# Accelerating the Creation of Custom, Production-Ready AI Models for Edge AI
## NVIDIA Tools, Part 1

Akhil Docca

Product Marketing, NVIDIA Corporation

# Building an AI Application Is Hard
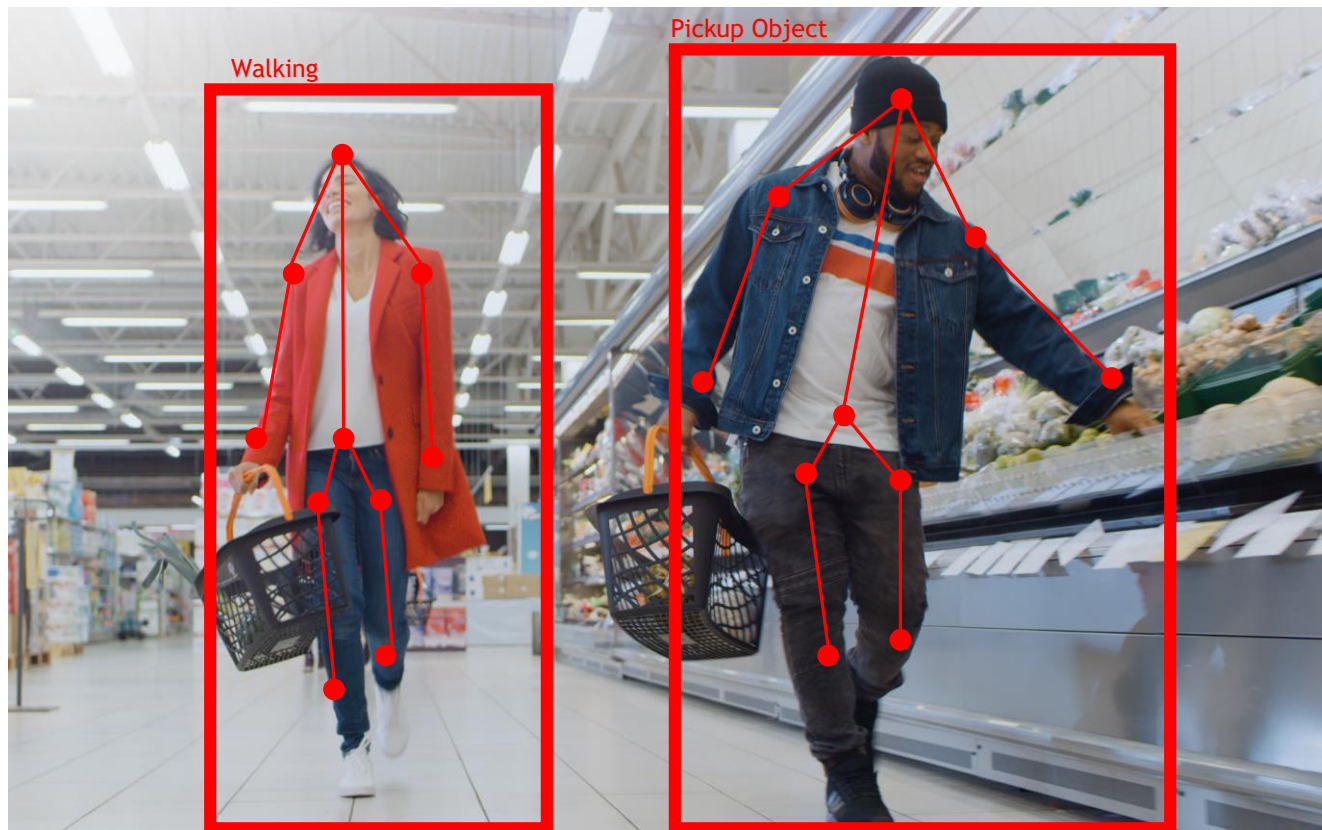
# Core Elements of an Autonomous Shopping App



| Customer enters the store | Picks up the item(s) of choice | Leaves the store and gets billed automatically |

