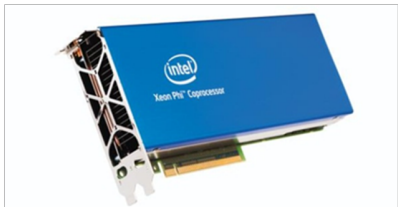# Performance Optimization of Deep Learning Frameworks on Modern Intel Architectures
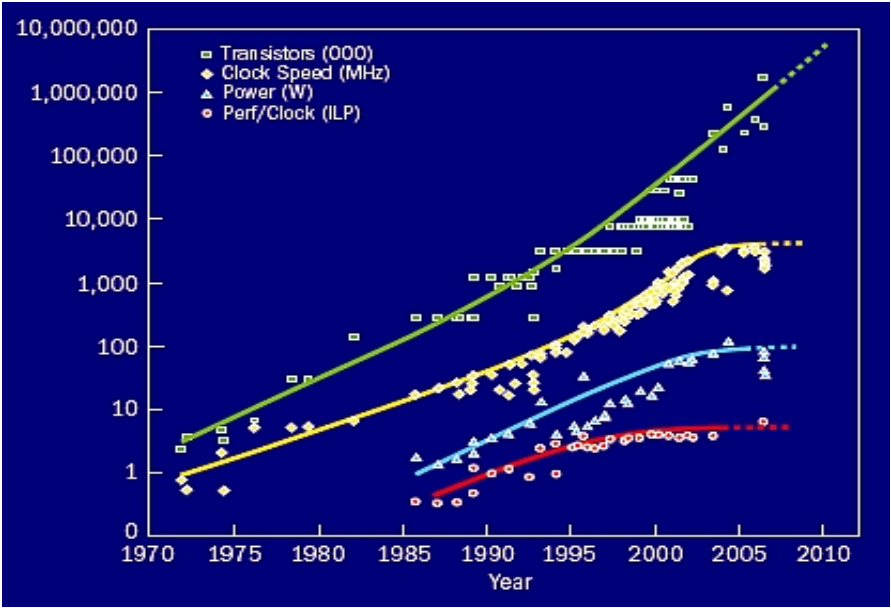
**ElMoustapha Ould-Ahmed-Vall, AG Ramesh, Vamsi Sripathi and Karthik Raman**

**Representing the work of many at Intel**

# Agenda

- Optimization matters on modern architectures

- Intel's recent Xeon and Xeon Phi products

- Introduction to Deep Learning

- Optimizing DL frameworks on IA
  - Key challenges
  - Optimization techniques
  - Performance data
  - DL scaling

# Moore's Law Goes on!



*Increasing clock speeds -> more cores + wider SIMD (Hierarchical parallelism)*

# Combined Amdahl's Law for Vector Multicores*

$$Speedup = (1/Serial_{frac} + 1 - Serial_{frac}/\textbf{NumCores}) * (1/Scalar_{frac} + 1 - Scalar_{frac}/\textbf{VectorLength})$$

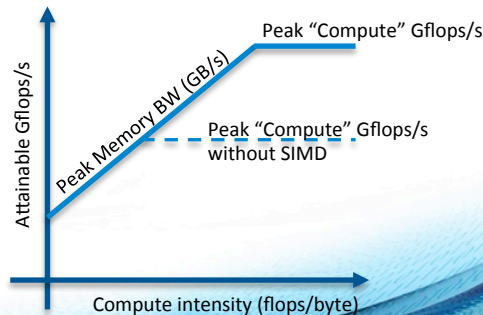Goal: *Reduce* **Serial Fraction** and *Reduce* **Scalar Fraction** of Code

Ideal Speedup: **NumCores*VectorLength  (requires zero scalar, zero serial work)**

