# S9391
# GstCUDA: Easy GStreamer and CUDA Integration

**Eng. Daniel Garbanzo**
**MSc. Michael Grüner**
**GTC March 2019**

**RidgeRun** Embedded Solutions

# Agenda

About RidgeRun

GStreamer Overview

CUDA Overview

GstCUDA Introduction

Application Examples

Performance Statistics

GstCUDA Demo on TX2

Q&A

# About Us

- **US Company - R&D Lab in Costa Rica**

- **15 years of experience**

- **Embedded Linux and GStreamer experts**

- **Custom multimedia solutions**

- **Digital signal/image processing**

- **AI and Machine Learning solutions**

- **System optimization: CUDA, GStreamer, OpenCL, OpenGL, OpenVX, Vulkan**

- **Support for embedded and resource constrained systems**

- **Professional services, dedicated teams and specialized tools**

# Multimedia Is Everywhere

**Medical Industry**

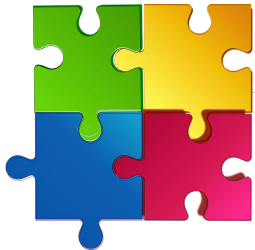**Automotive Industry**

**Smart Devices**

**Computer Vision**

- **Complex multimedia applications require a lot of processing resources**

- **GStreamer offers a flexible way for creating multimedia applications**

- **CUDA offers high performance accelerated processing capabilities**

**RidgeRun**
Embedded Solutions

**gstreamer**

- **Open source framework for audio and video applications**

- **Based on a pipeline architecture**

- **Extensible design based on plugins (more than 1000 freely available)**

- **Automatic format and synchronization handling**
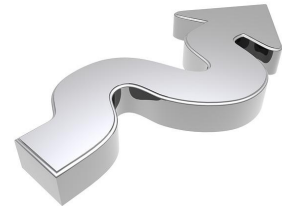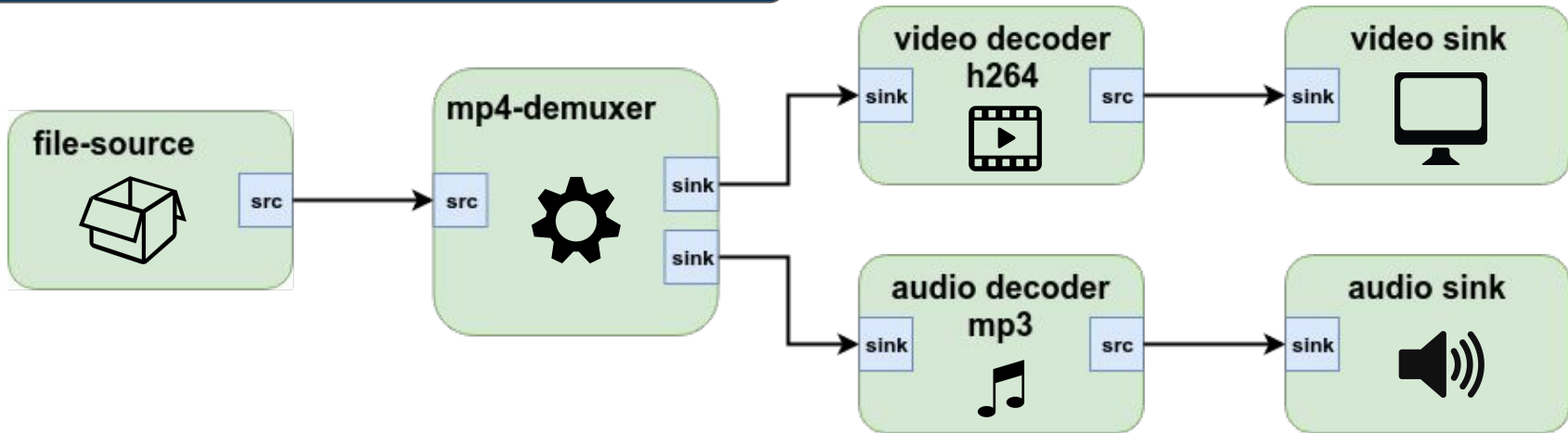
- **Tools for easy prototyping**

