# THE AI ERA INFRASTRUCTURE

## INDUSTRY TRANSFORMATION & INNOVATION OPPORTUNITIES

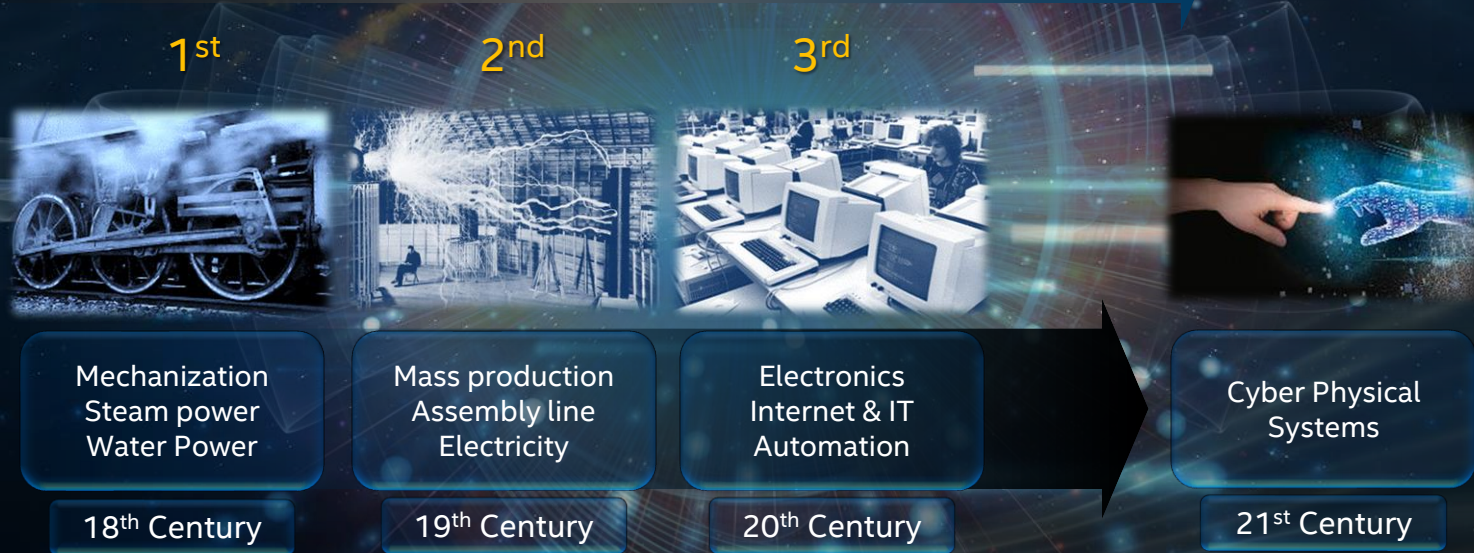**Meena Arunachalam**

Principal Engineer, Intel Architecture and Graphics Software

# INDUSTRIAL REVOLUTION -- FROM PHYSICAL TO DIGITAL

Fourth Industrial Revolution represents entirely new ways in which technology becomes embedded within industries, societies and even our human bodies

$4^{th}$

$1^{st}$      $2^{nd}$      $3^{rd}$



| Mechanization<br>Steam power<br>Water Power | Mass production<br>Assembly line<br>Electricity | Electronics<br>Internet & IT<br>Automation | Cyber Physical<br>Systems |
|---|---|---|---|
| $18^{th}$ Century | $19^{th}$ Century | $20^{th}$ Century | $21^{st}$ Century |

