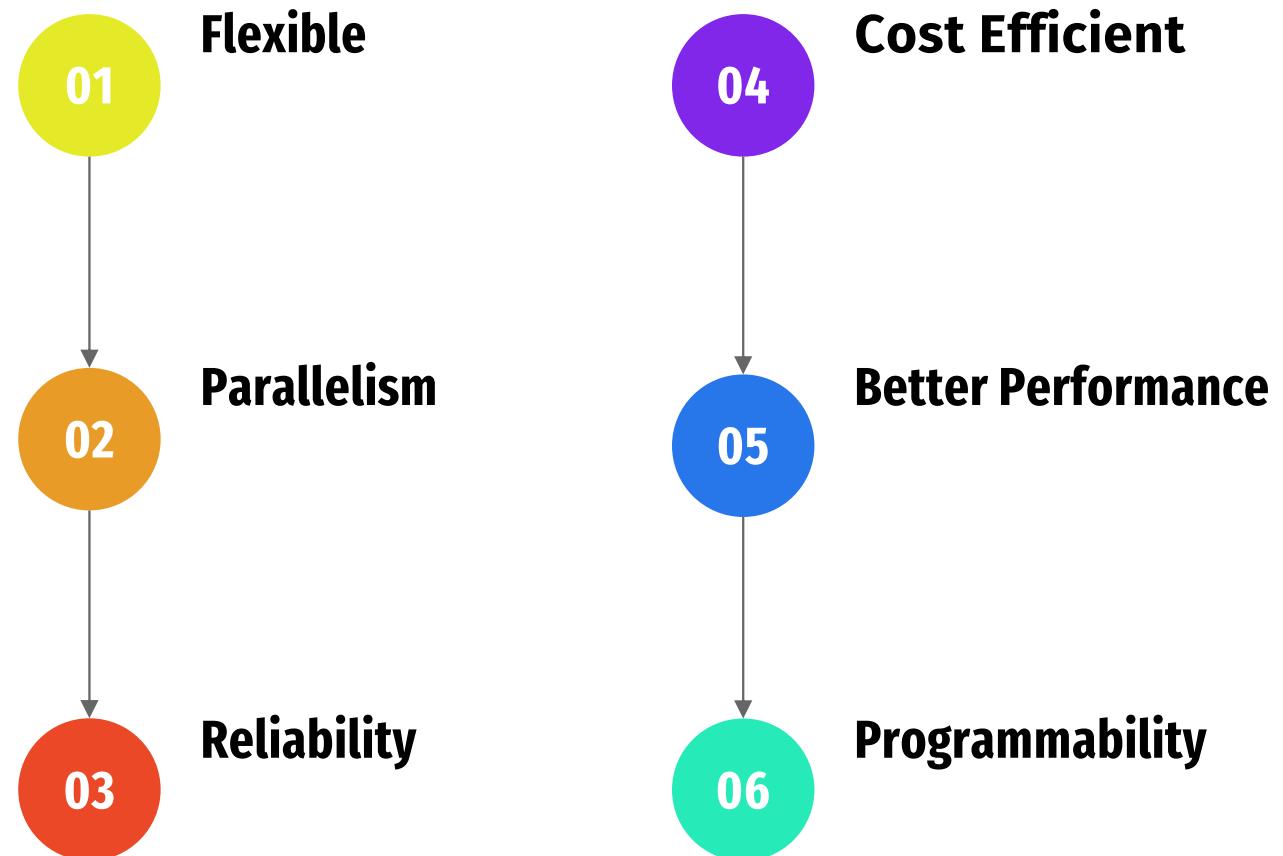
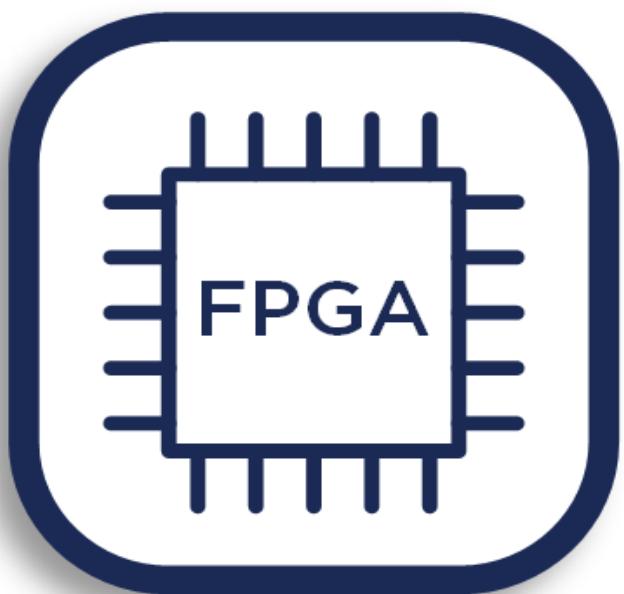




# Multi-Video Analytics System Based on FPGA

Name: Sabol Socare  
ID. : 62050104

# FPGA Infographics



# Applications on FPGA

## Some FPGA applications



**Medical**

Diagnostic and monitoring purposes

**Server & Cloud**

Data analytics and metadata and data processing

**Defense & Space**

Used for high-resolution optical data processing and radar imaging

**Video & Image Processing**

Accelerate image and video processing applications while meeting low-power constraints

# Multi-Video Analytics System Based on FPGA



## Project's objectives

### Multi-Video Analytics

02

Develop multi-video analytics system for Detection and Classification in security system

### Customized AI model

01

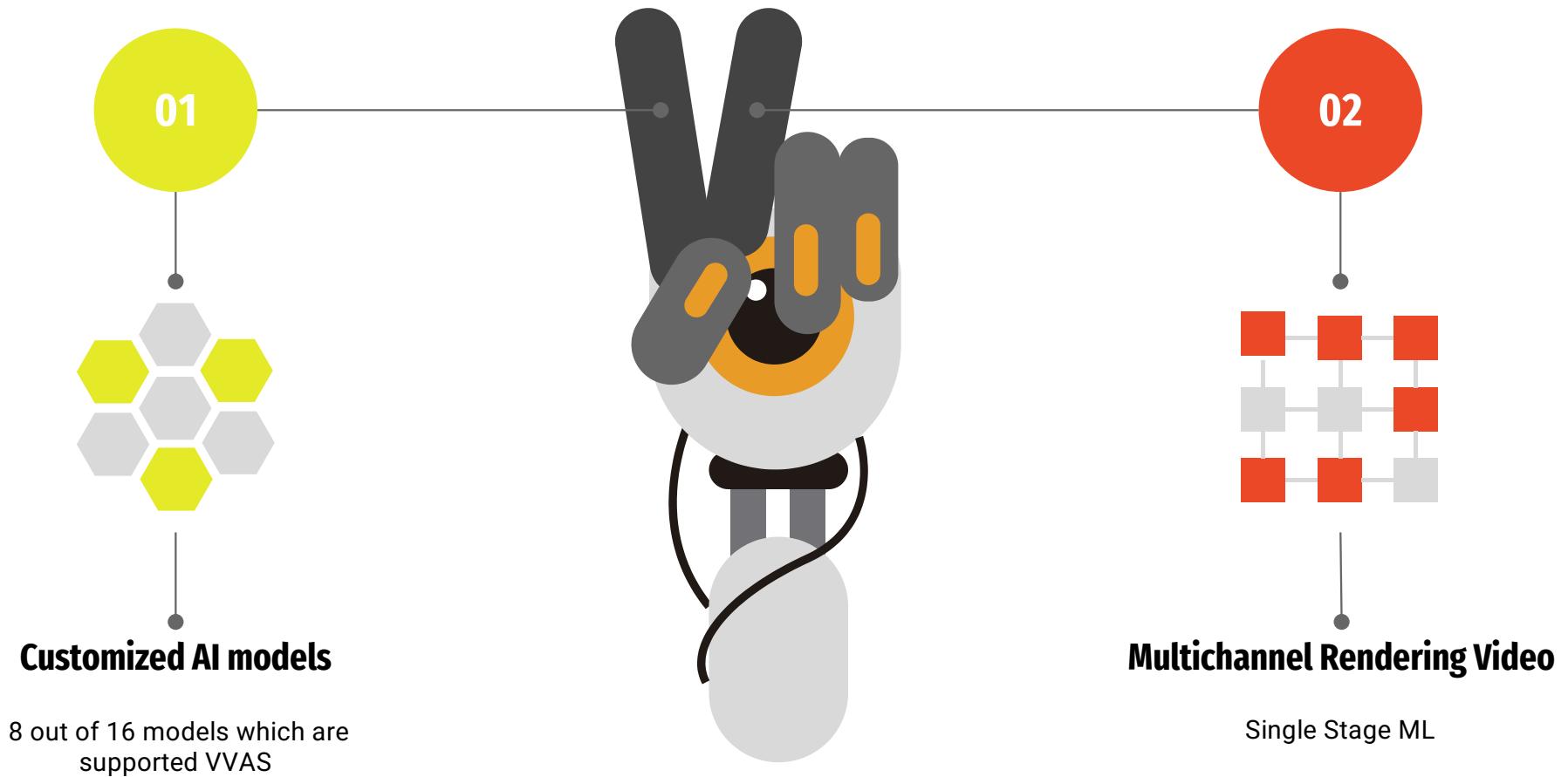
Compile all 8 models which are supported VVAS framework

