

Real-time Inference for the Real World

A USE CASE SERIES

[vectorize]

Vectorize aims to help developers and enterprises build fast, accurate, production-ready Gen AI applications in hours rather than weeks and months. Customers can turn unstructured data into perfectly optimized vector search indexes, purpose-built for retrieval augmented generation.

The Groq LPU™ Inference Engine delivers game-changing speed, proving that large language models (LLMs) can power applications with even the most demanding latency and performance requirements. However, even the fastest LLM generation can't overcome knowledge gaps when an application requires context that the LLM wasn't trained to understand.

That's why we're excited to announce a new collaboration between [Groq](#) and [Vectorize](#) to bring together the world's fastest real-time LLM inference and the world's most powerful Retrieval Augmented Generation (RAG) experimentation and pipeline platform.



The Groq LPU™ Inference Engine
delivers game-changing speed,
proving that LLMs can power
applications with even the most
demanding latency and
performance requirements.

Bridging the Knowledge Gap with RAG

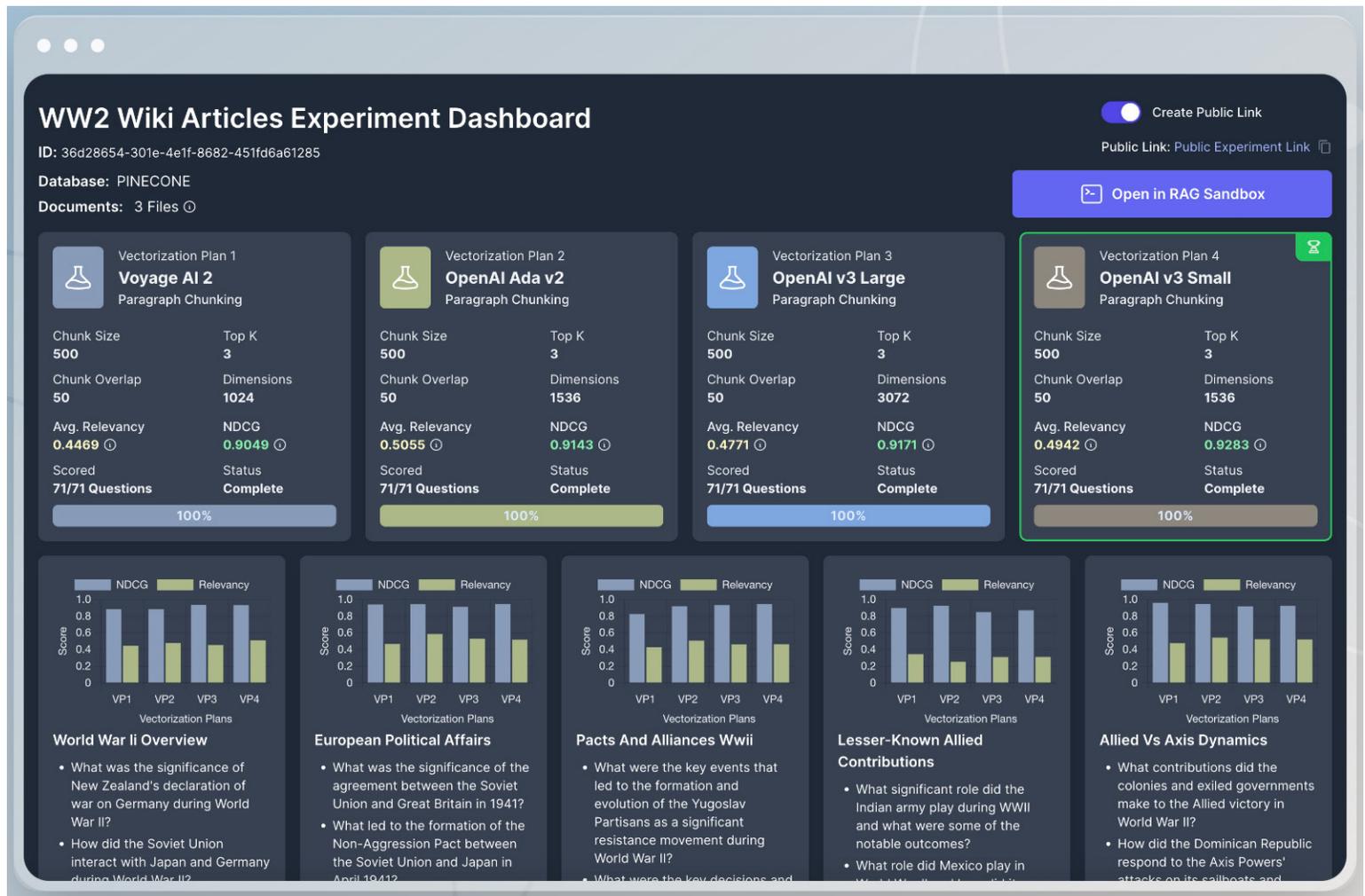
There's a large opportunity to generate contextually relevant responses for subjects outside of the LLM's training data. This is where [RAG](#) has emerged as the de facto solution. By incorporating relevant context into the LLM prompt, RAG addresses this inherent limitation of LLMs, enabling the LLM to generate responses even for subjects on which it wasn't trained.

