# AI-Assisted Patent Search & Analysis for Redox Flow Batteries

**Authors:** Georg Niess, Stefan Spirk, Roman Kern

**Contact:** georg.niess@tugraz.at

## Abstract

Patent search is essential to identify and avoid patent infringement. By January 2025, there are more than 26,000 patents which contain the term redox-flow battery. The sheer volume makes it very challenging to identify relevant technical or legal information. Language barriers and the use of obscure or vague terminology further reduce the effectiveness of traditional keyword-based searches.

**Our project**: Over the past 18 months, we have collected information on more than 26.000 patents related to redox flow batteries and used this dataset to develop an AI-powered patent-search assistant built on large language models (LLMs).
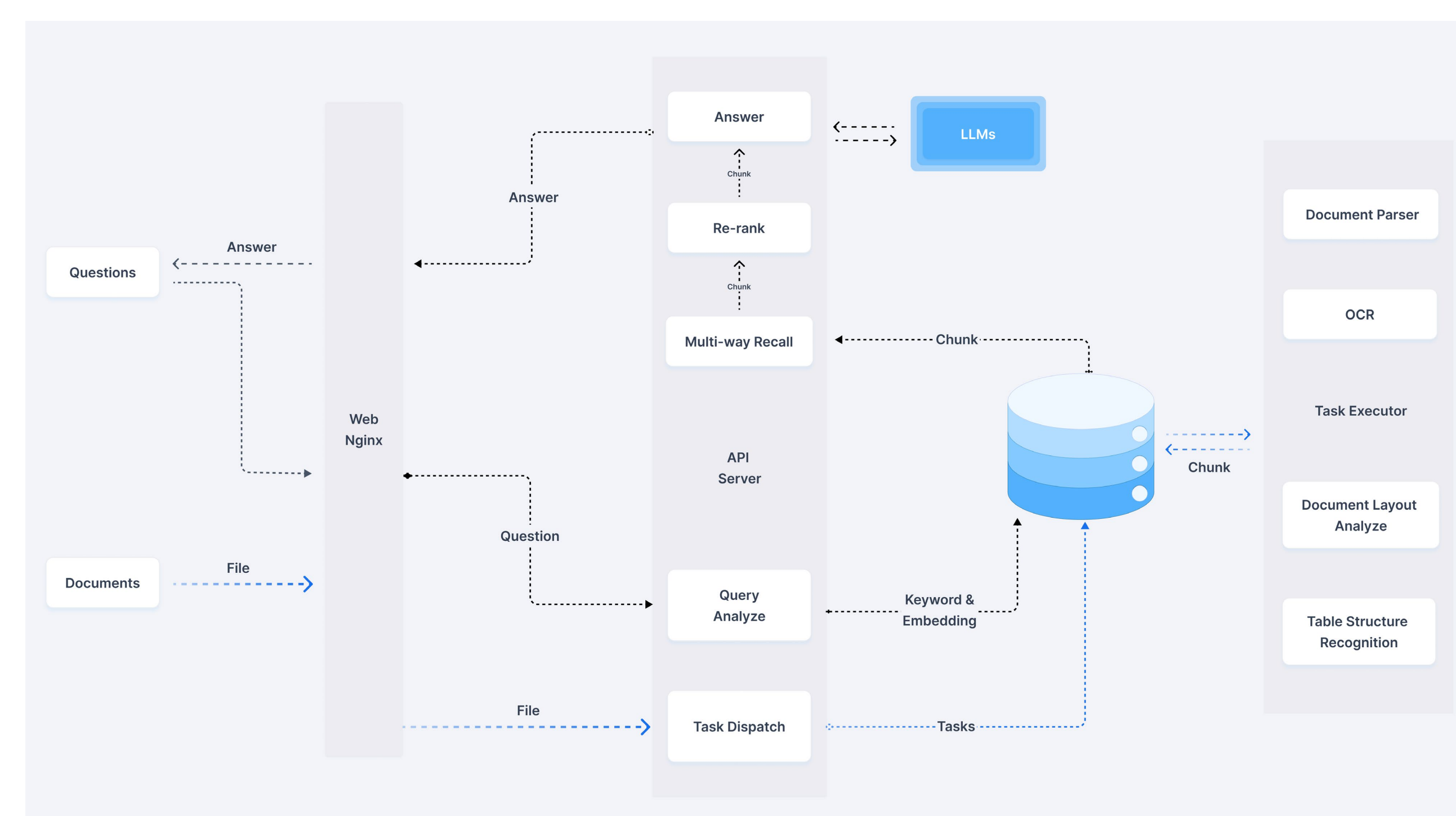
**Key features**:
1. **Embeddings-based Knowledge Base**: We embed all collected patents into a vector store, enabling semantic retrieval based on user queries.
2. **Retrieval-Augmented Generation (RAG):** Retrieved documents serve as context for the LLM, improving answer accuracy and minimizing hallucinations.
3. **Agentic Research Reports:** A series of specialized LLM agents can generate higher-level research reports (e.g., market trends) by leveraging the patent knowledge base.

