





COMPUTE AND THE BRAIN

OPERANT CONDITIONING 1938 SPIKING NEURON 1952 FIRST NEUROSCIENCE DEPARTMENT DEEP LEARNING PREVALENCE mid 2000s

THE TURING MACHINE

TRANSISTOR

FIRST COMPUTER SCIENCE DEPARTMENT

INTEL FOUNDED

FIRST 1
BILLION
TRANSISTOR
PROCESSOR
mid 2000s







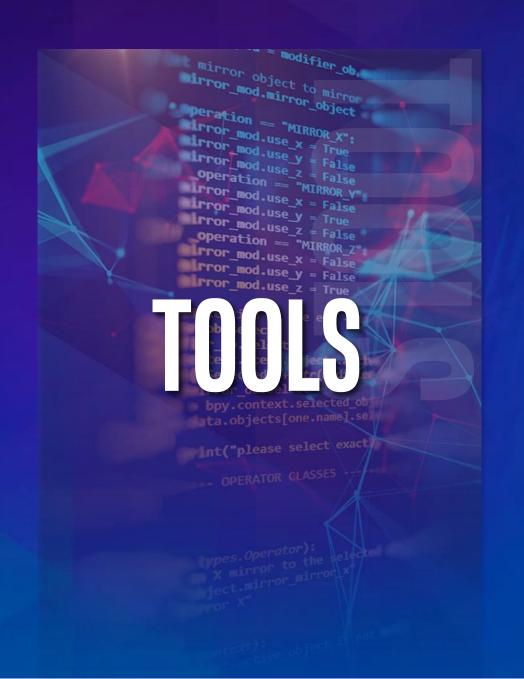




















JASON KNIGHT HEAD OF SOFTWARE PRODUCTS INTEL AI



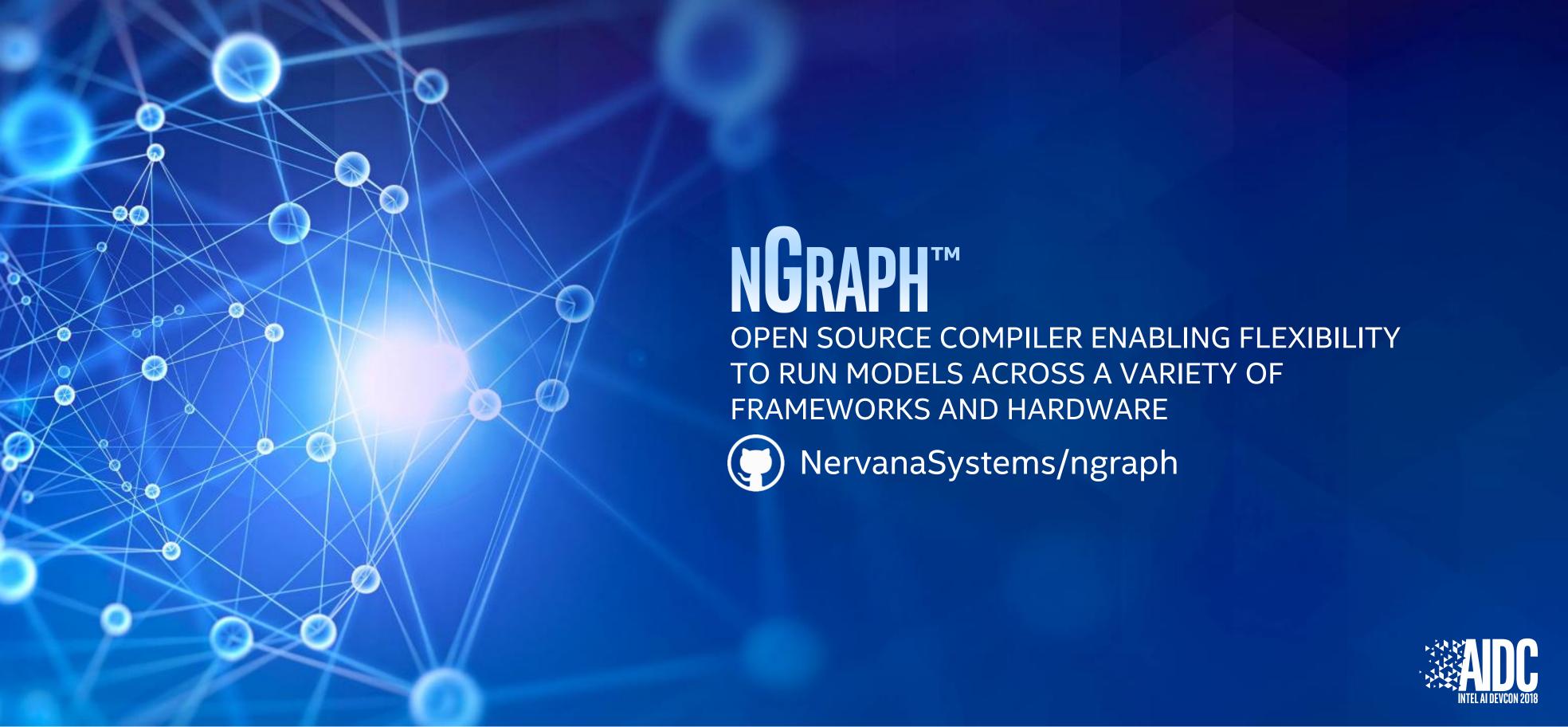


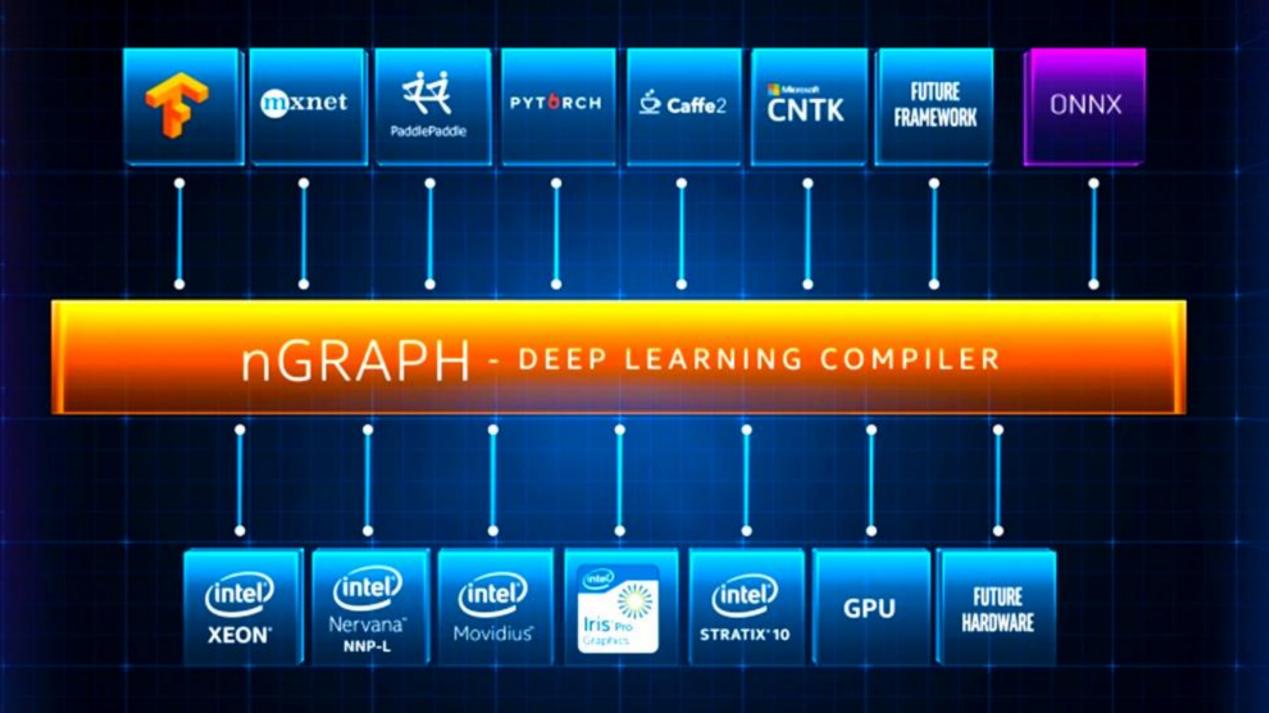


OPTIMIZING TENSORFLOW TO SUPERCHARGE AI WORKLOADS

RAJAT MONGA
TENSORFLOW ENGINEERING DIRECTOR
GOOGLE

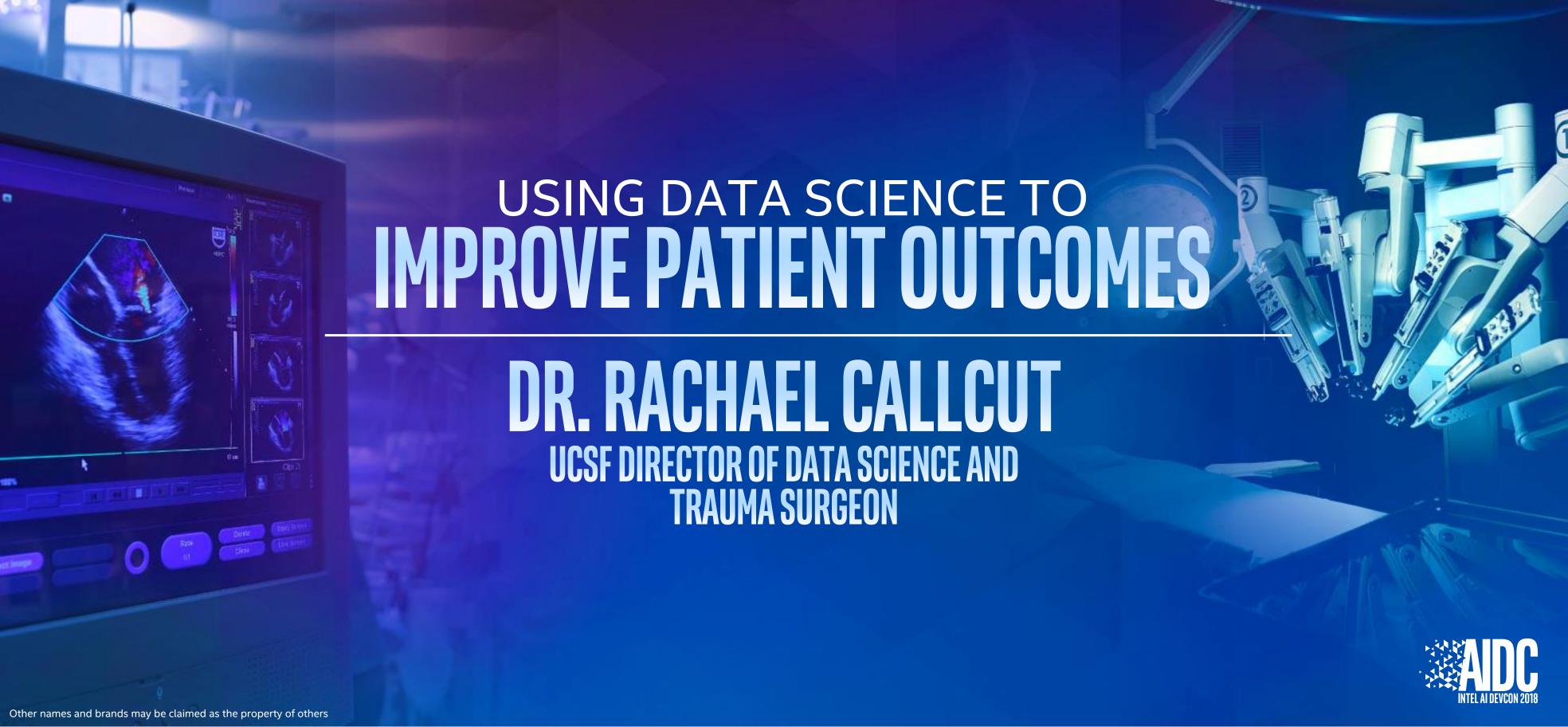














VIDEO IS THE ULTIMATE IOT SENSOR











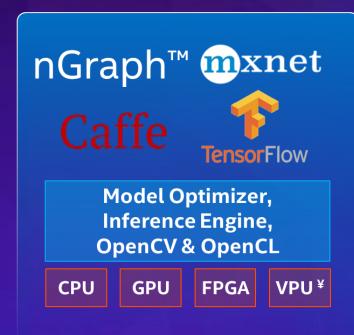








ANNOUNCING: OPENVINO SOFTWARE TOOLKIT VISUAL INFERENCING AND NEURAL NETWORK OPTIMIZATION



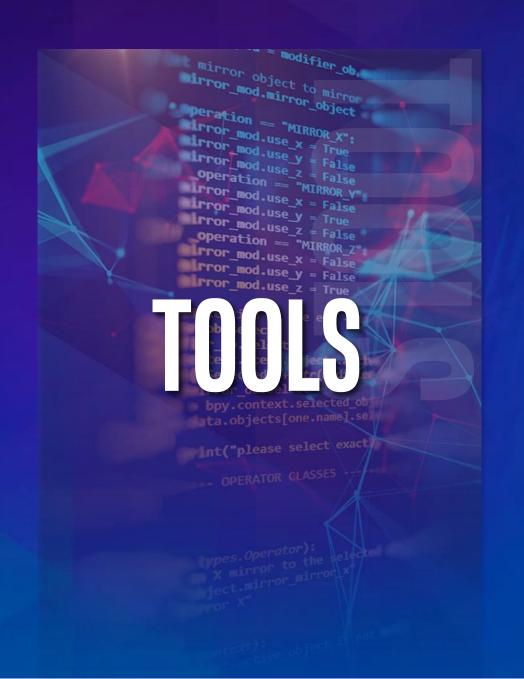
DEPLOY COMPUTER
VISION AND DEEP
LEARNING CAPABILITIES
TO THE EDGE



























FOUNDATIONAL FOR ARTIFICIAL INTELLIGENCE



Intel® Xeon® Processor

Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective

Kim Hazelwensd, Sarah Bird, David Brooks, Sammith Chunda, Uthu Dud, Davin Debulgder. Mohamed Fawzy, Bill Ita. Vangqing Ita. Aditya Kahu, James Law Kevin Lee, Issue Lu. Pacter Neutrillanis, Micha Smelyanskiy, Lung Xinng, Xinsdong Wing

Abstract - Machine learning sits at the core of many essential products and services at Facebook. This paper describes the hardware and software infrastructure that supports machine learning at global scale. Facebook's machine bearing workloads are extremely diverse; services require many different types of required by machine learning services present in data to be dead of the de models in practice. This diversity has implications at all layers in the global scale of Uncheck's datacenters beyond schanges the system stack. In addition, a sizable fraction of all data stored one used to efficiently feed