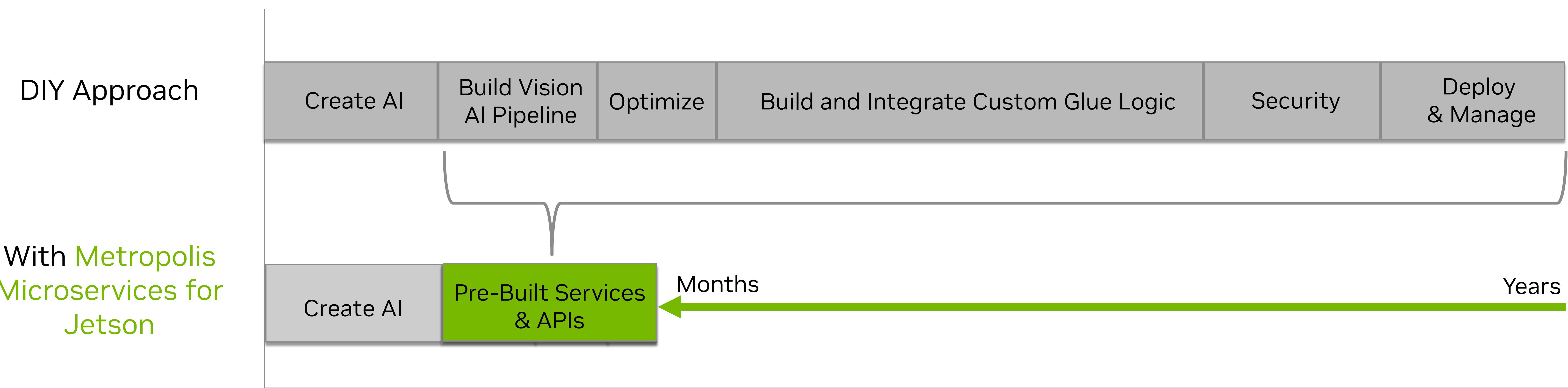




# Metropolis Microservices for Jetson

# Slashing the Cost of Bringing Edge AI Apps to Production



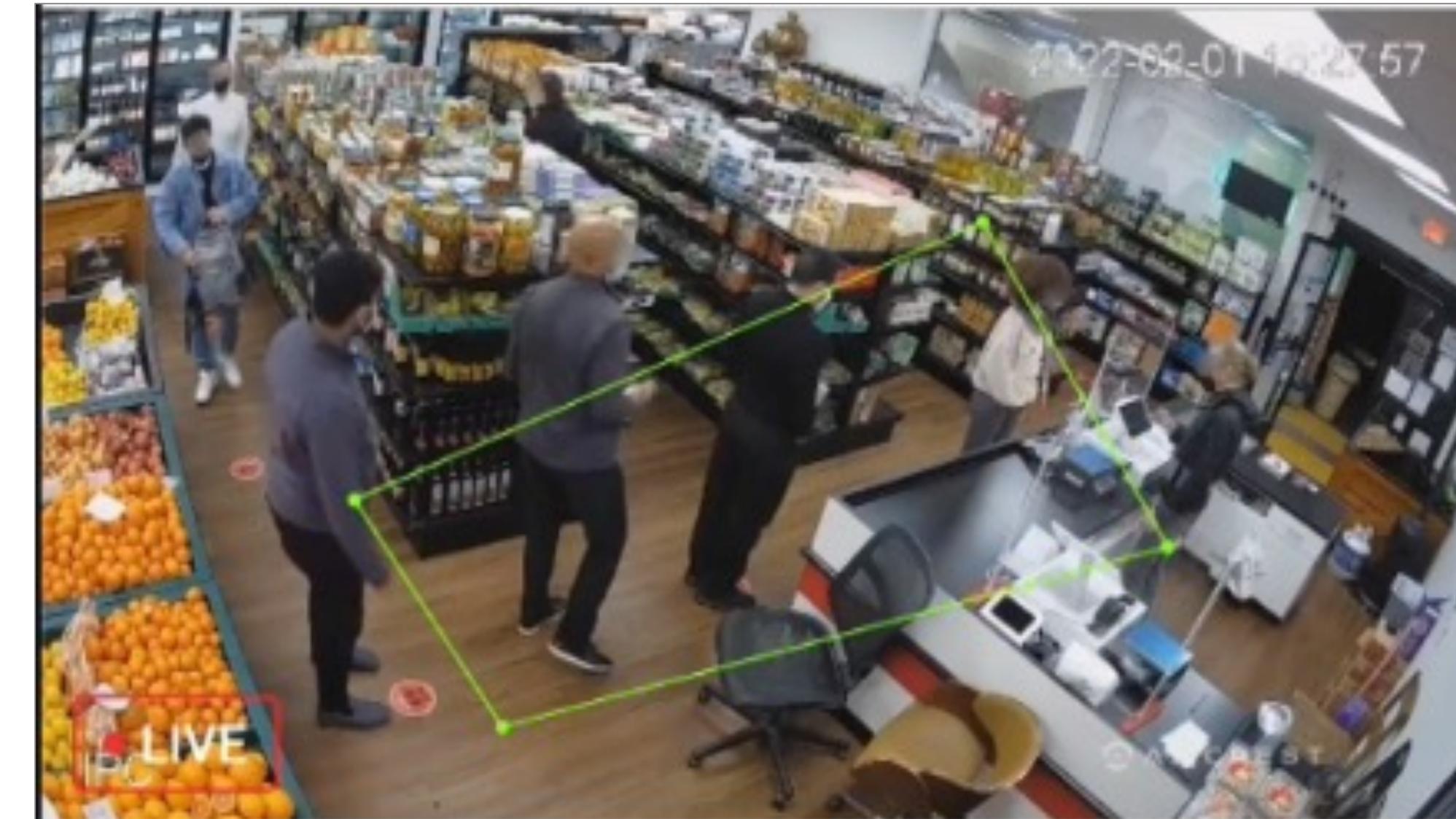
Years to Months = Huge Cost Savings

# Announcing Metropolis Microservices for Jetson



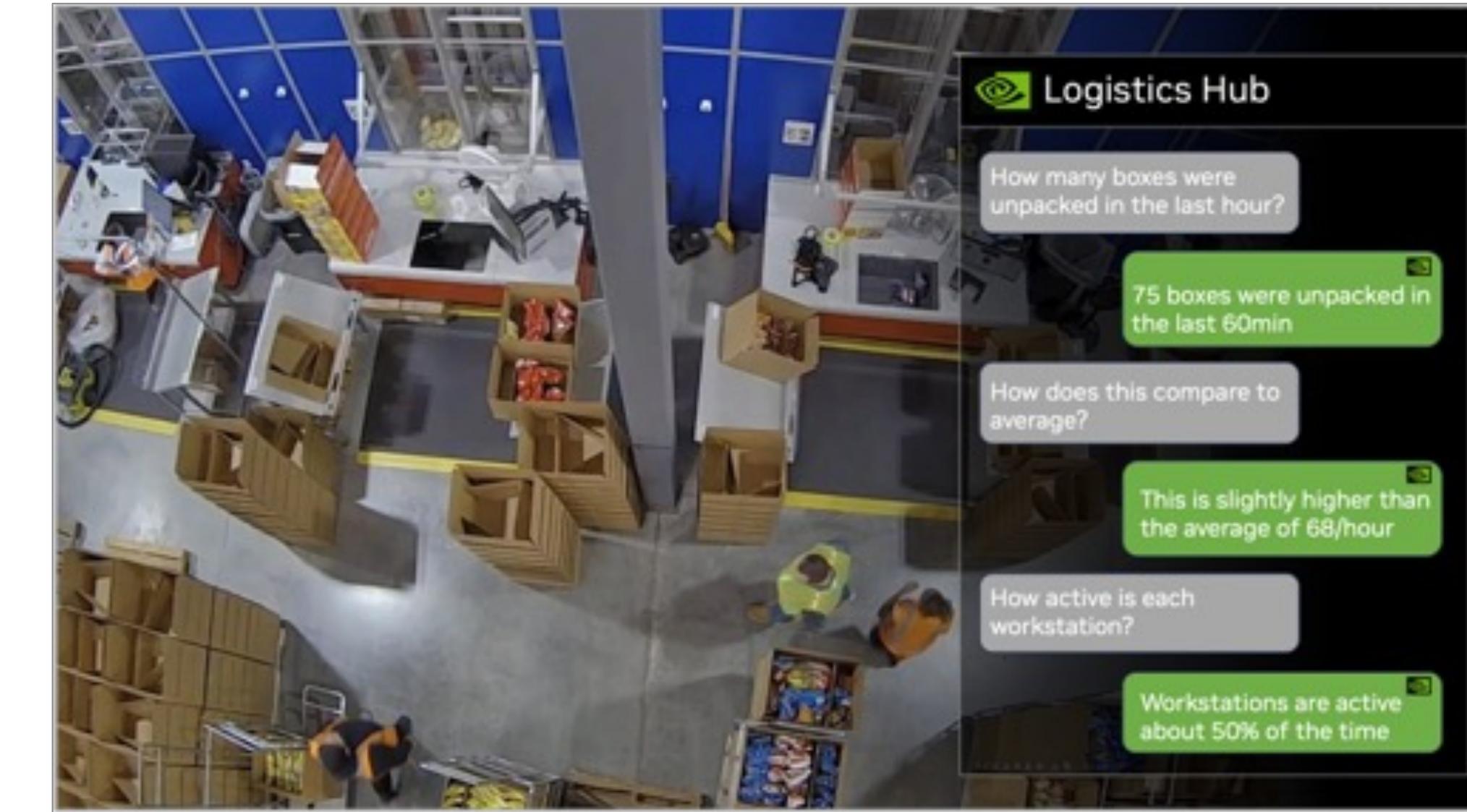
## Cloud-Native

- API-driven Microservices
- Fully Containerized
- Modular
- Extensible



## Suite of Pre-Built Services

- Sensor Storage & Management
- AI Perception Services
- IoT Gateway
- Monitoring, and more

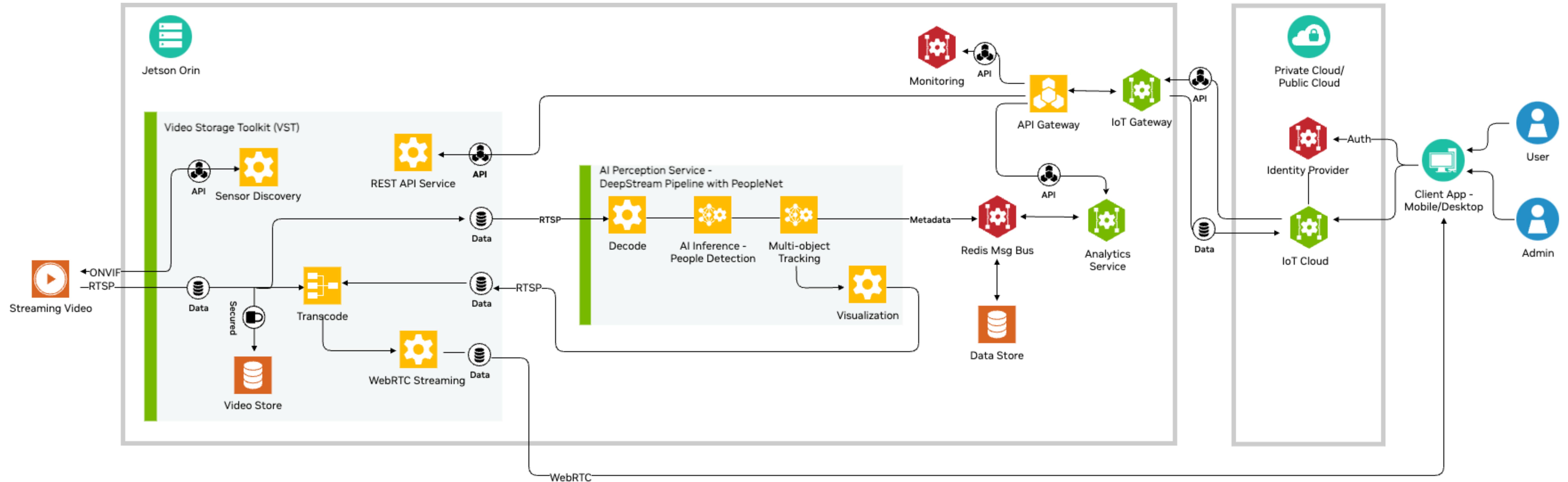


## Ready for Generative AI

- Flexible API-driven modules make prompting easy
- Powerful Jetson ORIN Compute capable of running multiple large models

# Cloud-Native Vision AI Workflows for the Edge

Use some or all the services



# Built for Maximum Performance

Freeing up Compute Resources for other Tasks

## Compute Intensive App

		Orin AGX 64	Orin NX 16
Input Streams		16	6
Resource Utilization	CPU	53%	61%
	GPU	42%	28%
	RAM	31GB	8 GB
	DLAx2	78%	33%
	PVA	34%	26%
	VIC	87%	55%

