

DNN Video Analysis OS

George Chen(陳澤世)

Division for AI Computing Platform(X000)

ICL/ITRI



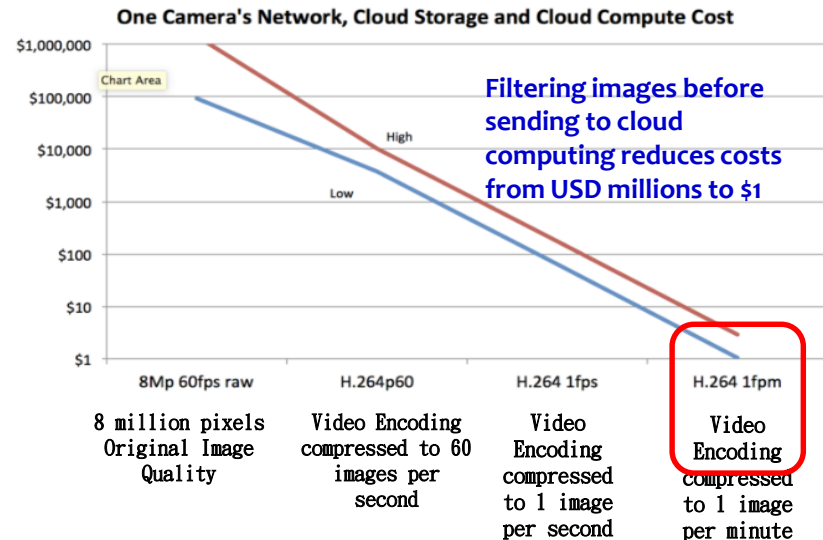
Cost of Image Streaming to Cloud Computing

- **Connecting Camera to Internet**
 - Optical Fiber : USD\$0.10/TB
 - Broadband : USD\$8-20/TB
 - T1 dedicated line : USD\$100/TB
- **Cloud Storage for Images**
 - Standard : USD\$12.50/TB per month
 - Long-term : USD\$2.5/TB per month
- **Image Analysis Cloud Computing**
 - YOLO on AWS GPU costs approx. USD\$0.58 per million images (equivalent to one day with 10 images per second)

Data Streaming	Network Costs		Storage Costs		Sub. Cost	3 yrs Costs
	OF	BB	Day	Month		
8Mp raw	9,500	950,000	79,000	476,000	3,300	92,000-1,430,000
H.264 p60	47	4,700	400	2,400	3,300	3,700-10,000
H.264 @ 1 fps	1	79	7	39	55	62-173
H.264 @ 1 fpm	0	1.30	0.11	0.66	0.91	1-3

Experience We learned

- **Video filtering on edge servers or front end devices greatly reduces the cost of network transmission/Storage/Computing Analysis**



Reference : <http://www.cogniteventures.com/2017/10/07/what-does-a-5-camera-cost/>

Video Recognition Speedup

- **AI inferences for real-time video:**
 - Inference video frames as **independent** images
 - Using open-source YOLO, one GPU card can inference 30 images per second
 - Thus, each 30 FPS camera costs one GPU card
 - **Very expensive**
- **Our video recognition improvement**
 - Using the similarity of neighboring frames to speedup AI
 - Current version can be applied for 1. DNN Object Detection 2. Classification + Localization
 - We are not developing for 1. DNN Object Classification and 2. Object Segmentation



Video AI Recognition PaaS Platform

- **DNN video recognition system**

- As a PaaS platform, integrating large number of cameras and scheduling AI jobs on multiple GPU hardware