### Retrain AI models

03

Utilize AI model for detection and classification for testing performance on CPU, GPU, FPGA

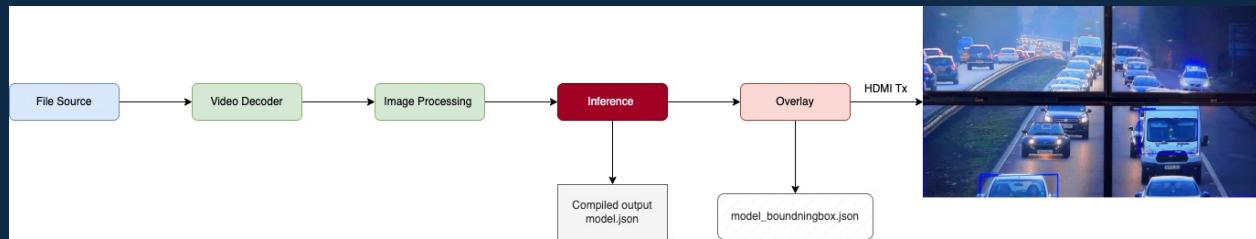
# Multi-Video Analytics System Based on FPGA



# What is Single stage ML ?

Single Stage ML (Machine Learning) is ML operation that involve only one ML operation on the input image.

So we can scale up to build 4 channels video analytics by process 4 streams in parallel.