RAG essentially creates a bridge between the LLM and an external set of knowledge in the form of data, but getting access to that data is not always straightforward. That's because it may live in file systems, knowledge bases, SaaS platforms, or other databases. Data extraction can be tricky, and once you have it extracted, figuring out how to make that data not only searchable but also ensure search results are relevant and accurate adds to the challenges.

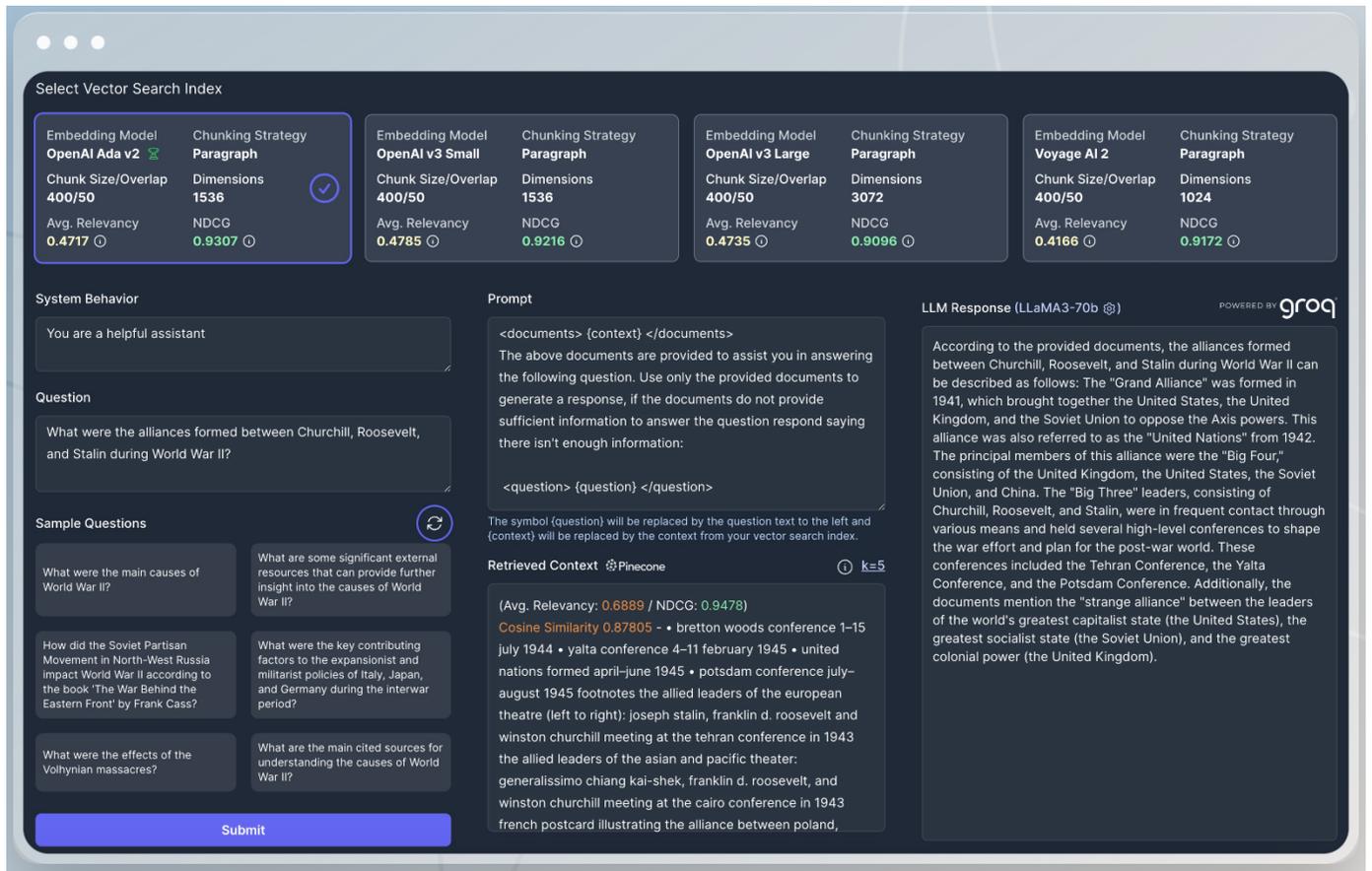
To build an optimized search solution, most AI engineers and developers rely on a vector database to provide semantic search on their data. Building a search index that is optimized for RAG involves finding the right chunking strategy and text embedding model to ensure accurate retrieval and relevant context.

