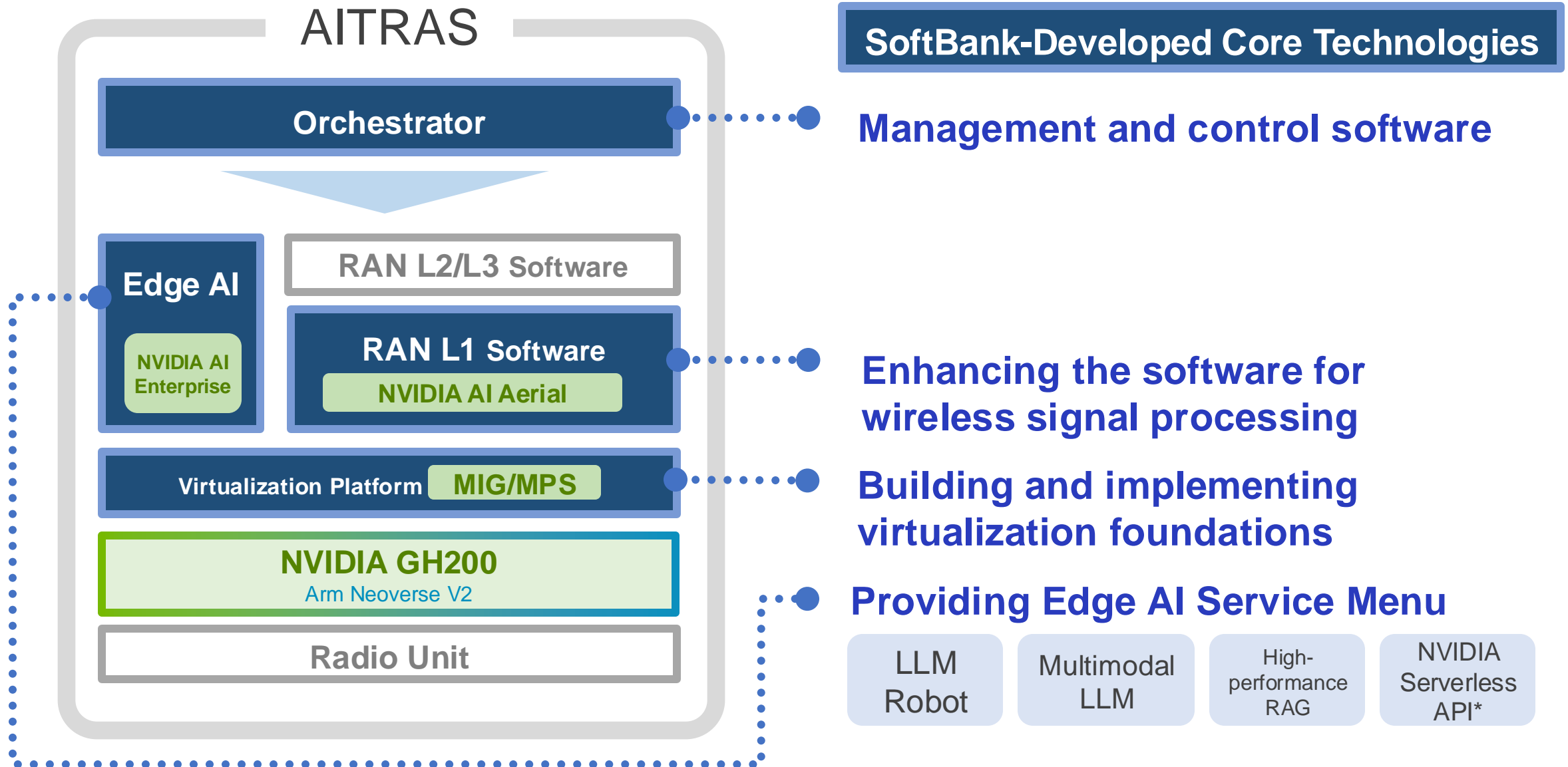