**Modularity**
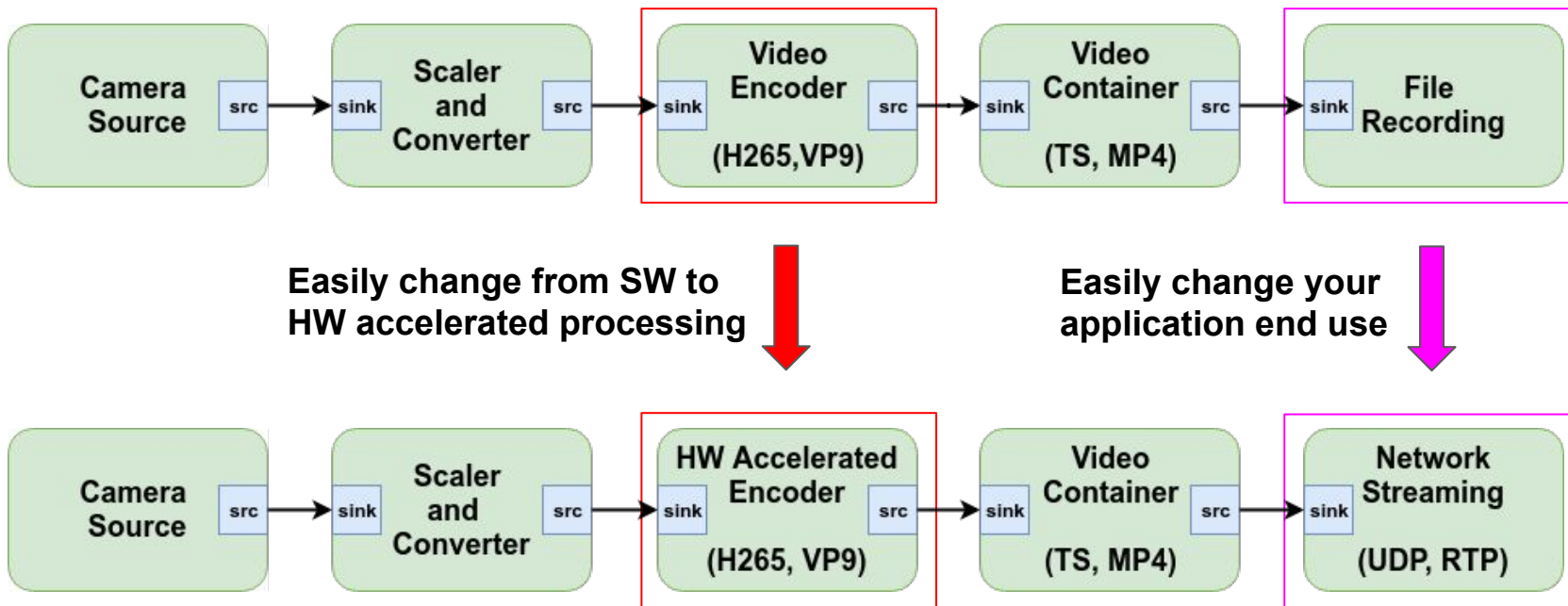
**Portability**

**Flexibility**

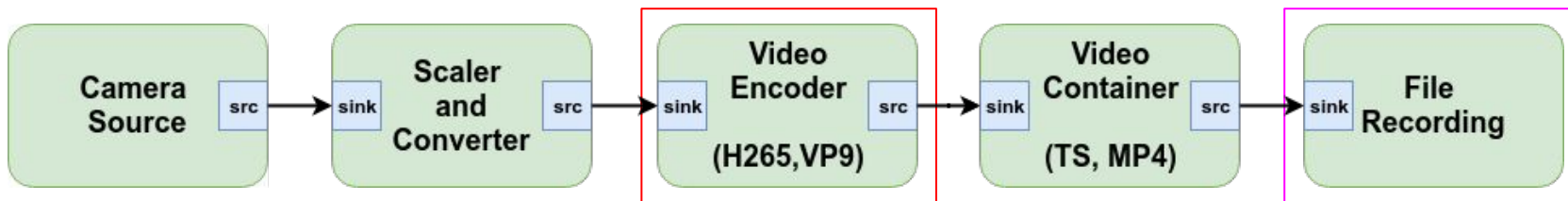**Basic MP4 player GStreamer Pipeline**

- **Each plugin represents a different processing module**

- **The plugins are linked and arranged in a pipeline**

- **Freedom to build arbitrary pipelines for different applications**

# Modular design lets you change your application easily!

Easily change from SW to HW accelerated processing
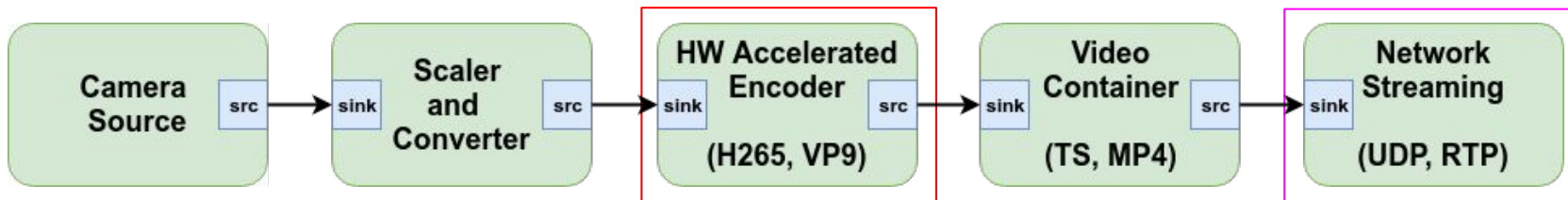
Easily change your application end use

# Modular design lets you change your application easily!

**Code equivalent :**

```
gst-launch v4l2src ! videoconverter ! x265enc ! mpegtsmux ! filesink
```
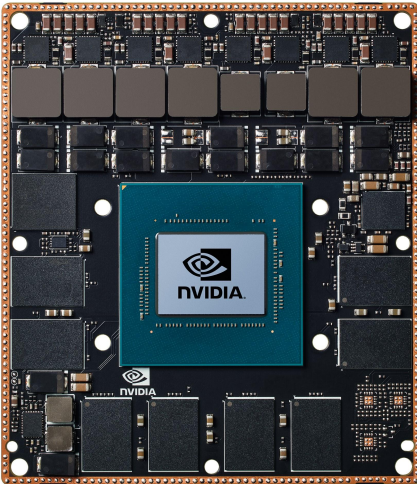
**Code equivalent :**

```
gst-launch v4l2src ! videoconverter ! omxh265enc ! mpegtsmux ! udpsink
```

