

The logo for the 2024 Embedded VISION Summit is centered on the left side of the slide. It features a white octagonal background with a colorful, multi-layered border in shades of purple, blue, green, yellow, and orange. The text "2024" is at the top, "embedded" is below it, "VISION" is in large, bold, dark blue letters with a gradient, and "SUMMIT" is at the bottom in a smaller, dark blue font.

2024
embedded
VISION
SUMMIT®

Optimized Vision Language Models for Intelligent Transportation System Applications

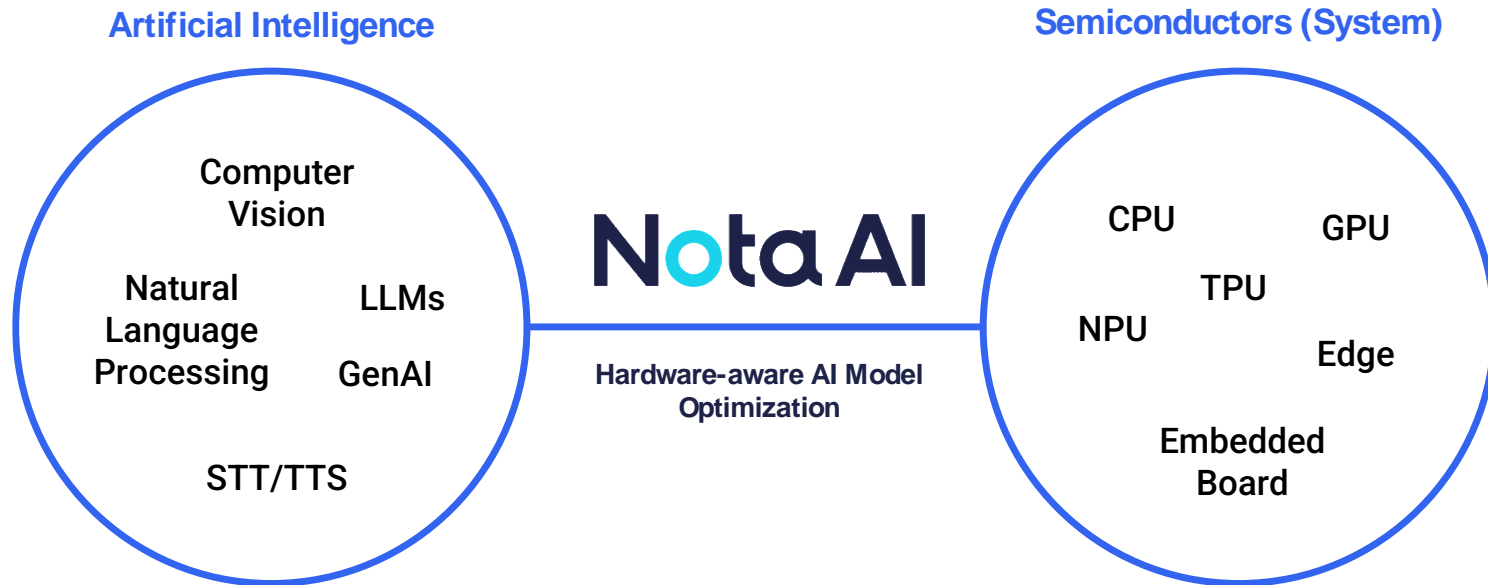
Tae-Ho Kim

CTO & Co-Founder

Nota Inc.