**Compute Bound Performance**
Most kernels of ML codes are compute bound
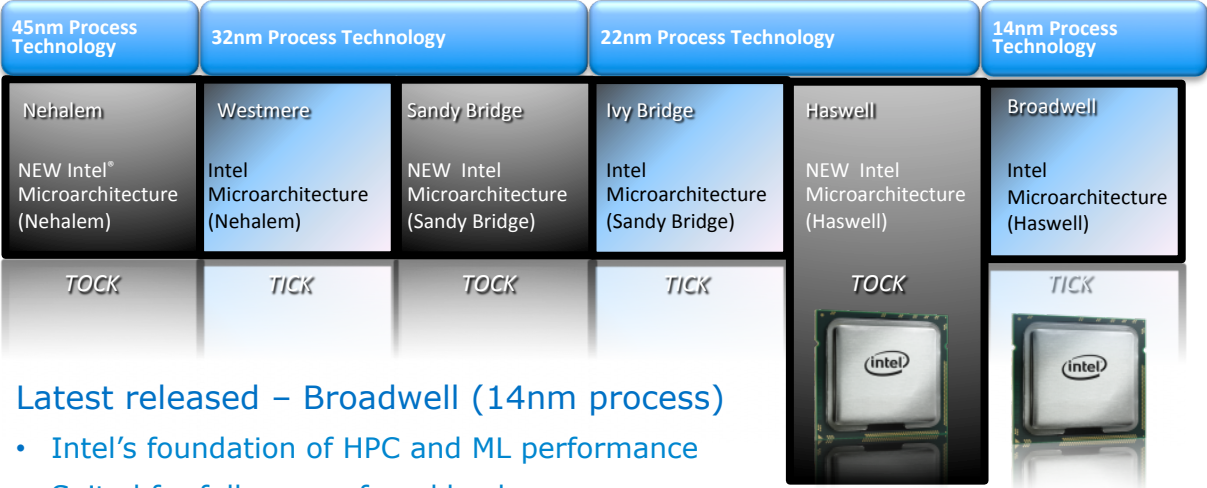i.e. raw FLOPS matter

Roofline Model
Gflops/s = min (Peak Gflops/s, Stream BW * flops/byte)

# Overview of Current Generation of Intel Xeon and Xeon Phi Products

# Current Intel® Xeon Platforms

| 45nm Process Technology | 32nm Process Technology | | 22nm Process Technology | | 14nm Process Technology |
|---|---|---|---|---|---|
| Nehalem<br><br>NEW Intel® Microarchitecture (Nehalem) | Westmere<br><br>Intel Microarchitecture (Nehalem) | Sandy Bridge<br><br>NEW Intel Microarchitecture (Sandy Bridge) | Ivy Bridge<br><br>Intel Microarchitecture (Sandy Bridge) | Haswell<br><br>NEW Intel Microarchitecture (Haswell) | Broadwell<br><br>Intel Microarchitecture (Haswell) |
| *TOCK* | *TICK* | *TOCK* | *TICK* | *TOCK* | *TICK* |

## Latest released – Broadwell (14nm process)

- Intel's foundation of HPC and ML performance
- Suited for full scope of workloads
- Industry leading performance/watt for serial & highly parallel workloads.
- Upto 22 cores / socket (Broadwell-EP) (w/ Hyper-Threading technology)

**Software optimization helps maximize benefit and adoption of new features**
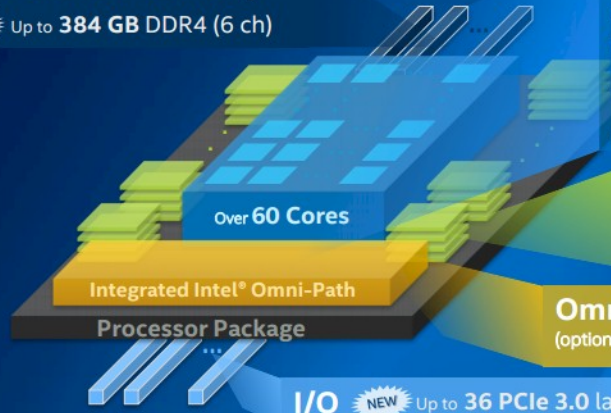
# 2<sup>nd</sup> Generation Intel® Xeon Phi™ Platform