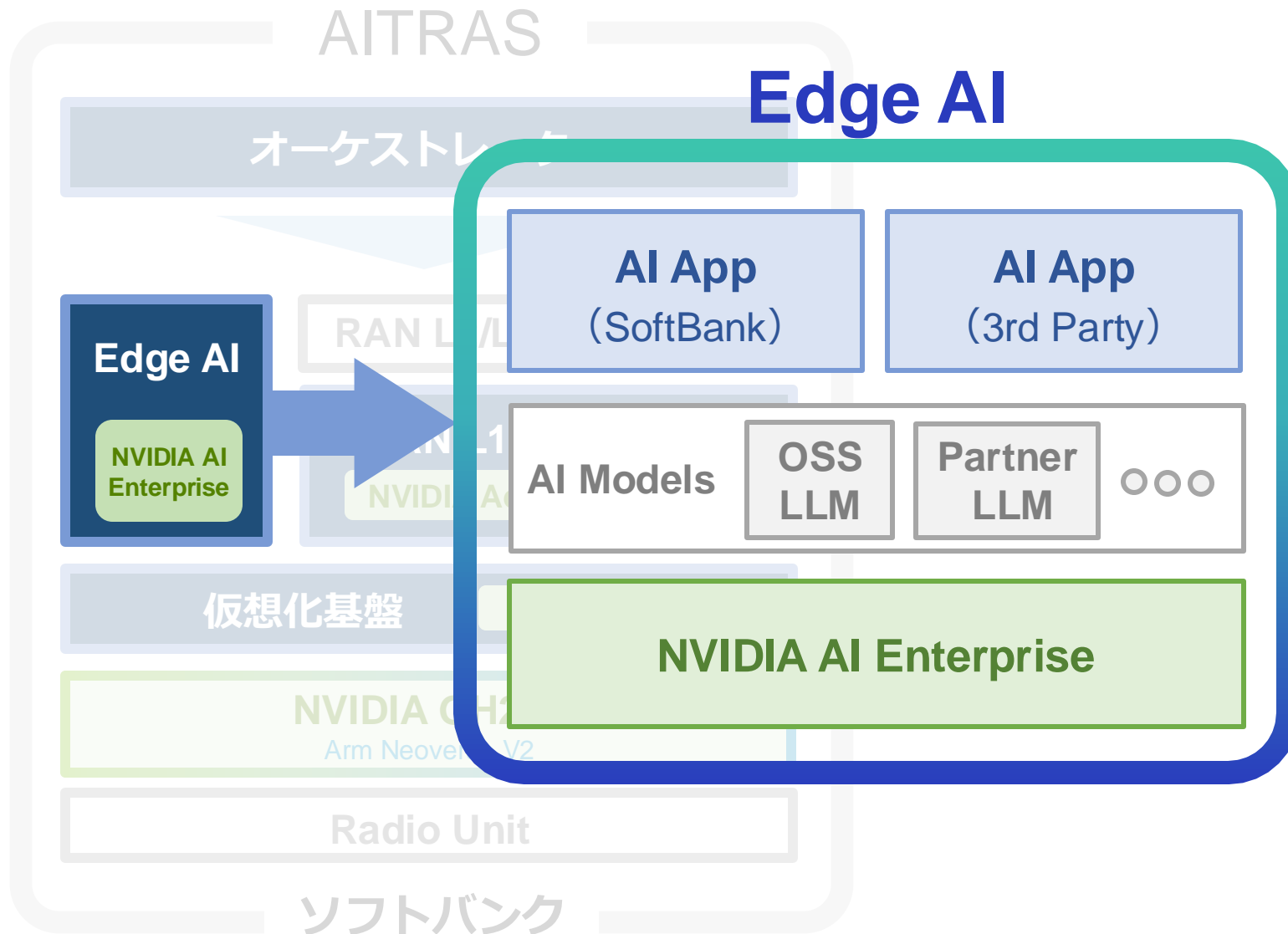