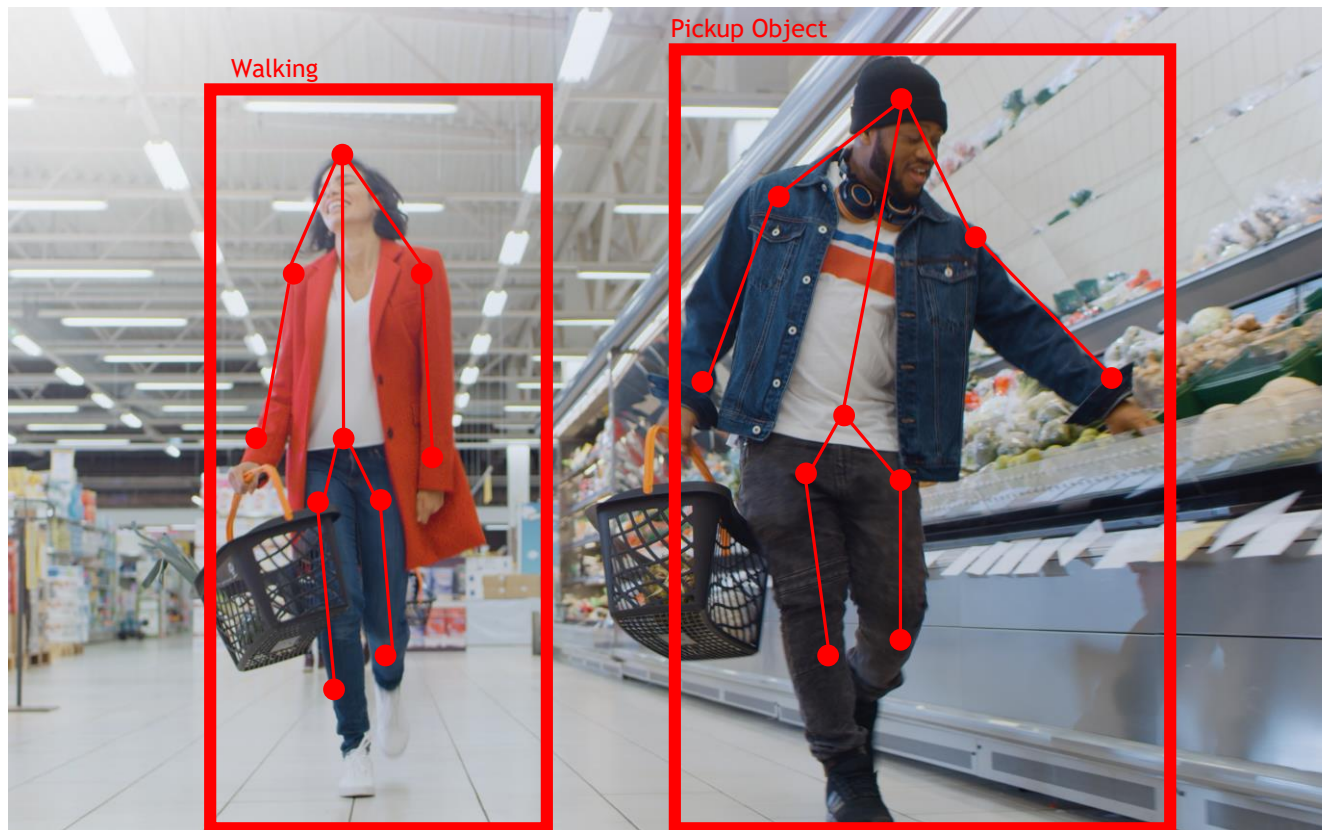# Core Elements of an Autonomous Shopping App



| Detect person | Detect pose based on key points | Recognize actions | Identify the item(s) picked |

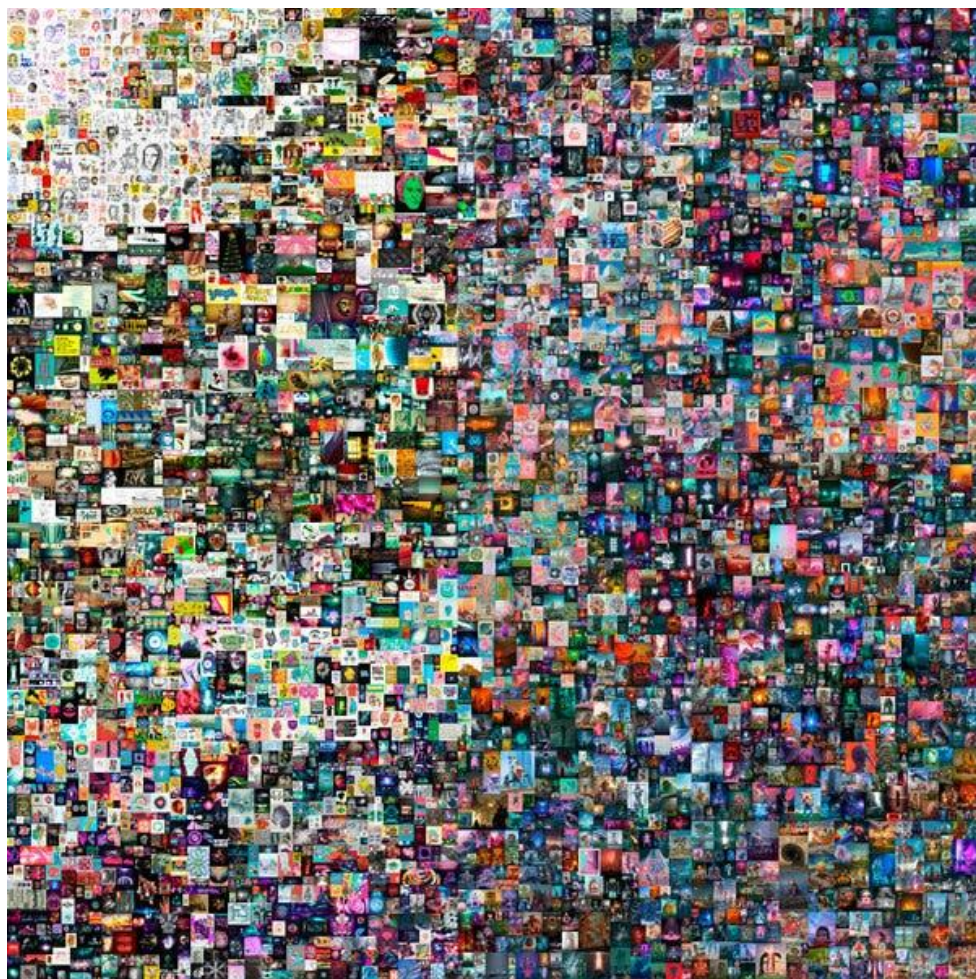# Core Elements of an Autonomous Shopping App



OPTION 1
Build and train a model from scratch

OPTION 2
Customize a model from a model zoo

# Option 1: Training from Scratch

First 5000 Days NFT (Image Courtesy: The New York Times)