data to the models including designificant challenges in delivering data to high-performance distributed training flows. Computational requirements are also intense, leveraging both GPU and CPU platforms for training and abundant CPU capacity for real-time inference. Addressing these and other emerging challenges continues to require diserve efforts that span machine learning algorithms, software, and hardware

I. INTRODUCTION

Facebook's mission is to 'Give people the power to build community and bring the world closer together." In support of that mission, Facebook connects more than two billion people as of December 2017. Meanwhile, the past several deploying the intrustructure for these services. While signifiyears have seen a revolution in the application of machine cant opportunities exist to optimize infrastructure on existing learning to real problems at this scale, building upon the virtuous cycle of machine learning algorithmic innovations, new bardware solutions while remaining cognizat of gameenormous 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 all aspects of user experience including services like tanking posts for News Feed, speech and text translations, and photo and real-time video classification [2], [3].

Facebook leverages a wide variety of machine learning algorithms 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 training frequencies at the required service latency. For machine learning inference. Facebook primarily relies on CPUs for all major services with neural network ranking services

Vocabawk terroch a large traction of all merch data formigh. machine learning papelines, and this teaction is mescaring over compling of data beed and training data/company co-incatang and networking optimizations. At the same time, Vacchook's scale provides unique opportunites. Diunal had cycles have a seguificant number of CPUs wealthic for distributed teaming algorithms during off-peak periods. With Facebook's company theer spread over ten datacenter because, scale also provides disaster recovery capability. Disaster recovery planning to essential as timely delivery of new machine bearing models. is important to Facebook's operations.