Cloud (LLM) Robot

AITRAS System Architecture



Positioning of this Demo



Original AI Service

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

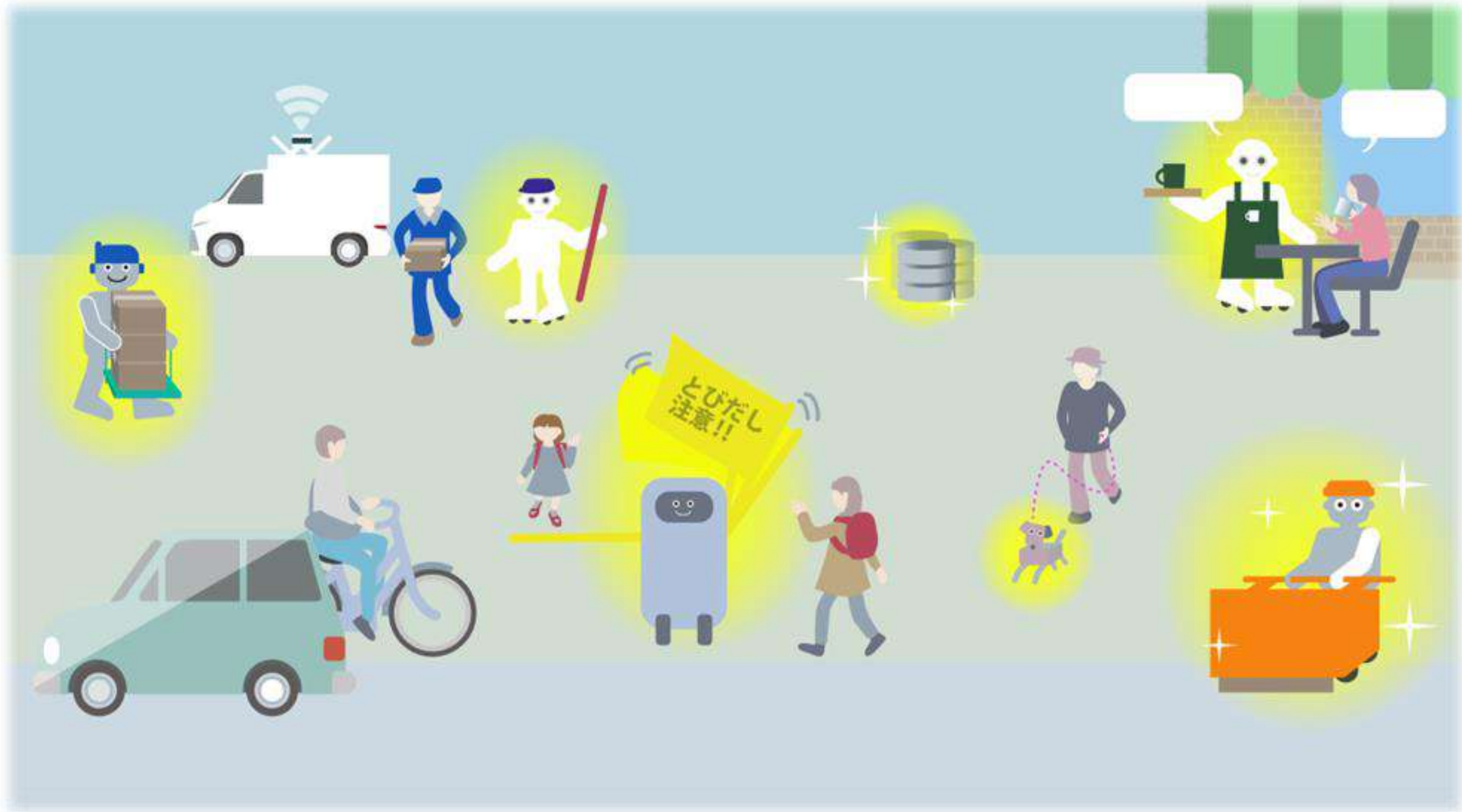
NVIDIA Serverless API

Allocates unused resources to meet customer demands on demand

NVIDIA AI Enterprise

Industry-standard AI framework

A Society Coexisting with Robots



Robots with advanced decision-making abilities are actively supporting daily life

Offloading Processing

Advanced decision-making requires high-performance computing.

