# GridKa School 2019

KIT, Campus North, FTU

# INTEL'S HARDWARE & SOFTWARE SOLUTIONS - DIRECTIONS FOR ARTIFICIAL INTELLIGENCE

Edmund Preiss

Business Development Manager, EMEA
Intel Computing Performance and Software Products (CPDP)

# Agenda

**Intel Xeon Family – Latest Status**

**Intel Software Development Tools**

- Intel Parallel Studio XE (IPS XE) Tool Suites

- News on the IPS XE 2020 Edition

- Upcoming **oneAPI** development tools concept

**Intel optimized AI Solutions and Directions**

# Agenda

**Intel Xeon Family Update**

**Intel Software Development Tools**

- Intel Parallel Studio XE (IPS XE) Tool Suites

- News on the IPS XE 2020 Edition

- Upcoming **oneAPI** development tools concept

**Intel optimized AI Solutions and Directions**

# Code Modernization

Stage 1: Use Optimized Libraries

Stage 2: Compile with Architecture-specific Optimizations

Stage 3: Analysis and Tuning

Stage 4: Check Correctness

# What's Inside Intel® Parallel Studio XE

## Comprehensive Software Development Tool Suite

### COMPOSER EDITION

**BUILD**
Compilers & Libraries

C / C++ Compiler
Optimizing Compiler

Intel® Math Kernel Library

Fortran Compiler
Optimizing Compiler

Intel® Integrated
Performance Primitives
Image, Signal & Data Processing

Intel® Threading
Building Blocks
C++ Threading Library

Intel® Data Analytics
Acceleration Library

Intel® Distribution for Python*
High Performance Scripting

### PROFESSIONAL EDITION

**ANALYZE**
Analysis Tools

Intel® VTune™ Amplifier
Performance Profiler

Intel® Inspector
Memory & Thread Debugger

Intel® Advisor
Vectorization Optimization
& Thread Prototyping

### CLUSTER EDITION

**SCALE**
Cluster Tools

Intel® MPI Library
Message Passing Interface Library

Intel® Trace Analyzer & Collector
MPI Tuning & Analysis

Intel® Cluster Checker
Cluster Diagnostic Expert System

Intel® Architecture Platforms

intel CORE inside · intel XEON inside · intel XEON PHI inside

Operating System: Windows*, Linux*, MacOS[1]*

**More Power for Your Code –** software.intel.com/intel-parallel-studio-xe

# SELECTED INTEL PARALLEL STUDIO XE HIGHLIGHTS
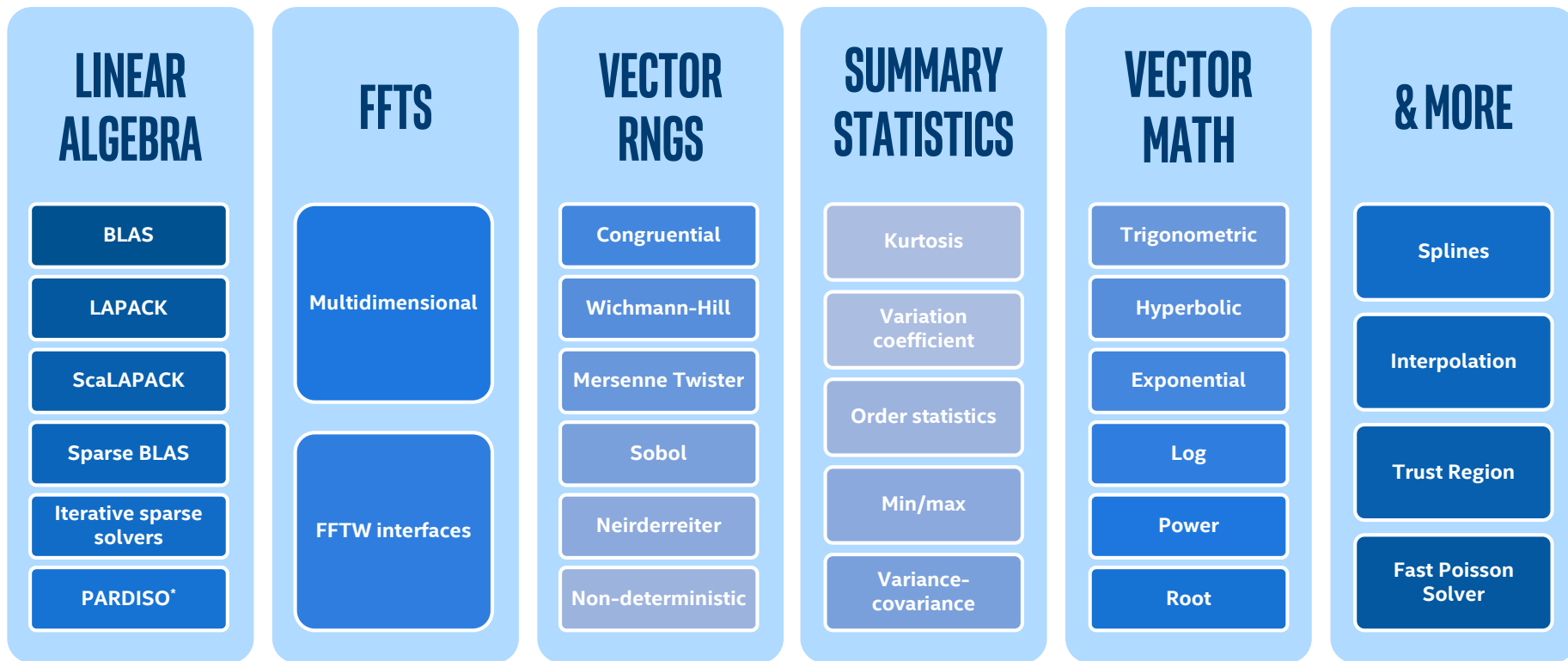
# What's Inside Intel® Math Kernel Library

| LINEAR ALGEBRA | FFTS | VECTOR RNGS | SUMMARY STATISTICS | VECTOR MATH | & MORE |
|---|---|---|---|---|---|
| BLAS | Multidimensional | Congruential | Kurtosis | Trigonometric | Splines |
| LAPACK | | Wichmann-Hill | Variation coefficient | Hyperbolic | Interpolation |
| ScaLAPACK | | Mersenne Twister | Order statistics | Exponential | |
| Sparse BLAS | FFTW interfaces | Sobol | Min/max | Log | Trust Region |
| Iterative sparse solvers | | Neiderreiter | Variance-covariance | Power | |
| PARDISO* | | Non-deterministic | | Root | Fast Poisson Solver |

¹Available only in Intel® Parallel Studio Composer Edition.