Resources free for other compute-intensive tasks

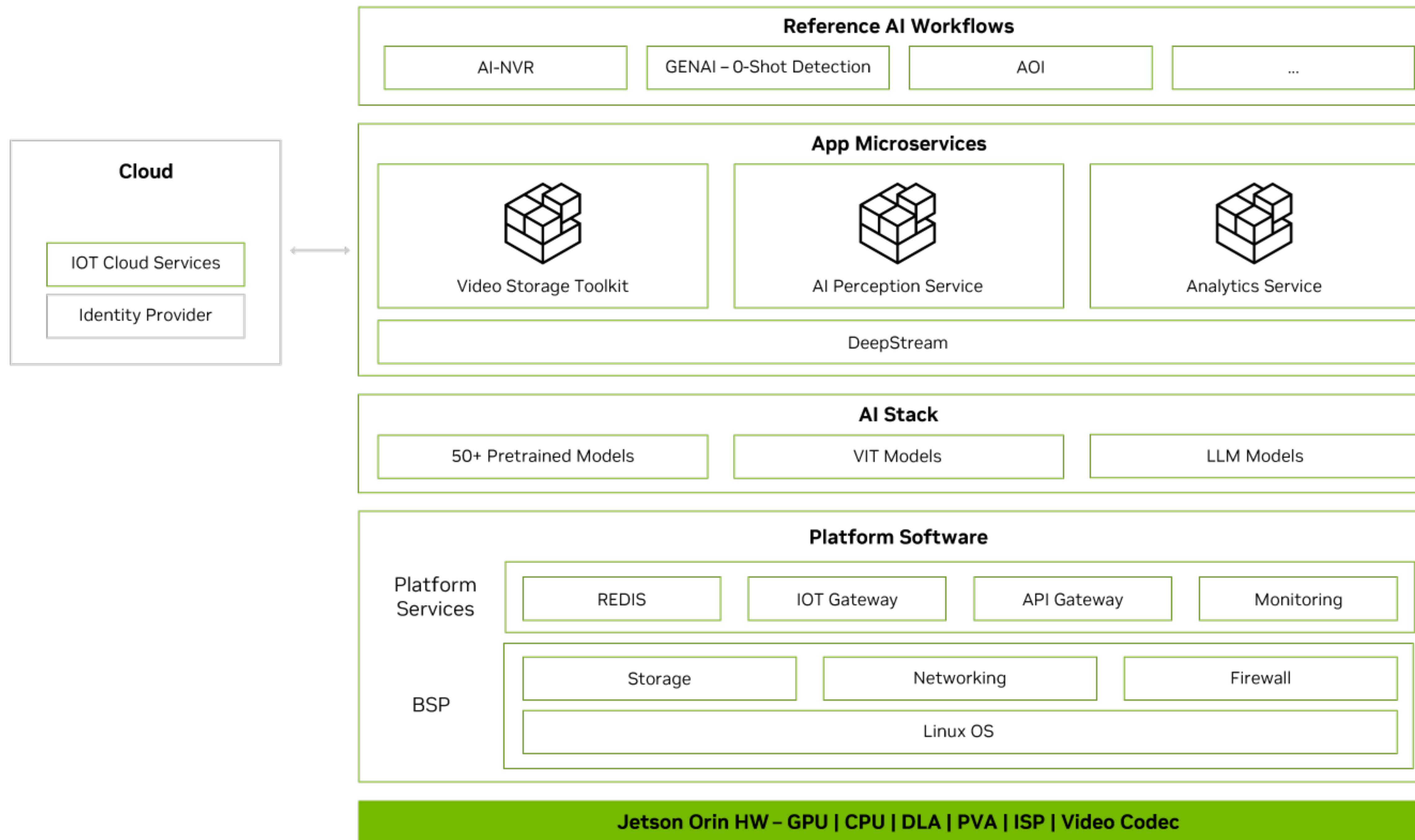
Maximized utilization of Orin's accelerators

- **Input:** 1080p @30fps
- **AI Perception Service**
  - 16x channels PeopleNet v2.6 (running on DLA)
  - Complex NvDCF Object Tracking (running on PVA)
- **Output:** 4 streams over WebRTC

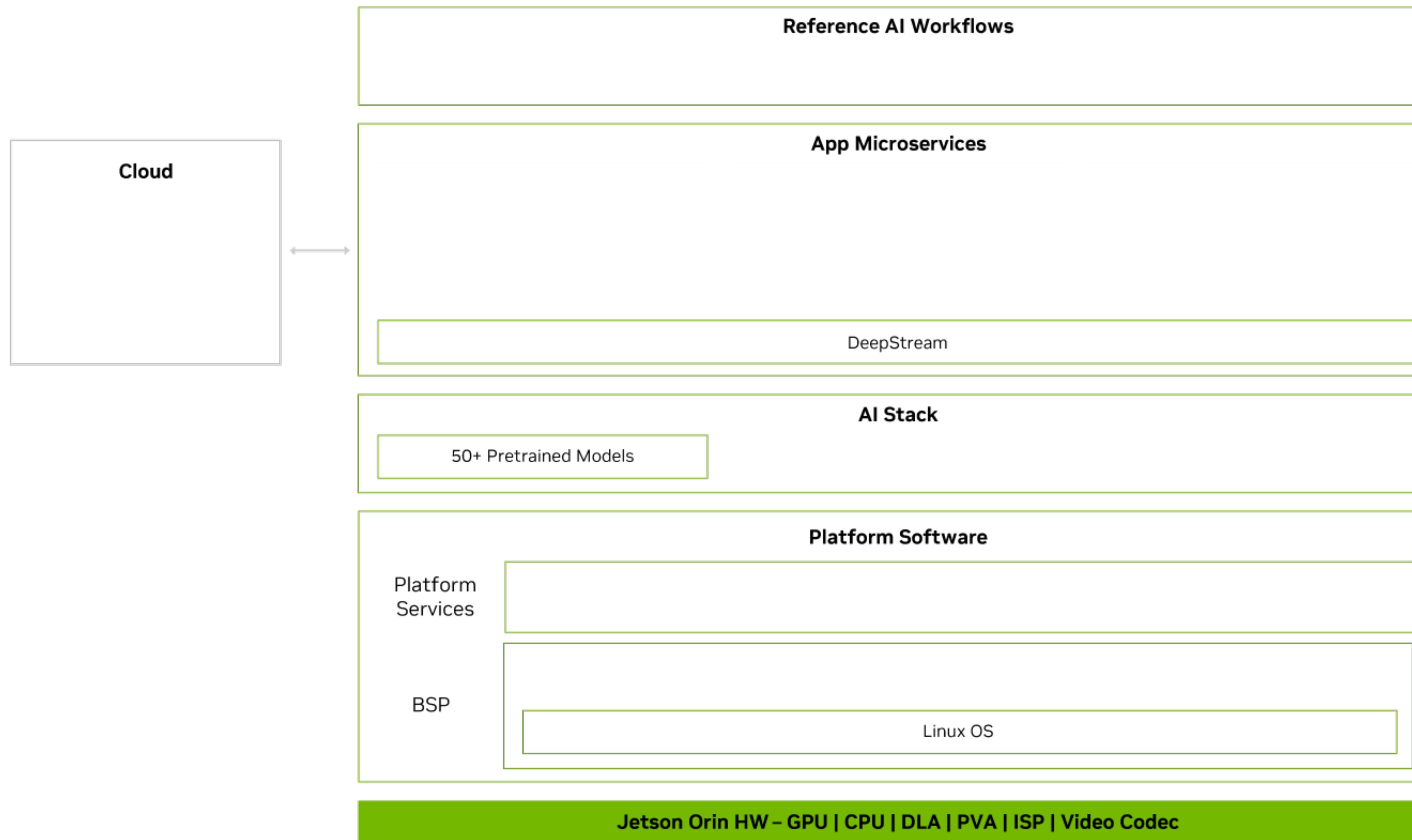
# Collection of APIs & Microservices

Category	SW Component	Delivery Mechanism (Debian/Container/other)	Hosted Where?
App Services	<ul style="list-style-type: none"><li>• Video Storage Toolkit (VST)</li><li>• AI Perception Service - DeepStream</li><li>• AI Perception Service - Generative AI</li><li>• Analytics Service - Line Crossing, ROIs, Counts</li><li>• Sensor Distribution &amp; Routing (SDR)</li></ul>	Containers	NGC
Cloud Services	<ul style="list-style-type: none"><li>• IoT Cloud Service</li></ul>	Containers with deployment scripts	NGC
Platform Services	<ul style="list-style-type: none"><li>• Redis</li><li>• API Gateway</li><li>• Monitoring</li><li>• IoT Gateway</li></ul>	Containers	NGC
Reference Workflow	<ul style="list-style-type: none"><li>• GenAI Sample App</li><li>• AI NVR Runtime app</li></ul>	Docker compose pkg	NGC
	Mobile app	APK	NGC

