

GADI SINGER VICE PRESIDENT, ARCHITECTURE GENERAL MANAGER, AI PRODUCTS GROUP - INTEL



THE AI REVOLUTION IS REALLY A COMPUTING EVOLUTION





ARTIFICIAL INTELLIGENCE

TYPES OF ANALYTICS/ML (PARTIAL LIST)



REGRESSION

CLUSTERING

FEATURE LEARNING **ANOMALY DETECTION**





ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

DEEP LEARNING





AI IS TRANSFORMING ['14/15]

SAMPLE-SCALE BENCHMARKS

PROGRAMMING MODEL C++, Proprietary APIs

EARLY ADOPTION

EARL)



COMPUTE RATIO Inference ~= Training





DL COMING OF AGE ['19/20]

AT-SCALE DEPLOYMENT OF SCIENCE AND REAL-LIFE SOLUTIONS





ML FRAMEWORKS 'Democratization of Data Science'









COMPUTE RATIO Inference >> Training

EARLY MAJORITY

LATE MAJORITY

LAGGARDS

HARDWARE ARCHITECTURE **CPU, DL Acceleration/NNPs, GPUs**



AI IN ACTION - LEADERSHIP AT SCALE

INTEGRATION WITH EXISTING WORKLOADS

BLENDED REAL-LIFE AI WORKLOADS

END DEVICES TO DATACENTER

PERFORMANCE VS PERF/W VS LATENCY

BUILT FOR AI. DESIGNED TO SCALE.

TCO-TOTAL COST OF OWNERSHIP

FULL PLATFORM OPTIMIZATION



ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION

mirror object to mirror irror_od_mirror_object Peration = "MIRROR_X": irror_mod.use_x = True irror_mod.use_y = False operation = "MIRROR_Y" irror_mod.use_y = True irror_mod.use_y = True irror_mod.use_x = False operation = "MIRROR_Z" irror_mod.use_x = True

SOFTWARE TOOLS

pint("please ser

- OPERATOR CLASSE

X mirror to the sel X mirror to the sel ject.mirror_mirror_x ror X

HARDWARE Platforms

ECOSYSTEM COMMUNITIES





AMR KHOSROWSHAHI Vice president, ai products group Chief technology officer - intel



ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION

mirror object to mirror irror_mod.mirror_object peration = "MIRROR_X": irror_mod.use_x = True irror_mod.use_y = False operation = "MIRROR_Y" irror_mod.use_x = False irror_mod.use_y = True irror_mod.use_x = False irror_mod.use_x = False irror_mod.use_x = False irror_mod.use_y = False irror_mod.use_y = False

SOFTWARE TOOLS

pint("please see

- OPERATOR CLASSE

types.Operator): X mirror to the sele ject.mirror_mirror_X ror X

ently object in

A COMPANY OF A COM





OPEN SOURCE INTEL[®] MATH KERNEL LIBRARY For Deep Neural Networks (Intel[®] MKL-DNN) HELPS REALIZE THE INCREDIBLE BENEFITS OF DIRECT OPTIMIZATION MATRIX MULTIPLICATION **BATCH NORM** POOLING NORMALIZATION

ACTIVATION

CONVOLUTION





TensorFlow

OPTIMIZING TENSORFLOW TO SUPERCHARGE A WORKLOADS

Other names and brands may be claimed as the property of others





ENTERPRISE DEEP LEARNING PLATFORM

VINAY KUMAR SANKARAPU Founder & Ceo, Arya.ai



HIGH PERFORMANCE DELIVERED ON INTEL® XEON® SCALABLE PROCESSOR

USE CASE: AUTOMATED CLAIMS PROCESSING for health insurance



CUSTOMER IMPACT

Inferencing Performance on Intel® Xeon® Scalable using Intel distribution for Python* & Intel-optimized TensorFlow* enabled the arya platform to reduce claim processing time from 48hrs to 0.2 seconds



PERFORMANCE WITH THE INTEL® XEON® SCALABLE PROCESSOR

7

6

5

3

2

0





NGRAPH

AND HARDWARE

) NervanaSystems/ngraph



OPEN SOURCE COMPILER ENABLING FLEXIBILITY TO RUN MODELS ACROSS A VARIETY OF FRAMEWORKS

BIGDL: DISTRIBUTED DEEP LEARNING

AKANKSHA BALANI Country Lead - Intel Software Developer Products





HIGH PERFORMANCE DEEP LEARNING FOR APACHE SPARK* ON CPU INFRASTRUCTURE



DESIGNED AND OPTIMIZED FOR INTEL® XEON® PROCESSOR

No need to deploy costly accelerators, duplicate data, or suffer through scaling headaches!