PeopleNet – An NVIDIA built people detection model

Data
- 3.5M images
- 16M+ people in the images
- 40 people, 5 years to collect, curate and label

Training and Optimizing
- 30 Days A100 8 GPUs
- Several months before reaching production-quality
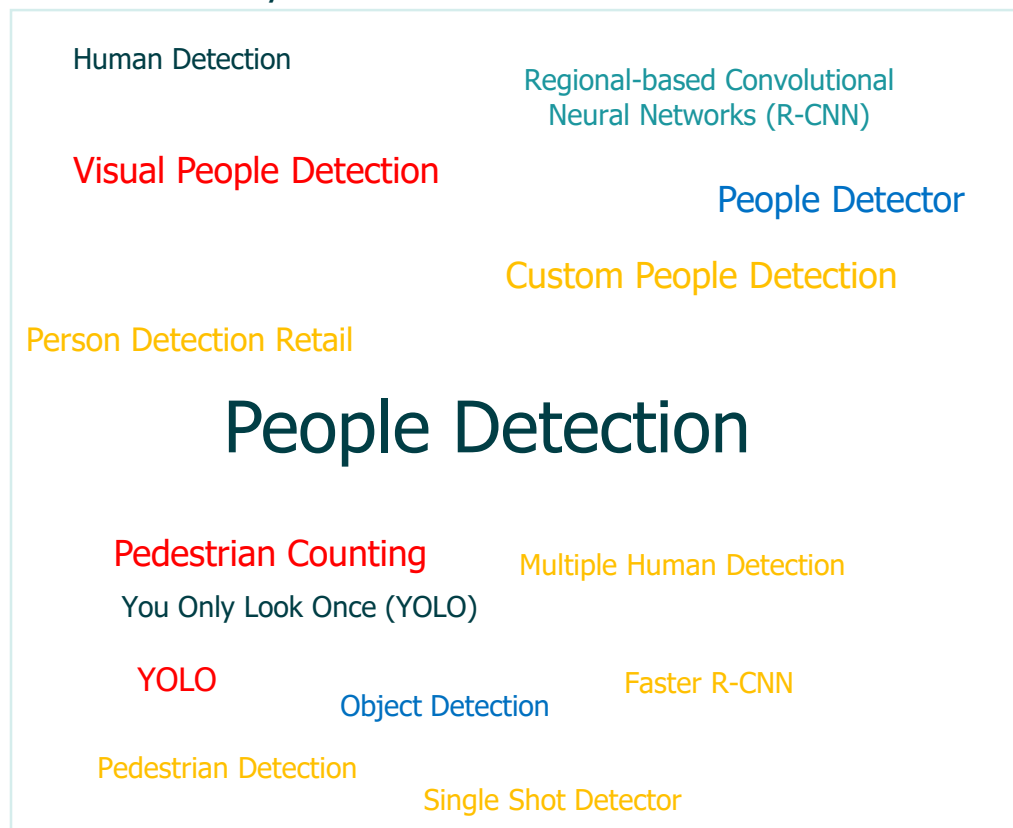
# Option 2: Train from Existing Models

Accuracy?     Dataset?     Performance?

Human Detection

Regional-based Convolutional Neural Networks (R-CNN)

Visual People Detection

People Detector

Custom People Detection

Person Detection Retail

## People Detection

Pedestrian Counting

Multiple Human Detection

You Only Look Once (YOLO)

YOLO

Object Detection

Faster R-CNN

Pedestrian Detection

Single Shot Detector

Write Code, Train, Iterate, Customize and Optimize

Caffe2   Chainer

K   mxnet

PaddlePaddle   PyTorch

TensorFlow

## Still requires lots of expertise and time!

NVIDIA.

# Industry Wide Challenge

LACK OF DATA

SKILLED PEOPLE

*IS THERE AN OPTION 3?*

# The NVIDIA TAO Framework



**NGC**

PRETRAINED MODELS

Pretrained Model

**TAO**

Your Data

Train    Adapt    Optimize

Your Custom Production Model

**MANY INDUSTRIES**

DeepStream      Riva

Triton

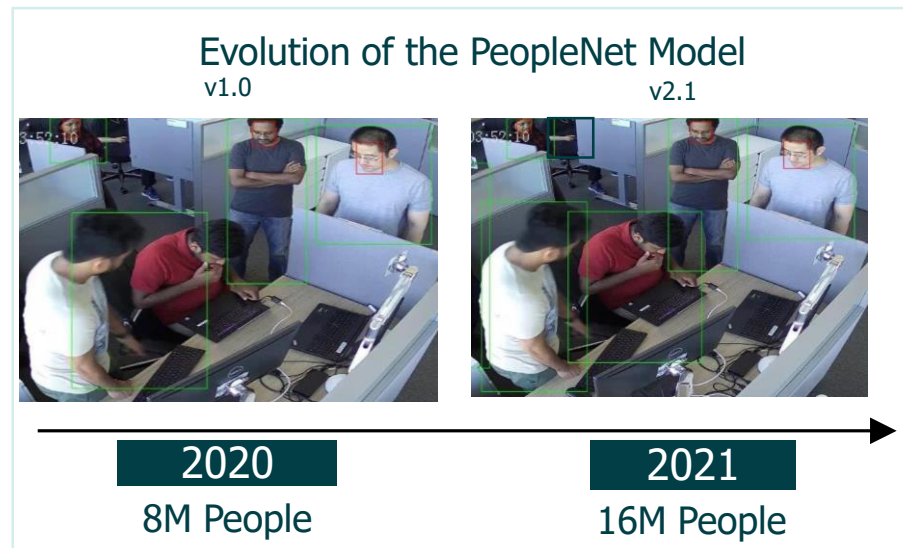| TRAIN EASILY | CUSTOMIZE FASTER | OPTIMIZE FOR DEPLOYMENT | INTEGRATE AND DEPLOY |
|---|---|---|---|
| Easy to use solutions that abstract away the AI framework complexity | Fine tune NVIDIA pre-trained AI models with fraction of the data as opposed to training from scratch | Optimize for low latency and high-throughput | Integrate the optimized models from TAO into DeepStream (Vision) and Riva (Speech) |