# Project's Scope



# Project Equipment

01 ZCU104 Evaluation Board rev 1.0

02 Micro USB wire

03 HDMI 2.0

04 HDMI 2.0 Wire

05 MicroSD Card

Vitis Unified Software platform version 2022.2

01

PetaLinux tool version 2022.2

02

Serial terminal emulator

03

Git

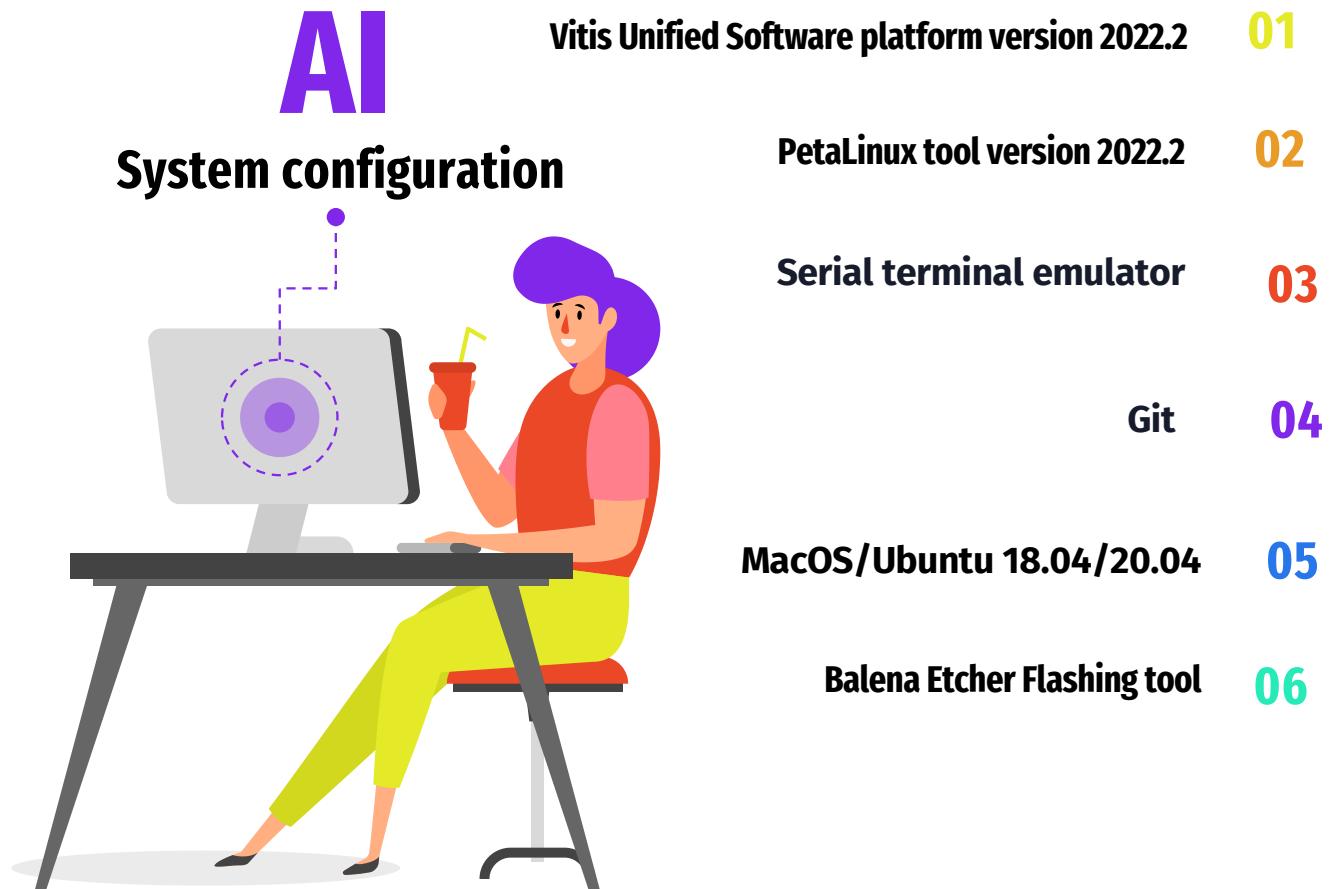
04

MacOS/Ubuntu 18.04/20.04

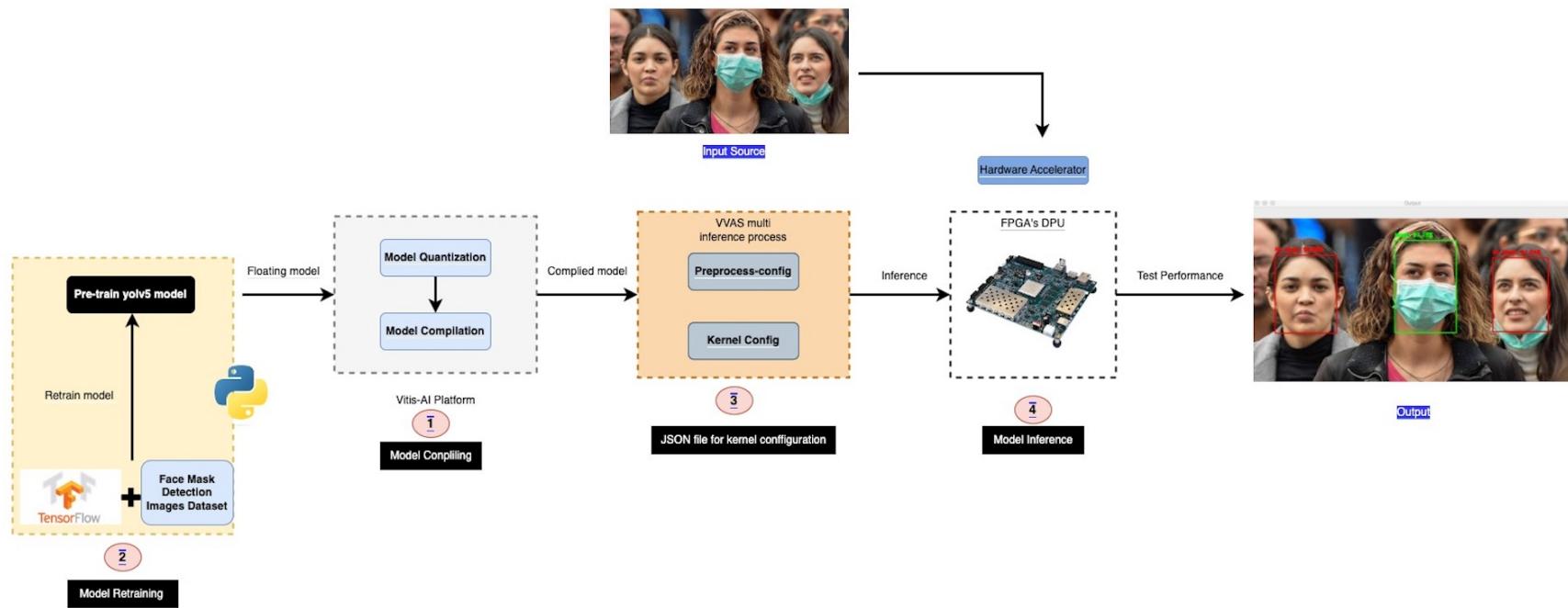
05

Balena Etcher Flashing tool

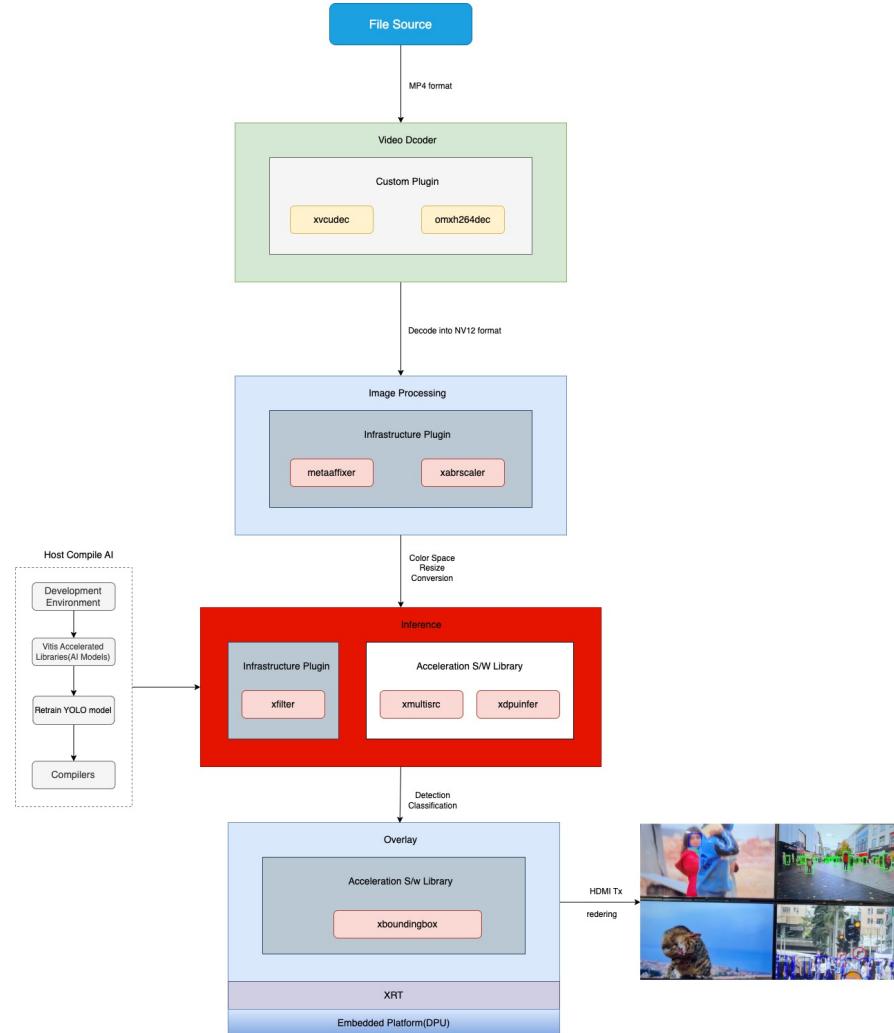
06



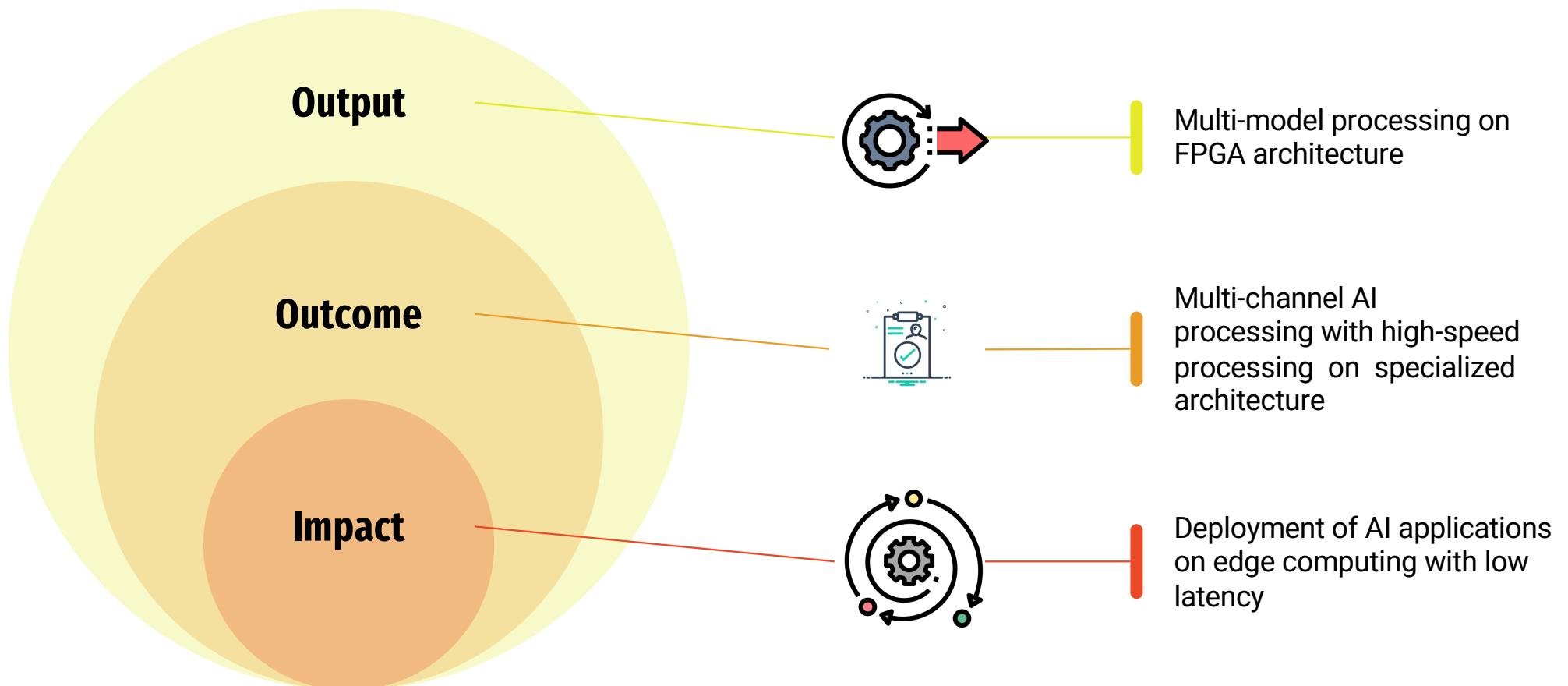
# Project Overall's diagram



# AIC's Infrastructure



# Output/Outcome/Impact



# Experiment test

01

## Compile AI models

Compiling AI model  
from model Zoo, Xilinx

02

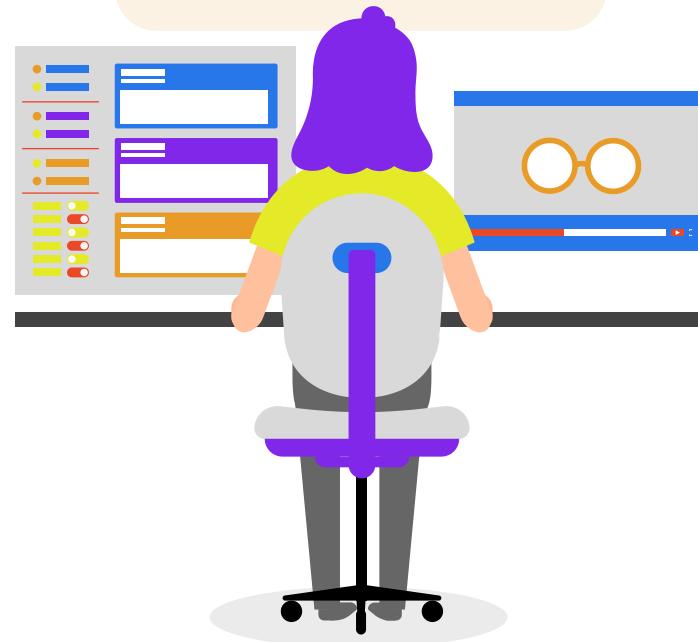
## Multi-video analytics

Analyze multi-video on  
FPGA

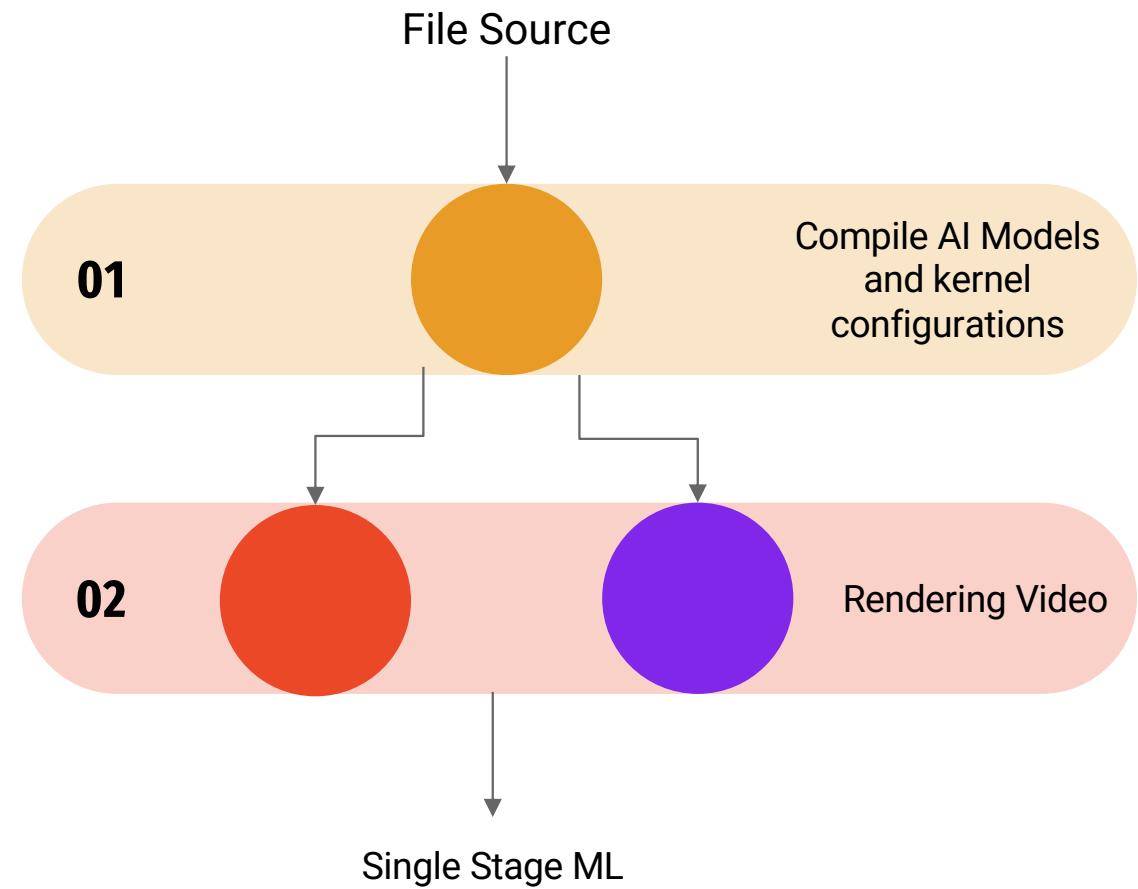
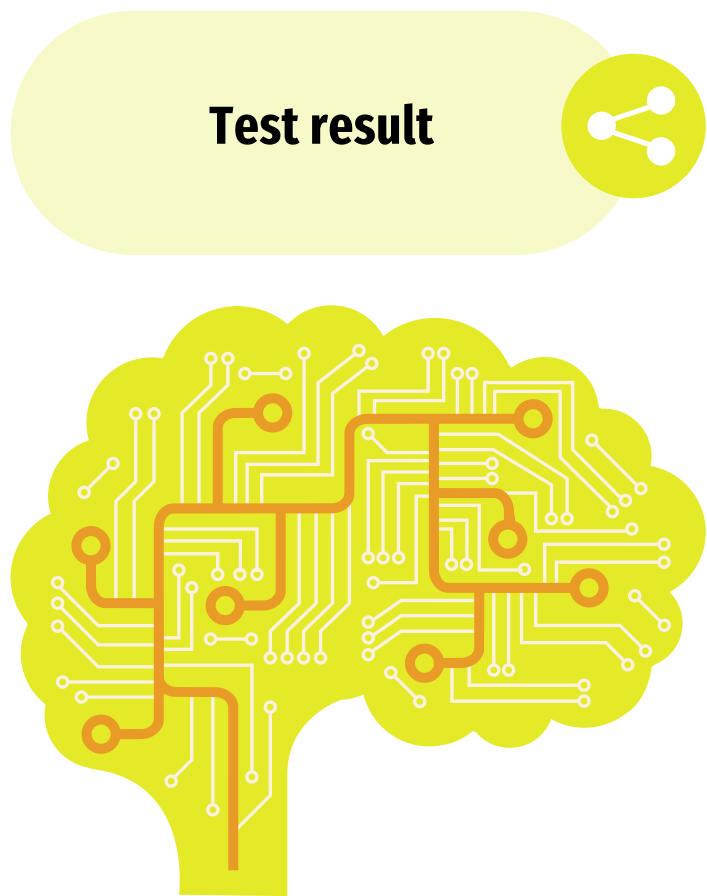
03

## Train AI models

Train and test performance  
with AI models



# Experiment Result



# 1. Compile AI model

## Cross compile

```
((vitis-ai-caffe) Vitis-AI /workspace/vvas_model > source compile_cf_model.sh ssd_traffic_pruned_0_9 cf_
ssdtraffic_360_480_0.9_11.6G_2.0
*****
* VITIS_AI Compilation - Xilinx Inc.
*****
[INFO] Namespace(batchsize=1, inputs_shape=None, layout='NCHW', model_files=['../models/AI-Model-Zoo/caffemodel(cf_ssdtraffic_360_480_0.9_11.6G_2.0/quantized/deploy.caffemodel)', model_type='caffe', named_inputs_shape=None, out_filename='/tmp/ssd_traffic_pruned_0_9.org.xmodel', proto='../models/AI-Model-Zoo/caffemodel(cf_ssdtraffic_360_480_0.9_11.6G_2.0/quantized/deploy.prototxt')
