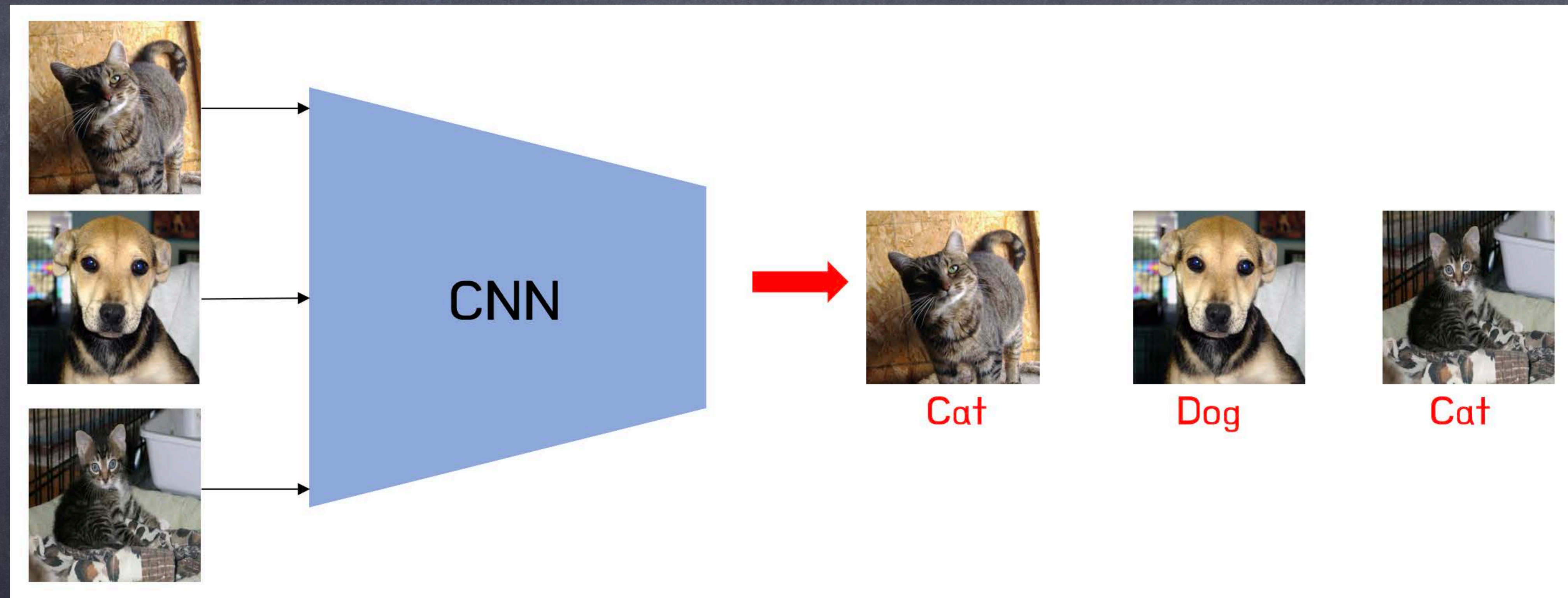


# Cascade models in Computer Vision

Argo Saakyan - Computer Vision Researcher



# Classification





# Object Detection task





# Difficulties

- Need high FPS for realtime processing
- 24/7 inference
- Accuracy/speed trade-off



# Cascade example

Grab the frame



**Detector - YOLOv5**

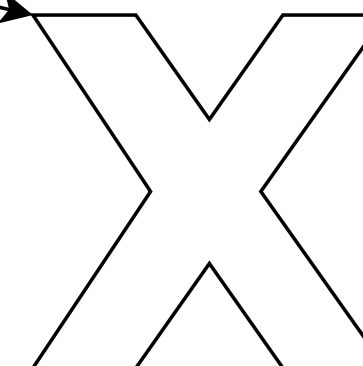
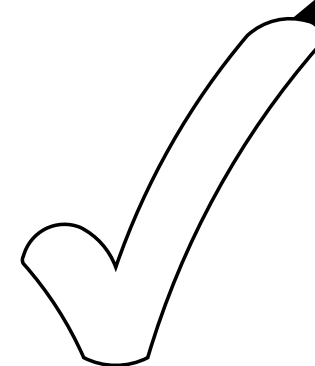
Can detect a gun,  
sometimes gives FP



