# Legal Disclaimers

Information in this document is provided in connection with Intel® products. No license, express or implied, by estoppels or otherwise, to any intellectual property rights is granted by this document. Except as provided in Intel's Terms and Conditions of Sale for such products, Intel assumes no liability whatsoever, and Intel disclaims any express or implied warranty, relating to sale and/or use of Intel® products including liability or warranties relating to fitness for a particular purpose, merchantability, or infringement of any patent, copyright or other intellectual property right. Intel products are not intended for use in medical, life saving, or life sustaining applications.

This document contains information on products in the design phase of development. The information here is subject to change without notice. Do not finalize a design with this information.

Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them. Intel may make changes to specifications and product descriptions at any time, without notice.

Intel processors and chipsets 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.

Intel® Advanced Vector Extensions (Intel® AVX)* provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause a) some parts to operate at less than the rated frequency and b) some parts with Intel® Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration and you can learn more at http://www.intel.com/go/turbo.

Intel® processors of the same SKU may vary in frequency or power as a result of natural variability in the production process.

Performance estimates or simulated results based on internal Intel analysis or architecture simulation or modeling are provided to you for informational purposes. Any differences in your system hardware, software or configuration may affect your actual performance

Intel, Intel Xeon, Intel Itanium, and Intel Netburst are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries.

# Legal Disclaimer & Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

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.

## 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

# Agenda

- Brief Deep Learning Training Overview
- AI Optimizations
- Why Multi-node Training?
- Summary

# ARTIFICIAL INTELLIGENCE



**ARTIFICIAL INTELLIGENCE**
A program that can sense, reason, act, and adapt

**MACHINE LEARNING**
Algorithms whose performance improve as they are exposed to more data over time

**DEEP LEARNING**
Subset of machine learning in which multilayered neural networks learn from vast amounts of data

# MACHINE LEARNING

## CLASSIC MACHINE LEARNING

How do you engineer the best features?

$N \times N$



$(f_1, f_2, ..., f_K)$

Roundness of face
Dist between eyes
Nose width
Eye socket depth
Cheek bone structure
Jaw line length
...etc.

**CLASSIFIER ALGORITHM**

SVM
Random Forest
Naïve Bayes
Decision Trees
Logistic Regression
Ensemble methods

**Arjun**

## DEEP LEARNING

How do you guide the model to find the best features?

$N \times N$



**NEURAL NETWORK**

input layer    hidden layer 1    hidden layer 2    hidden layer 3

output layer

~60 million parameters

**Arjun**

# INTEL® NERVANA™ PORTFOLIO

**EXPERIENCES**

**PLATFORMS**
Intel® Nervana™ Cloud & Appliance
Intel® Nervana™ DL Studio
Intel® Computer Vision SDK
Intel® Movidius™ MDK
(intel) Saffron™

**FRAMEWORKS**
Apache Spark™ MLlib bigDL
neon
TensorFlow
mxnet
Microsoft CNTK
torch
Caffe
Caffe2
Chainer
theano

**LIBRARIES**
python Intel Python Distribution
Intel® Data Analytics Acceleration Library (DAAL)
Intel® Nervana™ Graph*
Intel® Math Kernel Library (MKL, MKL-DNN)

**HARDWARE**
intel XEON inside | intel ARRIA 10 inside | intel CORE i7 inside | intel ATOM inside
Compute
Memory & Storage
Networking

INSIDE AI

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

(intel) Nervana™

9

# DATACENTER AI

Intel® Stratix® 10 FPGA

Intel® Xeon® Processor Scalable Family

Intel® Nervana™ Neural Network Processor*

## FLEXIBLE ACCELERATION

Accelerate the widest range of AI and other workloads & configurations

## FOUNDATION FOR AI

Begin your journey with the AI you need on the chip you know

## DEEP LEARNING BY DESIGN

Accelerate the most intensive deep learning deployments with this custom-built processor

*Future product that was formerly codenamed the "Crest family"
All performance positioning claims are relative to other processor technologies in Intel's AI datacenter portfolio
All products, computer systems, dates, and figures specified are preliminary based on current expectations, and are subject to change without notice.

# AI Optimizations

# Performance Drivers for AI Workloads

**Compute**

**Bandwidth**

**SW Optimizations**

# Hardware Optimized Libraries and Frameworks

**Scaling**

- Improve load balancing

- Reduce synchronization events, all-to-all comms

**Utilize all the cores**

- OpenMP, MPI

- Reduce synchronization events, serial code

- Improve load balancing

**Vectorize/SIMD**

- Unit strided access per SIMD lane

- High vector efficiency

- Data alignment

**Efficient memory/cache use**

- Blocking

- Data reuse

- Prefetching

- Memory allocation

Important to use optimized software frameworks and libraries for best AI workload performance

# Deep Learning Software Summary

Nervana Deep Learning Studio

Titanium: HW mgmt.