WE ARE ~~WITNESSING~~ THIS DIGITAL TRANSFORMATION

CREATING

# DATA DEFINES THE FUTURE



**COMPETITIVENESS AND BUSINESS GROWTH ARE INCREASINGLY DETERMINED BY THE POWER OF DATA.**

# INNOVATION ACROSS ALL INDUSTRIES

**AI / ML**

**Processing Power**

**"DIGITAL FUSION"**

Blending of Traditional & Digital Business Models

Smart Cities & Surveillance

Retail: Real-Time Pricing & Inventory

Enterprise/Consumer Analytics

Precision Medicine & Genomic Analytics

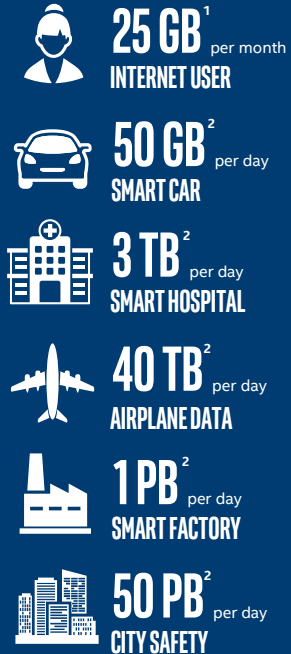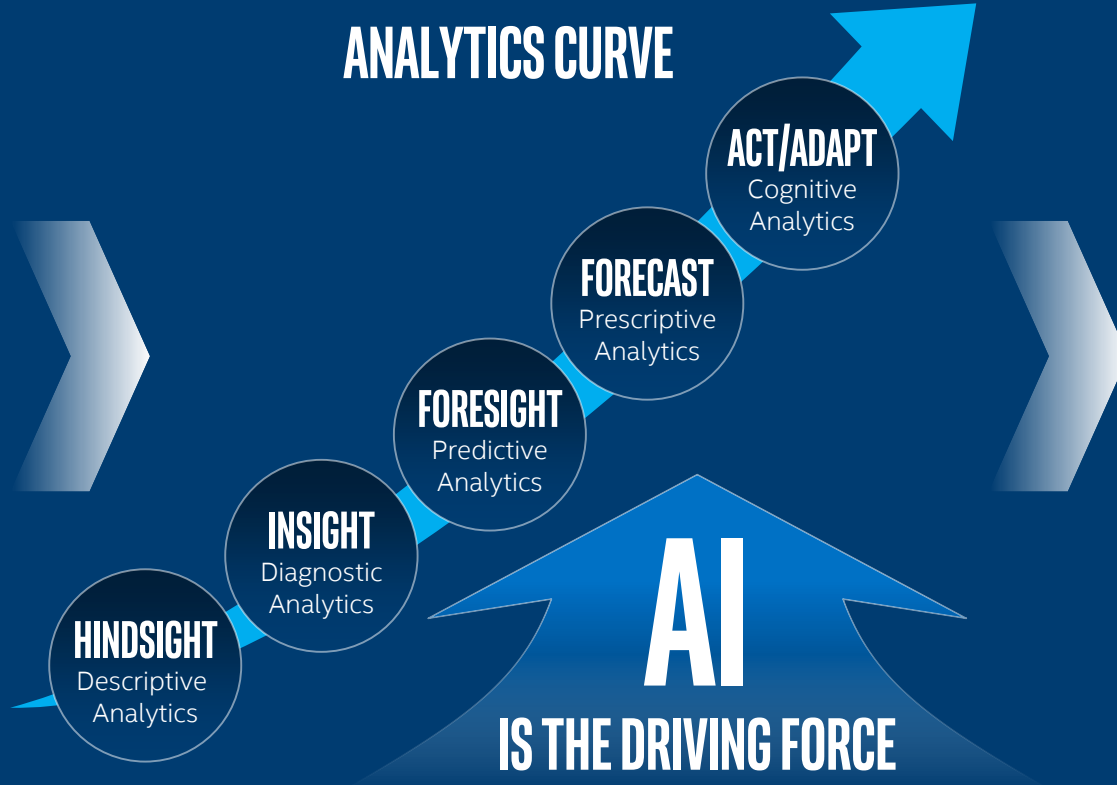Autonomous Cars & 5G Connectivity

Virtual/Augmented Reality

**Data**

# AI IS DRIVING ADVANCED ANALYTICS

## DATA DELUGE (2019)

25 GB[1] per month
**INTERNET USER**

50 GB[2] per day
**SMART CAR**

3 TB[2] per day
**SMART HOSPITAL**

40 TB[2] per day
**AIRPLANE DATA**

1 PB[2] per day
**SMART FACTORY**

50 PB[2] per day
**CITY SAFETY**

## ANALYTICS CURVE

**ACT/ADAPT**
Cognitive
Analytics

**FORECAST**
Prescriptive
Analytics

**FORESIGHT**
Predictive
Analytics

**INSIGHT**
Diagnostic
Analytics

**HINDSIGHT**
Descriptive
Analytics

**AI**
IS THE DRIVING FORCE

## INSIGHTS

**BUSINESS**

**OPERATIONAL**

**SECURITY**

# AI INSIDE INTEL

## REGULATORY
Audit & compliance automation

## HR
Diversity, recruiting & retention

## IT
Digital transformation with AI

## LOGISTICS
Supply chain optimization

## SALES
Info processing to improve efficiency

## HEALTH
Pharmaceutical analytics platform

## PRODUCTION
Factory process automation

## QUALITY
Automating visual defect detection
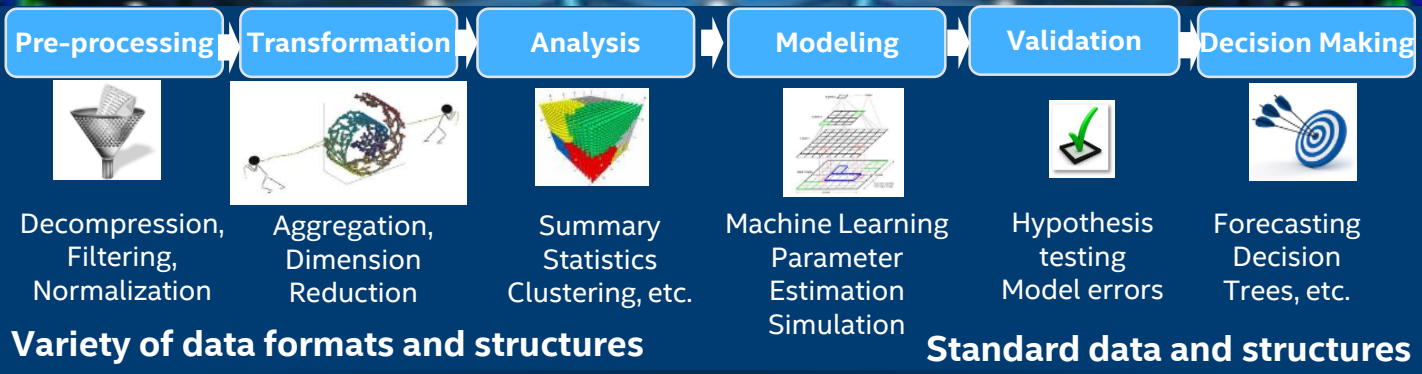
## RELIABILITY
Accelerating product validation