**Development environment for high performance GPU-accelerated applications**

**General purpose data processing via parallel algorithm execution on GPU**

**Extensive development, debugging and profiling set of tools**

# GstCUDA Integrates the Best of Both Worlds

gstreamer + GstCUDA + NVIDIA CUDA

RidgeRun
Embedded Solutions

RidgeRun
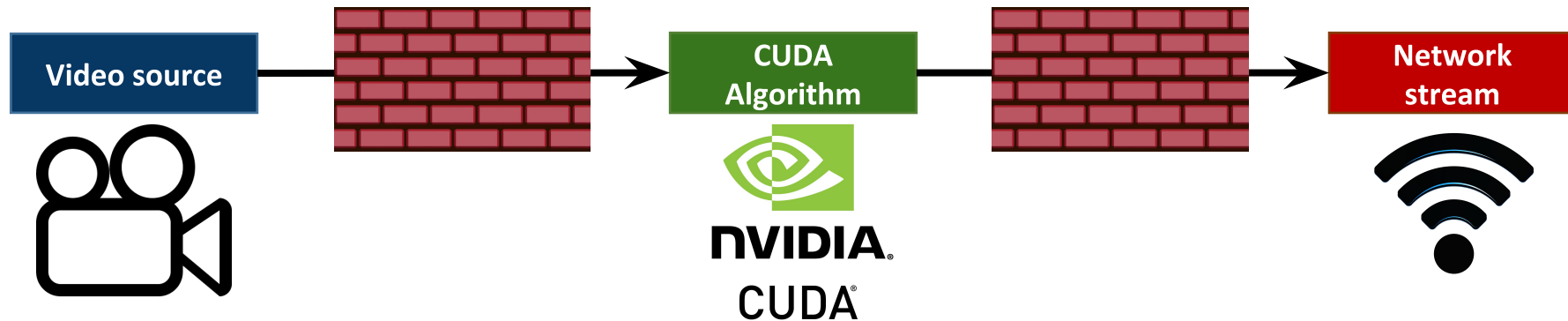Embedded Solutions

**Framework enabling easy integration of CUDA algorithms into GStreamer pipelines**

# GstCUDA

**Eliminates the need to learn GStreamer internals**

**Focus on your CUDA algorithm to reduce time to market!**
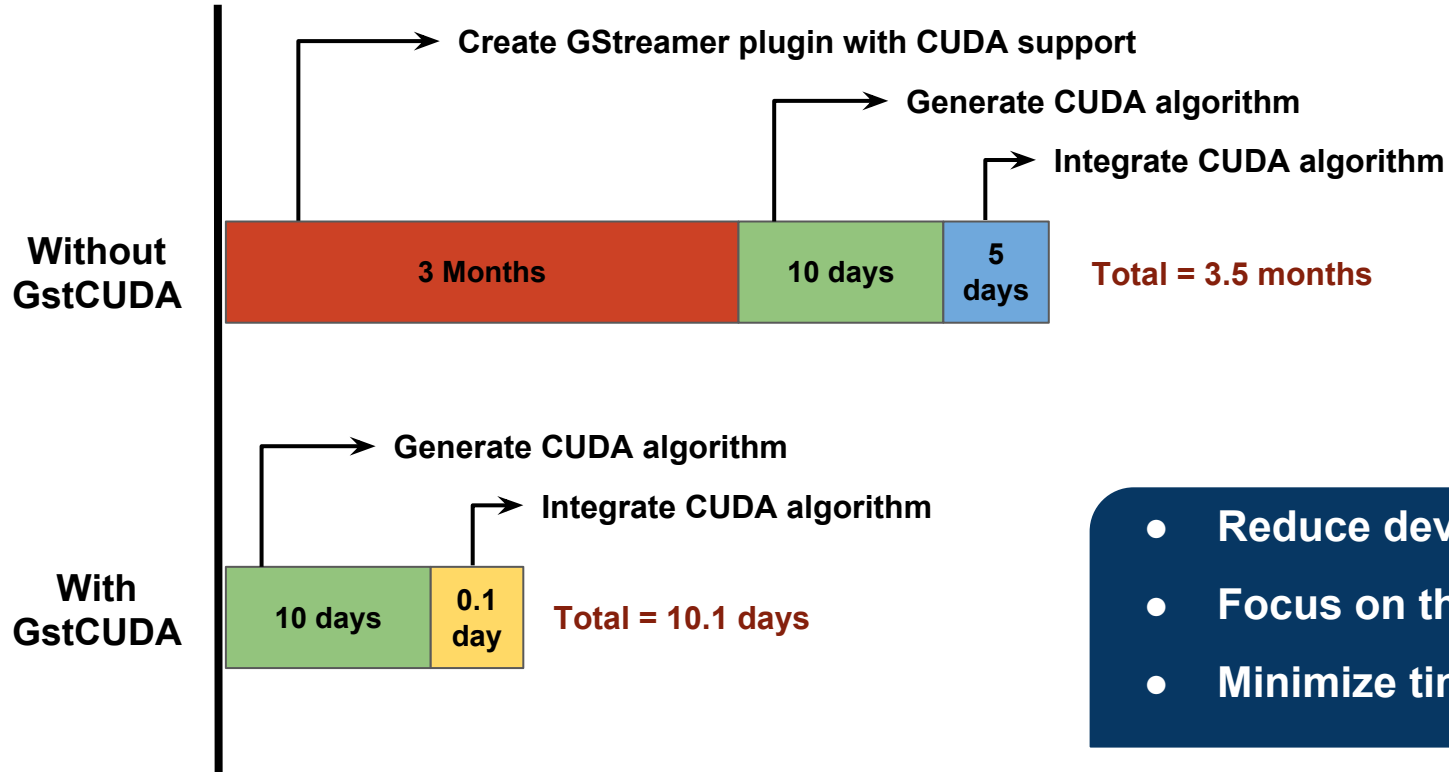
# What Does GstCUDA Solve?