- Data scientist and developer DL productivity tools
- Build and train models
- Compress and export modes to end points
- SaaS or open source

Frameworks



- Frameworks for developers
- Back end APIs to Nervana Graph

Nervana Graph

MKL-DNN ,other math libraries

HW Transformers, Non-x86 libraries

- Accelerate framework optimization on IA
- For framework developers & Intel
- Multi-node optimizations
- Open source

*Intel manages low level SW transition from one IA HW to another*

System SW: Drivers

- HW specific SW
- For system developers
- Delivered through HW sales

# Majority of Cycles Spent in Convolution

Percentage of time spent on individual layers of different network topologies running inference with Caffe* framework on Intel® Xeon® Processor E5-2699 v4



Legend: ■ Convolution  ■ InnerProduct  ■ ReLU  ■ LRN  ■ Pooling  ■ BatchNorm  ■ Eltwise  ■ Concat

## 60 to 85% Time Spent in Convolution - AI workloads such as image recognition are compute heavy

Batch Sizes AlexNet:256 VGG-19: 64 ResNet-50: 50 GoogleNet-V1: 96  Configuration Details on Slide: 27

# Building blocks (primitives) from mkl-dnn

1. Matrix multiplication

2. Convolution and direct batched convolution

3. Inner product

4. Pooling: maximum, minimum, average

5. Normalization

6. Activation: rectified linear unit (ReLU)

7. Data manipulation: multi-dimensional transposition

https://software.intel.com/en-us/articles/tensorflow-optimizations-on-modern-intel-architecture

# Intel® Nervana™ graph

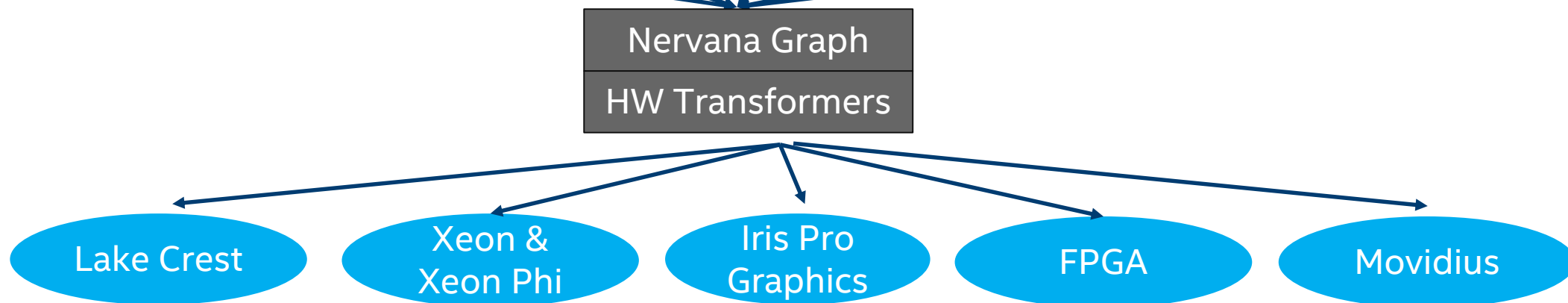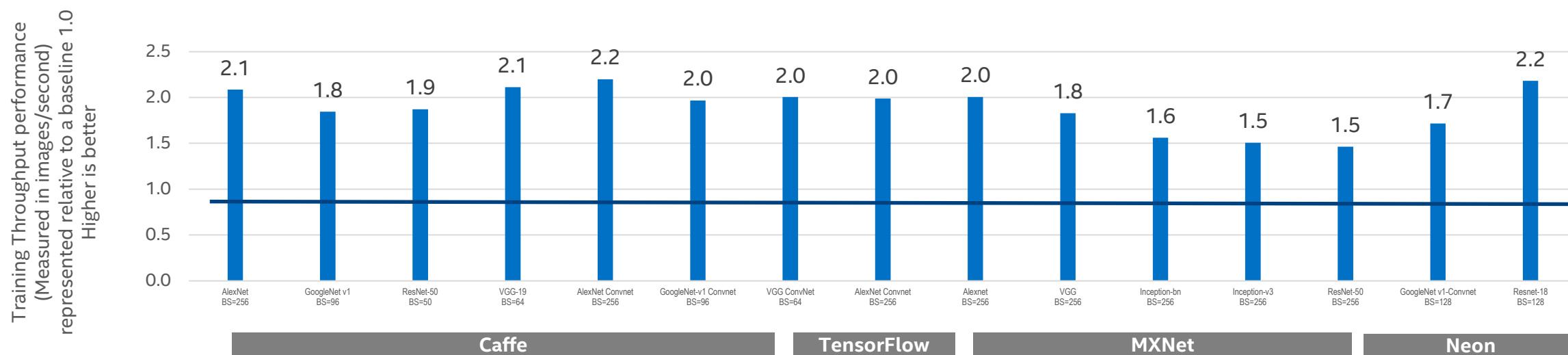## High-Performance Execution Graph for Neural Networks

Use Cases

Models

Frameworks

| neon | Caffe{2} | Torch | TensorFlow | CNTK | ... |

Nervana Graph

HW Transformers

Hardware

Lake Crest

Xeon & Xeon Phi

Iris Pro Graphics

FPGA

Movidius

# Up to 2.2x Higher Training Throughput
## on Intel® Xeon® Platinum 8180 Processor

Intel® Xeon® Platinum 8180 Processor Training throughput
over Intel® Xeon® Processor E5-2699 v4