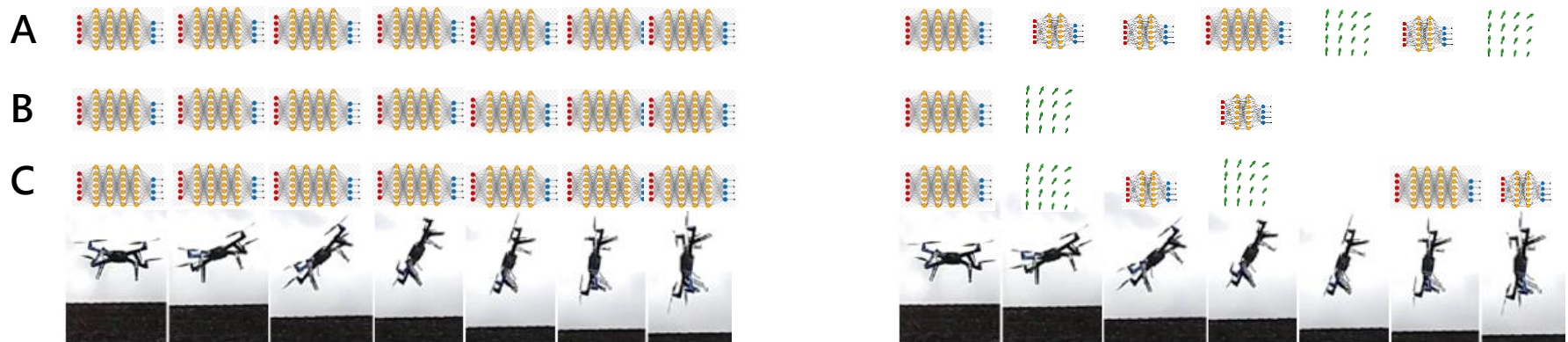
- **Basic version**

- Run M full-version AI model on every frames

- **Improved version:**

- Run full-version, partial version AI model, and/or speedup version, on every frames

models



Basic

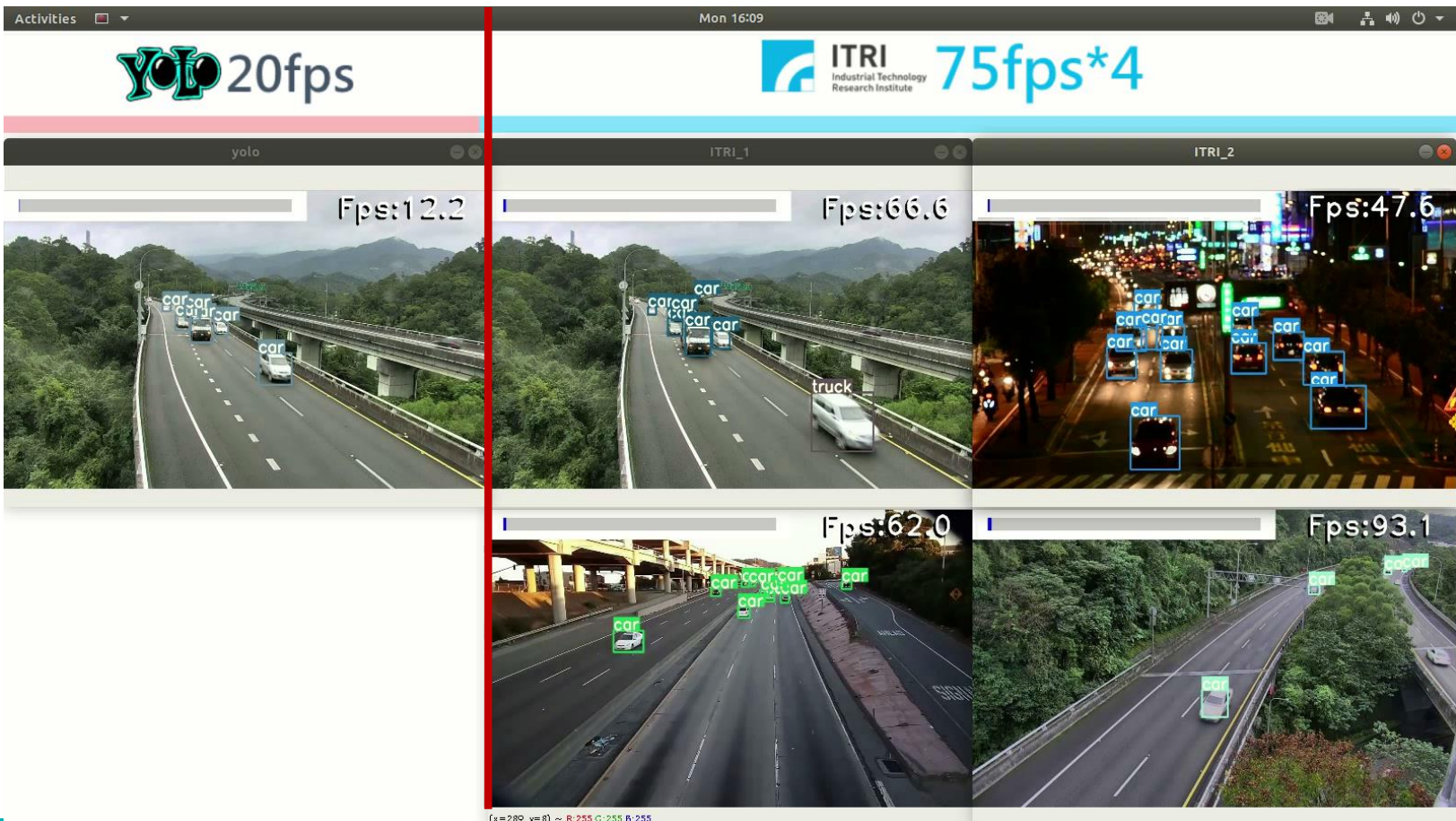
Improved

Performance Improvements

- Our technology speedup video inferencing, by interleaves using GPU and CPU, and supports multiple video streaming processes in parallel.

