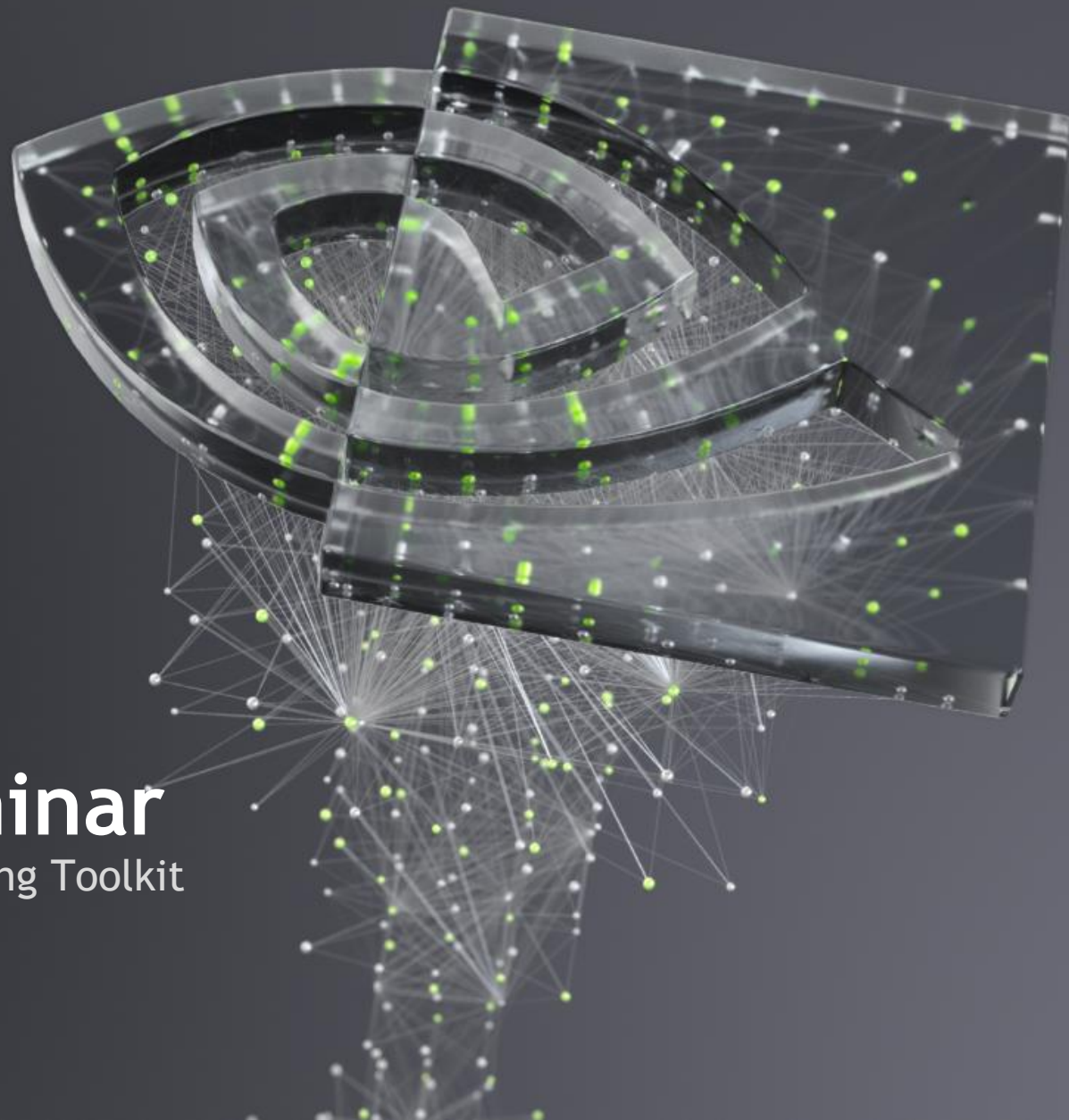




# NVIDIA Hands-on Seminar

Deepstream SDK, TensorRT, Transfer Learning Toolkit

Jonghwan Lee, Deep Learning Data Scientist, NVIDIA





# AGENDA

## NVIDIA Inference Technology

TensorRT/ Deepstream SDK/ Transfer Learning Toolkit

---

## NVIDIA TensorRT

High-performance deep learning inference platform

---

## NVIDIA Deepstream SDK

AI-powered Intelligent Video Analytics framework

---

## Hands-on with Jetson NX

Mainly Deepstream SDK, TensorRT and TLT

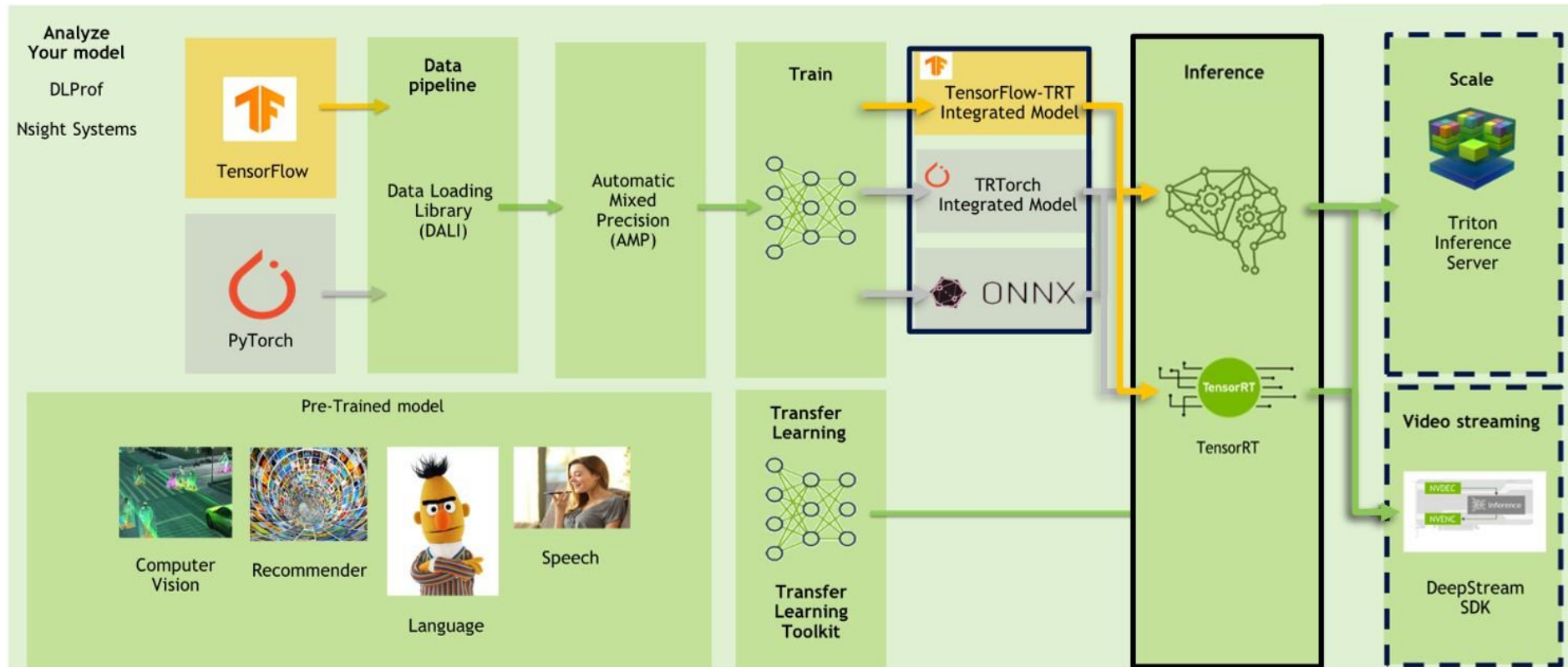
---



NVIDIA INFERENCE  
TECHNOLOGY

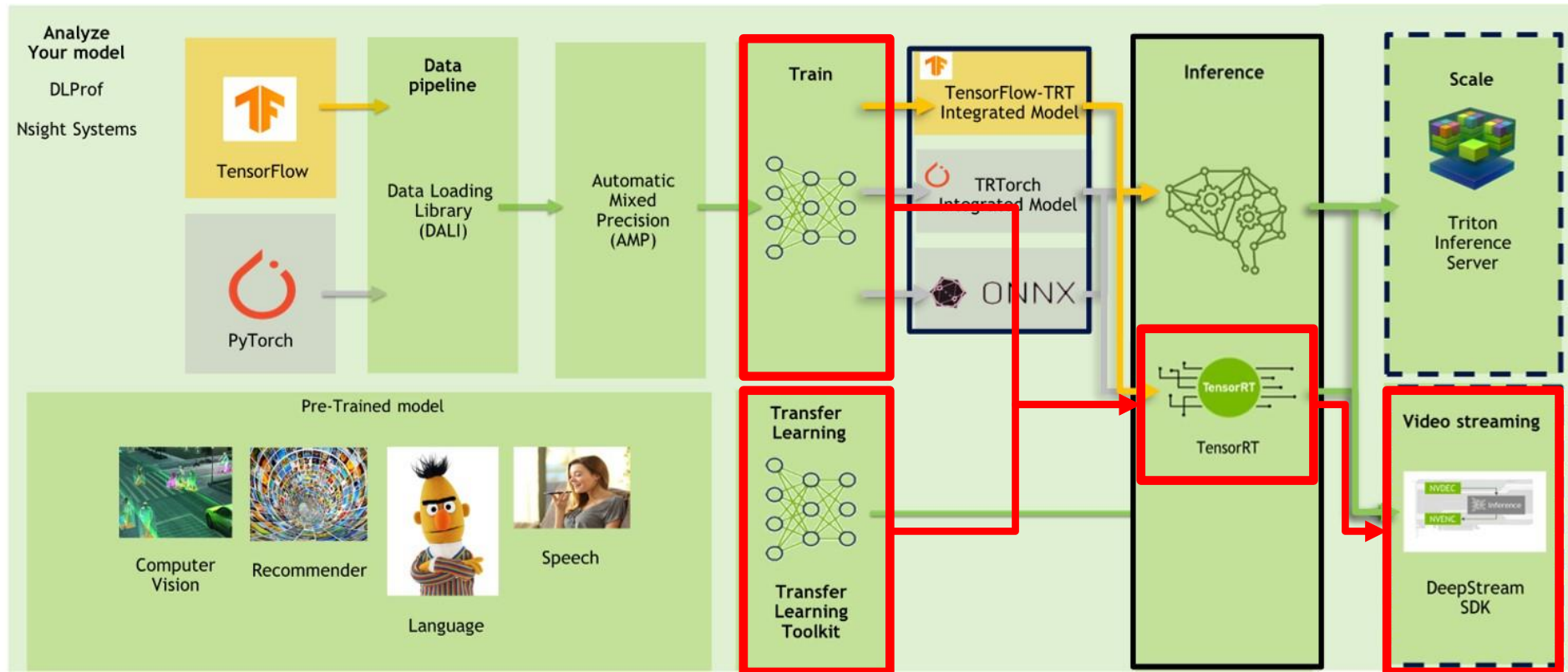
# NVIDIA INFERENCE SOLUTIONS

Easily deploy AI, Tools and Libraries to accelerate your AI model and Service



# NVIDIA INFERENCING SOLUTIONS

Easily deploy AI, Tools and Libraries to accelerate your AI model and Service

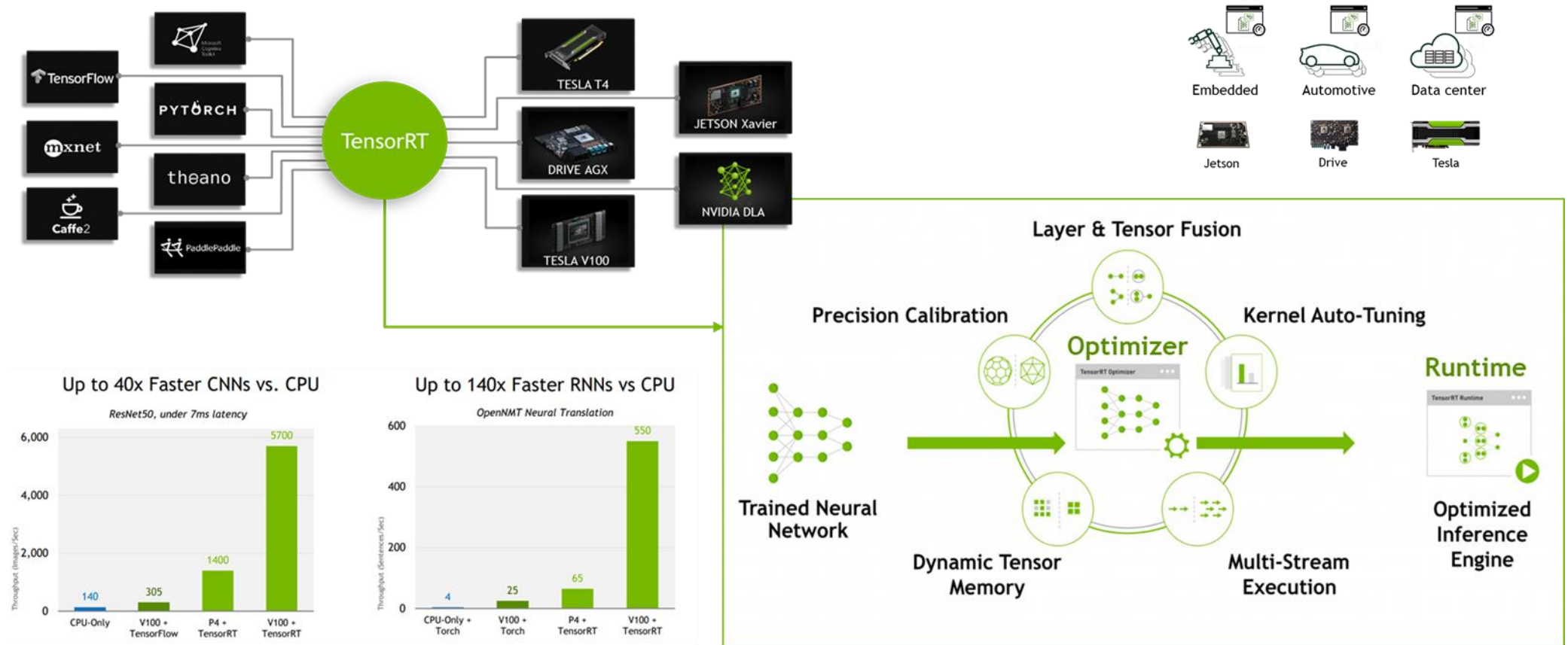


Today's focus!!!