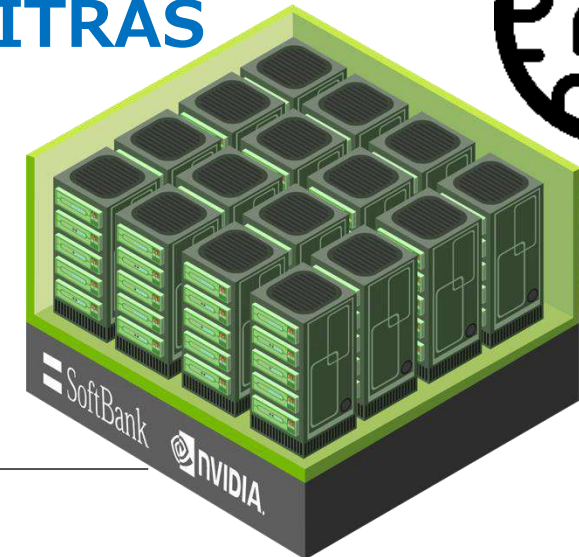
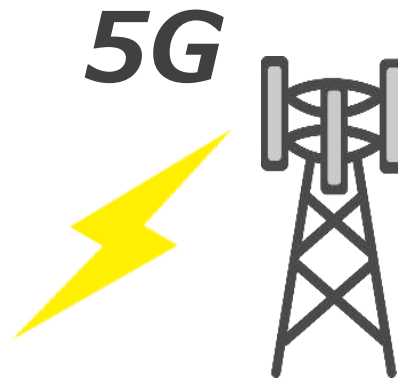
Computing Resources
Small



Computing Resources
Large

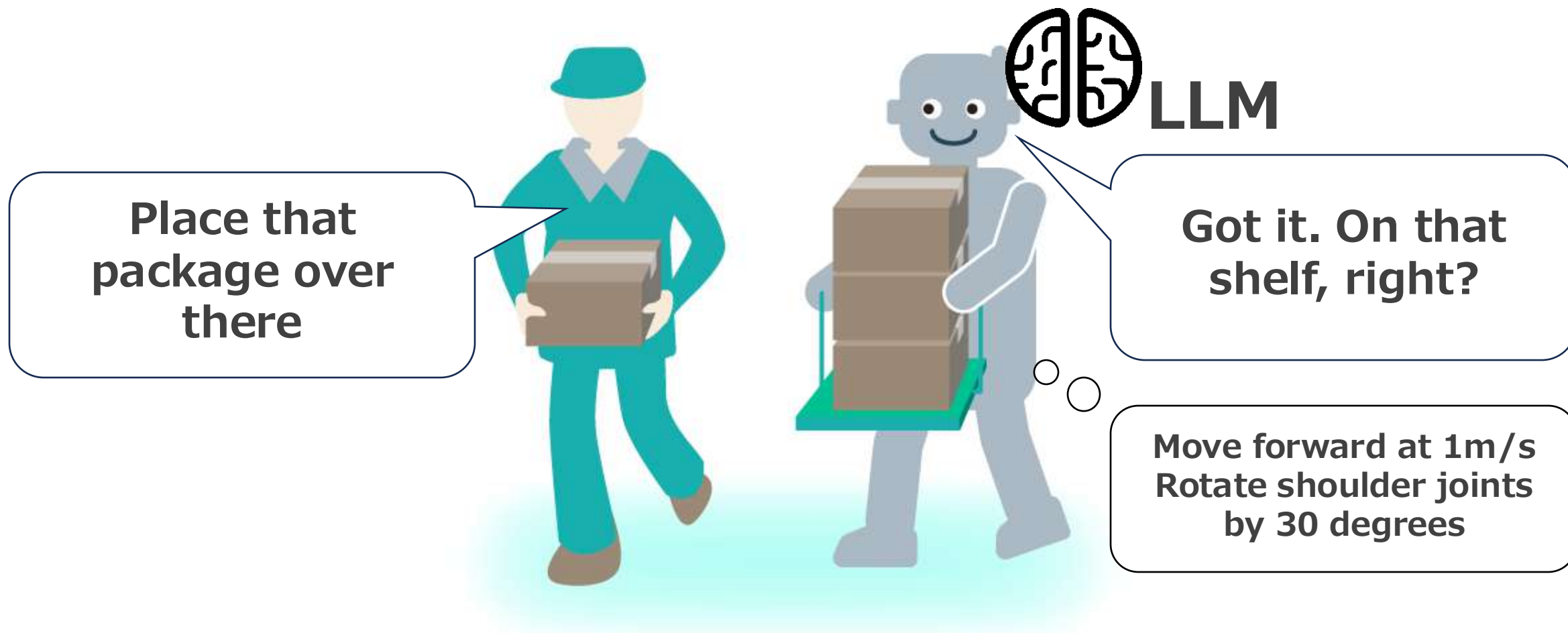


AITRAS



AITRAS can leverage more advanced computing power.

