



Harnessing Generative AI and Large Language Models With Vision AI Agents

Padmavathy Subramanian, Kaustubh Purandare, Subashree Radhakrishnan

GTC 2024



Agenda

- Why Multi-modal AI
- Introducing Visual Insight Agent
- Architecture
- Workflows
- Vision-Language Models – Customization
- VLM Challenges
- CTA

Why Multi-modal AI?

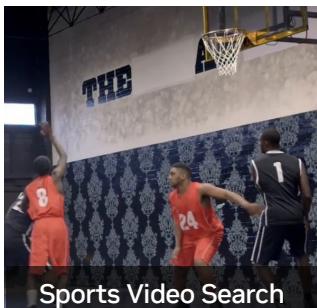
Multi-modal models autonomously provide video understanding and insights



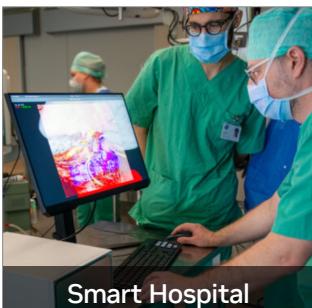
Automated Inspection



Process Automation



Sports Video Search



Smart Hospital

- AI automation for every industry
- Multiple sensors
- Petabytes of multi-modal data



Challenges

- Petabytes of videos, images
- Live streams
- Scene understanding
- Proactive alerts



Requirements

- Custom data
- Accuracy
- Context
- Video, audio, OCR
- Low latency



Are Pallets Being Moved Out?



Patient Emergency Alert

- Cloud native microservices
- Ability to customize models
- Ease of creating vision AI agents

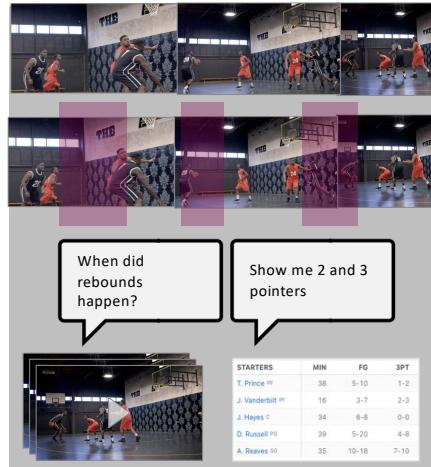
- Enterprise operational excellence
- Natural language interaction
- Contextual insights, proactive alerts

Generative AI use cases across Industries



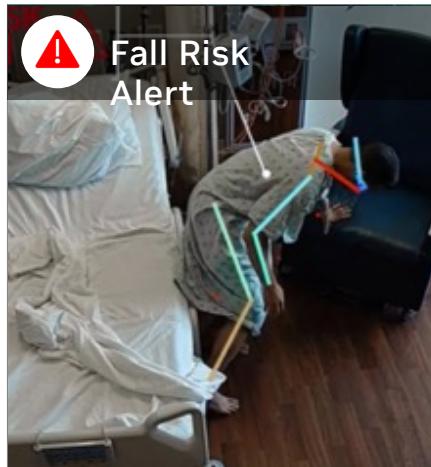
Contextual Understanding

Generalization



Spatio-temporal Localization

Insights correlating space and time, Q&A



Video Search

Highlights, proactive alerts

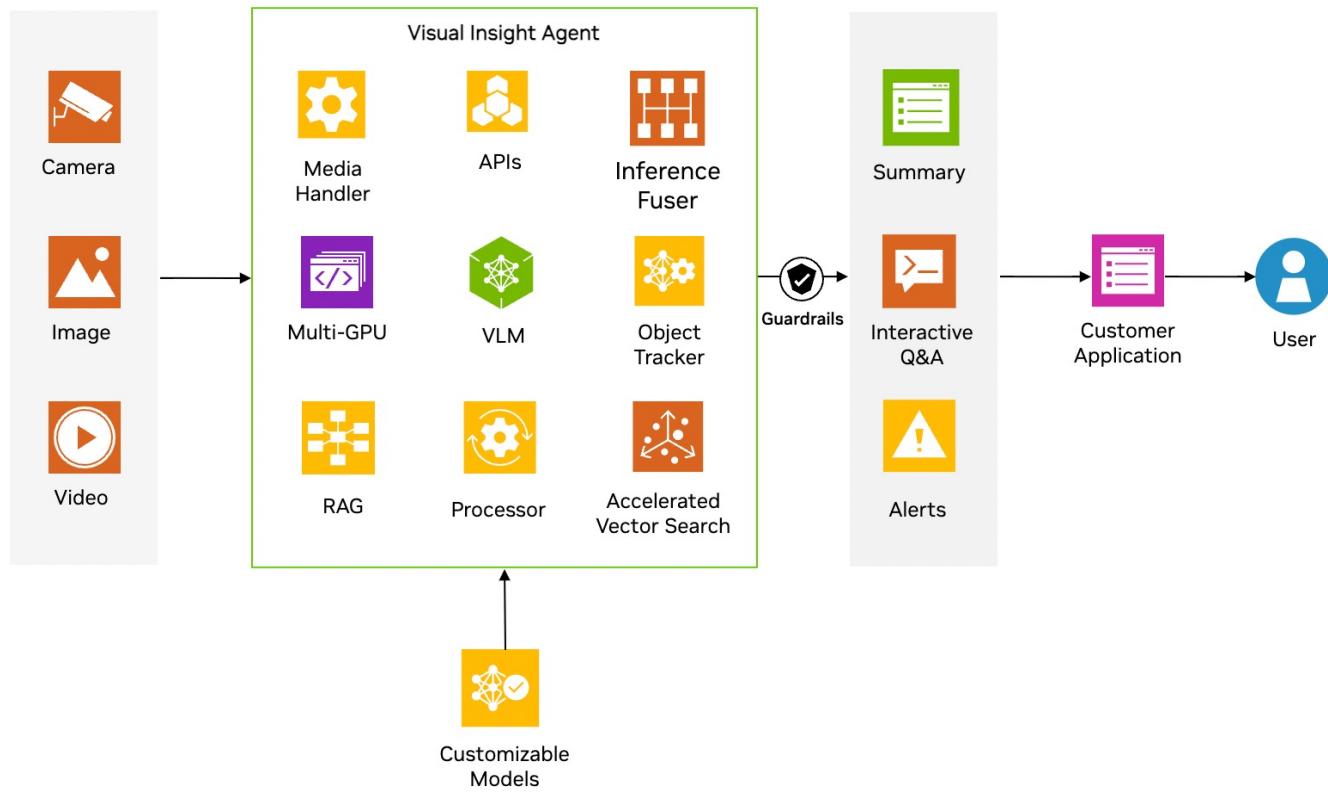


Video/Image Summarization

Complex reasoning

Introducing Visual Insight Agent

AI workflows to build Vision AI agents with VLMs



Build next wave of Vision AI agents using VIA

Leverage VIA AI workflows with disaggregated, accelerated inference for insights using VLMs

Build your AI agent using VIA Workflows!



Customize VLM with your data
— Optimize with TRT-LLM



Microservice for ingestion,
data disaggregation, scheduling



Microservice for inferencing,
object tracking.



Microservice for inferred multi-modal
data, metadata



Accelerated vector search
Spatio-temporal understanding

Accelerate inferencing on large amounts of videos, images and live streams!