## INVENTORY
Optimizing inventory management

## AND MORE...

# INTEL® IS INFUSING AI INTO EVERYTHING WE DO
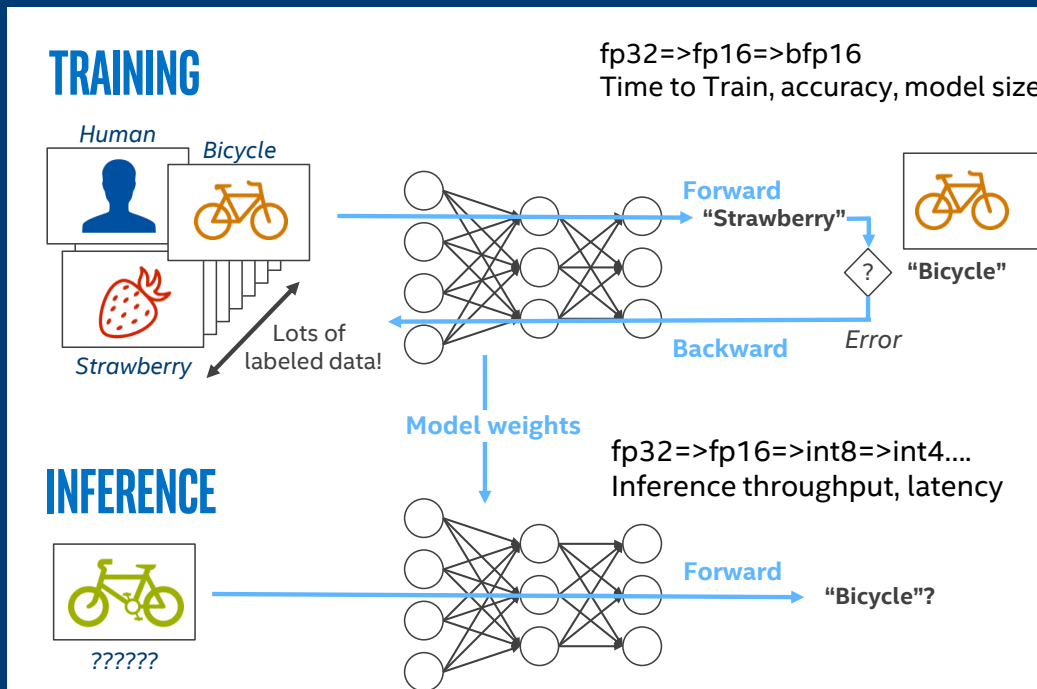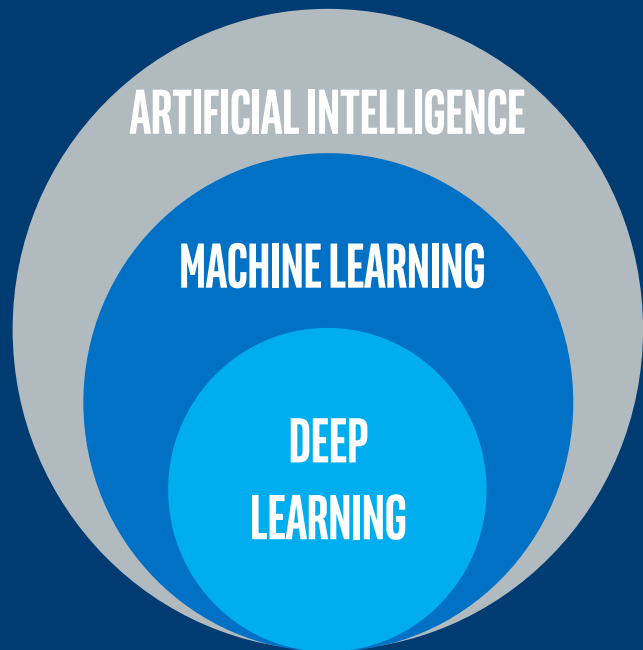
# End to End Data Analytics Flow



| Pre-processing | Transformation | Analysis | Modeling | Validation | Decision Making |
|---|---|---|---|---|---|
| Decompression, Filtering, Normalization | Aggregation, Dimension Reduction | Summary Statistics Clustering, etc. | Machine Learning Parameter Estimation Simulation | Hypothesis testing Model errors | Forecasting Decision Trees, etc. |