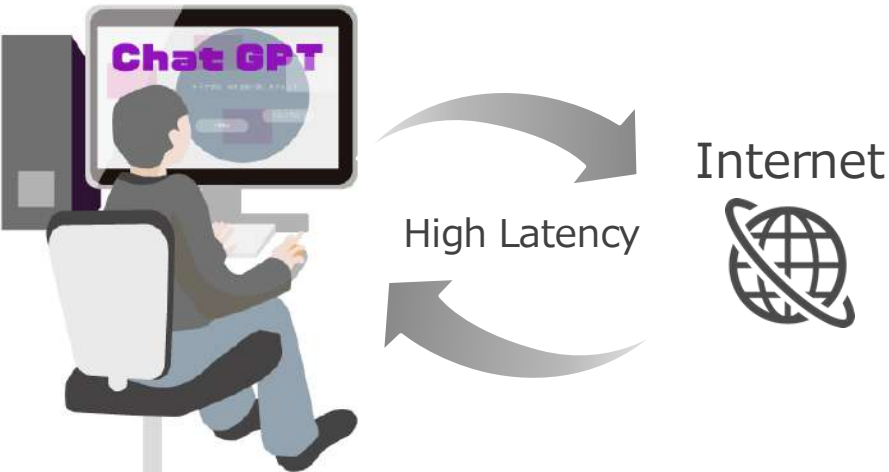
What is an LLM Robot?



Capable of responding to new situations with common sense.