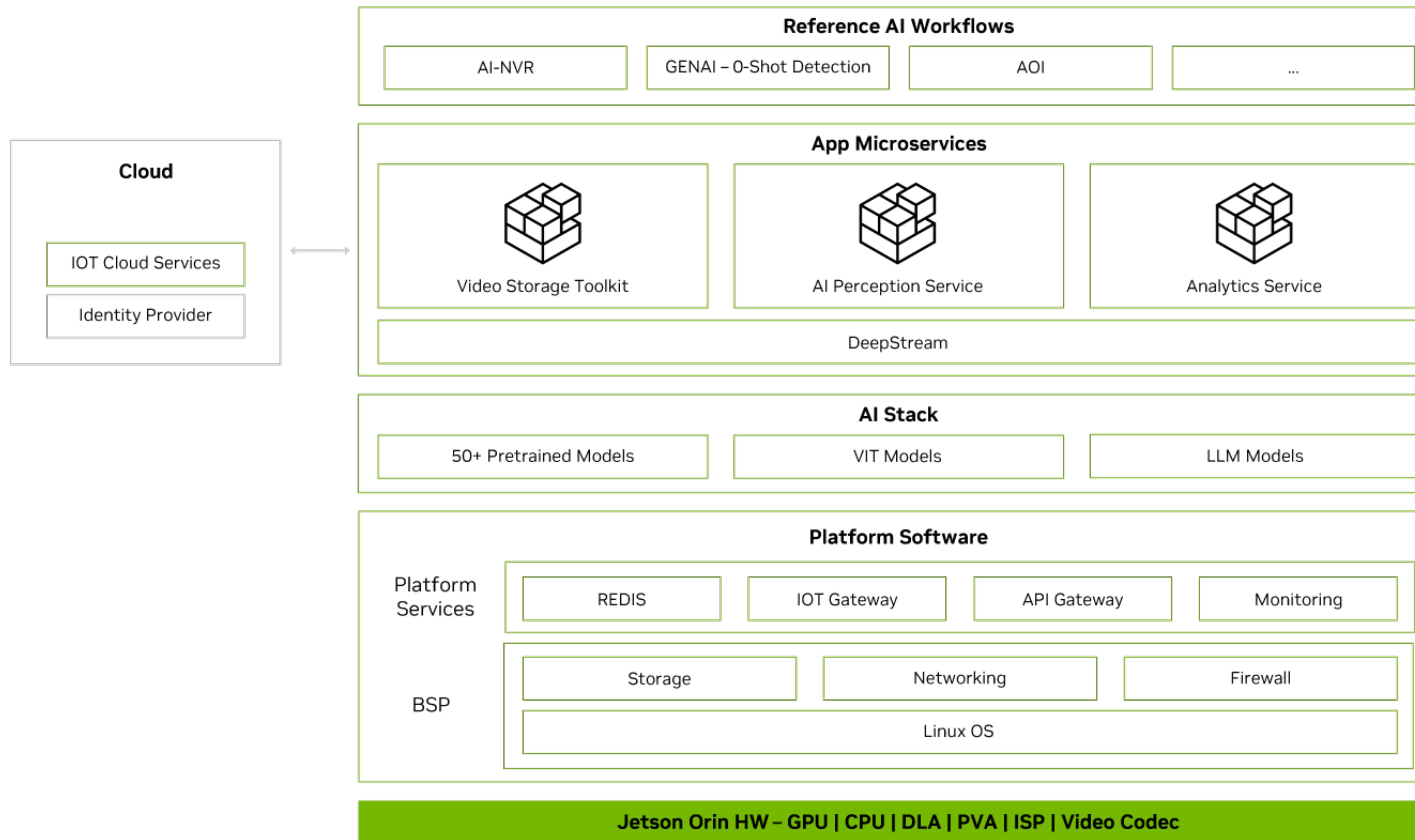
# Metropolis Microservices for Jetson



# Suite of Microservices and APIs (Before)

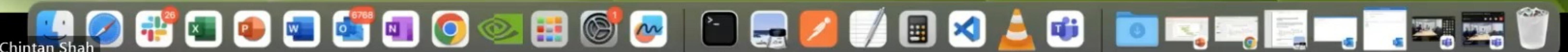


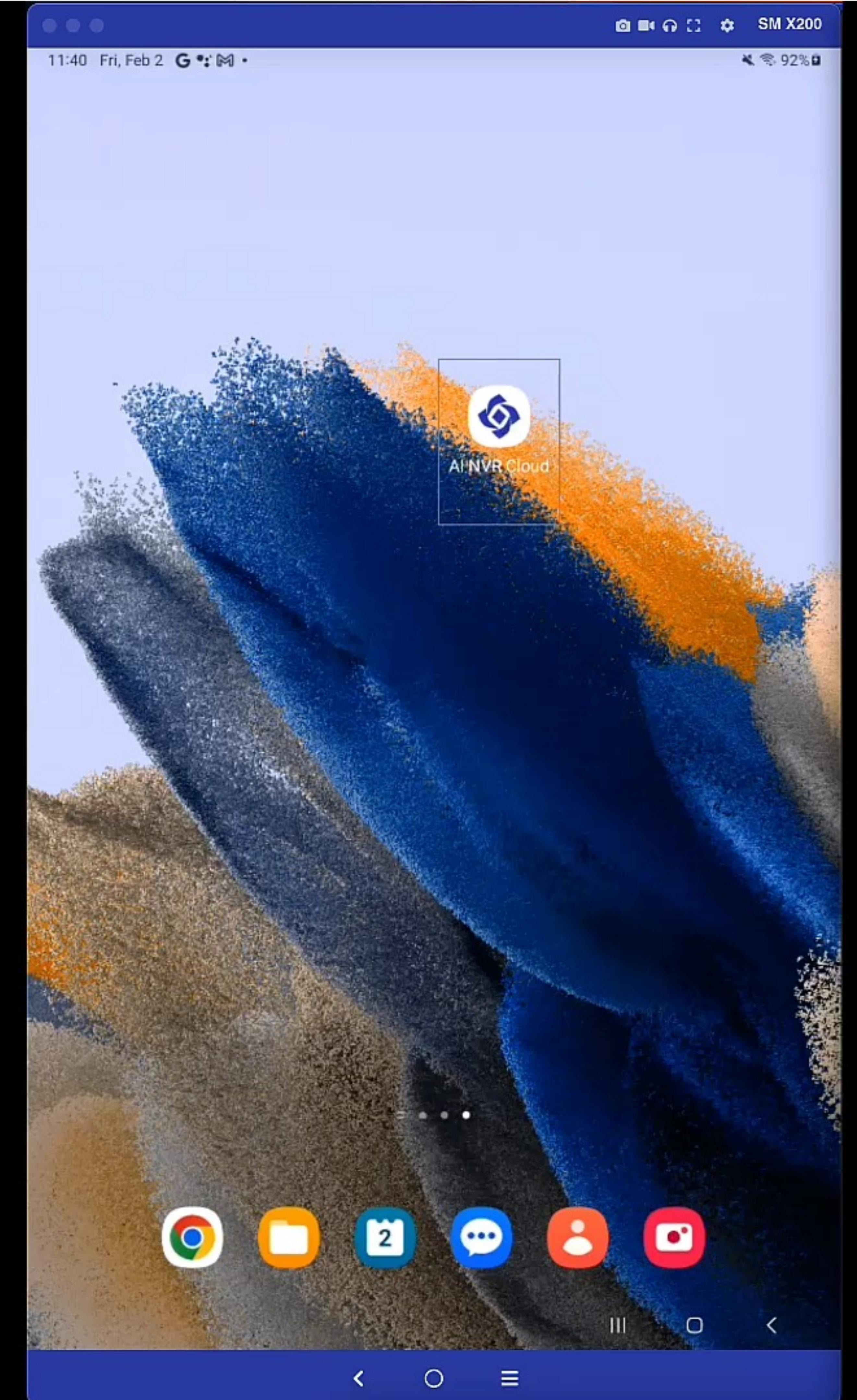
# Metropolis Microservices for Jetson



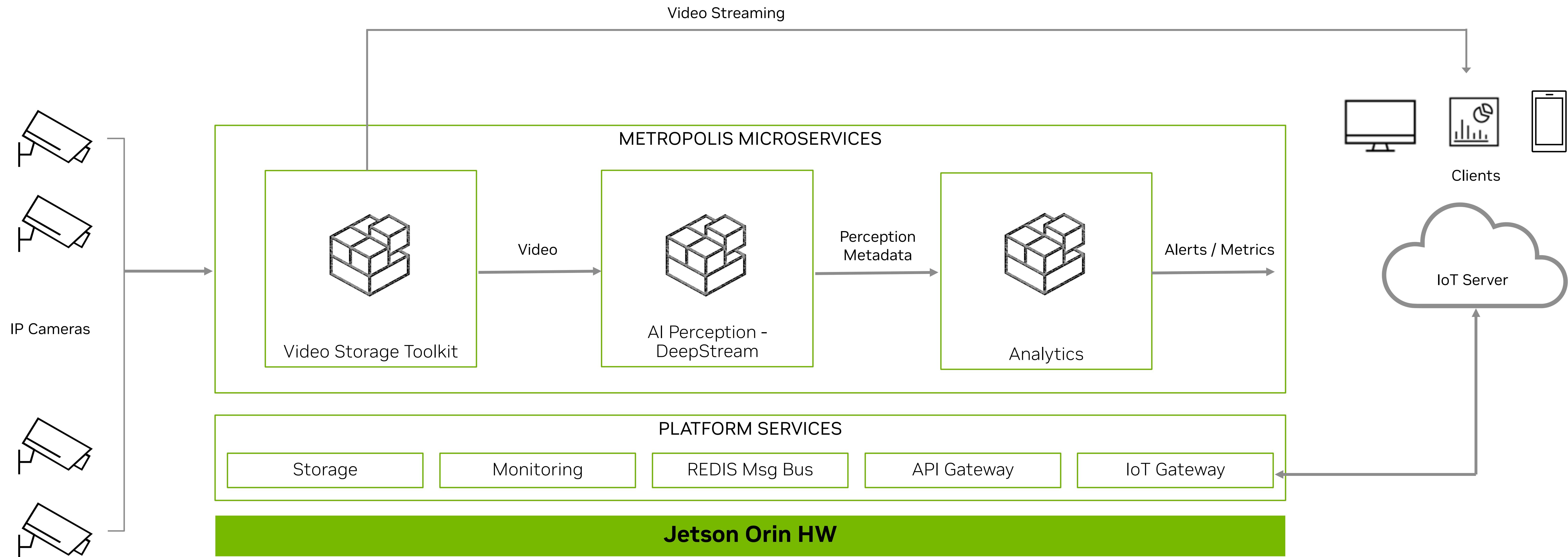
```
cshah — nvidia@tegra-ubuntu: ~/ai_nvr — ssh nvidia@172.17.170.143 — 245x64

base_compose.yaml compose_agx.yaml compose_nx.yaml config README.md
[nvidia@tegra-ubuntu:~/ai_nvr$ sudo docker compose -f compose_agx.yaml down --remove-orphans
[sudo] password for nvidia:
[+] Running 12/12
✓ Container deepstream                         Removed
✓ Container sdr-emdx                           Removed
✓ Container sdr                                Removed
✓ Container vst                                Removed
✓ Container emdx-webapi                         Removed
✓ Container ai_nvr-moj-http-based-init-sdr-emdx-1 Removed
✓ Container emdx-analytics-02                   Removed
✓ Container emdx-analytics-01                   Removed
✓ Container ai_nvr-moj-http-based-init-sdr-1   Removed
✓ Container ai_nvr-moj-init-vst-1              Removed
✓ Container ai_nvr-moj-init-ds-1              Removed
✓ Container ai_nvr-moj-http-based-init-emdx-analytics-1 Removed
nvidia@tegra-ubuntu:~/ai_nvr$ 
nvidia@tegra-ubuntu:~/ai_nvr$ 
nvidia@tegra-ubuntu:~/ai_nvr$ 
nvidia@tegra-ubuntu:~/ai_nvr$ sudo docker ps
[sudo] password for nvidia:
Sorry, try again.
[sudo] password for nvidia:
CONTAINER ID   IMAGE                               COMMAND             CREATED          STATUS           PORTS          NAMES
fb303a2ddaf0   prom/node-exporter                "/bin/node_exporter ..." 24 hours ago   Up 24 hours    nodeexporter
0c45f86b9c4c   grafana/grafana                  "/run.sh"          24 hours ago   Up 24 hours    grafana
290fb841f47d   nvcr.io/e7ep4mig3lne/release/its-monitoring:v1.6.0_arm64  "/root/its_monitorin..." 24 hours ago   Up 24 hours    its-monitoring
6a7f004e3920   prom/prometheus                 "/bin/prometheus --c..." 24 hours ago   Up 24 hours    prometheus
1fd2bdfde9c2   prom/alertmanager               "/bin/alertmanager --..." 24 hours ago   Up 24 hours    alert-manager
5d16583a7227   prom/pushgateway                "/bin/pushgateway"  24 hours ago   Up 24 hours    push-gateway
38c669eaa6bc   nvcr.io/e7ep4mig3lne/release/prov-agent:v1.1.0_arm64v8  "/opt/prov-agent/ent..." 24 hours ago   Up 24 hours    prov-agent
fa0c3e9874fe   nvcr.io/e7ep4mig3lne/release/tcpmux-client:v1.2.0_arm64v8  "/opt/tcpmux-client/..." 24 hours ago   Up 24 hours    tcpmux-client
0ca81fcfb1b1   nvcr.io/e7ep4mig3lne/release/ialpha-ingress-arm64v8:0.8  "sh -c '/nginx.sh 2>..." 24 hours ago   Up 24 hours    ingress
2a36814a55f8   redisfab/redistimeseries:master-arm64v8-jammy        "docker-entrypoint.s..." 24 hours ago   Up 24 hours    redis
3857eaaaadb3   nvcr.io/e7ep4mig3lne/release/vst:nvstreamer_v0.2.24_aarch64  "sh -c '/root/vst_re..." 25 hours ago   Up 24 hours    nvstreamer
[nvidia@tegra-ubuntu:~/ai_nvr$ sudo docker compose -f compose_agx.yaml up -d --force-recreate
[+] Running 12/12
✓ Container ai_nvr-moj-init-vst-1              Exited
✓ Container ai_nvr-moj-http-based-init-sdr-1   Exited
✓ Container ai_nvr-moj-init-ds-1              Exited
✓ Container ai_nvr-moj-http-based-init-sdr-emdx-1 Exited
✓ Container ai_nvr-moj-http-based-init-emdx-analytics-1 Exited
✓ Container emdx-webapi                      Started
✓ Container emdx-analytics-01                 Started
✓ Container sdr                            Started
✓ Container emdx-analytics-02                 Started
✓ Container deepstream                     Started
✓ Container vst                           Started
✓ Container sdr-emdx                     Started
nvidia@tegra-ubuntu:~/ai_nvr$ sudo docker compose -f compose_agx.yaml down --remove-orphans
[+] Running 12/12
✓ Container sdr-emdx                         Removed
✓ Container emdx-webapi                      Removed
✓ Container deepstream                       Removed
✓ Container sdr                             Removed
✓ Container vst                            Removed
✓ Container emdx-analytics-02                Removed
✓ Container emdx-analytics-01                Removed
✓ Container ai_nvr-moj-http-based-init-sdr-emdx-1 Removed
✓ Container ai_nvr-moj-http-based-init-sdr-1   Removed
✓ Container ai_nvr-moj-init-vst-1              Removed
✓ Container ai_nvr-moj-init-ds-1              Removed
✓ Container ai_nvr-moj-http-based-init-emdx-analytics-1 Removed
nvidia@tegra-ubuntu:~/ai_nvr$ 
```





# System Architecture for Retail Use Case



# Video Storage Toolkit (VST)

## WHAT DOES IT HELP WITH

Set of APIs to easily build complex multi-camera ingestion, storage, and streaming.

