

# #GOTOcph





BEST PRACTICES FOR REAL-TIME INTELLIGENT VIDEO ANALYTICS DR. EKATERINA SIRAZITDINOVA, NOVEMBER 2021



# WHY INTELLIGENT VIDEO ANALYTICS?











### WHY INTELLIGENT VIDEO ANALYTICS? Every industry is going through rapid innovation to bring intelligent insights





#### PUBLIC SPACES







#### HEALTHCARE & HOSPITALS



#### **TRANSPORTATION HUBS**







# CHALLENGES WITH INTELLIGENT VIDEO ANALYTICS



#### **THROUGHPUT** Achieving real-time, low latency results

















































-



# TIPS AND TRICKS FOR EFFICIENT AI VIDEO ANALYTICS



#### FEATURE TRANSFER





#### **TRANSFER LEARNING** Transferring learned features from one model to another



# Less data required to train accurately

**KEY BENEFITS** 



#### Reduce training time and cost





### DATA AUGMENTATION Extending the dataset by applying simple transformations

#### SPATIAL



#### COLOR

#### Color shift

#### Hue rotation





#### Saturation



# Contrast adjustment



🧼 NVIDIA.

Reduce memory requirements 

- Speed-up memory intensive operations by using half the bytes
- Speed-up math intensive operations by using Tensor Cores
- Train with half-precision while maintaining network accuracy as with single precision in order to:





# AUTOMATIC MIXED PRECISION (AMP)

Train larger models or larger batches









# Post Training Quantization (PTQ) for quantization after training is done

#### Quantization Aware Training (QAT) for modelling quantization error from weights and tensors during training



Quantize(x, r) = round(s \* clip(x, -r, r)) where r = |Max| and s = 127 / r

#### QUANTIZATION Reducing bits per weight

PeopleNet ResNet34

PeopleNet ResNet18

0



2x throughput Increase



#### 2-step process

- Reduce model size
- 2 Incrementally retrain model after pruning to recover accuracy

6 inputs, 6 neurons, 32 connections

### **NETWORK PRUNING** Reducing the number of weights

#### 6 inputs, 5 neurons, 24 connections



#### Network - ResNet18 4-class





#### Before layer fusion



### NETWORK GRAPH OPTIMIZATIONS Layer & tensor fusion





#### After horizontal and vertical layer fusion



- Target hardware platform
- Batch size
- Input dimensions
- Filter dimensions
- Tensor layout
- Specific algorithms implementation
- 100s of specialized kernels otimized for every GPU platform





TX2

Xavier NX

### **KERNEL AUTO-TUNING** Picking right algorithms depending on your deployment hardware





Xavier AGX



x86 GPUs



### **DYNAMIC TENSOR MEMORY UPON INFERENCE** Reduced network memory footprint & improved memory re-use

- Combining tensors into regions
  - Region lifetime is a section of network execution time
- Assigning regions to blocks
  - Regions assigned to a block have disjoint lifetimes



Similar to *register allocation in* compiling, the process of assigning a large number of target program variables onto a small number of registers



### Region 3

Tensor E



### **MULTISTREAM CONCURENT EXECUTION** For better GPU utilization and higher throughput

# Serial

#### Concurrent

H2D



#### Memory copy (D2H)

#### Performance improvement

Execution time



# FREE NVIDIA PRODUCTS DESIGNED TO MAKE YOUR AI APPLICATIONS EFFICIENT



# platform company."

Jensen Huang, CEO of NVIDIA

"NVIDIA is not a GPU company. It's a





#### PRE-TRAINED MODEL LIBRARY



**VEHICLE DETECTION** 

NLP + ASR



000

POSE ESTIMATION

LICENSE PLATES



FACE DETECT

NVIDIA GPU Cloud

### NVIDIA'S END-TO-END AI WORKFLOW Develop & deploy production ready solutions



#### TAO TOOLKIT

### DEVELOPMENT AND DEPLOYMENT

#### Turnkey apps

#### Development environment

#### **DEEPSTREAM SDK**



Jetson appliances

1	$\prec$	
	•	
	-!	

EGX servers





# TAO TOOLKIT

#### **CONVERSATIONAL AI**

Prune	
Те	nsorRT
NCE PLATFO	RMS
	F 2 2-1
	Ampore
4	Ampere

**T**4





### HIGH PERFORMANCE PRE-TRAINED VISION AI MODELS Download for free from <a href="https://ngc.nvidia.com/">https://ngc.nvidia.com/</a>

- Optimized for high throughput
- Trained for >80% accuracy
- **Production-ready**
- Adaptable with NVIDIA TAO



People detection



Gaze estimation





Heart rate estimation



Emotion recognition





Vehicle & pedestrian detection



License plate detection



License plate recognition

Facial landmark



People segmentation



Gesture recognition

Pose estimation



Face detect IR



Dash camera vehicle detection



Vehicle make net



Vehicle type net







### ENABLING BEYOND PRE-TRAINED AI MODELS 100+ combinations of model architectures and backbones

	Image Classification	Object Detection						Segmentation		
		DetectNet_V2	FasterRCNN	SSD	YOLOV3	YOLOV4	RetinaNet	DSSD	MaskRCNN	UNET
ResNet10/18 /34/50/101	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$
VGG16/19	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$
GoogLeNet	$\checkmark$			$\checkmark$	$\checkmark$	$\checkmark$				
MobileNet V1/V2	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
SqueezeNet	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
DarkNet 19/53	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
CSPDarkNet 19/53	$\checkmark$					$\checkmark$				
Efficient Net B0/B1	$\checkmark$		$\checkmark$	$\checkmark$			$\checkmark$	$\checkmark$		