High-speed processing is essential for controlling robots

Latency of General LLMs
Time from Question to Answer



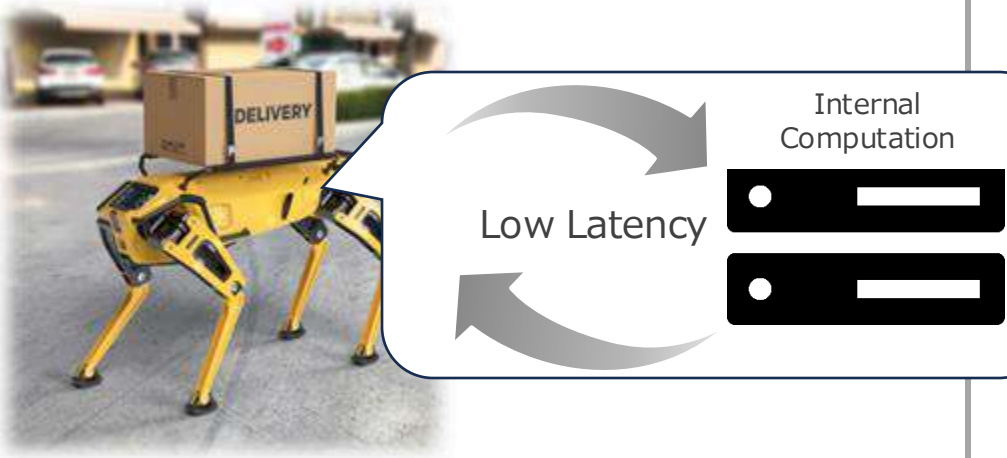
High Latency

Internet

More than 1 second



Robot Control Latency
Output of Control Information