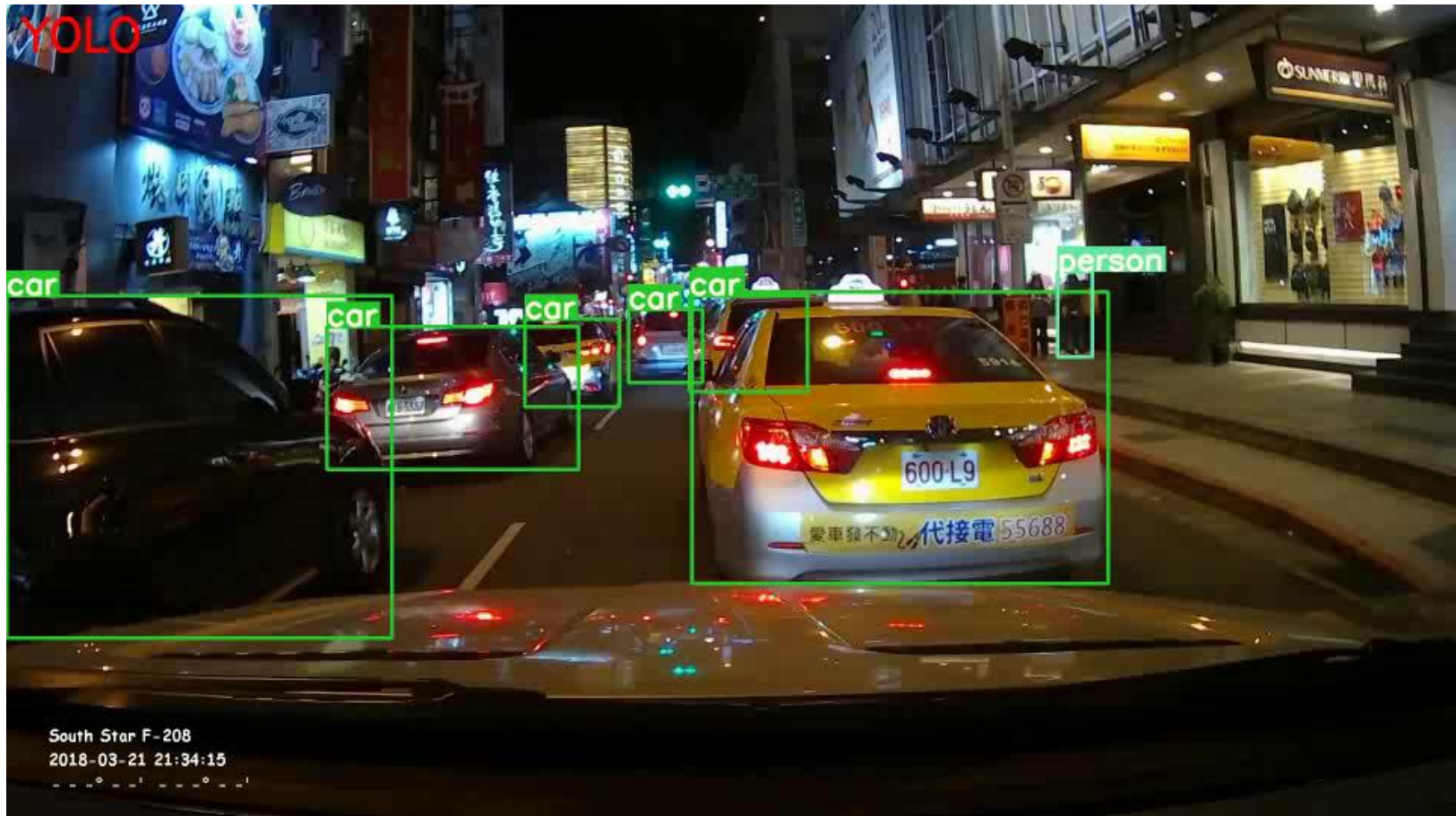
– Left: Open-source YOLO (20 FPS)