Y-axis: Training Throughput performance (Measured in images/second) represented relative to a baseline 1.0 Higher is better

Bar values and labels:
- AlexNet BS=256: 2.1
- GoogleNet v1 BS=96: 1.8
- ResNet-50 BS=50: 1.9
- VGG-19 BS=64: 2.1
- AlexNet Convnet BS=256: 2.2
- GoogleNet-v1 Convnet BS=96: 2.0
- VGG ConvNet BS=64: 2.0
- AlexNet Convnet BS=256: 2.0
- Alexnet BS=256: 2.0
- VGG BS=256: 1.8
- Inception-bn BS=256: 1.6
- Inception-v3 BS=256: 1.5
- ResNet-50 BS=256: 1.5
- GoogleNet v1-Convnet BS=128: 1.7
- Resnet-18 BS=128: 2.2

Frameworks: Caffe | TensorFlow | MXNet | Neon

**Intel® Xeon® Platinum Processor delivers high Training throughput performance across different frameworks**

# Up to 2.4x Higher Inference Throughput
## on Intel® Xeon® Platinum 8180 Processor

Intel® Xeon® Platinum 8180 Processor  Inference throughput
over Intel® Xeon® Processor E5-2699 v4



**Intel® Xeon® Platinum Processor delivers high Inference throughput performance across different frameworks**

INFERENCE using FP32 Configuration Details on Slide:24,25
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  Source: Intel measured as of June 2017

# INTEL® XEON® PROCESSOR PLATFORM PERFORMANCE

## Hardware plus optimized software
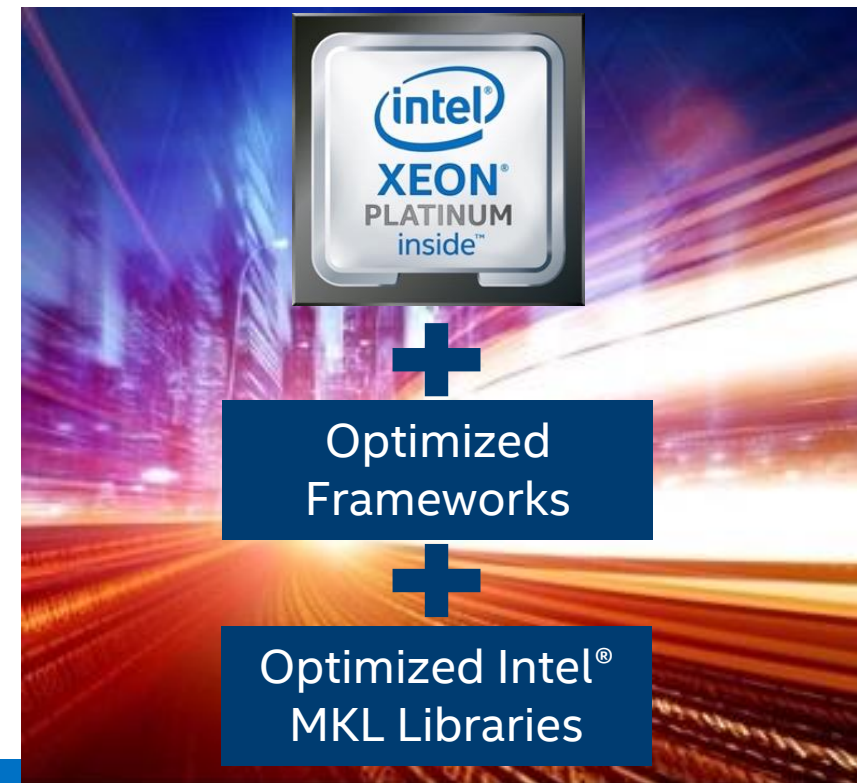
### INFERENCE THROUGHPUT

Up to
**138x**

Intel® Xeon® Platinum 8180 Processor
higher Intel optimized Caffe GoogleNet v1 with Intel® MKL
inference throughput compared to
Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

### TRAINING THROUGHPUT

Up to
**113x**

Intel® Xeon® Platinum 8180 Processor
higher Intel Optimized Caffe AlexNet with Intel® MKL
training throughput compared to
Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

Inference and training throughput measured with FP32 instructions. Inference with INT8 will be higher.

**+**
Optimized
Frameworks
**+**
Optimized Intel®
MKL Libraries

## Deliver significant AI performance with hardware and software optimizations on Intel® Xeon® Scalable Processors

Why Multi-node Training?

# Accuracy and reduced time to train:
## *The goal to drive innovation*

**Properties of Neural Networks[1]:**

**Results get better with**

- More Data

- Bigger models

- More computation.