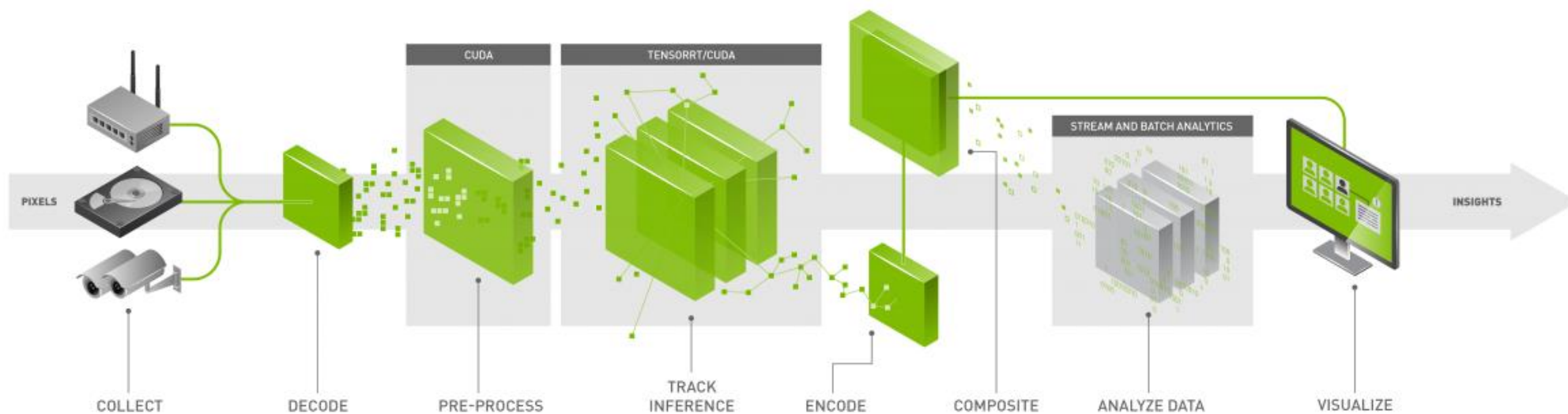
# NVIDIA TensorRT

From Every Framework, Optimized For each Target platform



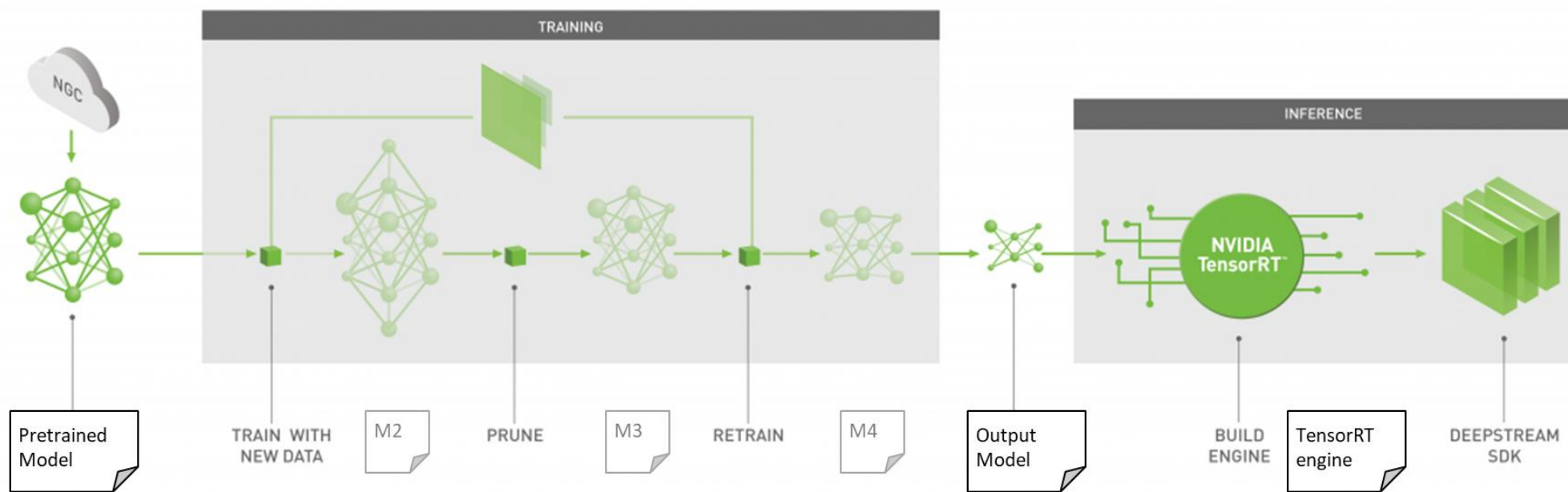
# NVIDIA Deepstream SDK

Analyze Data from Cameras, Sensors and IoT Gateways in Real-Time



# NVIDIA Transfer Learning Toolkit

## Accelerates Throughput on Leading Industry Platforms



1. Download docker container -> 2. Pull Model -> 3. Train with your data -> 4. Prune -> 5. Retrain -> 6. Export

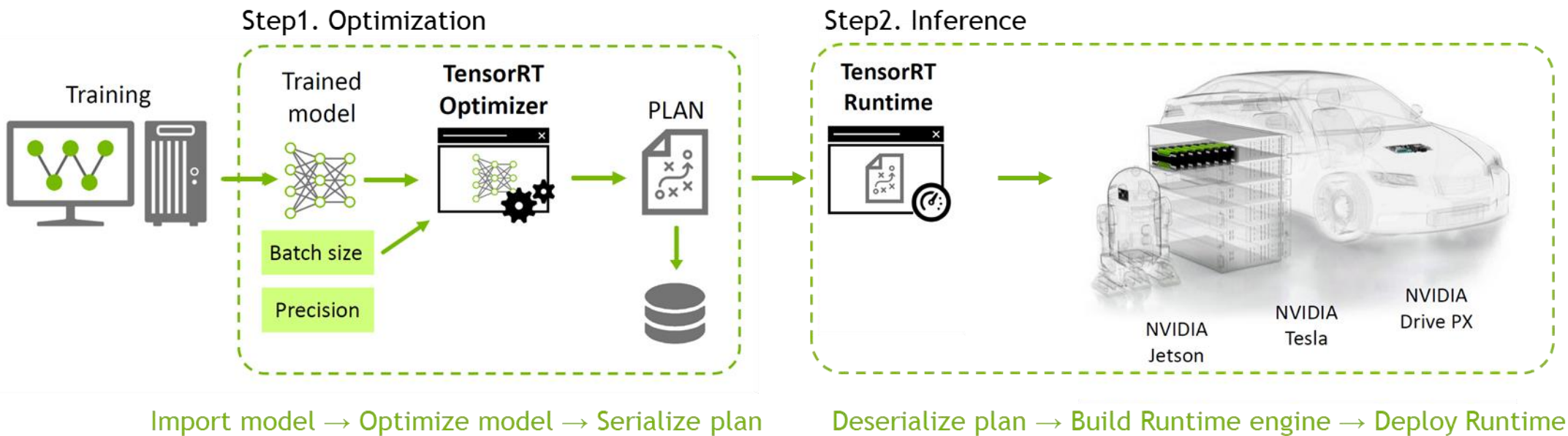
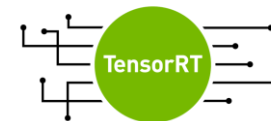


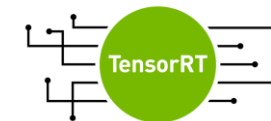


NVIDIA TensorRT

# TENSORRT WORKFLOW

How to accelerate DNN inference with TensorRT





# NVIDIA TENSOR CORE

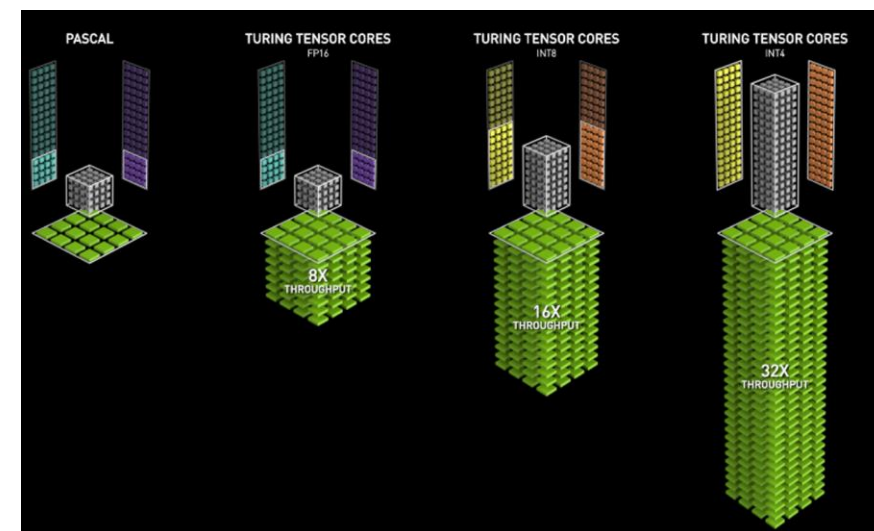
Mixed precision Matrix-Multiply-and-Add

Ampere has new 3<sup>rd</sup> TensorCore

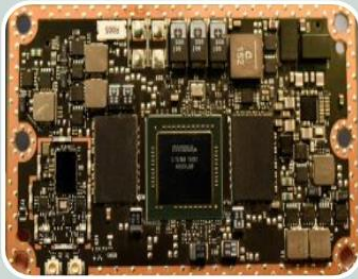
- ▶ Available in Volta and Turing architecture GPUs
- ▶ 125 Tflops in FP16 vs 15.7 Tflops in FP32 (8x speed-up)
- ▶ Optimized 4x4x4 dot operation (GEMM)

$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32                      FP16                      FP16                      FP16 or FP32

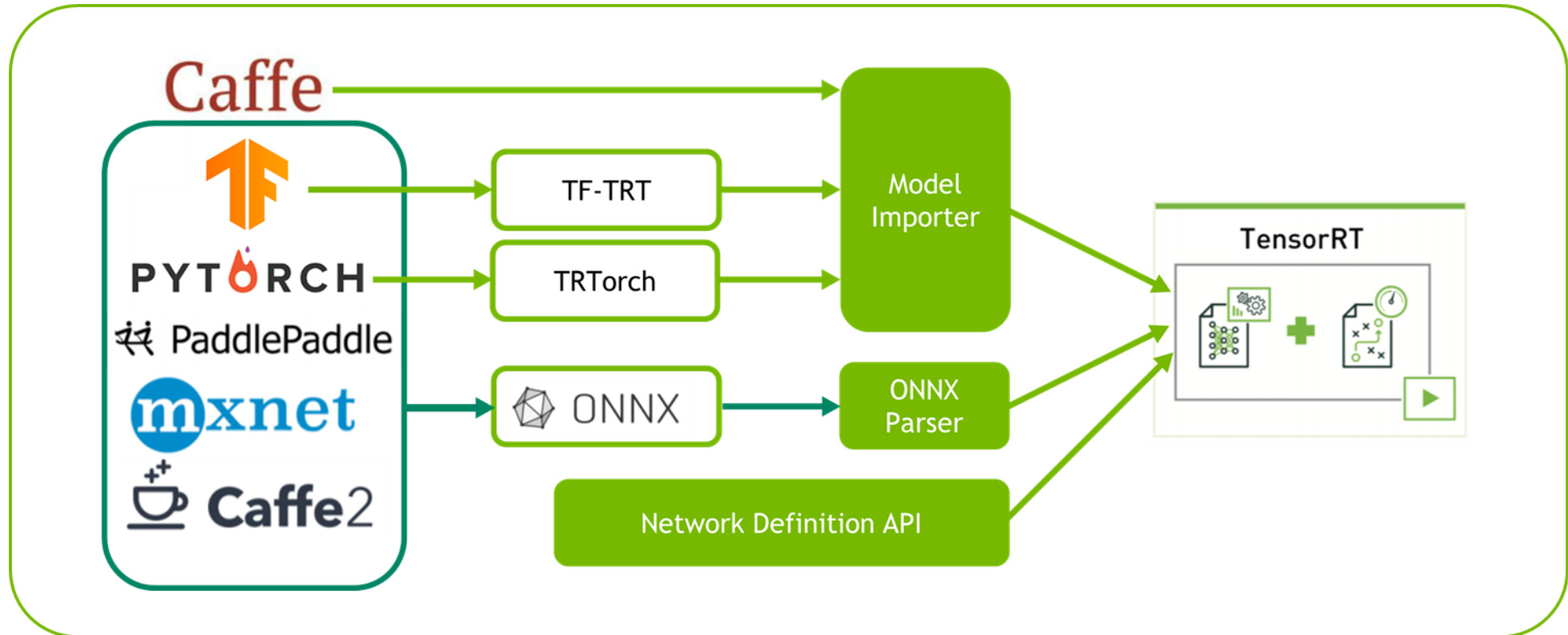
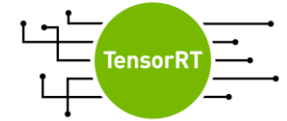


# NVIDIA Deep Learning Accelerator

[illegible]

# NVIDIA TENSORRT

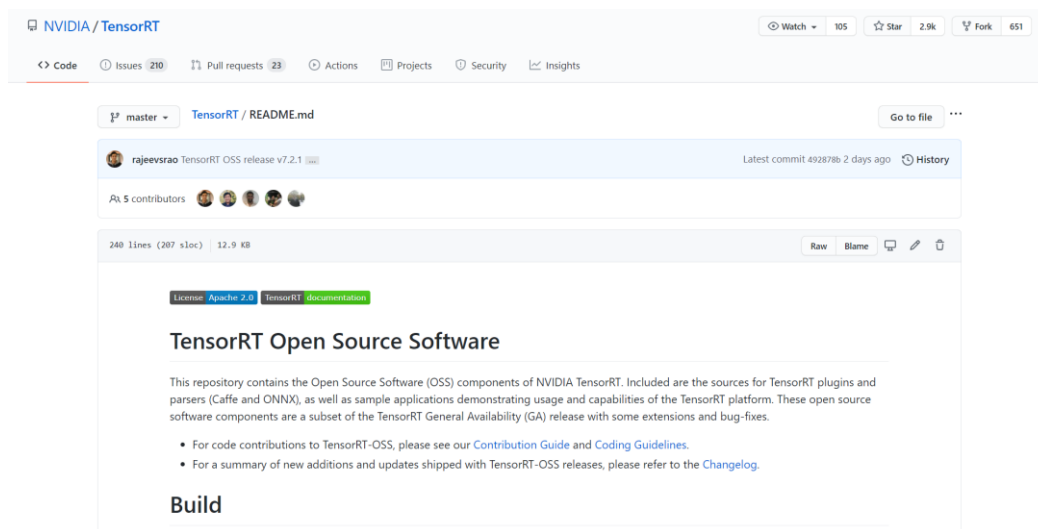
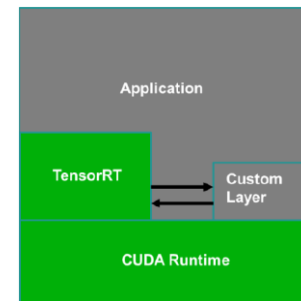
## TensorRT Model import flow