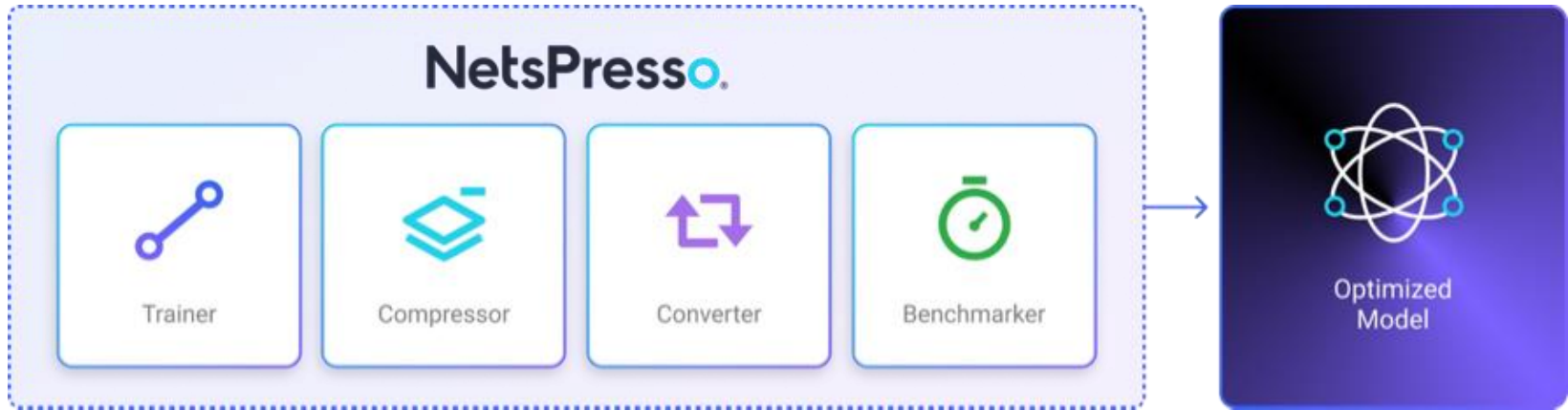
NotaAI

- In this talk, we explore the **challenges in ITS**.
- How **vision language model (VLM)** can solve these challenges.
- **Future work** will also be addressed.



Nota AI bridges the gap between AI & semiconductors.

Nota AI's Main Product: NetsPresso



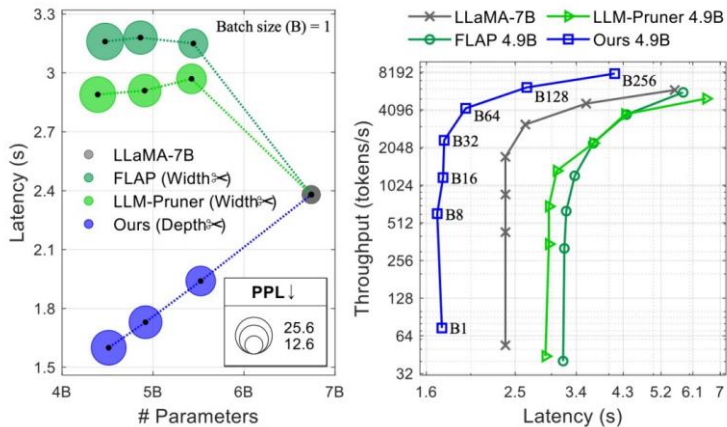
NetsPresso® simplifies AI model optimization for target devices with automated processes.

Platform & Edge Solutions: Elevating Excellence

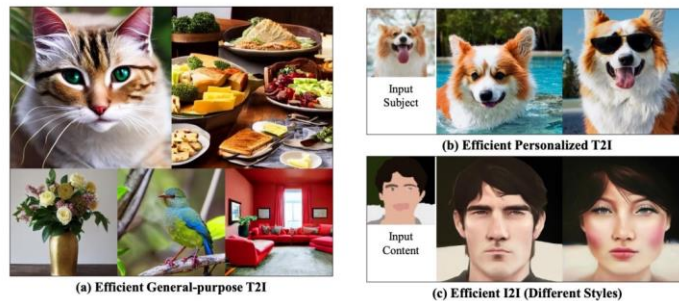


Using NetsPresso, Nota AI also has created a solution business in ITS

LLM Optimization



Stable Diffusion Optimization



Device	SDM-v1	Base [Ours]	Small [Ours]	Tiny [Ours]
AGX Orin 32GB	4.9s [53]	3.4s (-31%)	3.2s (-35%)	2.8s (-43%)
iPhone 14	5.6s [52]	4.0s (-29%)	3.9s (-29%)	3.9s (-29%)