Looking torward. Facebook expects rapid growth in machang learning across existing and new services (4). This growth will lead to growing scalability challenges for teams platforms, we continue to actively evaluate and prototype

The key contributions of this paper include the following major insights about machine learning at Facebook.

- Machine learning is applied personnely across seath, ill. services, and computer vision represents only a small fraction of the resource requirements.
- Facebook relies upon an incredibly diverse set of machine learning approaches including, but not humed to.
- Tremendous amounts of data are tunneled through our machine learning pipelines, and this creates engineering and efficiency challenges for beyond the compute modes.
- Facebook currently relies beavily on CPUs for interence. and both CPUs and GPUs for training, but constantly broughly, and evaluates new pardware reparent from a
- The worldwide scale of people on Facebook and corresponding diarnal activity patterns result in a buge number of machines that can be harnessed for machine learning tasks such as distributed training at scale.

SU LET'S GFT GROUNDED IN REALITY



fatebook.

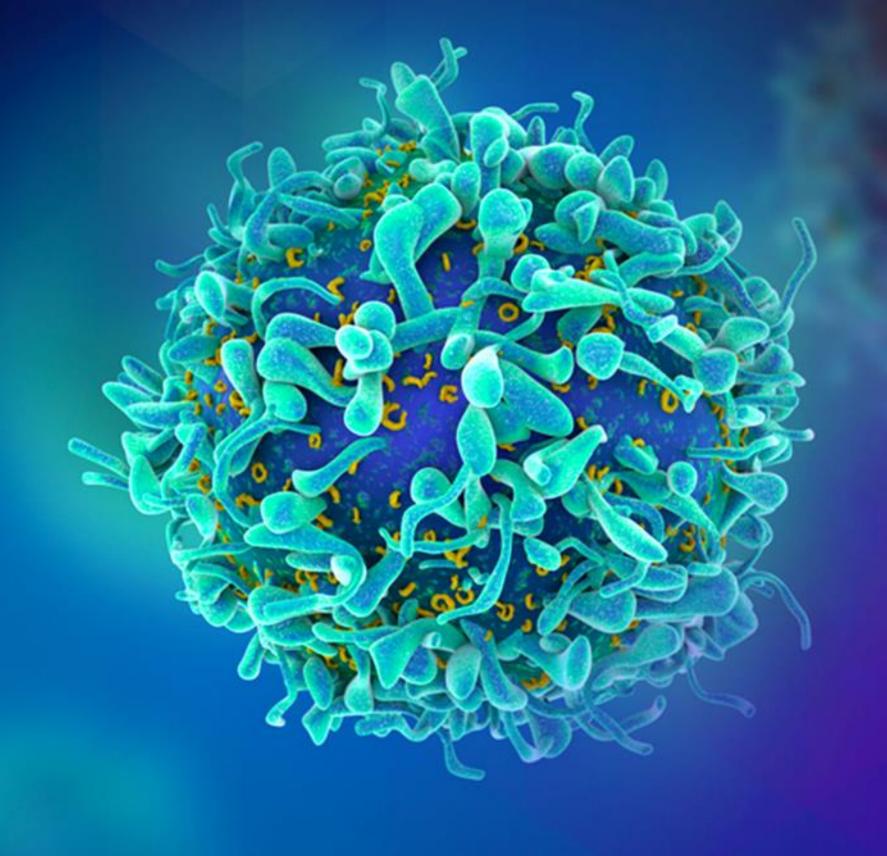
HEAD OF AI INFRASTRUCTURE FOUNDATION, FACEBOOK

Services	Relative Capacity	Compute	Memory
News Feed	100X	Dual-Socket CPU	High
Facer	10X	Single-Socket CPU	Low
Lumos	10X	Single-Socket CPU	Low
Search	10X	Dual-Socket CPU	High
Language Translation	1X	Dual-Socket CPU	High
Sigma	1X	Dual-Socket CPU	High
Speech Recognition	1X	Dual-Socket CPU	High

TABLE III

RESOURCE REQUIREMENTS OF ONLINE INFERENCE WORKLOADS.





