# REVOLUTIONIZING RETAIL WITH ARTIFICIAL INTELLIGENCE

Scott Brubaker, Paul Hendricks & Alex Sabatier

💽 NVIDIA.

### **INCEPTION PARTNERS & RETAIL ECOSYSTEM**





# **AI FOR RETAIL**



# **SHOPPING EXPERIENCE: STORE (IVA)**



### **TOP RETAIL IVA USE CASES**



# **SHOPPING EXPERIENCE: ONLINE**







# **RECOMMENDATION ENGINES ON GPU CLOUD**

### SONG RECOMMENDATIONS

### Spotify's Top Ten Most Popular Curated Playlists



### VIDEO RECOMMENDATIONS





#### TARGETED RECOMMENDATIONS











# **AI IN SUPPLY CHAIN**

#### WAREHOUSE OPTIMIZATION







#### FORECASTING AND REPLENISHMENT



# AI AT CORPORATE HQ



### GPU POWERED MACHINE LEARNING

### DATA SCIENCE IN RETAIL

- Supply Chain Replenishment
- Inventory Management
- Price Simulation & Management
- Prioritize Promotion Ad Targeting
- Marketing Optimization
- Personalized Recommendations
- Truck Routing
- **Online Delivery**



# THE STORE OF THE FUTURE

### Future-Proofed IVA Infrastructure

# DL-BASED IVA EDGE USE CASES

Loss Prevention Stock Out Reduction Store Analytics Security







# **NVIDIA VALUE**

### Comprehensive Platform for Retail IVA

NVIDIA DELIVERS	IVA Inference w/NVIDIA T4 GPU	Video Inference
peed Up	27*X CPU	27X
mages/second (1080P)	4400	25
Netropolis Platform optimized for IVA	DS Inference SDK TensorRT	20
GPU accelerated IVA Software Partners	70+	15
Deep Learning Education	Developer Blogs + IVA DLI	10 <u>10X</u>

GPU hardware accelerator engines for video decoding and encoding support faster than real-time video processing.

ResNet-50

CPU Server P4 T4

## ART OF THE POSSIBLE The State of AI in Retail

### 

Paul Hendricks Solutions Architect phendricks@nvidia.com



## INTRODUCTION

- Paul Hendricks is a Solutions Architect at NVIDIA, helping enterprise customers with their deep learning and AI initiatives
- Paul's background is primarily in retail, and has spent the past 5 years working with many Fortune 500 retail companies to implement data science and AI solutions.
- Prior to joining NVIDIA, Paul worked at Victoria's Secret as a Data Scientist building models to understand customer propensity to purchase and how to optimize assortment in stores.
- Currently, Paul's research at NVIDIA focuses on intelligent video analytics, machine leaning, recommendation systems, GANs, and reinforcement learning.





## INTRODUCTION

- Paul Hendricks is a Solutions Architect at NVIDIA, helping enterprise customers with their deep learning and AI initiatives
- Paul's background is primarily in retail, and has spent the past 5 years working with many Fortune 500 retail companies to implement data science and AI solutions.
- Prior to joining NVIDIA, Paul worked at Victoria's Secret as a Data Scientist building models to understand customer propensity to purchase and how to optimize assortment in stores.
- Currently, Paul's research at NVIDIA focuses on intelligent video analytics, machine leaning, recommendation systems, GANs, and reinforcement learning.

## **INTELLIGENT VIDEO ANALYTICS**

# **Image Classification**

### Problem Background

- Input Data: Images, Videos
- Goal: Given an input, identify the class that input belongs to



# **Object Detection**

### **Problem Background**

- Input Data: Images, Videos
- Goal: Given an input, identify objects and output bounding boxes around the objects and their classes



### Object Segmentation (Semantic Segmentation) Problem Background

- Input Data: Images, Videos
- Goal: Given an input, identify objects and output a mapping of pixels to their respective classes



Figure 1. The Mask R-CNN framework for instance segmentation.

## LOSS PREVENTION, STORE ANALYTICS, AND FRICTIONLESS CHECKOUT



https://www.standardcognition.com/

#### Single Stage Detectors

 These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass



#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference





#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference
- Examples: YOLOv3, SSD, RetinaNet, RetinaMask



#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference
- Examples: YOLOv3, SSD, RetinaNet, RetinaMask



#### **Two Stage Detectors**

 These algorithms generate a number of region proposals which are then passed to a CNN and classified



#### **R-CNN:** Regions with CNN features

#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference
- Examples: YOLOv3, SSD, RetinaNet, RetinaMask



#### **Two Stage Detectors**

- These algorithms generate a number of region proposals which are then passed to a CNN and classified
- Slower during inference since regions must be proposed and then evaluated (often redundant if overlaps)

#### **R-CNN:** Regions with CNN features



#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference
- Examples: YOLOv3, SSD, RetinaNet, RetinaMask