# NVIDIA TENSORRT OSS

Supports custom plug-ins by CUDA

- ▶ Allow users to express and provide implementations of novel layers
  - TensorRT provides APIs and dedicated implementations for most common layers
  - Use the custom layer API for infrequent or more innovative layers or your own confident implementation
  - Register custom implementations via a callback mechanism
  - Can be used in conjunction with reduced precision optimization



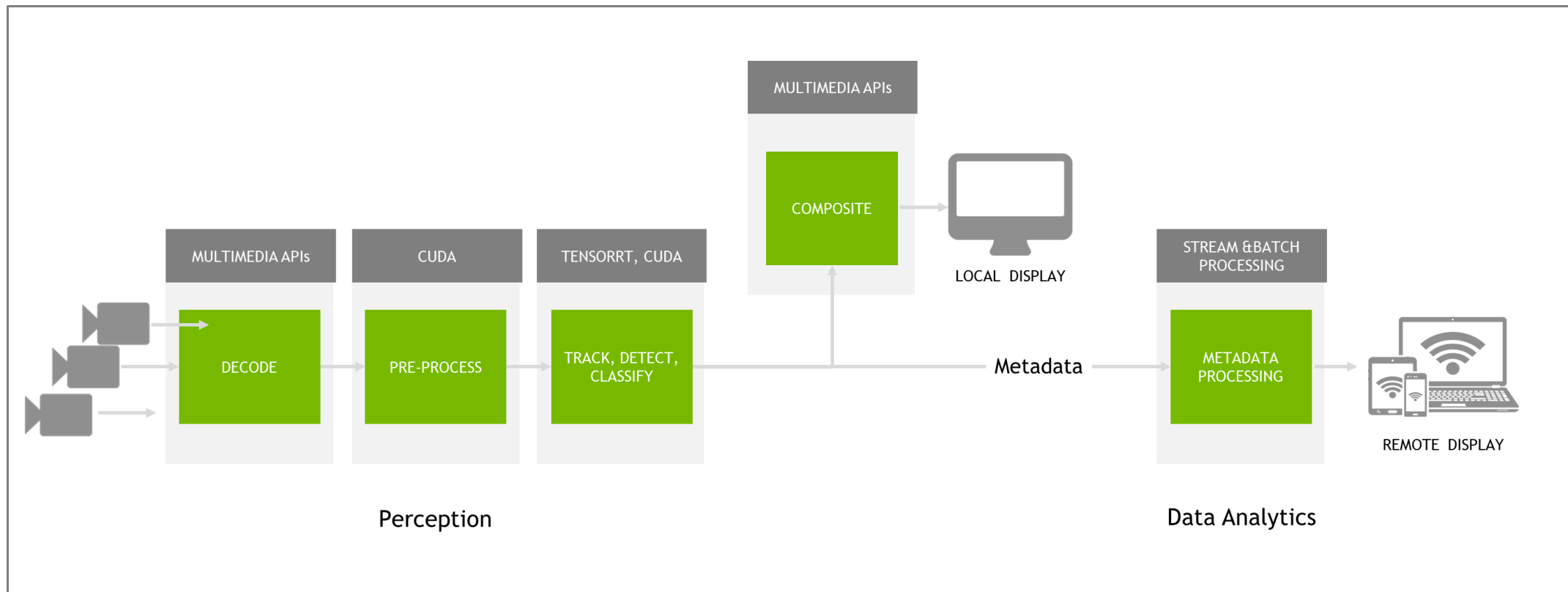




NVIDIA Deepstream SDK

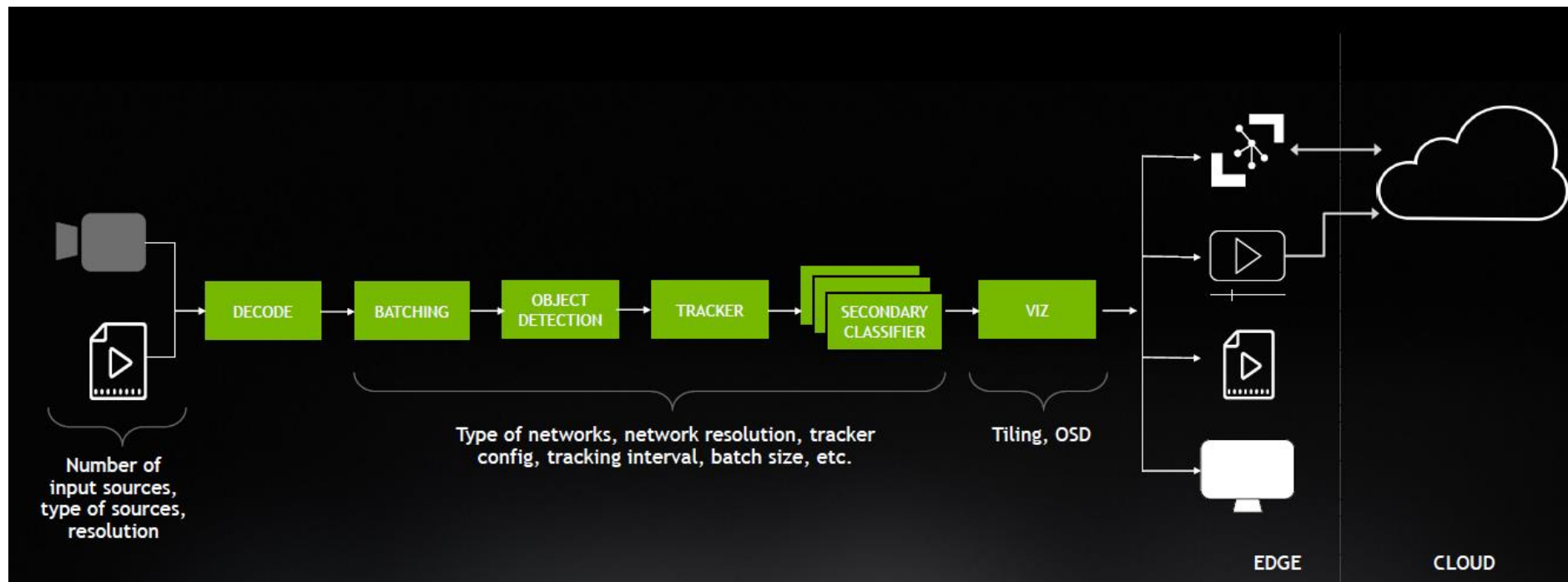
# GENERAL INTELLIGENCE VIDEO ANALYSIS

Framework for analyzing video



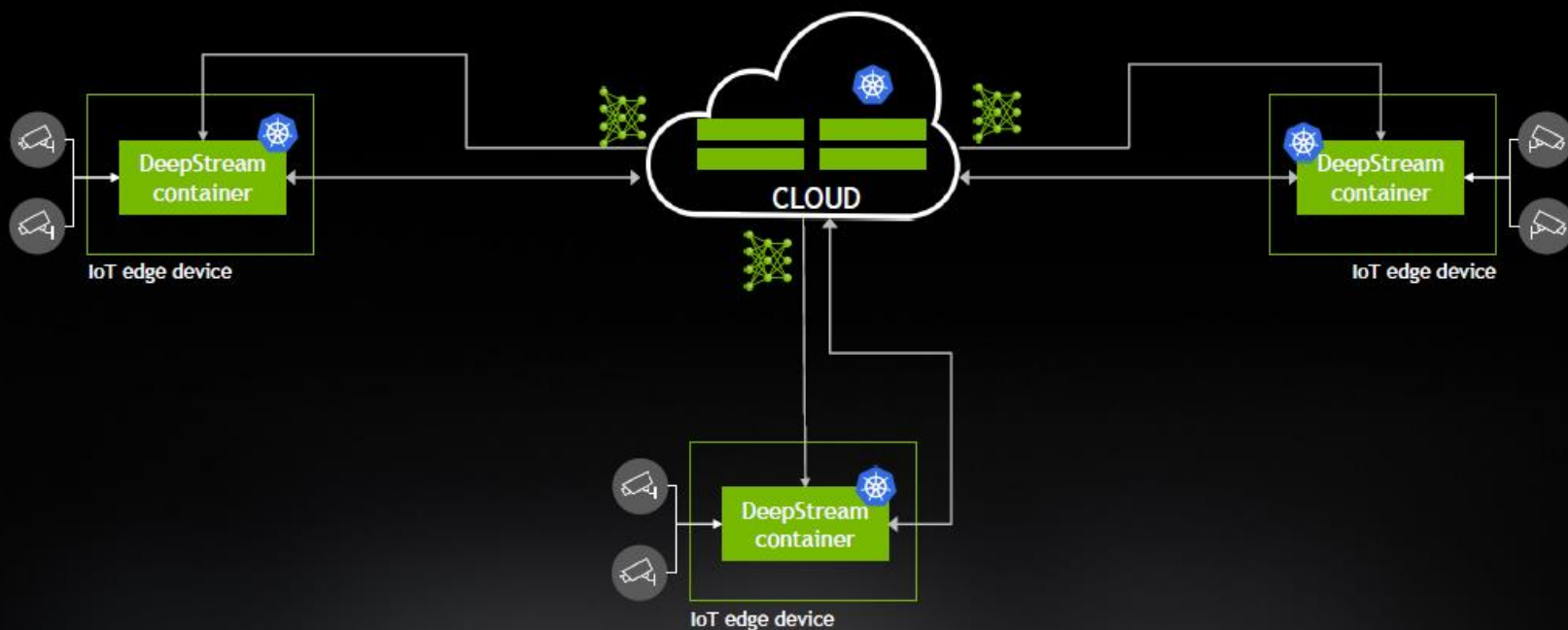
# DEEPSTREAM SDK

## End-to-End Deepstream Application - deepstream-test5



# DEEPSTREAM SDK

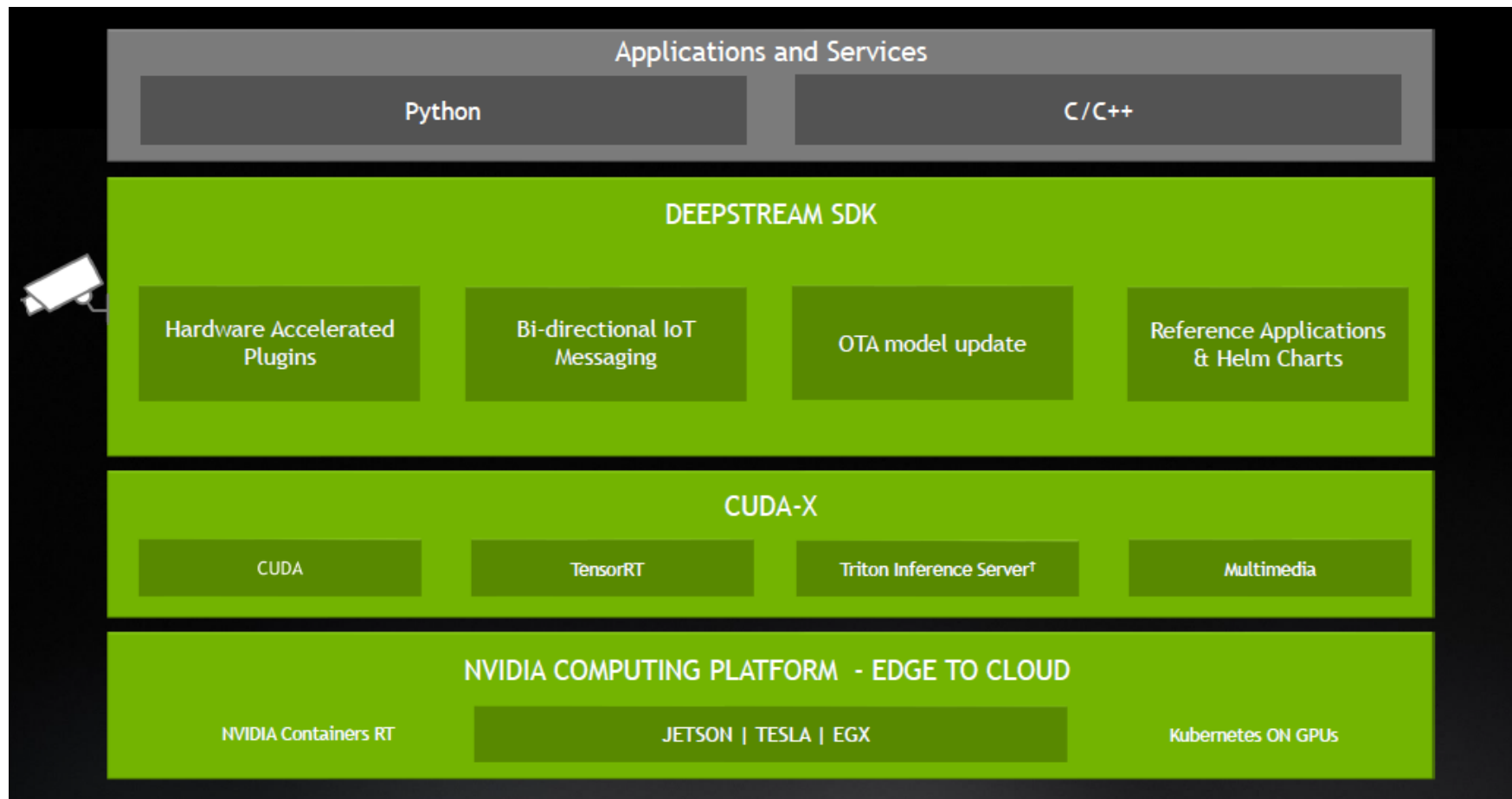
At scale deployment



Download the [Demo](#) to deploy DeepStream application using Helm charts and Kubernetes on NGC