# The Power of Pretrained Models



### Evolution of the PeopleNet Model

v1.0                         v2.1

| 2020 | 2021 |
|------|------|
| 8M People | 16M People |

### Inference Throughput - FPS

PeopleNet - Pruned & Quantized    4X

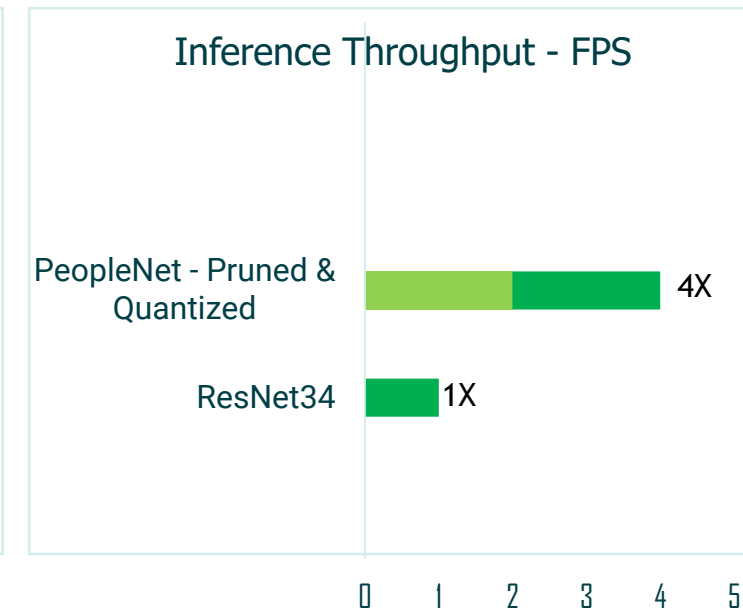ResNet34    1X

0  1  2  3  4  5

## WIDE RANGE OF USE CASES

100+ permutations of NVIDIA-optimized model architectures (EfficientDet, YOLOv3/v4)

Task based models - People Detection, Vehicle Detection, Gaze, Speech Recognition and Text to Speech
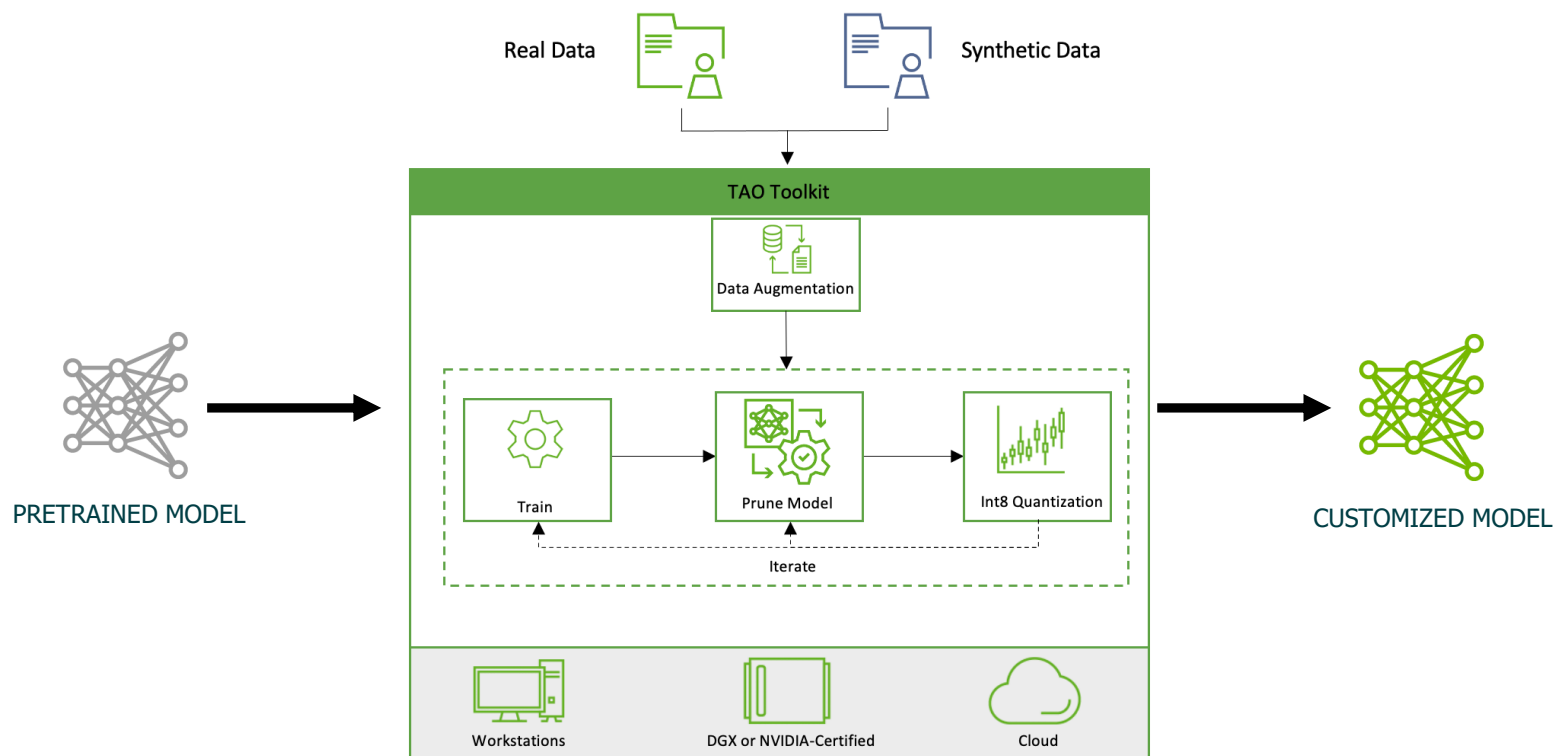
## HIGHLY ACCURATE