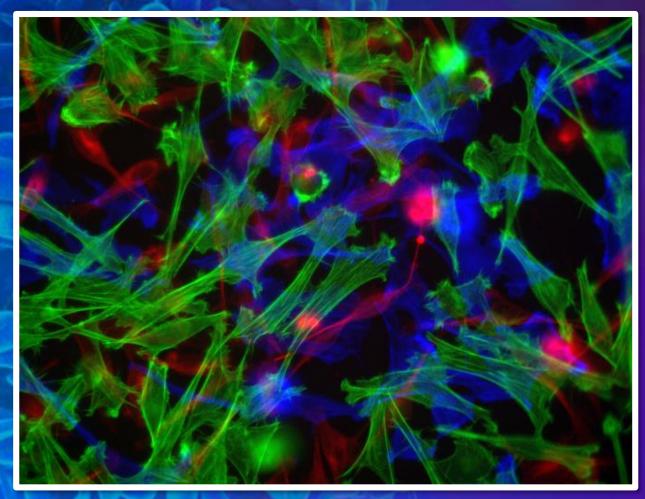
A SCIENTIFIC COLLABORATION BETWEEN INTEL AND NOVARTIS

KUSHAL DATTA, PHD
RESEARCH SCIENTIST
INTEL AI





26X LARGER



1024 X 1280 X 3

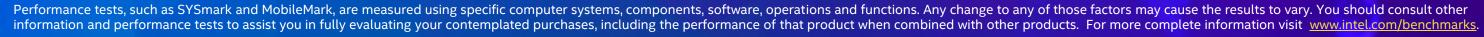
IMAGENET



224 X 224 X 3

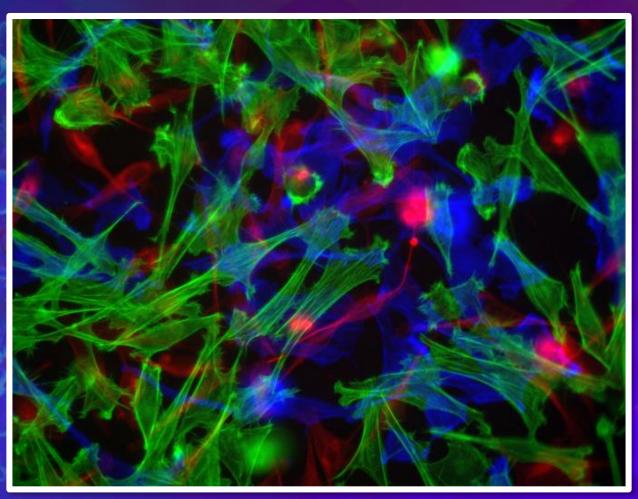


Software and workloads used in performance tests may have been optimized for performance only on Intel® microprocessors.





26X LARGER MULTIPLE OBJECTS



1024 X 1280 X 3

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

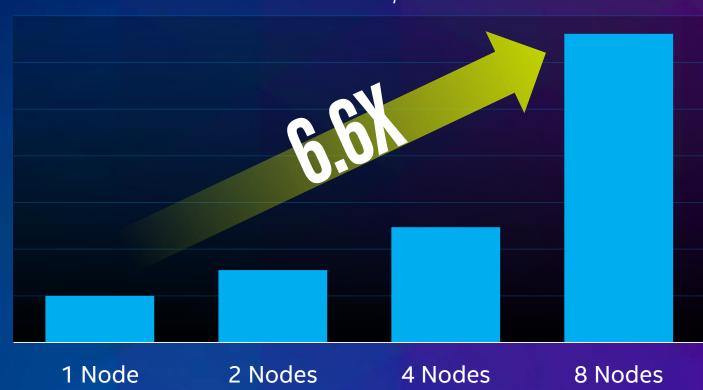
Software and workloads used in performance tests may have been optimized for performance only on Intel® microprocessors.

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 www.intel.com/benchmarks.