High performance — Parallelism
techniques



Latest chunking algorithm
— Aggregation of inference



Multi-GPU multi-node
disaggregation & acceleration



Low latency — High throughput
Intra/inter-node communication

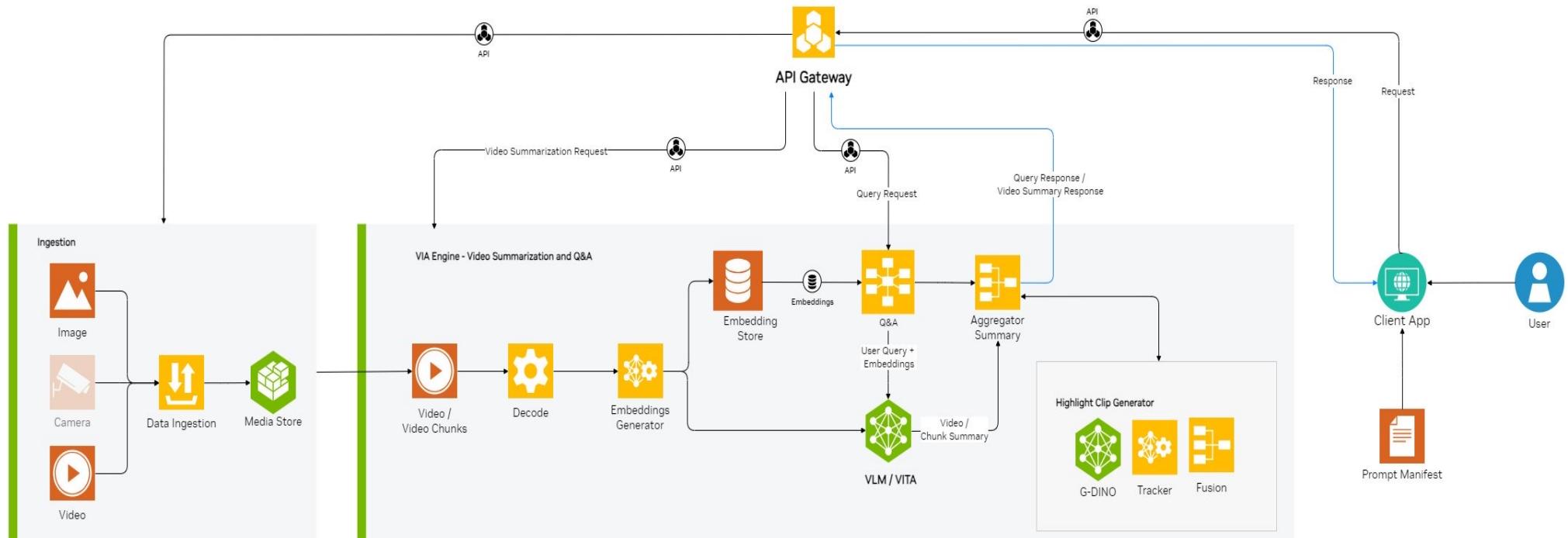


Create AI agent with VLM, RAG,
Guardrail, Analytics & more!



VIA – Workflows for building Vision AI agents

VIA Architecture



Summarization Workflow Demo

The screenshot shows the NVIDIA Visual Insight Agent interface, specifically the "INTERACTIVE Q&A" tab. The interface is designed for summarizing video content.

Left Panel (Interactive Q&A):

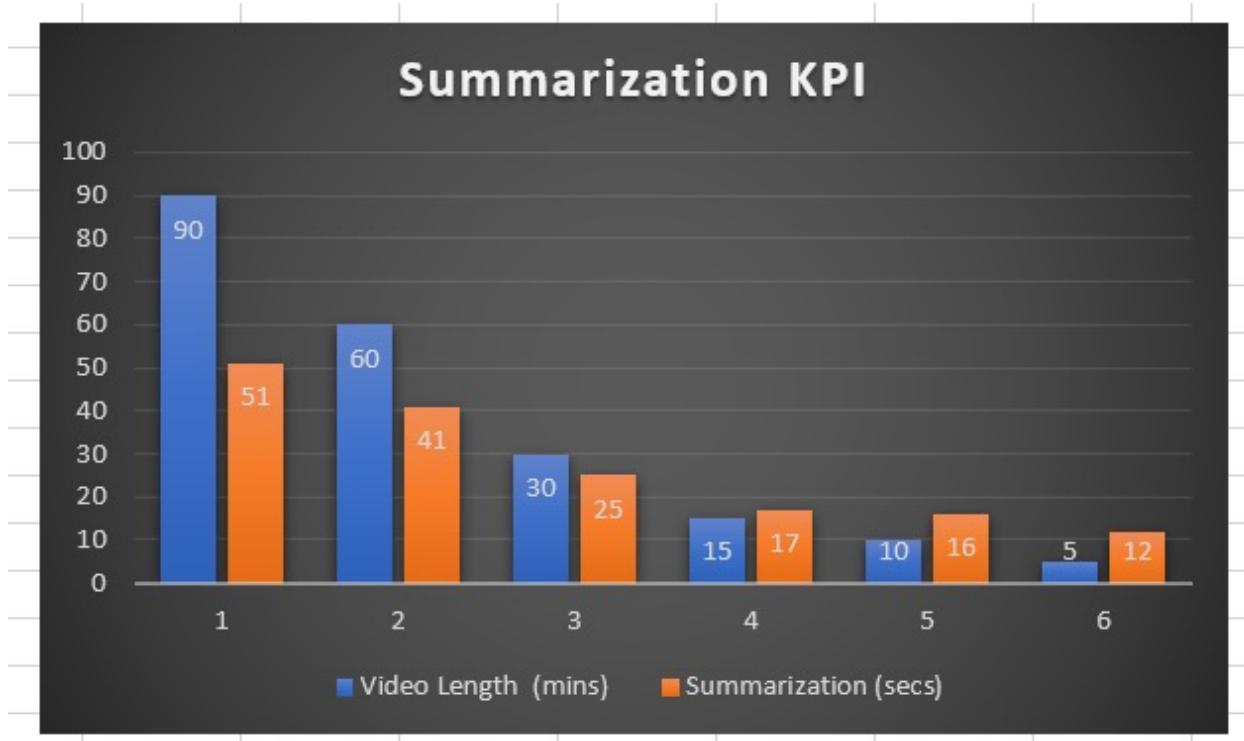
- Upload Area:** A large box with a green "Video" button and a "Drop Video Here" placeholder. Below it is a "CHUNK SIZE" dropdown set to "1 min".
- Select a Sample:** A section showing a thumbnail of a basketball game and the label "All-Star".
- Buttons:** "Upload & Start Chat" (green), "Show parameters", "Clear Chat", and "Restart App".
- Checkboxes:** "Enable Guardrails" and "Aggregate Chunk Responses".

Right Panel (Prompt Manifest):