Trained and continuously updated by experts so you can adapt to your domain or deploy as-is

## OPTIMIZED FOR INFERENCE

Deploy in the data center or at the edge

# A Closer Look at the NVIDIA TAO Toolkit



RUN ANYWHERE

Container based solution

TURNKEY JUPYTER NOTEBOOKS

Computer Vision, Speech and NLU

DATA AUGMENTATION

Spatial and color transformation

MULTI-GPU AND MULTI-NODE
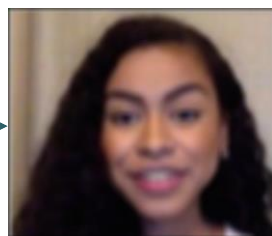
Accelerated model training

PRUNE AND QUANTIZE AWARE TRAINING

More than 4X speed up in inference

# Data Augmentation



**BLUR**

Gaussian Blur

**SPATIAL**
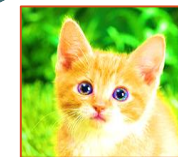
Vertical Flip | Horizontal Flip
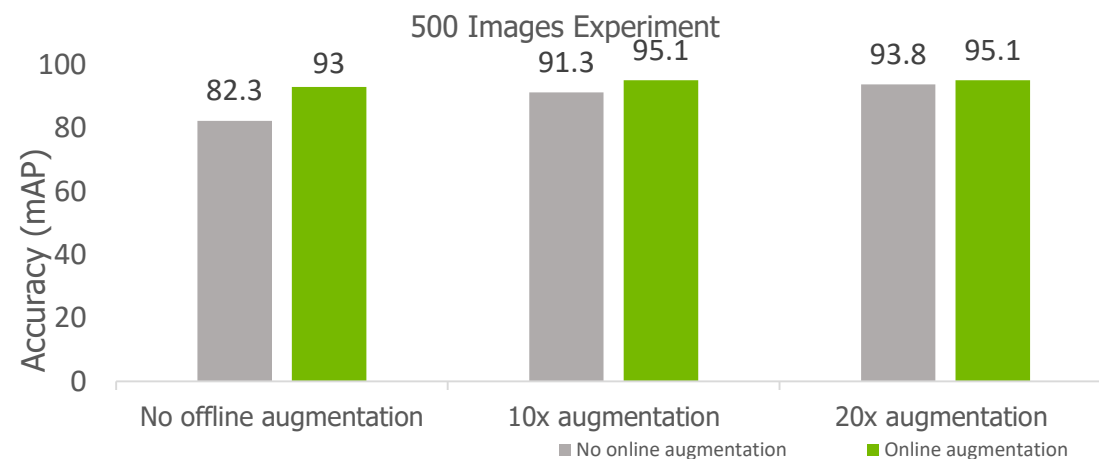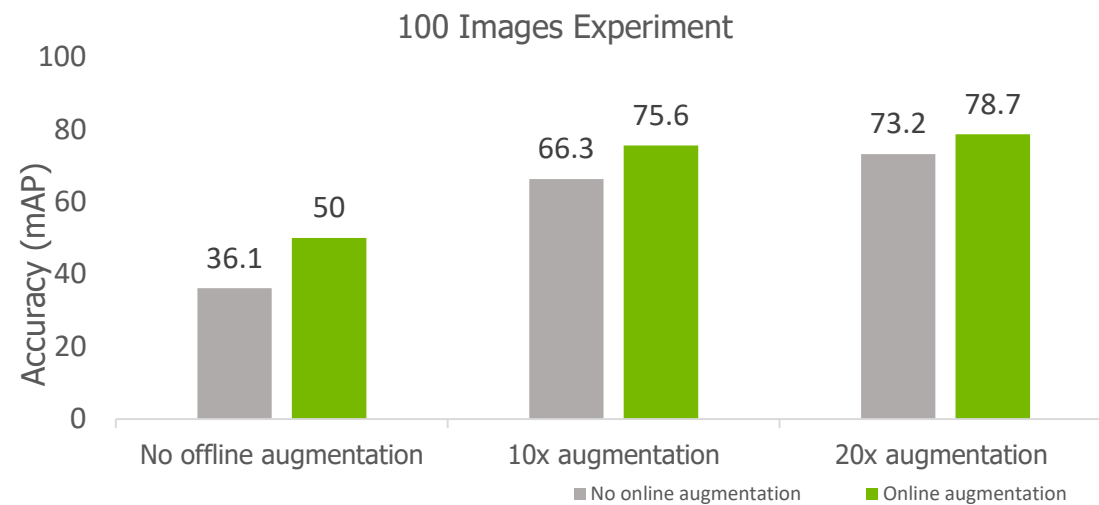
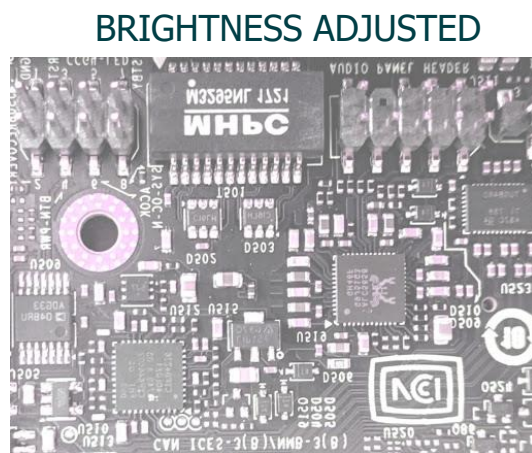Zoom | Shift