HIGH PERFORMANCE AT SCALE

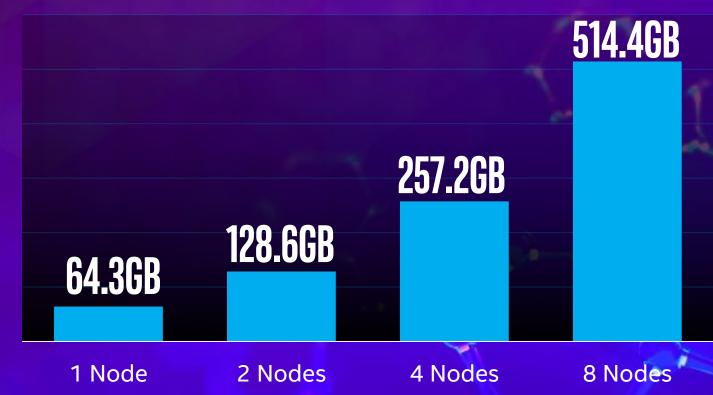
SCALING OF TIME TO TRAIN

Intel® Omni-Path Architecture, Horovod and TensorFlow®



TOTAL MEMORY USED

192GB DDR4 PER INTEL® 2S XEON® 6148 PROCESSOR



MULTISCALE CONVOLUTION NEURAL NETWORK

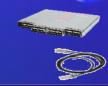


OPTIMIZED LIBRARIES

Intel® MKL/MKL-DNN, clDNN, DAAL

INTEL® OMNI-PATH ARCHITECTURE







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

[§] 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







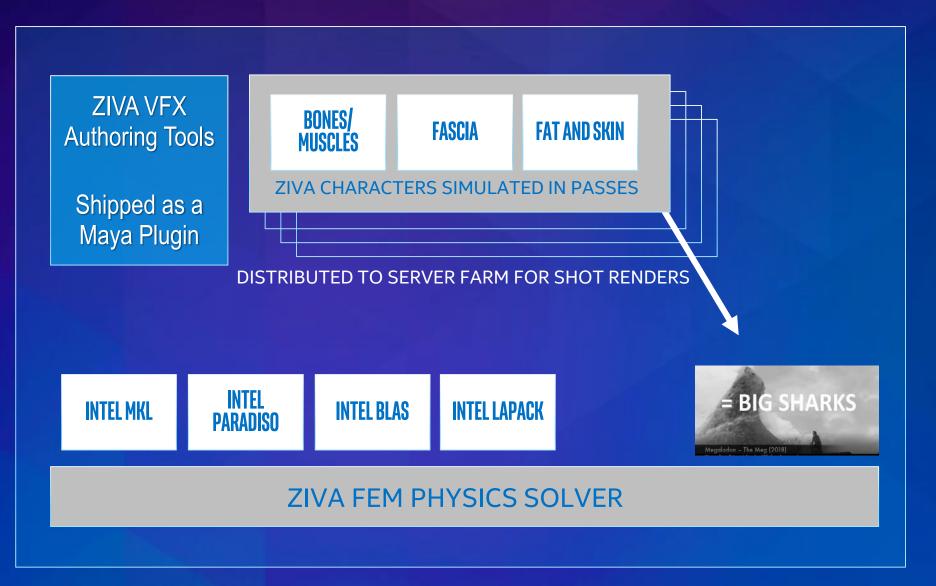


MACHINE LEARNING IMAGE RENDERING

JAMES JACOBS CEO OF ZIVA



THE POWER OF ZIVA SIMULATION AND COMPUTE IN TODAY'S MEDIA & ENTERTAINMENT



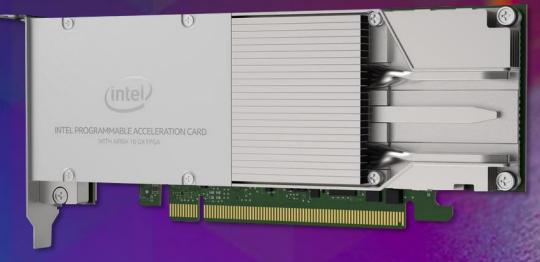






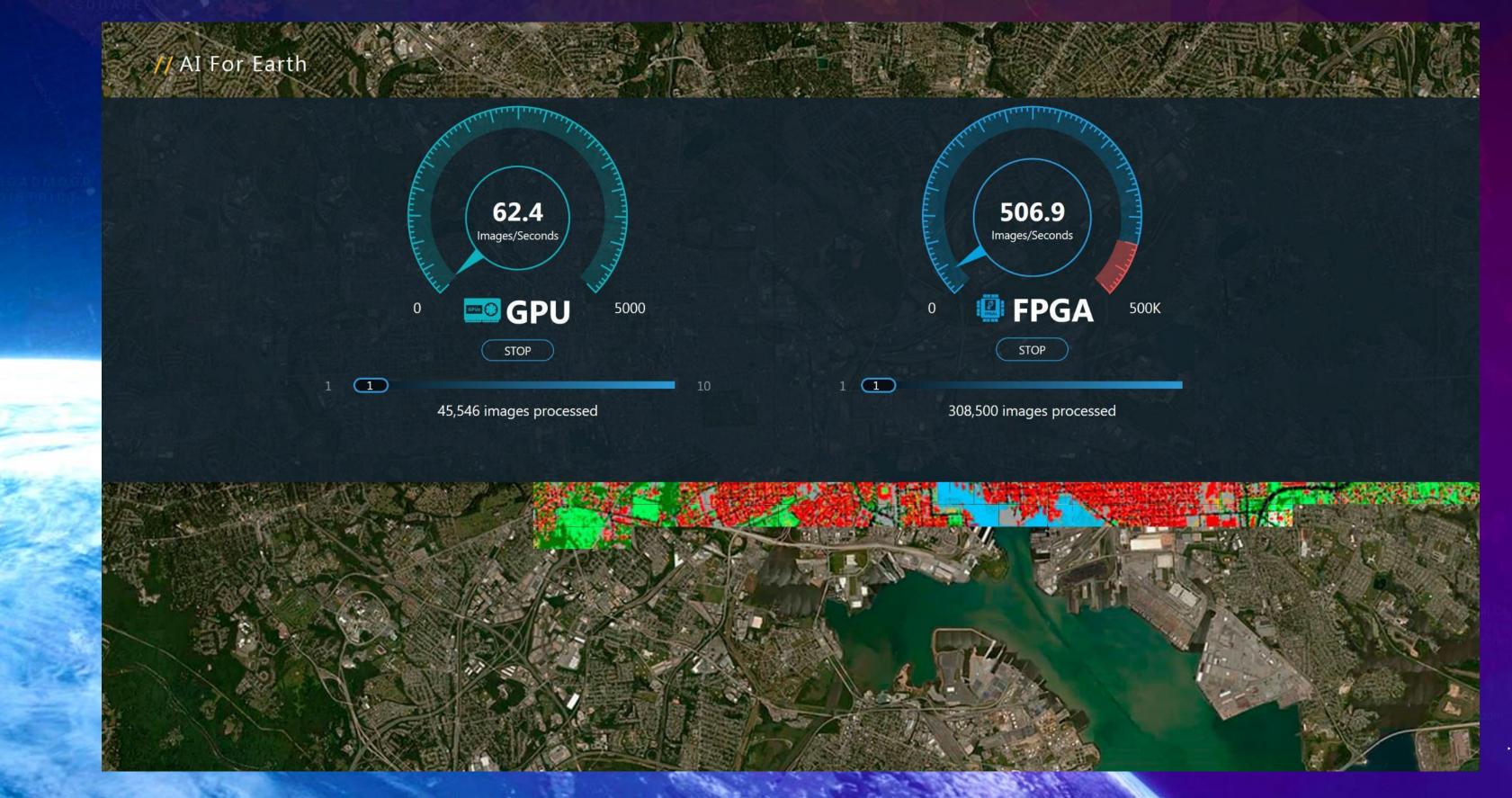






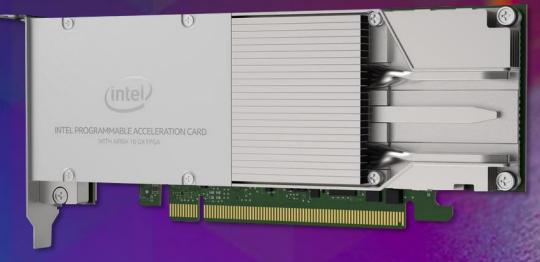
FLEXIBLE REAL-TIME INFERENCING FPGAPRODUCTS