## Knights Landing
*Holistic Approach to Real Application Breakthroughs*

(intel) inside
**XEON PHI**

**Platform Memory**
NEW Up to **384 GB** DDR4 (6 ch)

**Over 60 Cores**

Integrated Intel® Omni-Path

**Processor Package**

**Compute**
- Intel® Xeon® Processor Binary-Compatible
- **3+ TF**LOPS[1], **3X ST**[2] (single-thread) perf. vs KNC
- **2D Mesh** Architecture
- **Out-of-Order** Cores

**On-Package Memory**
- Over **5x** STREAM vs. DDR4[3]
- Up to **16 GB** at launch

**Omni-Path** (optional)  ▪ **1<sup>st</sup>** Intel processor to integrate

**I/O** NEW Up to **36 PCIe 3.0** lanes

# Intel® AVX Technology

| SNB/IVB | HSW/BDW | SKX & KNL |
|---|---|---|
| 256b AVX1<br>Flops/Cycle: 16 SP / 8 DP | 256b AVX2<br>Flops/Cycle: 32SP / 16 DP (FMA) | 512b AVX512<br>Flops/Cycle: 64SP / 32 DP (FMA) |

| AVX | AVX2 |
|---|---|
| 256-bit basic FP<br>16 registers<br>NDS (and AVX128)<br>Improved blend<br>MASKMOV<br>Implicit unaligned | Float16 (IVB 2012)<br>256-bit FP FMA<br>256-bit integer<br>PERMD<br>Gather |

| AVX512 |
|---|
| 512-bit FP/Integer<br>32 registers<br>8 mask registers<br>Embedded rounding<br>Embedded broadcast<br>Scalar/SSE/AVX "promotions"<br>Native media additions<br>HPC additions<br>Transcendental support<br>Gather/Scatter |

# Overview of Deep Learning and DL Frameworks

# Deep Learning – Convolutional Neural Network



Convolutional layer 1 — Convolutional layer 2 — Convolutional layer 3 — Convolutional layer 4 — Soft-max layer

Input layer
Max-pooling layer 1
Max-pooling layer 2
Max-pooling layer 3
Deep hidden identity features (DeepID)

Convolution Parameters:
Number of outputs/feature-maps: < 4 >
Filter size: < 3 x 3 >
Stride: < 2 >
Pad_size (for corner case): <1>

Feature maps

Filter = 3 x 3
Stride = 2
Pad_size = 1

Image

Convolved Feature

# Deep Learning: Train Once Use Many Times

# Deep Learning: Why Now?

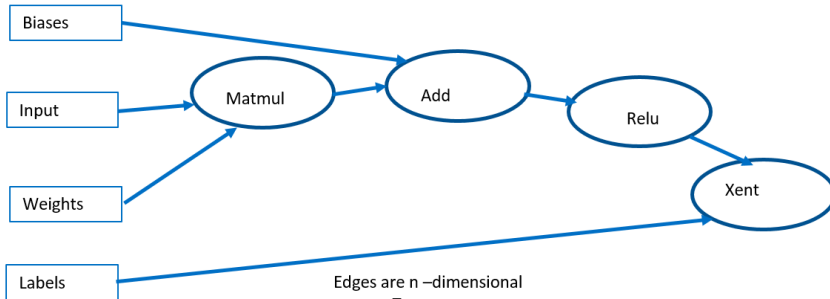| Bigger Data | Better Hardware | Smarter Algorithms |
|---|---|---|
| Image: 1000 KB / picture | Transistor density doubles every 18 months | Advances in algorithm innovation, including neural networks, leading to better accuracy in training models |
| Audio: 5000 KB / song | Cost / GB in 1995: $1000.00 | |
| Video: 5,000,000 KB / movie | Cost / GB in 2015: $0.03 | |