Crop goes to second neural net



**Classifier -  
EfficientNet**





# Why does it work?

## Speed:

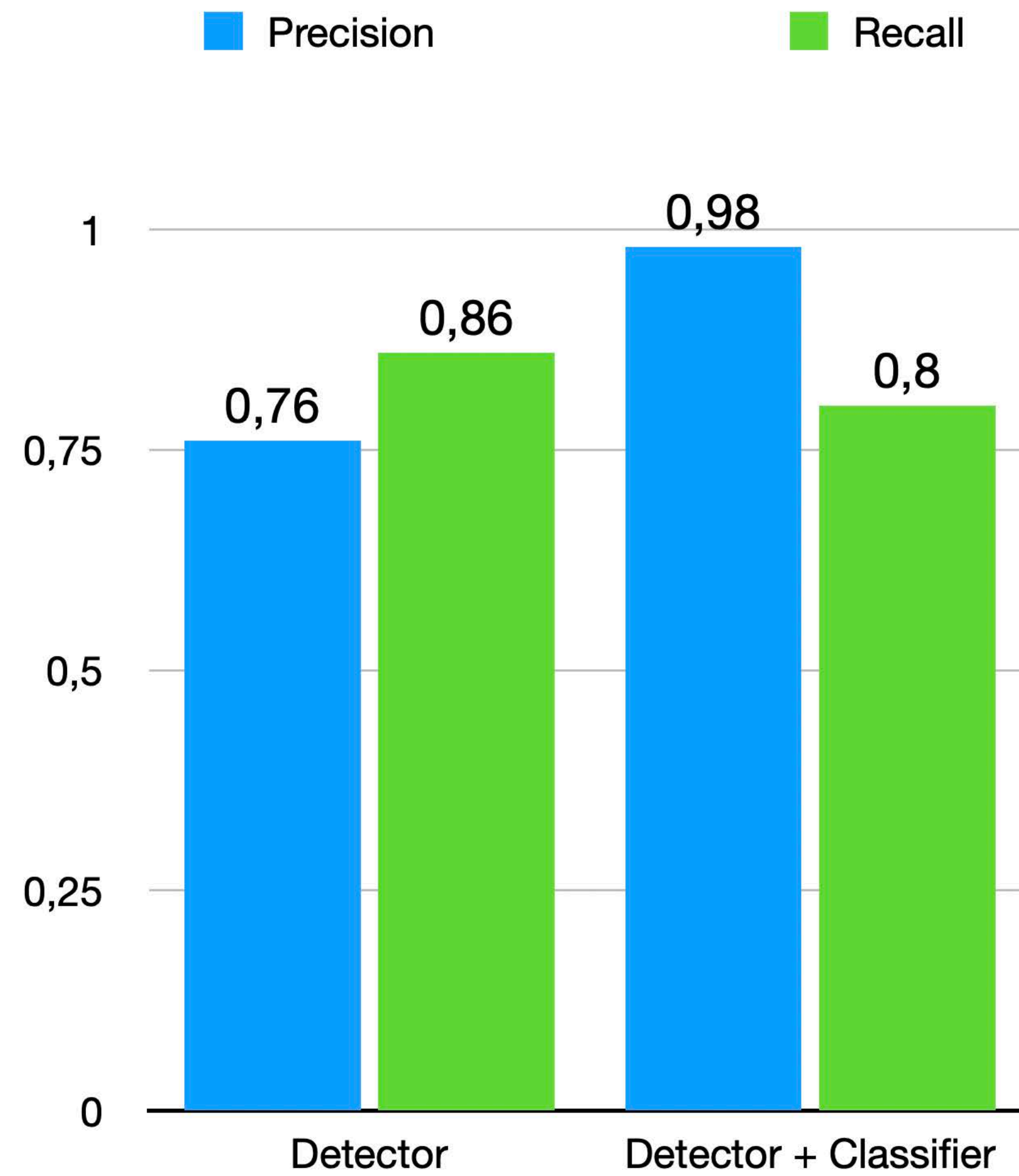
- Speed of the solution = speed of the detector
- Can choose fast detector and not lose accuracy

## Accuracy:

- Ensemble
- Higher res validation
- Different data
- Best of both worlds
- Retraining