[INFO] caffe model: /workspace/models/AI-Model-Zoo/caffemodel(cf_ssdtraffic_360_480_0.9_11.6G_2.0/quantized/deploy.caffemodel
[INFO] caffe model: /workspace/models/AI-Model-Zoo/caffemodel(cf_ssdtraffic_360_480_0.9_11.6G_2.0/quantized/deploy.prototxt
[INFO] parse raw model :100%[██████████] 123/123 [00:00<00:00, 332.52it/s]
[INFO] infer shape (NCHW) :100%[██████████] 123/123 [00:00<00:00, 1423.88it/s]
[INFO] infer shape (NHWC) :100%[██████████] 123/123 [00:00<00:00, 1619.24it/s]
[INFO] perform level-1 opt :100%[██████████] 3/3 [00:00<00:00, 206.10it/s]
[INFO] generate xmodel :100%[██████████] 123/123 [00:00<00:00, 2264.00it/s]
[INFO] dump xmodel: /tmp/ssd_traffic_pruned_0_9.org.xmodel
[UNILOG][INFO] Compile mode: dpu
[UNILOG][INFO] Debug mode: function
[UNILOG][INFO] Target architecture: DPUCZDX8G_ISA0_B3136_MAX_BG2
[UNILOG][INFO] Graph name: deploy, with op num: 263
[UNILOG][INFO] Begin to compile...
[UNILOG][INFO] Total device subgraph number 5, DPU subgraph number 1
[UNILOG][INFO] Compile done.
[UNILOG][INFO] The meta json is saved to "/workspace/vvas_model/.compiled_output(cf_ssdtraffic_360_480_0.9_11.6G_2.0/meta.json"
[UNILOG][INFO] The compiled xmodel is saved to "/workspace/vvas_model/.compiled_output(cf_ssdtraffic_360_480_0.9_11.6G_2.0/ssd_traffic_pruned_0_9.xmodel"
[UNILOG][INFO] The compiled xmodel's md5sum is f91e919eb6beba63de3fb34f42e5a907, and has been saved to "/workspace/vvas_model/.compiled_output(cf_ssdtraffic_360_480_0.9_11.6G_2.0/md5sum.txt"
(vitis-ai-caffe) Vitis-AI /workspace/vvas_model > 
```