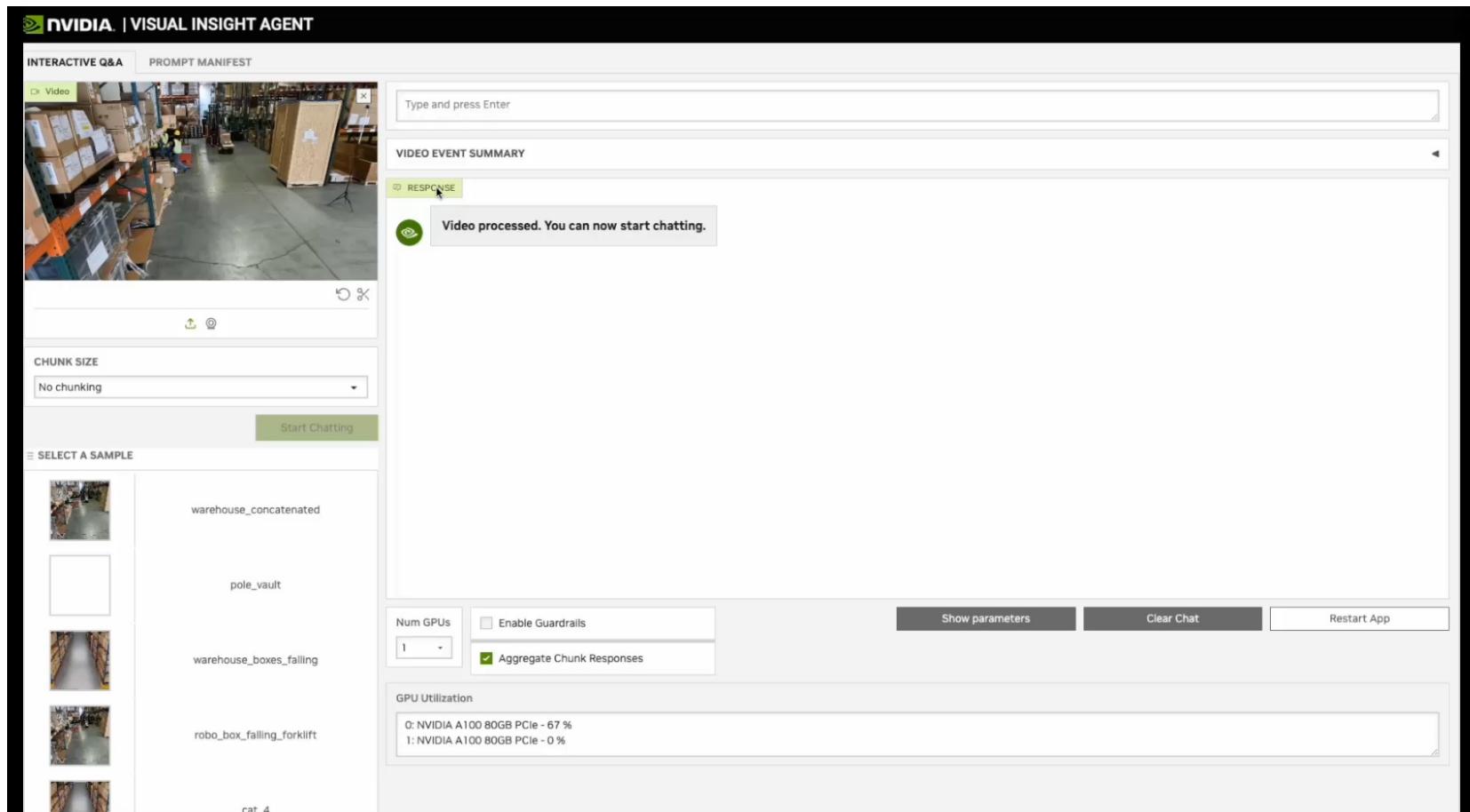
- VIDEO EVENT SUMMARY:** A large empty area labeled "RESPONSE".
- Text Placeholder:** "Upload your image/video first and then click on Start Chatting, or directly click the examples at the bottom of the page."

Performance on Long Video

End to End



Interactive Q&A Workflow Demo



Spatial Grounding on Video

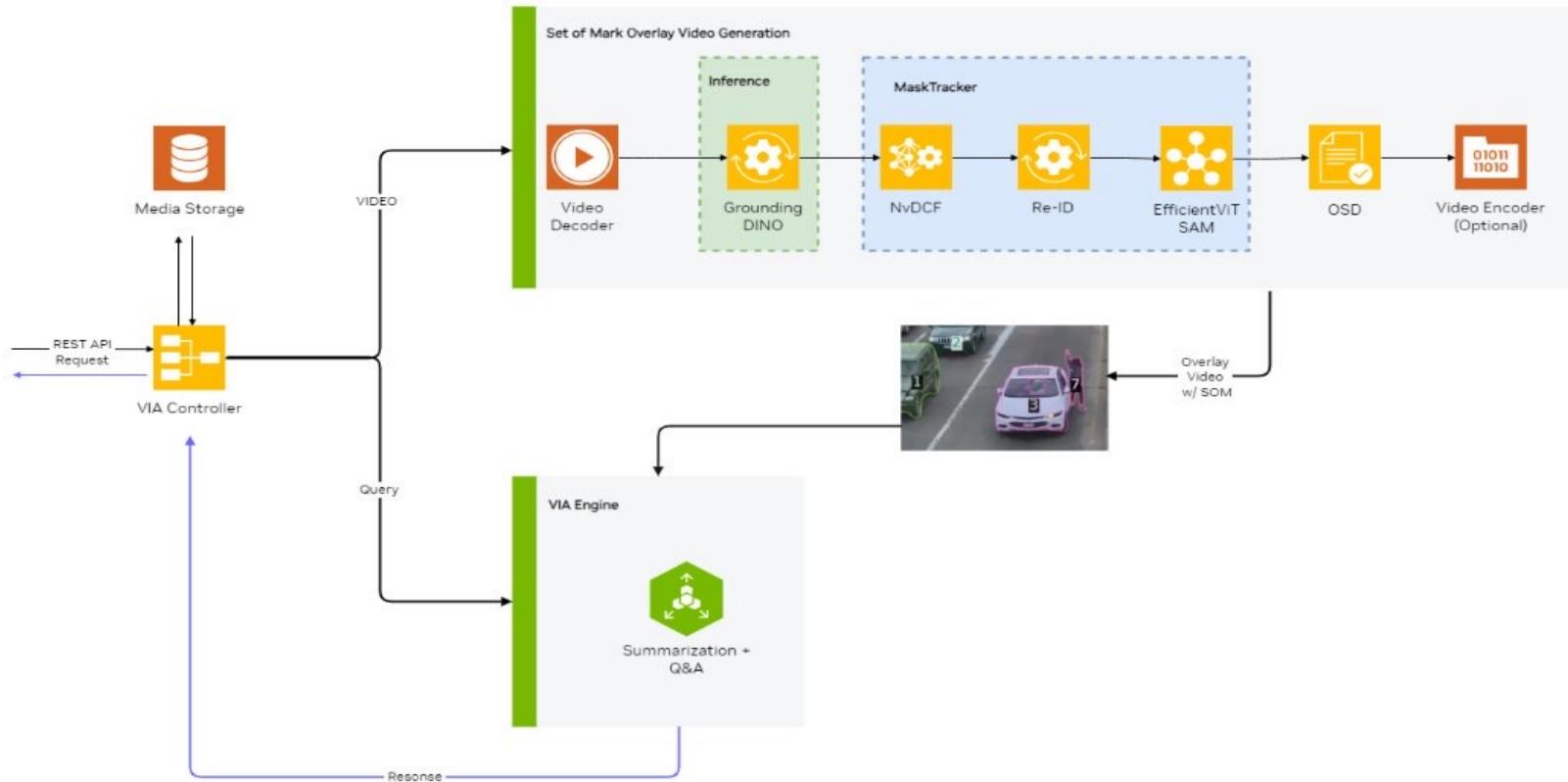


Key Challenges

- Accuracy Challenges - To help VLM to better understand and describe the scene in the video with localization info.
- Persistent ID over many frames throughout the entire video.
- Performance - spatial grounding is expensive along segment anything for every frame, especially for long Videos