Rotate | Shear

**COLOR**

Color Shift | Hue Rotation

Saturation | Contrast Adjustment

# Data Augmentation

ORIGINAL

ROTATED

BRIGHTNESS ADJUSTED

## 100 Images Experiment

Accuracy (mAP)

| | No offline augmentation | 10x augmentation | 20x augmentation |
|---|---|---|---|
| No online augmentation | 36.1 | 66.3 | 73.2 |
| Online augmentation | 50 | 75.6 | 78.7 |

■ No online augmentation  ■ Online augmentation

## 500 Images Experiment

Accuracy (mAP)

| | No offline augmentation | 10x augmentation | 20x augmentation |
|---|---|---|---|
| No online augmentation | 82.3 | 91.3 | 93.8 |
| Online augmentation | 93 | 95.1 | 95.1 |

■ No online augmentation  ■ Online augmentation

TAO Toolkit White Paper

Git Hub Project Repository

13

# Art of the Possible – Camera Types



Thermal IR Data Set

TAO TOOLKIT

Train     Adapt     Optimize

PeopleNet Model

### Accuracy on IR dataset



65.6   78.14   81.2   82.2   82.7
44.8   63.4   70.9   73.4   77.3

Accuracy (mAP)

Number of Images

—●— Accuracy with PeopleNet

# Art of the Possible – Adding New Classes



Goal: Add "Helmet" class to existing People detect model

Labeled Helmet Dataset

PeopleNet Model

**TAO TOOLKIT**
Inference

People and face annotation

Complete annotated dataset

**TAO TOOLKIT**
Train | Adapt | Optimize

80% AP over 100 epochs

Helmet Class Average Precision

Model with Helmet, People and Face Detection
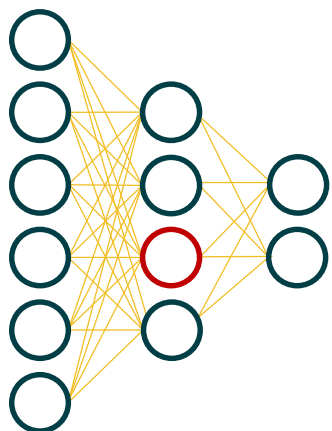
TAO Toolkit White Paper     Git Hub Project Repository
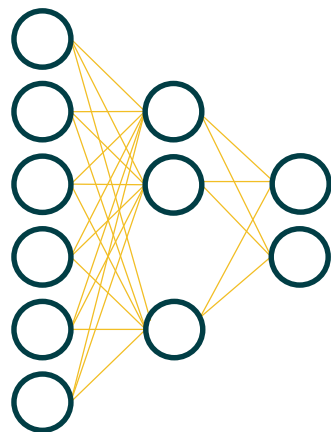
# Model Pruning

TrafficCamNet

## 2 STEP PROCESS

1. Reduce model size

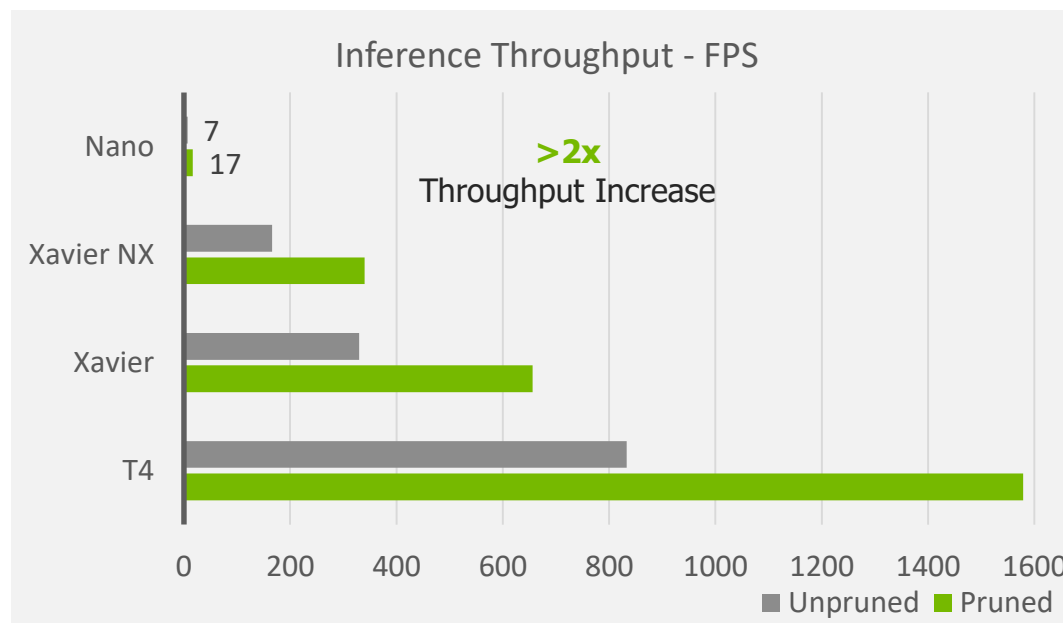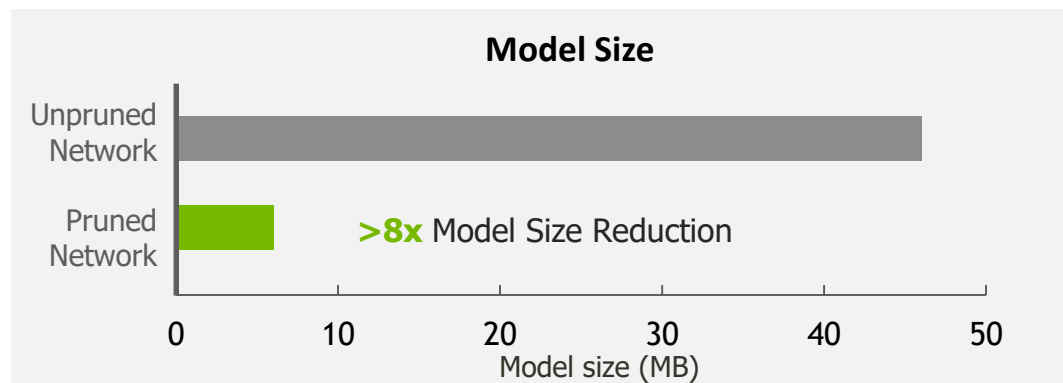2. Incrementally retrain model after pruning to recover accuracy