**Training is not a one time effort.**

- Many operational neural networks as part of different applications

- Each neural network may be trained with domain specific training sets

- Evolving input data sets drives the need for re-training

Google Brain's Jeff Dean quantifies benefits of reducing the time to train:

**Minutes, hours:**
- Interactive research
- Instant gratification of results

**1-4 Days:**
- Tolerable
- Interactivity replaced by running many experiments in parallel

**1-4 weeks:**
- High value experiments only
- Progress stalls

**>1 Month:**
- Don't even try.

*Jeff Dean: Large-Scale Deep Learning for Intelligent Computer Systems (GoogleBrain)*

**Striving for interactive research drives need for more computational power and multi-node training options.**

1 https://static.googleusercontent.com/media/research.google.com/en//people/jeff/BayLearn2015.pdf

# Progression towards Multi-node Training

Let's solve this problem using DL → Start with a server & GPU → Interesting results. Drive for better accuracy

- More data for better results
- Now training takes much longer
- How to scale and keep time-to-train manageable

*Options*

Single server takes too long to train now

- Scale out training infrastructure
- Multi-Node
- Fabric interconnected training

**Increased accuracy and reduced time to train**

# How can Deep Learning Training be Parallelized?



## Data parallelism

- Each system runs on its own dataset

- Communication occurs at each iteration, but less frequently as model parallelism

- Fast, non-blocking communication best to insure computation is not waiting on data.

## Model parallelism

- Share the Neural Network across many nodes

- Communication occurs for layers in each iteration; creates lots of communication

- Only as fast as slowest machine due to interactivity of code

# Intel® Omni-Path Architecture & AI

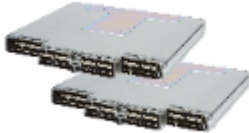## World-Class Interconnect Solution for Shorter Time to Train

**HFI Adapters & Integrated Processors**
*Single port*
x8 and x16

**Edge Switches**
*1U Form Factor*
24 and 48 port

**Director Switches**
*QSFP-based*
192 and 768 port

**Software**
*Open Source*
Host Software and Fabric Manager

**Cables**
*Third Party Vendors*
Passive Copper
Active Optical

*Fabric interconnect for breakthrough performance on scale-out workloads like deep learning training*

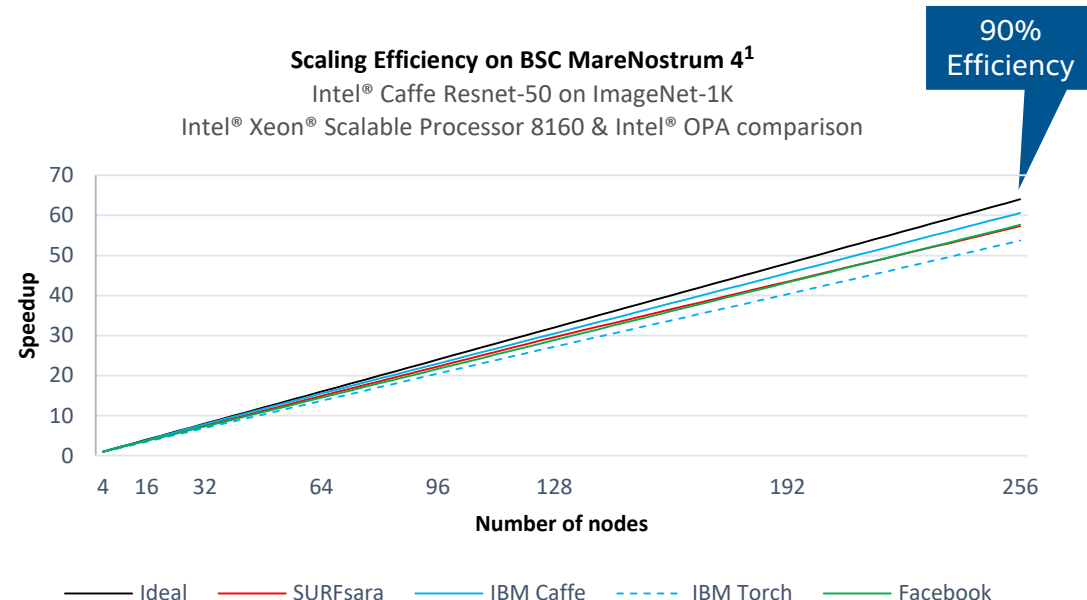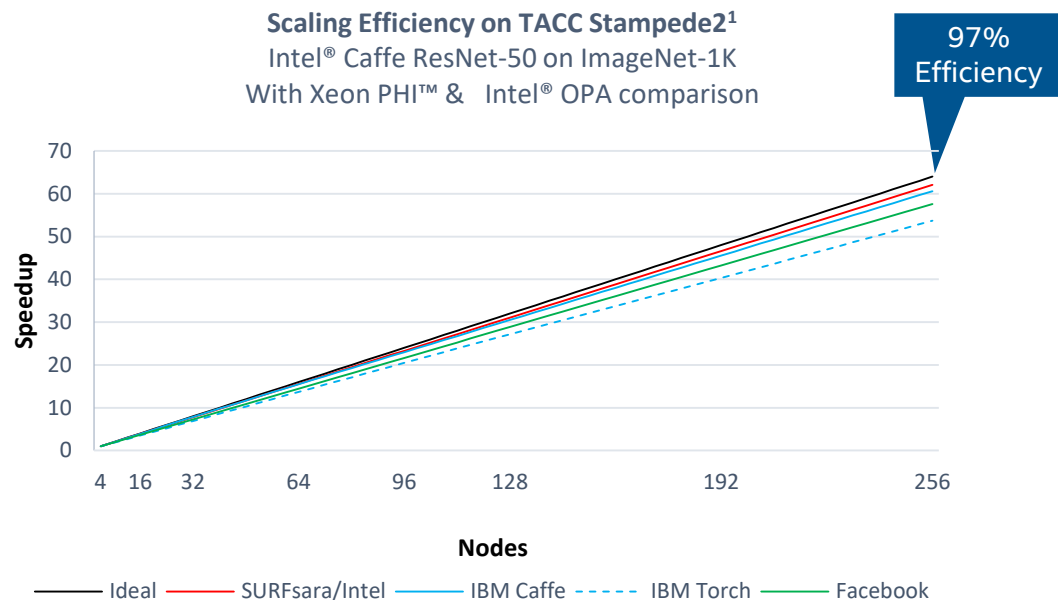## OMNI-PATH BENEFITS FOR AI