Nota AI also specializes in GenAI compression

Intelligent Transportation System (ITS)



Smart Intersection System

Real-time incident detection and analysis for safe road condition



VRU Safety Solutions

AI-based real-time hazardous situation screening and control system



Automatic Incident Management System

Real-time incident detection and analysis for safe road condition



AI Smart Parking

Real-time parking occupancy and parking facility utilization analysis

Use Case 1: Smart Intersection System

교차로 영상 AI 분석



The screenshot displays a real-time video feed of a city intersection. The video is overlaid with green bounding boxes and labels for various vehicles and pedestrians, indicating AI detection. On the right side of the interface, there are two data tables. The top table, titled '차량 카운팅 정보' (Vehicle Counting Information), shows a summary of vehicle counts by direction. The bottom table, titled '카운트 실시간 정보' (Count Real-time Information), provides a detailed log of individual vehicle counts, including time, direction, and ID.

차선	방향	총 카운트
1	좌회전	89
1	유턴	10
2	직진	82
2	좌회전	23
3	직진	150
4	직진	93
4	우회전	12

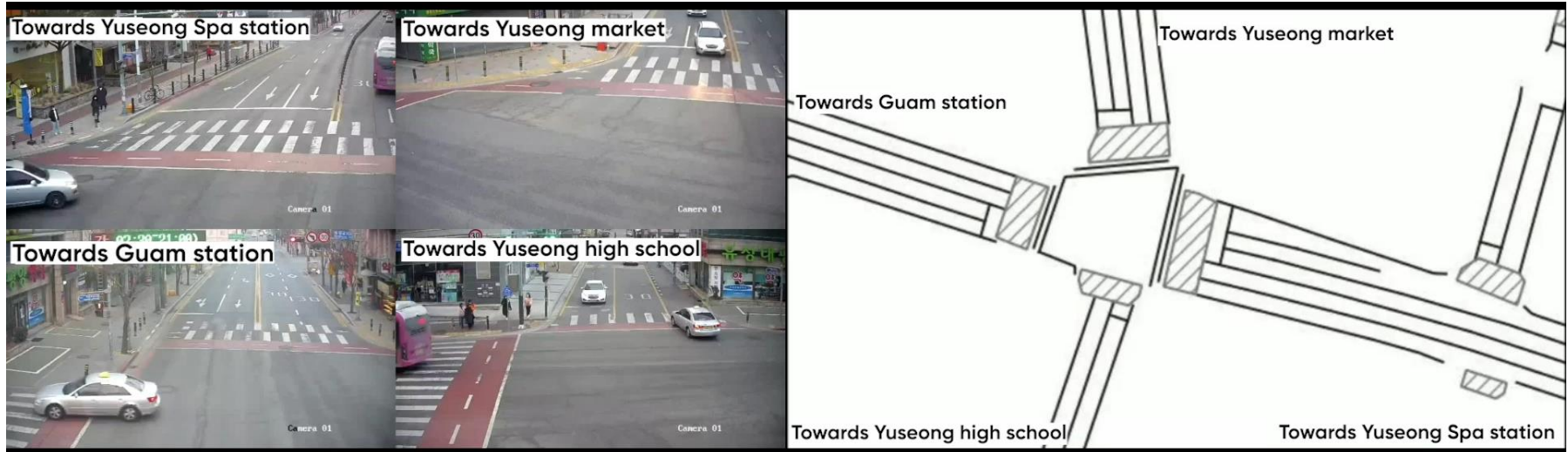
차종	시간	차선	방향	ID
🚗	17:00:50.351	3	직진	1629
🚗	17:00:49.518	4	직진	1586
🚗	17:00:47.887	3	직진	1623
🚗	17:00:47.319	3	직진	1581
🚗	17:00:45.953	3	직진	1572
🚗	17:00:45.367	4	직진	1570
🚗	17:00:43.421	4	직진	1583
🚗	17:00:42.889	3	직진	1576