Vectorize streamlines this process, thereby ensuring the LLM has the necessary context to generate an accurate response. It does this through a powerful experimentation framework that allows users to quickly analyze how well different vectorization strategies work for your unique set of data.

Vectorize automatically analyzes a representative sample of data and then handles the extraction and vectorization processes using a combination of approaches. The net result is concrete data that defines which embedding model, chunking strategy, and retrieval settings produce the most relevant context for a customer's [RAG pipelines](#).



What's more, users can immediately test newly vectorized data with [Llama 3 or any other Groq-powered LLM](#) in the RAG Sandbox. This combination of quantitative data with a qualitative end-to-end assessment in the RAG sandbox means users can confidently launch LLM-based applications knowing your customers will receive fast, accurate responses.



Fast Inference for RAG Applications

Developer velocity is critical for enterprises; it allows them to get the highest quality GenAI workloads to market faster and cheaper. **Vectorize becomes a huge asset in this quest as it compresses the development process of assimilating data from disparate sources into a RAG-based solution from days or weeks into minutes or hours.** This only works because the platform is fast. Speed encourages more interactivity and experimentation. Developers have the flexibility to test and fine-tune various model configurations until they get it just right.

Groq ultra-low latency, and the user experience it enables, is critical to Vectorize's performance. Groq powers the Vectorize sandbox and gives developers a tool to test various model configurations by experimenting with different queries and seeing instant responses. Combining Vectorize's data driven experiments with Groq ultra-low latency enables developers and businesses to continuously refine and optimize their RAG configurations.

Exceptional performance from Groq enhances a broad range of real-time AI applications like call center and sales co-pilots, real-time customer experience customization, customer service bots that actually help people, personalized recommendation systems, and dynamic content delivery. It significantly improves user experiences and accuracy, encouraging greater interactivity.

Likewise, this immediate feedback is crucial in areas like dynamic content delivery, where the speed of content updates can directly influence user engagement and satisfaction. Publishers and content managers can instantly adjust offerings based on real-time user behavior and feedback, ensuring that the content remains relevant and engaging. The ability to quickly process user data and update content accordingly helps to achieve a hyper-personalized user experience that keeps individuals engaged and returning.

Getting Started with Groq & Vectorize

Sign up for a Vectorize account at platform.vectorize.io and visit the quickstart documentation to see for yourself how Groq and Vectorize help you build better RAG applications. When you're ready to build your RAG application, sign up for an API key on GroqCloud™ at console.groq.com/login, and you'll be up and running in no time!