# Integration Complexities



- **A lot of roadblocks between CUDA and GStreamer**
- **These are complex and time consuming**
- **Time is money!**

# Development Time

Create GStreamer plugin with CUDA support

Generate CUDA algorithm

Integrate CUDA algorithm

**Without GstCUDA**

| 3 Months | 10 days | 5 days |

**Total = 3.5 months**

Generate CUDA algorithm

Integrate CUDA algorithm

**With GstCUDA**

| 10 days | 0.1 day |

**Total = 10.1 days**

- **Reduce development time**
- **Focus on the CUDA logic**
- **Minimize time to market**

# Performance Bottleneck

**Video source** → Memcpy → **CUDA Algorithm** (NVIDIA CUDA) → Memcpy → **Network stream**

- **Data transfers can be a bottleneck**

- **Memory copies can degrade performance**

- **Incompatibility between different memory types**

# Performance Bottleneck



## Without GstCUDA

## With GstCUDA

- **Data transfers bottleneck cause poor performance**

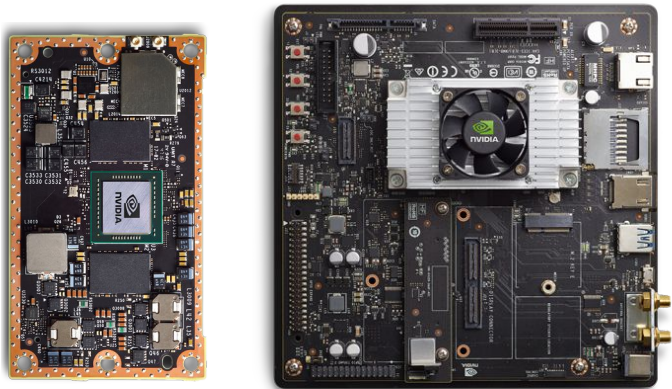- **Limited framerate at high resolutions**

- **Efficient memory handling improves performance**

- **Up to 2x 4K@60fps**

# Supported Platforms

- **Focused for NVIDIA Embedded Platforms**

**Jetson TX1, TX2, TX2i and Nano**

**Jetson AGX Xavier**

# GstCUDA Key Features

Allows CUDA algorithm easy integration into GStreamer pipelines

Out of the box quick prototyping tools

Fine-grained control of image memory layout (planes, strides, etc ...)

Automatic efficient memory handling

# GstCUDA Key Features

High performance for GStreamer/CUDA applications

Zero memory copy interface between CUDA and GStreamer

Direct handling of HW (NVMM) buffers

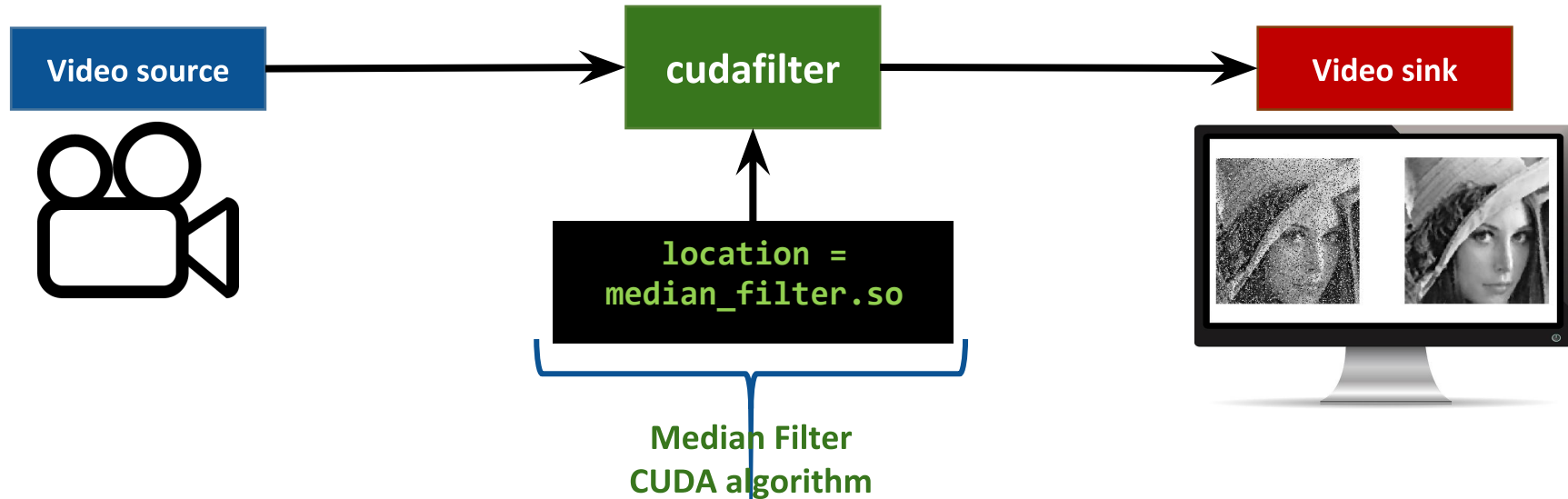Unified Memory allocation mechanism

# Framework Overview

# Quick Prototyping Elements

GStreamer elements for CUDA quick prototyping

Algorithms are loaded at runtime as plug-ins

Change the algorithm on the fly

# Cudafilter Element



**Video source**

**cudafilter**

```
location =
median_filter.so
```

**Median Filter
CUDA algorithm**

**Video sink**

**Single input / Single output configuration**

# Cudamux Element



**Multiple input / Single output configuration**

# CUDA Algorithm Interface

- **Make your CUDA algorithm compatible by implementing these interfaces**

## Cudafilter Interface

```
bool open();

bool close();

bool process (const GstCudaData &inbuf,
  GstCudaData &outbuf);

bool process_ip (const GstCudaData
  &inbuf, GstCudaData &outbuf);
```

## Cudamux Interface

```
bool open();

bool close();

bool process (vector<GstCudaData>
  &inbufs, GstCudaData &outbuf);

bool process_ip (vector<GstCudaData>
  &inbufs, GstCudaData &outbuf);
```

# Buffer Processing Methods

**process_ip**
**(In place)**

Algorithm outputs are written to the input buffer

**process**
**(Not in place)**

Algorithm inputs and outputs are different buffers

# Create Your Custom Element

- **Some applications may require specialized elements**
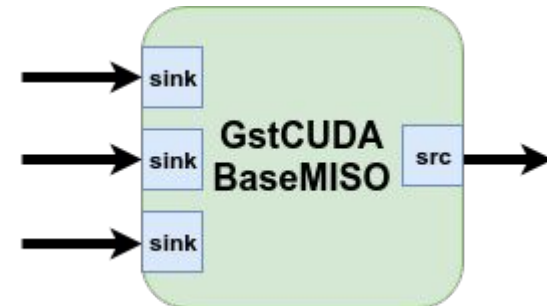- **GstCUDA provides bases classes to simplify development**