# Application Performance Snapshot (VTune Amplifier)

# Find Effective Optimization Strategies

Intel Advisor:  Cache-aware roofline analysis

## Roofs Show Platform Limits

Memory, cache & compute limits
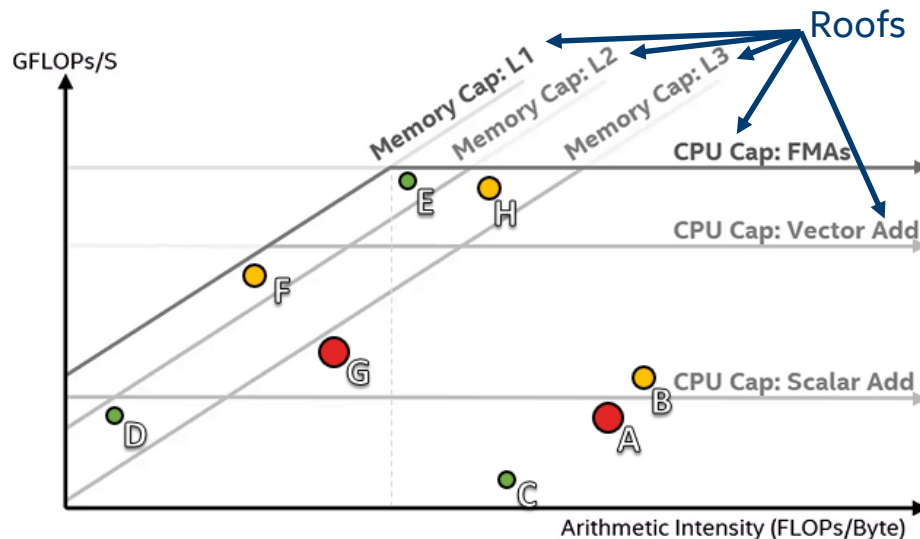
## Dots Are Loops

Bigger, red dots take more time so optimization has a bigger impact

Dots farther from a roof have more room for improvement

## Higher Dot = Higher GFLOPs/sec

Optimization moves dots up
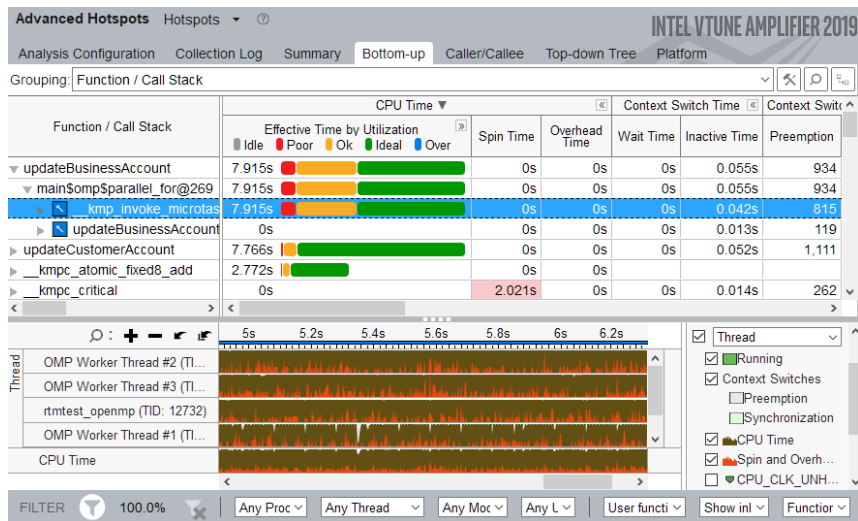
Algorithmic changes move dots horizontally



## Which loops should we optimize?

- A and G are the best candidates
- B has room to improve, but will have less impact
- E, C, D, and H are poor candidates

Roofline tutorial video

# Analyze & Tune Application Performance & Scalability with Intel® VTune™ Amplifier—Performance Profiler



## Fast, Scalable Code, Faster

- Accurately profile C, C++, Java*, Python*, Go*, or any mix
- Optimize CPU/GPU, threading, memory, cache, storage & more
- Save time: rich analysis leads to insight

## What's New in 2019 Release (Highlights)

- Simplified workflow for easier tuning
- I/O Analysis—Tune SPDK storage & DPDK network performance
- New Platform Profiler—Get insights into platform-level performance, identify memory & storage bottlenecks & imbalances

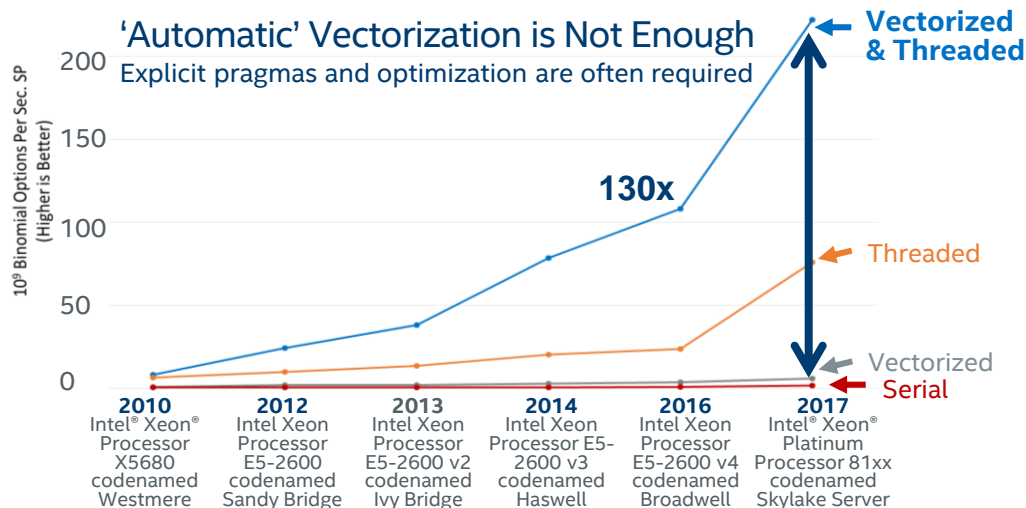OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos

# Optimize Vectorization & Threading with Intel® Advisor

**Performance Increases Scale with Each New Hardware Generation**

### 'Automatic' Vectorization is Not Enough
Explicit pragmas and optimization are often required



**Vectorized & Threaded**

**130x**

Threaded

Vectorized

Serial

Y-axis: $10^9$ Binomial Options Per Sec. SP (Higher is Better) — 0, 50, 100, 150, 200

| 2010 Intel® Xeon® Processor X5680 codenamed Westmere | 2012 Intel Xeon Processor E5-2600 codenamed Sandy Bridge | 2013 Intel Xeon Processor E5-2600 v2 codenamed Ivy Bridge | 2014 Intel Xeon Processor E5-2600 v3 codenamed Haswell | 2016 Intel Xeon Processor E5-2600 v4 codenamed Broadwell | 2017 Intel® Xeon® Platinum Processor 81xx codenamed Skylake Server |

## Modern Performant Code

- Vectorized for Intel® Advanced Vector Extensions (Intel® AVX-512 & Intel® AVX)
- Efficient memory access
- Threaded

## Capabilities

- Adds & optimizes vectorization
- Analyzes memory patterns
- Quickly prototypes threading

Benchmark: Binomial Options Pricing Model software.intel.com/en-us/articles/binomial-options-pricing-model-code-for-intel-xeon-phi-coprocessor

**Learn More: http: intel.ly/advisor**

12

# Get Breakthrough Vectorization Performance

Intel® Advisor—Vectorization Advisor

## Faster Vectorization Optimization

- Vectorize where it will pay off most
- Quickly ID what is blocking vectorization
- Tips for effective vectorization
- Safely force compiler vectorization
- Optimize memory stride

## Data & Guidance You Need

- Compiler diagnostics + Performance Data + SIMD efficiency
- Detect problems & recommend fixes
- Loop-Carried Dependency Analysis
- Memory Access Patterns Analysis



Optimize for Intel® Advanced Vector Extensions 512 (Intel® AVX-512) with or without access to Intel AVX-512 hardware

# Agenda

**Intel Xeon Family Update**

**Intel Software Development Tools**

- Intel Parallel Studio XE (IPS XE) Tool Suites
- News on the IPS XE 2020 Edition
- Upcoming **oneAPI** development tools concept

**Intel optimized AI Solutions and Directions**

# Key Updates for Intel® Parallel Studio XE 2020

**Speed Artificial Intelligence Inferencing**  -  Intel® Compiler and analyzer support for Vector Neural Instructions (VNNI) in Cascade Lake/AP platform

**512GB DIMMs with Persistent Memory** – Identify, Optimize & Tune Platforms for Intel® Optane™ Persistent Memory with Intel® VTune™ Amplifier

**Extended Coarse Grain Profiling** – Platform level collection and analysis in Intel® VTune™ Amplifier

**Cache Simulation Insights for Vectorization -** Roofline analysis for L1, L2 L3, DRAM in Intel® Advisor

**Expanded standard support** — More Fortran 2018 features & Expanded support of C++17 with initial C++20 support

**Latest Processor Support** - Intel® Xeon® Scalable Processors (codenamed Cascade Lake / Cascade Lake AP)

**New OS Support** – Clear Linux & Amazon Linux 2*

 * Supported features of tools and libraries may vary by instances and configurations

(intel)

# Intel® VTune™ Profiler Server Architecture

## Just launch your browser and go

Users         VTune Profiler         Target Systems
                   Server



## Easier profiling

- **Access with a web browser** – no install required by users
- **Share results** – all results available to all users with server access
- **Profile any system on the network** – server installs collector on the target

# What's Coming in Intel® Parallel Studio XE 2020

Coming Q3'2019

- **Intel® Math Kernel Library**
  - Increased AVX512 Optimizations for Complex Vector Math Functions
  - Strided Vector Math API
- **Intel® Data Analytics Acceleration Library**
  - Performance improvements and feature improvements such as
    - Gradient Boosted Trees for large dimensional data sets
    - Extended Z-score support for PCA algorithm
    - XGBoost accelerated with DAAL

# Agenda

**Intel Xeon Family Update**

**Intel Software Development Tools**

- Intel Parallel Studio XE (IPS XE) Tool Suites

- News on the IPS XE 2020 Edition

- Upcoming **oneAPI** development tools concept

**Intel optimized AI Solutions and Directions**

# DIVERSE WORKLOADS REQUIRE DIVERSE ARCHITECTURES

**The future** is a **diverse** mix of scalar, vector, matrix, and spatial **architectures** deployed in CPU, GPU, AI, FPGA and other accelerators



