



**OCP**  
SUMMIT

March 20-21  
**2018**  
San Jose, CA

**OPEN. FOR BUSINESS.**



# INTEL DEEP LEARNING PLATFORM

**Jordan Plawner, Sr. Director, Product Management**  
Artificial Intelligence Product Group, Intel

**OPEN. FOR BUSINESS.**



# AI Adoption in the Enterprise is Just Beginning

In a recent Forrester Research survey...

58%

of business and technology professionals said they're researching AI, but **only...**

12%

said they are currently using AI systems.

OPEN. FOR BUSINESS.

# The AI Opportunity: Democratize AI

**NARROW AI** *DRIVEN BY INDUSTRY*



*Dedicated to assist with  
or take over specific tasks*

**STRONG AI** *DRIVEN BY IT*

## **PLANNING**

Smart R&D and forecasting

## **PRODUCTION**

Optimize production and maintenance

## **PROMOTION**

Targeted sales and marketing

## **PROVIDE**

Enhanced user experiences

## **IT TECHNICAL SERVICES**

Image classification, language translation, facial recognition, speech-2-text, speech generation, time series...

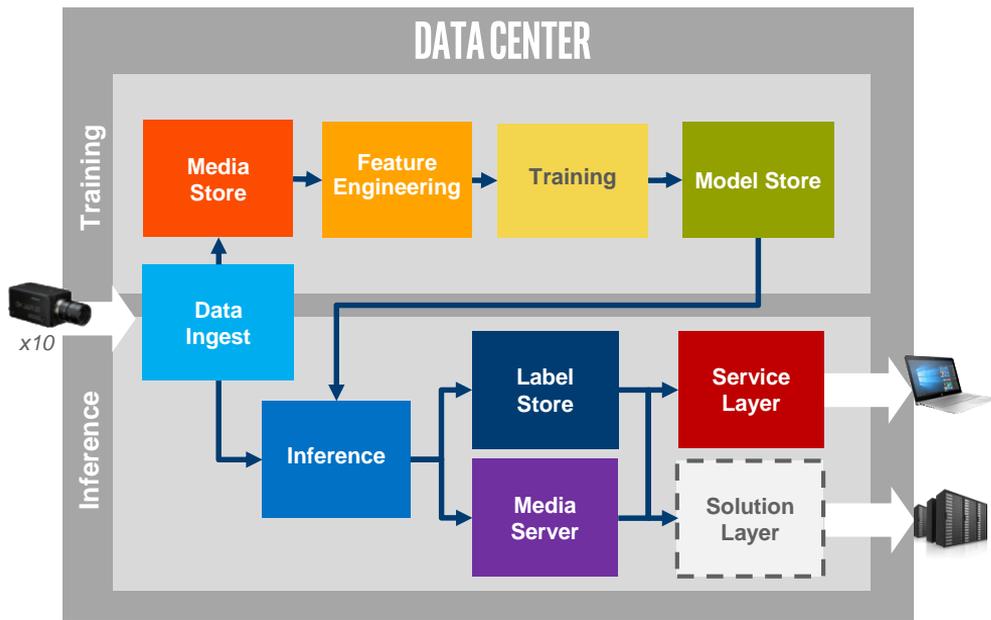
## **PRODUCTION OPERATIONS**

Data pipeline and workflow, logging and monitoring, feature store, life cycle management...

**OPEN. FOR BUSINESS.**



# Deploying Deep Learning



## Media Storage

- Media Store
- Media Store
- Media Store

110 Nodes  
 8 TB/day per camera  
 10 cameras  
 3x replication  
 1-year video retention  
 4 mgmt nodes

- Media Store
- Media Store
- Media Store

Per Node  
 1x Intel Xeon E5-2680v4  
 20x 4TB SSD

## Multi-Purpose Cluster

- Data Ingestion 4 nodes
- Data Ingestion One ingestion per day, one-day retention
- Data Ingestion
- Data Ingestion
- Inference 4 nodes
- Inference 20M frames per day
- Inference

Tagged Datasets 2 nodes  
 Tagged Datasets Infrequent op

- Service Layer 3 nodes
- Service Layer Simult users

- Media Server 3 nodes
- Media Server 10k clips stored
- Media Server

Per Node  
 1x Intel Xeon E5-2680v4  
 20x 4TB SSD

## Data Storage

- Model Store 4 nodes
- Model Store 1-year of history
- Model Store
- Model Store
- Label Store 4 nodes
- Label Store Labels for 20M frames/day
- Label Store

