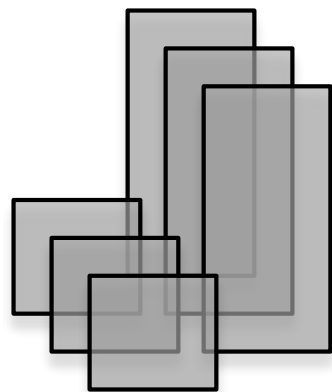# Compute for Everyone

Intel has a long tradition of democratizing compute

*by*

Making it **easier**

Making it **powerful**

Making it **accessible**

Mainframes
a few thousand users

# Compute for Everyone

Intel has a long tradition of democratizing compute

*by*

Making it **easier**

Making it **powerful**

Making it **accessible**

Mainframes
a few thousand users

Personal Computers
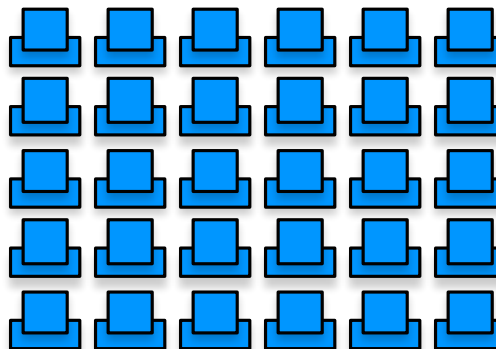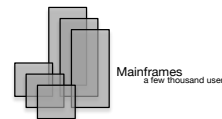millions and millions of users

# Compute for Everyone

Intel has a long tradition of democratizing compute

*by*

Making it **easier**

Making it **powerful**

Making it **accessible**

Mainframes
a few thousand users

Personal Computers
millions and millions of users
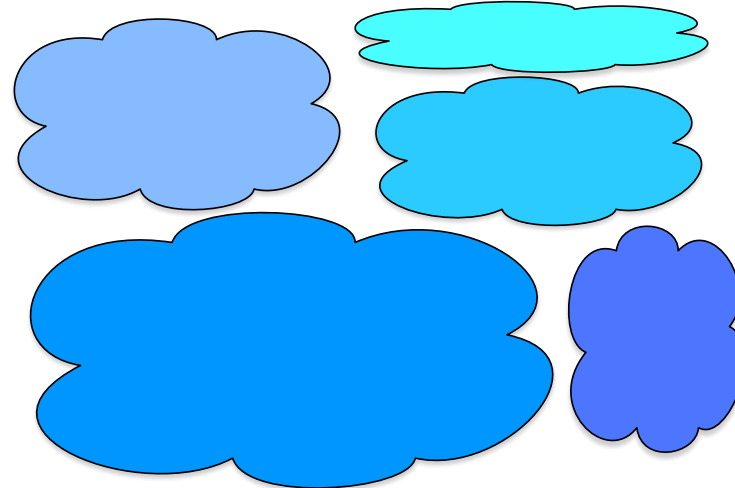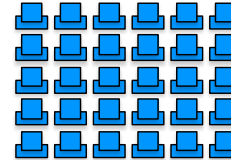
Cloud Computing
billions of users

# Compute for Everyone
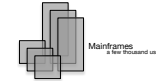
Intel has a long tradition of democratizing compute

*by*

Making it **easier**

Making it **powerful**

Making it **accessible**

Mainframes
a few thousand users

Personal Com...
millions and...