**Variety of data formats and structures**

**Standard data and structures**

**Data Prep, ETL, Dimension Reduction**

Start with Data

ETL → Feature Engr, Classical ML → Data Prep
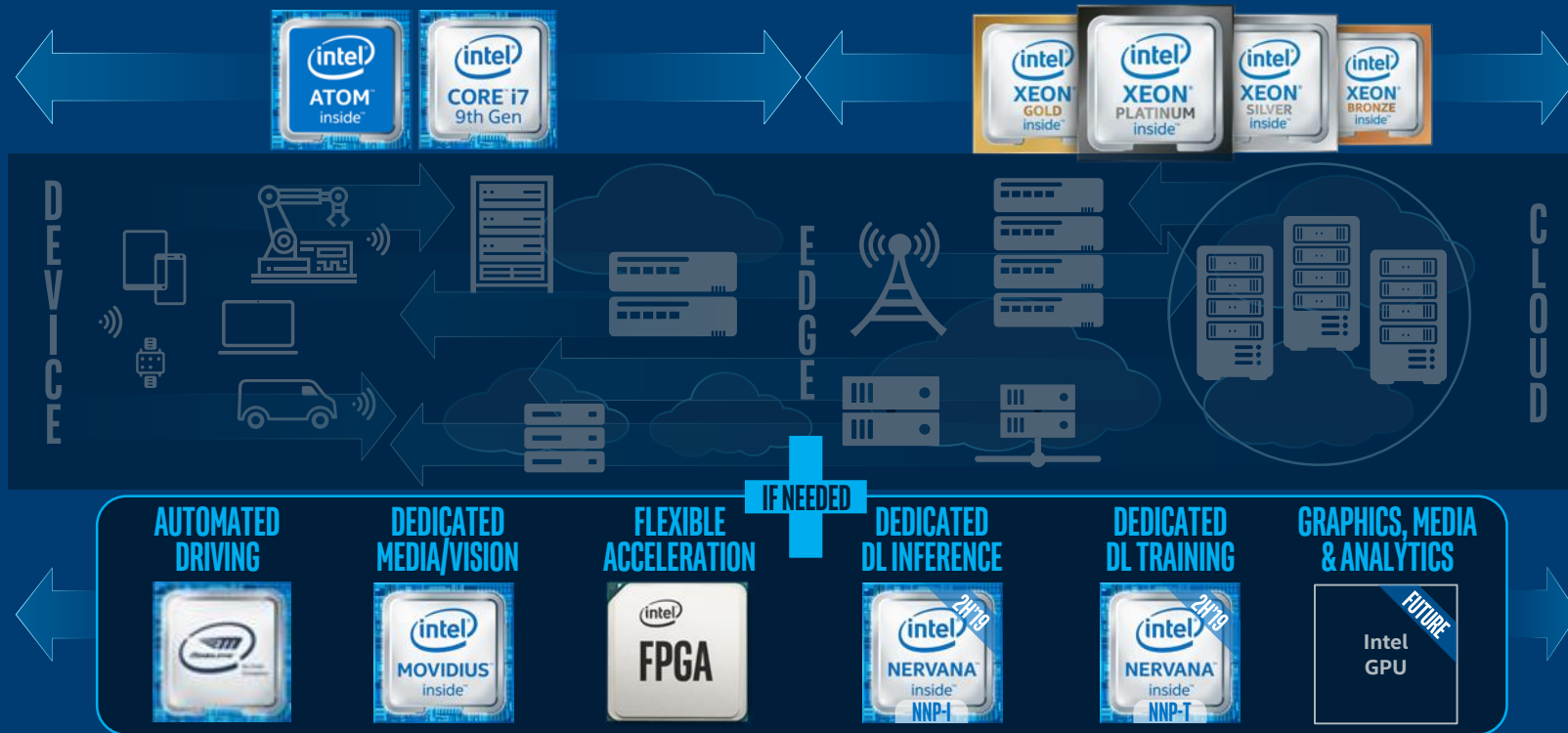
# DEEP LEARNING FLOW



Artificial Intelligence > Machine Learning > Deep Learning

TRAINING

Human
Bicycle
Strawberry

Lots of labeled data!

fp32=>fp16=>bfp16
Time to Train, accuracy, model size

Forward
"Strawberry"

?
"Bicycle"

Backward
Error

Model weights

INFERENCE

??????

fp32=>fp16=>int8=>int4....
Inference throughput, latency

Forward
"Bicycle"?