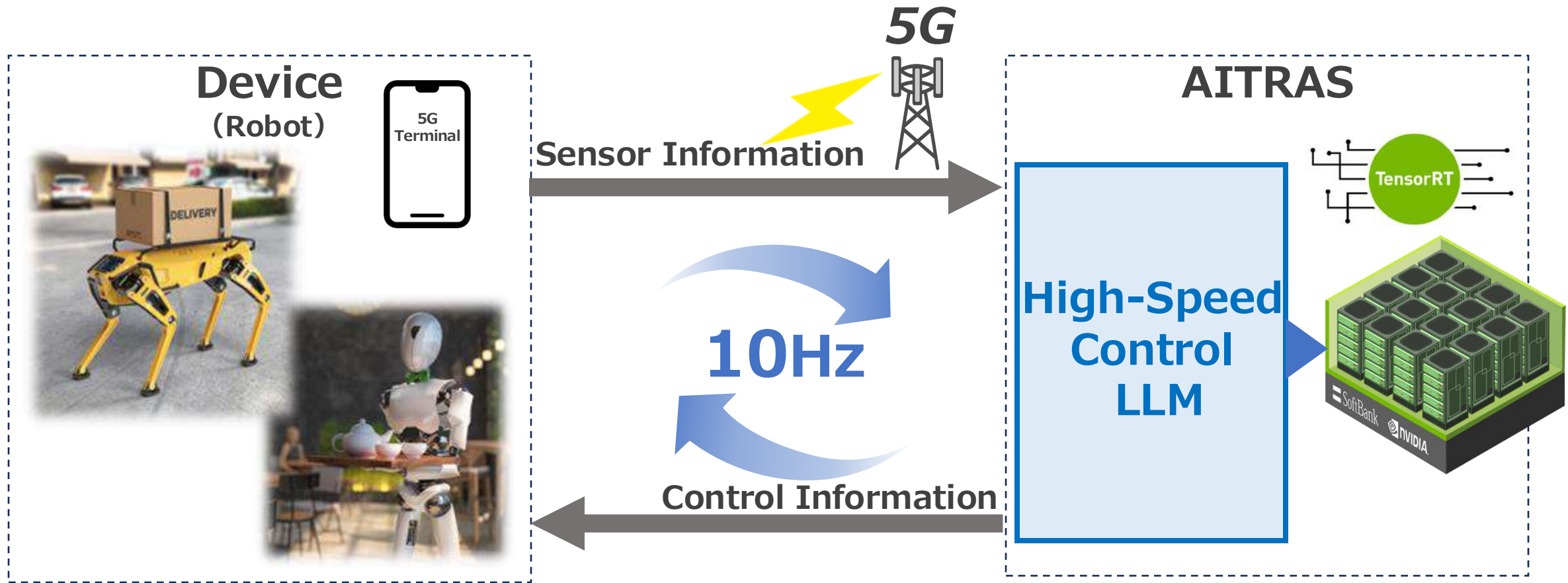


Low Latency

Internal Computation

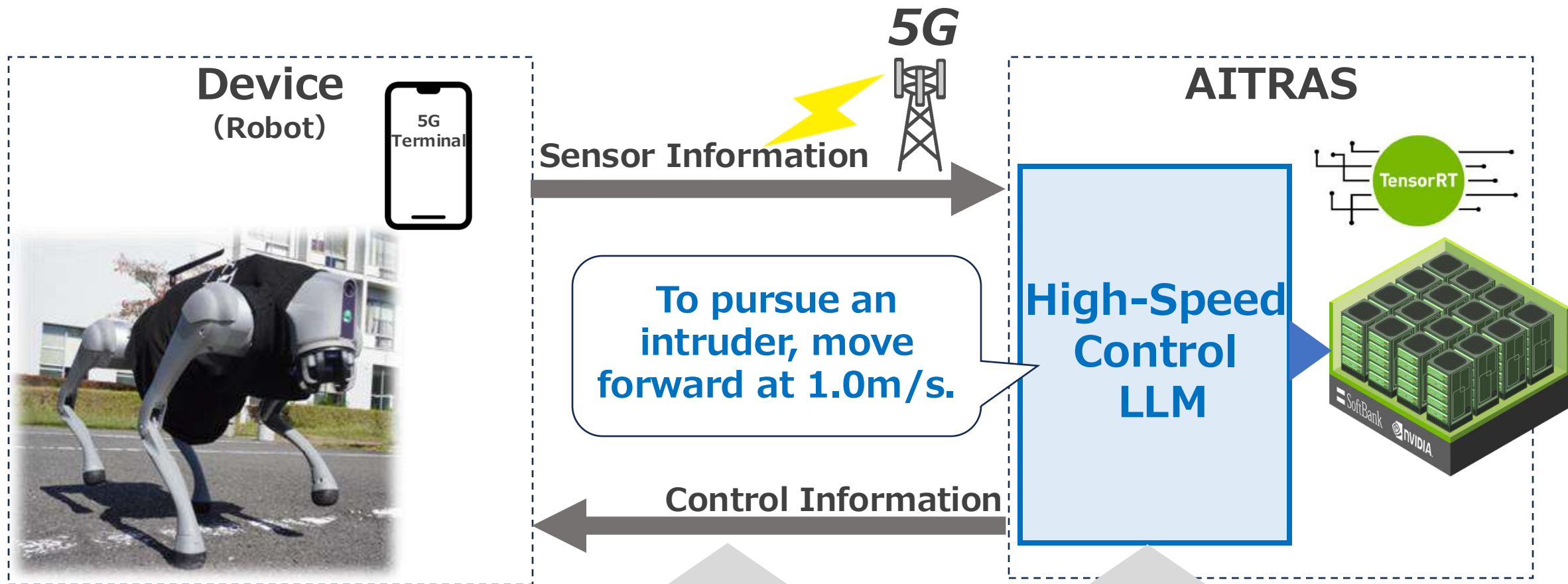
Around 0.1 seconds

Low-Latency LLM Robots in Action

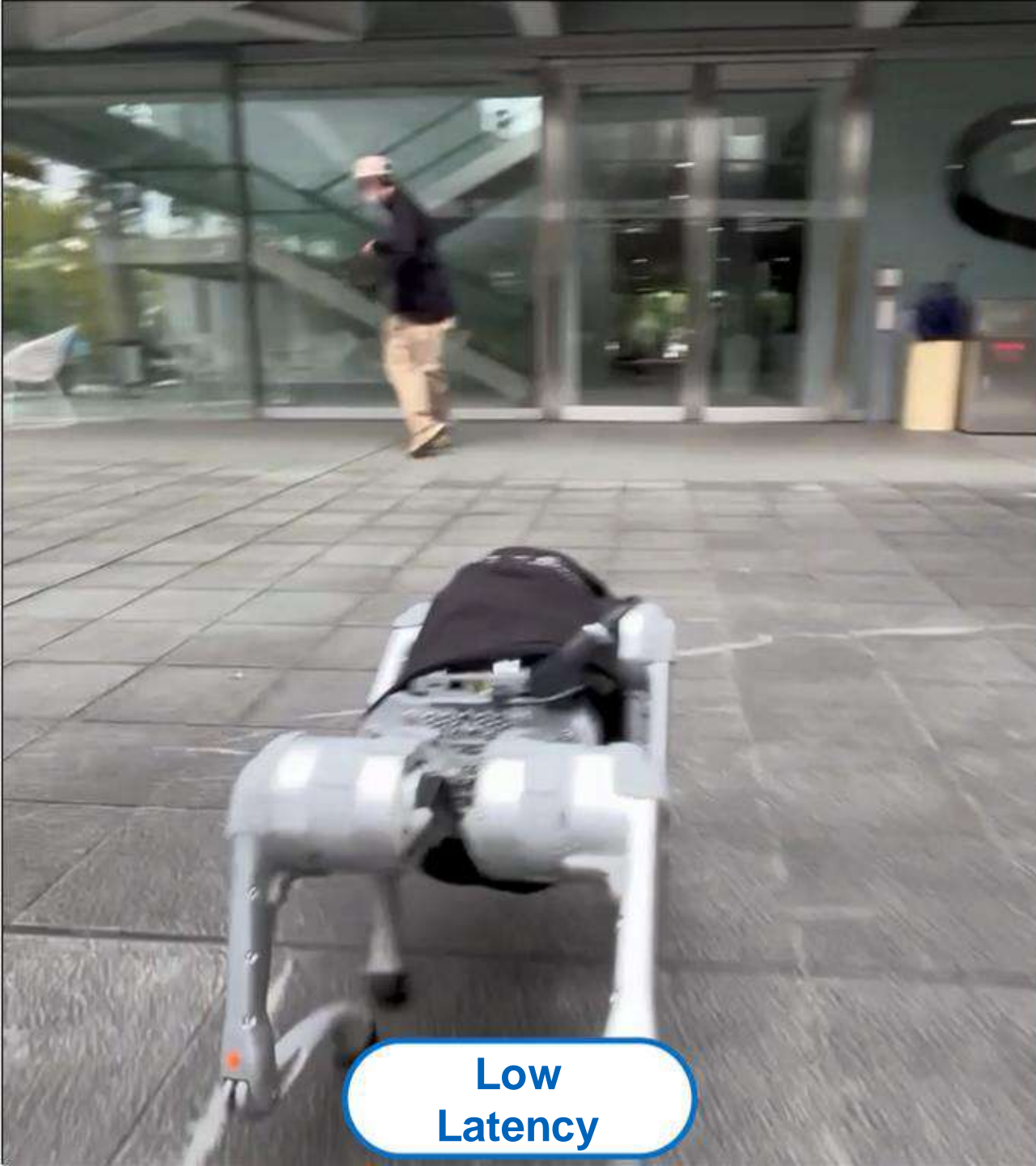