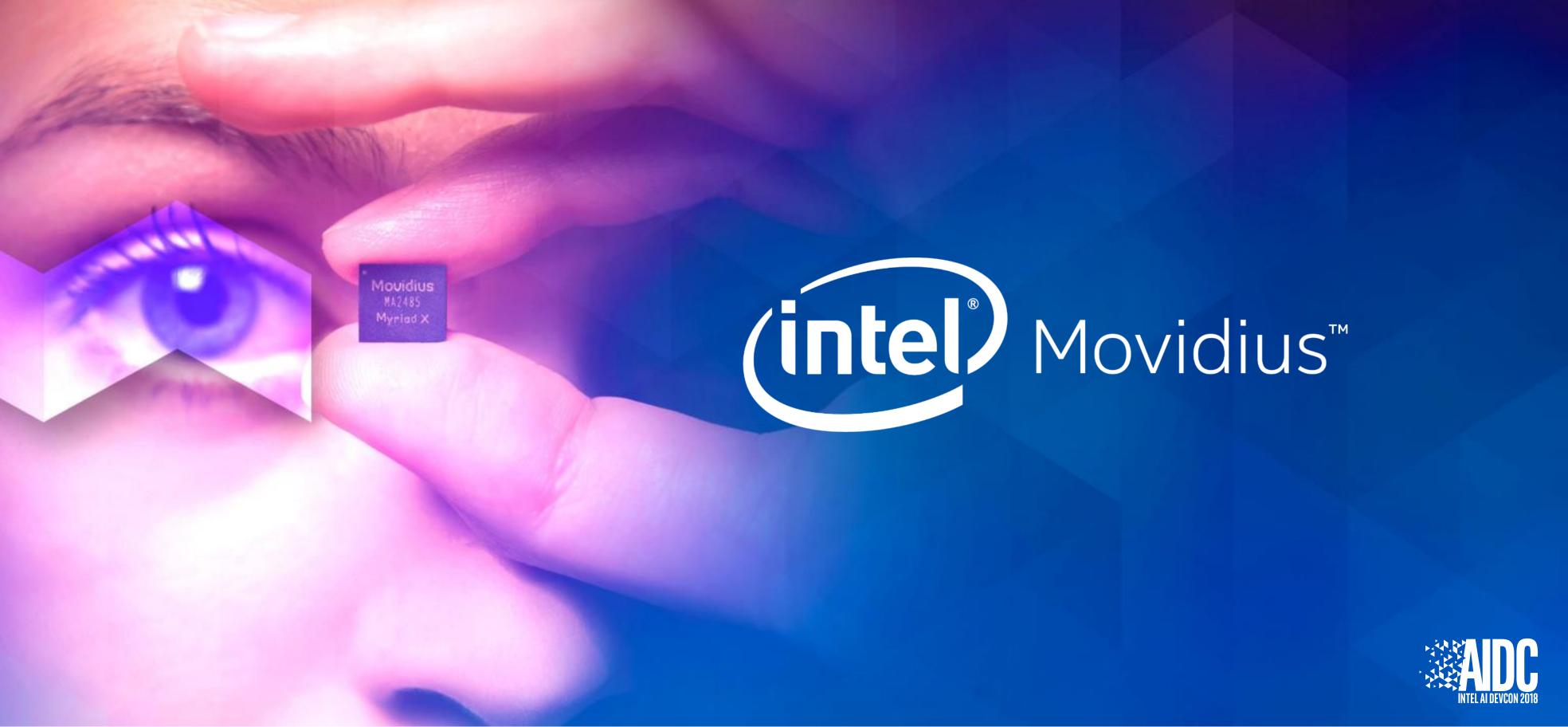




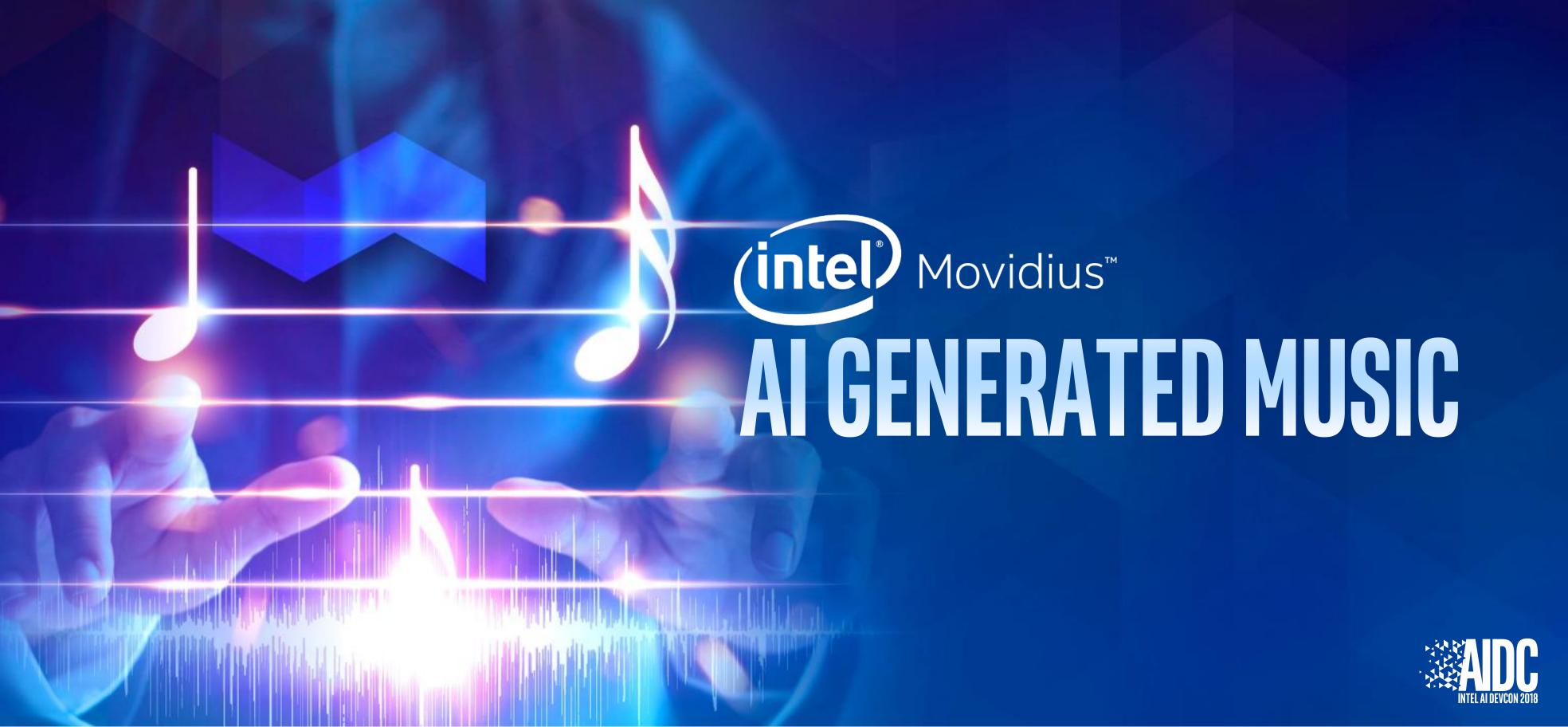


FLEXIBLE REAL-TIME INFERENCING FPGAPRODUCTS







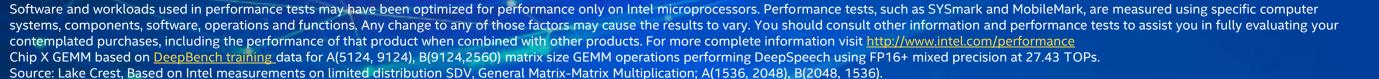




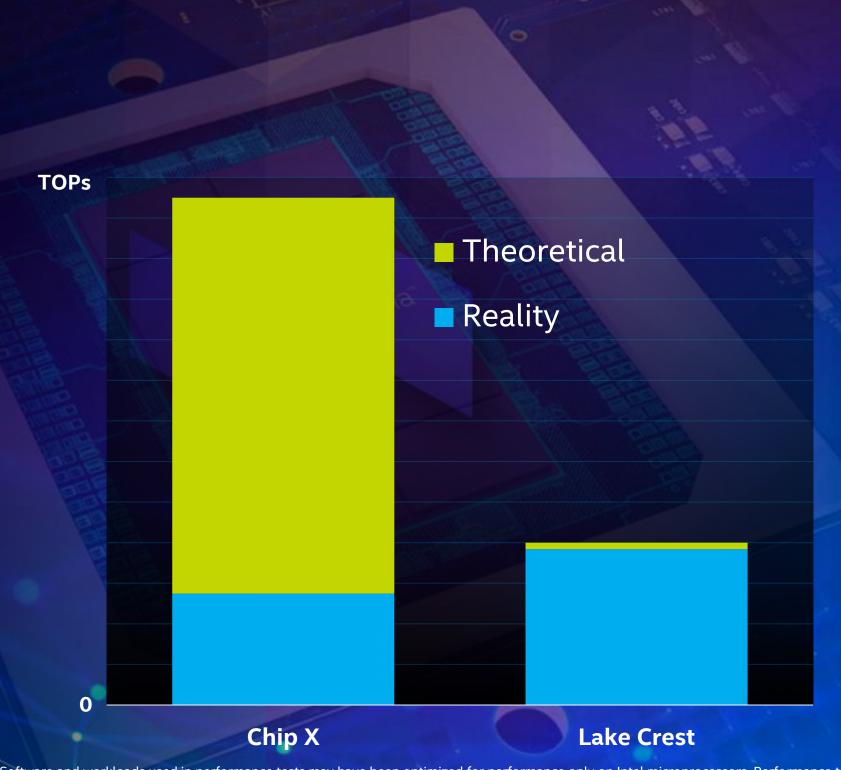














HIGH UTILIZATION & MODEL PARALLELISM

GEMM OPERATION UTILIZATION¹

96.4%

A(1536, 2048) B(2048, 1536) **MULTI-CHIP** SCALING²

96.2%

