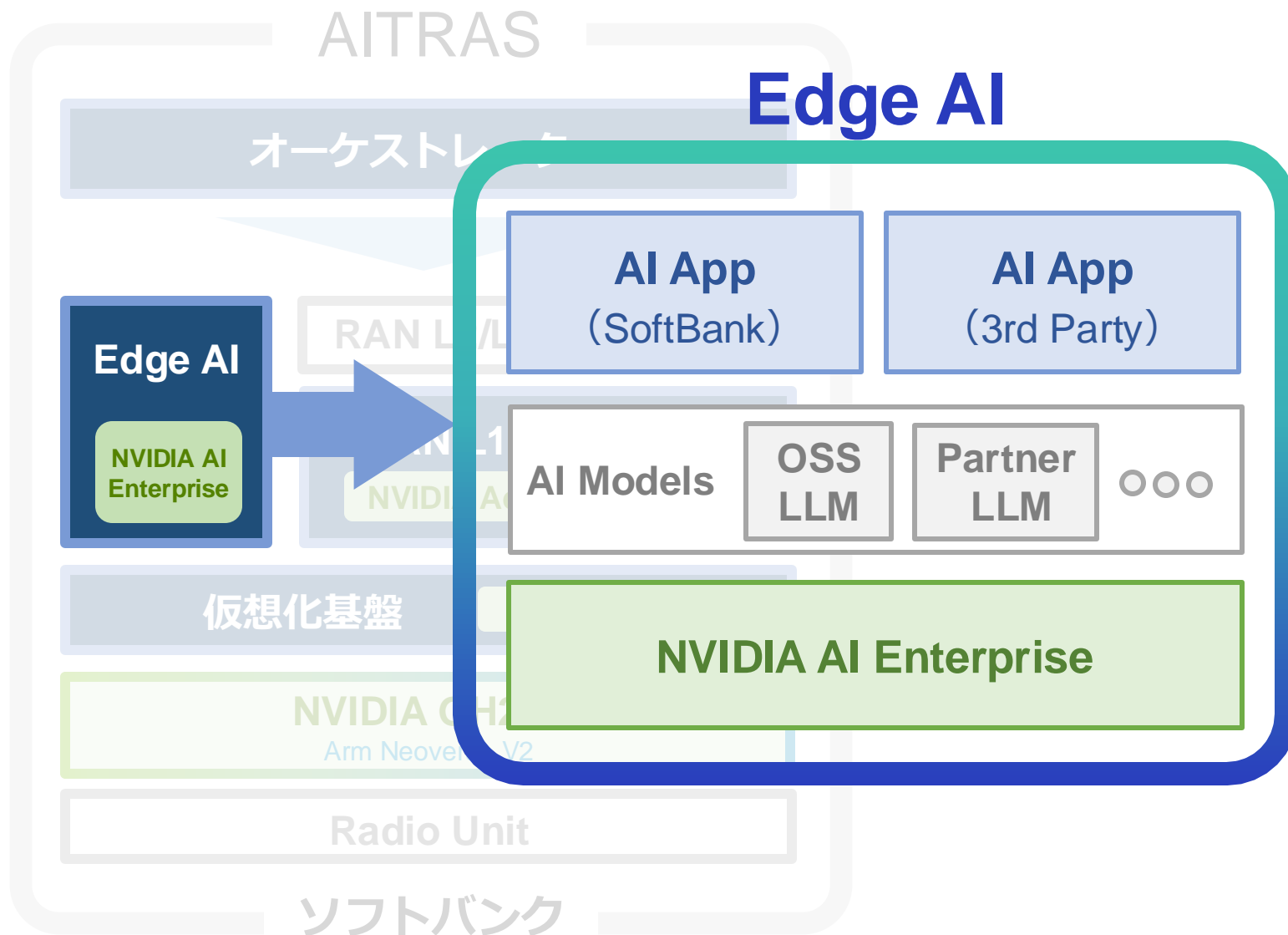


RAG Menu @Edge

Key Technology of AITRAS : AI Edge



Original AI Service

- Autonomous Driving Monitoring App
- Cloud Robot
- **RAG Menu @edge**

NVIDIA Serverless API

Allocates unused resources to meet customer demands on demand

NVIDIA AI Enterprise

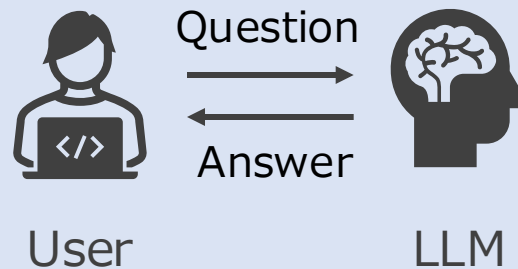
Industry-standard AI framework

What is RAG?

(Retrieval-Augmented Generation)

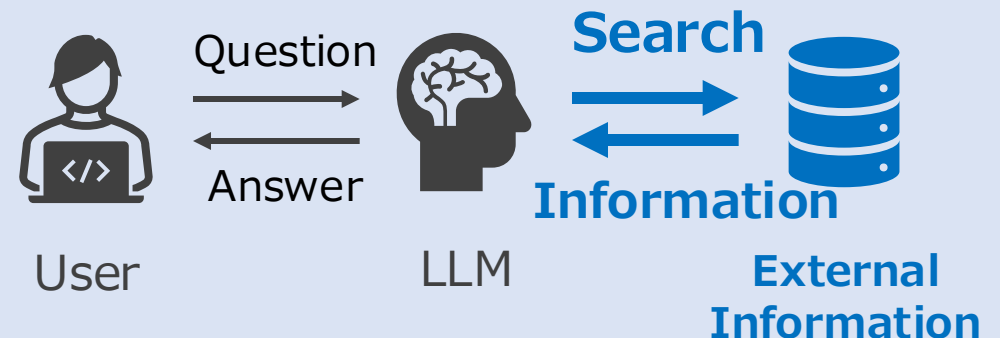
Allows answering information that LLMs do not know, such as internal confidential information