Lower TCO, improved ease of use



Efficient Scale-Out

Powered by Intel® MKL-DNN



CASE STUDY: IMAGE RECOGNITION JD.京东 JD.com

JD.京东 JD.Com JD.Com, 2nd largest online retailer in China, ~ 250 M users. S CHALLENGE: Building of

Building deep learning applications such as image similarity search without moving data.



https://software.intel.com/en-us/articles/building-large-scale-image-feature-extraction-with-bigdl-at-jdcom

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.

*Other names and brands may be claimed as the property of others



Intel[®] Xeon[®] CPU

SOLUTION:

Switched from GPU to CPU cluster. Using Apache Spark* with BigDL, running on Intel® Xeon® processors



VIDEO IS THE ULTIMATE IOT SENSOR



















ANNOUNCING: OPENVINO[™] SOFTWARE TOOLKIT VISUAL INFERENCING AND NEURAL NETWORK OPTIMIZATION





DEPLOY COMPUTER VISION AND DEEP LEARNING CAPABILITIES TO THE EDGE





ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION

mirror object to mirror irror_od.use_x = false operation = "MIRROR_X": irror_mod.use_y = false operation = "MIRROR_Y" irror_mod.use_y = false irror_mod.use_y = false operation = "MIRROR_Y" irror_mod.use_x = false operation = "MIRROR_Y"

lection at the end -nd ob.select+1 er_ob.selecte1 ntext.scene.dijd.ts.action "selected" -- str(modifier pror ob.select + 0 bpy.context.selected_ob nte.objects[one-name].selecte1

pint("please select exac

-- OPERATOR GLASSES

types.Operator):
 X mirror to the sele
 X mirror_mirror_*
 ject.mirror_mirror_*

HARDWARE Platforms







BUILT FOR AI. DESIGNED TO SCALE.









FOUNDATIONAL FOR ARTIFICIAL INTELLIGENCE



ARTIFICIAL INTELLIGENCE AND INTEL® XEON® PROCESSORS





Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective

assumials Chantalla, Uthen U

Kim Hazeboood, Sarah Bird, David Bro Mohamed Fawyy, Bill Jia, Vanjajing Parter Normalians, May

Abstract-Machine learning sits at the core of ma products and services at Facebook. This paper a hardware and software infrastructure that supplies learning at global scale. Facebook's machine tearning workloads are extremely diverse; services require many different types of founded by membre learning services process chillings at models in practice. This diversity has implications at all inversion the global scale of Facebook's data entry broad to tamps the system stack. In addition, a sizable fraction of all data stored are used to efficiently leed data to the models matching deat I acchook flows through machine learning pipetines, presenting coupling of data teed and warming darkompac colocation. distributed training flows. Computational requirements are also intense, leveraging both GPU and CPU platforms for training and scale provides unager opportunities. Doned had creleviewe abundant CPU capacity for real-time inference. Addressing these a significant number of CPU- southing for distributed number and other emerging challenges continues to require diverse efforts algorithms during off-peak periods. With Exclosek's compar-

1. INTRODUCTION

Facebook's mission is to "Give people the power to build Looking toward, Facebook expects and proofs in ma community and bring the world closer together." In support chang learning across eviding and new services (4). This of that mission, Eacebook connects more than two billion growth will lead to growing scalability chillenges for teams people as of December 2017. Meanwhile, the past several deploying the infrastructure for these services. While sentyears have seen a resolution in the application of machine cant opportunities exist to optimize infrastructure on existing learning to real problems at this scale, building upon the platforms, we continue to actuely evaluate and policype virtuous cycle of machine learning algorithmic innovations, new hardware solenous while temaning constrait of ganeenormous amounts of training data for models, and advances changing algorithmic innovations in high-performance computer architectures [1]. At Facebook, machine learning provides key capabilities in driving nearly major insights about machine learning at Vacebook. all aspects of user experience including services like ranking posts for News Feed, speech and text translations, and photoand real-time video classification [2], [3]. Facebook leverages a wide variety of machine learning al-

gorithms in these services including support vector machines. gradient boosted decision trees, and many styles of neural networks. This paper describes several important aspects of datacenter infrastructure that supports machine learning at Facebook. The infrastructure includes internal "ML-as-a-Service" flows, open-source machine learning frameworks, and distributed training algorithms. From a hardware point of view. Facebook leverages a large fleet of CPU and GPU platforms for training models in order to support the necessary Other names and brands may be claimed as the property of other for all major applied-machine-learning

and networking optimizations. At the same time, Facebook's that span machine learning algorithms, software, and hardware their spread over ten datacenter because, sede abso provides disaster receivery capitality. Display receivery planame in essential as timely delayery of new machine learning models. is important to Facebook's operations.