CPU    GPU    AI    FPGA

SCALAR    VECTOR    MATRIX    SPATIAL

SVMS

# PROGRAMMING CHALLENGE

Diverse set of data-centric hardware

No common programming language or APIs

Inconsistent tool support across platforms

Each platform requires unique software investment



CPU    GPU    AI    FPGA

SCALAR    VECTOR    MATRIX    SPATIAL

SVMS

# INTEL'S ONE API CORE CONCEPT

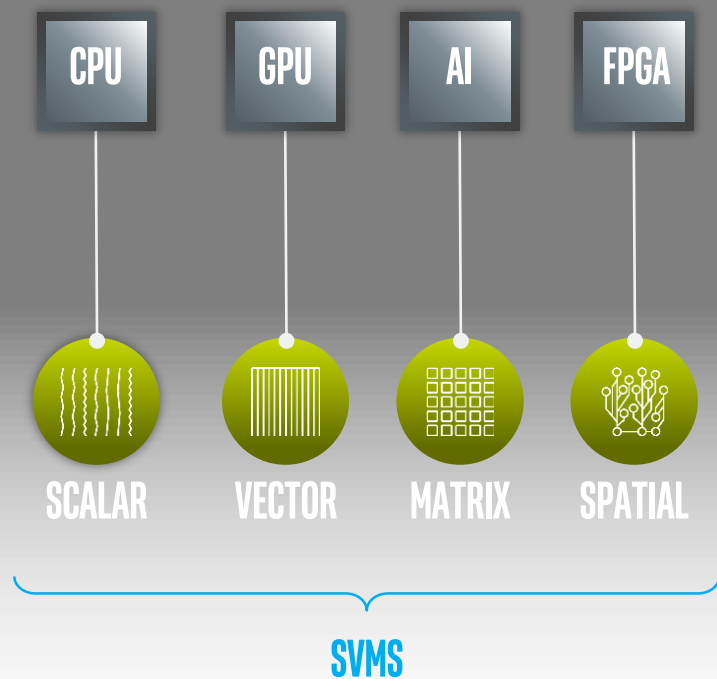**Project oneAPI** delivers a unified programming model to simplify development across diverse architectures

Common developer experience across Scalar, Vector, Matrix and Spatial (SVMS) architecture

Unified and simplified language and libraries for expressing parallelism

Uncompromised native high-level language performance

Support for CPU, GPU, AI and FPGA

Based on industry standards and open specifications

**One API Tools**

Optimized Applications

Optimized Middleware / Frameworks

**One API** Language & Libraries

| CPU | GPU | AI | FPGA |
| SCALAR | VECTOR | MATRIX | SPATIAL |

(intel)

# ONE API FOR CROSS-ARCHITECTURE PERFORMANCE

Optimized Applications

Optimized Middleware & Frameworks

**One API** Product

**Direct Programming**

Data Parallel C++

**API-Based Programming**

Math
Threading
DPC++ Library
Analytics/ML
DNN
ML Comm
Video Processing
Rendering

**Porting Tool**

**Analysis & Debug Tools**
VTune™
Advisor
Debugger

CPU

GPU

AI

FPGA

Some capabilities may differ per architecture.

# Agenda

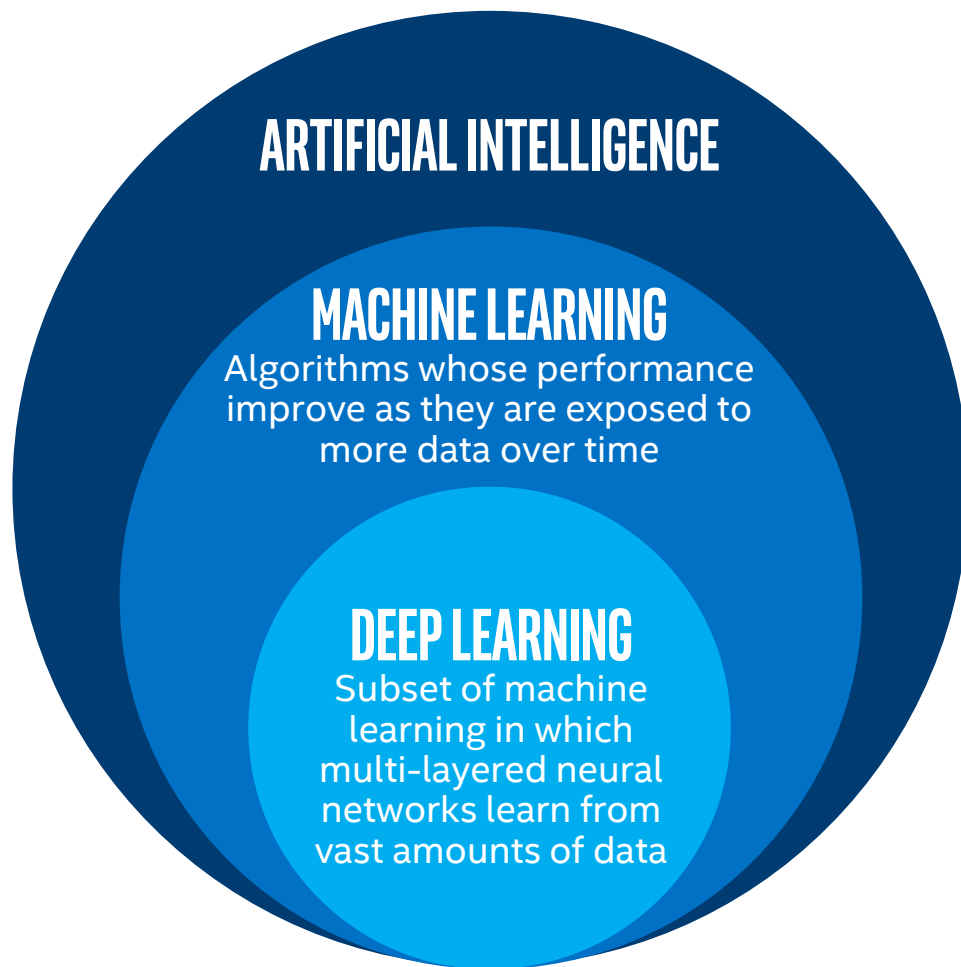**Intel Xeon Family Update**

**Intel Software Development Tools**

- Intel Parallel Studio XE (IPS XE) Tool Suites

- News on the IPS XE 2020 Edition

- Upcoming **oneAPI** development tools concept

**Intel optimized AI Solutions and Directions**

# ARTIFICIAL INTELLIGENCE

is the ability of machines to learn from experience, without explicit programming, in order to perform cognitive functions associated with the human mind

## ARTIFICIAL INTELLIGENCE

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

### DEEP LEARNING
Subset of machine learning in which multi-layered neural networks learn from vast amounts of data

(intel)

**ARTIFICIAL INTELLIGENCE**

A R T I F I C I A L   I N T E L L I G E N C E

## SOLUTIONS
*Solution Architects*

**ARTIFICIAL INTELLIGENCE**

AI Solutions Catalog (<u>Public</u> & <u>Internal</u>)

Platforms | Finance | Healthcare | Energy | Industrial | Transport | Retail | Home | More…