6 inputs, 6 neurons, 32 connections

6 inputs, 5 neurons, 24 connections

**Model Size**

Unpruned Network

Pruned Network

**>8x** Model Size Reduction

| 0 | 10 | 20 | 30 | 40 | 50 |

Model size (MB)

Inference Throughput - FPS

Nano   7
       17

**>2x**
Throughput Increase

Xavier NX

Xavier

T4

| 0 | 200 | 400 | 600 | 800 | 1000 | 1200 | 1400 | 1600 |

■ Unpruned  ■ Pruned

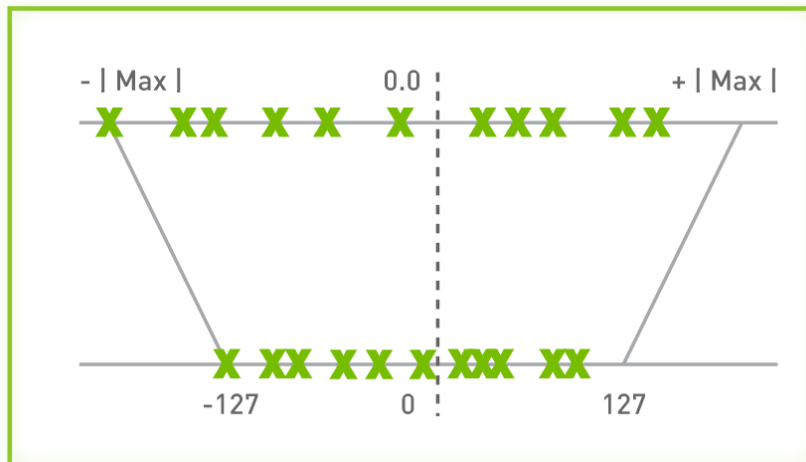*TrafficCamNet is DetectNet_v2 + ResNet18 based architecture (available on NGC)*
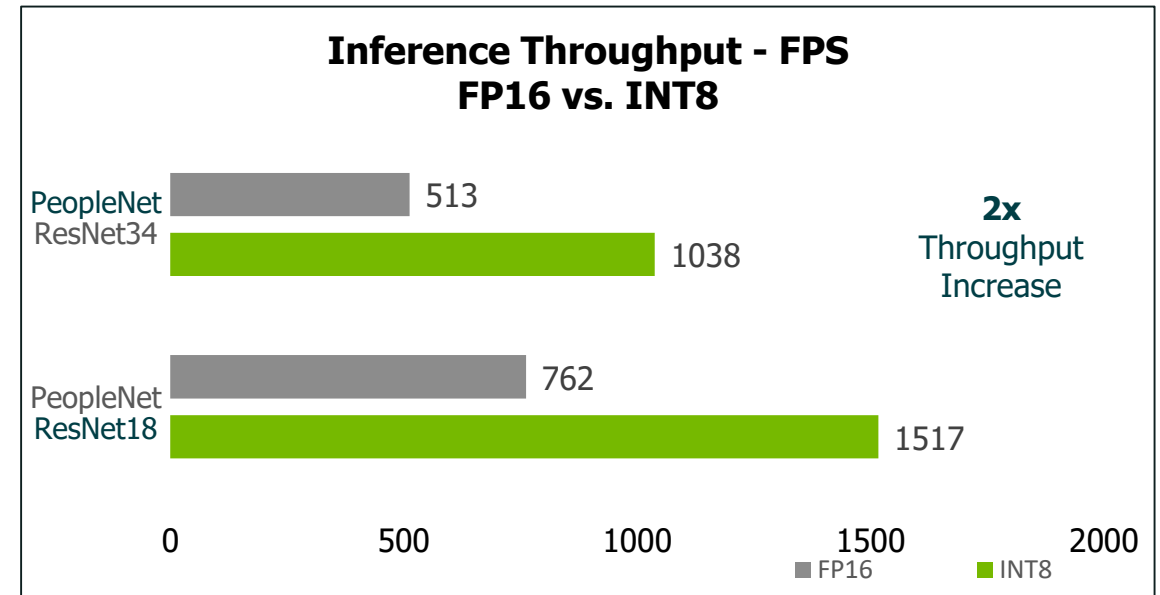
# Quantization