- High bandwidth (100Gbps), low latency and high message rate
- Improve performance, reliability & QoS:
  - Traffic Flow Optimization to maximize QoS in mixed traffic
  - Packet Integrity Protection for rapid and transparent recovery of transmission errors
  - Dynamic lane scaling to maintain link continuity
  - Switches offered with redundancy and hot swappable FRUs
- Heterogeneous cluster support
- Excellent price/performance & price/port, 48 radix
- Dispersive routing
- Multi-Rail support to provide high injection rate
- CPU/OPA integration for increased performance, lower power & cost

## BREAKTHROUGH PERFORMANCE

- ImageNet-1K training in **less than** 40 minutes with Intel® Caffe[2]
- 90%+ efficiency with Resnet-50 training to 256 nodes and beyond on Xeon® and Xeon® Phi™
- Reduce time to train
- #1 Green500 (June 2017) cluster designed with AI in mind
- Reduced communication latency compared to InfiniBand EDR[1]:
  - Up to **21%** Higher Performance, lower latency at scale
  - Up to **17%** higher messaging rate
  - Up to **9%** higher application performance

# Multi-Node Intel® Caffe Resnet-50 on ImageNet-1K:
## *High Scaling Efficiency with Intel® OPA*

**Scaling Efficiency on TACC Stampede2[1]**
Intel® Caffe ResNet-50 on ImageNet-1K
With Xeon PHI™ &  Intel® OPA comparison

**97% Efficiency**



Legend: Ideal — SURFsara/Intel — IBM Caffe — IBM Torch (dashed) — Facebook

**Scaling Efficiency on BSC MareNostrum 4[1]**
Intel® Caffe Resnet-50 on ImageNet-1K
Intel® Xeon® Scalable Processor 8160 & Intel® OPA comparison

**90% Efficiency**



Legend: Ideal — SURFsara — IBM Caffe — IBM Torch (dashed) — Facebook

- **TACC Stampede 2**
  - **97% scaling efficiency** from 4 to 256 Intel® Xeon Phi™ 7250 nodes interconnected with Intel® OPA
  - Convergence with Top1/5 > **74%/92%**
  - 4 - 256 node runs:  batch size of 16 per node, scaling efficiency of 97% in **63 minutes**

- **BSC MareNostrum 4**
  - Convergence with Top1/5 > **74%/92%**
  - 4 - 256 node runs:  Batch size of 32 per node, 90% scaling efficiency, Total time to train: 70 Minutes

**Strong multi-node training, with high accuracy with Intel® OPA**

1. https://blog.surf.nl/en/imagenet-1k-training-on-intel-xeon-phi-in-less-than-40-minutes/
More Information
  • https://www.bsc.es/user-support/mn4.php
  • http://portal.tacc.utexas.edu/user-guides/stampede2
  • Goyal, Priya, et al. "Accurate, Larg Minibatch SGD: Training ImagNetin 1 Hour." arXiv preprint arXiv:1706.02677 (2017)
  • Cho, Minsik, et al. "powerAI DDL." arXiv preprint arXiv:1708.02188 (2017)
  • IBM claims 95% scaling efficiency and Facebook claims 89%

# Multi-Node Intel® Caffe Resnet-50 on ImageNet-1K:
## *Time to Train on Intel® Xeon® and Xeon Phi™ with Intel® OPA*

**Time to Train (Minutes)[1]**
**TACC Stampede2 with Intel® Xeon Phi™ and Intel® OPA**



**Time to Train (Minutes)[1]**
**BSC MareNostrum 4 with Intel® Xeon® Scalable Processor 8160 & Intel® OPA**



- Scaling continues, up to 768 nodes
- Efficiency reduced due to number of iterations (epochs) needed to achieve 74%/92% Top1/Top5 accuracy