Cloud Computing
billions of users

Fog and IoT Computing
so many billions of users

# Compute for Everyone

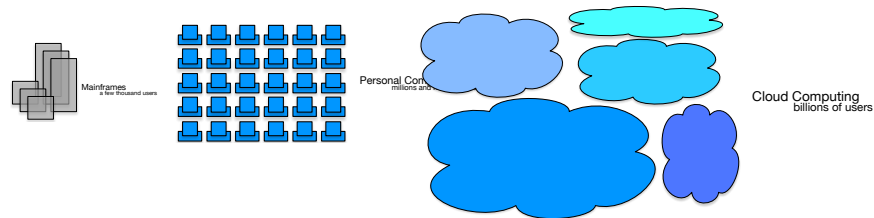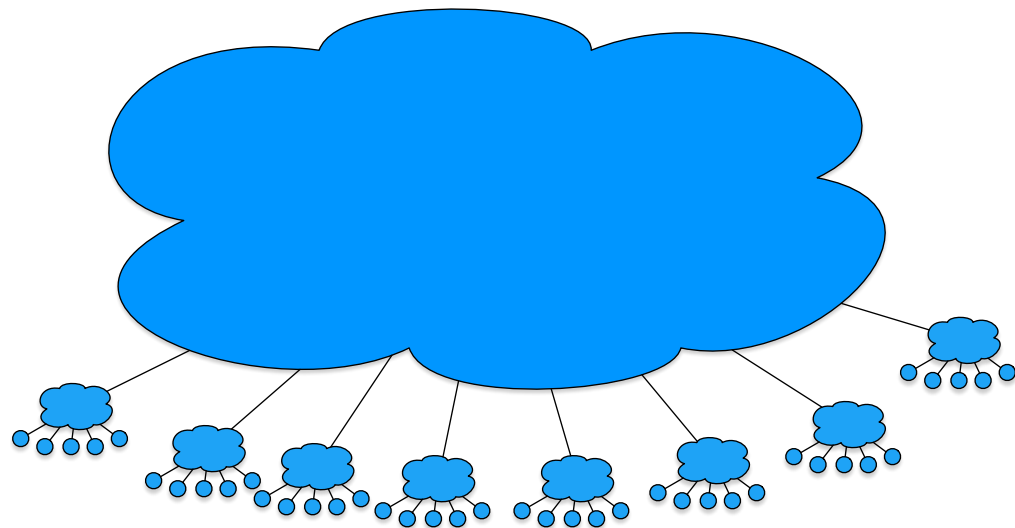Intel has a long tradition of democratizing compute

*by*

Making it **easier**

Making it **powerful**

Making it **accessible**

*but*

What does

## Democratizing AI

actually mean?

# Democratizing AI

## What does that actually mean?

Making it **easier**   *by*   **Automating** and abstracting anything that is not AI

Making it **powerful**

Making it **accessible**

# Democratizing AI

What does that actually mean?

Making it **easier**

Making it **powerful**

*by*

Making it **accessible**

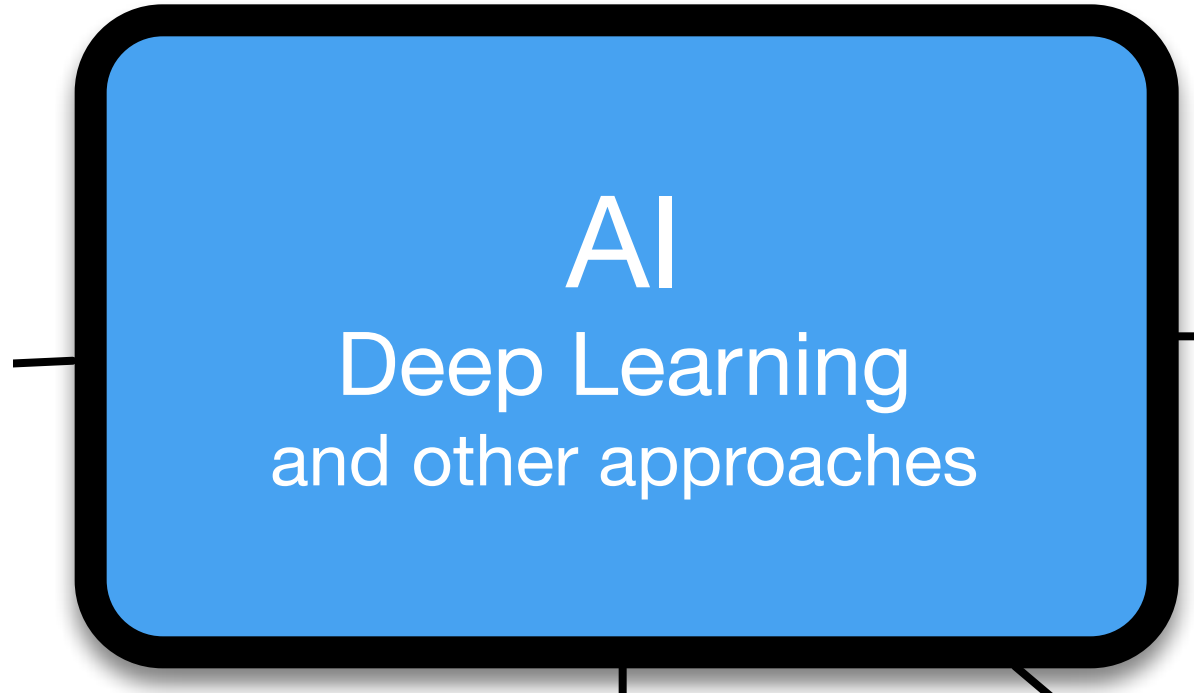**Automating** and abstracting anything that is not AI