# Metrics





# Dataset

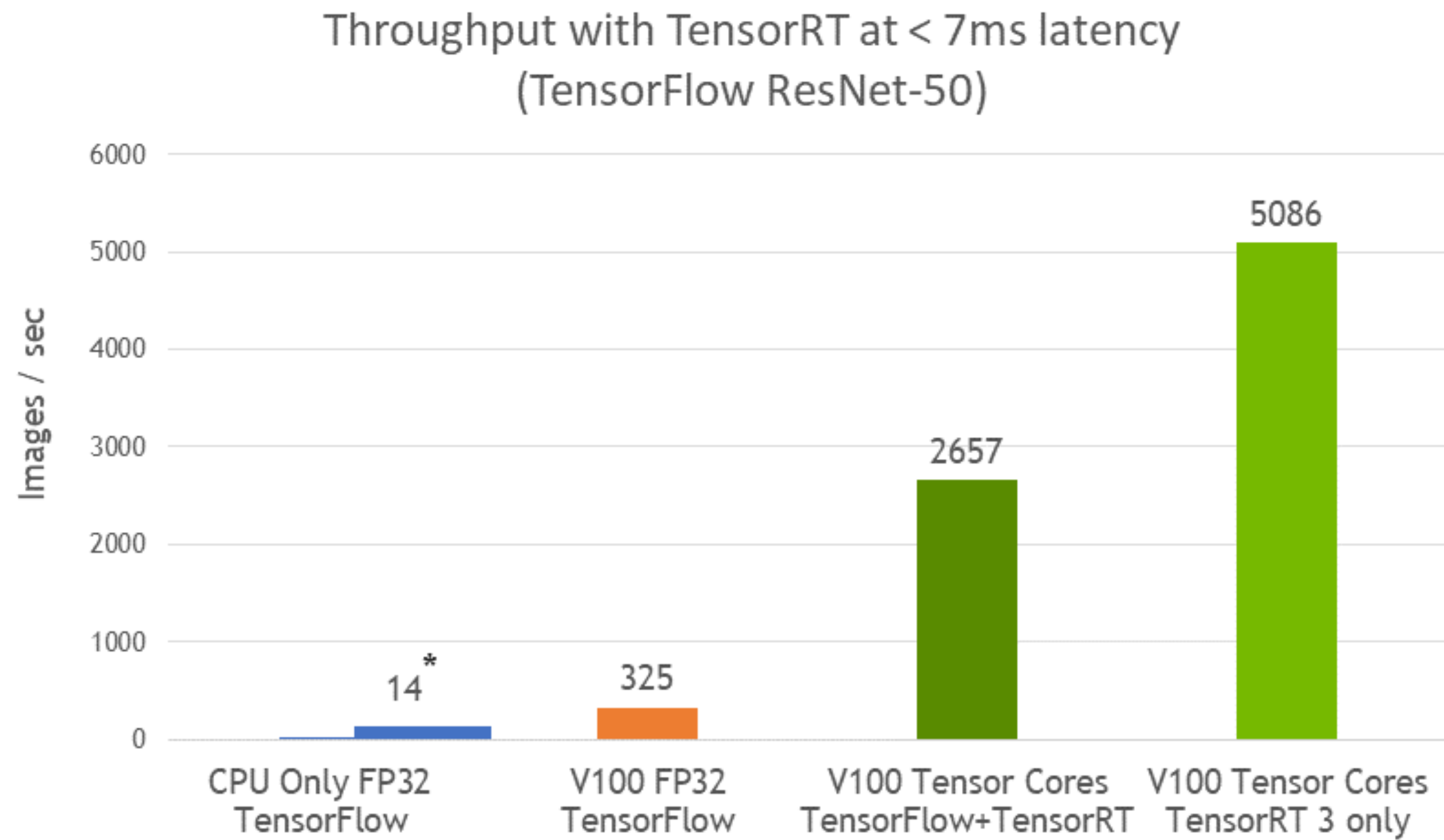


- Detector - close to inference
- Classifier - crops from detector





# Deployment

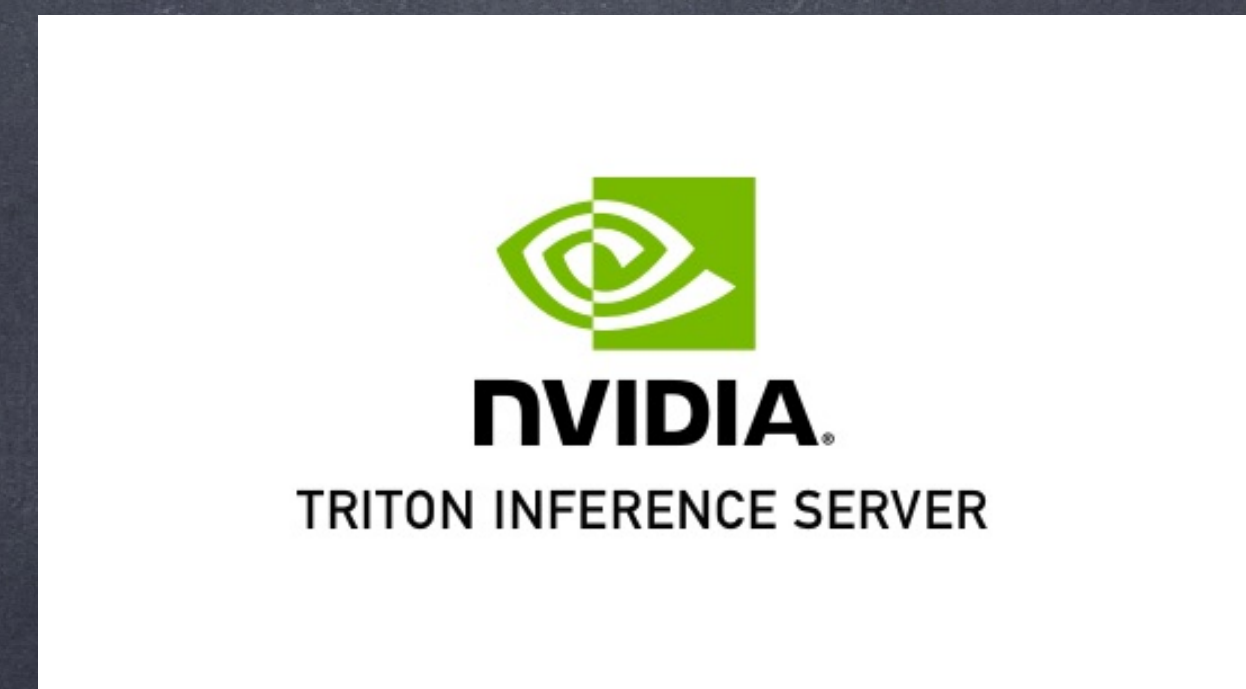
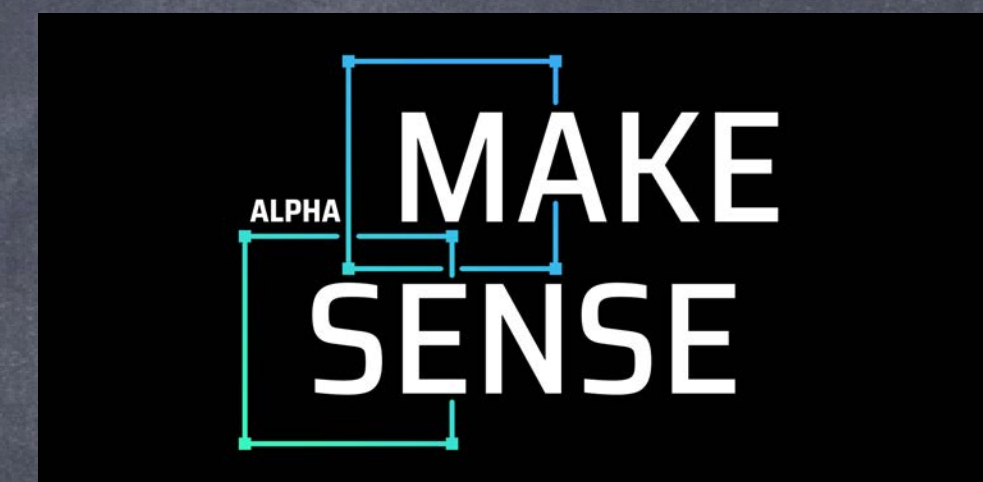


Updated 3/28/2018. \* Min CPU latency measured was 70 ms. It is not < 7ms.  
CPU: Skylake Gold 6140, Ubuntu 16.04, 18 CPU threads. Volta V100 SXM; CUDA (384.111;V9.0.176);  
Batch sizes: CPU=1;V100\_FP32=2; V100\_TensorFlow\_TensorRT=16; V100\_TensorRT=32; Latency=6ms. TensorRT 3.  
Latest results at: <https://developer.nvidia.com/deep-learning-performance-training-inference>



# Examples

- YOLOv8 - Detector
- Makesense.ai - Labeling
- EfficientNet - Classifier
- Deploy with Triton Server





# Thank you!

• [linkedin.com/in/argo-saakyan/](https://www.linkedin.com/in/argo-saakyan/)