The key contributions of this paper include the following

- Machine learning is applied pervavisely across nearly all services, and computer vision represents only a small
- fraction of the resource requirements · Facebook relies upon an incredibily diverse set of machine learning approaches including, but we hunch to
- Tremendous amounts of data are tunneled through our
- machine learning pipelines, and this creates engineering and efficiency challenges far beyond the compare nodes Facebook currently relies heavily on CPUs far inference.
- and both CPUs and GPUs for training, but constantly protocypes and evaluates new hundware solutions from a
- . The worldwide scale of people on Facebook and corresponding diurnal activity patterns result in a bage number
- of machines that can be harnessed for machine learning tasks such as distributed training at scale

Trainin Data Features FBLearner FBLearn Feature Flow Store CPU+GP CPU Fig. 1. Example of Facebook's Machine Learning Flow and Infrastructure.

Rela Services News Feed 1002 10X Facer 10X Lumos 10X Search 1XLanguage Translation 1X Sigma Speech Recognition 1X

g Model	Inference	Predictions	HOME AR	G, Sensch Suit Protosi Comment	CONNUNTY ¢ Tians
er	FBLearner Predictor		August 15 Have an HPCA- Is Vienna in Pea September 3.	HIPCA 2018	mun? Wart to go posai ky
U	CPU	Download		Recipional And	(POSALS

tive Capacity	Compute	Memory
x	Dual-Socket CPU	High
	Single-Socket CPU	Low
	Single-Socket CPU	Low
	Dual-Socket CPU	High

TABLE III RESOURCE REQUIREMENTS OF ONLINE INFERENCE WORKLOADS.



wipro holmes

FROM CHIP TO APPLIED AI

TAPATI BANDOPADHYAY General Manager and Practice Head Wipro Holmes



AI-OPTIMIZED CHIP TO BUSINESS SERVICES



SOLUTION DEPLOYED

wipro holmes



Performance Benchmarking in Image Detection & Recognition

Leveraging CNNs on Intel[®] Xeon[®] Platinum

OUTCOME



320+ Accounts



Analyst Reports Leadership **20 FY18**



Automation Arbitrage **19,000** people equivalent productivity



3,200+ Bots Deployed FY18



A SCIENTIFIC COLLABORATION BETWEEN INTEL AND NOVARTIS



HIGH PERFORMANCE AT SCALE

SCALING OF TIME TO TRAIN

Intel[®] Omni-Path Architecture, Horovod and TensorFlow[®]



Other names and brands may be claimed as the property of others § Configuration: CPU: Xeon 6148 @ 2.4GHz, Hyper-threading: Enabled. NIC: Intel® Omni-Path Host Fabric Interface, TensorFlow: v1.7.0, Horovod: 0.12.1, OpenMPI: 3.0.0. OS: CentOS 7.3, OpenMPU 23.0.0, Python 2.7.5

Time to Train to converge to 99% accuracy in model 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 c 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

Speedup compared to baseline 1.0 measured in time to train in 1 nodes

TOTAL MEMORY USED 192GB DDR4 PER INTEL® 2S XEON® 6148 PROCESSOR





FLEXIBLE REAL-TIME INFERENCING FPGA PRODUCTS



G

P104300033

// AI For Earth







DEEP NEURAL NETWORKS - PROJECT BRAINWAVE

DR. VIVEK SESHADRI Researcher – Microsoft Labs



PROJECT BRAINWAVE



A SCALABLE FPGA-POWERED SERVING PLATFORM FOR DEEP NEURAL NETWORKS

INDUSTRY-LEADING LATENCY (<2MS) AT ULTRA LOW COST

DEPLOYABLE TO INTELLIGENT CLOUD OR INTELLIGENT EDGE THROUGH AZURE IOT



http://aka.ms/aml-real-time-ai brainwave-edge@microsoft.com









Movidius MA2485 Myriad)