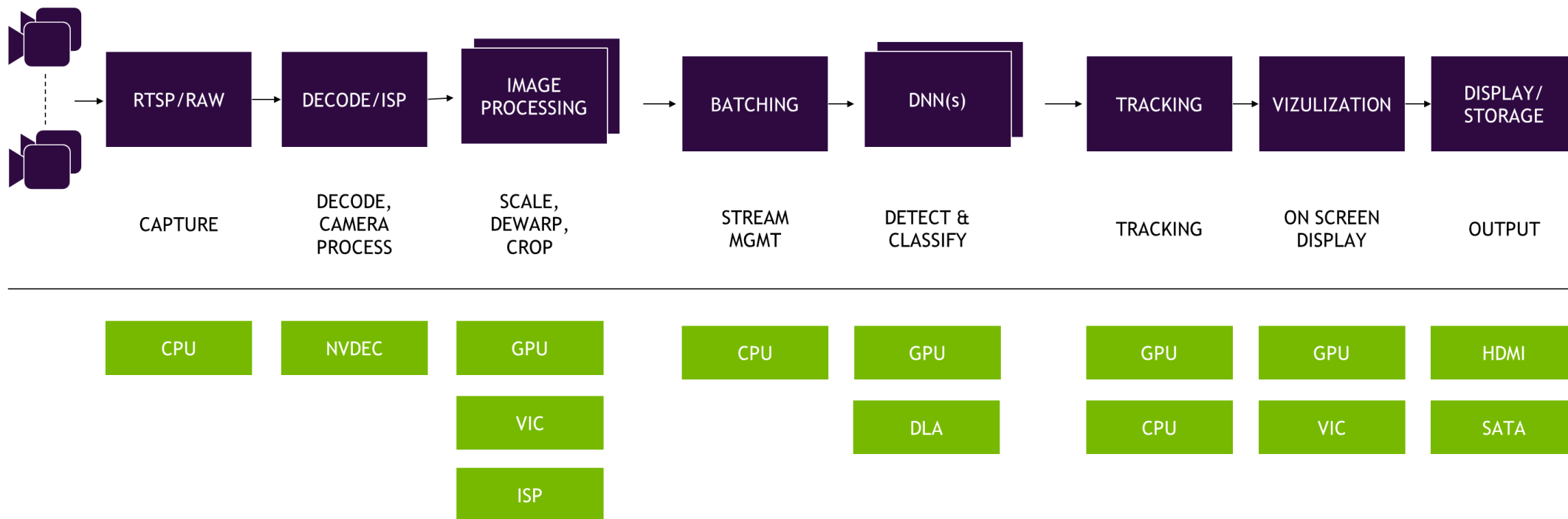
# WHAT IS DEEPSTREAM?

AI-powered Intelligent Video Analytics framework



# DEEPSTREAM GRAPH ARCHITECTURE

AI-powered Intelligent Video Analytics framework





# CONCEPTS OF DEEPSTREAM SDK

## Gstreamer Foundations

The **DeepStream SDK** is based on the open source [GStreamer multimedia framework](#). There are a few key concepts in GStreamer that we need to touch on before getting started. These include Elements, Pads, Buffers, and Caps. We will be describing them at a high level, but encourage those who are interested in the details to read the [GStreamer](#) documentation to learn more.



# CONCEPTS OF DEEPSTREAM SDK

## Example of Gstreamer pipeline

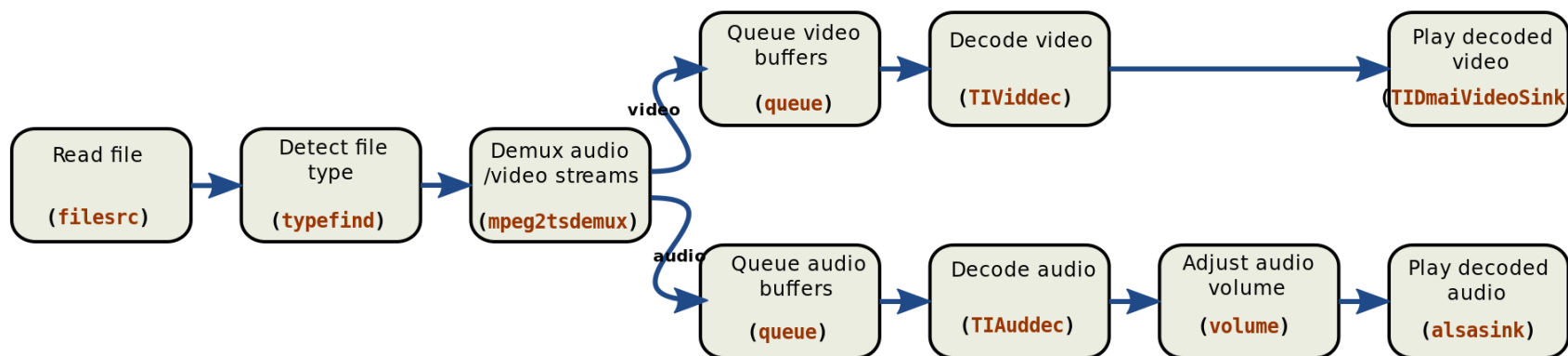
The goals of GStreamer are to separate the application (e.g. Video player, Video editor, etc.) from the streaming media complexity (e.g. hardware acceleration, remoteness)

GStreamer – streaming media

D-Bus – inter process communication

Use **gst-launch** command to create the GStreamer pipeline

```
gst-rtsp-server - v4l2src ! video/x-raw,width=1280,height=720 ! omxh264enc ! video/x-h264,profile=baseline ! h264parse config-interval=1 ! rtph264pay na
```

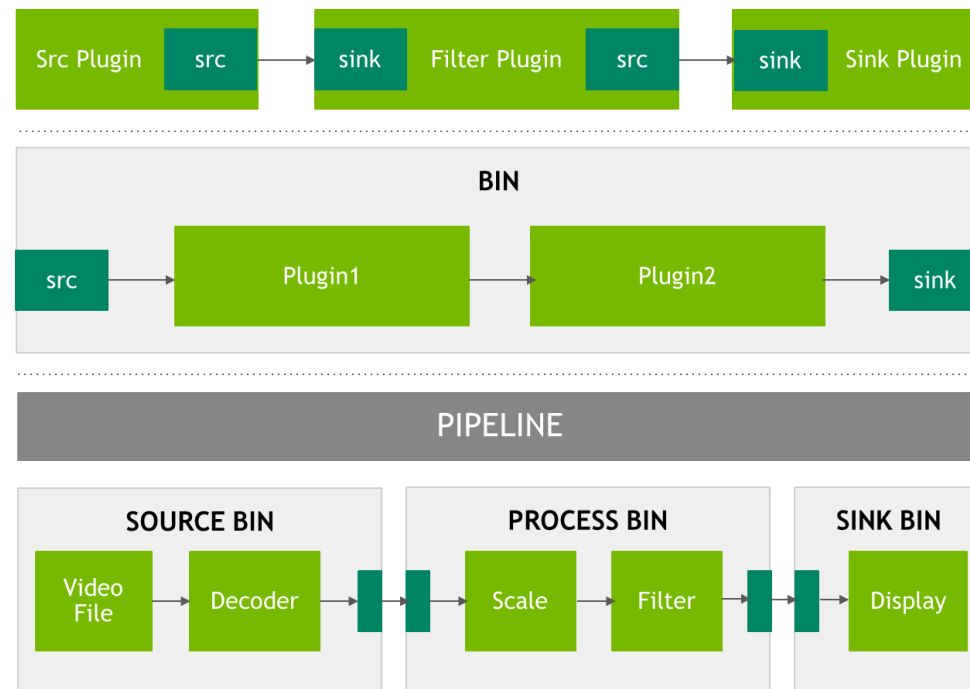


```
gst-launch filesrc location="video.ts" ! typefind ! mpeg2tsdemux name=demux \
demux. ! 'video/x-h264' ! queue ! TIViddec ! TIDmaVideoSink \
demux. ! 'audio/mpeg' ! queue ! TIAuddec ! volume volume=5 ! alsasink
```

# CONCEPTS OF DEEPSTREAM SDK

## Gstreamer pipeline

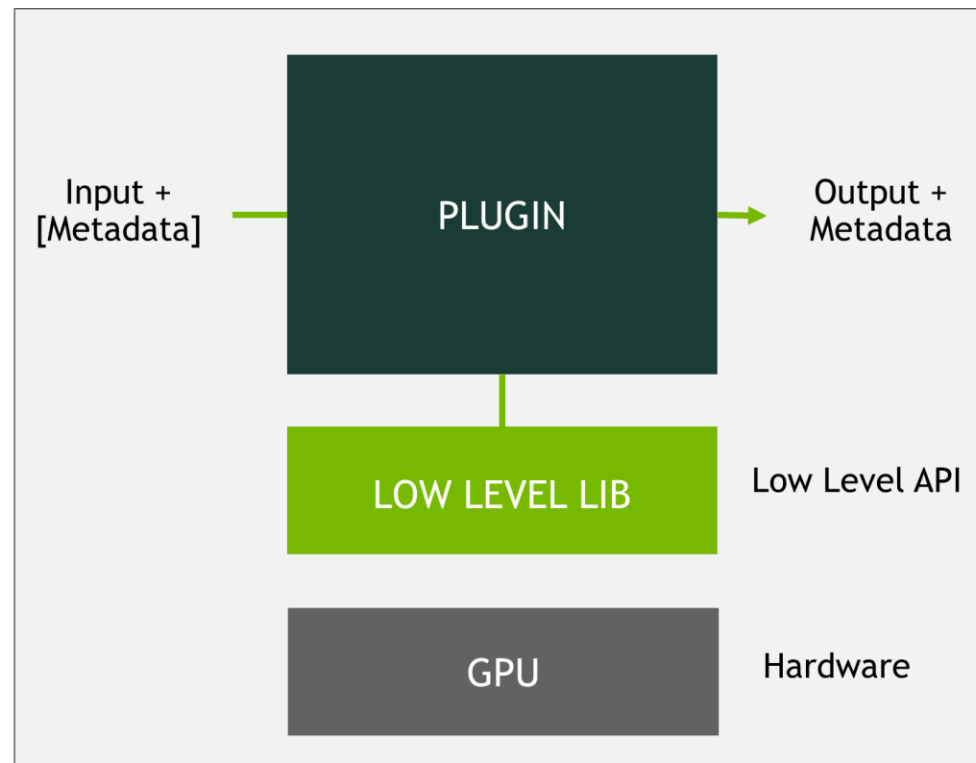
Level	Component	Function
1	PLUGINS	Basic building block connected through PADs
2	BINS	A container for a collection of plugins
3	PIPELINE	Top level bin providing a bus and managing the synchronization



# CONCEPTS OF DEEPSTREAM SDK

## Deepstream Building block

- A plugin model-based pipeline architecture
- Graph-based pipeline interface to allow high-level component interconnect
- Heterogenous processing on GPU and CPU
- Hides parallelization and synchronization under the hood
- Inherently multi-threaded



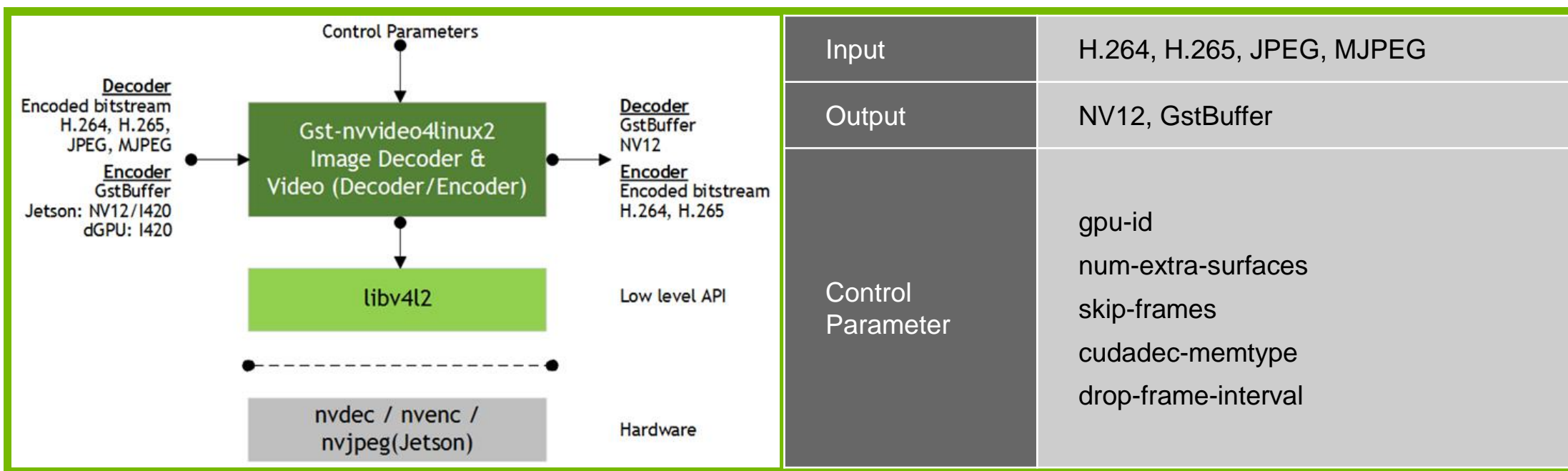
# DEEPSTREAM SDK

## Deepstream accelerated plugin

Plugin Name	Functionality
Gst-nvinfer	DL inference for detection, classification and segmentation with TRT engine
Gst-nvinferserver	DL inference for detection, classification and segmentation with any other native framework
Gst-nvtracker	Reference object trackers; KLT, IOU, NvDCF
Gst-nvstreammux	Stream aggregation, multiplexing and batching
Gst-nvstreamdemux	Demux batched frames into individual buffers
Gst-nvmultistreamtiler	Composites a 2D tile from batched buffers
Gst-nvdsosd	Draw boxes and text overlay
Gst-nvvideoconvert	Scaling, format conversion and rotation
Gst-nvdewarper	Dewarping for fish-eye degree cameras
Gst-nvof	Fast optical flow calculation on special HW accelerator for it
Gst-nvofvisual	Visualize input optical flow vector on frames
Gst-nvsegvisual	Visualize segmentation results on frames
Gst-nvvideo4linux2	Hardware accelerated decode and encode(NVDEC / NVENC)
Gst-nvjpegdec	Decoding JPEG images
Gst-nvmsgconv	Metadata generation
Gst-nvmsgbroker	Messaging to cloud
Gst-nvdsanalytics	ROI filtering, Overcrowding Detection, Direction Detection and Line crossing