**GstCUDABaseFilter:**

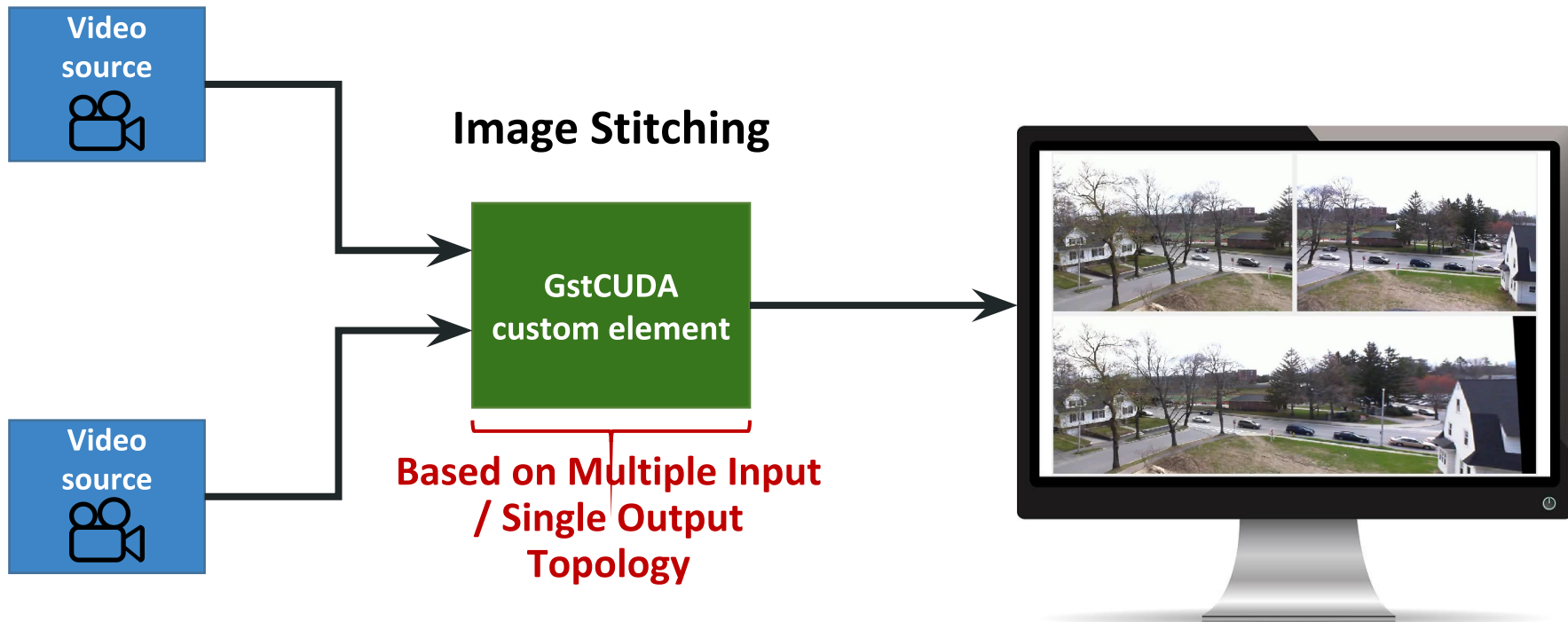- **Single-input / Single-output topology**

**GstCUDABaseMISO:**

- **Multiple-input / Single-output topology**

# GstCUDA Framework Usage Example

- **Inherit parent classes and focus on the algorithm!**

**Video source**

**Image Stitching**

**GstCUDA custom element**

**Based on Multiple Input / Single Output Topology**

**Video source**

# GstCUDA Framework Summary

- **The framework includes:**

### GstCUDA API

- **Utils to handle memory interfaces**
- **GStreamer Unified Memory allocators**
- **Parent classes for different topologies**
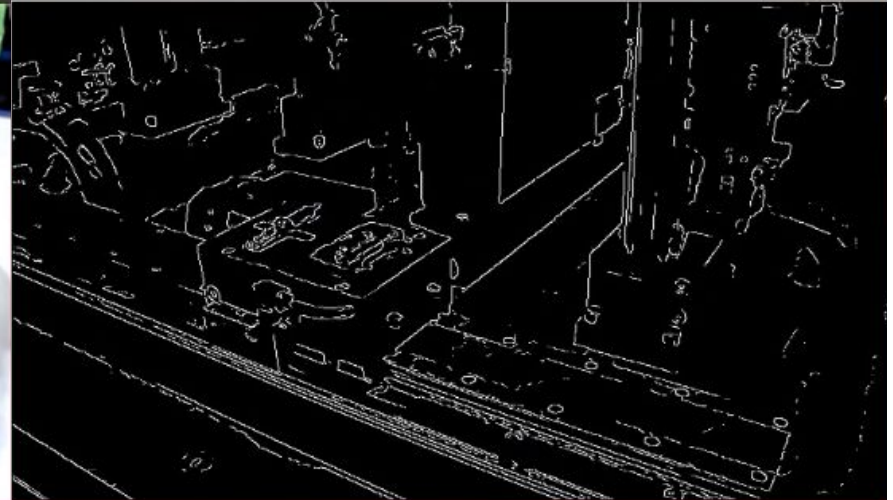
### Quick prototyping elements