**Outcome:** Our prototype shows promise in enhancing patent search and classification for redox flow battery R&D; early feedback indicates it helps project members more quickly identify relevant patents, gain insight into technology trends, and generate concise research summaries.
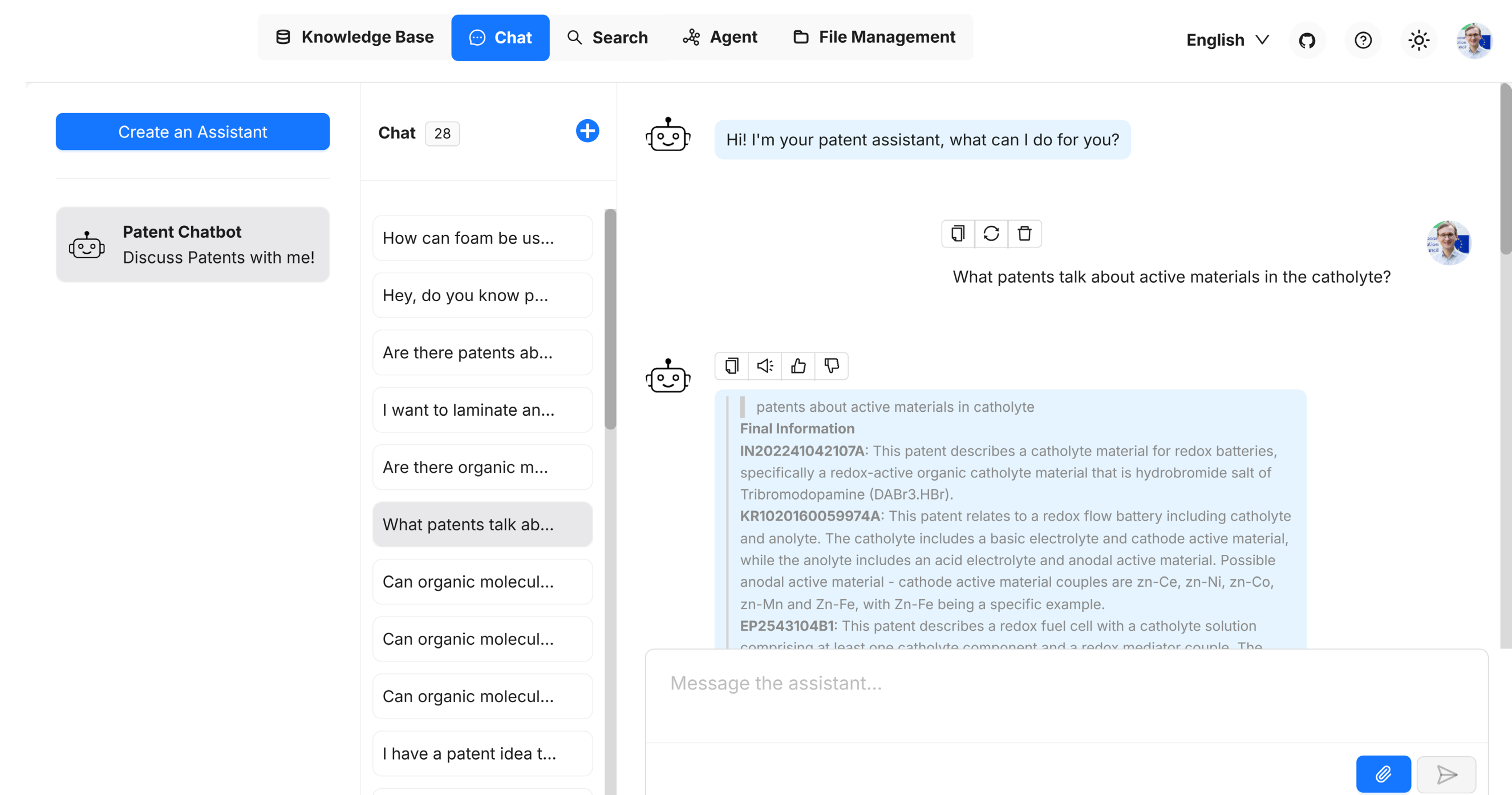
## Architecture