1. https://blog.surf.nl/en/imagenet-1k-training-on-intel-xeon-phi-in-less-than-40-minutes/
More information
- https://www.bsc.es/user-support/mn4.php
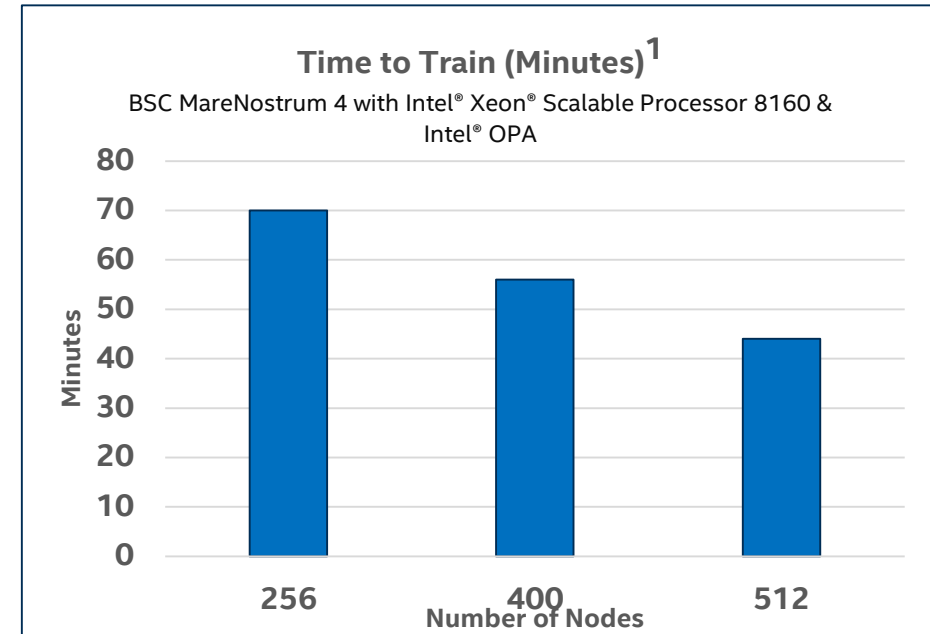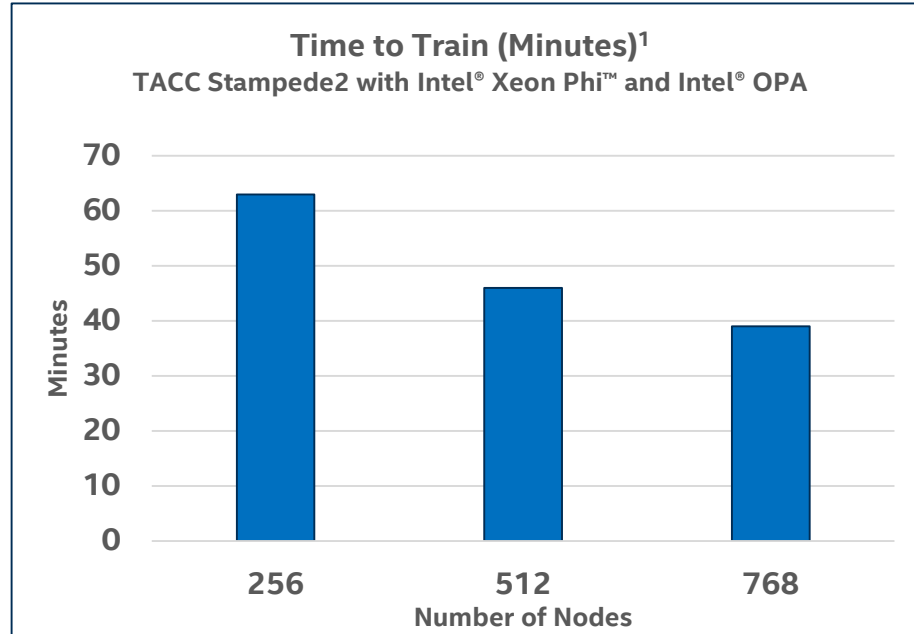- http://portal.tacc.utexas.edu/user-guides/stampede2

# TensorFlow on OPA

- OPA based systems support TensorFlow with any of gRPC, Verbs and MPI communications interfaces [1]

- Preliminary multi-node measurements with Resnet-50 show good scaling with TensorFlow (using gRPC & sockets)

- Measurement and tuning efforts with various models and MPI are making good progress

- MPI Multi-End Point capability for TensorFlow in PoC stage

**ResNet-50 Performance Scaling on TensorFlow 1.2 [1]
(gRPC sockets on KNL/OPA, Batch Size=128 per node)**

Worker Nodes - Parameter Server Nodes

Images/sec    Scaling Efficiency

Source: Intel, Aug 2017

# Multi-node GoogLeNet comparison with Intel® Caffe on Intel® Omni-Path Architecture (Intel® OPA)

- Intel® Caffe and MLSL have been optimized for Intel® Xeon® processors & Intel® OPA

- Parallelization using Machine Learning Scaling Library (MLSL)
  - MLSL abstracts communication patterns and supports data/model/hybrid parallelism
  - Leverages Intel MPI optimizations for OPA for communication but could use other runtimes/message layers
  - MLSL API is being designed to be applicable to variety of popular Frameworks (e.g. Caffe, Torch, Theano, etc)
  - MLSL also provides statistical data collection to monitor time spent on different operations; including computation and communication.