## TOOLKITS
*App Developers*

**DEEP LEARNING DEPLOYMENT**

**OpenVINO™** †
*Open Visual Inference & Neural Network Optimization toolkit for inference deployment on CPU, processor graphics, FPGA & VPU using TF, Caffe\* & MXNet\**

**Intel® Movidius™ SDK**
*Optimized inference deployment for all Intel® Movidius™ VPUs using TensorFlow\* & Caffe\**

**DEEP LEARNING** *COMING SOON !*
Intel® Deep Learning Studio‡
Open-source tool to compress deep learning development cycle

## LIBRARIES
*Data Scientists*

**MACHINE LEARNING LIBRARIES**

**Python**
• Scikit-learn
• Pandas
• NumPy

**R**
• Cart
• Random Forest
• e1071

**Distributed**
• MlLib (on Spark)
• Mahout

**DEEP LEARNING FRAMEWORKS** *COMING SOON !*

**Now optimized for CPU**

TensorFlow\* | MXNet\* | Caffe\* | BigDL/Spark\*

**Optimizations in progress**

Caffe2\* | PyTorch\* | PaddlePaddle\*

## FOUNDATION
*Library Developers*

**ANALYTICS, MACHINE & DEEP LEARNING PRIMITIVES**

**Python**
*Intel distribution optimized for machine learning*

**DAAL**
*Intel® Data Analytics Acceleration Library (for machine learning)*

**MKL-DNN   clDNN**
*Open-source deep neural network functions for CPU, processor graphics*

**DEEP LEARNING GRAPH COMPILER**

**Intel® nGraph™ Compiler** (Alpha)
*Open-sourced compiler for deep learning model computations optimized for multiple devices (CPU, GPU, NNP) using multiple frameworks (TF, MXNet, ONNX)*

## HARDWARE
*IT System Architects*

**AI FOUNDATION**

**DEEP LEARNING ACCELERATORS**

ATOM inside | Iris Graphics | CORE 8th Gen | Iris Graphics | XEON PLATINUM inside

Data Center
Edge
Device

NERVANA inside *COMING 2019* | STRATIX 10 inside | ARRIA 10 inside | MOVIDIUS inside | Mobileye | GNA inside

NNP L-1000

← Inference →

† *Formerly the Intel® Computer Vision SDK*
\**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.*

**AI.INTEL.COM**

# Components Comparison : Intel MKL-DNN vs Intel MKL

## MKL-DNN    (Open Source)

- Convolution
- Pooling

- ReLU
- Inner Product

- Normalization

## MKL (Math Kernel Library)

| Linear Algebra | Fast Fourier Transforms | Vector Math | Summary Statistics | And More... |
|---|---|---|---|---|
| • BLAS<br>• LAPACK<br>• ScaLAPACK<br>• Sparse BLAS<br>• Sparse Solvers<br>• Iterative<br>• PARDISO*<br>• Cluster Sparse Solver | • Multidimensional<br>• FFTW interfaces<br>• Cluster FFT | • Trigonometric<br>• Hyperbolic<br>• Exponential<br>• Log<br>• Power<br>• Root<br>• Vector RNGs | • Kurtosis<br>• Variation coefficient<br>• Order statistics<br>• Min/max<br>• Variance-covariance | • Splines<br>• Interpolation<br>• Trust Region<br>• Fast Poisson Solver |

# Data Transformation & Analysis Algorithms

Intel® Data Analytics Acceleration Library

| Basic Statistics for Datasets | Correlation & Dependence | Matrix Factorizations | Dimensionality Reduction | Outlier Detection |
|---|---|---|---|---|
| Low Order Moments | Cosine Distance | SVD | PCA | Univariate |
| Quantiles | Correlation Distance | QR | Association Rule Mining (Apriori) | Multivariate |
| Order Statistics | Variance-Covariance Matrix | Cholesky | Optimization Solvers (SGD, AdaGrad, lBFGS) | Math Functions (exp, log,…) |

☐ Algorithms supporting batch processing

☐ Algorithms supporting batch, online and/or distributed processing

# Machine Learning Algorithms

Intel® Data Analytics Acceleration Library



Regression
- Logistic Regression
- Ridge Regression
- Linear Regression

Supervised Learning
- Neural Networks

Classification
- Decision Forest
- Decision Tree
- Boosting (Ada, Brown, Logit)
- Naïve Bayes — Weak Learner
- k-NN
- Support Vector Machine

Unsupervised Learning
- K-Means Clustering
- EM for GMM

Collaborative Filtering — Alternating Least Squares

Algorithms supporting batch processing

Algorithms supporting batch, online and/or distributed processing

# Faster Python* with Intel® Distribution for Python

**Intel® Distribution for Python\* Performance Speedups
for Select Math Functions on Intel® Xeon™ Processors**

■ Speedup with Intel Python vs pip/numpy

Up to **440X speedup** versus stock
*NumPy from pip*

Speedup Factor (Y-axis: 0 to 500)

| Math function | Speedup |
|---|---|
| array-array | 16X |
| array-scalar | 17X |
| array*array | 16X |
| array*scalar | 17X |
| array+array | 16X |
| array+scalar | 17X |
| erf | 51X |
| exp | 258X |
| invsqrt | 77X |
| log10 | 442X |

Math functions (Array size = 1M)

**Configuration:** Hardware: Intel® Xeon® CPU E5-2699 v4 @ 2.20GHz (2 sockets, 22 cores per socket, 1 thread per core – HT is off), 256GB DDR4 @ 2400MHz.
Software: Stock: CentOS Linux release 7.3.1611 (Core), python 3.6.2, pip 9.0.1, numpy 1.13.1, scipy 0.19.1, scikit-learn 0.19.0. Intel® Distribution for Python* 2018 Gold: mkl 2018.0.0 intel_4, daal 2018.0.0.20170814, numpy 1.13.1 py36_intel_15, openmp 2018.0.0 intel_7, scipy 0.19.1 np113py36_intel_11, scikit-learn 0.18.2 np113py36_intel_3

Learn More: software.intel.com/distribution-for-python

# Intel-Optimized AI Frameworks

## Popular DL Frameworks are now optimized for CPU!

**CHOOSE YOUR FAVORITE FRAMEWORK**

TensorFlow™ *   Caffe *   mxnet *   BigDL FOR Apache Spark™ *   PYTORCH

See installation guides at  ai.intel.com/framework-optimizations/

*More under optimization:*   Caffe2 *   PYTORCH *   PaddlePaddle *

SEE ALSO: Machine Learning Libraries for Python (Scikit-learn, Pandas, NumPy), R (Cart, randomForest, e1071), Distributed (MlLib on Spark, Mahout)
*Limited availability today
Other names and brands may be claimed as the property of others.

