Building Agentic RAGwith ADK + Vector Search



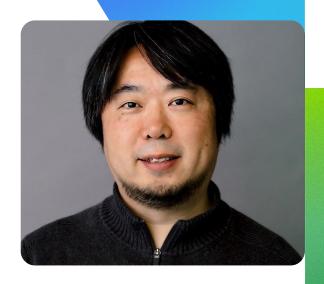


Kaz Sato

Google Cloud Cloud Al Developer Advocate







01. Agent dev stack on Google Cloud

To build an agentic app, you'd need...

Toolkit

for standardized & efficient development

Protocols

for connecting LLMs, tools & agents

Agent platform

for managing deployment and operations

Agent dev stack on Google Cloud







Model Context Protocol



Vertex Al Agent Engine



Agent2Agent (A2A) protocol

Open-source, code-first toolkit for building, evaluating, and deploying Al agents.

Open protocol that standardizes how applications provide context to LLMs.

Managed platform to deploy, manage, and scale Al agents in production.

Open standard designed to enable communication & collaboration between Al agents.

Agent Development Kit (ADK)



Multi-agent framework by Google

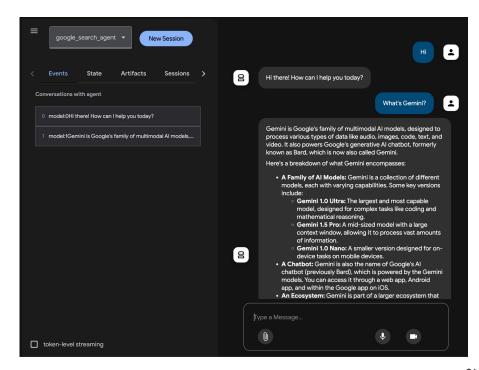
Actions driven by gen Al reasoning

Supports **Gemini** and **third party**

models with LiteLLM

Live audio, video, screen bi-directional streaming

adk web tool



MCP: Model Context Protocol

MCP is an open protocol that enables seamless integration between Al apps & agents and your tools & data sources.

APIs

Standardize how **web applications** interact with the **backend**:

- Servers
- Databases
- Services

LSP

Standardizes how **IDEs** interact with **language-specific tools**:

- Code navigation
- Code analysis
- Code intelligence

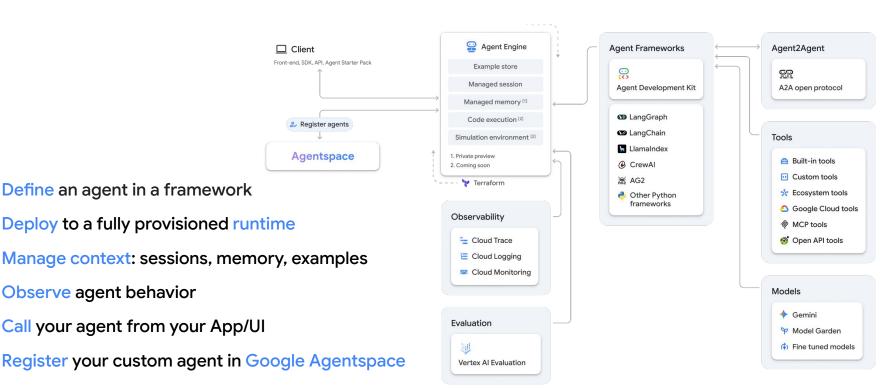
MCP

Standardizes how AI applications interact with external systems:

- Prompts
- Tools
- Resources

Vertex Al Agent Engine

Observe agent behavior



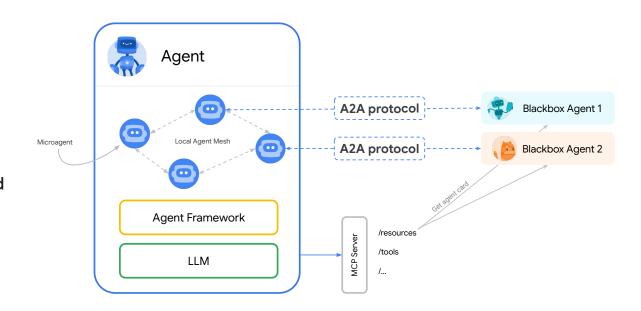
A2A (Agent2Agent)

Open protocol to handle agent collaboration

Enables agent discovery using Agent Card

Enterprise-ready for a secure and governed experience.