Deploy DNN and Computer Vision at the Edge

Native FP16 and Fixed Point 8 bit support

4 TOPS with 1 TOPS of DNN Compute at 1W



Novidius" A CONTRACT NOVIDIUS



(intel) Nervana[®] Nervana[®]

intel

Nervana





GOALS - HIGH UTILIZATION & MODEL PARALLELISM



SPRING CREST PURPOSE BUILT DESIGN OPTIMIZED ACROSS MEMORY BANDWIDTH, UTILIZATION, AND POWER

INTEL[®] NERVANA[™] **NNP L-1000** in 2019

3-4x training performance of first generation Lake Crest product





ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION

mirror object to mirror irror_mod.mirror_object Peration = "MIRROR_X": irror_mod.use_X = True irror_mod.use_Y = False operation = "MIRROR_Y irror_mod.use_X = False irror_mod.use_X = False operation == "MIRROR_Z irror_mod.use_X = False irror_mod.use_X = False irror_mod.use_Y = False irror_mod.use_Y = False

lection at the endob.select=1 or_ob.select=1 ntext.scene.ddjetts.au "Selected" = str(modif or of select = 8

It a objects [one-name] .sn

NUMBER OF STREET

- OPERATOR CLASSES

types.Operator): X mirror to the sele ject.mirror_mirror_X rject.X

TOTAL OF THE



ECOSYSTEM COMMUNITIES





PRAKASH MALLYA Managing Director, Sales & Marketing Group Intel India



THE PROMISE OF DATA DRIVEN INDIA



- 1B+ citizen data digitized; payment digitization
- 4G infra momentum, content explosion



- 3rd largest footprint of Al startups*
- 2nd largest developer ecosystem in the world



- •

COMMUNITY BUILDING CRITICAL TO UNLEASH THE TRUE POTENTIAL OF AI



OPPORTUNITIE

Our demographics

• High friction verticals: Healthcare, Agriculture and Education







SUPERCHARGE YOUR ML MODEL

ATANU ROY AI SPECIALIST SOLUTIONS ARCHITECT - AMAZON

Other names and brands may be claimed as the property of others



MACHINE LEARNING AT AMAZON: **A LONG HERITAGE**



amazon go

Inventing entirely new customer experiences



AMAZON SAGEMAKER







One-click deployment



Fully managed hosting with autoscaling





PHILIPS BONE-AGE DL INFERENCE RAVI RAMASWAMY PHILIPS HEALTHCARE - SENIOR DIRECTOR



PHILIPS BONE-AGE DL INFERENCE MODEL ON INTEL® XEON® PROCESSOR PLATFORM

BIOMEDICAL IMAGE BASED ANALYTICS FOR DISEASE SCREENING/DIAGNOSIS

1



HEALTHCARE ANALYTICS WLs **OPTIMIZATION ON INTEL® XEON® PROCESSOR**







40X IMPROVEMENT IN INFERENCE PERFORMANCE





AI FOR SAFE MOBILITY

DR. C. V. JAWAHAR Amazon Chair Professor, CVIT, IIIT-Hyderabad





INDIA DRIVING DATA SET AN INTEL & IIIT-HYDERABAD INITIATIVE

DATA MAKES AI POSSIBLE



Enabling AI Research in Unstructured Driving **Conditions & Innovations** around Safe Mobility



WORLD'S FIRST PUBLIC DATA SET OF INDIAN DRIVING CONDITIONS



ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION

mirror object to mirro irror_mod.mirror_object Peration = "MIRROR_X": irror_mod.use_x = True irror_mod.use_y = False operation = "MIRROR_Y" irror_mod.use_y = True irror_mod.use_y = True irror_mod.use_y = True irror_mod.use_x = False irror_mod.use_y = False irror_mod.use_y = False

SOFTWARE TOOLS

Pint("please set

- OPERATOR CLASSE

types.Operator): X mirror to the sel ject.mirror_mirror_x ror X

HARDWARE Platforms

ECOSYSTEM COMMUNITIES



HELPING MAKE YOUR **AIVISION A REALITY**

Stratix°10

SG280LN3F43E3VGS1



-

Movidius

Nervana

Intel* Xeon* Processor E5 V4



THE ADDRESS



RISK FACTORS

Today's presentation contains forward-looking statements. All statements made that are not historical facts are subject to a number of risks and uncertainties, and actual results may differ materially. Please refer to our most recent earnings release, Form 10-Q and 10-K filing available on our website for more information on the risk factors that could cause actual results to differ.