# AI IS EXPANDING

Deploy AI anywhere
with unprecedented hardware choice



DEVICE — EDGE — CLOUD

intel ATOM inside
intel CORE i7 9th Gen

intel XEON GOLD inside
intel XEON PLATINUM inside
intel XEON SILVER inside
intel XEON BRONZE inside

IF NEEDED

| AUTOMATED DRIVING | DEDICATED MEDIA/VISION | FLEXIBLE ACCELERATION | DEDICATED DL INFERENCE | DEDICATED DL TRAINING | GRAPHICS, MEDIA & ANALYTICS |
|---|---|---|---|---|---|
| MobilEye | intel MOVIDIUS inside | intel FPGA | intel NERVANA inside NNP-I 2H'19 | intel NERVANA inside NNP-T 2H'19 | Intel GPU FUTURE |

All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice.

intel

# CPU FOUNDATION FOR ARTIFICIAL INTELLIGENCE

**MATRIX OPERATIONS**
*Intel® Advanced Vector Extensions*

**LOWER & MIXED PRECISION**
*Intel® Deep Learning Boost*

**LARGER CACHES, MEMORY LATENCY & BANDWIDTH**

**OPTIMAL DATA MOVEMENT & TRANSFORMATIONS**

**OPTIMIZED LIBRARIES AND FRAMEWORKS**

| KTI | PCU | Global Hub | PCIe | PCIe | UPI |
|-----|-----|-----|-----|-----|-----|
| CORE | CORE | CORE | CORE | CORE | CORE |
| CORE | CORE | CORE | CORE | CORE | CORE |
| IMC | CORE | CORE | CORE | CORE | IMC |
| CORE | CORE | CORE | CORE | CORE | CORE |
| CORE | CORE | CORE | CORE | CORE | CORE |

MESH | LCC

**DATA CENTER PROCESSOR CORE**

MESH:
Intel® Mesh Architecture
LCC:
Processor Cache

Illustration: Intel® Xeon® Scalable Processor

**INTEL® XEON® SCALABLE PROCESSOR: ENABLES INFRASTRUCTURE-WIDE AI READINESS**

intel XEON PLATINUM inside

# INTEL DLBOOST – VNNI EXAMPLE

3X peak operations providing significant improvement in inferencing performance

# REINVENTING XEON FOR AI

**Intel® Optimization for Caffe ResNet-50[1] Inference Throughput**

Relative Inference Throughput (images/sec) (Higher is better)

**Intel® Deep Learning Boost**

**14X[3]**

Introducing new INT8 VNNI instruction

**30X[3]**

2S Intel® Xeon® Platinum 8280 processor (28 cores/S)

2S Intel® Xeon® Platinum 9282 processor (56 cores/S)

**INT8**

**5.7X** [2]

Enabling Lower precision & system optimizations for higher throughput
March 6th 2019

**FP32**

**2.8 X** [2]

With new library and framework optimizations
Jan 19th 2018

**FP32**

**1.0** [2]

Intel® Optimized Caffe
At launch, July 11th 2017

**Intel® Xeon® Platinum 8180 Processor (Codenamed: Skylake)**

**2nd Generation Intel® Xeon® Scalable Processor (Cascade Lake)**

# SOFTWARE IS ESSENTIAL

| | | | |
|---|---|---|---|
| **TOOLKITS** Application Developers | | **OPENVINO™ TOOLKIT** | **INTEL® MOVIDIUS™ SDK** |

**LIBRARIES** Data Scientists

**MACHINE LEARNING LIBRARIES**

Scikit-Learn    NumPy    MLlib

**DEEP LEARNING FRAMEWORKS**

TensorFlow    mxnet    Caffe    BigDL on Spark

Caffe2    PYT🔥RCH    Microsoft CNTK    PaddlePaddle

**FOUNDATION** Library Developers

**ANALYTICS, MACHINE & DEEP LEARNING PRIMITIVES**

MKL-DNN    clDNN    MLSL    Python    DAAL

**DEEP LEARNING GRAPH COMPILER**

Intel® nGraph™ Compiler

intel ATOM x / intel Iris Graphics    intel CORE i7 8th Gen    intel XEON inside    intel NERVANA inside    intel STRATIX 10 inside    intel ARRIA 10 inside    intel MOVIDIUS inside    INTEL GNA (IP)    MOBILEYE An Intel Company

**HARDWARE**

intel AI

# INTEL® DISTRIBUTION FOR PYTHON*

software.intel.com/intel-distribution-for-python

## FOR DEVELOPERS USING THE MOST POPULAR AND FASTEST GROWING PROGRAMMING LANGUAGE FOR AI

### EASY, OUT-OF-THE-BOX ACCESS TO HIGH PERFORMANCE PYTHON

- Prebuilt, optimized for numerical computing, data analytics, HPC
- Drop in replacement for your existing Python (no code changes required)

### DRIVE PERFORMANCE WITH MULTIPLE OPTIMIZATION TECHNIQUES

- Accelerated NumPy/SciPy/Scikit-Learn with Intel® MKL
- Data analytics with pyDAAL, enhanced thread scheduling with TBB, Jupyter* Notebook interface, Numba, Cython
- Scale easily with optimized MPI4Py and Jupyter notebooks

### FASTER ACCESS TO LATEST OPTIMIZATIONS FOR INTEL® ARCHITECTURE

- Distribution and individual optimized packages available through conda and Anaconda Cloud
- Optimizations upstreamed back to main Python trunk

## ADVANCING PYTHON* PERFORMANCE CLOSER TO NATIVE SPEEDS

# INTEL® DATA ANALYTICS ACCELERATION LIBRARY (INTEL® DAAL)

BUILDING BLOCKS FOR ALL DATA ANALYTICS STAGES, INCLUDING DATA PREPARATION, DATA MINING & MACHINE LEARNING

| Pre-processing | → | Transformation | → | Analysis | → | Modeling | → | Validation | → | Decision Making |

**Open Source | Apache\* 2.0 License**

**Common Python, Java and C++ APIs across all Intel hardware**

Optimized for large data sets including streaming and distributed processing

Flexible interfaces to leading big data platforms including Spark\* and range of data formats (CSV, SQL, etc.)