## Test jpeс

```
root@xlnx-zcu104-vDec-hdmiTx-2021-2:~/Vitis-AI-Library/samples/ssd# ./test_jpeg_ssdsd_traffic_pruned_0_9 sample_ssds.jpg
WARNING: Logging before InitGoogleLogging() is written to STDERR
I0919 07:38:46.888861 2247 demo.hpp:1183] batch: 0 image: sample_ssds.jpg
I0919 07:38:46.889010 2247 process_result.hpp:431] RESULT: 1 469.372 319.82 518.265 360.897 0.964248
I0919 07:38:46.889141 2247 process_result.hpp:431] RESULT: 1 549.768 338.249 649.877 383.833 0.953987
I0919 07:38:46.889193 2247 process_result.hpp:431] RESULT: 1 627.717 318.577 797.316 543.785 0.953987
I0919 07:38:46.889242 2247 process_result.hpp:431] RESULT: 1 499.001 338.743 561.493 368.852 0.951322
I0919 07:38:46.889389 2247 process_result.hpp:431] RESULT: 1 8.27212 387.261 284.645 529.78 0.77812
I0919 07:38:46.889372 2247 process_result.hpp:431] RESULT: 1 569.494 297.068 737.173 365.71 0.6773
I0919 07:38:46.889427 2247 process_result.hpp:431] RESULT: 1 175.256 328.477 296.489 413.135 0.62172
root@xlnx-zcu104-vDec-hdmiTx-2021-2:~/Vitis-AI-Library/samples/ssd# 
```

## Test Performance

```
root@xlnx-zcu104-vDec-hdmiTx-2021-2:~/Vitis-AI-Library/samples/ssd# ./test_performance_ssdsd_traffic_pruned_0_9 ./test_performance_ssdsd_traffic_pruned_0_9.list
WARNING: Logging before InitGoogleLogging() is written to STDERR
I0928 02:46:15.769080 3112 benchmark.hpp:184] writing report to <STDOUT>
I0928 02:46:15.769458 3112 benchmark.hpp:211] waiting for 0/30 seconds, 1 threads running
I0928 02:46:15.769737 3112 benchmark.hpp:211] waiting for 10/30 seconds, 1 threads running
I0928 02:46:15.769747 3112 benchmark.hpp:211] waiting for 20/30 seconds, 1 threads running
I0928 02:46:15.770086 3112 benchmark.hpp:219] waiting for threads terminated
FPS=46.569
E2E_MEAN=21454.4
DPU_MEAN=21348.1

```

## \*\* Compiling AI model (Model Zoo)

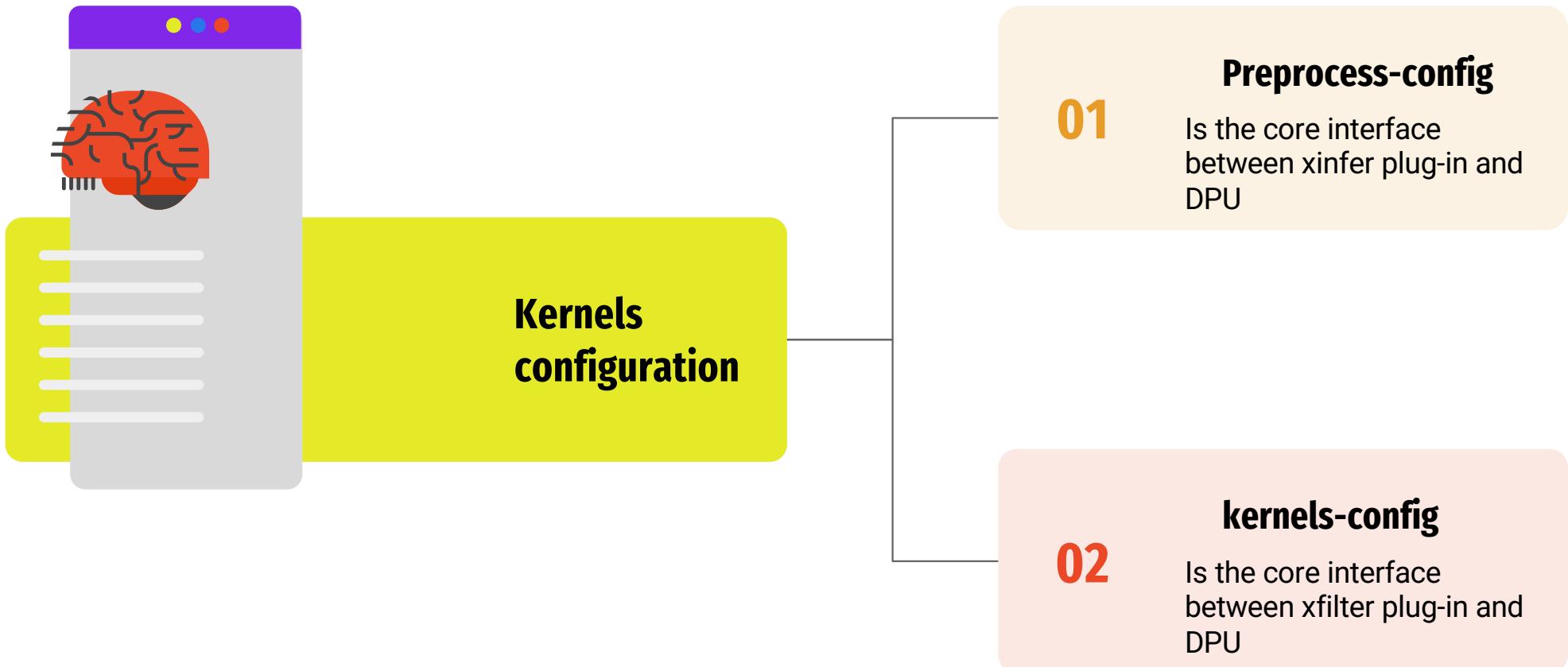
# Compile AI model

ตารางโมเดลที่หลังจากปรับแต่งทดสอบประสิทธิภาพของแต่ละโมเดล

ลำดับที่	ชื่อโมเดล	ทดสอบภาพ	ทดสอบประสิทธิภาพ
1	resnet18	ผ่าน	FPS= 68.8985
2	mobilenet_v2	ผ่าน	FPS= 239.659
3	inception_v1	ผ่าน	FPS= 148.714
4	ssd_adas_pruned_0_95	ผ่าน	FPS= 77.8735
5	ssd_mobilenet_v2	ผ่าน	FPS= 22.9255
6	ssd_pedestrian_pruned_0_97	ผ่าน	FPS= 78.0783
7	yolov3_voc_tf	ผ่าน	FPS= 10.0492
8	densebox_640_360	ผ่าน	FPS= 179.286

## 2. Multi-video analytics

Customize JSON file for rendering video



# Customize JSON file

```
1 {
2   "vvas-library-repo": "/usr/lib/",
3   "inference-level": 1,
4   "attach-ppe-outbuf": false,
5   "kernel": {
6     "library-name": "libvvas_xdpuninfer.so",
7     "config": {
8       "batch-size" : 0,
9       "model-name" : "densebox_320_320",
10      "model-class" : "FACEDECTECT",
11      "model-format" : "BGR",
12      "model-path" : "/usr/share/vitis_ai_library/models/",
13      "run_time_model" : false,
14      "need_preprocess" : false,
15      "performance_test" : false,
16      "debug_level" : 1,
17      "max-objects":3
18    }
19  }
20}
21
```

Preprocess-config

```
1 ∨{
2   "ivav-library-repo": "/usr/lib/",
3   "element-mode": "inplace",
4   "kernels": [
5     {
6       "library-name": "libvvas_xboundingbox.so",
7       "config": {
8         "model-name" : "densebox_320_320",
9         "display_output" : 1,
10        "font_size" : 0.5,
11        "font" : 3,
12        "thickness" : 3,
13        "debug_level" : 0,
14        "label_color" : { "blue" : 0, "green" : 0, "red" : 0 },
15        "label_filter" : [ "class", "probability" ],
16        "classes" : [
17        ]
18      }
19    }
20  ]
21}
```