# Intel Caffe – ML Framework
## Optimized for Xeon and Xeon Phi Products

- ❑ Fork of BVLC Caffe by Intel to optimize for IA
- ❑ Leverages Intel MKL Deep Neural Network (DNN) API's
- ❑ Optimized for BDW (AVX2) and KNL (MIC_AVX512)
- ❑ https://github.com/intel/caffe

# Tensorflow ™ : Open Source ML Framework (Google)

- **Computation is a Dataflow Graph with Tensors**
- General computing mathematical framework – widely used for
  - Deep Neural Networks
  - Other machine learning algorithms
  - HPC applications
- Key computational kernels, extendable user operations
- Core in C++, front end wrapper in python
- Multi node support using GRPC
  - Google Remote Procedural Calls



Edges are n –dimensional arrays : Tensors

Example from Jeff Dean's presentation

# Optimizing Deep Learning Frameworks

# Performance Optimization on Modern Platforms

## Hierarchical Parallelism

### Coarse-Grained / multi-node
**Domain decomposition**

### Fine-Grained Parallelism / within node
**Sub-domain: 1) Multi-level domain decomposition (ex. across layers)**
**2) Data decomposition (layer parallelism)**

### Scaling

Improve load balancing

Reduce synchronization events, all-to-all comms

### Utilize all the cores

OpenMP, MPI, TBB…

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

# Intel Strategy: Optimized Deep Learning Environment



| | |
|---|---|
| Fuel the development of vertical solutions |
| Intel® Deep Learning SDK | Accelerate design, training, and deployment |
| Caffe · theano · n · TensorFlow · torch · CNTK Microsoft · dmlc mxnet | Drive optimizations across open source deep learning frameworks |
| Intel® Math Kernel Library (Intel® MKL) · Intel® MKL-DNN | Maximum performance on Intel architecture |
| intel XEON PHI inside + Intel® Omni-Path Architecture (Intel® OPA) · intel XEON inside + ALTERA Arria 10 FPGA+SoC | Deliver best single node and multi-node performance |

**Training**          **Inference**

# Example Challenge 1: Data Layout Has Big Impact on Performance

- Data Layouts impacts performance
  - Sequential access to avoid gather/scatter
  - Have iterations in inner most loop to ensure high vector utilization
  - Maximize data reuse; e.g. weights in a convolution layer
- Converting to/from optimized Layout is some times less expensive than operating on unoptimized Layout

| 21 | 18 | 32 | 6 | 3 | |
|----|----|----|----|----|----|
| 1 | 8 | 92 | 37 | 29 | 44 |
| 40 | 11 | 9 | 22 | 3 | 26 |
| 23 | 3 | 47 | 29 | 88 | 1 |
| 5 | 15 | 16 | 22 | 46 | 12 |
| | 29 | 9 | 13 | 11 | 1 |

| 21 | 18 | ... | 1 | .. | 8 | 92 | .. |
|----|----|----|----|----|----|----|----|

Better optimized for
some operations

vs

| 21 | 8 | 18 | 92 | .. | 1 | 11 | .. |
|----|----|----|----|----|----|----|----|

# Example Challenge 2: Minimize Conversions Overhead

- End to end optimization can reduce conversions
- Staying in optimized layout as long as possible becomes one of the tuning goals
- Minimize the number of back and forth conversions
  - Use of graph optimization techniques



Native to MKL layout    Convolution    MKL layout to Native    Max Pool    Native to MKL layout    Convolution    MKL layout to Native

# Example Challenge 3: Ensuring Enough Parallelism to Leverage all Cores

- Maximize parallelism to use all cores efficiently

- Intra operation/layer parallelism within operators (OpenMP)