Pre-trained weights trained on OpenImage dataset



# ACHIEVING STATE OF THE ART ACCURACY FOR PUBLIC DATASETS





#### NVIDIA TENSORRT TensorRT Optimizer & TensorRT Runtime

#### Layer and Tensor Fusion

#### Weight and Activation **Precision Calbration**

#### **Kernel Auto-Tuning**

#### **Dynamic Tensor Memory**

#### **Multi-Stream Execution**

![](_page_22_Picture_9.jpeg)

# Step 1: Optimize trained model

![](_page_23_Picture_1.jpeg)

**Trained Neural** Network

# Step 2: Deploy optimized plans with runtime

![](_page_23_Picture_4.jpeg)

#### **Optimized Plans**

### **TENSORRT WORKFLOW** Optimize and deploy

![](_page_23_Picture_7.jpeg)

TensorRT Optimizer

![](_page_23_Picture_9.jpeg)

### **Optimized Plans**

![](_page_23_Picture_12.jpeg)

![](_page_24_Figure_1.jpeg)

### **TRITON INFERENCE SERVER** Open-source software for scalable, simplified inference serving

![](_page_24_Figure_4.jpeg)

![](_page_24_Picture_6.jpeg)

![](_page_25_Figure_0.jpeg)

### DEEPSTREAM SDK

![](_page_25_Picture_2.jpeg)

### DEEPSTREAM APPLICATION ARCHITECTURE End-to-end hardware accelerated pipeline

![](_page_26_Figure_1.jpeg)

![](_page_26_Picture_4.jpeg)

# PIPELINE EFFICIENCY WITH ZERO MEMORY COPIES

![](_page_27_Figure_1.jpeg)

![](_page_27_Figure_2.jpeg)

MEMORY)

![](_page_27_Picture_5.jpeg)

![](_page_28_Figure_0.jpeg)

#### NVIDIA GRAPH COMPOSER Drag + drop development environment

![](_page_28_Figure_2.jpeg)

150 building blocks

#### **Extensions:** collection of components

![](_page_29_Picture_0.jpeg)

![](_page_29_Picture_1.jpeg)

#### DEEPSTREAM VIDEO DEMO PeopleNet: detecting people

#### https://ngc.nvidia.com/catalog/models/nvidia:tlt\_peoplenet

![](_page_29_Picture_4.jpeg)

![](_page_29_Picture_5.jpeg)

![](_page_29_Picture_6.jpeg)

![](_page_30_Picture_0.jpeg)

![](_page_30_Picture_1.jpeg)

### DEEPSTREAM VIDEO DEMO DashCamNet and VehicleTypeNet in action

![](_page_30_Picture_3.jpeg)

Car 2805 largevehicle Ci Car Car 295 ( Car 2 Car 2303 sedan

https://ngc.nvidia.com/catalog/models/nvidia:tlt\_dashcamnet https://ngc.nvidia.com/catalog/models/nvidia:tlt\_vehiclemakenet

![](_page_30_Picture_6.jpeg)

![](_page_30_Picture_8.jpeg)

### DEEPSTREAM VIDEO DEMO Hand detector adapted with TAO Toolkit and GestureNet is used out the box

![](_page_31_Picture_1.jpeg)

https://github.com/NVIDIA-AI-IOT/gesture\_recognition\_tlt\_deepstream

![](_page_31_Picture_5.jpeg)

#### TRAIN

Transfer learning, AMP, multi GPU, data augmentation, pruning and quantization aware training with TAO Toolkit

![](_page_32_Picture_2.jpeg)

![](_page_32_Picture_3.jpeg)

#### **SUMMARY** For efficient Al video analytics

OPTIMIZE

Layer fusion, kernel auto-tuning, dynamic tensor memory, etc. with TensorRT

![](_page_32_Picture_7.jpeg)

#### DEPLOY

### Concurrent execution, zero copies, integrated encoders and decoders with Deepstream SDK

![](_page_32_Picture_10.jpeg)

🧼 NVIDIA.

TAO Toolkit <a href="https://developer.nvidia.com/tao-toolkit">https://developer.nvidia.com/tao-toolkit</a> DeepStream SDK <a href="https://developer.nvidia.com/deepstream-getting-started">https://developer.nvidia.com/deepstream-getting-started</a> TensorRT <a href="https://developer.nvidia.com/tensorrt">https://developer.nvidia.com/tensorrt</a> Triton Inference Server: <a href="https://github.com/triton-inference-server">https://github.com/triton-inference-server</a> NVIDIA GPU Cloud (NGC) https://ngc.nvidia.com/ Developer forums <a href="https://forums.developer.nvidia.com/">https://forums.developer.nvidia.com/</a> NVIDIA Deep Learning Institute https://www.nvidia.com/en-us/training/

#### **DEVELOPER RESOURCES** Powerful end-to-end AI video analytics made easy

![](_page_33_Picture_4.jpeg)

![](_page_33_Picture_6.jpeg)

![](_page_34_Picture_0.jpeg)

# EARNING DEEP LEARNING

Theory and Practice of Neural Networks, Computer Vision, Natural Language Processing, and Transformers Using TensorFlow

MAGNUS EKMAN

### THEORY AND PRACTICE OF NEURAL NETWORKS, COMPUTER VISION, NATURAL LANGUAGE PROCESSING, AND TRANSFORMERS USING TENSORFLOW

- Full-colour guide
- Illuminates both the core concepts and the hands-on programming techniques needed to succeed
- Shows how to build advanced architectures, including the Transformer
- Includes concise, well-annotated code examples using TensorFlow with Keras; corresponding PyTorch examples are provided online

Available at the GOTO Copenhagen Conference Bookstore

![](_page_34_Picture_11.jpeg)

![](_page_35_Picture_0.jpeg)

![](_page_35_Picture_1.jpeg)