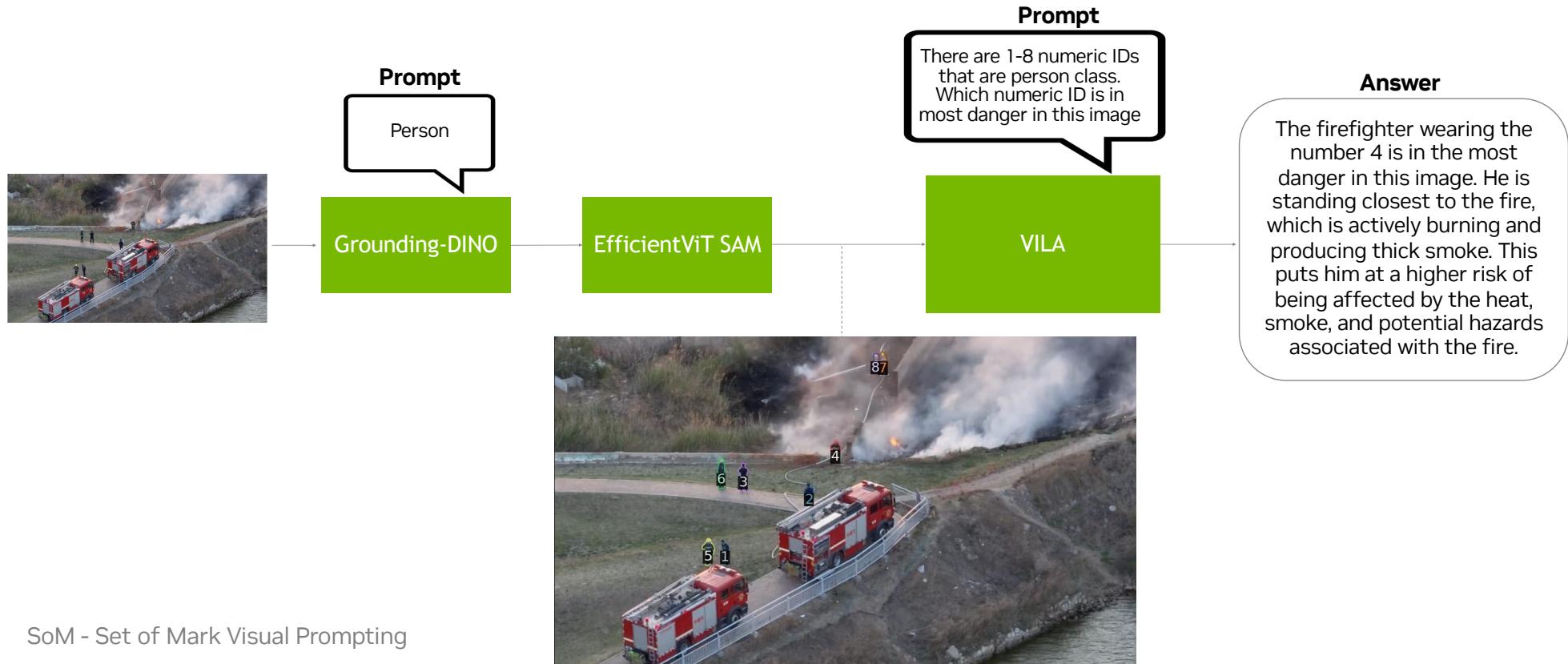
Interactive Q&A with Spatial Grounding on Video

Workflow Architecture



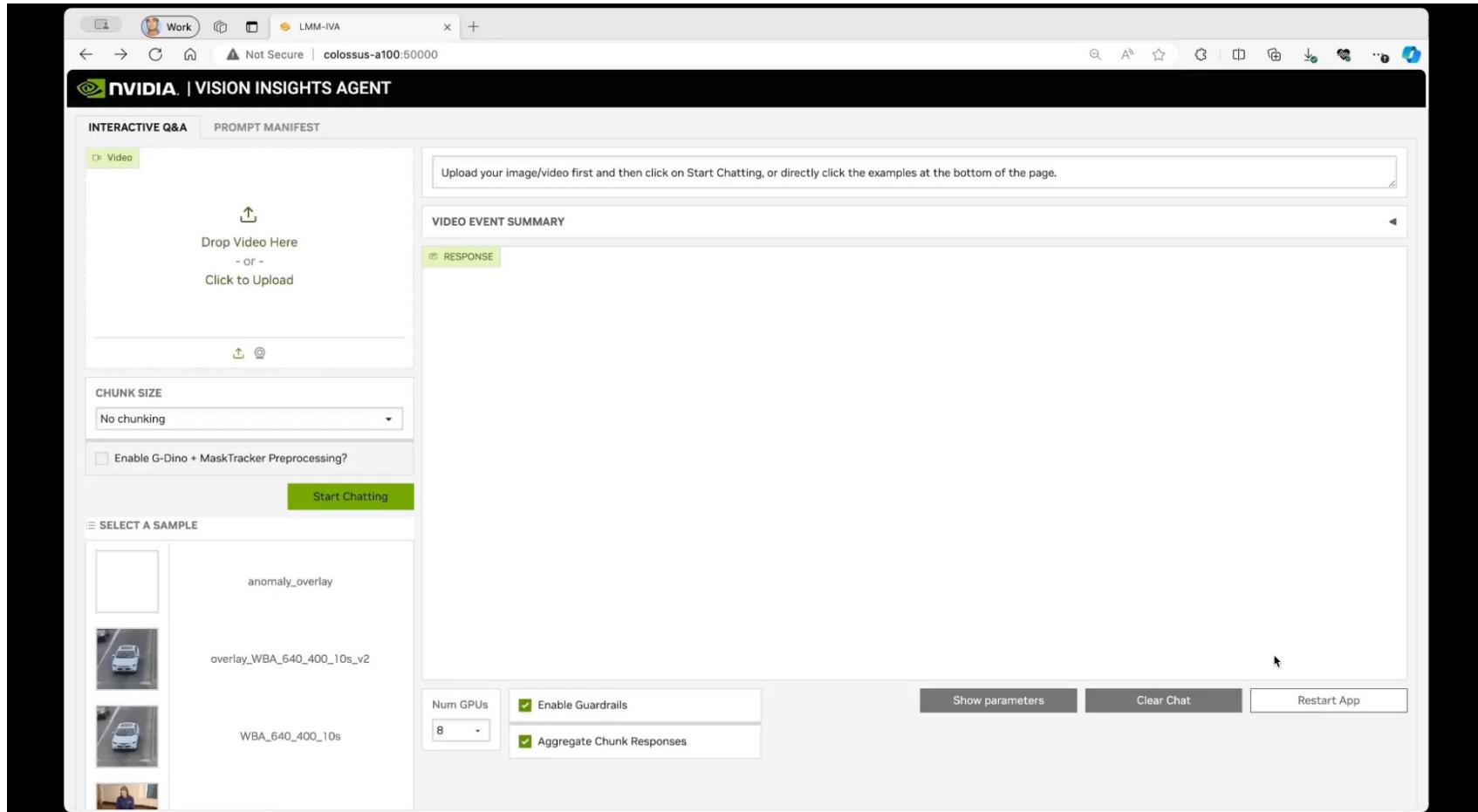
Improving Spatial and Contextual Awareness using SoM Prompting