(intel)

# INTEL® XEON® PROCESSOR PLATFORM PERFORMANCE

## Hardware plus optimized software
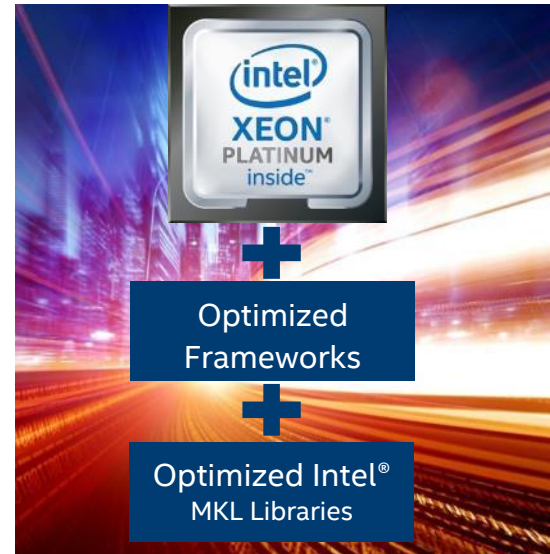
### INFERENCE THROUGHPUT

**198x**

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

Inference and training throughput uses FP32 instructions

### TRAINING THROUGHPUT

**127x**

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



Optimized
Frameworks

+

Optimized Intel®
MKL Libraries

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

Up to 191X Intel® Xeon® Platinum 8180 Processor higher Intel optimized Caffe Resnet50 with Intel® MKL inference throughput compared to Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe
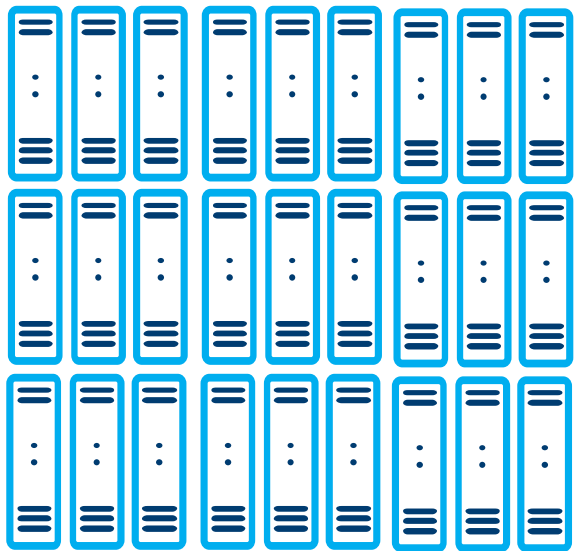Up to 93X Intel® Xeon® Platinum 8180 Processor higher Intel optimized Caffe Resnet50 with Intel® MKL training throughput compared to Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

# Use your Intel Architecture based HPC Infrastructure ---> Shorten Deep Learning Training Time

## HPC

## AI

# AI (ML & DL) SOFTWARE STACK FOR INTEL® PROCESSORS



## Distributed DL Training

Uber's open source Distributed training framework for TensorFlow

**Intel MKL**

**Intel MKL-DNN**

Intel Processors

**Intel MPI**

**Intel MLSL**

Intel Processors

# More details on Intel and AI at GridKa School

**Join the Intel hands-on Training :  Thursday afternoon at 13:15**

**Title : Enhance Machine Learning Performance with Intel® Software tools**

1st  Session:

- Introduction of Intel tools for ML and DL

  - DLBoost, VNNI instructions ; Intel distribution for Python

2nd session:

- Classical ML

  - Numpy and MKL, K-Means, clustering and DAAL4PY Distributed ML algorithms

3rd session:

- Intel MKL for Deep Neural Network Intel optimized Framework and Tensorflow

- Distributed Tensorflow with Horovod

# THANK YOU

# Legal Disclaimer & Optimization Notice

Performance results are based on testing as of August 2017 to September 2018 and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure.

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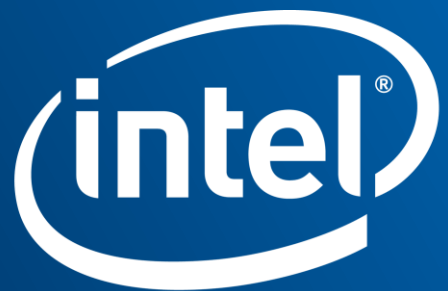
INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING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

**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.
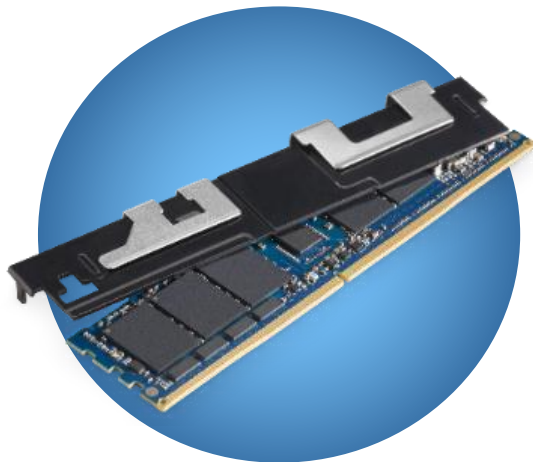Notice revision #20110804

# INTEL® OPTANE™ DC PERSISTENT MEMORY

# THE BEST OF BOTH WORLDS WITH INTEL® OPTANE™ DC PERSISTENT MEMORY

## DRAM ATTRIBUTES

Performance comparable to DRAM at *low latencies*[1]

## NAND SSD ATTRIBUTES

Data persistence with higher capacity than DRAM[2]

Optane™ Media

1. "*Fast performance comparable to DRAM*" – Intel persistent memory is expected to perform at latencies near DDR4 DRAM. Benchmarks and proof points forthcoming. "*low latencies*" – Data transferred across the memory bus causes latencies to be orders of magnitude lower when compared to transferring data across PCIe or I/O bus' to NAND/Hard Disk. Benchmarks and proof points forthcoming.
2. Intel persistent memory offers 3 different capacities – 128GB, 256GB, 512GB. Individual DIMMs of DDR4 DRAM max out at 256GB.