- **Generic elements to evaluate custom algorithms**
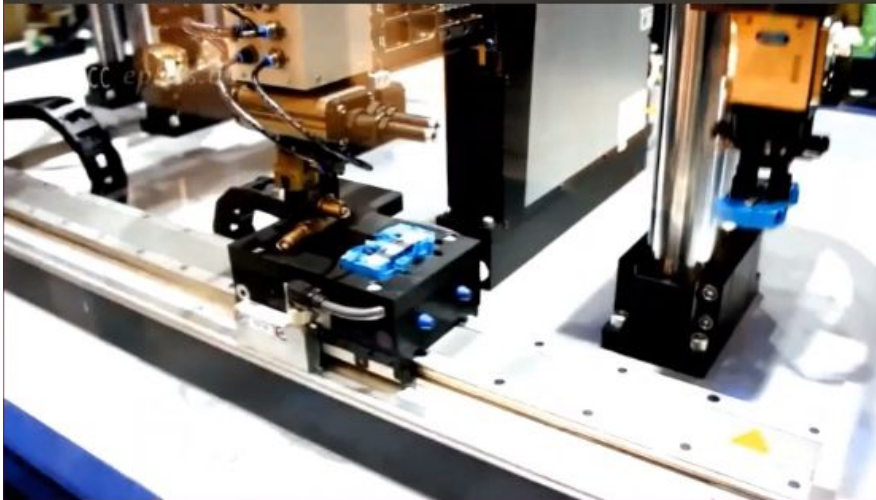- **Runtime loading of CUDA algorithms**

### Set of examples

- **Complete GstCUDA element boilerplate**
- **CUDA algorithms for the prototyping elements**

# GstCUDA Application Areas Examples Video

# Industrial Applications: Border Enhancement

# Automation Applications: Hough Transform

# Security Applications: Motion Detection/Estimation

# Performance Statistics

# Varying Algorithm / Fixed Image Size

## Test Conditions

- **Image convolution algorithm**

$$(f * g)[n] = \sum_{m=-\infty}^{\infty} f[m]g[n-m]$$

- **Stressing compute capabilities**

- **Variable convolution kernel size**

- **1080p@240fps / 1080p@60fps stream input**

- **Cudafilter element**

- **Unified Memory allocator**

- **Jetson TX2 platform**

- **Not In-place**

```
location =
convolution.so
```

nvcamerasrc | src → sink | cudafilter | src → sink | fakesink

# Varying Algorithm / Fixed Image Size

## Framerate Stats

### Average maximum framerate at different convolution kernel sizes
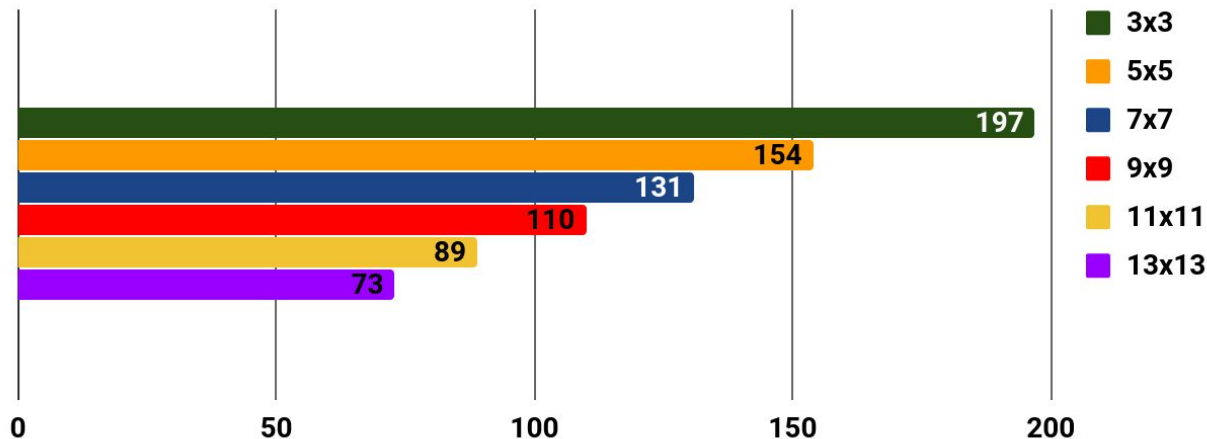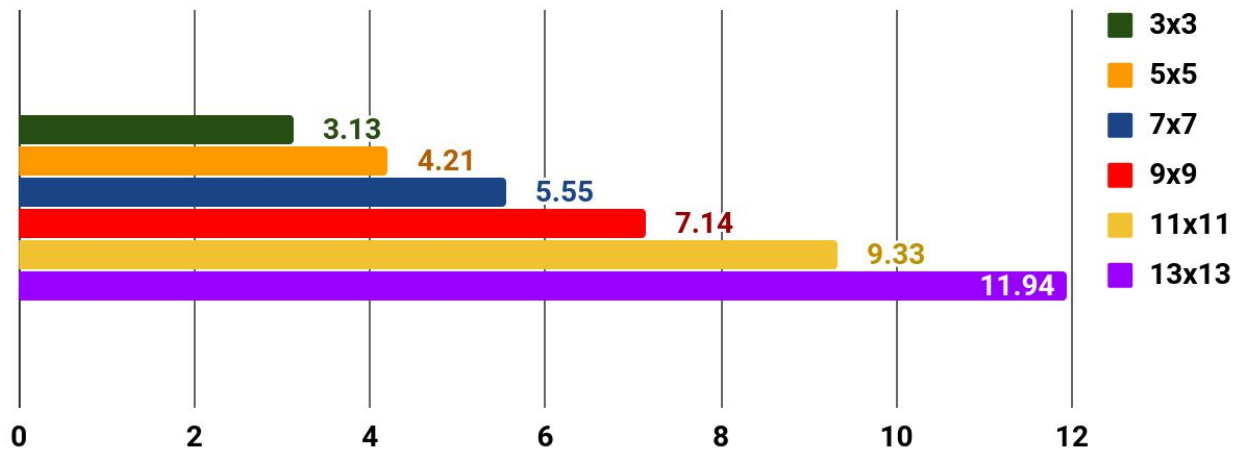
**1080p@240ps input**