AI framework for Image Insights - Grounding DiNO + Efficient SAM with ViLA

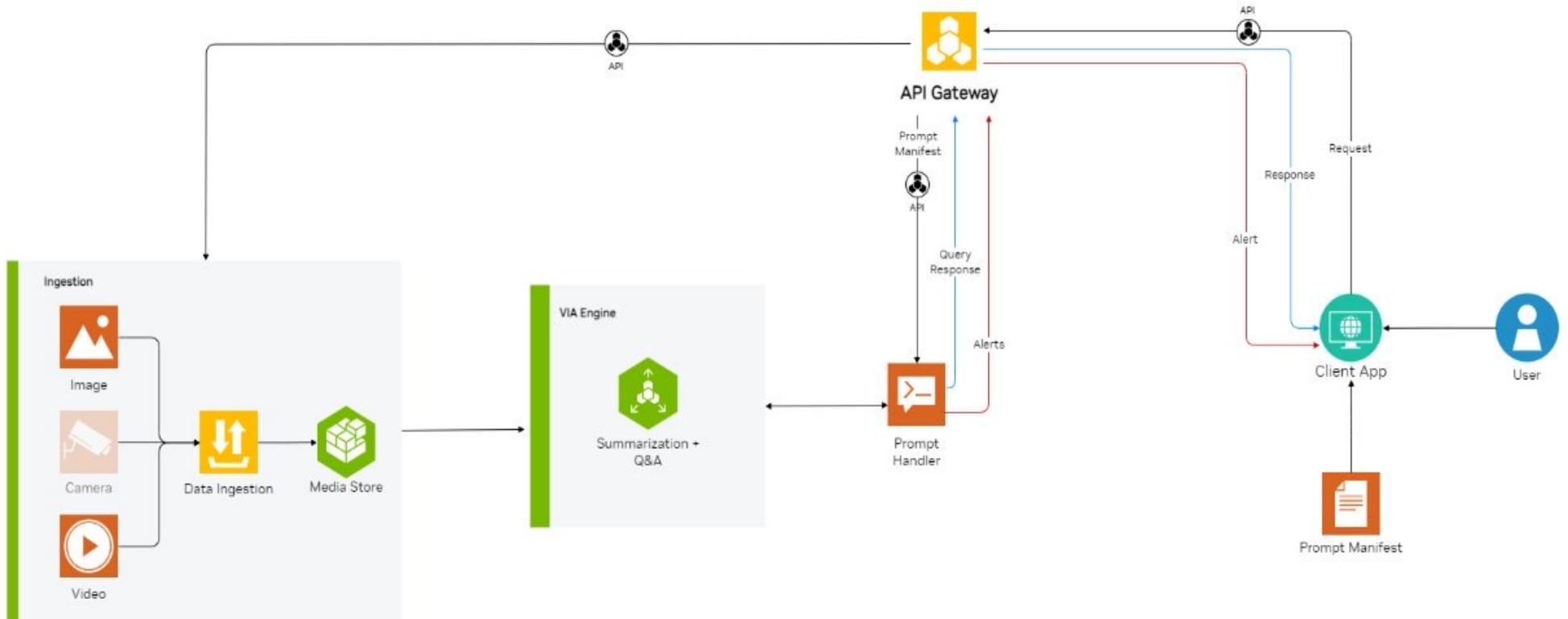


VIA Workflow with Spatial Grounding on Video

Video Demo



Autonomous Alerts Workflow Architecture



Autonomous Alerts Workflow Demo

The screenshot shows the NVIDIA Visual Insight Agent interface running in a browser window. The title bar indicates the page is "Not Secure" and the URL is "colossus-a100:50000". The main header reads "NVIDIA. | VISUAL INSIGHT AGENT". Below it, there are two tabs: "INTERACTIVE Q&A" and "PROMPT MANIFEST", with "PROMPT MANIFEST" being active.

The left sidebar contains two upload sections:

- Image:** A large area with a placeholder image of a drone and the text "Drop Image Here - OR - Click to Upload".
- Prompt Manifest:** A smaller area with a placeholder image and the text "Drop File Here - OR - Click to Upload".

Below these sections is a green button labeled "Upload & Start Chat".

The right side of the interface is divided into two main sections:

- Video Event Summary:** A large white area where the application displays its findings.
- RESPONSE:** A smaller white area below the summary.

At the bottom of the interface, there are several controls:

- "SELECT A SAMPLE" section with a table:

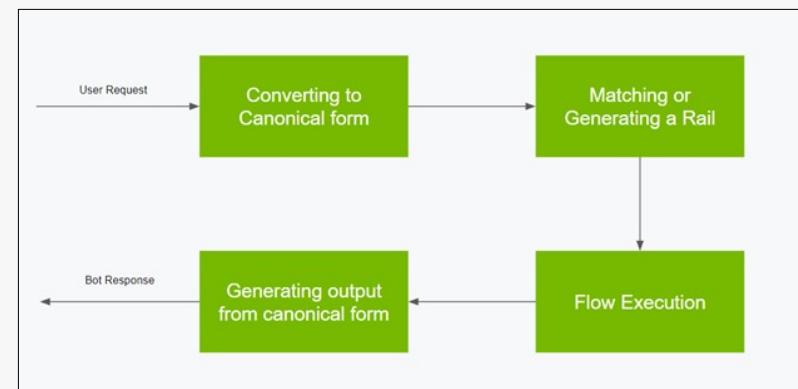
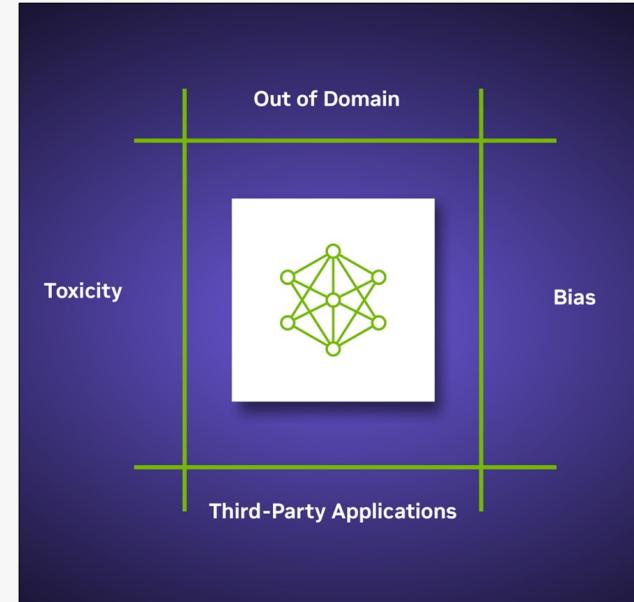
Image	drone2
	drone2
- "Num GPUs" dropdown set to 8.
- "Enable Guardrails" checkbox checked.
- "Show parameters" button.
- "Clear Chat" button.
- "Restart App" button.

At the very bottom, there are links: "Use via API" and "Built with Gradio".

Mitigating Hallucinations with Guardrails

Building trustworthy, safe, and secure LLM conversational systems

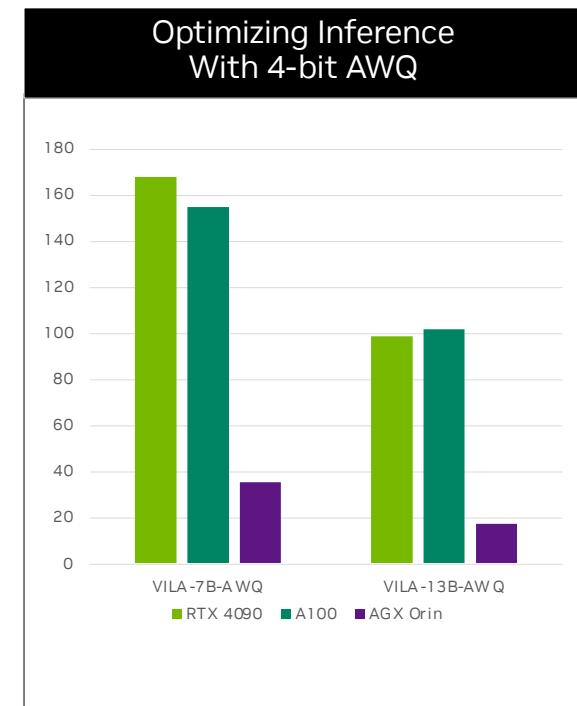
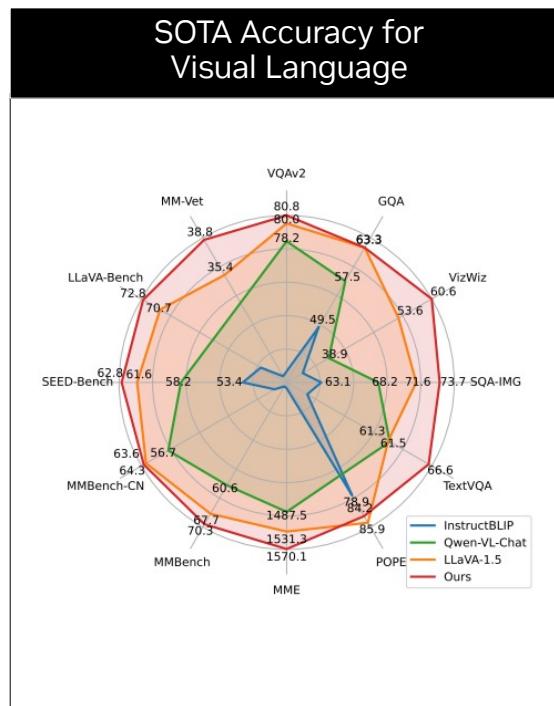
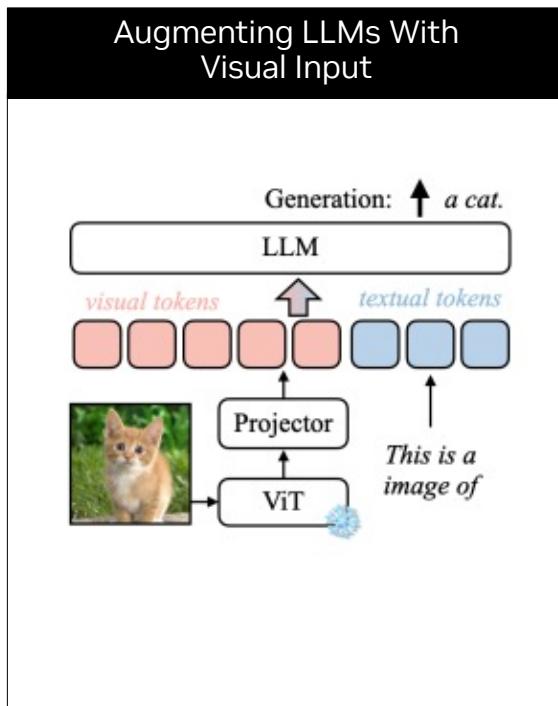
- When LLMs are not supplied with factual actual information, they hallucinate and provide faulty, but convincing responses.
- RAG reduces the likelihood of hallucinations by providing relevant facts
- NeMo Guardrails helps enterprises keep applications built on LLM aligned with their safety and security requirements.
- Today, NeMo Guardrails supports three broad categories of guardrails:
 - Topical
 - Safety
 - Security