# DEEPSTREAM SDK

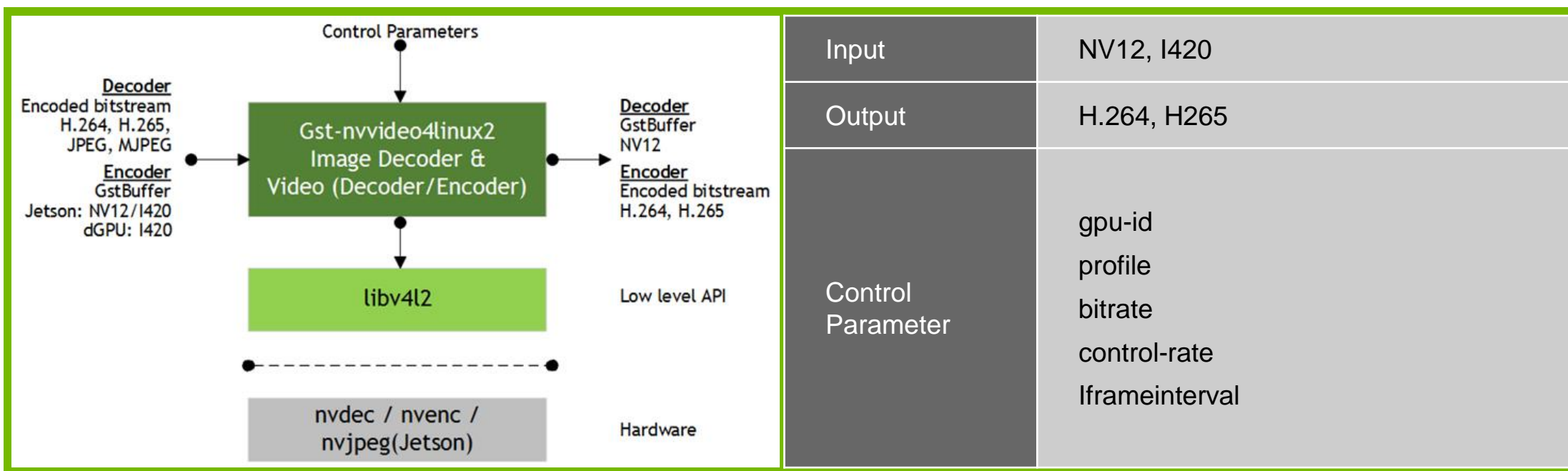
## Deepstream accelerated plugin - Gst-nvvideo4linux2 (Decode)





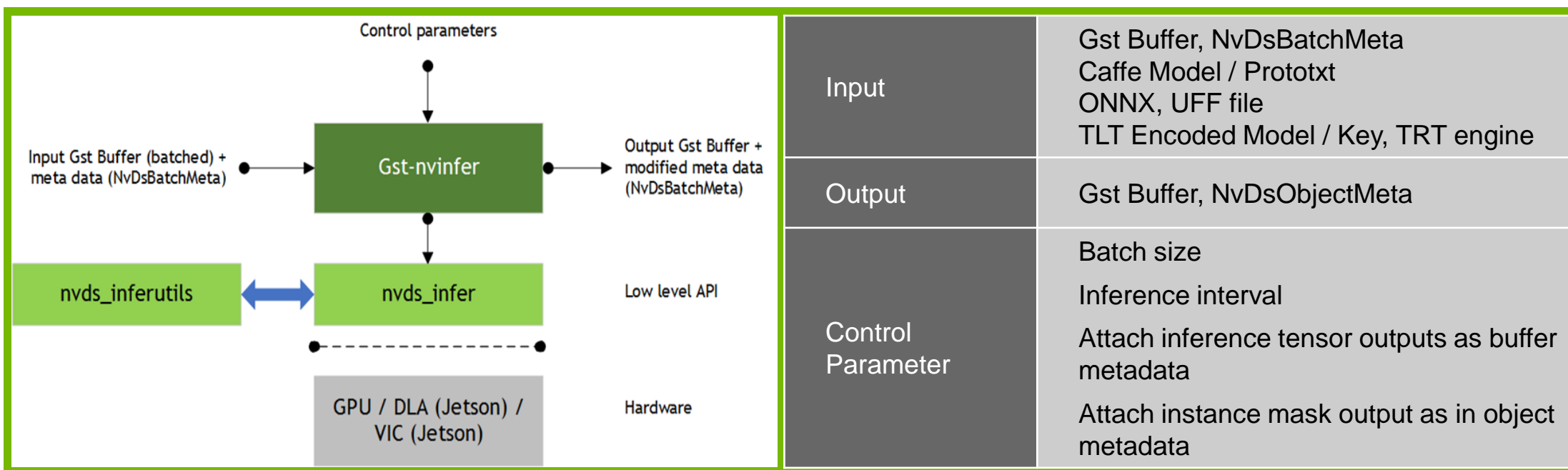
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvvideo4linux2 (Encode)



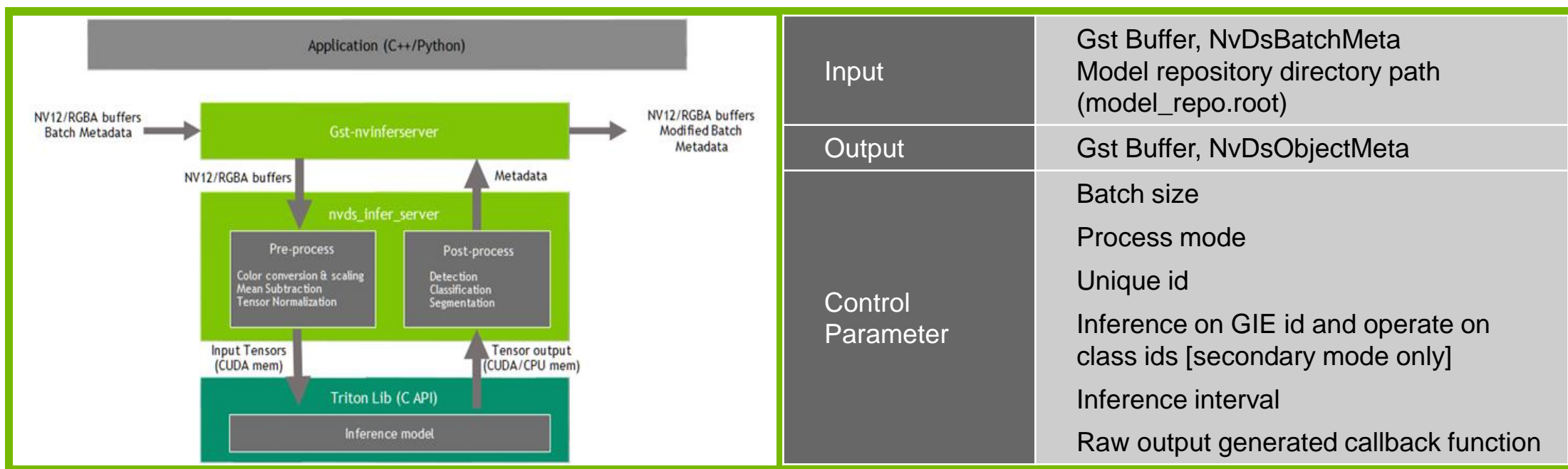
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvinfer



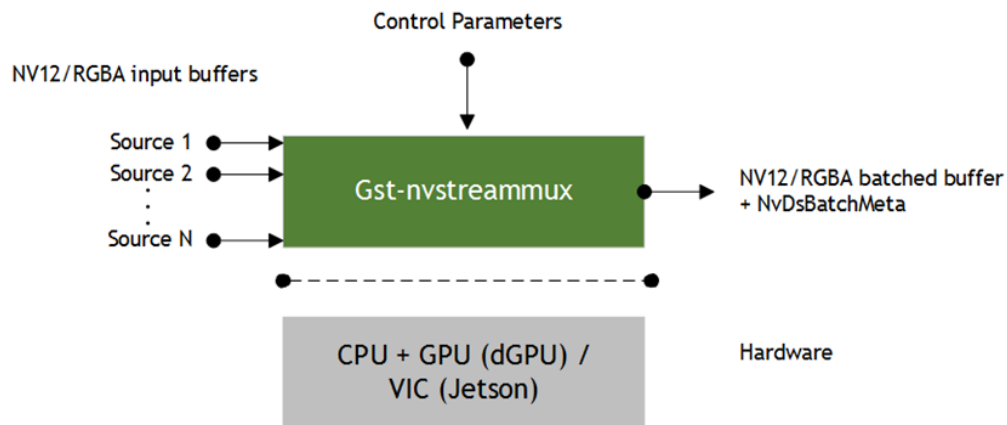
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvinferserver



# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvstreammux



Input

NV12/RGBA buffers

Output

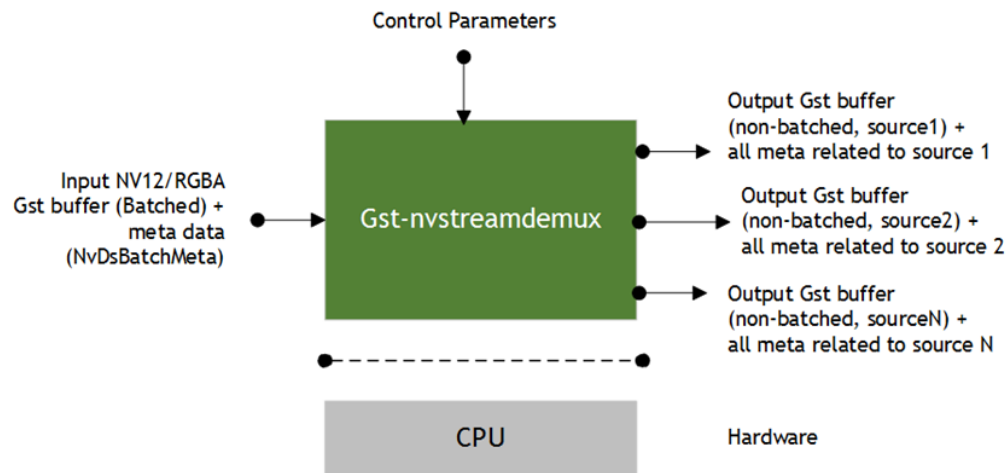
NV12/RGBA batched buffer  
GstNvBatchMeta

Control  
Parameter

batch-size, batched-push-timeout  
Width, height, enable-padding  
gpu-id, live-source, nvbuf-memory-type  
num-surfaces-per-frame  
buffer-pool-size  
attach-sys-ts

# DEEPSTREAM SDK

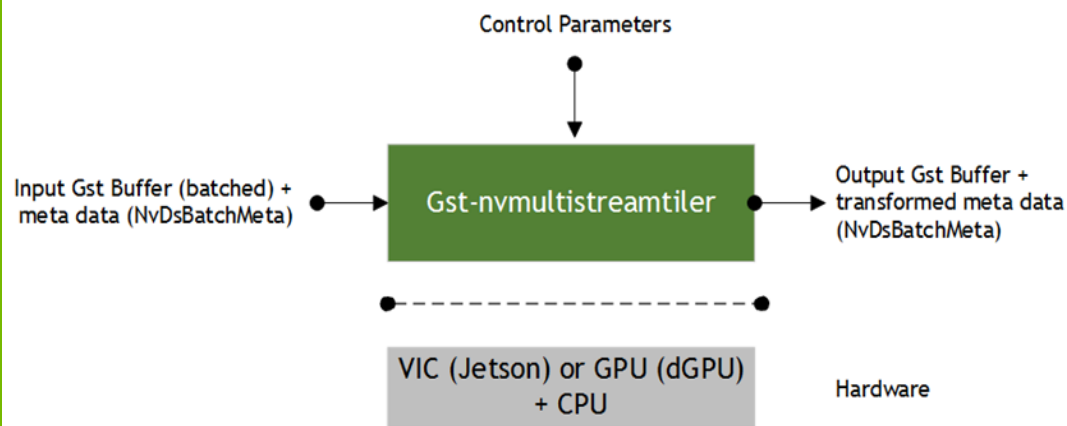
## Deepstream accelerated plugin - Gst-nvstreamdemux



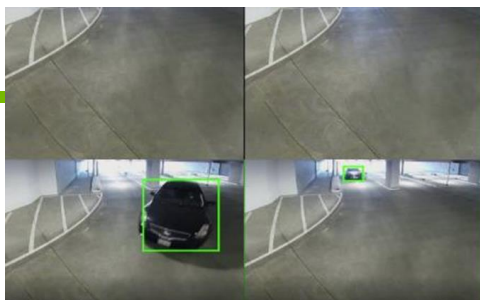
Input	Gst Buffer(batched), NvDsBatchMeta
Output	Gst Buffer (non-batched, single source) Meta related to each Gst Buffer source
Control Parameter	N/A

# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvmultistreamtiler



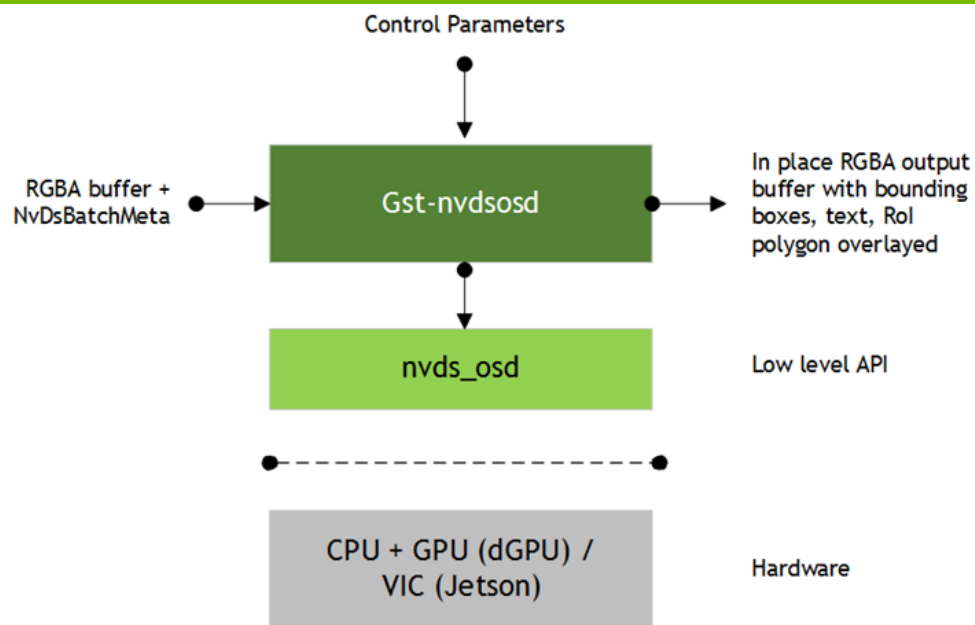
Input	Gst Buffer batched buffer NvDsBatchMeta
Output	Gst Buffer with composited input frames Transformed NvDsBatchMeta
Control Parameter	rows, columns width, height gpu-id (dGPU only) show-source nvbuf-memory-type custom-tile-config