**Enabling** scale up, scale out and novel AI techniques for *everyone*

# Democratizing AI

What does that actually mean?

Making it **easier** — **Automating** and abstracting anything that is not AI

Making it **powerful** — **Enabling** scale up, scale out and novel AI techniques for *everyone*

Making it **accessible** *by* **Bringing it** to the compute platform you already have

# Making AI Easier *by*

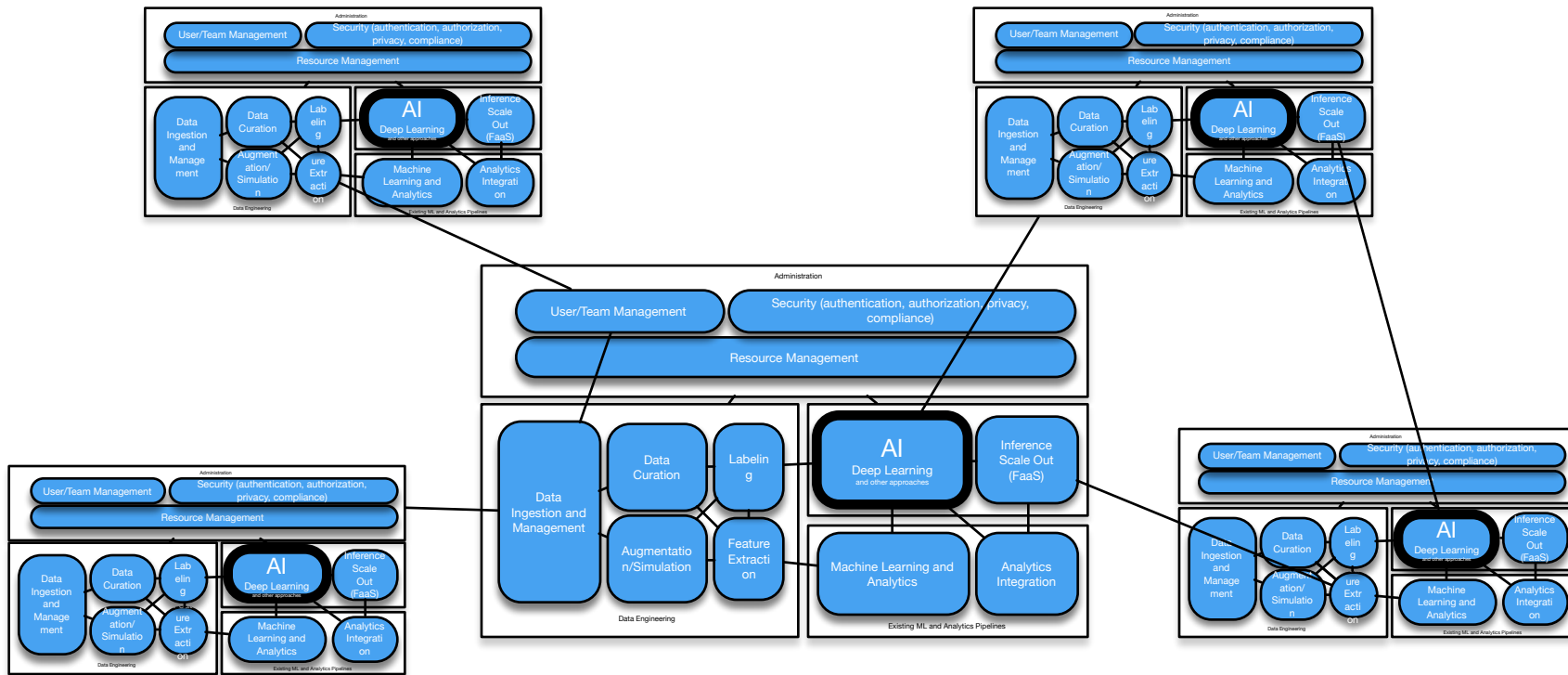## Automating and abstracting anything that is not AI