Vision-Language Models

Vision-Language Model for Images

ViLA (Visual Language Assistant)



<https://arxiv.org/pdf/2312.07533.pdf>

ViLA Capabilities



Pred: Home to the greatest pizza.

prediction



Home to the best burgers and fried chicken.

Home to unbeatable fish and chips.

Home to outstanding ramen.

context

In-Context Learning



Q: Can the vehicle proceed through the traffic now?

A: Based on the image, the vehicle cannot proceed through the traffic yet. There are multiple people and bicycles in the crosswalk, and the traffic light is red. The vehicle must wait for the traffic light to turn green before proceeding.

Zero-shot Generalization and Complex Reasoning

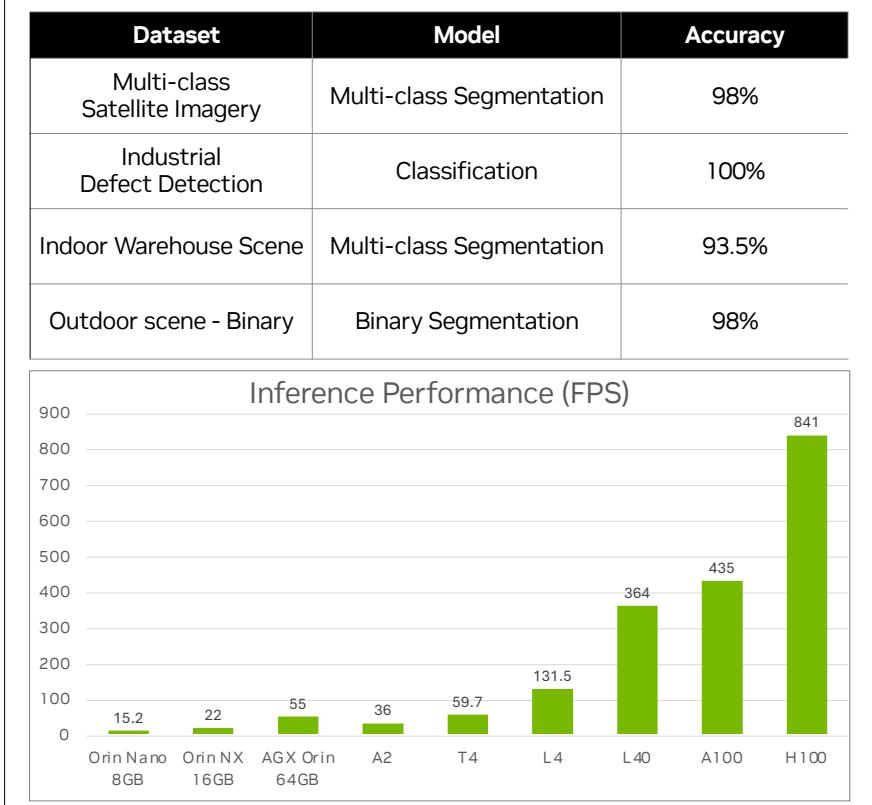
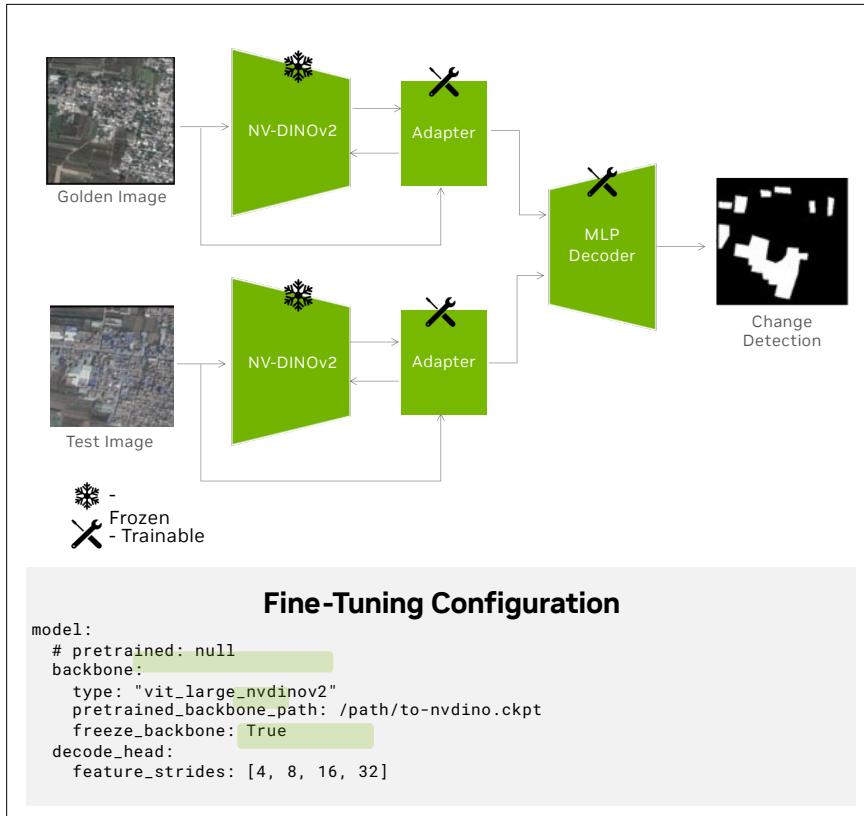


Prompt: \nI am getting just two custom pizza for me and my friend. How much should I pay in total according to the image + *Let's think step-by-step*.

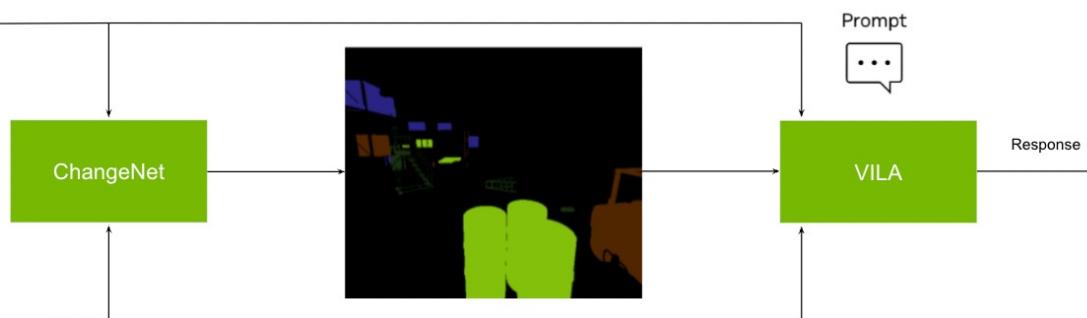
Answer: The image shows a chalkboard menu with a special for today, which is to create your own pizza. The price for this special is \$9.99. Since you are getting two custom pizzas, you should pay $\$9.99 \times 2 = \19.98 for the two pizzas. So, the answer is **\$19.98**.

