la experiencia de usuario para el flujo integrado.

Puede que José necesite exponer alguna opción en la UI para su sistema. Por ejemplo, un botón “Consultar Grafo” por fuera del flujo normal, o una forma de visualizar las relaciones encontradas. Dado que el frontend es modular, él podría crear un componente React nuevo que presente, por decir, un subgrafo cuando el usuario selecciona un resultado. Aunque esto es opcional, añadir visualizaciones de los mapas mentales (ya que los tiene) sería un *bonus* enorme – quizás integrando una librería de graph drawing en el frontend para mostrar cómo la respuesta se conectó con los nodos del grafo (sería un nivel de explainability extra).

Paso 8: Documentar y iterar.

Como último paso, conviene que José documente su configuración personalizada (por su propio bien y para la comunidad). Podría anotar en un README interno cómo levantó el modelo Qwen3 con Ollama, cómo agregó la llamada al grafo, etc. Contribuir esos cambios de vuelta al proyecto principal (si son de

aplicación general) sería ideal, aunque si su sistema grafo es propietario tal vez no lo haga público. En cualquier caso, la documentación le ayudará a mantener su fork.

Anticipación de desafíos técnicos:

- *Compatibilidad de Qwen3*: Asegurarse de que el modelo Qwen3 esté bien soportado en el runtime elegido. Si es Ollama, verificar que Qwen está en formato compatible (p. ej., Qwen-7B tiene versión int4 para Ollama, si Qwen3 es similar tendría que portarlo). Si no, otra opción es usar LocalAI or llama.cpp to run Qwen weights, configurando LiteLLM para huggingface.
- *Ajuste de parámetros de chunking*: Qwen3 podría tener un contexto máximo distinto a GPT-3.5/4. SurfSense usa por defecto *LateChunker* que adapta el tamaño de fragmentos al máximo de tokens del modelo ¹²². Es probable que necesite especificar en la config cuántos tokens maneja Qwen3 para que LateChunker lo use. Revisar `LateChunker` en la config (quizá se autodetecta vía nombre del modelo, no seguro). Si no se hace, podría haber fragmentos demasiado grandes y el modelo truncarlos o tardar.
- *Rendimiento del modelo*: Si Qwen3 es grande, las inferencias pueden ser lentas. Para mitigar, José podría: usar `FAST_LLM` con un modelo más pequeño para dar respuestas provisionales rápidas ⁴³; ajustar parámetros de generación (menos max tokens, más temperatura si quiere respuestas concisas).
- *Integración del grafo - sincronización*: Debe decidir si la información del grafo se actualiza en SurfSense en tiempo real o con alguna periodicidad. Por ejemplo, si el mindmap se actualiza, ¿se reindexa en SurfSense? Podría automatizarlo: cada vez que cambie el grafo, llamar al API de ingesta de SurfSense para añadir la nueva info. O, más elegante, usar la herramienta del agente para siempre consultar el grafo de fuente viva en vez de duplicarlo.
- *Decisiones de arquitectura del agente*: Incorporar una herramienta extra en LangChain puede hacer las conversaciones más largas o complejas. Quizás deba aumentar `max_turns` del agente o ajustar los prompt templates para que no malinterprete la salida del grafo. Estas finezas se ajustarán con prueba y error.
- *Mantenimiento al actualizar SurfSense*: Al haber personalizado, cada vez que quiera actualizar a una nueva versión del proyecto tendrá que mergear sus cambios. Para reducir dolores, intentar hacer los cambios de manera aislada (p. ej., en archivos nuevos o configuraciones externas) para que un pull del upstream no cause muchos conflictos. Por ejemplo, en lugar de modificar fuertemente `server.py`, podría envolver su llamada al grafo en un plugin externo.