Right: Our method (Total 300 FPS)



Accelerated Results of Image Recognition and Tracking

- Processing Speed of Videos is 4.1 times faster than original Yolov4
- Accuracy improves from 61.27% to 61.02%

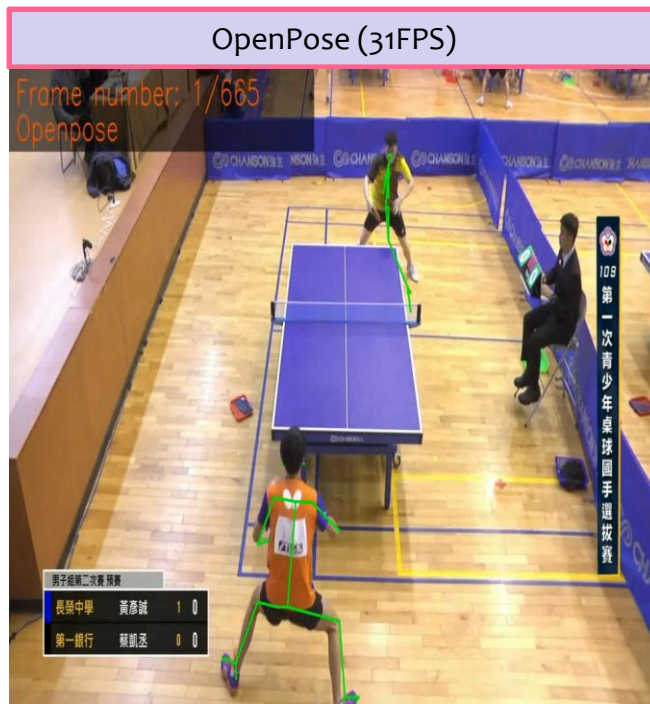


Remarks : DEMO video played with 30 FPS

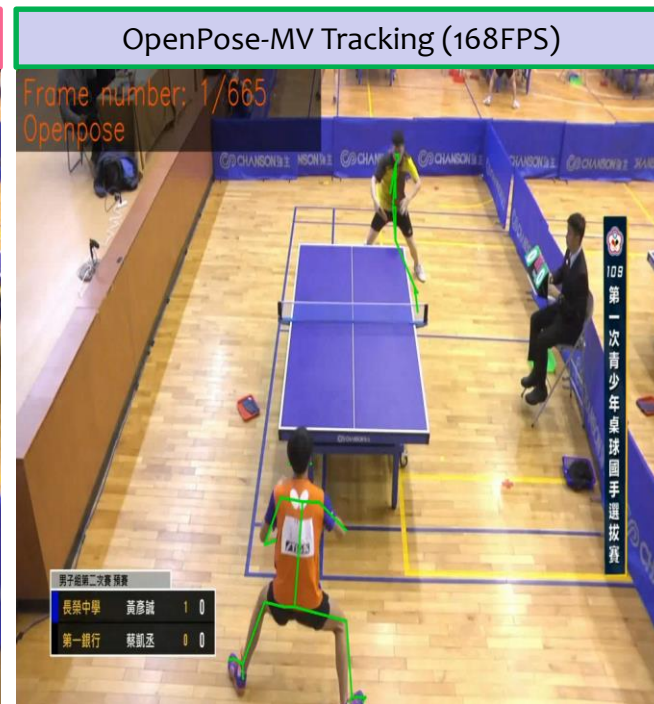
Human Body Joint Point Test

- Table tennis competition as example, processing speed from 31 FPS increased to 168 FPS which is 5.4 faster.
 - Accuracy increased from 54.8% to 56.3%

- **AI Body Joint Points**



- **AI Body Joint Point +Tracking**



ITRI DNN Farm

機器學習 & 深層學習框架 (Framework)

FY108-110



FY109-110



theano

NVIDIA DIGITS



虛擬化資源管理 (Docker / Resource Management)

FY108



NFS

FY109



FY110



硬體設施 / 網路環境 (Hardware / Network Infrastructure)

- Infra
- Compute
- GPU

Commodity Server
 - Tesla V100 x4
 - Tesla P100 16GB x2
 - Tesla P40 24GB x2
 - Titan V x16
 - Titan XP x4

ISO 27001

 ≥ 28 GPUs

- Infra
- Compute
- GPU
- Infra
- Compute
- GPU

Accelerated Server with V100

Accelerated Server with V100

基本資安控管
 DNN Farm
 ≥ 64 GPUs



George Chen
Tschen@itri.org.tw