- Daejeon metropolitan city ITS construction project
- Intersection CCTV AI video analysis (600 ch)
- 98% accuracy in traffic volume counting in night and rainy conditions

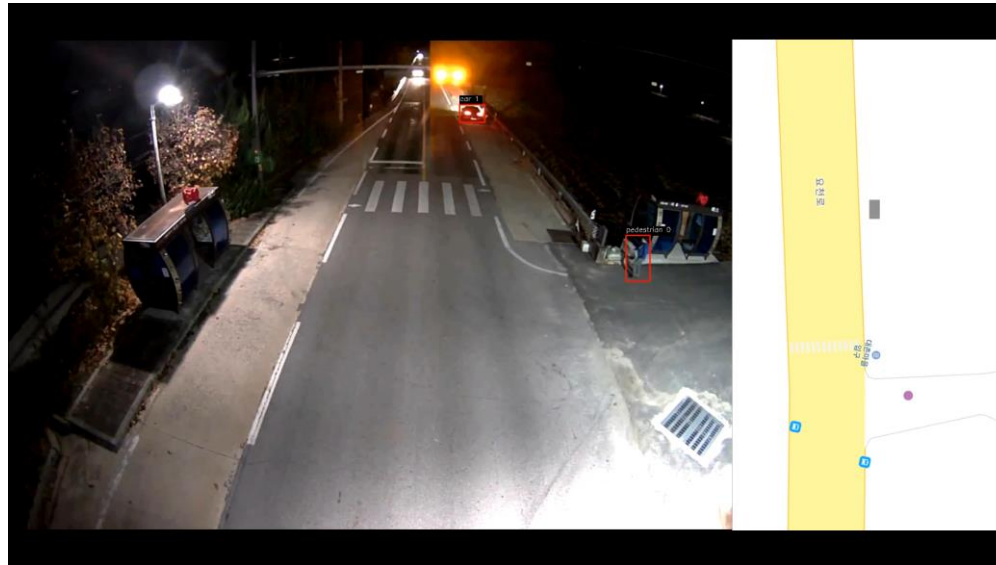
Use Case 2: AI Safe Crossing

Daejeon City | Implementing LDM using vehicle trajectory prediction information

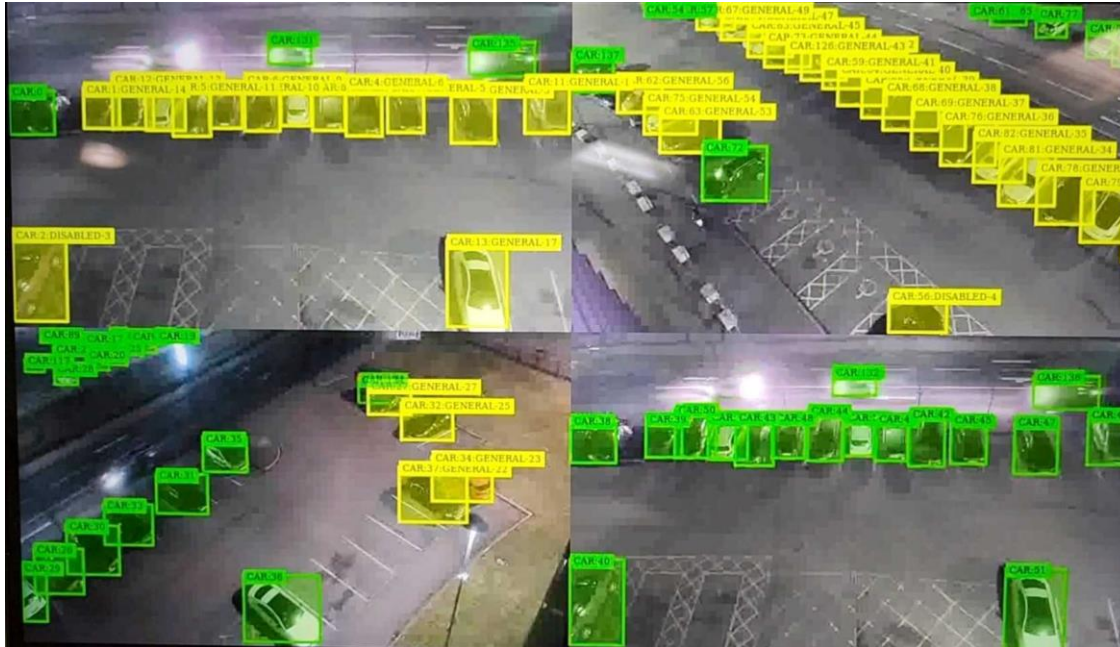


Use Case 3: VRU Safety Solutions

National Highways | C-ITS: AI video analysis for smart crosswalk safety transmitted through V2X communication



Use Case 4: AI Smart Parking



The City of