Siguiendo esta hoja de ruta, José debería lograr **clonar e instalar SurfSense**, configurarlo con su LLM personalizado Qwen3, e **integrar su potente base de conocimiento en grafo** para enriquecer las respuestas. El resultado será un asistente único que combina la recuperación de información contextual con inferencias simbólicas de alto nivel – justo el híbrido que motivó su investigación. Vale la pena avanzar paso a paso, probando cada componente antes de complicarlo más, y aprovechar la comunidad (Discord, GitHub) en caso de encontrarse con algún escollo inesperado. Con SurfSense como plataforma base, José puede concentrarse en lo novedoso (su motor de grafos y fine-tuning incremental) sin reinventar toda la infraestructura de un asistente conversacional, obteniendo así lo mejor de ambos mundos.

Enlaces relevantes a código/discusiones:

- Repositorio SurfSense (GitHub): conexión con funcionalidades NotebookLM/Perplexity ³ ⁸.
- Documentación oficial (surfsense.net/docs): guía de instalación, variables de entorno y limitaciones Docker ⁵⁹ ¹⁰².
- Post Reddit del autor sobre integración de LiteLLM (LLMs soportados ilimitadamente) ⁴².

- Devpost de SurfSense: detalles de configuración de múltiples LLM y proveedores ⁶⁰ , así como estado de características (podcasts, etc.) ⁵⁷ .
- Discusiones Reddit con la comunidad: YouTube frames vs transcript ²¹ , integración con herramientas externas (Zapier, Lyrz) ⁶⁸ – para entender la visión comunitaria.
- Investigación relacionada: TEA-GLM (alineamiento de GNN y LLM) ⁹⁴ , NOCL (LLM resolviendo tareas de grafo sin MPNN) ⁹⁶ , GraphGPT (construcción de grafos a partir de texto) ⁹¹ – potencial teórico que puede aprovecharse en la integración.

Estas referencias respaldan las afirmaciones de este informe y sirven como punto de partida si se desea profundizar en algún aspecto específico. En conclusión, SurfSense se presenta como una **plataforma robusta y de vanguardia** que, con la estrategia adecuada, puede fusionarse con sistemas de conocimiento simbólico para lograr un asistente de investigación verdaderamente potente y personalizado.

8 44

-
- ¹ ² ¹³ ¹⁴ ¹⁸ ¹⁹ ⁴² ⁹⁹ SurfSense - Personal AI Assistant for World Wide Web Surfers. : r/selfhosted
https://www.reddit.com/r/selfhosted/comments/1f158vh/surfsense_personal_ai_assistant_for_world_wide/
- ³ ⁷ ¹⁰ ¹⁶ ¹⁷ ²³ ²⁴ ²⁷ ²⁸ ³¹ ³² ³³ ³⁴ ³⁷ ⁴¹ ⁴⁴ ⁴⁵ ⁴⁶ ⁴⁷ ⁴⁸ ⁴⁹ ⁵⁰ ⁵³ ⁵⁴ ⁵⁵ ⁶¹ ⁶⁴ ⁶⁶ ¹⁰⁴
¹²² GitHub - Decentralised-AI/SurfSense-Open-Source-Alternative-to-NotebookLM: Open Source Alternative to NotebookLM / Perplexity / Glean, connected to external sources such as search engines (Tavily, Linkup), Slack, Linear, Notion, YouTube, GitHub and more.
<https://github.com/Decentralised-AI/SurfSense-Open-Source-Alternative-to-NotebookLM>
- ⁴ ¹² ¹⁵ ²⁹ ³⁰ ⁴⁰ ⁴³ ⁵⁷ ⁵⁸ ⁶⁰ ⁷¹ ¹¹¹ ¹¹⁶ ¹¹⁷ ¹¹⁸ ¹²⁰ SurfSense | Devpost
<https://devpost.com/software/surfsense>
- ⁵ ⁶ ²⁵ ²⁶ ⁹³ ¹¹⁵ ¹¹⁹ SurfSense: A Knowledge Graph Brain for your Web Browsing Sessions - DEV Community
https://dev.to/rohan_a426b7a9d85310208cf/surfsense-a-knowledge-graph-brain-for-your-web-browsing-sessions-3pi0
- ⁸ ⁹ ¹¹ ³⁵ ³⁶ ³⁸ ³⁹ ⁷⁸ ⁹⁸ ¹⁰⁵ SurfSense - The Open Source Alternative to NotebookLM / Perplexity / Glean : r/LocalLLM
https://www.reddit.com/r/LocalLLM/comments/1kag0nx/surfsense_the_open_source_alternative_to/
- ²⁰ ²¹ ²² ⁶⁸ ⁶⁹ Perplexity like LangGraph Research Agent : r/LangChain
https://www.reddit.com/r/LangChain/comments/1kag9yc/perplexity_like_langgraph_research_agent/
- ⁵¹ ⁸⁶ ⁸⁷ ⁸⁹ ⁹⁰ ¹⁰³ ¹⁰⁶ ¹⁰⁷ ¹²¹ Prerequisites
<https://www.surfsense.net/docs>
- ⁵² ⁶⁵ ¹⁰⁰ CHANGELOG.md - MODSetter/SurfSense - GitHub
<https://github.com/MODSetter/SurfSense/blob/main/CHANGELOG.md>
- ⁵⁶ ⁵⁹ ¹⁰¹ ¹⁰² ¹¹² ¹¹³ Docker Installation
<https://www.surfsense.net/docs/docker-installation>
- ⁶² GitHub Trending JS/TS (@github-trending-js.bsky.social) - Bluesky
<https://bsky.app/profile/github-trending-js.bsky.social>
- ⁶³ Issues · MODSetter/SurfSense - GitHub
<https://github.com/MODSetter/SurfSense/issues>