Administration

| User/Team Management | Security (authentication, authorization, privacy, compliance) |

Resource Management

**Data Engineering**

- Data Ingestion and Management
- Data Curation
- Labeling
- Augmentation/Simulation
- Feature Extraction

**AI**
Deep Learning
and other approaches

Inference Scale Out (FaaS)

**Existing ML and Analytics Pipelines**

- Machine Learning and Analytics
- Analytics Integration

# Making AI Easier *by* **Automating** and abstracting anything that is not AI

# Making AI Easier *by* **Automating** and abstracting anything that is not AI

# How do we solve messy problem?

# Making AI Easier *by* **Automating** and abstracting anything that is not AI

How do we solve messy problem?

The open source community, with Intel's support, is converging on solutions.

DLaaS offerings are flourishing

Kubernetes is the API

# Democratizing AI

What does that actually mean?

Making it **easier**

Making it **powerful**

Making it **accessible**

*by*

**Automating** and abstracting anything that is not AI

**Enabling** scale up, scale out and novel AI techniques for *everyone*

**Bringing it** to the compute platform you already have

# Making it Powerful

## Xeon Democratizes AI

*by*

## Enabling scale up, scale out and novel AI techniques for *everyone*

Intel® – SURFsara* Research Collaboration – Multi-Node Intel® Caffe ResNet-50
Scaling Efficiency on 2S Intel® Xeon® Platinum 8160 Processor Cluster



90% Scaling Efficiency

- MareNostrum4 Barcelona Supercomputing Center
- ImageNet-1K
- 256 nodes
- 90% scaling efficiency
- Top-1/Top-5 > 74%/92%
- Batch size of 32 per node
- Global BS=8192
- Throughput: 15170 Images/sec

## Time-To-Train: 70 minutes (50 Epochs)

# Xeon Democratizes AI

## INFERENCE THROUGHPUT

Up to
## 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

## TRAINING THROUGHPUT

Up to
## 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

Intel® Xeon® Platinum 8180 Processor higher Intel optimized Caffe Resnet50 with Intel® MKL inference throughput 133X and training throughput 73X compared to Intel® Xeon® Processor E5-2699 v3 with BVLC-Caffe

Inference and training throughput measured with FP32 instructions. Inference performance with INT8 is expected to be higher

**AI performance is constantly improving with hardware and software optimizations on Intel® Xeon® Scalable Processors**

# Xeon Democratizes AI: Case Study

GE Healthcare

**Intel's Solution Stack includes**

Intel® Xeon® Scalable processors

Intel® Solid State Drives

Intel Deep Learning Deployment Toolkit

Intel® Math Kernel Library for Deep Neural Networks

**OPTIMIZED MODEL**
Exceeds GE Inferencing Target

**14x** FASTER | **5.9x** ABOVE TARGET

# Democratizing AI

What does that actually mean?

Making it **easier**

Making it **powerful**

Making it **accessible**

*by*

**Automating** and abstracting anything that is not AI

**Enabling** scale up, scale out and novel AI techniques for *everyone*

**Bringing it** to the compute platform you already have

# Democratizing AI

Making it **accessible**

*by*

**Bringing it** to the compute platform you already have

**Optimizing Xeon AI**

**Augmenting Xeon with a broad compute portfolio**

**Enabling End-to-end AI**

*And most importantly*

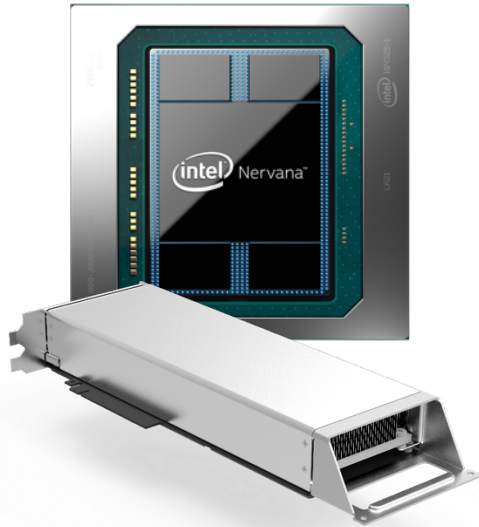**Making it easier to leverage the full stack**