In our system, user-submitted questions first pass through a reverse proxy (Web/Nginx) and are routed to the API server's "Query Analyze" component, which converts each question into both keyword queries and dense embeddings. These embeddings retrieve semantically relevant "chunks" from the central vector-indexed knowledge base, while traditional keyword queries can run in parallel for a hybrid recall approach. Retrieved chunks are then re-ranked, and the top passages are passed (with the original question) to a LLM to generate a precise answer, which is sent back through the proxy to the user. When a user uploads a document, Nginx forwards the file to the "Task Dispatch" component, which creates OCR, document parsing, layout analysis, and table-structure recognition jobs. The Task Executor processes each job, splits the document into smaller text or table chunks, and indexes those chunks (with both keyword indices and embeddings) into the same central knowledge base. This unified design allows the LLM to answer queries using both previously ingested and freshly uploaded files, ensuring that all information is semantically searchable and available for RAG-based retrieval.
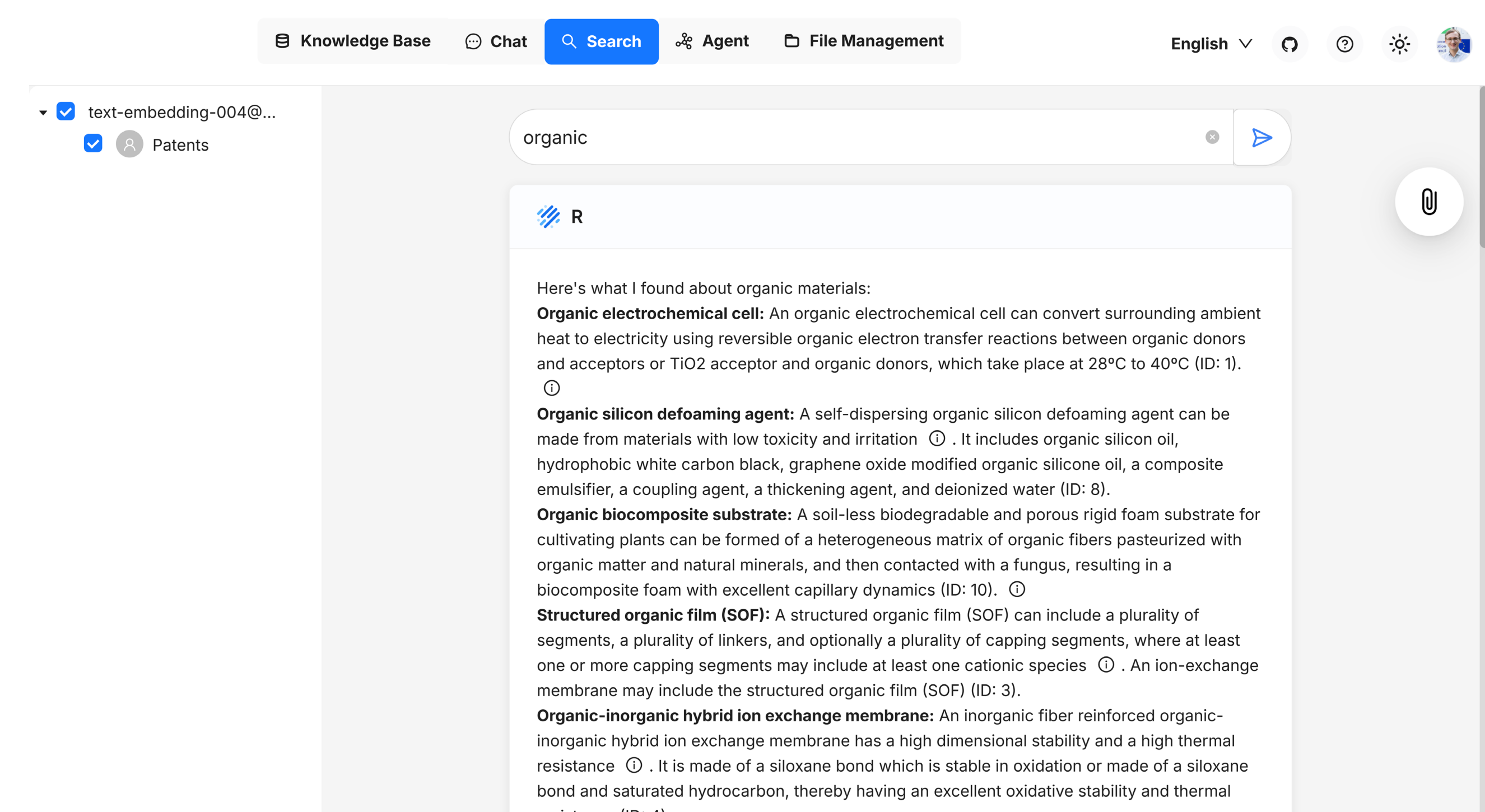