# DEEPSTREAM SDK

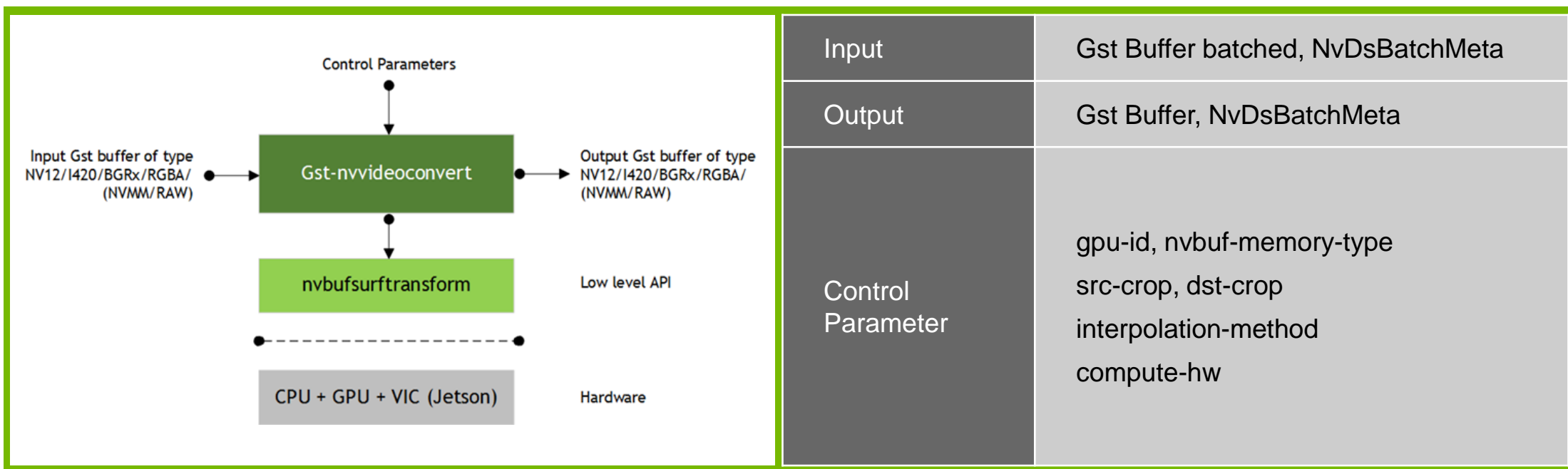
## Deepstream accelerated plugin - Gst-nvdsosd



Input	RGBA buffer, NvDsBatchMeta
Output	RGBA buffer modified in place to overlay
Control Parameter	gpu-id, display-clock, display-text clock-font, clock-font-size x-clock-offset, y-clock-offset clock-color, process-mode hw-blend-color-attr, display-bbox, display-mask

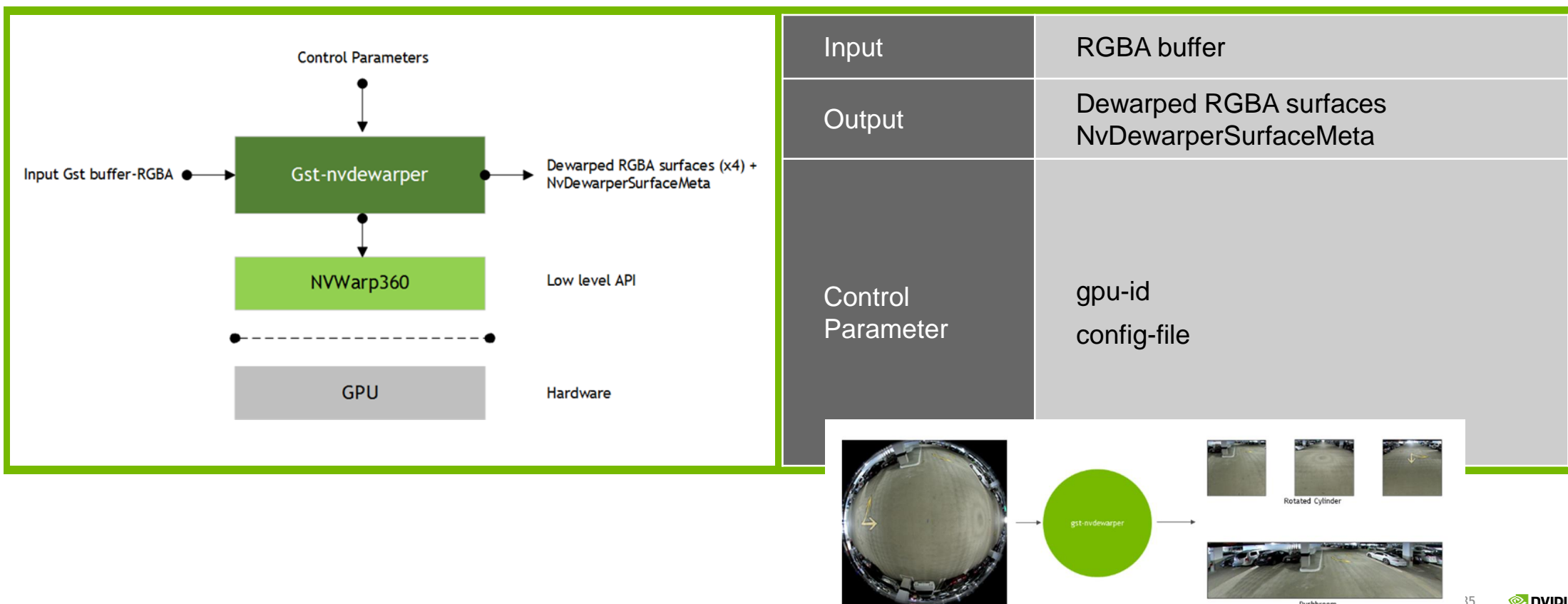
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvvideoconvert



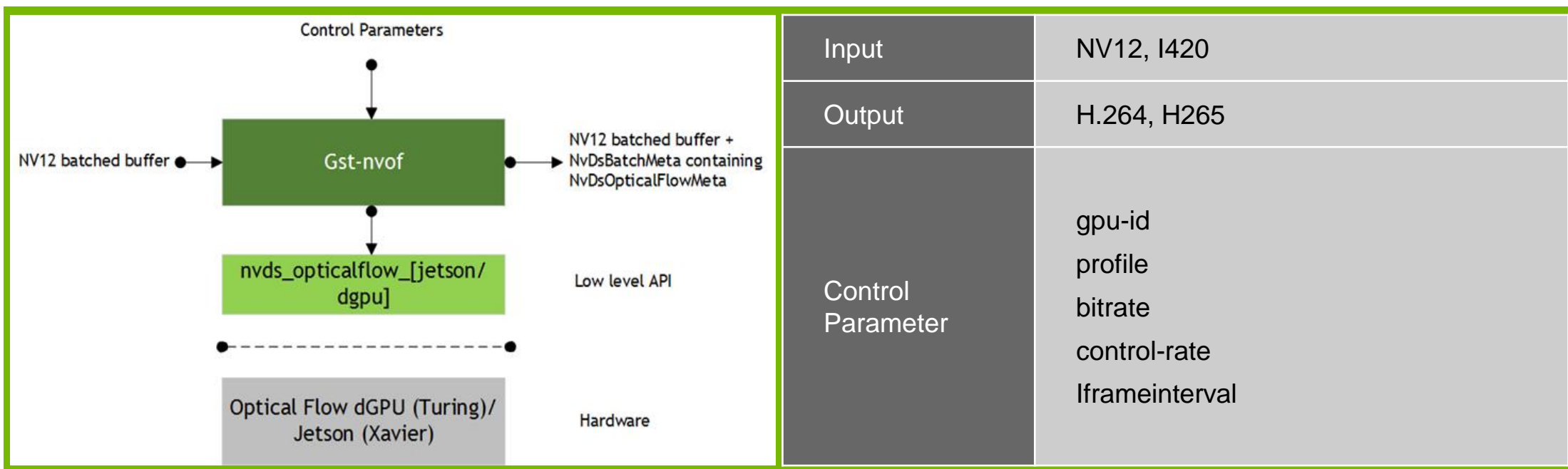
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvdewarper



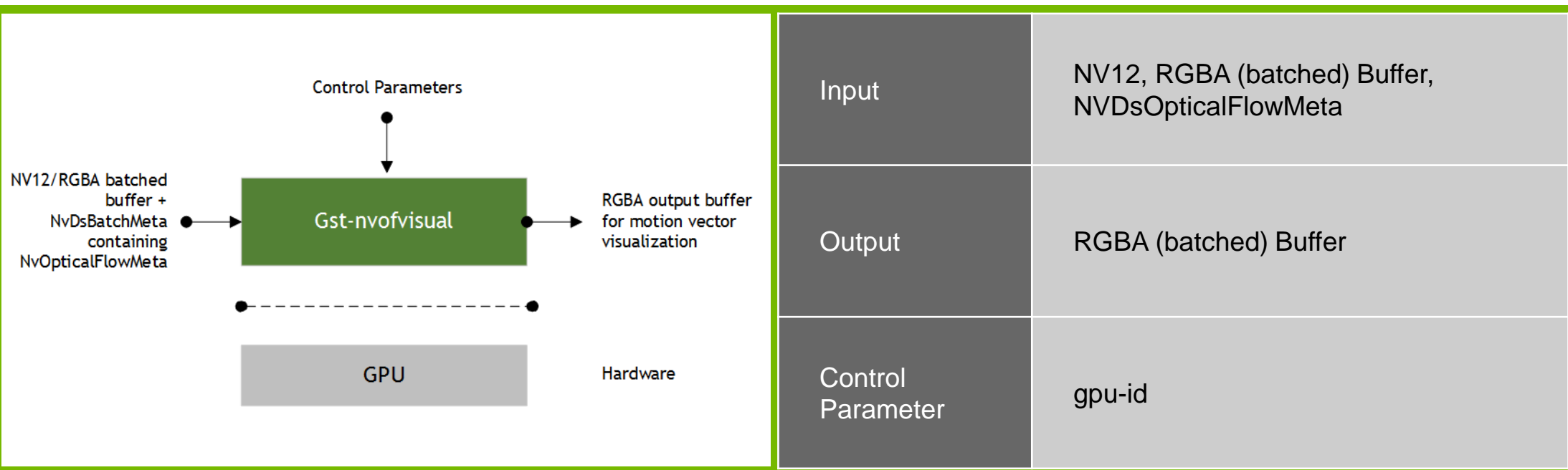
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvof



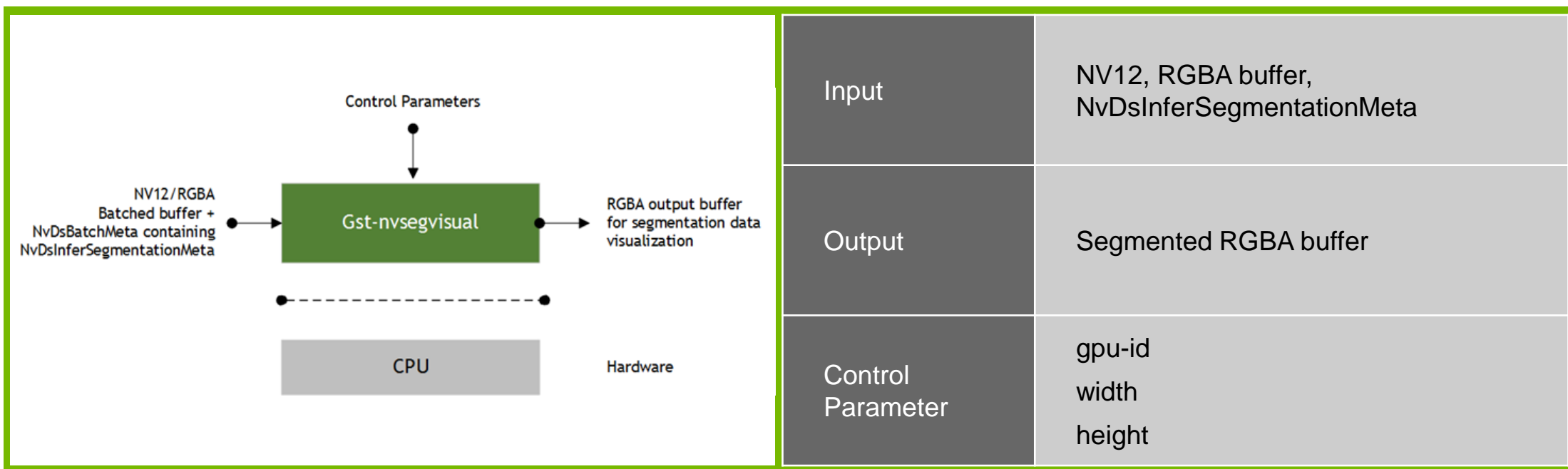
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvofvisual



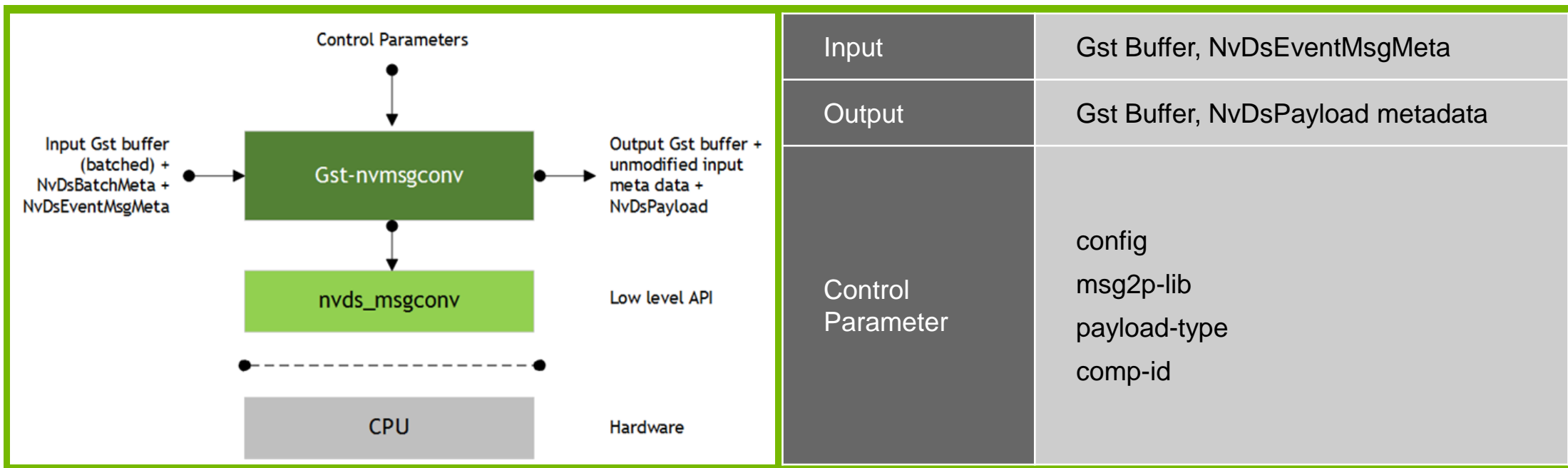
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvsegvisual



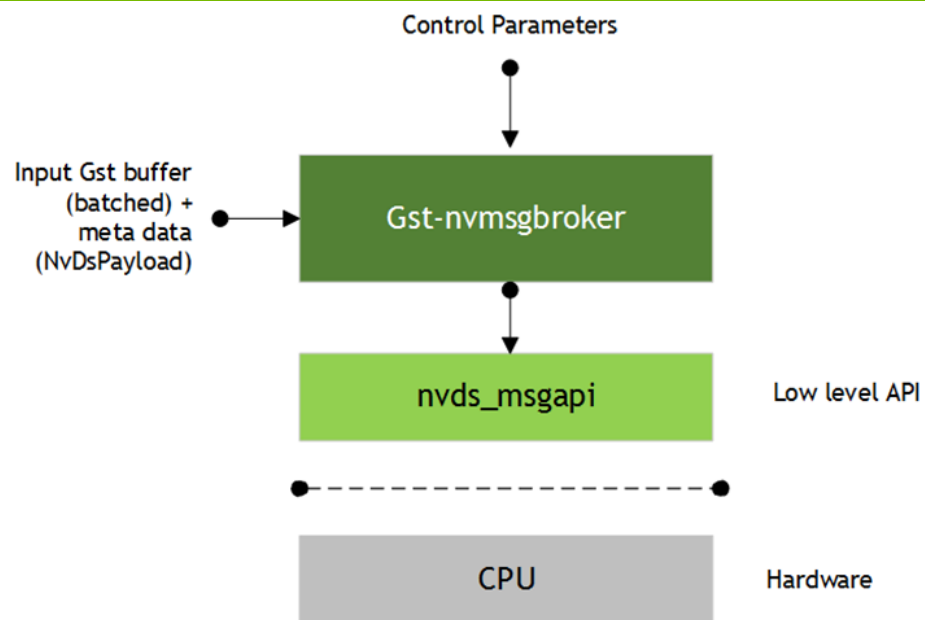
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvmsgconv



# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvmsgbroker

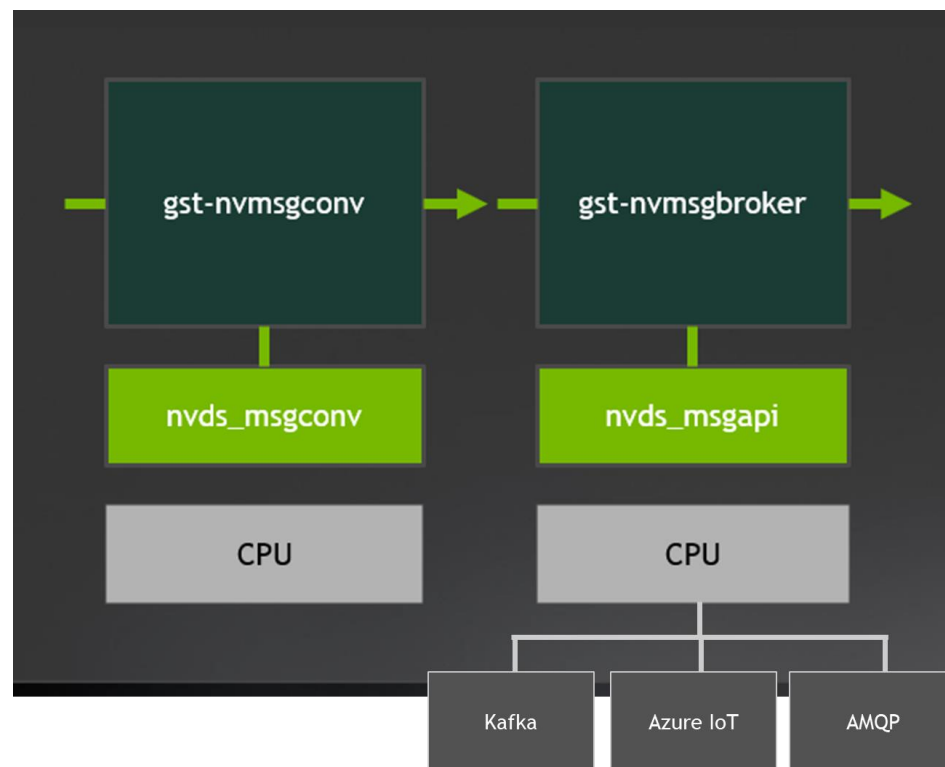


Input	Gst Buffer, NvDsPayload
Output	N/A
Control Parameter	Config conn-str proto-lib comp-id Topic



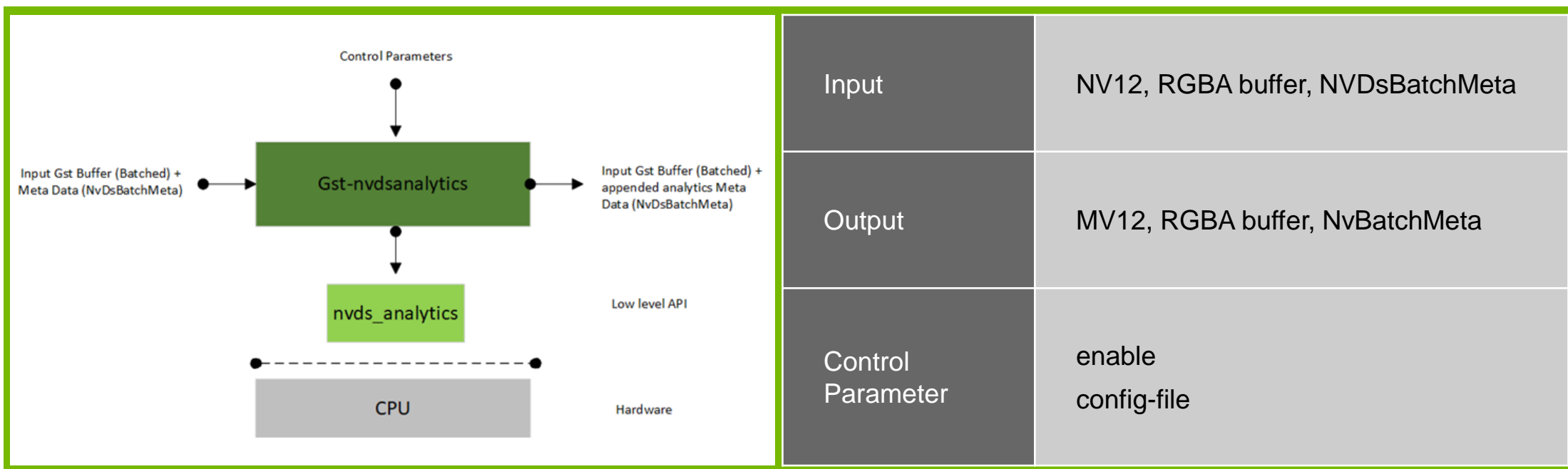
# DEEPSTREAM SDK

Deepstream accelerated plugin - Gst-nvmsgconv & Gst-nvmsgbroker




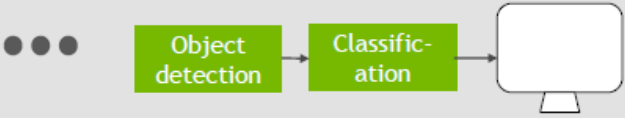

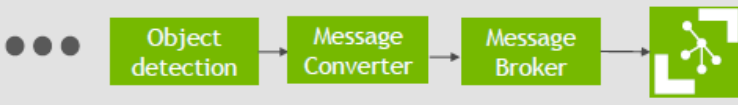
# DEEPSTREAM SDK

## Deepstream accelerated plugin - Gst-nvdsanalytics



# DEEPSTREAM SDK

Get started applications (available both C and Python)

Name	Function	
deepstream-test1	DeepStream Hello world. Single video from file to on screen display with bounding box	
deepstream-test2	Builds on test1 and adds secondary object classification on detected objects	
deepstream-test3	Builds on test1 and adds multiple video inputs	
deepstream-test4	Builds on test1 and adds connections to IoT services thru the nvmsgbroker plugin	

Native-C apps: [sources/apps/sample\\_apps/](https://github.com/NVIDIA-AI-IOT/deepstream_python_apps)

Python apps: [https://github.com/NVIDIA-AI-IOT/deepstream\\_python\\_apps](https://github.com/NVIDIA-AI-IOT/deepstream_python_apps)



# HANDS-ON FOR DEEPSTREAM

# DEEPSTREAM SDK HANDS-ON

## Setting up for the Hands-on

- Prerequisites

- **Host machine**

- > SSH terminal
    - > VLC player
    - > (For TLT exercise) GPU machine is required & docker runtime environment
    - > `$ docker run --gpus all -it -v $(pwd):/workspace -w /usr/local/src -p 8888:8888 cycoslee/nv-deepstream-sdk-handson:tlt_host_210127`

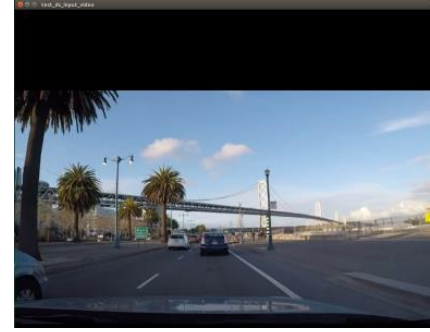
- **Jetson NX (based on Jetpack 4.4.1)**

- > Docker container pull
    - > `$ sudo docker run -ti --runtime=nvidia --rm --net=host -e DISPLAY=$DISPLAY -w /opt/nvidia/deepstream/deepstream-5.0 --device /dev/video0 -v /tmp/.X11-unix:/tmp/.X11-unix -v $(pwd):/workspace cycoslee/nv-deepstream-sdk-handson:ds_nx_210127`

- **Assets**

- > IP camera(supports RTSP)
    - > USB camera(webcam)