# Latency Estimates for different Memory and Storage Devices



- Volatile Memory
- Load/Store Instructions
- Cache Line Granularity

- Non-Volatile Storage
- Load/Store Instructions
- Cache Line Granularity

- Non-Volatile Storage
- I/O Commands
- Bock Granularity

Cost ($/GiB)

| | Latency |
|---|---|
| CPU Registers | ~0.1ns |
| CPU Caches (L1, L2, L3, L4) | 1-10ns |
| DDR DRAM | ~80-100ns |
| Persistent Memory | <1us* |
| NAND SSD | 10-100us |
| Hard Disk Drives (HDD) | ~10ms |
| Tape | ~100ms |

Capacity

# Distributed TensorFlow™ Compare

**Distributed Tensorflow with Parameter Server**

With Parameter Server →

Averages All the Gradients

Parameter Server

Worker A    Worker B    Worker C

or

Each Averages Portion of the Gradients

Parameter Server A    Parameter Server B    Parameter Server C
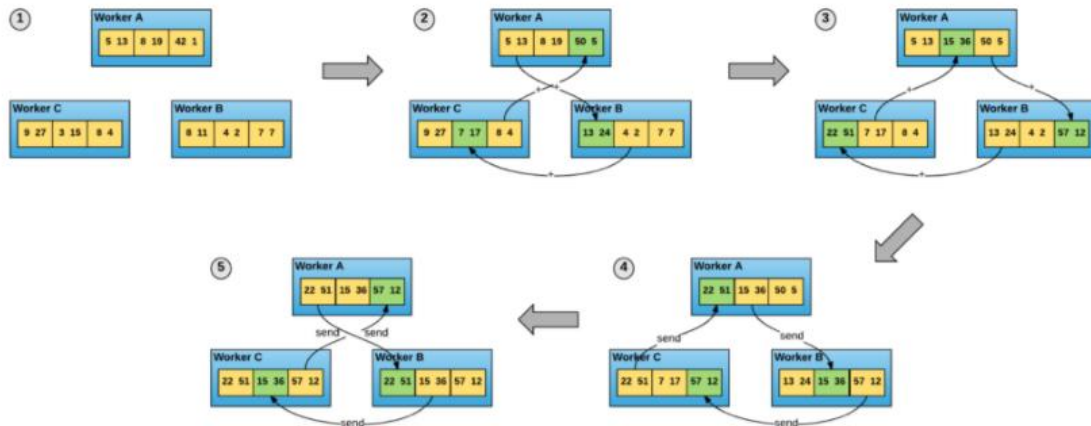
Worker A    Worker B    Worker C

The parameter server model for distributed training jobs can be configured with different ratios of parameter servers to workers, each with different performance profiles.

HOROVOD

No Parameter Server →

① Worker A — 5 13 8 19 42 1
Worker C — 9 27 3 15 8 4
Worker B — 8 11 4 2 7 7

② Worker A — 5 13 8 19 50 5
Worker C — 9 27 7 17 8 4
Worker B — 13 24 4 2 7 7

③ Worker A — 5 13 15 36 50 5
Worker C — 22 51 7 17 8 4
Worker B — 13 24 4 2 57 12

⑤ Worker A — 22 51 15 36 57 12
Worker C — 22 51 15 36 57 12    send    send
Worker B — 22 51 15 36 57 12

④ Worker A — 22 51 15 36 50 5
Worker C — 22 51 7 17 57 12    send    send
Worker B — 13 24 15 36 57 12    send

Uber's open source Distributed training framework for TensorFlow

The ring all-reduce algorithm allows worker nodes to average gradients and disperse them to all nodes without the need for a parameter server.

Source: https://eng.uber.com/horovod/

# UNSUPERVISED LEARNING EXAMPLE

**MACHINE LEARNING**

- Regression
- Classification
- Clustering
- Decision Trees
- Data Generation
- Image Processing

**DEEP LEARNING**

- Speech Processing
- Natural Language Processing
- Recommender Systems
- Adversarial Networks
- Reinforcement Learning

## MACHINE LEARNING (CLUSTERING)
*An 'Unsupervised Learning' Example*



K-MEANS

Revenue / Purchasing Power

Choose the right AI approach for your challenge

(intel)