**System Architecture:** End-to-end workflow illustrating how user questions and uploaded documents are handled by the API server, knowledge base, and task executor modules.
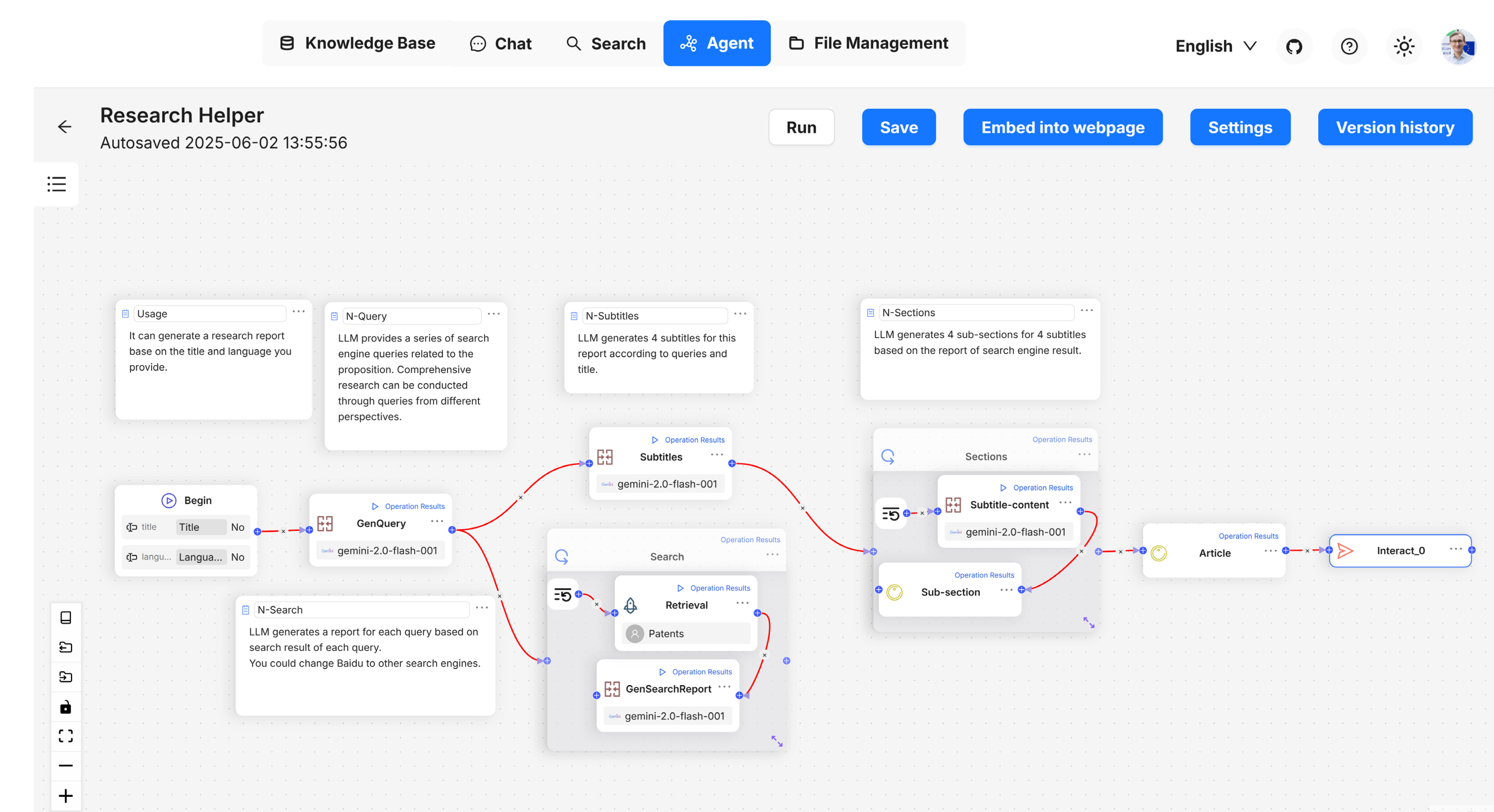
## Examples



*Chat Mode: Interactive Q&A interface over the entire patent corpus.*



*Search Mode: Keyword/phrase search augmented with semantic ranking.*



*Research-Agent Mode: Automated agents that can process user-level "research questions" (e.g., "What are the recent material innovations in redox flow electrolytes?") and generate structured reports.*



Want to see some example reports? **Scan me!**

**TU Graz – Institute of Machine Learning and Neural Computation**
8010 Graz, Inffeldgasse 16b/I, Austria, Tel.: +43 316 873-5811
georg.niess@tugraz.at