SAN DIEGO

- USA San Diego outdoor parking management (Caltrans)

MK Milton Keynes
City Council

- UK Milton Keynes stadium outdoor parking management

- Ill-posed problem: How can we define “road debris”?
- Contextual problem: How can we define “accidents” with object detection?
- Rare dataset: How can we obtain dataset?

Ill-posed Problem: Road Debris



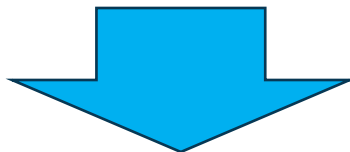
Can it be detected by legacy AI models?

Contextual Problem: Accident Detection



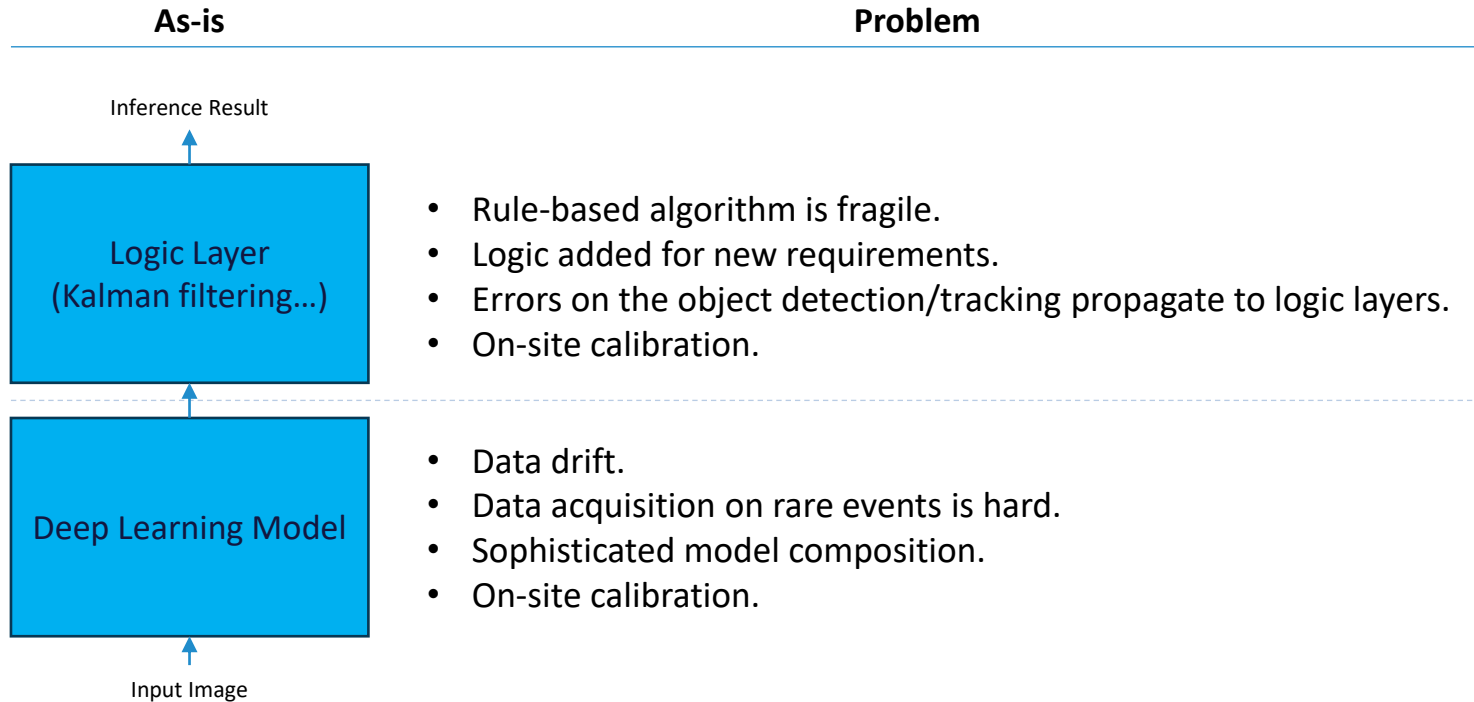
Can it be detected by legacy AI models?