### GoogLeNet Performance of Intel Caffe + MLSL on Broadwell/OPA



*Y-axis: Throughput Scaling Efficiency %, X-axis: Nodes (1, 2, 4, 8, 16, 24, 32, 48)*

Legend: GoogLeNet-V1, GoogLeNet-V2

# Naïve Multi-Node Implementation

- Global averaging (usually Allreduce) of weights for all layers performed across nodes either :
  - At beginning of iteration (e.g. before $F_1$) .
  - After weight update $U_k$

- No overlapping of computation & communication



Iteration i

Allreduce

Iteration i+1

$F_k$ – Forward propagation computation for layer k
$B_k$ – Back propagation computation for layer k
$U_k$ – Local weight update for layer k
$S_{1-N}$ – Allreduce of weights for all layers

# MLSL – Multi-node communication technique

- Overlap computation & communication by using Non-blocking Allreduce (MPI_Iallreduce) of weights after back propagation for each layer

- Note for first layer ($F_1$) communication cannot be hidden but is partially/totally hidden for subsequent layers



Iteration i

Iteration i+1

MPI_Iallreduce

$F_k$ – Forward propagation computation for layer k
$B_k$ – Back propagation computation for layer k
$U'_k$ – Weight update & MLSL servers perform non-blocking Allreduce of weights for layer k
$S'_k$ – Synchronization (Waits) for completion of MPI_Iallreduce from $U'_k$

# RESOURCES

# INTEL® NERVANA™ AI ACADEMY

## For developers, students, instructors and startups

### LEARN

- Online tutorials
- Webinars
- Student kits
- Support forums

### DEVELOP

- Intel Optimized Frameworks
- Exclusive access to Intel® Nervana™ DevCloud

### TEACH

- Comprehensive courseware
- Hands-on labs
- Cloud compute
- Technical Support

### SHARE

- Project showcase opportunities at
- Intel Developer Mesh
- Industry & Academic events

## software.intel.com/ai

# INTEL® NERVANA™ DEVCLOUD

## Free AI cloud access for Intel® Nervana™ AI Academy members

✓ Get started for FREE

✓ 4 weeks access to remote cluster of Intel® Xeon® Scalable processors

✓ 200 GB file storage

✓ Pre-configured libraries & frameworks†

## software.intel.com/ai/DevCloud

†neon™ framework, Intel® Optimization for Theano*, Intel® Optimization for TensorFlow*, Intel® Optimization for Caffe*, Intel® Distribution for Python* (including NumPy, SciPy, and scikit-learn*), Keras* library
*Other names and brands may be claimed as the property of others.

# FIND OUT MORE

**LEARN**

**Find out more at www.intelnervana.com**

**EXPLORE**

**Use Intel's performance-optimized libraries & frameworks**

**ENGAGE**

**Contact your Intel representative for help and POC opportunities**

# BACKUP

# Configuration Details

Platform: 2S Intel® Xeon® Platinum 8180 CPU @ 2.50GHz (28 cores), HT disabled, turbo disabled, scaling governor set to "performance" via intel_pstate driver, 384GB DDR4-2666 ECC RAM. CentOS Linux release 7.3.1611 (Core), Linux kernel 3.10.0-514.10.2.el7.x86_64. SSD: Intel® SSD DC S3700 Series (800GB, 2.5in SATA 6Gb/s, 25nm, MLC).

**Performance measured with**: Environment variables: KMP_AFFINITY='granularity=fine, compact', OMP_NUM_THREADS=56, CPU Freq set with cpupower frequency-set -d 2.5G -u 3.8G -g performance

**Deep Learning Frameworks:**

- Caffe: (http://github.com/intel/caffe/), revision f96b759f71b2281835f690af267158b82b150b5c. Inference measured with "caffe time --forward_only" command, training measured with "caffe time" command. For "ConvNet" topologies, dummy dataset was used. For other topologies, data was stored on local storage and cached in memory before training. Topology specs from https://github.com/intel/caffe/tree/master/models/intel_optimized_models (GoogLeNet, AlexNet, and ResNet-50), https://github.com/intel/caffe/tree/master/models/default_vgg_19 (VGG-19), and https://github.com/soumith/convnet-benchmarks/tree/master/caffe/imagenet_winners (ConvNet benchmarks; files were updated to use newer Caffe prototxt format but are functionally equivalent). Intel C++ compiler ver. 17.0.2 20170213, Intel MKL small libraries version 2018.0.20170425. Caffe run with "numactl -l".

- TensorFlow: (https://github.com/tensorflow/tensorflow), commit id 207203253b6f8ea5e938a512798429f91d5b4e7e. Performance numbers were obtained for three convnet benchmarks: alexnet, googlenetv1, vgg(https://github.com/soumith/convnet-benchmarks/tree/master/tensorflow) using dummy data. GCC 4.8.5, Intel MKL small libraries version 2018.0.20170425, interop parallelism threads set to 1 for alexnet, vgg benchmarks, 2 for googlenet benchmarks, intra op parallelism threads set to 56, data format used is NCHW, KMP_BLOCKTIME set to 1 for googlenet and vgg benchmarks, 30 for the alexnet benchmark. Inference measured with --caffe time -forward_only -engine MKL2017option, training measured with --forward_backward_only option.

- MxNet: (https://github.com/dmlc/mxnet/), revision 5efd91a71f36fea483e882b0358c8d46b5a7aa20. Dummy data was used. Inference was measured with "benchmark_score.py", training was measured with a modified version of benchmark_score.py which also runs backward propagation. Topology specs from https://github.com/dmlc/mxnet/tree/master/example/image-classification/symbols. GCC 4.8.5, Intel MKL small libraries version 2018.0.20170425.