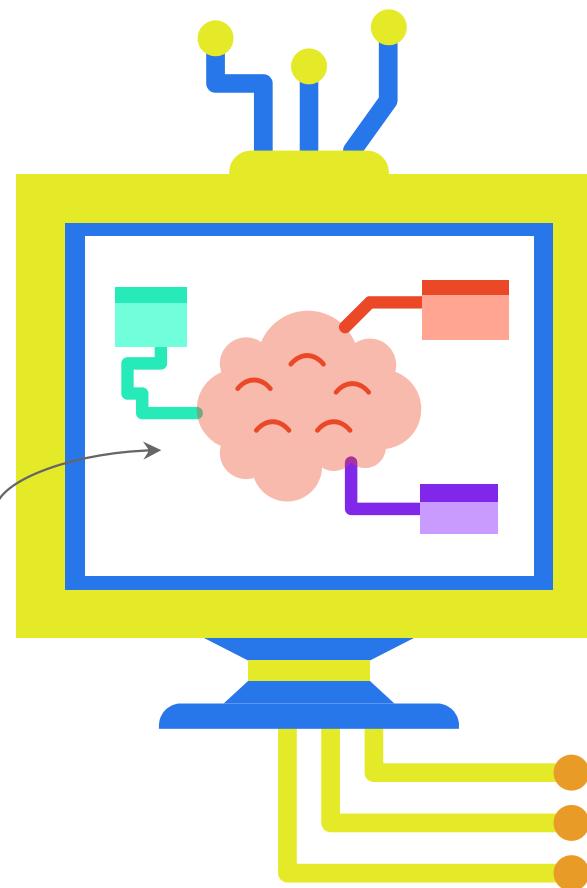
kernels-config

# Compile AI model

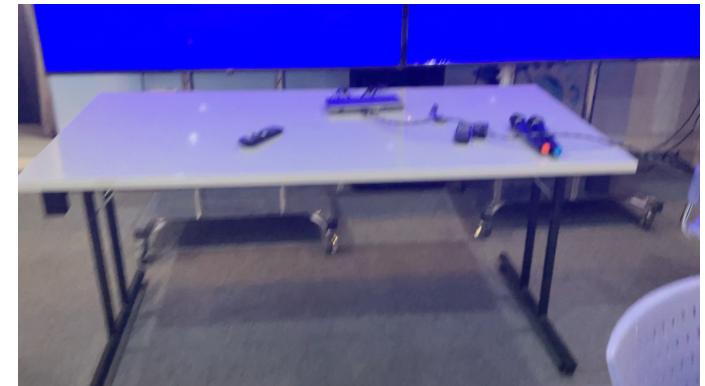


# Multichannel Machine Learning

Single Stage

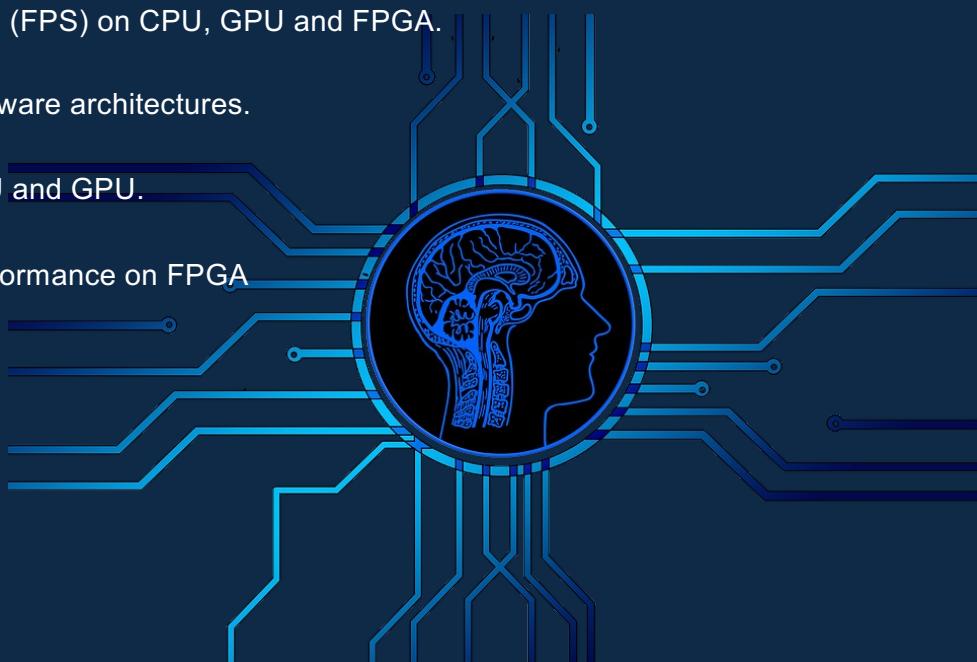


Cascaded

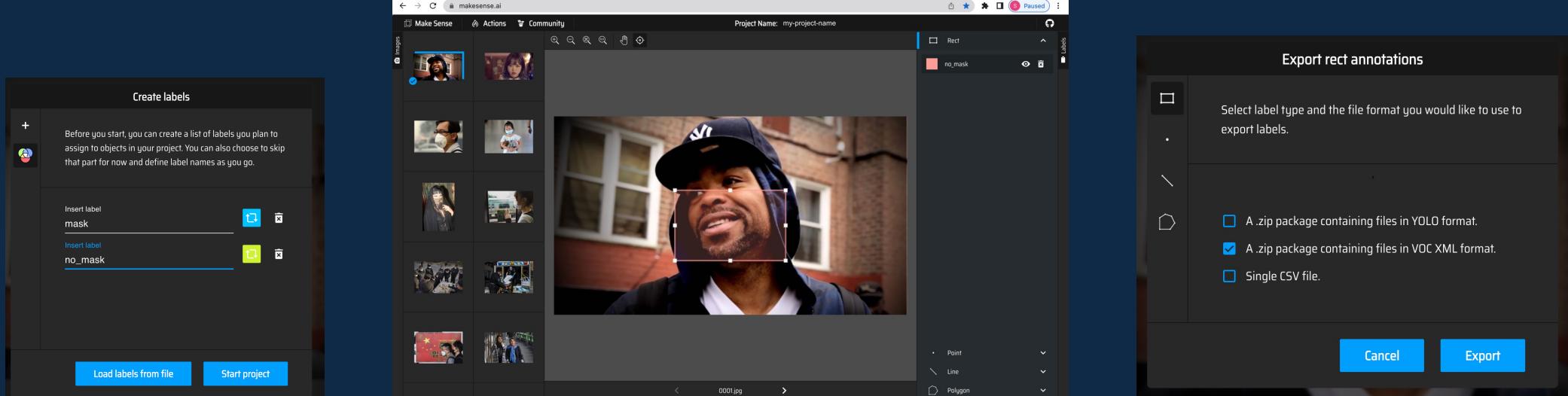


### 3. Train and test performance with AI models

- We present a comparative analysis of the performance of Frames Per Second (FPS) on CPU, GPU and FPGA.
- The FPS performance was measured and compared across the different hardware architectures.
- Train Yolov5 model using PyTorch framework for testing performance on CPU and GPU.
- Train Object classification model using Tensorflow2 framework for testing performance on FPGA



# Prepare Dataset (Yolov5)



# Prepare Dataset

```
(base) Socares-MacBook-Pro:train_data_voc socaresabol$ tree -d  
.  
└── images  
    ├── train  
    └── val  
── labels  
    ├── train  
    └── val
```

Note:

- Train images = 1295 images
- Val images = 254 images

# Training with Yolov5 model

Yolov5.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Comment Share S

RAM Disk

Files

41m 96/99 5.36G 0.01354 0.007992 0.0002363 44 640: 100% 81/81 [00:20<00:00, 3.86it/s] Class Images Instances P R mAP50 mAP50-95: 100% 8/8 [00:01<00:00, 4.06it/s] all 254 443 0.936 0.917 0.931 0.574

Epoch 97/99 GPU\_mem box\_loss obj\_loss cls\_loss Instances Size 640: 100% 81/81 [00:23<00:00, 3.48it/s] 5.36G 0.01378 0.008237 0.000329 54 mAP50 mAP50-95: 100% 8/8 [00:01<00:00, 4.08it/s] Class Images Instances P R all 254 443 0.935 0.909 0.929 0.578