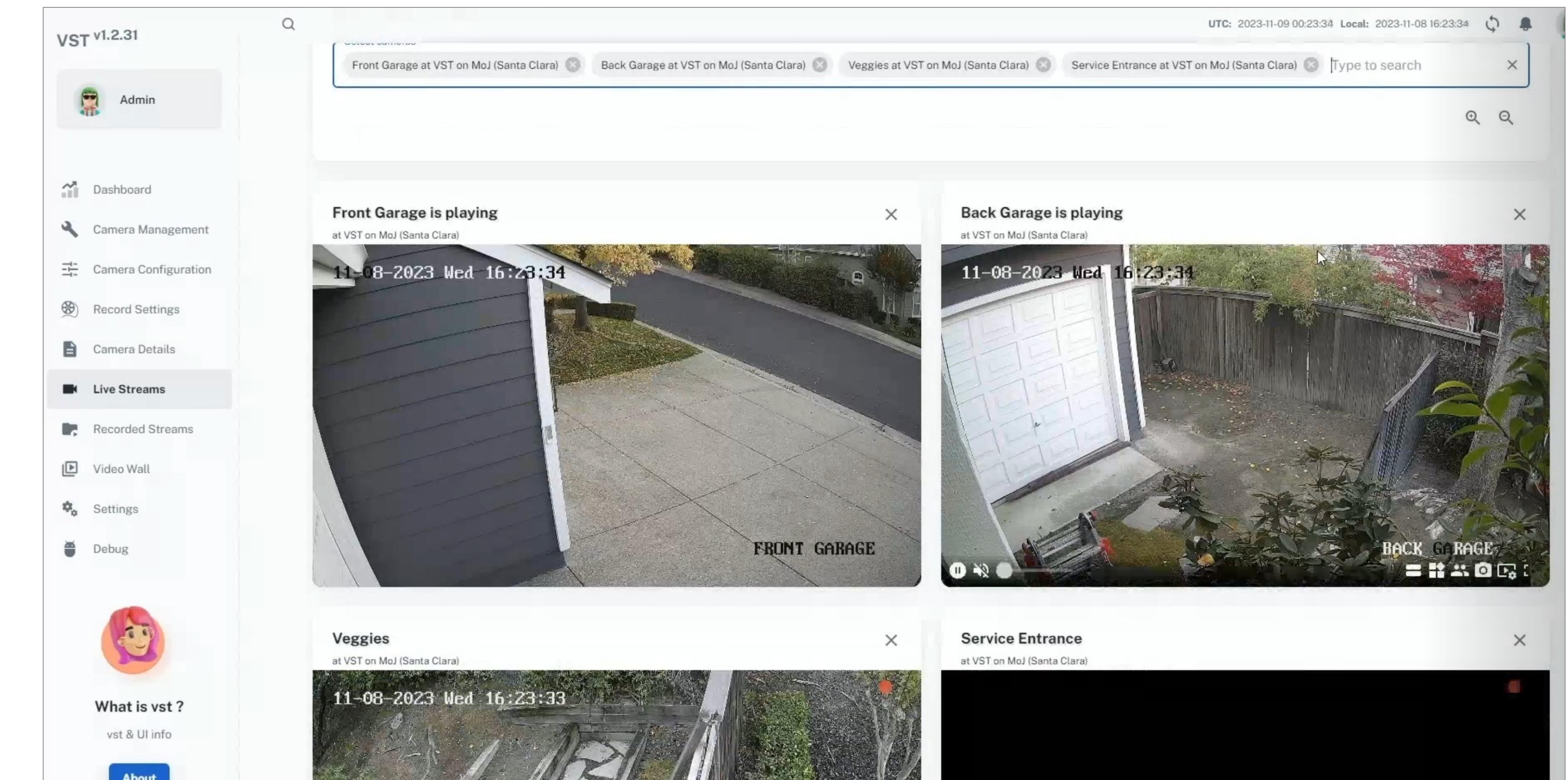
## WHAT'S INCLUDED

Collection of 20+ APIs including:

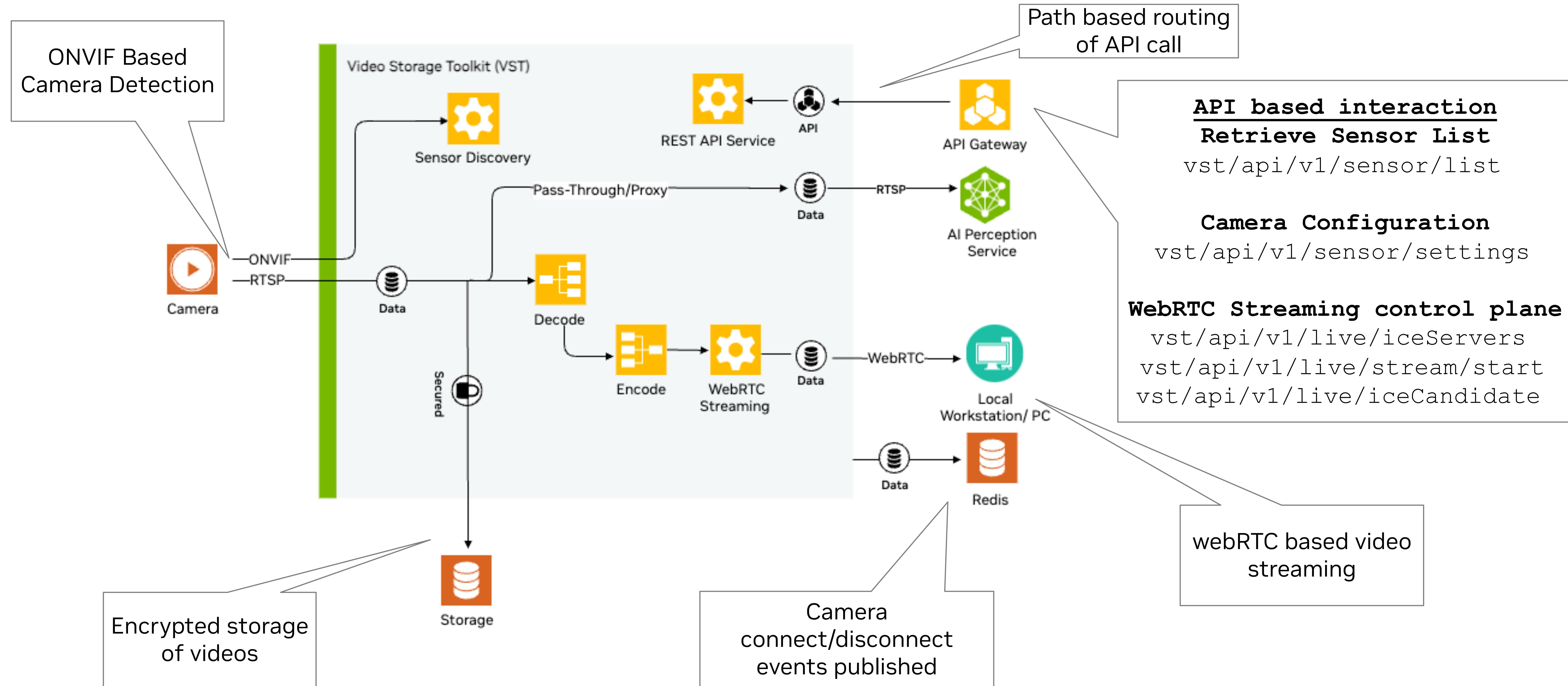
- ONVIF camera discovery, monitoring, & management
- Video recording / storage
- Media server, WebRTC streaming
- WebUI for setup and visualization

## HOW TO INTEGRATE INTO YOUR APPLICATION

REST APIs



# Video Storage Toolkit Functionality



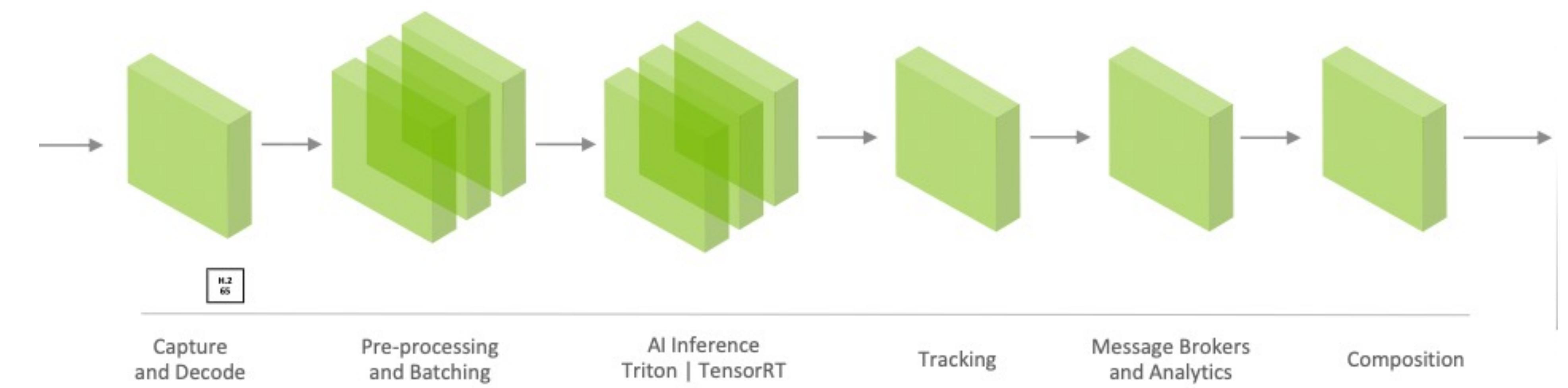
# AI Perception Service - Pre-Built DeepStream Pipelines

## WHAT DOES IT HELP WITH

Creating an efficient pipeline to process pixels to metadata leveraging DeepStream SDK

## WHAT'S INCLUDED

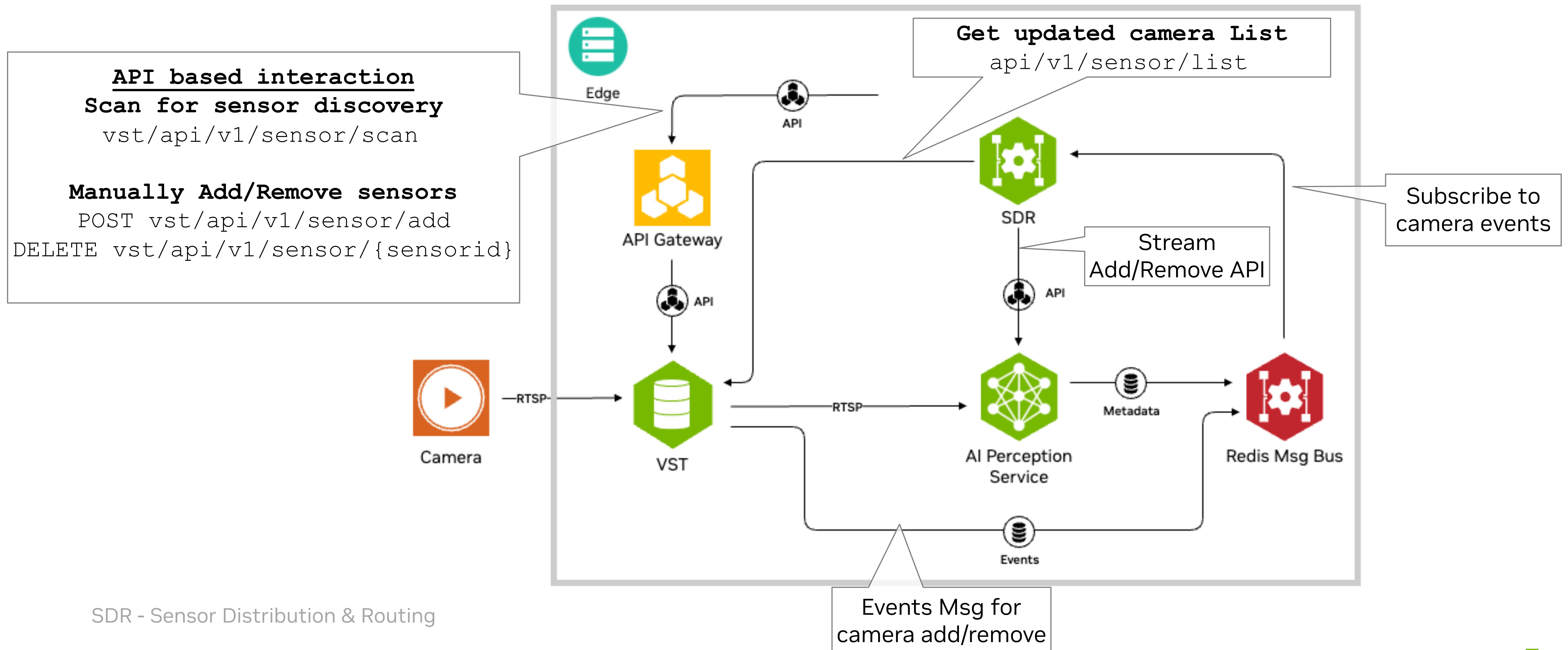
- Up to 16 streams across 2 DLAs (Orin) for PeopleNet
- World-class NvDCF tracker running on PVA
- Dynamic stream discovery, add/remove, reconnection
- Metadata generation and publishing on to Redis
- Integration of app KPIs with Monitoring services