Chain of Thought

ViLA for Change Detection



ViLA for Change Detection

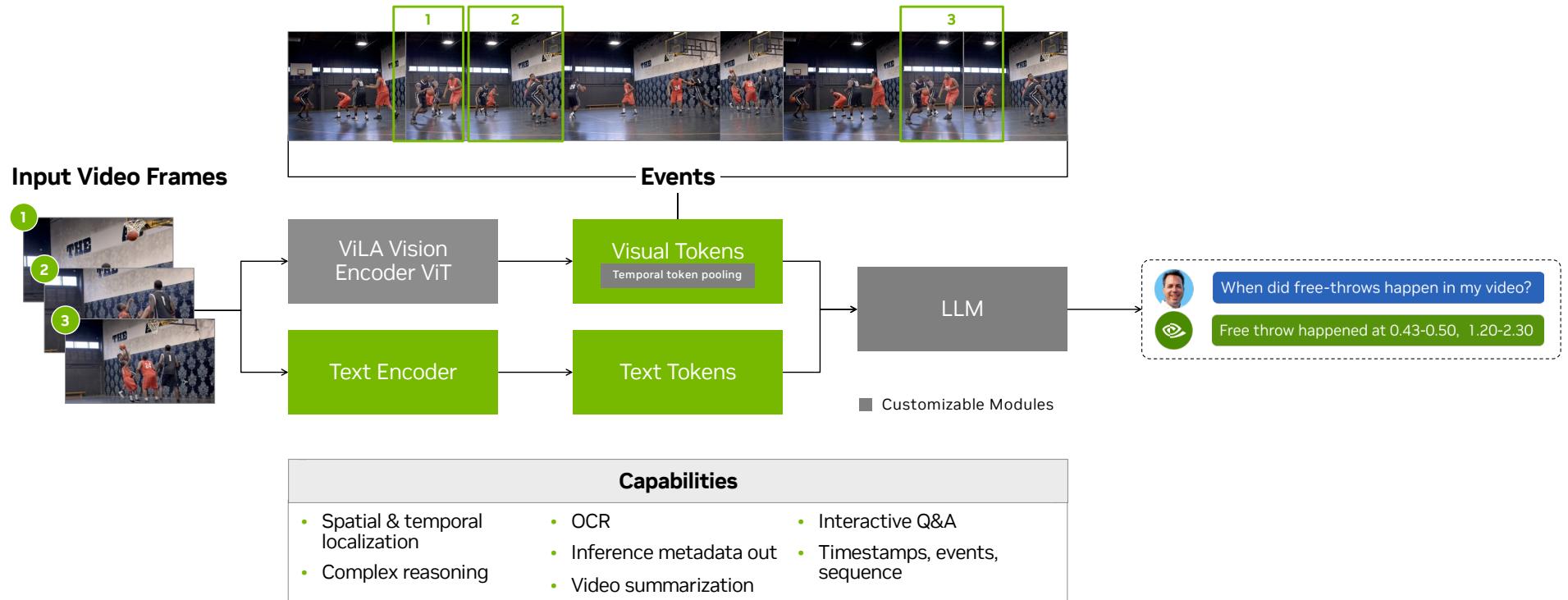


The warehouse has undergone significant changes.

The most noticeable transformation is the addition of a large number of barrels.

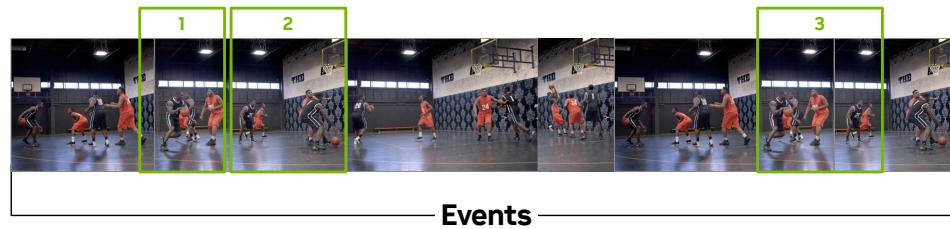
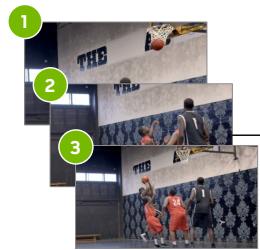
Vision-Language Models for Videos

LiTA (Language-instructed Temporal Assistant)



LiTA Capabilities

Input Video Frames



Events

NVIDIA | VISUAL INSIGHT AGENT

INTERACTIVE Q&A **PROMPT MANIFEST**

Type and press Enter

VIDEO EVENT SUMMARY

RESPONSE

- Video processed. You can now start chatting.
- The environment is an indoor basketball court.
- No, the video does not show any spectators.
- The players are wearing black and red jerseys.

CHUNK SIZE
No chunking

SELECT A SAMPLE

- warehouse_concatenated
- pole Vault
- warehouse_boxes_falling
- robo_box_falling_forklift