## HIGH PERFORMANCE MACHINE LEARNING AND DATA ANALYTICS LIBRARY

# INTEL® MATH KERNEL FOR DEEP LEARNING NEURAL NETWORKS (INTEL® MKL-DNN)

FOR DEVELOPERS OF DEEP LEARNING FRAMEWORKS FEATURING OPTIMIZED PERFORMANCE ON INTEL HARDWARE

## DISTRIBUTION DETAILS

- Open Source
- Apache* 2.0 License
- Common DNN APIs across all Intel hardware.
- Rapid release cycles, iterated with the DL community, to best support industry framework integration.
- Highly vectorized & threaded for maximal performance, based on the popular Intel® MKL library.

github.com/01org/mkl-dnn

**EXAMPLES:**

| Direct 2D Convolution | Local response normalization (LRN) | Rectified linear unit neuron activation (ReLU) | Maximum pooling | Inner product |

## Accelerate Performance of Deep Learning Models

Optimization Notice

# INTEL® DISTRIBUTION OF OPENVINO TOOLKIT

*OPTIMIZE EXISTING MODELS, RUN INFERENCE WHERE YOU NEED IT*

## Build high performance deep learning inference and computer vision

A toolkit to accelerate development of **high performance computer vision** & **deep learning inference into vision/AI applications** from edge to cloud. It enables deep learning on hardware accelerators and easy deployment across multiple types of Intel® platforms (CPU, GPU/Intel® Processor Graphics, FPGA, VPU).

**Who needs it?**
- Computer vision, deep learning developers
- Data scientists
- OEMs, ISVs, system integrators

**Usages**
Security surveillance, robotics, retail, healthcare, AI, office automation, transportation, non-vision use cases (speech, text) & more.

**HIGH PERFORMANCE AI EDGE TO CLOUD**

**STREAMLINED & OPTIMIZED DEEP LEARNING INFERENCE**

**HETEROGENEOUS, CROSS-PLATFORM FLEXIBILITY**

**Free Download** ▶ software.intel.com/openvino-toolkit
**Open Source version** ▶ 01.org/openvinotoolkit

# HPC ←→ AI: IMAGE ANALYSIS FOR DRUG DISCOVERY

## RESULTS

## NOVARTIS

### Joint Intel & Novartis collaboration

*Processing 1024x1280 large image dataset, reducing the training time to 31 minutes to >99% accuracy on 2S Intel® Xeon® processor based cluster.*



224 X 224 X 3
**ImageNet**

**26X**
**Larger**

1024 X 1280 X 3
**Microscopic Images**

(Highest Resolution Pathway)

Activation Size (in GB)

20.7GB

2.6GB   5.2GB   10.4GB

batchsize=8  batchsize=16  batchsize=32  batchsize=64

(Lowest Resolution Pathway)

High Content Screening/M-CNN Training on 8 Node Intel® 2S Xeon® 6148 processor HPC cluster
TensorFlow 1.7, Horovod, OpenMPI, BS=8/Worker, 4 Worker/Node, GBS=256, OPA Fabric

**Large Memory Usage Per Node**

**Time To Train**

31 Mins

■ TensorFlow variables
■ Σ(size of activation)

64.3GB

Speedup:
6.6x
Eff: 82.5%

16GB

17.5GB

3.4 Hrs

10.4

2.6

local batchsize=8     local batchsize=8     1 node    2 nodes    4 nodes    8 nodes
1 training worker     4 training workers

**Customer:**
Novartis Inst. of Biomedical Research (Switzerland) is one of the largest pharmaceutical companies in the world

**Challenge:** High content screening of cellular phenotypes is a fundamental tool supporting early stage drug discovery. While analyzing whole microscopic images are desirable, these images are 26X larger than benchmark dataset such as ImageNet*-1K. As a result, the high computational workload with high memory requirement would be prohibitive for deep learning model training

**Solution:** Intel and Novartis teams were able to scale and train the model with 32 TensorFlow* workers in 31 minutes.

(intel) XEON GOLD inside™

TensorFlow

http://aidc.gallery.video/detail/video/5790618241001/deep-learning-based-classification-of-high-content-cellular-images-on-intel-architecture?autoStart=true&q=Datta