Epoch 98/99 GPU\_mem box\_loss obj\_loss cls\_loss Instances Size 640: 100% 81/81 [00:21<00:00, 3.79it/s] 5.36G 0.01351 0.007794 0.0002291 38 mAP50 mAP50-95: 100% 8/8 [00:02<00:00, 3.46it/s] Class Images Instances P R all 254 443 0.939 0.916 0.939 0.579

Epoch 99/99 GPU\_mem box\_loss obj\_loss cls\_loss Instances Size 640: 100% 81/81 [00:22<00:00, 3.62it/s] 5.36G 0.01361 0.007939 0.0002859 44 mAP50 mAP50-95: 100% 8/8 [00:01<00:00, 4.07it/s] Class Images Instances P R all 254 443 0.934 0.921 0.934 0.577

100 epochs completed in 0.668 hours.  
Optimizer stripped from runs/train/exp2/weights/last.pt, 14.4MB  
Optimizer stripped from runs/train/exp2/weights/best.pt, 14.4MB

Validating runs/train/exp2/weights/best.pt...  
Fusing layers...  
Model summary: 157 layers, 7015519 parameters, 0 gradients, 15.8 GFLOPs

Class	Images	Instances	P	R	mAP50	mAP50-95
all	254	443	0.939	0.917	0.938	0.581
mask	254	357	0.965	0.95	0.976	0.644
no_mask	254	86	0.913	0.884	0.901	0.518

Results saved to runs/train/exp2

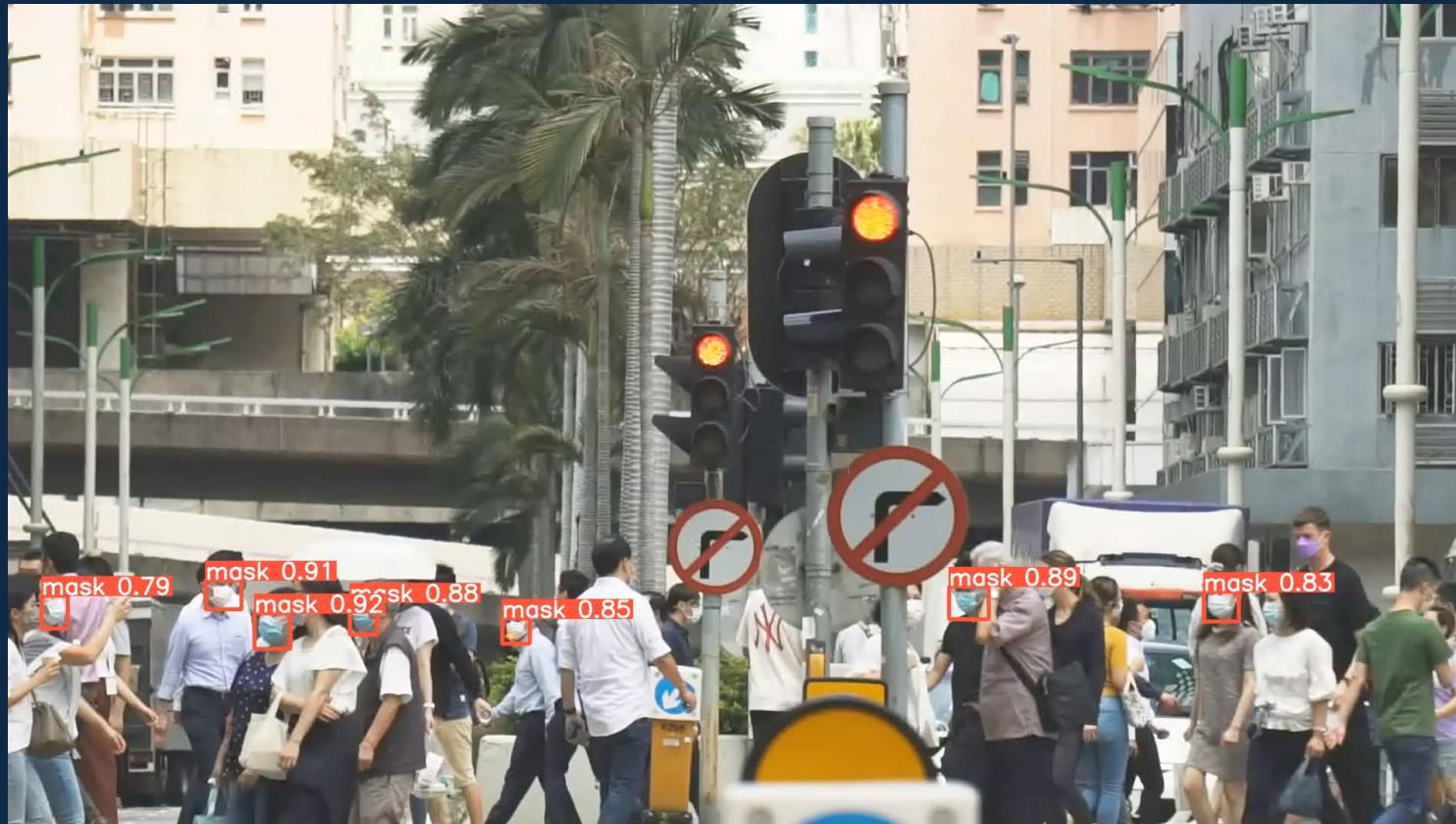
```
!python detect.py --source /content/drive/MyDrive/MacDocument/videos/video.mp4 --weights runs/train/exp2/weights/best.pt --conf 0.!
```

```
[ ] detect.py --source /content/drive/MyDrive/MacDocument/videos/REFINEDET.mp4 --weights runs/train/exp2/weights/best.pt --conf 0.5
```

```
[ ]
```

41m 26s completed at 20:57

# Detect Video



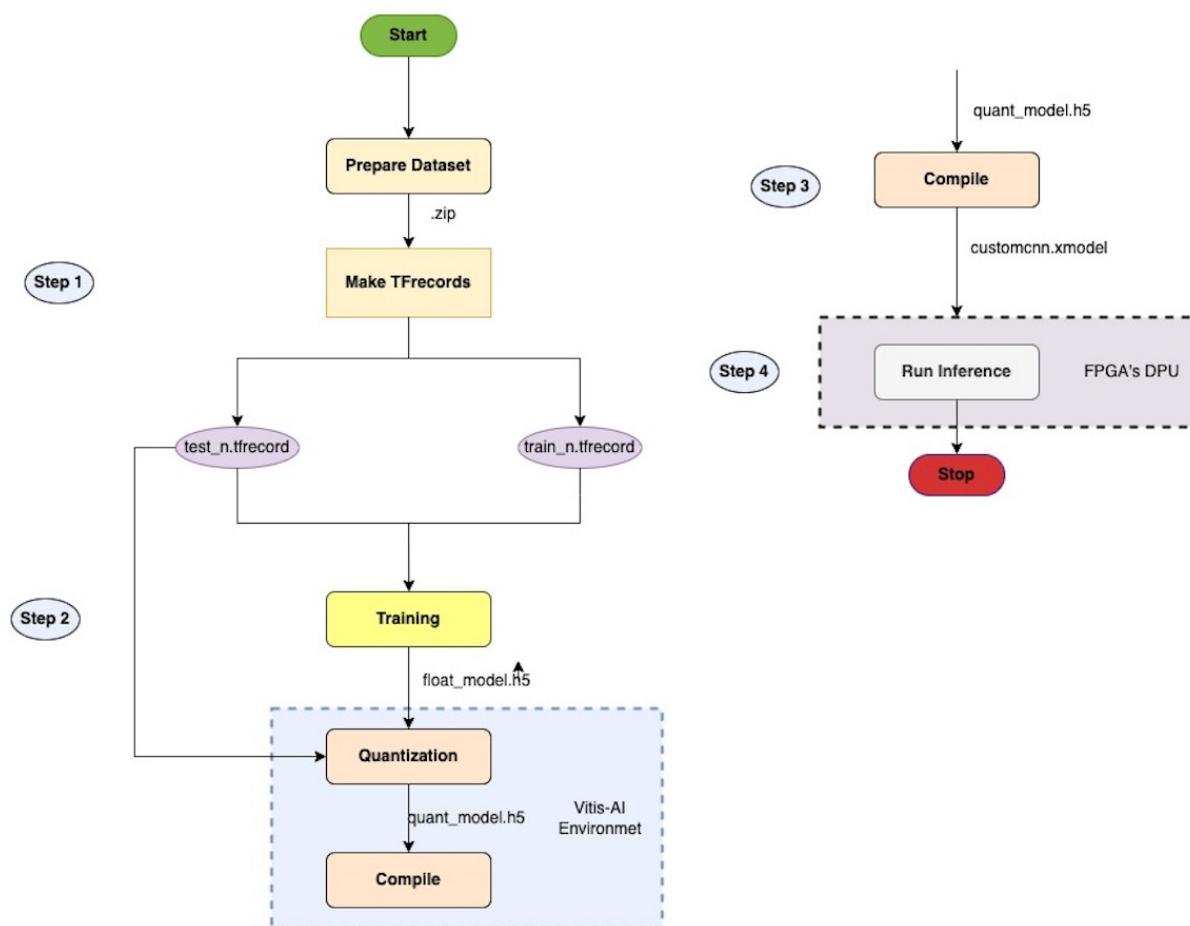
# Test performance on CPU and GPU

- ผลลัพธ์จากการนำโมเดล Yolov5 (detection) มาทดสอบประสิทธิภาพบน CPU กับ GPU

ลำดับที่	สถาปัตยกรรม hardware	FPS(fps)
1	Intel(R) Xeon(R) (Colab)	3.6193
2	NVIDIA V100 (Colab)	64.1026