## HOW TO INTEGRATE INTO YOUR CODE

Deploy using Docker container

# Automatic Stream Addition & Removal with SDR



# Analytics Service - Line Crossing, ROIs, Count

## WHAT DOES IT HELP WITH

Using APIs to configure insights and alerts such as line-crossing, ROIs and FOV

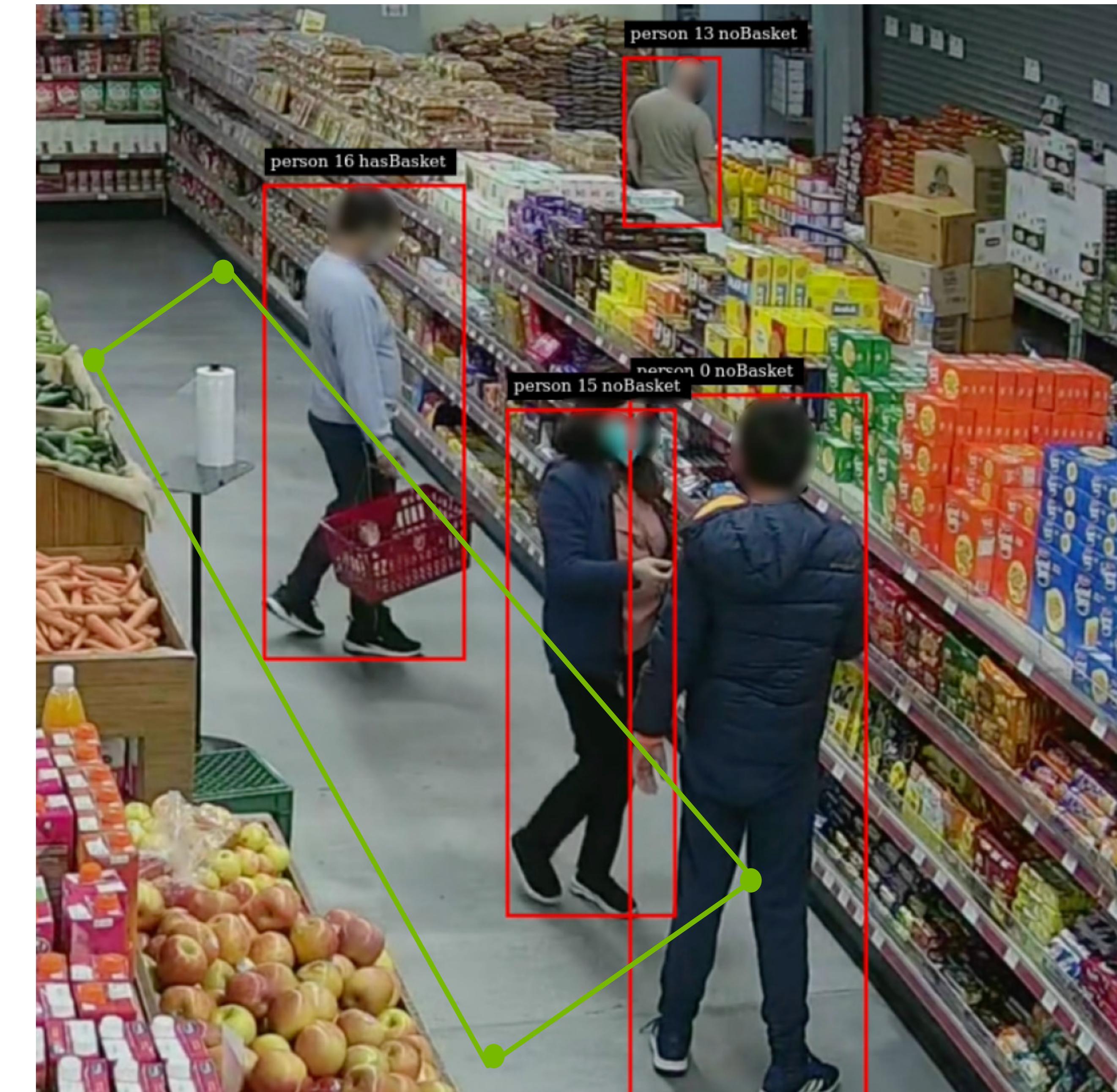
Can help identify traffic flow patterns, entry/exit tracking, encroachment into specific areas

## WHAT'S INCLUDED

- Integration into other platform services such as Redis - publishing system event, Monitoring and Storage
- Visualization through overlays
- REST APIs to get insights, set conditions, etc.

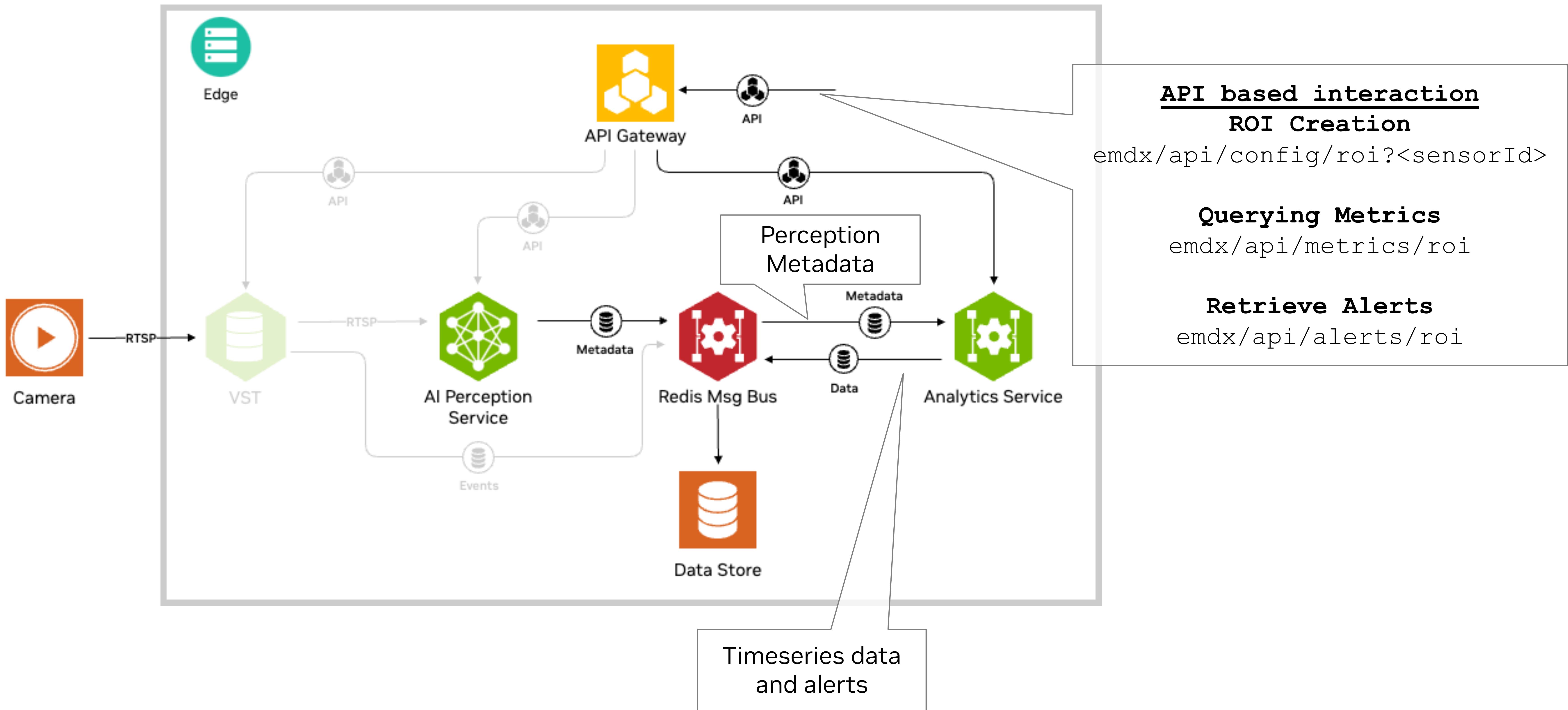
## HOW TO INTEGRATE INTO YOUR APPLICATION

REST APIs



ROI - Region of Interest  
FOV - Field of View

# Streaming Analytics for Spatio-Temporal Understanding



# IoT and Cloud Service

## WHAT DOES IT HELP WITH

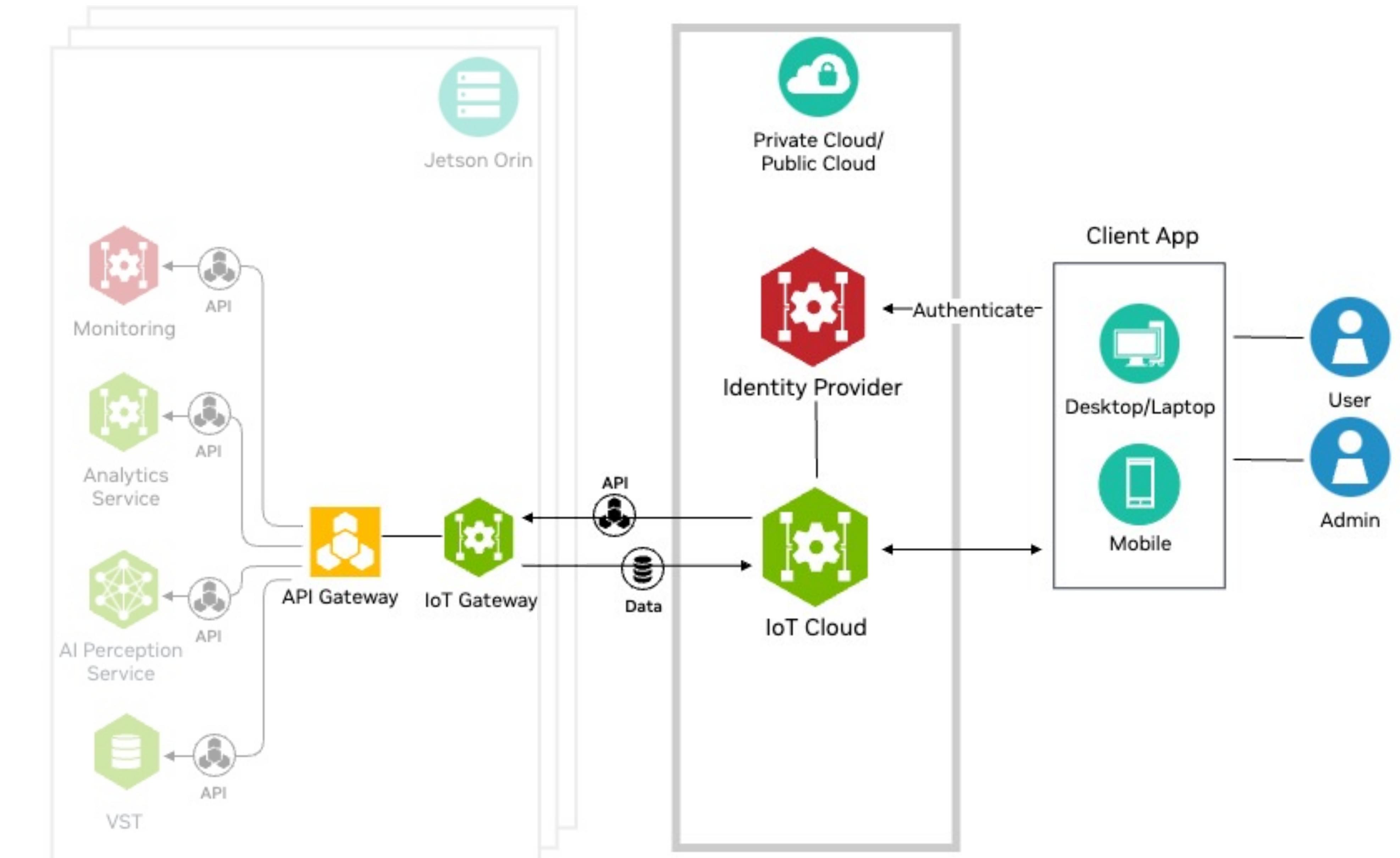
Remote, secure invocation of microservices APIs on devices from client applications

## WHAT'S SUPPORTED

- Reference cloud implementation with recipe to deploy on any Cloud
- Setup flows: device securely connecting to cloud through OTP; user getting access authorization through claim code
- “Always on” encrypted, bi-directional communication link between Edge & Cloud
- Public Proxy endpoint in the cloud using which clients can invoke device APIs
- Authentication, authorization, user mgmt., device claim in the cloud