Legend:
- 3x3
- 5x5
- 7x7
- 9x9
- 11x11
- 13x13

Bar values:
- 3x3: 197
- 5x5: 154
- 7x7: 131
- 9x9: 110
- 11x11: 89
- 13x13: 73

X-axis: 0, 50, 100, 150, 200

**Average Maximum Framerate [fps]**

# Varying Algorithm / Fixed Image Size

## CPU Load Stats

**Average CPU load at different convolution kernel sizes**

1080p@60fps input



Legend: 3x3, 5x5, 7x7, 9x9, 11x11, 13x13, 1080p@60fps baseline*

Values: 10.1, 10.3, 10.1, 10.2, 10.3, 10.5, 6.9

Average CPU load [%]

## GPU Load Stats

**Average GPU load at different convolution kernel sizes**

1080p@60fps input



Legend: 3x3, 5x5, 7x7, 9x9, 11x11, 13x13

Values: 7.1, 13.2, 21.5, 30.8, 44.8, 61.2

Average GPU load [%]

*baseline = simple capture pipeline (without GstCUDA)

# Fixed Algorithm / Varying Image Size

## Test Conditions

- **Memory copy algorithm**

$$y_{(n,m)} = x_{(n,m)}$$

- **Stressing data transfer**

- **Variable input resolution**

- **Cudafilter element**

- **Unified Memory allocator**

- **Jetson TX2 platform**

- **In-place vrs not In-place**

```
location = memcpy.so
```

nvcamerasrc  src → sink  cudafilter  src → sink  fakesink

# Fixed Algorithm / Varying Image Size



**Framerate Stats**

## Average maximum framerate at different resolutions

**Note:** Maximum Framerate limited to 245 fps by the video source

Legend:
- 4K Not-IP
- 4K IP
- 1080p Not-IP
- 1080p IP
- 720p Not-IP
- 720p IP
- 1280x540 Not-IP
- 1280x540 IP

Values:
- 4K Not-IP: 135
- 4K IP: 165
- 1080p Not-IP: 215
- 1080p IP: 244
- 720p Not-IP: 245
- 720p IP: 245
- 1280x540 Not-IP: 245
- 1280x540 IP: 245

X-axis: 0, 50, 100, 150, 200, 250

**Average Maximum Framerate [fps]**

# Fixed Algorithm / Varying Image Size

## Processing Time Stats

### Average processing time at different resolutions



**Legend:**
- 4K@60fps Not-IP
- 4K@60fps IP
- 1080p@60fps Not-IP
- 1080p@60fps IP
- 720p@60fps Not-IP
- 720p@60fps IP

Values:
- 4.195
- 0.856
- 2.595
- 0.723
- 1.599
- 0.727

**Average Processing Time [ms]**

# Fixed Algorithm / Varying Image Size

## CPU Load Stats

**Average CPU load at different resolutions**

- 4K@60fps Not-IP: 15.6
- 4K@60fps IP: 11.2
- 1080p@60fps Not-IP: 10
- 1080p@60fps IP: 7.8
- 720p@60fps Not-IP: 7.8
- 720p@60fps IP: 6.2
- 4K@60fps baseline*: 10
- 1080p@60fps baseline*: 6.9
- 720p@60fps baseline*: 5.2

**Average CPU load [%]**

## GPU Load Stats

**Average GPU load at different resolutions**

- 4K@60fps Not-IP: 10.2
- 4K@60fps IP: 1
- 1080p@60fps Not-IP: 3
- 1080p@60fps IP: 0
- 720p@60fps Not-IP: 1
- 720p@60fps IP: 0

**Average GPU load [%]**

*baseline = simple capture pipeline (without GstCUDA)