A(6144, 2048) B(2048, 1536)

MULTI-CHIP COMMUNICATION³

2.4 TB/S

OFF CHIP BANDWIDTH <790ns LATENCY

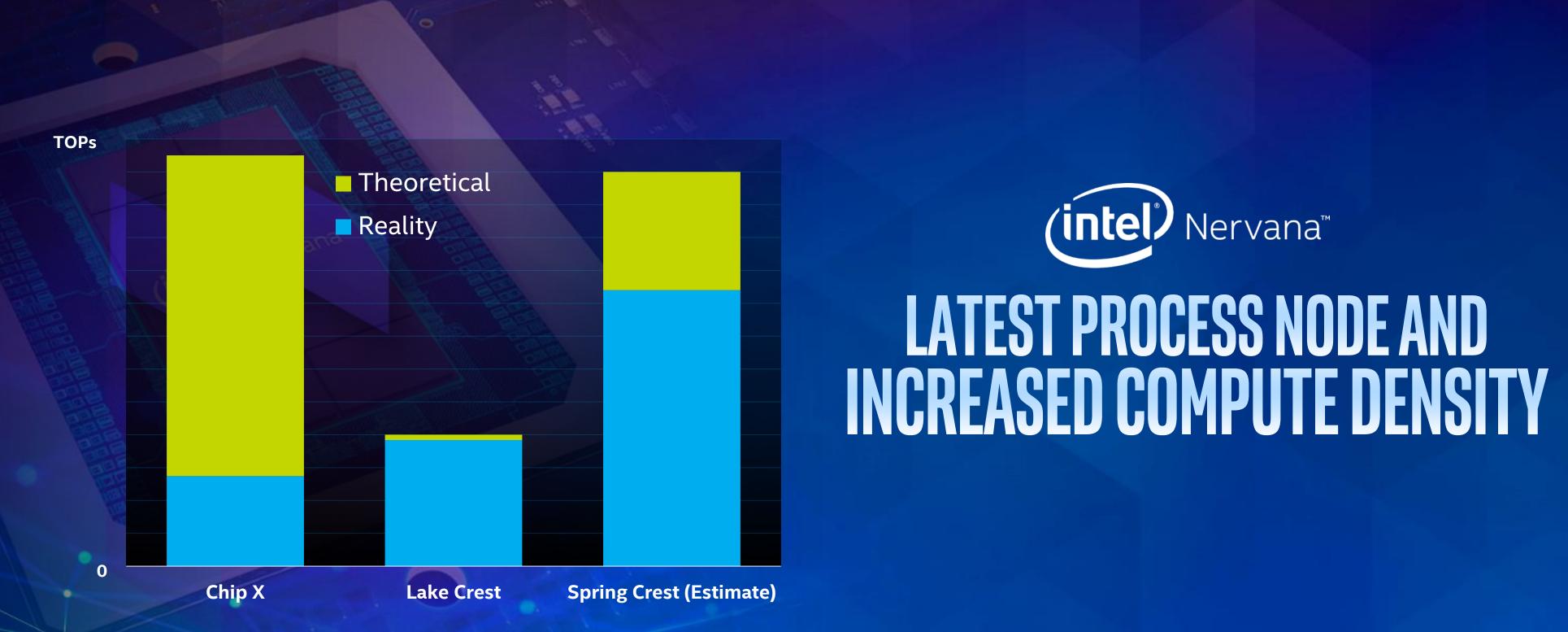
POWER < 210W

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 Chip X GEMM based on DeepBench training data for A(5124, 9124), B(9124,2560) matrix size GEMM operations performing DeepSpeech using FP16+ mixed precision at 27.43 TOPs.