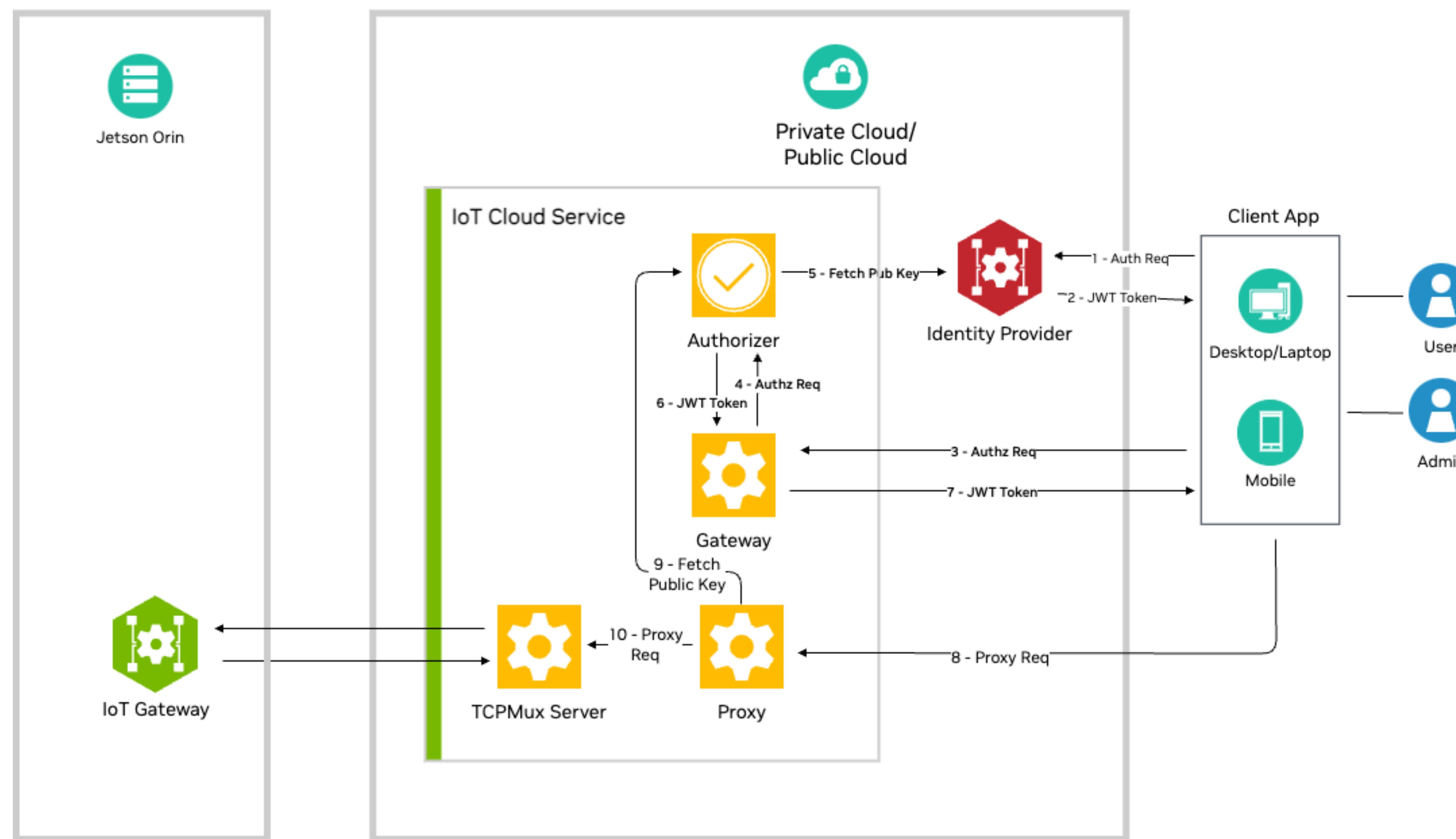
## HOW TO INTEGRATE INTO YOUR CODE

REST APIs



# 3-Step Process for Device API Calls by Client App

1. Authenticate with IDP using username/password to obtain ID token
2. Retrieve JWT token (containing authorization scopes) from Gateway by using ID Token
3. Invoke device APIs while presenting the JWT authorization tokens



# Application Deployment Using Docker Compose

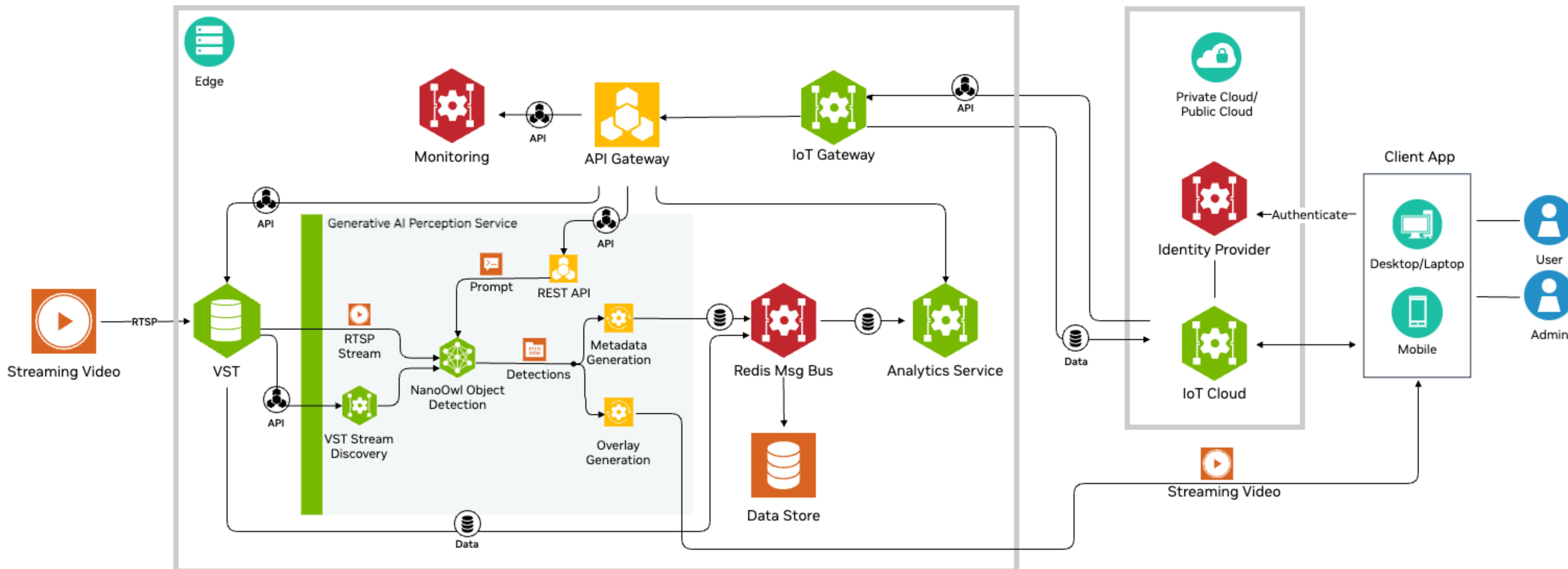
List all the microservice containers along with config, storage, startup ordering and deploy!

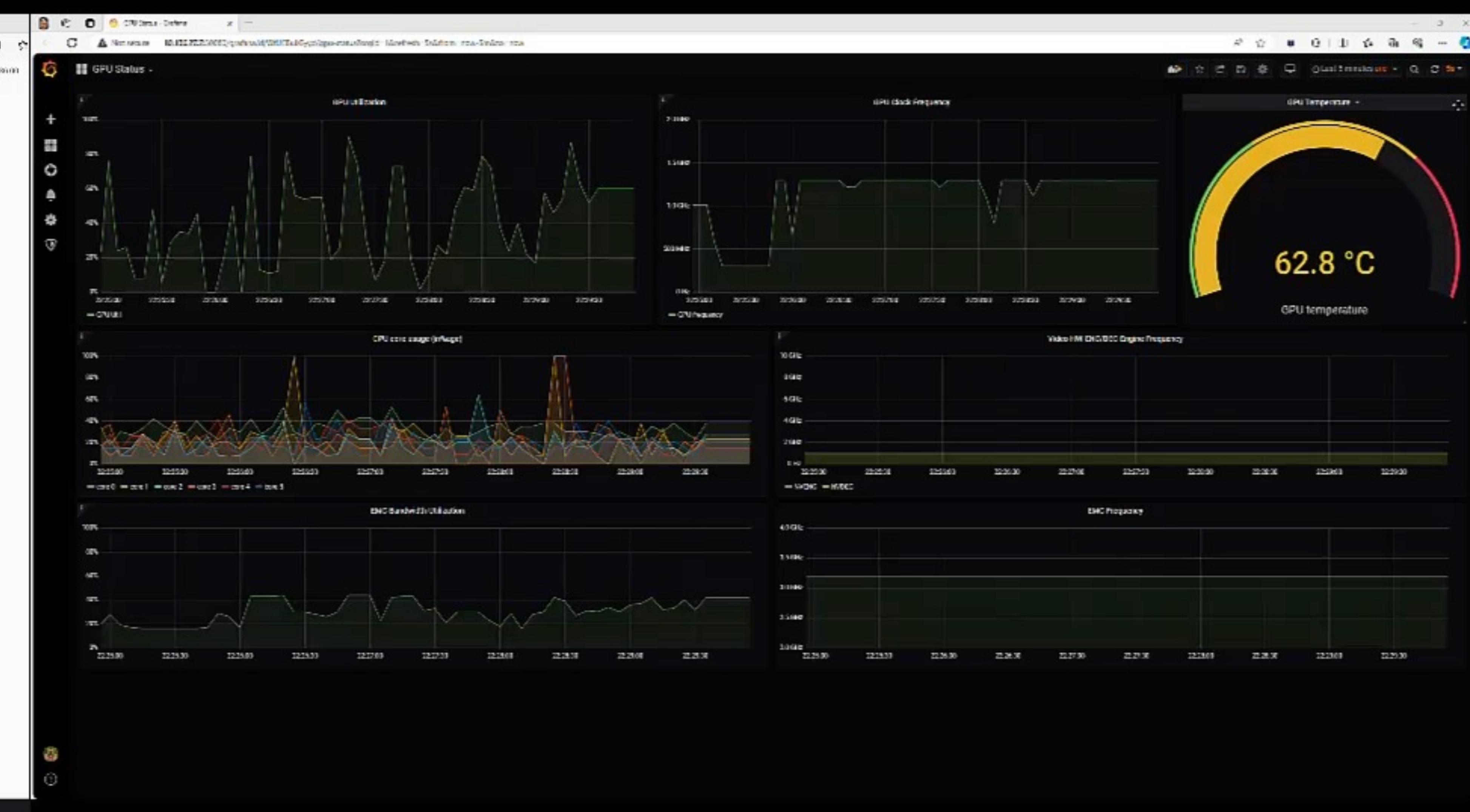
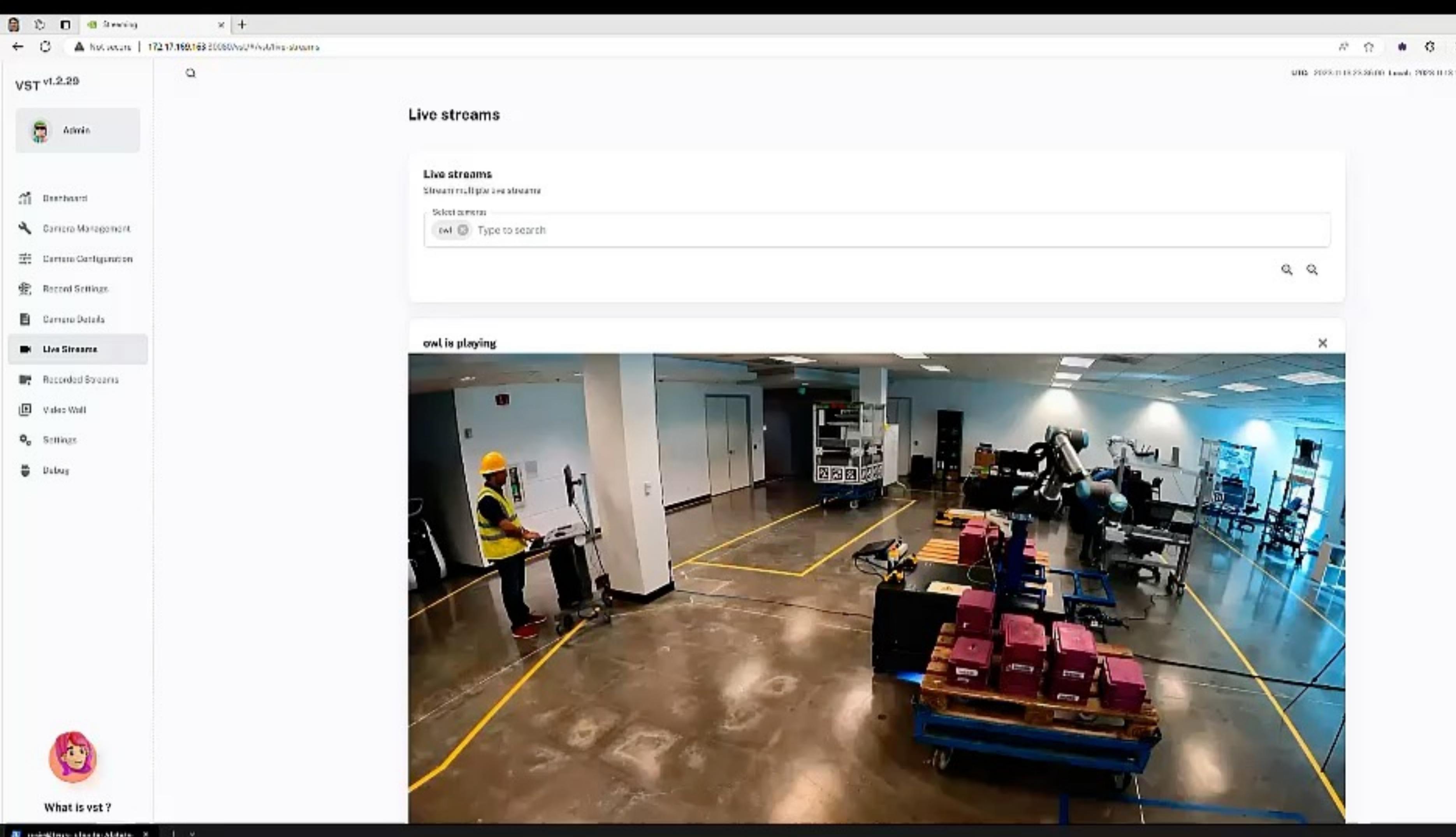
```
vst:  
  image: nvcr.io/e7ep4mig3lne/release/vst:v1.2.37_aarch64  
  .....  
  container_name: vst  
  entrypoint: sh -c '/root/vst_release/launch_vst --debug-level 3 2>&1  
  volumes:  
    - ./config/vst/vst_config.json:/root/vst_release/configs/vst_config.json  
    - ./config/vst/vst_storage.json:/root/vst_release/configs/vst_storage.json  
  depends_on:  
    moj-init-vst:  
      condition: service_completed_successfully  
  deploy:  
    resources:  
      limits:  
        memory: 5600M  
    restart_policy:  
      condition: always  
emdx-analytics-01:  
  image: nvcr.io/e7ep4mig3lne/release/emdx-analytics:mmj_v1  
  environment:  
    CONFIG_LOCATION: "/config"  
    PORT: 6001  
    INSTANCE_ID: emdx-analytics-01  
  .....
```