✓ General LLM



Answers based only on pre-trained information

✓ RAG

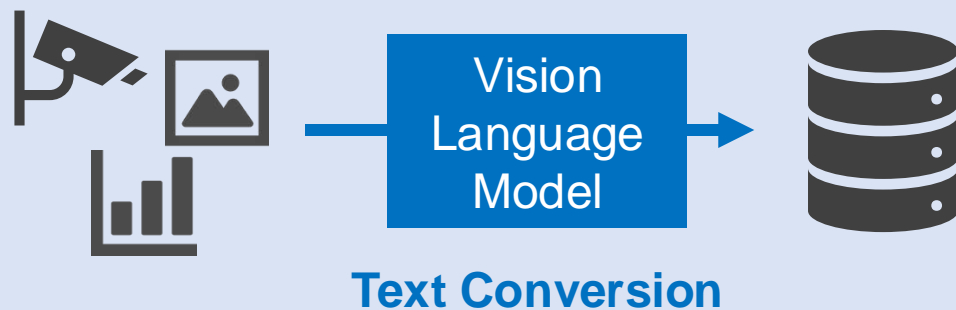


Answers based on pre-trained information + **External Information**

Multimodal RAG

Handling various data such as images and audio with RAG

Visual Information



Example: Diagrams and graphs

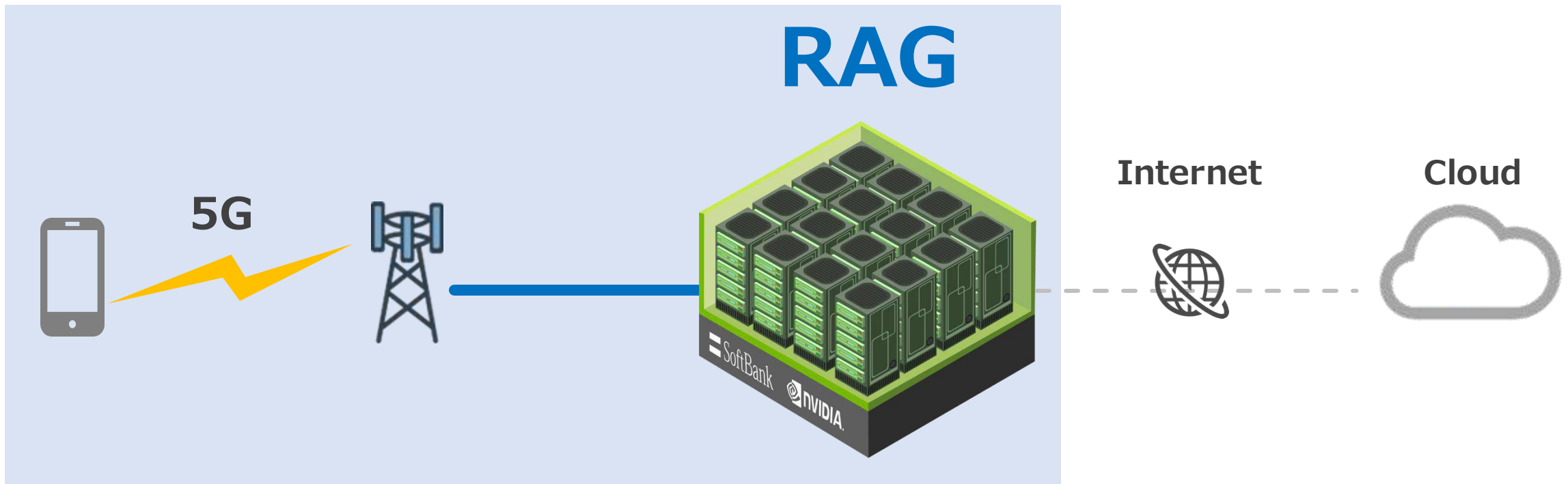
Audio Information



Example: Meetings, Call Centers

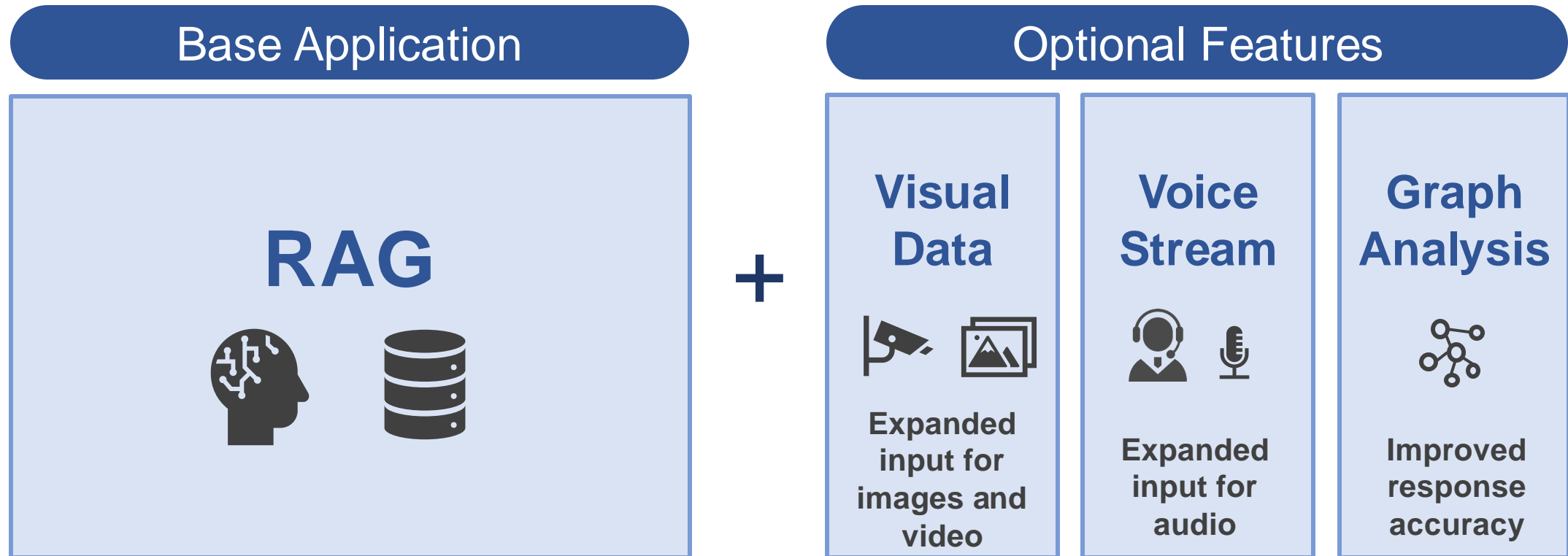
RAG Deployment at the Edge

RAN is a closed network
→ **Effective for ensuring the security of RAG using confidential information**