# SUPERVISED LEARNING EXAMPLE

**MACHINE LEARNING**
- Regression
- Classification
- Clustering
- Decision Trees
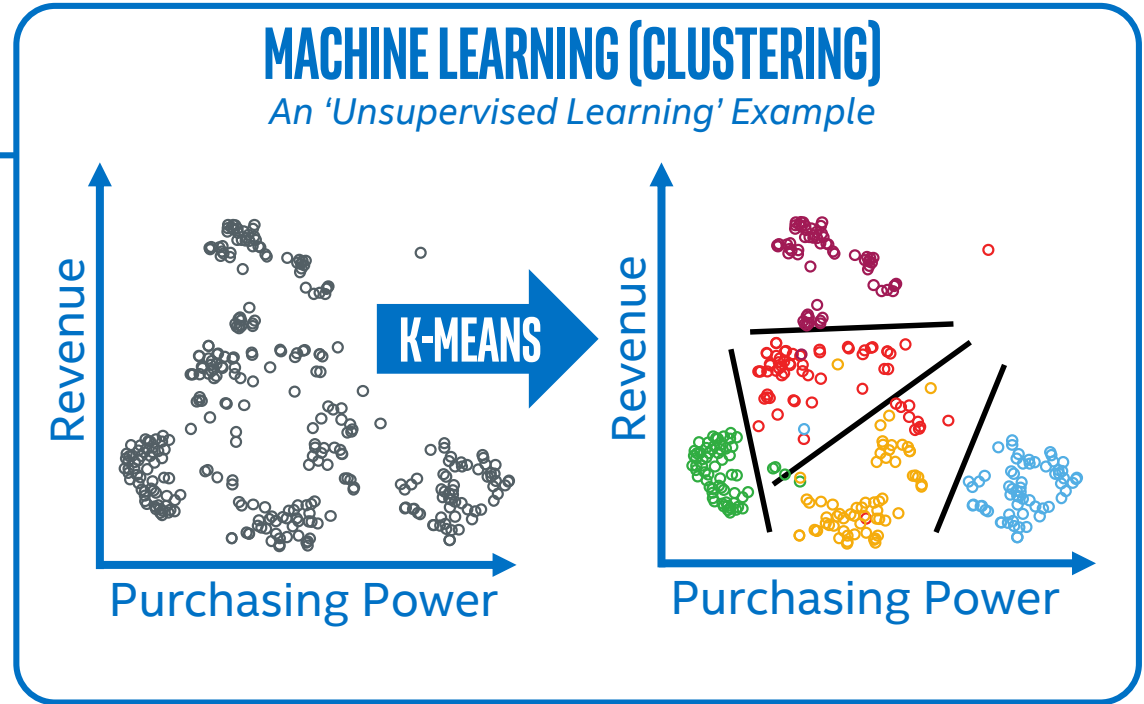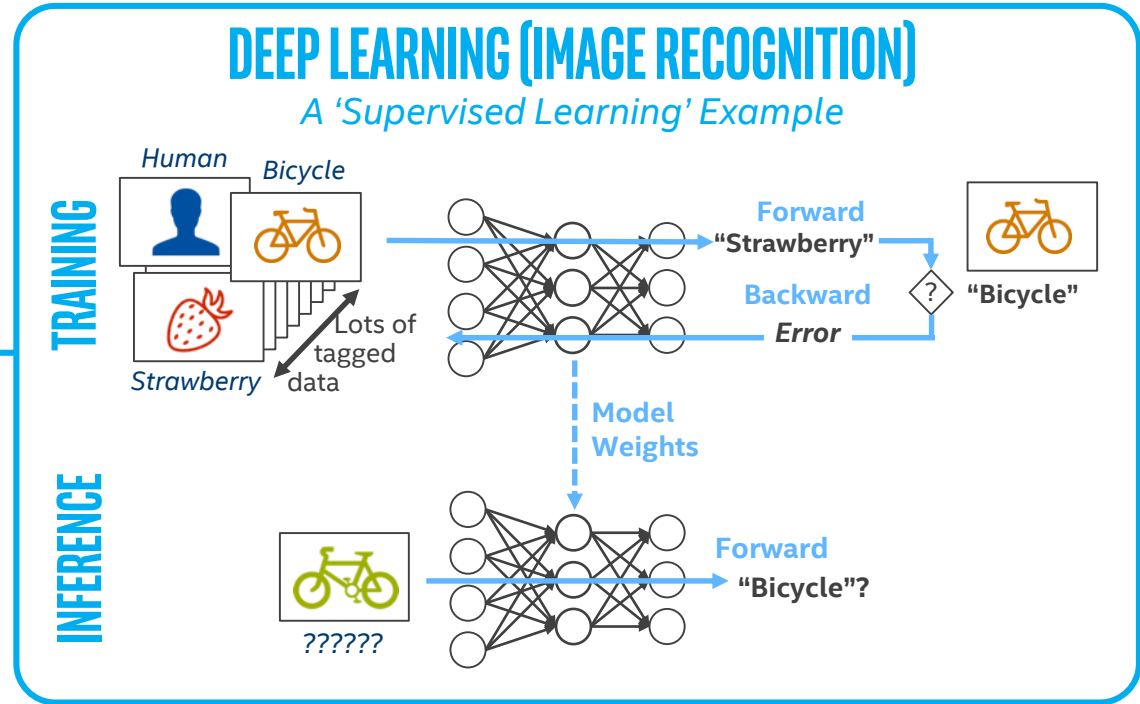- Data Generation
- Image Processing

**DEEP LEARNING**
- Speech Processing
- Natural Language Processing
- Recommender Systems
- Adversarial Networks
- Reinforcement Learning

## DEEP LEARNING (IMAGE RECOGNITION)
*A 'Supervised Learning' Example*

**TRAINING**

Human    Bicycle

Lots of tagged data

Strawberry

Forward "Strawberry"

Backward *Error*

? "Bicycle"

Model Weights

**INFERENCE**

?????? 

Forward "Bicycle"?

## Choose the right AI approach for your challenge

# REINFORCEMENT LEARNING EXAMPLE



Choose the right AI approach for your challenge

# Intel-Optimized Frameworks: How To Get?

## Check out our intel.ai for the framework optimizations page

### INTEL® OPTIMIZATION FOR TENSORFLOW*

This Python*-based deep learning framework is designed for ease of use and extensibility on modern deep neural networks and has been optimized for use on Intel® Xeon® processors.

→ Learn More   → GitHub
→ Documentation   → Resources

### MXNET*

The open-source, deep learning framework MXNet* includes built-in support for the Intel® Math Kernel Library (Intel® MKL) and optimizations for Intel® Advanced Vector Extensions 2 (Intel® AVX2) and Intel® Advanced Vector Extension 512 (Intel® AVX-512) instructions.

→ Learn More   → GitHub
→ Documentation   → Resources

### PYTORCH

Intel continues to accelerate and streamline PyTorch on Intel architecture, most notably Intel® Xeon® Scalable processors, both using Intel® Math Kernel Library for Deep Neural Networks (Intel® MKL-DNN) directly and making sure PyTorch is ready for our next generation of performance improvements both in software and hardware through the nGraph Compiler.

→ Learn More

### INTEL® OPTIMIZATION FOR CAFFE*

The Intel® Optimization for Caffe* provides improved performance for of the most popular frameworks when running on Intel® Xeon® processors.

→ Learn More   → GitHub
→ Documentation   → Resources

Installation example for optimized Tensorflow / Keras for a Medial Sample application
https://docs.google.com/document/d/1AAsWLSfYBx-pYfvzFQxOYi4Z1zx3TAEB5534s61Lf8w/edit?usp=sharing

Topology Examples : https://www.intel.ai/framework-optimizations