Convolution of tiles in parallel

Inter operation parallelism across operators

Parallel execution

# Example Challenge 4: Optimizing the Data Layer



- Data Layer comprises 3 major ops
  - Read data
  - Decode data: e.g. JPEG decode, decompression
  - Transform data
- Result of read, decode & transform is input to DNN layers
- Reduce number of cores dedicated to feed DNN
  - IO optimization: consider compression
  - Decode: consider LMDB instead of JPEG
  - Resizing/data processing: consider pre-processing
  - Then vectorize, parallelize

# Optimizing Deep Learning Frameworks for Intel® Architecture

- Leverage high performant compute libraries and tools
  - e.g. Intel® Math Kernel Library, Intel® Python, Intel® Compiler etc.
- Data Format/Shape:
  - Right format/shape for max performance: blocking, gather/scatter
- Data Layout:
  - Minimize cost of data layout conversions
- Parallelism:
  - Use all cores, eliminate serial sections, load imbalance
- Other Functions/Primitives (un-optimized in libraries):
  - Optimize via compiler knobs, improve existing implementations
- Memory allocation
  - unique characteristics and ability to reuse buffers
- Data layer optimizations:
  - parallelization, vectorization, IO
- Optimize hyper parameters:
  - e.g. batch size for more parallelism
  - learning rate and optimizer to ensure accuracy/convergence

# AlexNet Optimization Progression

# VGG Optimization Progression



Chart: Cumulative Speedup across optimization stages for Broadwell and Knights Landing

| Stage | Broadwell | Knights Landing |
|---|---|---|
| Baseline | 1.00x | 1.00x |
| MKL Integration | 3.15x | 15.80x |
| Thread Optimization | 5.40x | 23.60x |
| Compiler Knobs Tuning | 10.18x | 122.50x |
| Matrix Transpose/Data Transformations | 13.29x | 164.95x |
| Memory Allocations | 14.65x | 171.20x |
| Conversions Optimization | 19.27x | 273.50x |

■ Broadwell  ■ Knights Landing

# Configuration details

Intel® Xeon™ processor E5-2699v4 (22 Cores, 2.2 GHz), 128GB DDR memory, Centos 7.2 based on Red Hat* Enterprise Linux 7.2

Intel® Xeon Phi™ processor 7250 (68 Cores, 1.4 GHz, 16GB MCDRAM: Flat mode), 96GB DDR memory, Centos 7.2 based on Red Hat* Enterprise Linux 7.2

AlexNet and VGG benchmarks:

https://github.com/soumith/convnet-benchmarks

# Multi-Node Distributed Training

- Model Parallelism
  - Break the model into *N* nodes
  - The same data is in all the nodes


- Data Parallelism
  - Break the dataset into *N*  nodes
  - The same model is in all the nodes
  - Good for networks with few weights, e.g. GoogLeNet


- You can use either model or data parallelism or a hybrid of both

# Data Parallelism

# Scaling Efficiency: Intel® Xeon Phi™ Processor

## Deep Learning Image Classification <u>Training</u> Performance : MULTI-NODE Scaling

# Multi-node Challenges

- Need to optimize both compute (iteration) and communication (weight updates)
- More nodes mean higher batch per iteration
  - Enough work for each node
- Optimized hyper parameters (e.g. Batch Size)
  - Time to Train: increases with batch size
  - Accuracy: batch size impacts convergence and accuracy
- Communication overheads if small per node batch
  - e.g. Total batch size = 1024
    - 1024 nodes : Batch size = 1 per node – **communication dominates**
    - 64 nodes each : Batch size = 16 per node – **computation dominates**



$Time\ To\ Train\ (TTT)$

sweet spot

$batch\ size$

# Summary

- Don't be fooled by performance of DL workloads when using unoptimized frameworks

- Significant performance headroom from optimization on Xeon and Xeon Phi
  - Close to 300x speedup in certain topologies

- Traditional vectorization and parallelization strategies apply

- Other unique performance challenges: hyper parameters, data layer, inter/intra layer parallelization, etc.

- Call to action:
  - Try Intel optimized frameworks available today, more to come soon

# Legal Disclaimers

# Optimization Notice