**Intel AI PORTFOLIO**

| SOLUTIONS | | | |
|---|---|---|---|
| **PLATFORMS** | Intel® AI DevCloud | Intel® Deep Learning — Cloud / System˝ | Intel Saffron™ REASONING |
| **TOOLS** | Intel® Deep Learning Studio^ | Intel® Deep Learning Deployment Toolkit | Intel® Computer Vision SDK / Intel® Movidius™ SDK |
| **FRAMEWORKS** | TensorFlow Caffe mxnet Caffe2* PYTORCH* CNTK* PaddlePaddle* BigDL ON Spark neon | | |
| **LIBRARIES** | Intel® MKL/MKL-DNN, clDNN, DAAL, Intel Python Distribution, etc. DIRECT OPTIMIZATION | Intel® nGraph™ Library — CPU Transformer / NNP Transformer† / More... | |
| **TECHNOLOGY** | DATACENTER | EDGE/GATEWAY | SYSTEMS & COMPONENTS |

**SOLUTIONS:** Data Scientists · Technical Services · Reference Solutions

*Future product
†Beta available
˝Available in the Intel® Deep Learning Cloud, coming to other platforms later
Other names and brands may be claimed as the property of others.

# INTEL® NERVANA™ NEURAL NETWORK PROCESSOR (NNP)¥

Scalable acceleration with best performance for intensive deep learning

## PARALLELISM

Massively-parallel compute

Specialized on-die fabrics

Optimized numerics - Flexpoint

## SCALABILITY

Large on-die memory

High speed interconnects

Massive inter-chip data transfer

## UTILIZATION

Direct SW control for best on-chip memory usage

Managed data-flow paths

## ROADMAP

First silicon in 2017

Product roadmap on track to exceed performance goal[1]

# PROJECT BRAINWAVE FOR REAL-TIME AI

"A major leap forward in both performance and flexibility for cloud-based serving of deep learning models."

Doug Burger
Distinguished Engineer
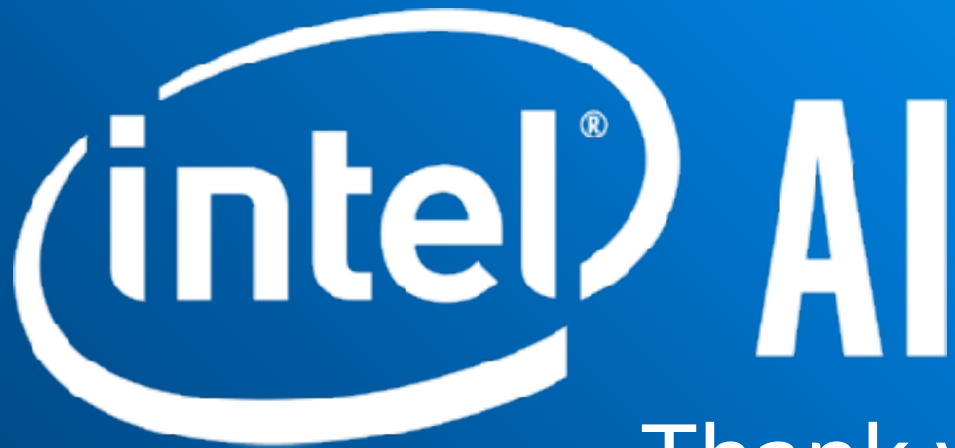
# INTEL IS DEMOCRATIZING AI

# INTEL IS DEMOCRATIZING AI

*by*

Offering edge-to-edge AI compute solutions

Developing key AI software with the open source community

*and*

## Making it work better together

# Notices and 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 complete information visit: **http://www.intel.com/performance**.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit **www.intel.com/benchmarks**.

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.

The products described 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, the Intel logo, Xeon, Xeon Phi and Nervana are trademarks of Intel Corporation in the U.S. and/or other countries.

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