INTEL AI DEVCON 2018



ARTIFICIAL INTELLIGENCE CASE STUDY

intel

Hongwei Yi 2018/08/08

LEGAL NOTICES & DISCLAIMERS

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at intel.com, or from the OEM or retailer. No computer system can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit <u>http://www.intel.com/performance</u>.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Statements in this document that refer to Intel's plans and expectations for the quarter, the year, and the future, are forward-looking statements that involve a number of risks and uncertainties. A detailed discussion of the factors that could affect Intel's results and plans is included in Intel's SEC filings, including the annual report on Form 10-K.

The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

Intel, the Intel logo, Pentium, Celeron, Atom, Core, Xeon, Movidius, Saffron and others are trademarks of Intel Corporation in the U.S. and/or other countries. *Other names and brands may be claimed as the property of others.









«Alpha available [†]Beta available

[‡] Future

*Other names and brands may be claimed as the property of others. All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.



CASE 1: TRANSACTION FRAUD DETECTION WITH DEEP LEARNING & MACHINE LEARNING



TRANSACTION FRAUD DETECTION WITH DEEP LEARNING & MACHINE LEARNING



Problem: no single algorithm delivers high-enough accuracy **Solution:** ensemble models of classical machine learning with deep learning achieved industry recorded high detection accuracy



The COMPLETE End-to-End detection runs on Intel[®] Xeon[®] Scalable Processor



THE FRAUD CREDIT CARD DETECTION USAGE OF POC





THE <u>within-between-within(wbw)</u> sandwich-structured sequence learning architecture

Fraud Credit Card Transaction Detection Flow Diagram



The Within-Between-Within(WBW) Architecture



- WBW model achieved the record high >25% precision on >60% recall rate, with latency of <200 ms
- GBDT & RF don't run well on GPU due to irregular computing in the algorithms



"WBW" (WITHIN-BETWEEN-WITHIN) SANDWICH-STRUCTURED SEQUENCE LEARNING ARCHITECTURE



CASE 2 - HIGH-CONTENT IMAGING IN DRUG Discovery



HIGH-CONTENT IMAGING (HCI) IN DRUG DISCOVERY

High-content imaging (HCI) or screening (HCS) has seen increased application in systems biology & drug discovery. Images acquired through microscopy-based assays provide visual information to investigate cellular phenotypes induced by genetic or chemical treatments.



Conventional HCS analysis pipeline

HCS analysis using Multi-Scale Convolutional Network¹





Multi-scale convolutional neural network



Phenotype probability





MCNN HAS A LARGE MEMORY FOOTPRINT

High content images are 26x larger than images from ImageNet dataset



MCNN size grows linearly with batch size used in each step of training







INTEL XEON® ENABLES LARGE BATCH TRAINING WITH LARGE MEMORY

Training Time Improves with large batches

Peak memory utilization in MCNN training can scale well beyond 16GB



¹ Workload: Image set <u>BBBC021</u>: Human MCF7 Cells – compound profiling experiment. Configuration details in backup

² Godinez et al, A multi-scale convolutional neural network for phenotyping high-content cellular images. Bioinformatics, 2017

tel[®]AI 13

LARGE BATCH TRAINING M-CNN REACHES SOTA ON 8 XEON® SERVERS

High Content Screening/M-CNN Training on 8 Node Intel[®] 2S Xeon[®] 6148 processor cluster TensorFlow 1.7, Horovod, OpenMPI, BS=32/Node, GBS=256, OPA Fabric



Workload: Image set <u>BBBC021</u>: Human MCF7 Cells – compound profiling experiment.

TensorFlow: 1.7.0, Python: 2.7.5, Horovod: 0.12.1: OMP_NUM_THREADS=10

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit http://www.intel.com/performance. *Other names and brands may be claimed as the property of others



HIGH PERFORMANCE AT SCALE WITH INTEL® XEON® SCALABLE PROCESSOR

Time To Train



Workload: Image set <u>BBBC021</u>: Human MCF7 Cells – compound profiling experiment.

TensorFlow: 1.7.0, Python: 2.7.5, Horovod: 0.12.1, OpenMPI 3.0.0

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit http://www.intel.com/performance. *Other names and brands may be claimed as the property of oth

15

'inte

CASE 3: DEEP LEARNING FOR DRUG DISCOVERY



DEEP LEARNING FOR DRUG DISCOVERY

Problem: Need high efficiency platform for clinical genome analysis on DeepChem Topologies: GCN

Solution: Optimized data feeding and computing architecture on Xeon Processor Scalable Family faster than GPGPU





Multi-task Graph Convolutional Network for PCBA and ChEMBL





Drug discovery (predicting compound-protein interactions)

CASE 4: DIAGNOSTIC IMAGING - BODY PART Classification in CT Scans



DIAGNOSTIC IMAGING - BODY PART CLASSIFICATION IN CT SCANS

Problem: Need fast inferencing to classify images coming off of a CT scanner

Solution:

High performance inferencing at ~600 images/sec. Exceeded customer goals by 6x and provided 10x performance boost over unoptimized version

CT scan images



Inference Throughput v/s Core Count







CASE 5: MEDICAL IMAGE ANALYSIS Application



AI INFERENCE OPTIMIZATION FOR USER EXPERIENCE & TCO

- <u>**Customer</u>**: A leading Medical Image solution company in China wants to port an inference application from NVidia GPGPU to Intel[®] Xeon[®] platform for production deployment of the solution on cloud. The motivation is to reduce the deployment cost and simplify the cloud infrastructure requirements.</u>
- <u>Challenge</u>: The Medical Image Analysis application requires fast responsiveness for end user experience. The original latency of inference on Xeon[®] is too high.
- <u>Support Requirement</u>: Inference latency on Xeon[®] can be equal to NVidia[®] GeForce[®] GTX 1080
- <u>**Result</u>**: Inference latency on Xeon[®] Gold 2S system is equal or lower than Nvidia[®] GPU</u>







WHAT CUSTOMER SHOWED AT INTEL AI WORKSHOP

Transfer to IA Platform

AMD GPU - Rendering

Nvidia GPU - AI computation

CPU - Business processing

Xeon

Low cost、high density、large clustering、Seamless transfer

Performance comparison Framework : Caffe & Tensorflow Model : RFCN Single image in Dicom format Data : (Chest radiograph) 68ms GTX 1080 Xeon 6148

Good performance and Seamless Migration



CASE 6: AI AS A SERVICE - AI CLOUD



UCLOUD UAI-INFERENCE

- UCloud Top5 public cloud service providers in PRC, served more than 50,000 enterprises
- UCloud launched AI online service UAI-inference using Intel[®] Xeon[®] SKX based servers to build up best TCO.
- UCloud has introduced a better performing AI framework: Intel Caffe with 43x higher performance

UCLOUD UAI-Service: The first commercial deployment of Intel optimized Caffe solution as AI differentiated service on SKX in public cloud provider

UCloud Object Detection CTPN Inference Optimization(Lower is better)







FIND OUT MORE