Built on top of existing, popular standards including HTTP, SSE, JSON-RPC



02. Vector Search advanced practices for Agentic RAG

A typical RAG scenario for e-commerce

Can you find Google Pixel 9?

Sure! I will use my Vector Search to find it!













User

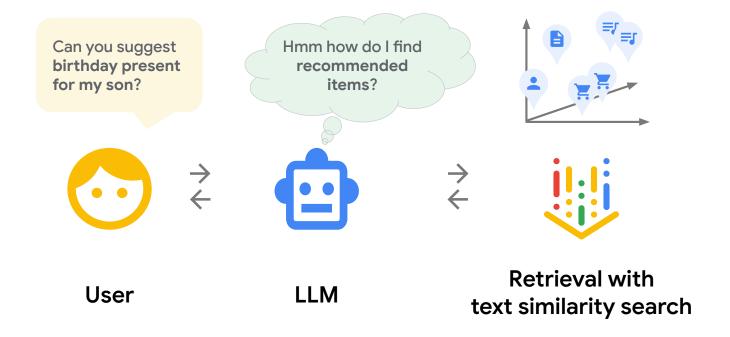
LLM

Retrieval with text similarity search

Challenge: Multimodal / Keyword search

Can you find cups Hmm how do I with dancing handle the image / figures? The SKU is keyword search..? #123-ABC. Retrieval with User LLM text similarity search

Challenge: Recommendations



What makes search quality better





Multimodal search

Use images, audio and video semantics



Hybrid search

Use both semantic and keyword search



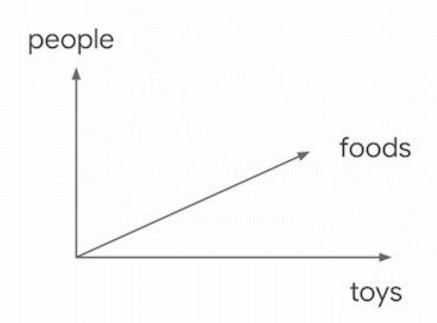
Task type embedding

To find by relevance, not by similarity

Multimodal Search

Using an embedding model that understands text and image





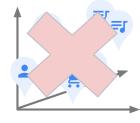
Hybrid Search

The weakness of the semantic search

Proprietary

Can you find item# 123-4567? **ASDFGHJ?**

Hmm... what do they mean???

















User

LLM

Vector Search

Contents

Hybrid Search

Combining the strength of semantic and keyword search

Can you find items matching "Kids"?

Here's the results with the **Hybrid Search!**

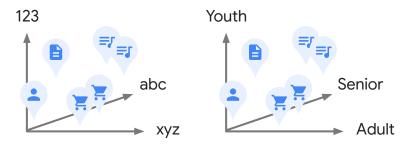








LLM



Google Blue Kids Sunglasses Google Red Kids Sunglasses YouTube Kids Coloring Pencils YouTube Kids Character Sticker Sheet

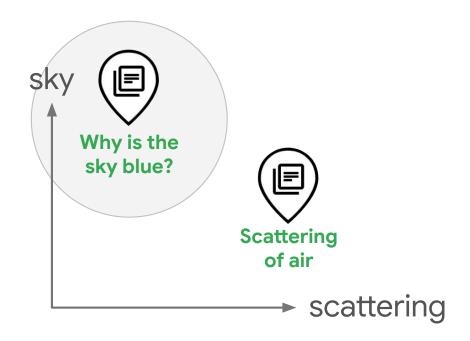
Google White Classic Youth Tee
Google Doogler Youth Tee
Google Indigo Youth Tee
Google Black Classic Youth Tee
Chrome Dino Glow-in-the-Dark Youth Tee
Google Bike Youth Tee