# Fixed Algorithm / Varying Image Size

### Test Conditions

- **Simple image mixing algorithm**

$$y_{(n,m)} = 0.5(x_1(n, m) + x_2(n, m))$$

- **Stressing data transfer**

- **Variable input resolution**

- **Cudamux element**

- **Unified Memory allocator**

- **In-place=True**

- **Jetson TX2 platform**



```
location =
mixer.so
```

# Fixed Algorithm / Varying Image Size

## Framerate Stats
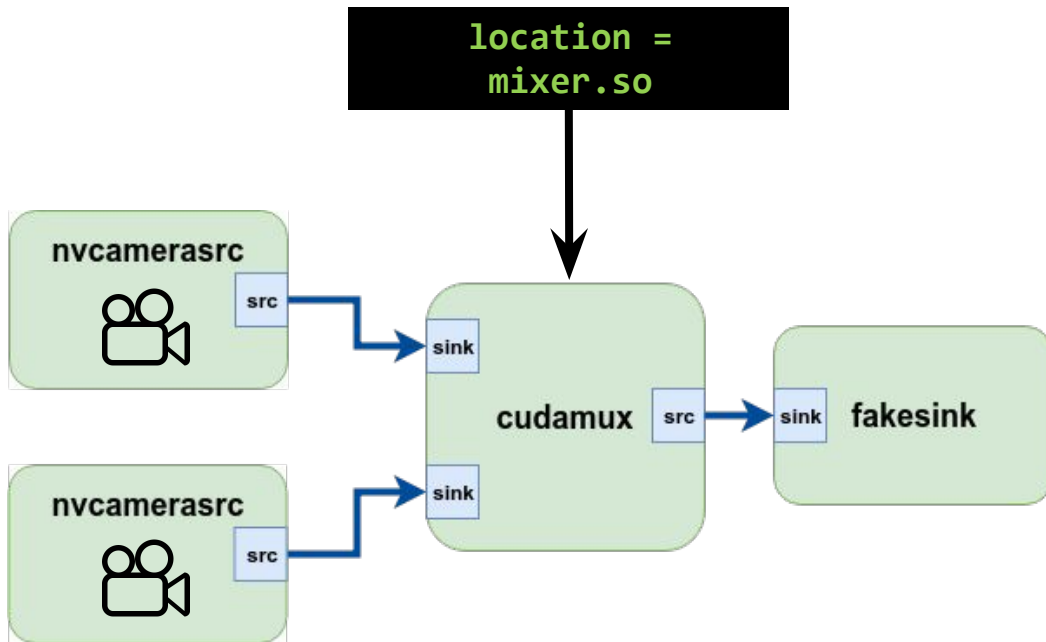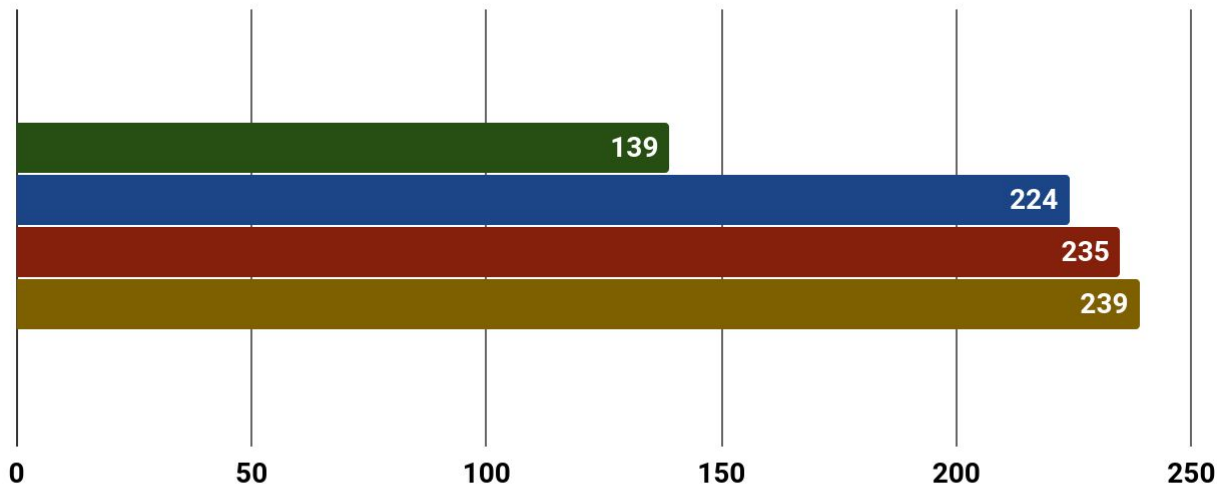
### Average maximum framerate at different resolutions

■ 4K  ■ 1080p  ■ 720p  ■ 1280x540
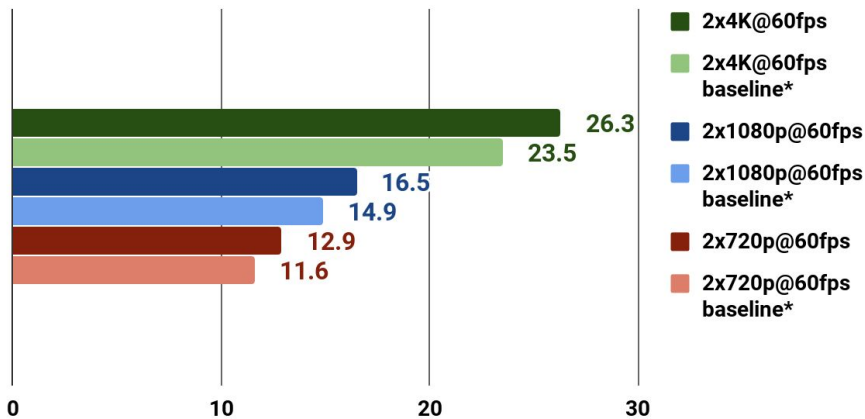
**Note:** Maximum Framerate limited to 240fps by the video source

| Resolution | Average Maximum Framerate [fps] |
|---|---|
| 4K | 139 |
| 1080p | 224 |
| 720p | 235 |
| 1280x540 | 239 |

Average Maximum Framerate [fps]

# Fixed Algorithm / Varying Image Size

RidgeRun
Embedded Solutions

## CPU Load Stats

**Average CPU load at different resolutions**

Legend:
- 2x4K@60fps
- 2x4K@60fps baseline*
- 2x1080p@60fps
- 2x1080p@60fps baseline*
- 2x720p@60fps
- 2x720p@60fps baseline*

- 26.3
- 23.5
- 16.5
- 14.9
- 12.9
- 11.6

0    10    20    30

**Average CPU load [%]**

## GPU Load Stats

**Average GPU load at different resolutions**

Legend:
- 4K@60fps
- 1080p@60fps
- 720p@60fps

- 16.5
- 3.1
- 1.1

0    5    10    15    20

**Average GPU load [%]**

**\*baseline = simple capture pipeline (without GstCUDA)**

# GstCUDA Live Demo on Jetson TX2
# Sobel Filter 1080p60fps

## Code equivalent :

```
gst-launch-1.0  nvcamerasrc  sensor-id=2  fpsRange=60,60  !
"video/x-raw(memory:NVMM),width=1920,height=1080,framerate=6
0/1,format=I420"  !  nvvidconv  !  "video/x-raw"  !  queue  !
cudafilter  in-place=false  location=/borders.so  !  queue  !
nvoverlaysink
```

# Resources

- **GstCUDA wiki page:**

  - **gstcuda.ridgerun.com**

- **RidgeRun Website:**

  - **ridgerun.com**

- **RidgeRun Contact:**

  - **ridgerun.com/contact**