# HOW TO CONTROL INPUTS



1\_test\_ds\_input\_video



filesrc

h264parse

nvv4l2decoder

nvstreammux

nveglglessink



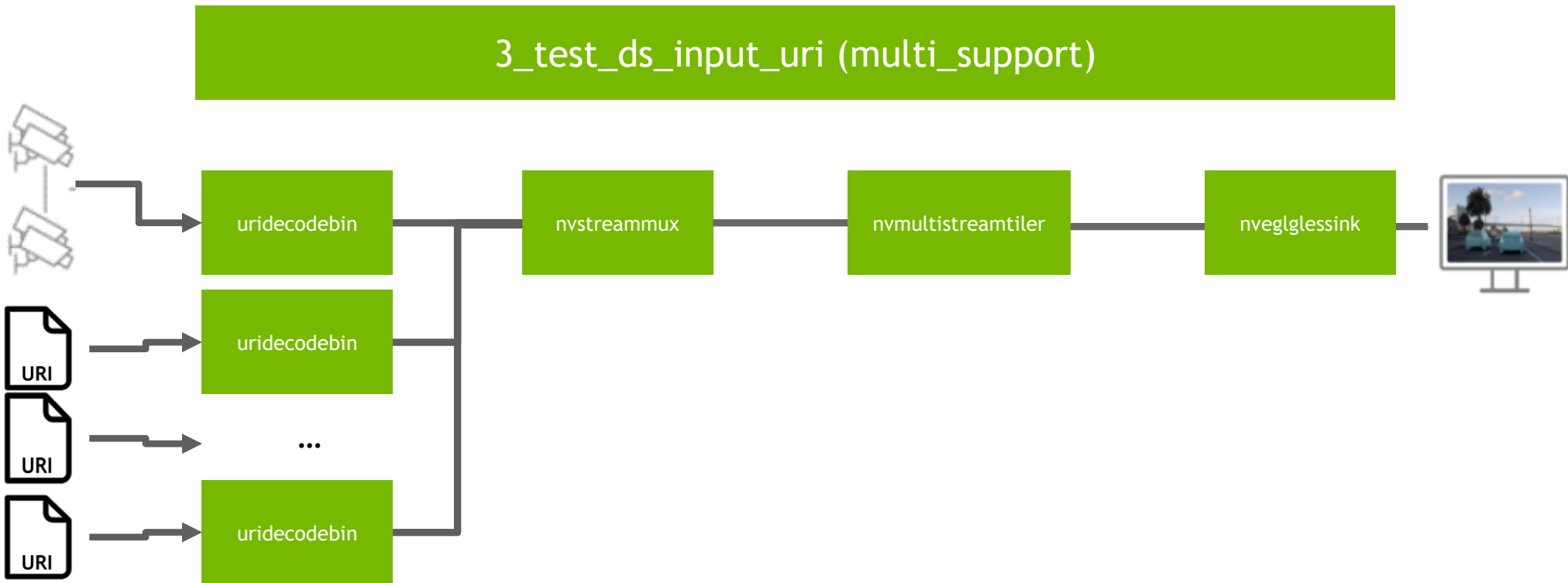
# HOW TO CONTROL INPUTS



2\_test\_ds\_input\_webcam

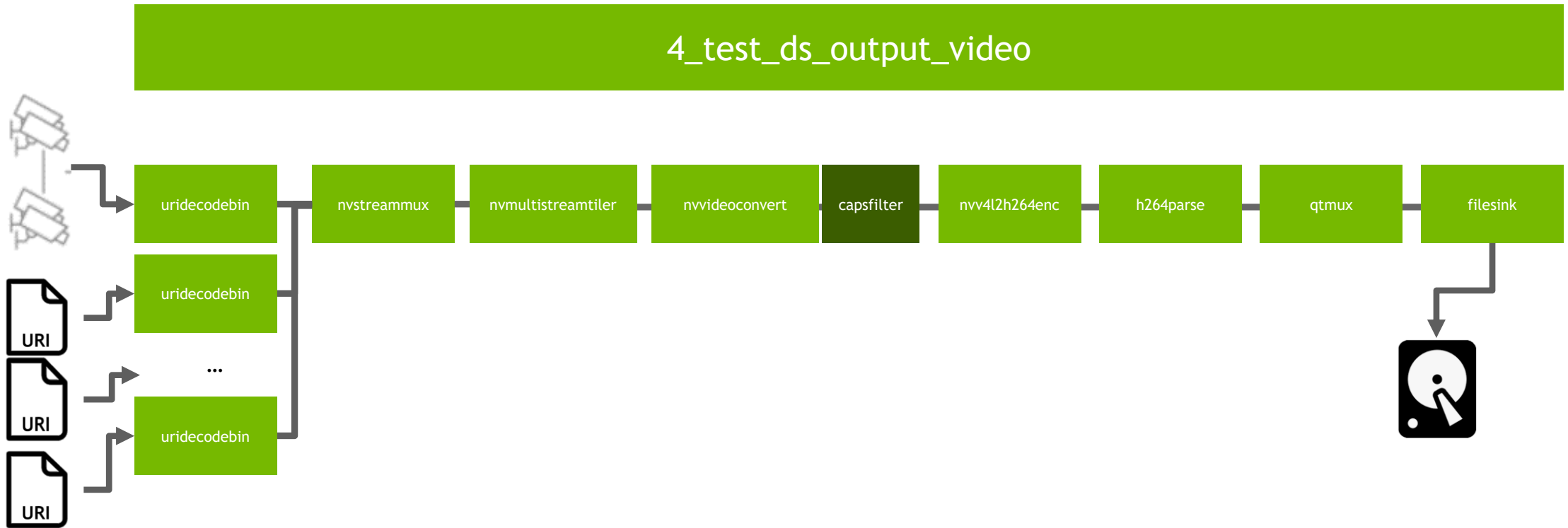


# HOW TO CONTROL INPUTS





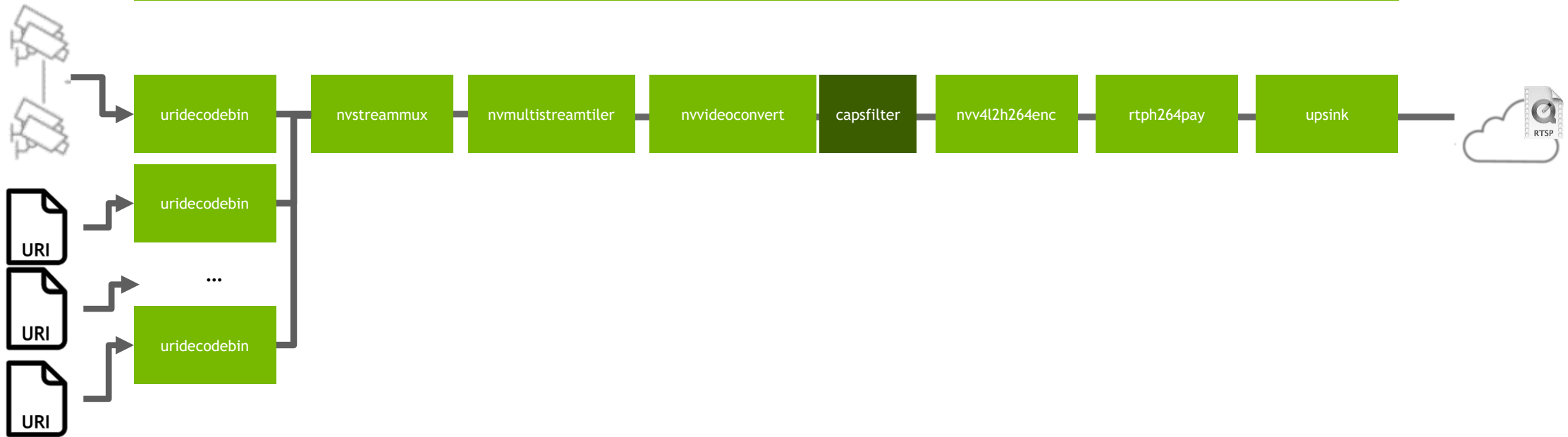
# HOW TO CONTROL OUTPUTS



# HOW TO CONTROL OUTPUTS



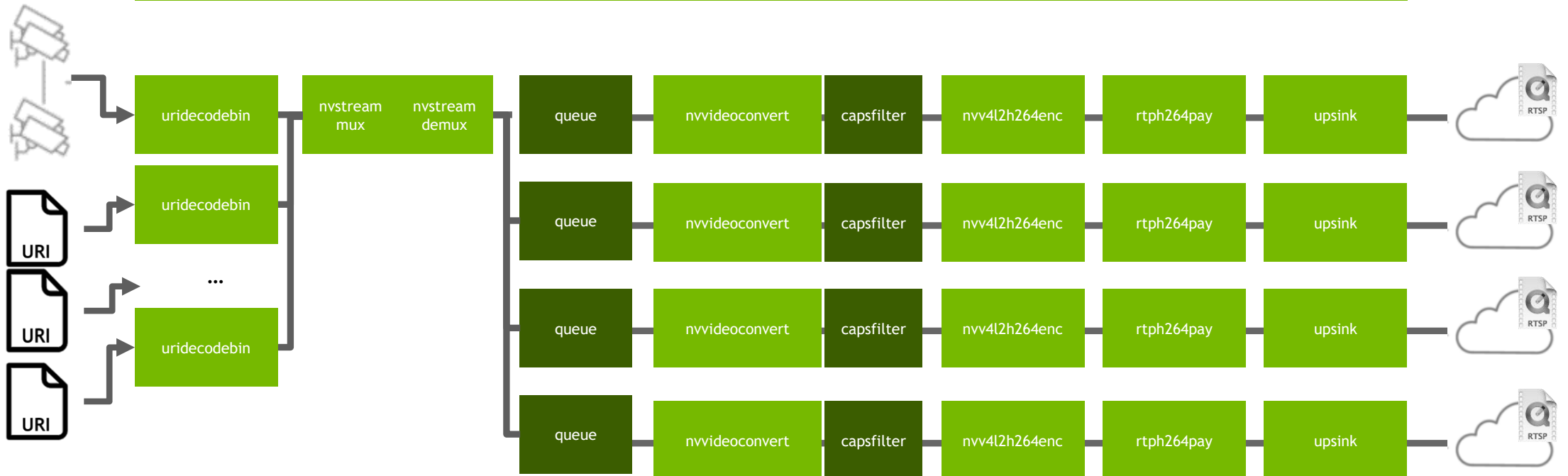
5\_test\_ds\_output\_rtsp



# HOW TO CONTROL OUTPUTS



6\_test\_ds\_output\_video



# HOW TO DEPLOY DL

How to use TensorRT easily by using trtexec

	Model Name	FPS
0	inception_v4	317.924917
1	super_resolution_bsd500	152.118646
2	unet-segmentation	148.245192
3	yolov3-tiny-416	556.410331
4	ResNet50_224x224	866.609455
5	ssd-mobilenet-v1	887.422532

## 7\_test\_trt\_benchmark

### Benchmarks Targeted for Jetson Xavier NX (Using GPU+2DLA)

The script will run following Benchmarks:

- Names : Input Image Resolution
- Inception V4 : 299x299
- ResNet-50 : 224x224
- OpenPose : 256x456
- VGG-19 : 224x224
- YOLO-V3 : 608x608
- Super Resolution : 481x321
- Unet : 256x256

For benchmark results on all NVIDIA Jetson Products; please have a look at [NVIDIA jetson\\_benchmark webpage](#)

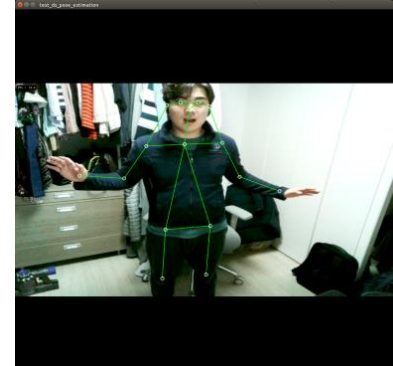
Following scripts are included:

1. Install Requirements for running benchmark script (install\_requirements.sh)
2. CSV files containing parameters (benchmark\_csv folder)
3. Download Model (utils/download\_models.py)

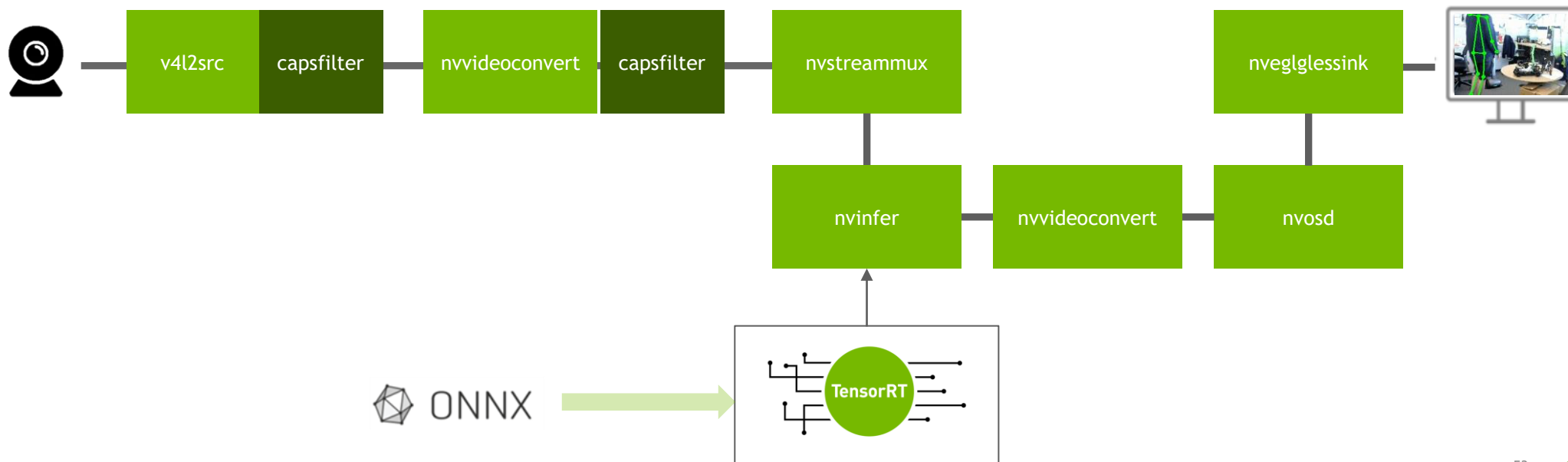
Model Name	FPS
inception_v4	311.73
vgg19_N2	66.43
super_resolution_bsd500	150.46
unet-segmentation	145.42
pose_estimation	237.1
yolov3-tiny-416	546.69
ResNet50_224x224	824.02
ssd-mobilenet-v1	887.6

# HOW TO DEPLOY DL

By using TensorRT



8\_test\_ds\_pose\_estimation

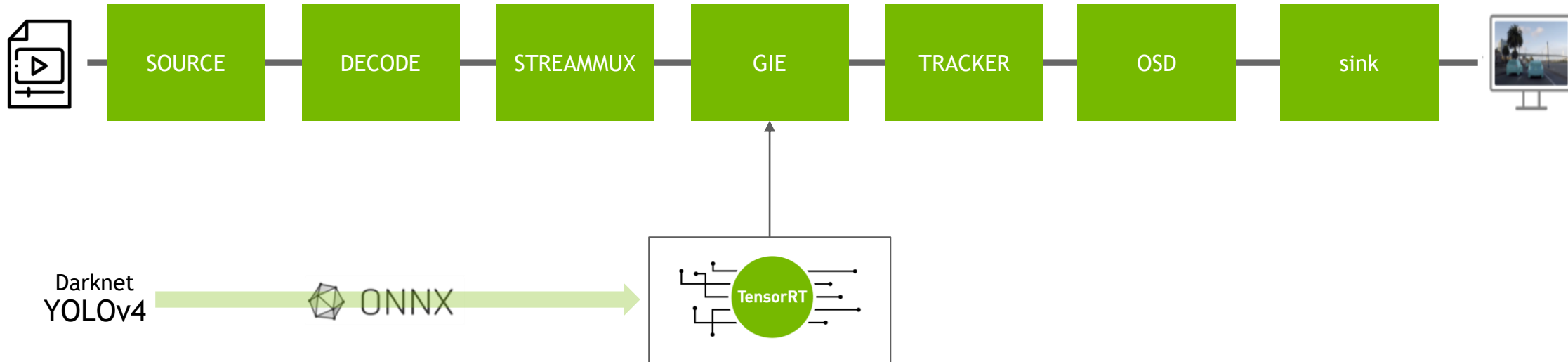


# HOW TO DEPLOY DL

By using TensorRT



9\_test\_ds\_yolov4

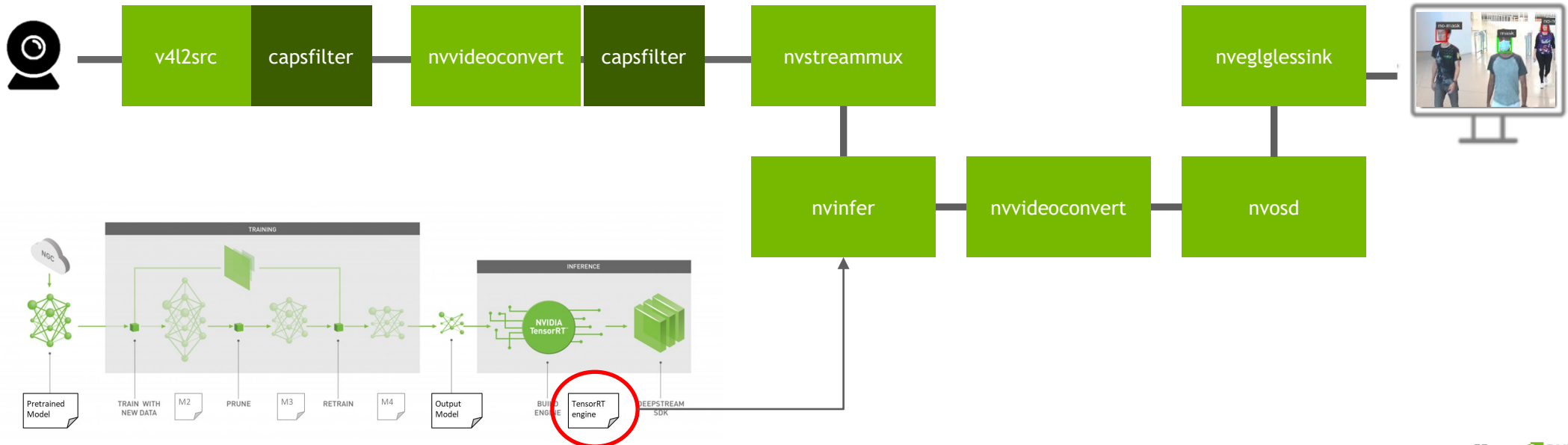


# HOW TO DEPLOY DL

By using Transfer Learning Toolkit



10\_test\_ds\_facemask



# Q&A

You can send an E-MAIL me ([jonghwanl@nvidia.com](mailto:jonghwanl@nvidia.com))  
if you have any questions for this