Robots Controlled 10 Times Per Second with LLM

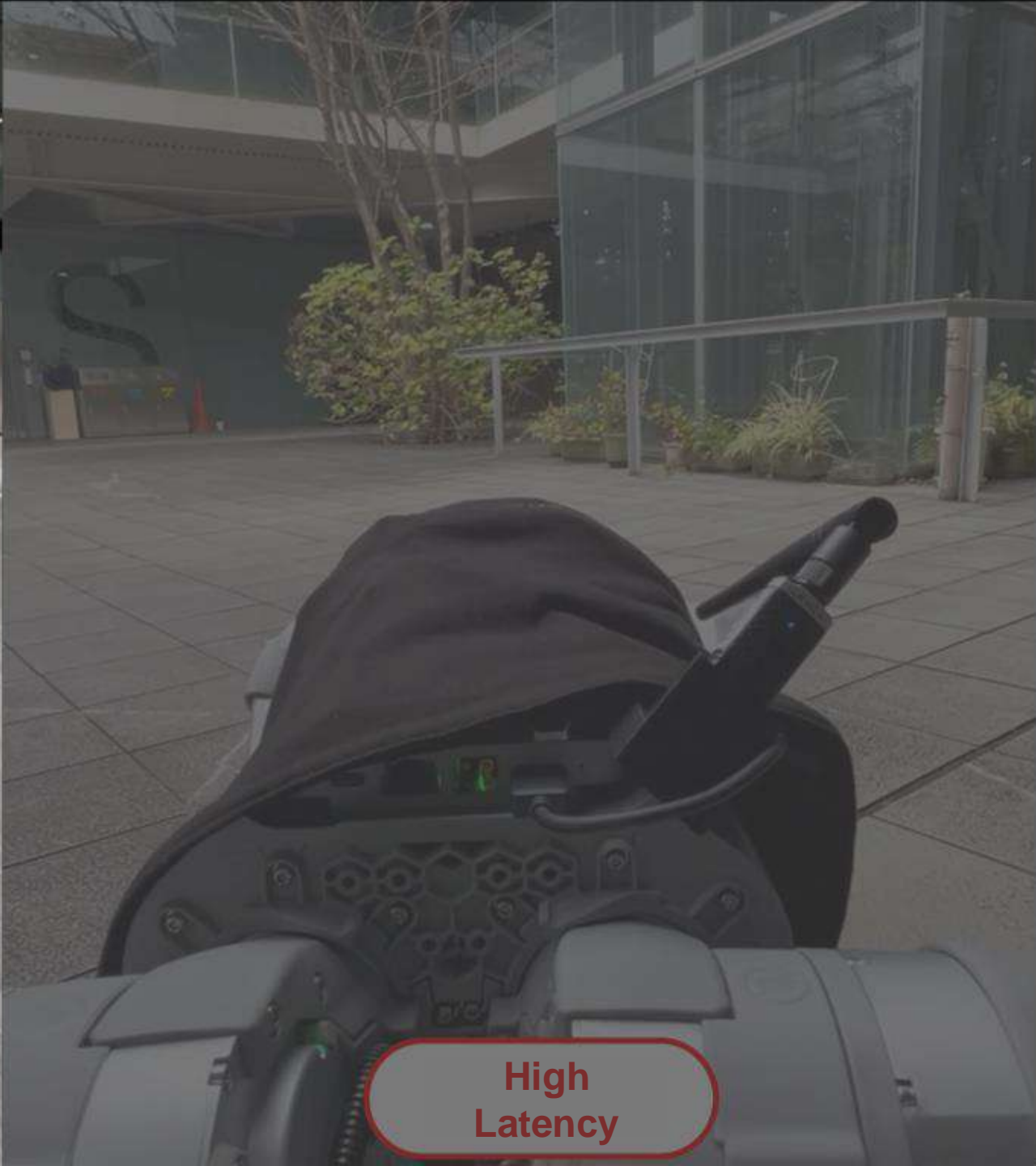
Low-Latency LLM Robots in Action



Total Latency: **Appx. 100ms** = **RAN Latency: Appx. 40ms** + **LLM Latency: Appx. 60ms**



**Low
Latency**



**High
Latency**