- 67 **SurfSense – Customizable AI Research & Knowledge Management Assistant**
<https://www.surfsense.net/>
- 70 **SurfSense-Open-Source-Alternative-to-NotebookLM/surfsense_backend at main · Decentralised-AI/SurfSense-Open-Source-Alternative-to-NotebookLM · GitHub**
https://github.com/Decentralised-AI/SurfSense-Open-Source-Alternative-to-NotebookLM/tree/main/surfsense_backend
- 72 73 74 75 76 77 79 81 82 **GitHub - theaiautomators/insights-lm-public: Open-source, self-hosted alternative to NotebookLM. Chat with your documents, generate audio summaries, and ground AI in your own sources—built with Supabase and N8N on a React frontend.**
<https://github.com/theaiautomators/insights-lm-public>
- 80 **I Built a NotebookLM Clone That You Can Sell (n8n + Loveable ...**
<https://www.theaiautomators.com/notebooklm-clone-that-you-can-sell/>
- 83 **Goekdeniz-Guelmez/Local-NotebookLM - GitHub**
<https://github.com/Goekdeniz-Guelmez/Local-NotebookLM>
- 84 **An Open Source Alternative to Google's NotebookLM - It's FOSS**
<https://itsfoss.com/open-notebooklm/>
- 85 **meng shao on X: "Local-NotebookLM：开源+ 本地版NotebookLM ...**
https://x.com/shao_meng/status/1898183580360360255
- 88 **LocalAI-An open-source, self-hosted alternative to OpenAI ... - Aibase**
<https://www.aibase.com/tool/29580>
- 91 **varunshenoy/GraphGPT: Extrapolating knowledge graphs ... - GitHub**
<https://github.com/varunshenoy/GraphGPT>
- 92 **GraphGPT**
<https://graphgpt.github.io/>
- 94 **[2408.14512] LLMs as Zero-shot Graph Learners: Alignment of GNN ...**
<https://arxiv.org/abs/2408.14512>
- 95 **Alignment of GNN Representations with LLM Token Embeddings**
<https://neurips.cc/virtual/2024/poster/96777>
- 96 **NOCL: Node-Oriented Conceptualization LLM for Graph Tasks ...**
<https://arxiv.org/html/2506.10014v1>
- 97 **[PDF] NOCL: Node-Oriented Conceptualization LLM for Graph Tasks ...**
<https://www.arxiv.org/pdf/2506.10014>
- 108 109 110 114 **Manual Installation**
<https://www.surfsense.net/docs/manual-installation>