#### **Two Stage Detectors**

- These algorithms generate a number of region proposals which are then passed to a CNN and classified
- Slower during inference since regions must be proposed and then evaluated (often redundant if overlaps)
- Often are more accurate than single stage detectors, especially when trained on semantic segmentations

#### **R-CNN:** Regions with CNN features



#### Single Stage Detectors

- These algorithms regress the bounding boxes as well as classify the object within that bounding box in a single pass
- Computationally efficient and can be very fast during inference
- Examples: YOLOv3, SSD, RetinaNet, RetinaMask



#### **Two Stage Detectors**

- These algorithms generate a number of region proposals which are then passed to a CNN and classified
- Slower during inference since regions must be proposed and then evaluated (often redundant if overlaps)
- Often are more accurate than single stage detectors, especially when trained on semantic segmentations
- Examples: Faster RCNN, Mask RCNN

**R-CNN:** Regions with CNN features



# **GETTING STARTED**

#### **DLI Courses**

Introduction to Object Detection with TensorFlow – <u>https://courses.nvidia.com/courses/course-v1:DLI+L-AV-04+V1</u>

#### Papers

- YOLOV3 <u>https://pjreddie.com/publications/</u>
- Faster RCNN <u>https://arxiv.org/pdf/1506.01497</u>
- Mask RCNN <u>https://arxiv.org/abs/1703.06870</u>
- RetinaNet <u>https://arxiv.org/abs/1708.02002</u>
- RetinaMask <u>https://arxiv.org/abs/1901.03353</u>

#### Libraries

- DarkNet <u>https://github.com/pjreddie/darknet</u>
- TensorFlow's Object Detection API <u>https://github.com/tensorflow/models/tree/master/research/object\_detection</u>
- Facebook's Mask RCNN Benchmark <u>https://github.com/facebookresearch/maskrcnn-benchmark</u>

#### Datasets

- ImageNet <u>https://www.kaggle.com/c/imagenet-object-detection-challenge</u>
- Pascal VOC <u>http://host.robots.ox.ac.uk/pascal/VOC/</u>
- COCO <u>http://cocodataset.org/</u>
- Open Images <u>https://storage.googleapis.com/openimages/web/index.html</u>

## **MACHINE LEARNING**

# DATA SCIENCE IN RETAIL

Supply Chain Replenishment

Inventory Management

Price Management / Markdown Optimization

Prioritize Promotion And Ad Targeting

Marketing Optimization

Personalized Recommendations

**Truck Routing** 

Online Delivery



## **ML WORKFLOW STIFLES INNOVATION**



# DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



#### **DATA PREPARATION**

GPUs accelerated compute for in-memory data preparation Simplified implementation using familiar data science tools Python drop-in Pandas replacement built on CUDA C++. GPU-accelerated Spark (in development)

# DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



#### **MODEL TRAINING**

GPU-acceleration of today's most popular ML algorithms XGBoost, Random Forest, Linear Regression, PCA, K-means, k-NN, DBScan, tSVD ...

## DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



#### VISUALIZATION

Effortless exploration of datasets, billions of records in milliseconds Dynamic interaction with data = faster ML model development Data visualization ecosystem (Graphistry & OmniSci), integrated with RAPIDS

### RAPIDS – OPEN GPU DATA SCIENCE Software Stack

**Data Preparation** Model Training Visualization ► **PYTHON** DEEP LEARNING FRAMEWORKS RAPIDS DASK CUDF CUML **CUGRAPH** CUDNN CUDA **APACHE ARROW** 

# **GETTING STARTED**

#### **DLI Courses**

• Accelerating Data Science Workflows with RAPIDS – <u>https://courses.nvidia.com/courses/course-v1:DLI+L-DS-01+V1</u>

#### Resources

RAPIDS GitHub – <u>https://github.com/rapidsai</u>

- cuDF <u>https://github.com/rapidsai/cudf</u>
- cuML <u>https://github.com/rapidsai/cuml</u>
- cuGraph <u>https://github.com/rapidsai/cugraph</u>
- XGBoost <u>https://github.com/rapidsai/xgboost</u>
- Dask cuDF <u>https://github.com/rapidsai/dask-cudf</u>
- Dask cuML <u>https://github.com/rapidsai/dask-cuml</u>
- Dask XGBoost <u>https://github.com/rapidsai/dask-xgboost</u>
- Notebooks <u>https://github.com/rapidsai/notebooks</u>
- Notebooks Extended- <u>https://github.com/rapidsai/notebooks-extended</u>

# **NVIDIA HARDWARE**

## TESLA V100 TENSOR CORE GPU

World's Most Advanced Data Center GPU

5,120 CUDA cores 640 Tensor cores 7.8 FP64 TFLOPS | 15.7 FP32 TFLOPS | 125 Tensor TFLOPS 20MB SM RF | 16MB Cache 32 GB HBM2 @ 900GB/s | 300GB/s NVLink