Note: FPS = 1/(Inference time)

# Training Object Classification



# Prepare Dataset (Object Classification)

1. Download dogs-vs-cats dataset from Kaggle.
2. dogs-vs-cat dataset consists of 25000 images of varying dimensions, divided into the two classes of cat and dog
3. Each image is intrinsically labelled or classified by its filename

Example: the image with filename cat.12.jpg is obviously of class cat

# Training Object Classification model

The screenshot shows a terminal window with a dark blue background. On the left is a file explorer sidebar titled "Files". The main area displays a log of training output. The log includes several epochs (96, 97, 98, 99, 100) with learning rate scheduler updates, loss values, accuracy, and validation metrics. It also mentions saving a model to "build/float\_model/f\_model.h5". A message at the bottom indicates TensorBoard can be opened with the command: tensorboard --logdir=build/tb\_logs --host localhost --port 6006.

```
+ Code + Text
Epoch 96/100
350/350 [=====] - ETA: 0s - loss: 0.2420 - accuracy: 0.9391
Epoch 96: val_accuracy did not improve from 0.94660
350/350 [=====] - 116s 331ms/step - loss: 0.2420 - accuracy: 0.9391 - val_loss: 0.6577 - val_accuracy: 0.7
Epoch 97: LearningRateScheduler setting learning rate to 0.0001.
Epoch 97/100
350/350 [=====] - ETA: 0s - loss: 0.2359 - accuracy: 0.9418
Epoch 97: val_accuracy did not improve from 0.94660
350/350 [=====] - 126s 359ms/step - loss: 0.2359 - accuracy: 0.9418 - val_loss: 0.2629 - val_accuracy: 0.9
Epoch 98: LearningRateScheduler setting learning rate to 0.0001.
Epoch 98/100
350/350 [=====] - ETA: 0s - loss: 0.2397 - accuracy: 0.9410
Epoch 98: val_accuracy improved from 0.94660 to 0.94740, saving model to build/float_model/f_model.h5
350/350 [=====] - 117s 335ms/step - loss: 0.2397 - accuracy: 0.9410 - val_loss: 0.2269 - val_accuracy: 0.9
Epoch 99: LearningRateScheduler setting learning rate to 0.0001.
Epoch 99/100
350/350 [=====] - ETA: 0s - loss: 0.2332 - accuracy: 0.9448
Epoch 99: val_accuracy did not improve from 0.94740
350/350 [=====] - 116s 332ms/step - loss: 0.2332 - accuracy: 0.9448 - val_loss: 0.5571 - val_accuracy: 0.8
Epoch 100: LearningRateScheduler setting learning rate to 0.0001.
Epoch 100/100
350/350 [=====] - ETA: 0s - loss: 0.2396 - accuracy: 0.9401
Epoch 100: val_accuracy did not improve from 0.94740
350/350 [=====] - 116s 331ms/step - loss: 0.2396 - accuracy: 0.9401 - val_loss: 0.2451 - val_accuracy: 0.9
TensorBoard can be opened with the command: tensorboard --logdir=build/tb_logs --host localhost --port 6006
```

# Train Classification model

Command line options:

```
--input_height : 200  
--input_width  : 250  
--input_chan   : 3  
--tfrec_dir   : build/tfrecords  
--batchsize   : 50  
--learnrate   : 0.001  
--epochs      : 100  
--chkpt_dir   : build/float_model
```

---

Output :

- f\_model.h5
- accuracy: 0.9401
- val\_accuracy: 0.9390

## Quantization Classification model

```
-----  
Command line options:  
--float_model    : build/float_model/f_model.h5  
--quant_model    : build/quant_model/q_model.h5  
--batchsize       : 50  
--tfrec_dir       : build/tfrecords  
--evaluate        : True  
-----
```

\*\* Quantized model was saved as: q\_model.h5

# Compile Classification model

```
-----  
Command line options:  
--target_dir   : ./build/target_zcu104  
--image_dir    : build/dataset/test  
--input_height : 200  
--input_width  : 250  
--num_images   : 1000  
--app_dir      : application  
--model        : ./build/compiled_zcu104/customcnn.xmodel  
-----
```

\*\* Compiled model was saved as: customcnn.xmodel

# Inference on FPGA board

```
Command line options:  
--image_dir : images  
--threads    : 1  
--model      : customcnn.xmodel
```

```
-----  
Pre-processing 1000 images...
```

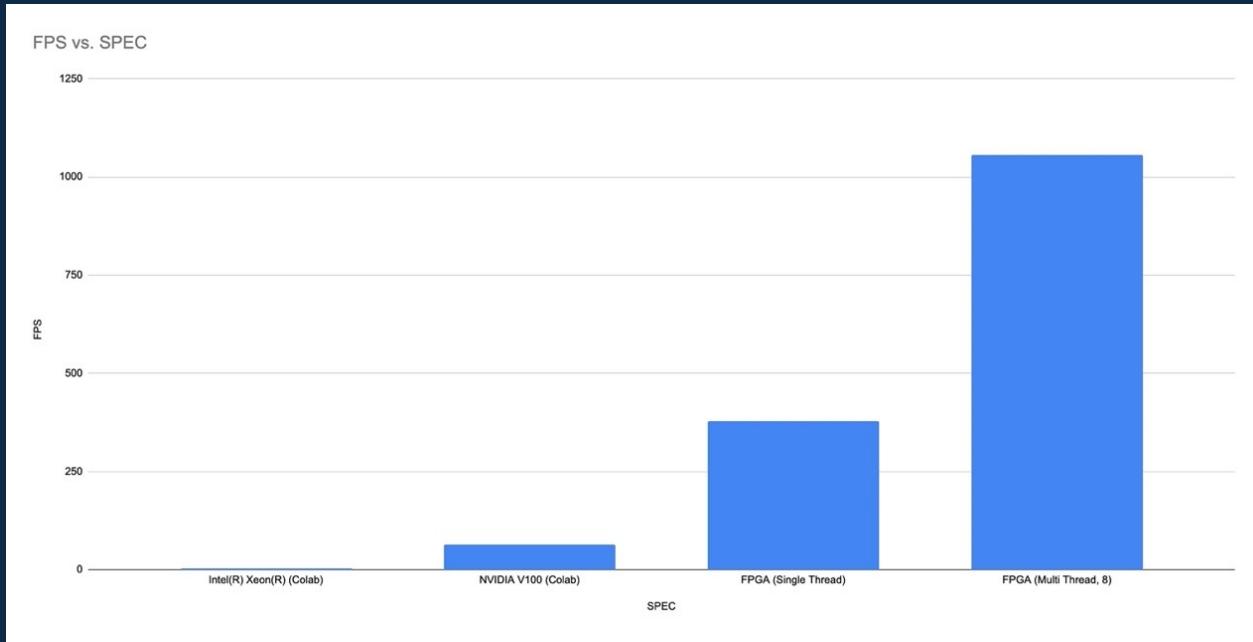
```
Starting 1 threads...
```

## \*\* Output

- ผลลัพธ์จากการนำโมเดล Classification มาทดสอบประสิทธิภาพบน FPGA

ลำดับที่	สถาปัตยกรรมฮาร์ดแวร์	FPS(fps)
1	FPGA (Single Thread)	377.79
2	FPGA (Multi Thread, 8)	1055.84

# Performance Comparative graph



\*\* Comparing each models on a graph

# THANKS!

Do you have any questions?

sabolsocare1028@gmail.com

+66 805535324

62050104@go.buu.ac.th



CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.