(intel) AI | 19

# HPC ←→ AI: IMAGE ANALYSIS FOR DRUG DISCOVERY

High Content Screening Training with 313K Images on 64-Node
Intel® 2S Xeon® Scalable processor 6148, TensorFlow*, "horovod*",
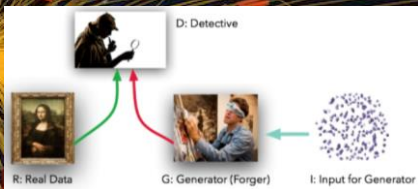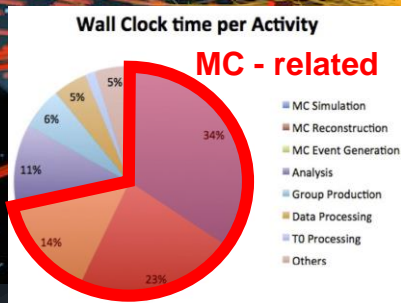OpenMPI*, Batch Size=32/Node, Intel® Omni-Path™ Fabric

# HPC ⟵⟶ AI: DIS/REPLACING MONTE CARLO SIM.
## CERN HIGH ENERGY PHYSICS
*JOINT COLLABORATION WITH INTEL AND SURFSARA*

# RESULT

*94% scaling efficiency up to 128 nodes, with a significant reduction in training time per epoch for 3D-GANs & >2500X Inference*



Wall Clock time per Activity

**MC - related**

- MC Simulation 34%
- MC Reconstruction 23%
- MC Event Generation 14%
- Analysis 11%
- Group Production 6%
- Data Processing 5%
- T0 Processing 5%
- Others

WLCG Wall Clock time for the ATLAS experiment



D: Detective
R: Real Data    G: Generator (Forger)    I: Input for Generator



GENERATOR

DISCRIMINATOR

**3D-Generative Adversarial Networks(GANs)**

### Time to create an electron shower

| Method | Machine | Time/Shower (msec) |
|---|---|---|
| **Full Simulation (geant4)** | 2S Intel® Xeon® Platinum 8180 | 17000 |
| **3D GAN (batch size 128)** | 2S Intel® Xeon® Platinum 8180 | **7** |

**Inference Perf: >2500X**

**Customer:** CERN, the European Organization for Nuclear Research, which operates the Large Hadron Collider (LHC), the world's largest and most powerful particle accelerator

**Challenge:** CERN currently uses Monte Carlo simulations for complex physics and geometry modeling, which is a heavy computational load that consumes up to >50% of the Worldwide LHC (Large Hadron Collider) Computing Grid (WLCG) power for electron shower simulations.

**Solution:** Distributed training using 128 nodes of the TACC Stampede 2 cluster (Intel® Xeon® Platinum 8160 processor, Intel® OPA) and a **3D Generative Adversarial Network (3D GAN)**. Performance was first optimized on a single node then scaled using TensorFlow* optimized with Intel® MKL-DNN, using 4 workers/node and an optimized number of convolutional filters.
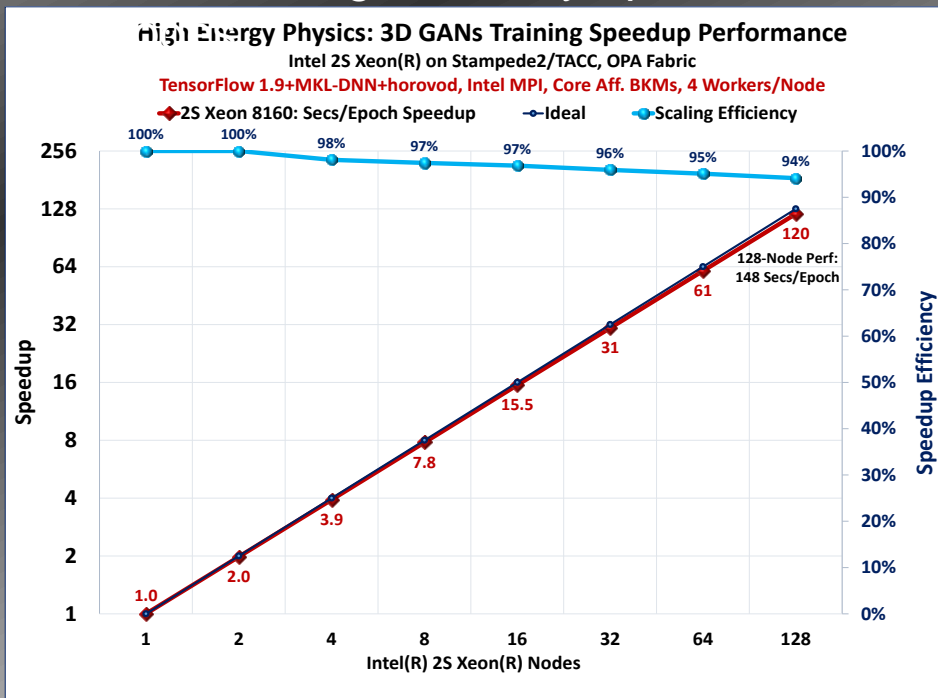
intel XEON PLATINUM inside™

TensorFlow

https://www.rdmag.com/article/2018/11/imagining-unthinkable-simulations-without-classical-monte-carlo

# Multi-Node Training Performance & Accuracy (2018)

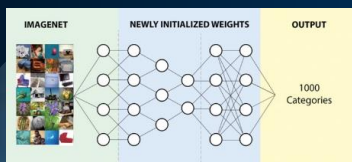**Distributed training using data parallelism**