Demo

Multimodal search and Hybrid Search

Task type embedding

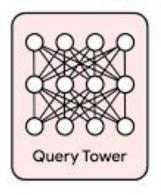
Questions and their answers are not semantically similar

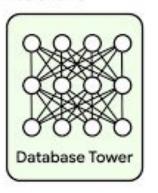


Task type embedding

Learns the relationship between questions and answers



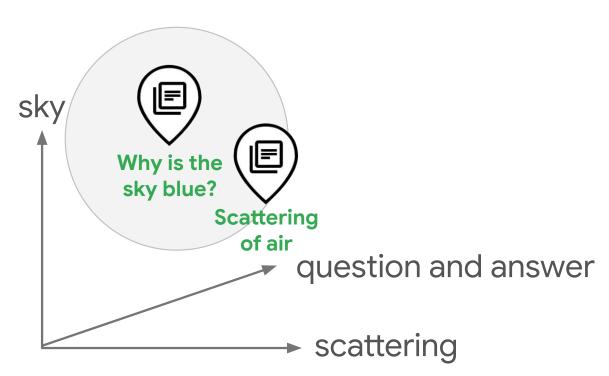




Google Cloud Next Tc.., 2

Proprietary

Maps the question and answer closer in the embedding space



Demo

Task type embeddings

03. Building Agentic RAG with ADK + Vector Search

Challenge: *Smart* Recommendations

Hmm how do I make Can you suggest birthday present for my son? recommendations What's the latest trends? with latest trends? Retrieval with User LLM text similarity search

Solution: Al Agents +

Vector Search

Any birthday present for my 10 yrs son?

Hey Search Agent, run 20 queries for "Science Toys"

Got it! Running 20 queries in parallel...

"STEM toys age 10"

"science kits 10 year old boy"

"experiments for kids 10 years old"















User

UI Agent

Search Agent

Google Search & **Vector Search**

Demo

Shopper's Concierge

Shopper's Concierge

Research with Google Search to find 5 item categories (~500 ms)



Any birthday present for my 10 yrs son?

Hey Google Search, what do people buy for birthday presents?









Sport goods, Outdoor toys, Science Toys

User

UI Agent

what's the title of the consumer at the highest level of a food chain



Google Search grounding

Food Chain | National Geographic Society

Jan 21, 2011 - The second trophic level consists of

organisms that eat the producers. These are called primary

Shopper's Concierge

Generate and run 20 queries for each category (From 3 million items to a few hundreds, in ~500 ms)

Proprietary

Hey Search Agent, run 20 queries for "Science Toys" Got it! Running 20 queries in parallel...

"STEM toys age 10"

"science kits 10 year old boy"

"experiments for kids 10 years old"

...















User

UI Agent

Search Agent

Vector Search

Shopper's Concierge

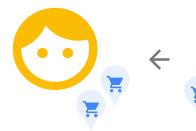
Personalized multimodal item curation for the user (From a few hundreds to tens, in ~4 sec)

Proprietary

Thanks!
These look amazing!

Hi User, Here's the result! Looking at the images and descs, choosing relevant items to the user intent and preference...















User

UI Agent

Search Agent

Vector Search

Multimodal item curation with Gemini 2.0 Flash

User intent: "decorating a living room", searching for "Wall Art & Mirrors"



04. Summary

Multimodal, multi-agent RAG with Generative Recommendations





Multimodal,
Multi-agent search
with ADK

Real-time multimodal communication



Generative Recommendations

Google Search grounding
Query generation

Multimodal item curation



Vector Search

Multimodal search

Hybrid search

Task-type embeddings

Ranking API

Get Started



Vector Search docs





ADK recsys sample





@kazunori279