The diagram illustrates various Docker Compose configuration blocks with corresponding annotations:

- Container name:** Points to the `image` field in the `vst` service definition.
- config:** Points to the `volumes` section of the `vst` service, specifically the two volume definitions.
- Storage mounts:** Points to the `volumes` section of the `vst` service.
- Init container for startup ordering:** Points to the `depends_on` section of the `vst` service, specifically the dependency on `moj-init-vst`.
- Failure restart:** Points to the `restart_policy` section of the `vst` service, specifically the `condition: always` setting.
- Environment setting:** Points to the `environment` section of the `emdx-analytics-01` service, specifically the `CONFIG_LOCATION` variable.

# Integrating Generative AI Services with Metropolis Microservices





# Summary

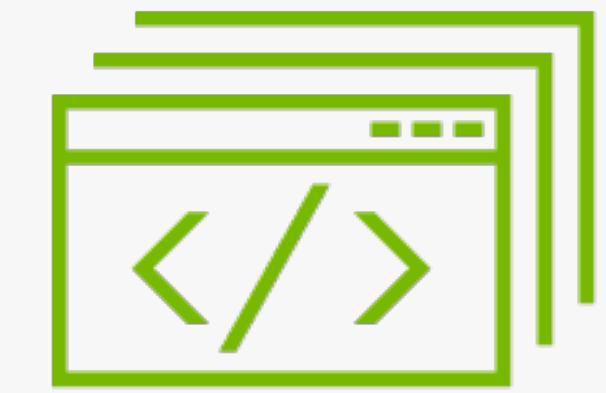
- Announced availability of Metropolis Microservices for Jetson
- Cloud-native Microservice Architecture
- API-driven platform for building AI apps at the Edge
- Accelerate development and reduce time to market

## Get Started

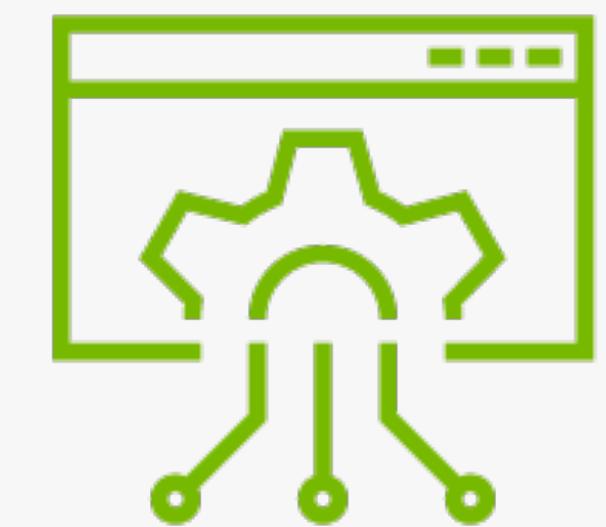
<https://developer.nvidia.com/metropolis-microservices/jetson-get-started>



Build complex applications quickly with APIs



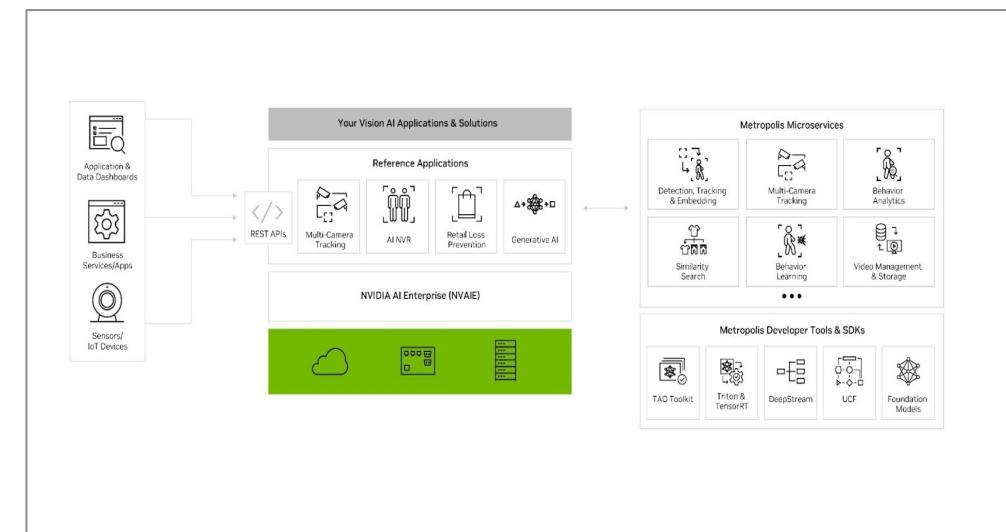
Maximize your system resources



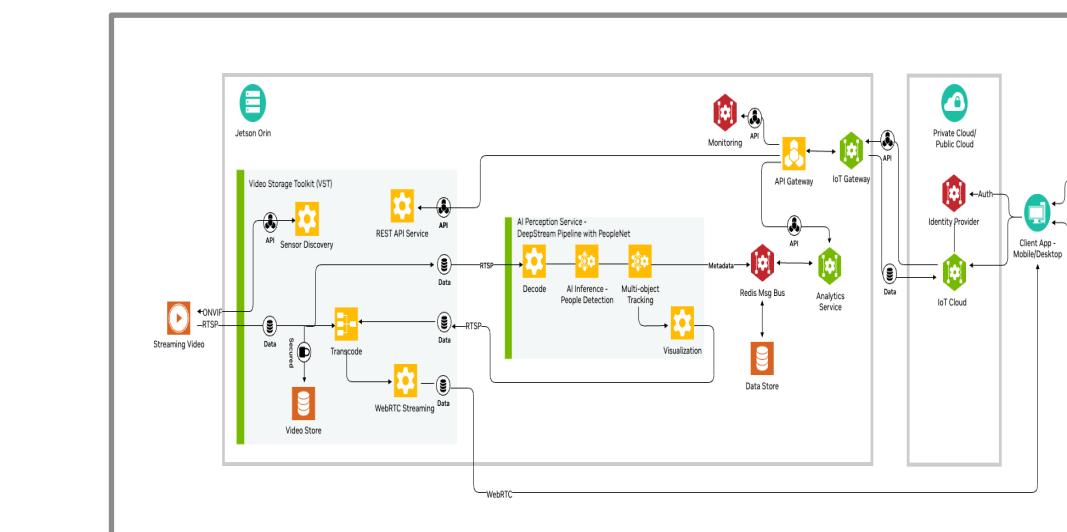
Integrate your own custom microservices

# Metropolis Microservices for Jetson - Resources

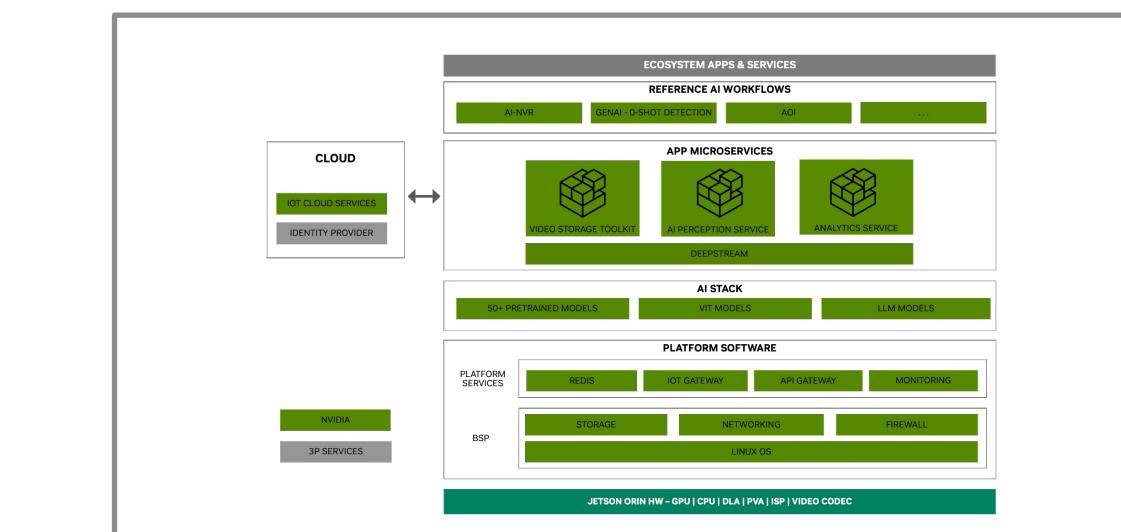
## Web



[Product Page](#)

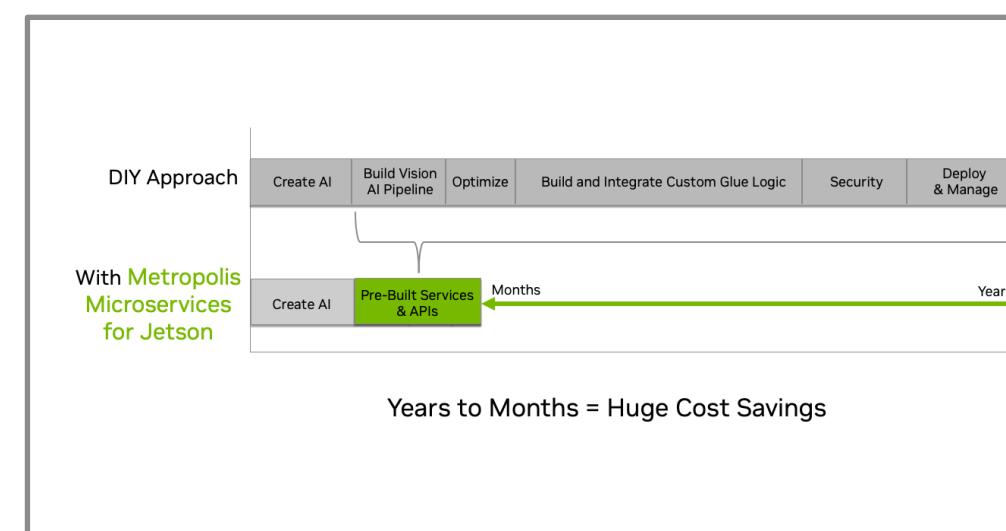


[Get Started Page](#)