Source: Lake Crest: Based on Intel measurements on limited distribution SDV 1 General Matrix-Matrix Multiplication; A(1536, 2048), B(2048, 1536)

2 Two chip vs. single chip GEMM performance; A(6144, 2048), B(2048, 1536) 3 Full Chip MRB-CHIP MRB data movement using send/recv, Tensor size = (1, 32), average across 50K iterations





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
Chip X GEMM based on DeepBench training data for A(5124, 9124), B(9124,2560) matrix size GEMM operations performing DeepSpeech using FP16+ mixed precision at 27.43 TOPs.
Source: Lake Crest - Based on Intel measurements on limited distribution SDV

Source: Spring Crest - Intel measurements on simulated product



FIRST COMMERCIAL NNP INTEL® NERVANA NNP L-1000 in 2019

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

SPRING CREST

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



ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION









ENABLE DEVELOPERS TO ACHIEVE THEIR AI VISION











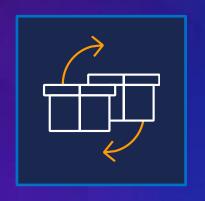




Machine Learning at Amazon: a long heritage



Personalized recommendations



Fulfillment automation / inventory management



Cargo



Voice driven interactions



Inventing entirely new customer experiences

achine Learning Platform





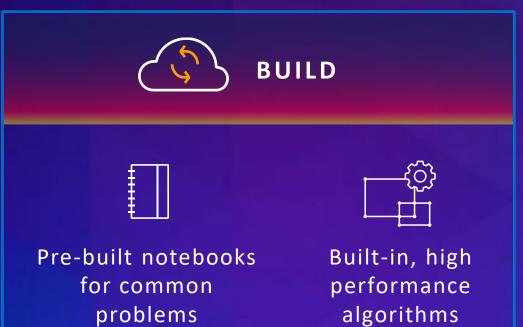






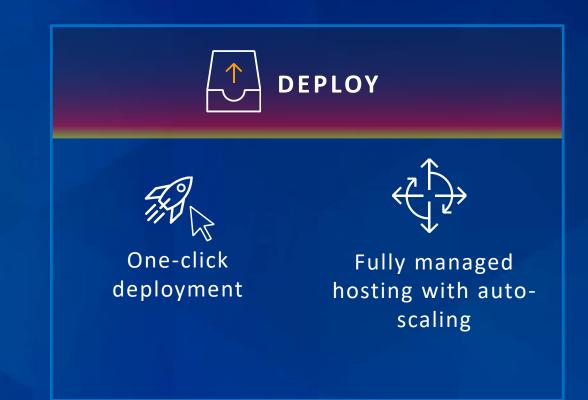


Amazon SageMaker











AWS Machine Learning Stack





AWS DEEPLENS



INFRASTRUCTURE



EC2 GPUs



EC2 CPUs



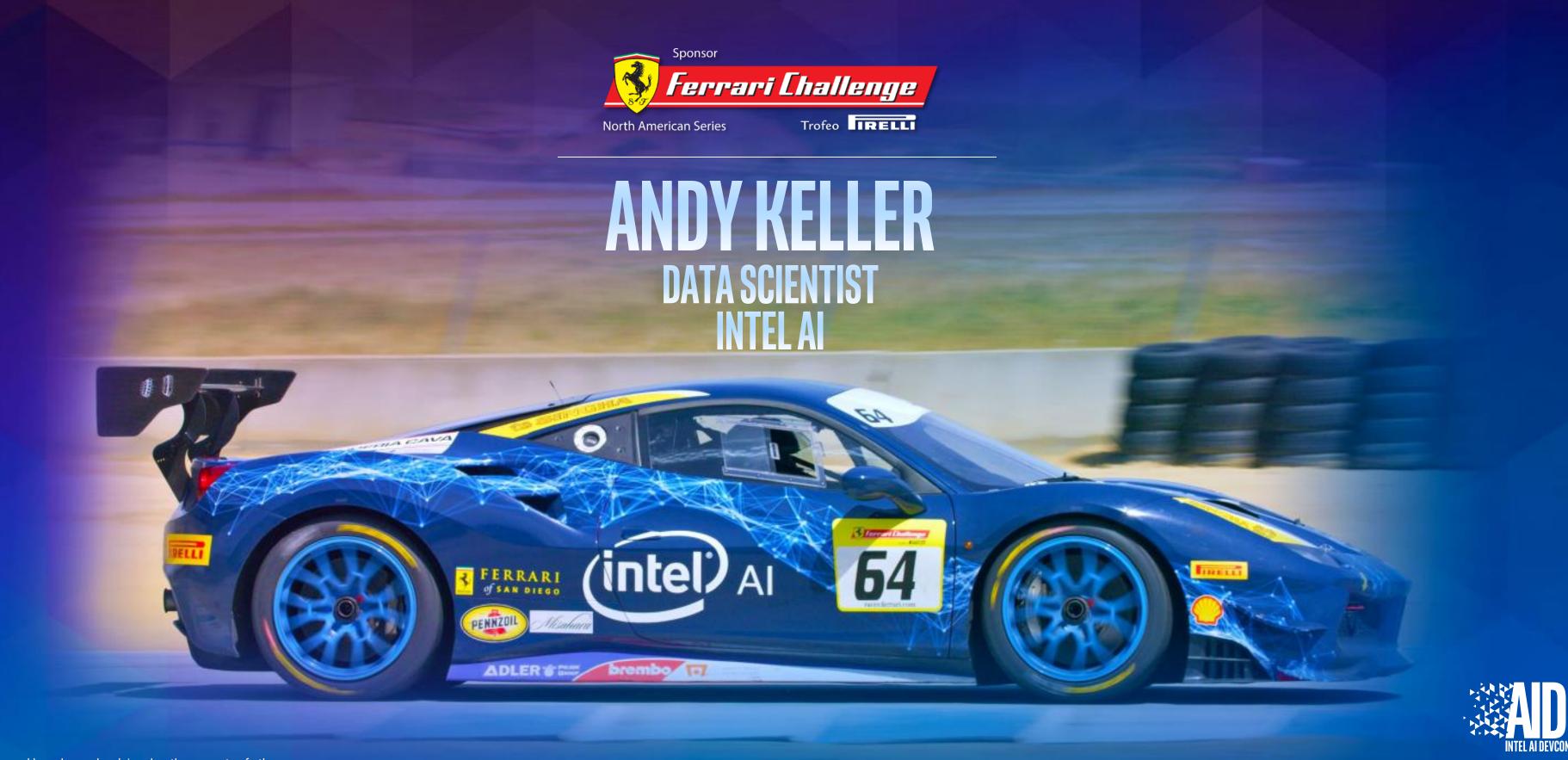
IoT



Edge









YINYIN LIU

HEAD OF DATA SCIENCE INTEL AI

OPEN SOURCE LIBRARIES

NLP ARCHITECT

NervanaSystems/nlp-architect

RLCOACH



NEURAL NETWORK DISTILLER











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.



DISCLAIMERS

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 information go to www.intel.com/benchmarks.

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice Revision #20110804 Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at [intel.com].

Results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling, and provided to you for informational purposes. Any differences in your system hardware, software or configuration may affect your actual performance.

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, Xeon, Intel Nervana and Intel Movidius are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

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

© Intel Corporation.