- Ill-posed problem: How can we define “Road debris”?
- Contextual problem: How can we define “Accidents” with Object detection?
- Rare dataset: How can we obtain dataset?



Requires super-generalized model: Foundation Model

Foundation Model on the Edge: Challenges in Industrial AI



To-be

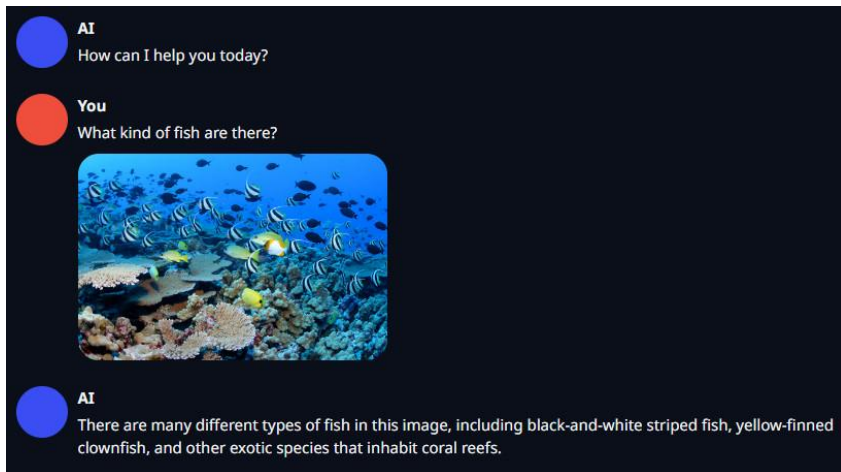
Features



- VLM is capable of comprehending complex scene.
- VLM already contains various logic.
- VLM is robust on data drift.
- VLM is aware of rare events.
- VLM needs less or no calibration.
- Still, VLM is not understanding video.

Vision Language Model: LLaVA

LLaVA



Live LLaVA



Source: jetson-ai-lab.com

Working Prototype of LLaVA on Accidents Detection



Benchmark on Models

Llava-13B (Jetson AGX Orin)	Quantization	Tokens/sec	Memory
<u><code>text-generation-webui</code></u>	4-bit (GPTQ)	2.3	9.7 GB
<u><code>llava.serve.cli</code></u>	FP16 (None)	4.2	27.7 GB
<u><code>llama.cpp</code></u>	4-bit (Q4_K)	10.1	9.2 GB
<u><code>local_llm</code></u>	4-bit (MLC)	21.1	8.7 GB

Source: jetson-ai-lab.com

- More advanced optimization required
- VLM needs to comprehend temporal consistency
- Domain adaptation might be required for user specific scenario
- Interface for product required
- Prompt engineering is required for higher performance

Conclusion

- Industrial AI is already a widely used technology, but technology is limited when the problem is complex and underdetermined.
- Among GenAI models, a Vision Language Model (VLM) can understand complex scenes, so it could analyze complex queries and events.
- For example, in ITS, road debris and car accidents are severe problems, but this couldn't be solved by legacy AI models.
- Using VLMs, these problems can be solved as well.
- However, VLMs are still computationally expensive, so lightweight VLMs are required for the next step.

Visit Nota AI @318 !



Thank you for your attention!