- Neon: ZP/MKL_CHWN branch commit id:52bd02acb947a2adabb8a227166a7da5d9123b6d. Dummy data was used. The main.py script was used for benchmarking , in mkl mode. ICC version used : 17.0.3 20170404, Intel MKL small libraries version 2018.0.20170425.

# Configuration Details

Platform: 2S Intel® Xeon® CPU E5-2699 v4 @ 2.20GHz (22 cores), HT enabled, turbo disabled, scaling governor set to "performance" via acpi-cpufreq driver, 256GB DDR4-2133 ECC RAM. CentOS Linux release 7.3.1611 (Core), Linux kernel 3.10.0-514.10.2.el7.x86_64. SSD: Intel® SSD DC S3500 Series (480GB, 2.5in SATA 6Gb/s, 20nm, MLC).

**Performance measured with**: Environment variables: KMP_AFFINITY='granularity=fine, compact,1,0', OMP_NUM_THREADS=44, CPU Freq set with cpupower frequency-set –d 2.2G –u 2.2G –g performance

**Deep Learning Frameworks:**

- Caffe: (http://github.com/intel/caffe/), revision f96b759f71b2281835f690af267158b82b150b5c. Inference measured with "caffe time --forward_only" command, training measured with "caffe time" command. For "ConvNet" topologies, dummy dataset was used. For other topologies, data was stored on local storage and cached in memory before training. Topology specs from https://github.com/intel/caffe/tree/master/models/intel_optimized_models (GoogLeNet, AlexNet, and ResNet-50), https://github.com/intel/caffe/tree/master/models/default_vgg_19 (VGG-19), and https://github.com/soumith/convnet-benchmarks/tree/master/caffe/imagenet_winners (ConvNet benchmarks; files were updated to use newer Caffe prototxt format but are functionally equivalent). GCC 4.8.5, Intel MKL small libraries version 2017.0.2.20170110.

- TensorFlow: (https://github.com/tensorflow/tensorflow), commit id 207203253b6f8ea5e938a512798429f91d5b4e7e. Performance numbers were obtained for three convnet benchmarks: alexnet, googlenetv1, vgg(https://github.com/soumith/convnet-benchmarks/tree/master/tensorflow) using dummy data. GCC 4.8.5, Intel MKL small libraries version 2018.0.20170425, interop parallelism threads set to 1 for alexnet, vgg benchmarks, 2 for googlenet benchmarks, intra op parallelism threads set to 44, data format used is NCHW, KMP_BLOCKTIME set to 1 for googlenet and vgg benchmarks, 30 for the alexnet benchmark. Inference measured with --caffe time -forward_only -engine MKL2017option, training measured with --forward_backward_only option.

- MxNet: (https://github.com/dmlc/mxnet/), revision e9f281a27584cdb78db8ce6b66e648b3dbc10d37. Dummy data was used. Inference was measured with "benchmark_score.py", training was measured with a modified version of benchmark_score.py which also runs backward propagation. Topology specs from https://github.com/dmlc/mxnet/tree/master/example/image-classification/symbols. GCC 4.8.5, Intel MKL small libraries version 2017.0.2.20170110.

- Neon: ZP/MKL_CHWN branch commit id:52bd02acb947a2adabb8a227166a7da5d9123b6d. Dummy data was used. The main.py script was used for benchmarking , in mkl mode. ICC version used : 17.0.3 20170404, Intel MKL small libraries version 2018.0.20170425.

# Configuration Details

Platform: 2S Intel® Xeon® CPU E5-2697 v2 @ 2.70GHz (12 cores), HT enabled, turbo enabled, scaling governor set to "performance" via intel_pstate driver, 256GB DDR3-1600 ECC RAM. CentOS Linux release 7.3.1611 (Core), Linux kernel 3.10.0-514.21.1.el7.x86_64. SSD: Intel® SSD 520 Series 240GB, 2.5in SATA 6Gb/s, 25nm, MLC.

**Performance measured with**: Environment variables: KMP_AFFINITY='granularity=fine, compact,1,0', OMP_NUM_THREADS=24, CPU Freq set with cpupower frequency-set –d 2.7G –u 3.5G –g performance

**Deep Learning Frameworks:**

– Caffe: (http://github.com/intel/caffe/), revision b0ef3236528a2c7d2988f249d347d5fdae831236. Inference measured with "caffe time --forward_only" command, training measured with "caffe time" command. For "ConvNet" topologies, dummy dataset was used. For other topologies, data was stored on local storage and cached in memory before training. Topology specs from https://github.com/intel/caffe/tree/master/models/intel_optimized_models (GoogLeNet, AlexNet, and ResNet-50), https://github.com/intel/caffe/tree/master/models/default_vgg_19 (VGG-19), and https://github.com/soumith/convnet-benchmarks/tree/master/caffe/imagenet_winners (ConvNet benchmarks; files were updated to use newer Caffe prototxt format but are functionally equivalent). GCC 4.8.5, Intel MKL small libraries version 2017.0.2.20170110.