Start Chatting

Num GPUs: 1 Enable Guardrails Aggregate Chunk Responses

Show parameters **Clear Chat** **Restart App**

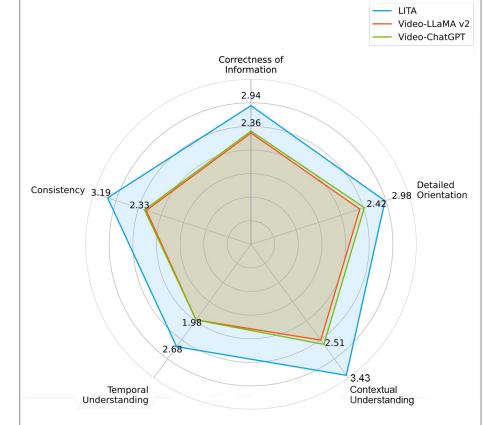
GPU Utilization

0: NVIDIA A100 80GB PCIe - 0 %
1: NVIDIA A100 80GB PCIe - 0 %

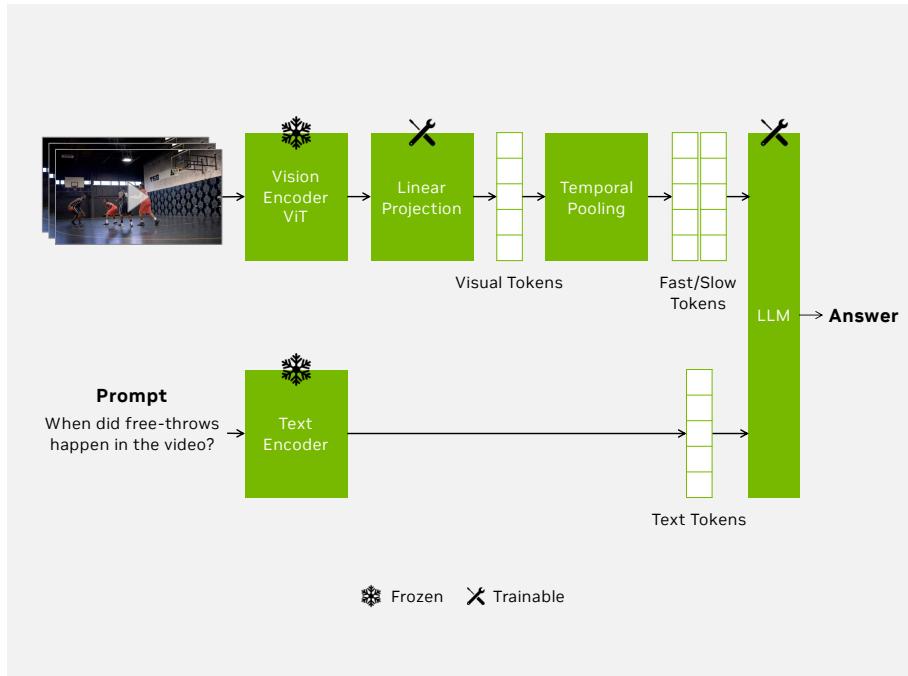
LiTA TRT-LLM Perf with TRT-LLM



LiTA SOTA Vision-LLM Benchmark

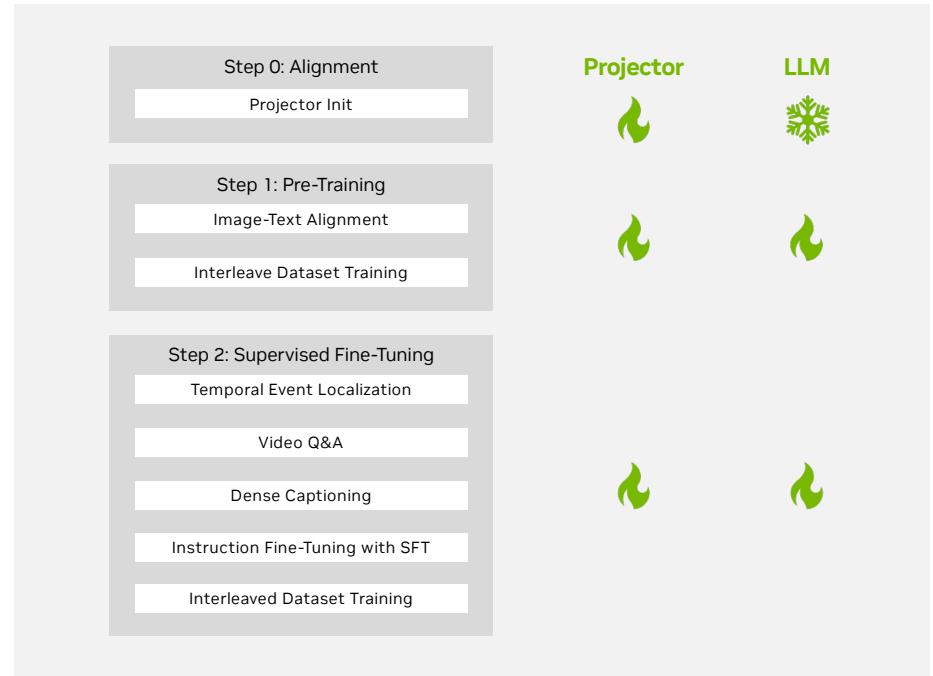


Fine-Tuning LITA for Custom Use Case



Supervised Fine-Tuning with LoRA

Model Size: 13.3B parameters



Training Workflow

Trainable Parameters: 50M (with LoRA)

*LoRA: Low-Rank Adaptation of LLMs <https://arxiv.org/pdf/2106.09685.pdf>

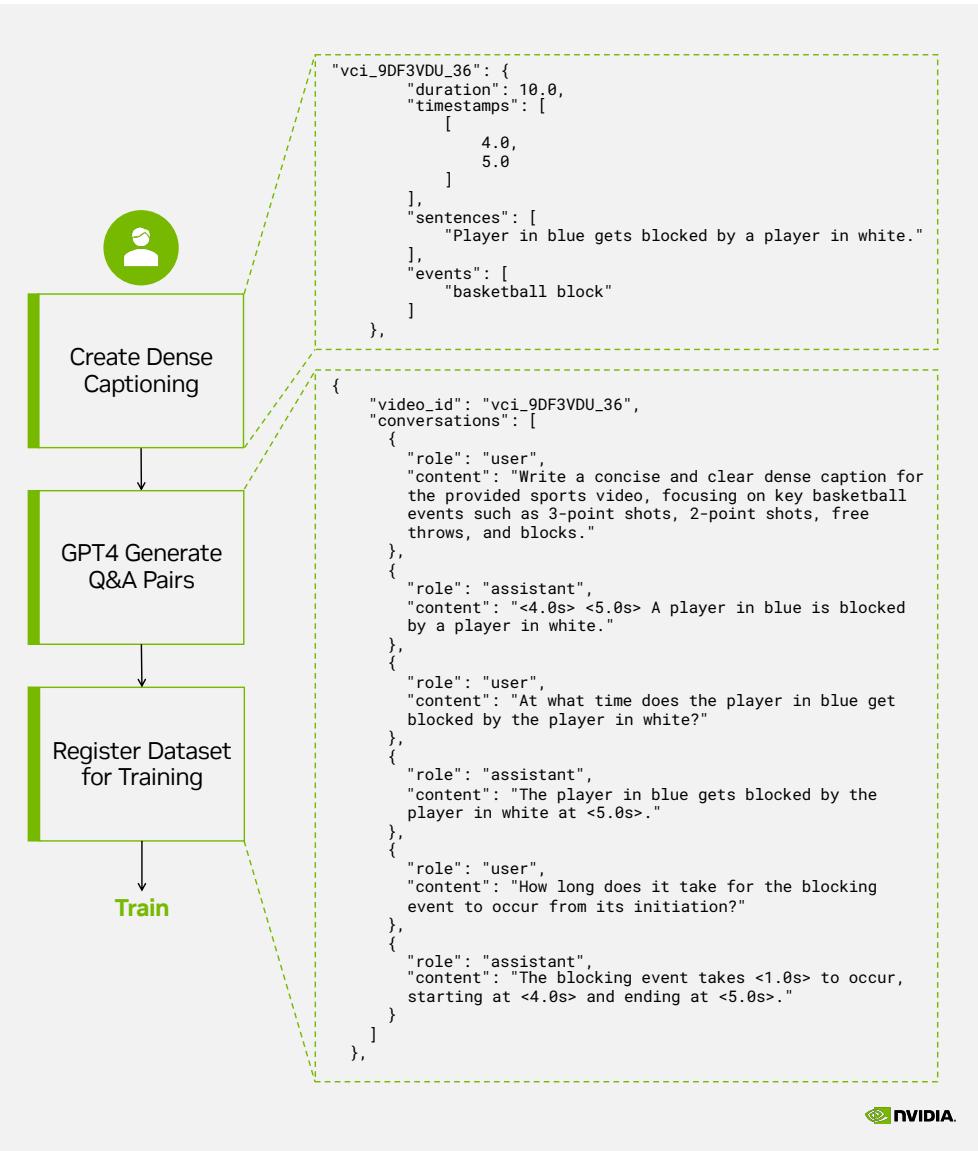
Steps to Customize LiTA

Dataset Recommendation

- 1000 videos for good fine-tuning
- 1 – 2 min video clips
- Events between 1% – 10% of clip duration

Training KPIs

- Sports (Basketball) action dataset
- Trained ~1K videos in 16 hours on 8xA100



After Fine-Tuning LiTA

Example: basketball use case



Question: What type of shot is this?

Original model: This is a basketball shot.

Fine-tuned model: This is a 2-pointer shot from the short corner.

Challenges with finetuning VLM for temporal localization

- Dynamic input number of frames
- Long context video
- Improving Spatial Reasoning
- Model resolution for OCR
- Improving temporal localization
- Adding Audio modality
- Fine-grained temporal localization data
- Inference performance

Apply for Early Access

<https://developer.nvidia.com/visual-insight-agent-early-access>

VIA workflows will be offered:

- To ecosystem ISVs, SDPs and enterprises partners in various verticals
- With workflows illustrating the ability to customize its microservices and build AI agents.

VIA will:

- Have AI workflows that include microservices to accelerate building next wave of vision AI agents
- Enable workflows for video summarization, interactive Q&A and automated alerts by leveraging VLMs

Leverage VIA AI workflows:

- Finetune and customize large multi-modal models
- Vision AI foundation models
- GPU training performance
- Acceleration in ingestion and streaming of multimedia
- Cloud-native deployment at enterprise scale