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

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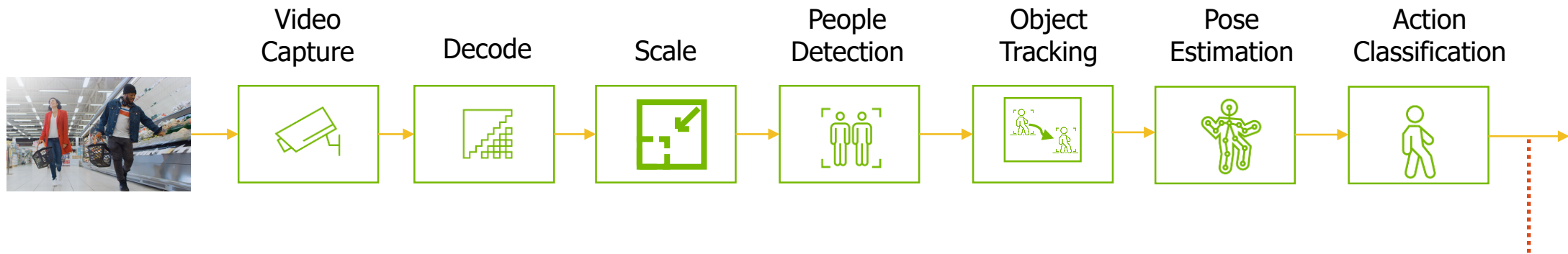


Transformation of floating-point weights to integer



**Inference Throughput - FPS FP16 vs. INT8**

PeopleNet ResNet34: FP16 513, INT8 1038 — **2x** Throughput Increase

PeopleNet ResNet18: FP16 762, INT8 1517

<1% loss in accuracy between FP16 and INT8

https://developer.nvidia.com/blog/improving-int8-accuracy-using-quantization-aware-training-and-the-tao-toolkit/
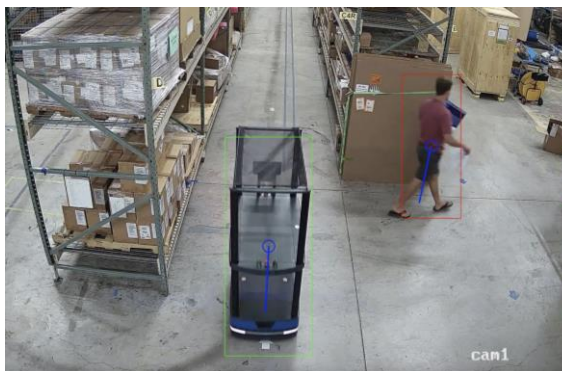
# Putting it all Together



Video Capture → Decode → Scale → People Detection → Object Tracking → Pose Estimation → Action Classification

# NVIDIA TAO – Accelerating AI Across Industries

**One CupAI**
Vision AI for Precision Agriculture

Remote tracking of animal health, growth and phenotype

Deployed an optimized solution that can run **inference on 1 TB** of data per day**, in just a few weeks**



**6River Systems**
Warehouse Logistics

Track objects in warehouse to optimize robot path planning and increase picking efficiency

Trained and deployed their model in application in just **weeks running with 30 parallel video streams instantaneously**
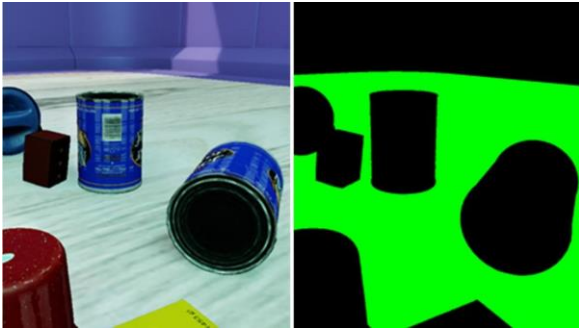


**Mavenir**
Quality Inspection

Detect defective bottles, labels before the final packaging step

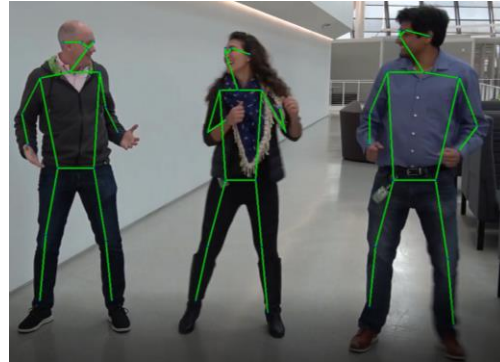Scaled across various use cases and **speed up development time by 3x**

# Summary

- The TAO Toolkit makes it easy for you to create custom, production-ready models for your speech and vision applications

  - Built on TensorFlow and PyTorch

- Diverse selection of pre-trained AI models

  - Removes the need for large training datasets

- Turnkey model optimization for inference

- Deploy easily with DeepStream for Vision AI applications at the edge

# Resources



AI-Powered Robots with Synthetic Data



2D Pose Estimation Model with NVIDIA TAO Toolkit
Part 1 | Part 2



Train and Deploy Action Recognition Model



Supercharge your AI workflow with TAO Toolkit Whitepaper

https://developer.nvidia.com/tao-toolkit