### **TENSOR CORE BUILT FOR AI**

### Delivering 125 TFLOPS of DL Performance



Frameworks

**VOLTA-OPTIMIZED cuDNN** 



VOLTA TENSOR CORE 4x4 matrix processing array D[FP32] = A[FP16] \* B[FP16] + C[FP32] Optimized For Deep Learning



ALL MAJOR FRAMEWORKS

# **NVIDIA DGX**

### Al Supercomputer-in-a-Box



1000 TFLOPS | 8x Tesla V100 32GB | NVLink Hybrid Cube Mesh 2x Xeon | 8 TB RAID 0 | Quad IB 100Gbps, Dual 10GbE | 3U - 3200W

### NVIDIA DGX-2 THE WORLD'S MOST POWERFUL DEEP LEARNING SYSTEM FOR THE MOST COMPLEX DEEP LEARNING CHALLENGES

- First 2 PFLOPS System
- 16 V100 32GB GPUs Fully Interconnected
- NVSwitch: 2.4 TB/s bisection bandwidth
- 24X GPU-GPU Bandwidth
- 0.5 TB of Unified GPU Memory
- 10X Deep Learning Performance



### TESLA T4 WORLD'S MOST ADVANCED SCALE-OUT GPU

2,560 CUDA Cores 320 Turing Tensor Cores 65 FP16 TFLOPS | 130 INT8 TOPS | 260 INT4 TOPS 16GB | 320GB/s 70 W



## **NEW TURING TENSOR CORE**

MULTI-PRECISION FOR AI INFERENCE & ENTRY LEVEL TRAINING 65 TFLOPS FP16 | 130 TeraOPS INT8 | 260 TeraOPS INT4



THROUGHPUT

### WORLD'S MOST PERFORMANT INFERENCE PLATFORM

Up To 36X Faster Than CPUs | Accelerates All AI Workloads



For all three graphs:

Dual-Socket Xeon Gold 6140 @ 3.6GHz with single GPU as shown 18.11-py3 | TensorRT 5.0 | CPU FP32, P4 & T4: INT8 | Batch Size = 128

### WORLD'S FASTEST INFERENCE PERFORMANCE

### **NVIDIA GPUs Set New Performance Records**



### THE JETSON FAMILY





JETSON TX1 7 - 15W 1 TFOPS (FP16) 50mm x 87mm JETSON TX2 7 - 15W 1.3 TOPS (FP16) 50mm x 87mm JETSON AGX XAVIER 10 - 30W 10 TFLOPS (FP16) | 32 TOPS (INT8) 100mm x 87mm

UAVs • AI subsystems • AI Cameras

- Fully autonomous machines

Factory automation • Logistics • Delivery robots

Multiple devices • Unified software

## **NVIDIA SOFTWARE**

# **CHALLENGES WITH DEEP LEARNING**

Current DIY deep learning environments are complex and time consuming to build, test and maintain

Requires high level of expertise to manage driver, library, framework dependencies

Development of frameworks by the community is moving very fast



# **NVIDIA GPU CLOUD**

### Deep Learning Everywhere, For Everyone

Innovate in minutes, not weeks Removes all the DIY complexity of deep learning software integration

Always up to date Monthly updates by NVIDIA to ensure maximum performance

**Deep learning across platforms** Containers run locally on DGX Systems and TITAN PCs, or on cloud service provider GPU instances



NVIDIA GPU Cloud integrates GPU-optimized deep learning frameworks, runtimes, libraries, and OS into a ready-to-run container, available at no charge

## COMMON SOFTWARE STACK ACROSS DGX FAMILY





# **TENSORRT DEPLOYMENT WORKFLOW**

### Step 1: Optimize trained model



### Step 2: Deploy optimized plans with runtime





# **NVIDIA TENSORRT**

### From Every Framework, Optimized For Each Target Platform



Frameworks

Platforms

# TensorRT 5 & TensorRT Inference Server

Turing Support • Optimizations & APIs • Inference Server



Up to 40x faster perf. on Turing Tensor Cores

New INT8 workflows, Win & CentOS support

Maximize GPU utilization, run multiple models on a node

Free download to members of NVIDIA Developer Program soon at developer.nvidia.com/tensorrt

# TensorRT Inference Server

Containerized Microservice for Data Center Inference

Multiple models scalable across GPUs

Supports all popular AI frameworks

Seamless integration into DevOps deployments leveraging Docker and Kubernetes

Ready-to-run container, free from the NGC container registry



## TRANSFER LEARNING TOOLKIT

# End to End NVIDIA Deep Learning Workflow

Pre-Trained model access from NGC \* Training & adaptation \* Applications ready to integrate with DeepStream



Accelerate time to market and save on compute resources!



### **NVIDIA DEEPSTREAM**

### Zero Memory Copies