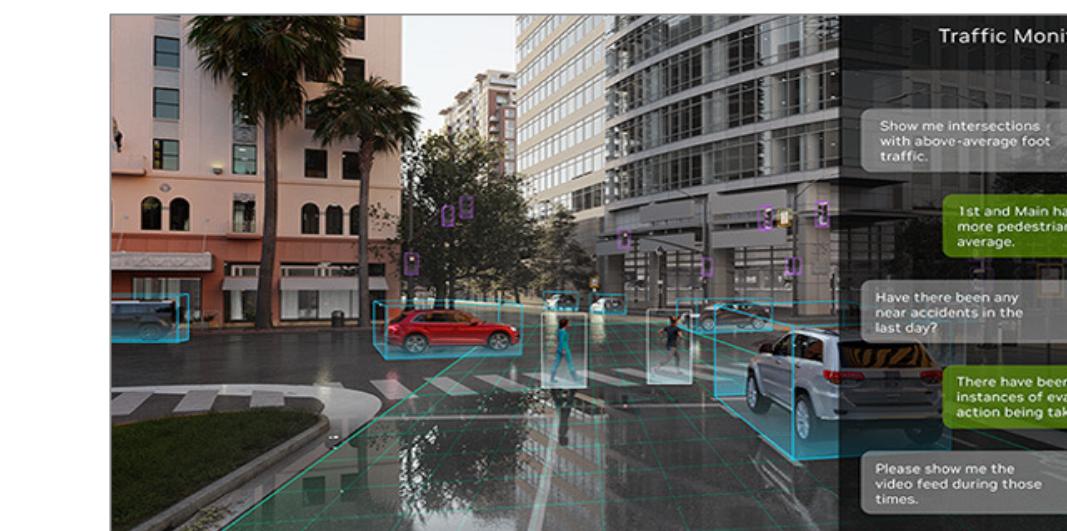


[SW Download](#)

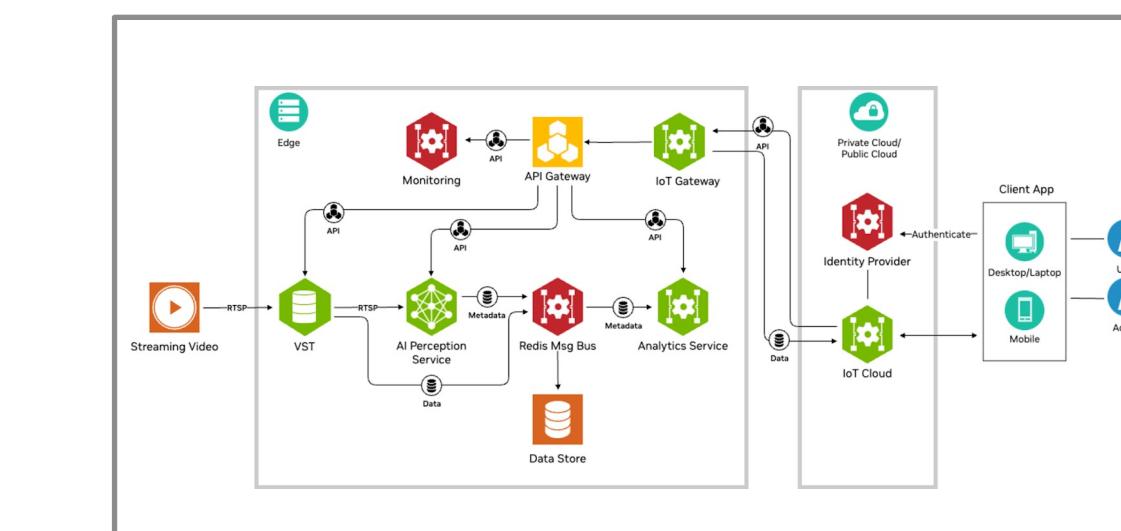
## Blogs



[Dev News](#)

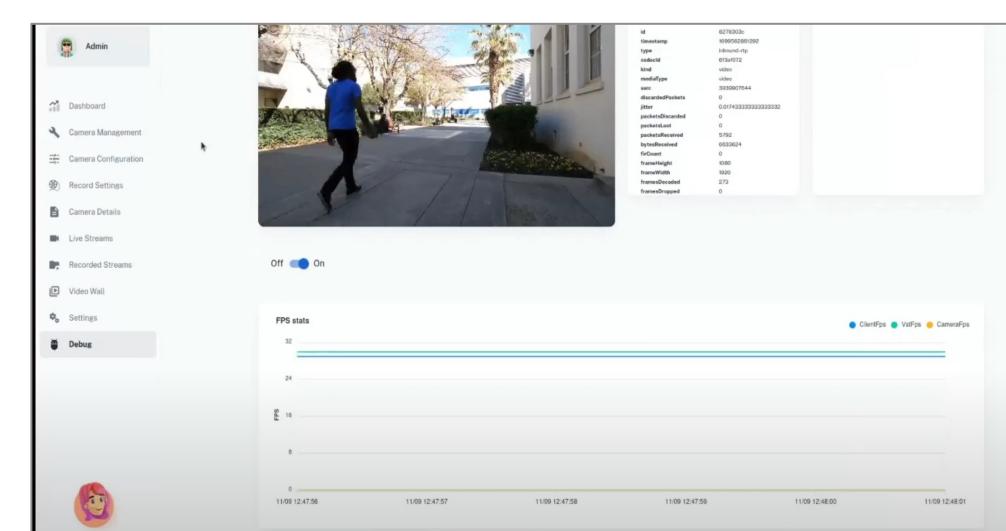


[Gen AI](#)

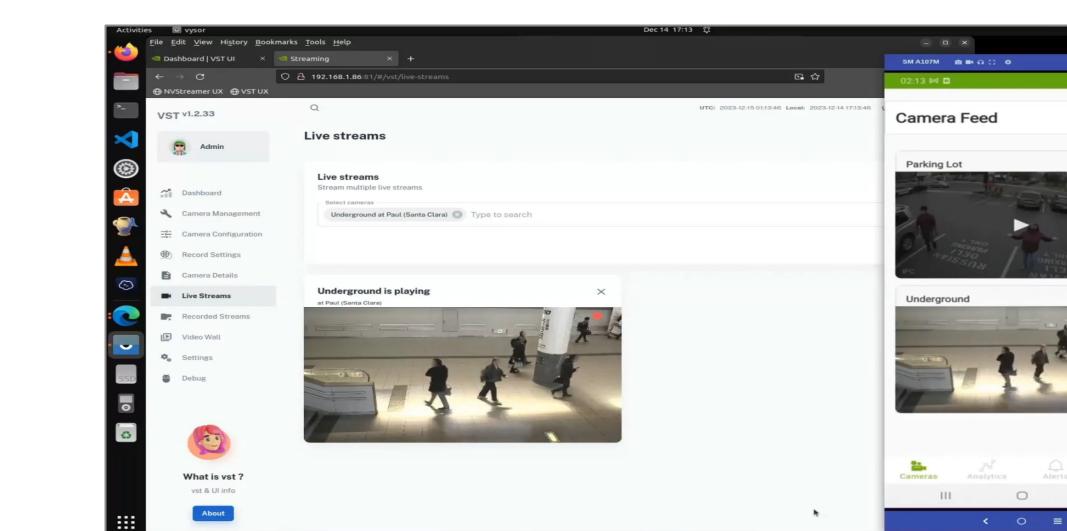


[API Workflow](#)

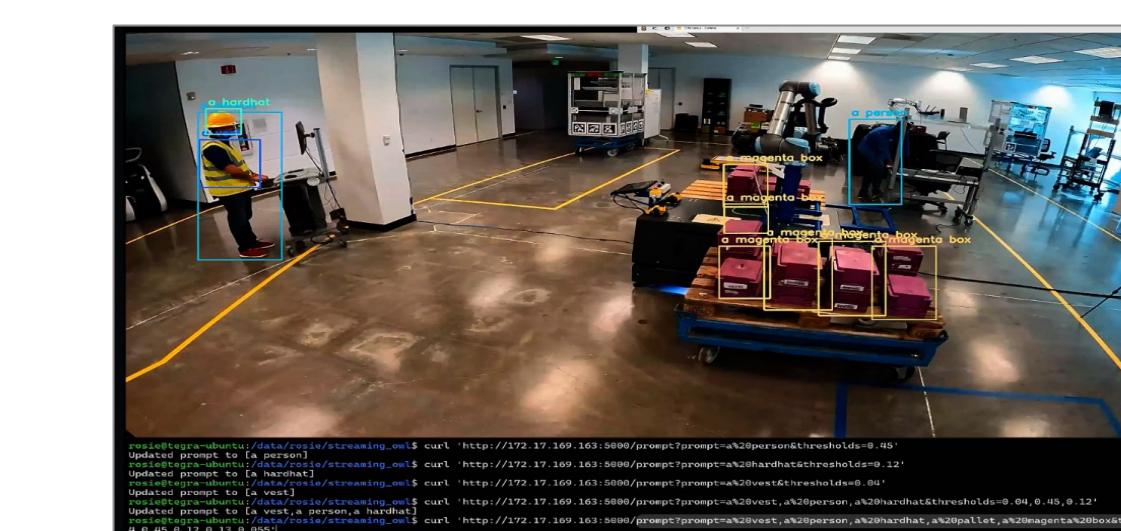
## Videos



[VST Demo](#)

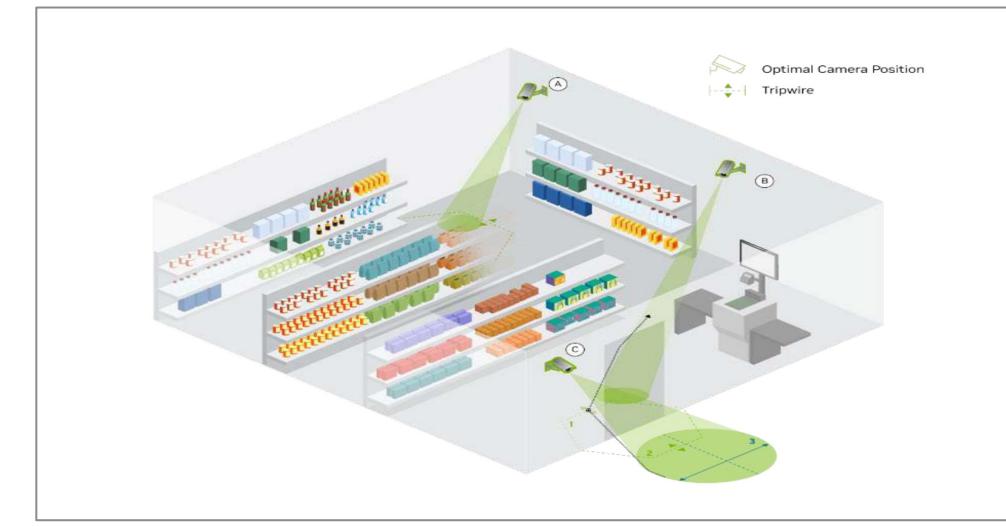


[AI NVR Demo](#)



[Gen AI Demo](#)

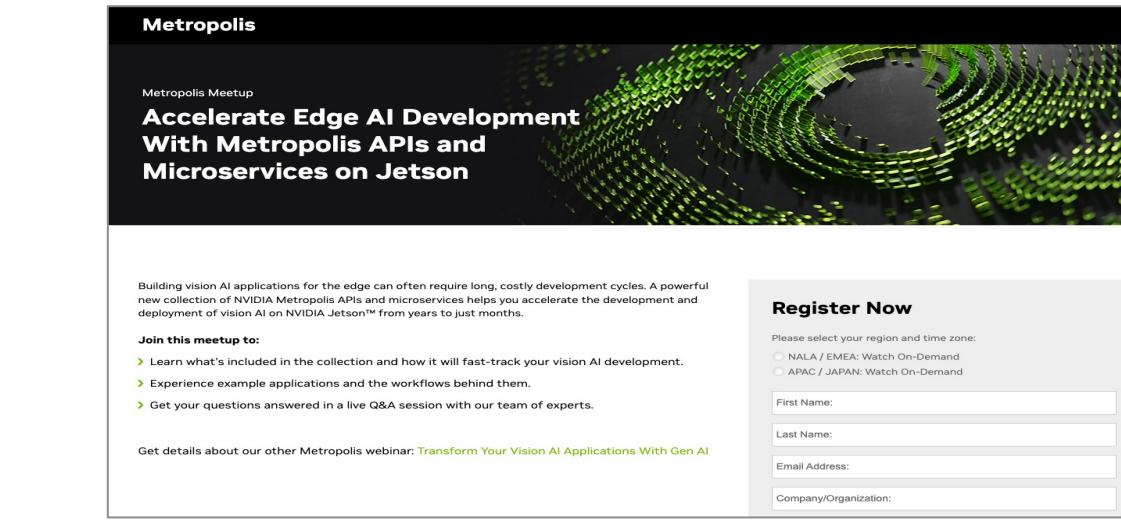
## Assets



[Whitepaper](#)



[Solution One-Pager](#)



[Webinar – Part 1](#)



[Webinar – Part 2](#)

# The In-Person GTC Experience Is Back

Come to GTC—the conference for the era of AI—to connect with a dream team of industry luminaries, developers, researchers, and business experts shaping what's next in AI and accelerated computing.

From the highly anticipated keynote by NVIDIA CEO Jensen Huang to over 600 inspiring sessions, 200+ exhibits, and tons of networking events, GTC delivers something for every technical level and interest area.

Be sure to save your spot for this transformative event. You can even take advantage of early-bird pricing when you register by February 7.

**March 18-21, 2024 | [www.nvidia.com/gtc](http://www.nvidia.com/gtc)**