94% Scaling efficiency up to 128



High Energy Physics: 3D GANs Training Speedup Performance
Intel 2S Xeon(R) on Stampede2/TACC, OPA Fabric
TensorFlow 1.9+MKL-DNN+horovod, Intel MPI, Core Aff. BKMs, 4 Workers/Node



Ratio of Ecal and Ep

# HPC ← → AI: CHEST X-RAY IMAGE CLASSIFICATION

## DellEMC
### *Joint collaboration with SURFsara, & DellEMC*



Identifying thoracic pathologies from the NIH ChestXray14 dataset
**Emphysema** affects more than: 3 Mil U.S & 65 Mil Worldwide
**Pneumonia** affects more than: 1 Mil US & 450 Mil Worldwide



**Transfer Learning**

# RESULT

*Training time reduced to 11 mins while increasing the accuracy across 10 categories & 4% (>90%) better relative to the existing DenseNet-121 model*



Categorical Accuracies in identifying diseases using ResNet-50 vs CheXNet-121

Better accuracy in 10 categories with scaled out trained ResNet50!

**Customer:** DellEMC* Research on AI applications on Intel® Xeon® CPUs: Medical, Cloud, HPC, etc.

**Challenge:** Train a chest X-ray model that delivers highly-efficient scaling performance on Intel® Xeon® processor nodes, while also delivering higher accuracy than the existing ChexNet-121 model

**Solution:** 256-node cluster consisting of dual Intel® Xeon® Gold 6148 processor, Intel® Omni-Path fabric, and ResNet-50 topology. ResNet50 tests performed with TensorFlow* and Horovod*.

https://ai.intel.com/diagnosing-lung-disease-using-deep-learning/

# TRAINING PERFORMANCE: RESNET-50 ON CHESTXRAY14

## INTEL® 2S XEON® GOLD 6148F PROCESSOR BASED DELLEMC* POWEREDGE C6420 ZENITH* CLUSTER ON OPA™ FABRIC
## TENSORFLOW* + "HOROVOD*", IMPI

Relative Training Throughput (images/sec) (Higher is Better)

200
180
160
140
120
100
80
60
40
20
0

1

104

**104x faster** using 128 Intel® Xeon® nodes!

**Training Time 8 MINUTES** to reach a solution with **256 2Skt Intel® Xeon® Gold 6148 processor**

187

**187x faster** using 256 Intel® Xeon® nodes!

TensorFlow ResNet-50 Node=1, Workers=4...

TensorFlow ResNet-50 Nodes=128, Workers=512...

TensorFlow ResNet-50 Nodes=256, Workers=1024...

(intel) AI | 24

# Case Study: Image Recognition

# World Bank

## RESULT

*High accuracy results using an AWS Databricks* platform to train a dataset consisting of almost 1 million images in 69 categories, with near linear scaling on a partial dataset*

**Client:** The International Comparison Program (ICP) in the World Bank Development Data Group

**Challenge:** The World Bank team needed to automate the process of confirming that the crowd-sourced photos, gathered from cellphone contributors from 15 countries, were accurately classified into one of 162 categories ranging from food to footwear, and to remove personally identifiable information (PII) from the photos.

**Solution:** Utilized Intel's BigDL framework (a distributed deep-learning library for Apache Spark*) and an AWS Databricks* platform running on Intel® Xeon® Processors (AWS R4.8xlarge instance with 20 nodes) to help classify more than 1 million crowdsourced photos before sharing the dataset with the public.

https://databricks.com/session/using-crowdsourced-images-to-create-image-recognition-models-with-bigdl
https://itpeernetwork.intel.com/artificial-intelligence-world-bank-image-recognition/

# CENTER FOR DIGITAL HEALTH INNOVATION (CDHI) AT UCSF

**RESULT**

*"Dataset, model development and training [...] implementing 3D CNN in BigDL to analyze MRI scans and classify OA (osteoarthritis) [...] provides rich 3D imaging support [...] on the same cluster where data is stored"*

(intel XEON inside)   hadoop*   Spark*   BigDL*

**Client:** Center for Digital Health Innovation (CDHI) at UCSF, leveraging new digital health technologies to transform healthcare.

**Challenge:** Projected by 2040 – 78M adults with doctor-diagnosed OA & 35M with arthritis-attributable activity limitations. Need automated system that classifies menisci based on presence/absence of lesions, provides immediate objective results at MRI scan, & eliminates intra-user variability.

**Solution:** Apache Spark* with BigDL on CDH 5.9*, on Intel® Xeon® servers from Dell*. With 3D image convolution in BigDL, the CDHI team built a MRI classification system & deployed it on their CDH Dell cluster.

https://cdn.oreillystatic.com/en/assets/1/event/269/Automatic%203D%20MRI%20knee%20damage%20classification%20with%203D%20CNN%20using%20BigDL%20on%20Spark%20Presentation.pdf

Artificial Intelligence will empower

**TRANSFORMATIVE INNOVATIONS**

WE ARE WITNESSING THIS DIGITAL TRANSFORMATION

CREATING