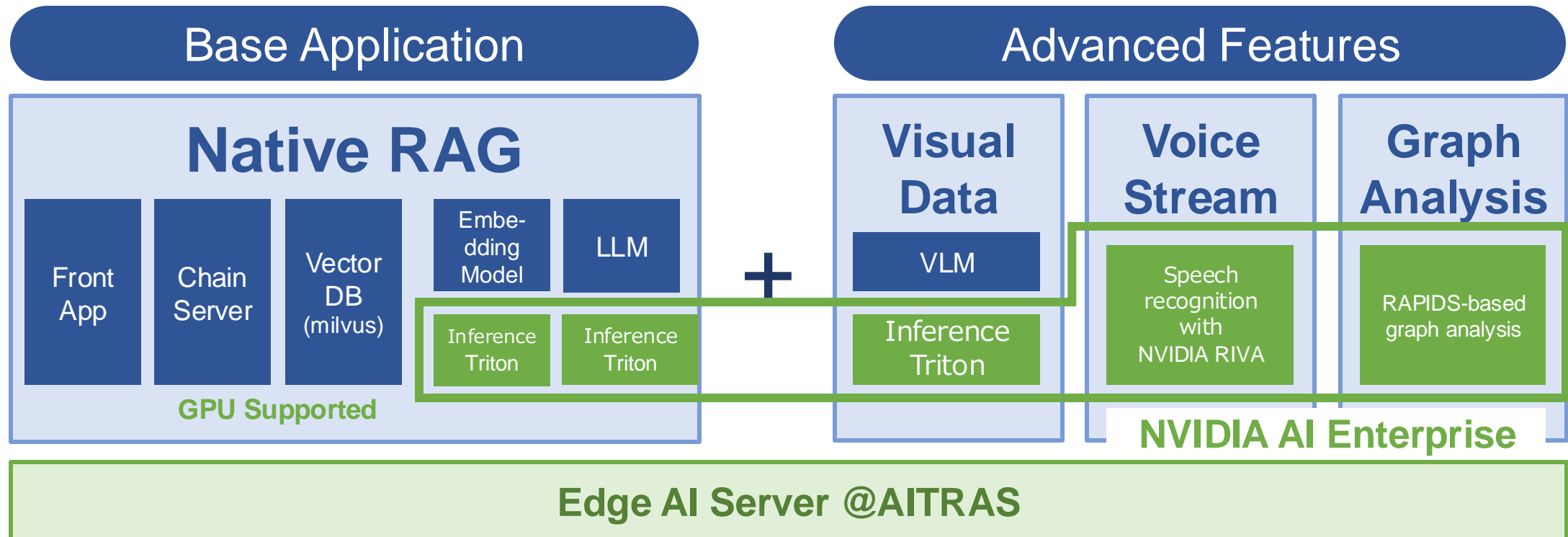
RAG Menu @Edge

Additional Options to Meet Customer Needs

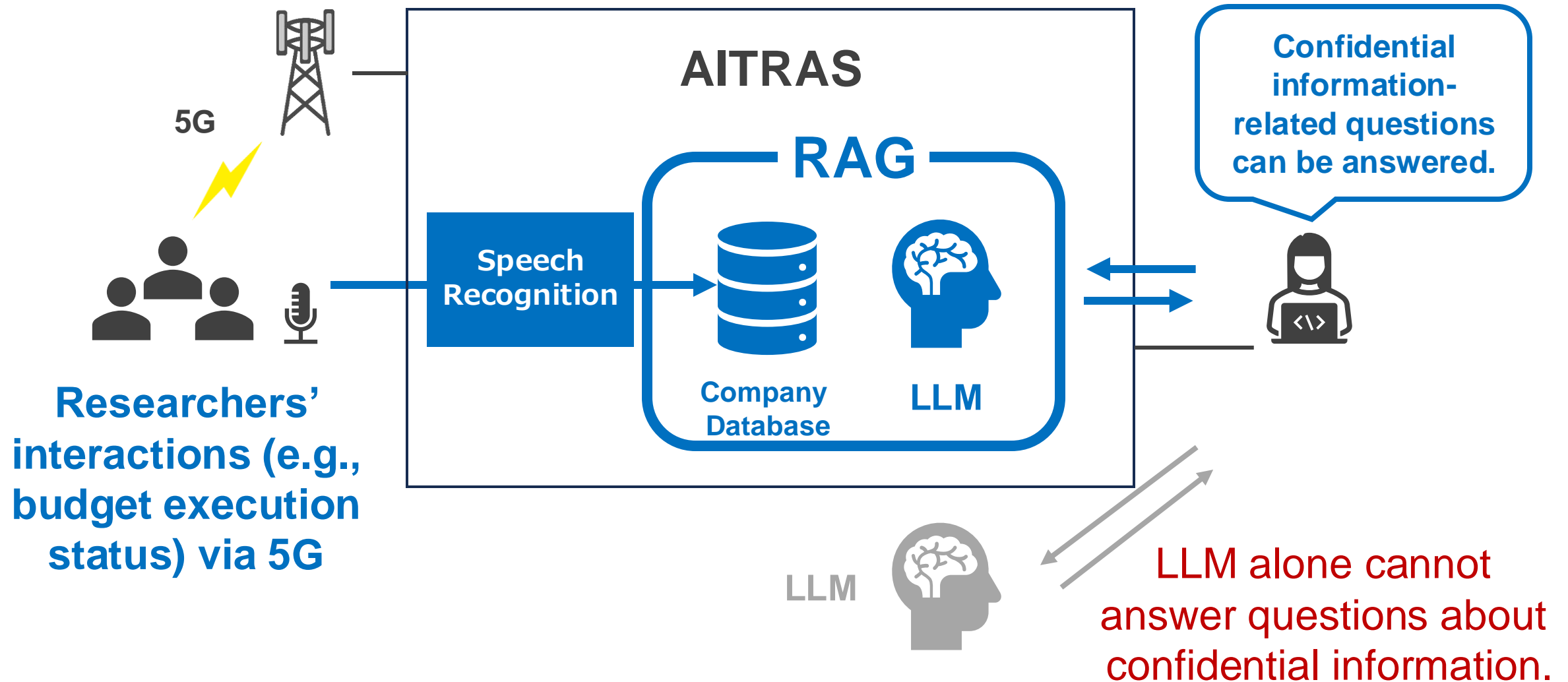


Provided Features

Delivering optimal GPU-enabled feature sets implemented with NVIDIA AI Enterprise



Demo : Research Institute of Advanced Technology RAG



RAG DEMO Converse

オーディオファイルを選択

SB_Sentan_Budget_Status.wav

Play

Input the audio data into the RAG

Clear

LLM Chatbox

RAG

テキストを入力してENTERを押してください

Example : Enterprise Scenario

Improving Factory Operations Efficiency with RAG