Per Node  
 1x Intel Xeon E5-2680v4  
 5x 4TB SSD

## Training

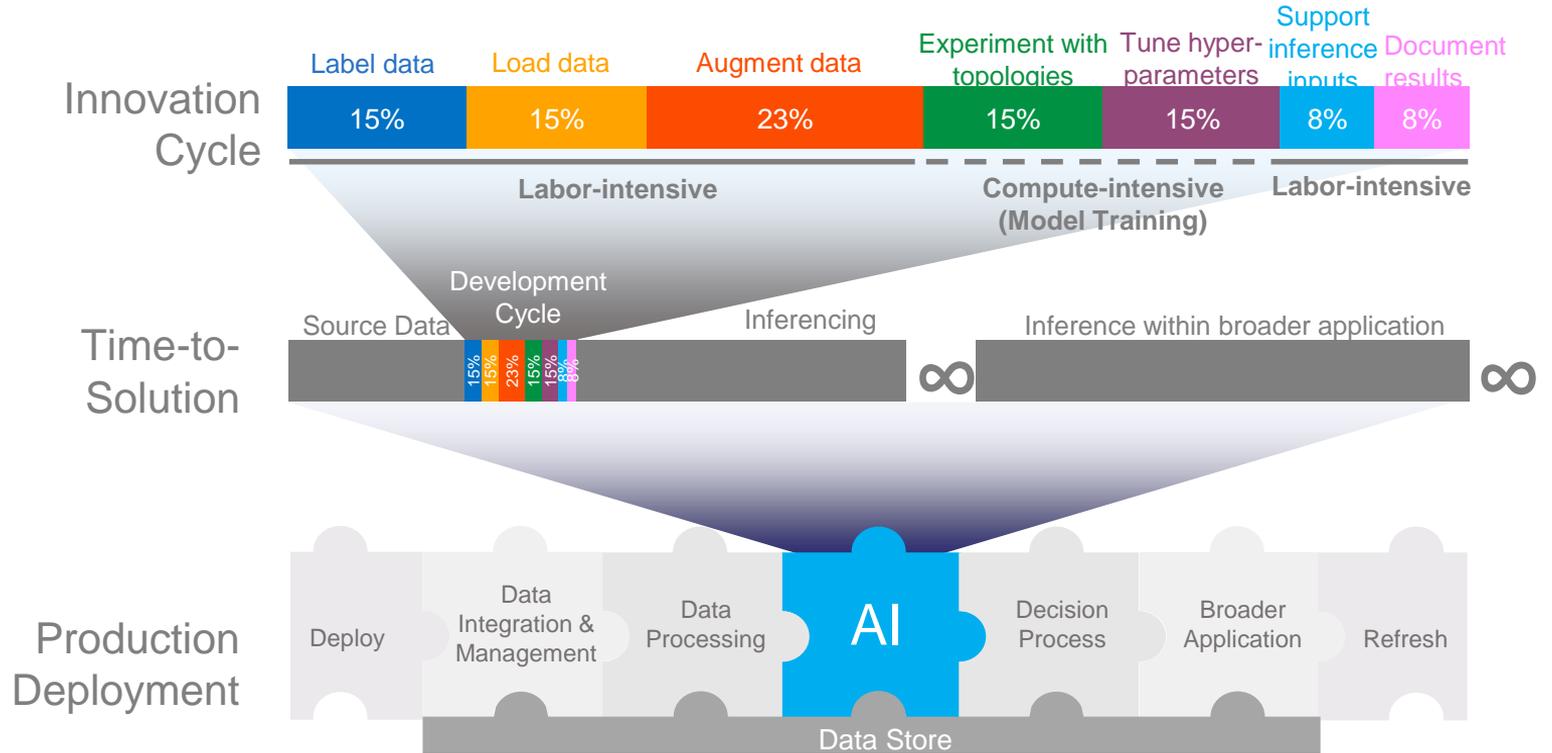
- Training 16 nodes < 10 hours TTT
- Training

Per Node  
 1x Intel Xeon E5-2680v4  
 1x 4TB SSD

OPEN. FOR BUSINESS.



# Deploying Deep Learning Time-to-Solution



OPEN. FOR BUSINESS.

# Intel AI Portfolio

*If it computes and is connected it will do AI*

**GENERAL  
AI**



Mainstream AI



Flexible Acceleration

**TRAINING**

**DATA CENTER/  
WORKSTATION**



Mainstream Training



Intensive Training

**INFERENCE**

**DATA CENTER/  
WORKSTATION**



Mainstream Inference



To be announced

Intensive Inference



Real-time Inference

**GATEWAY/EDGE**



Mainstream Inference



Higher Inference Throughput



Vision 1-20W



Speech/Audio 1-100+mW



Autonomous driving



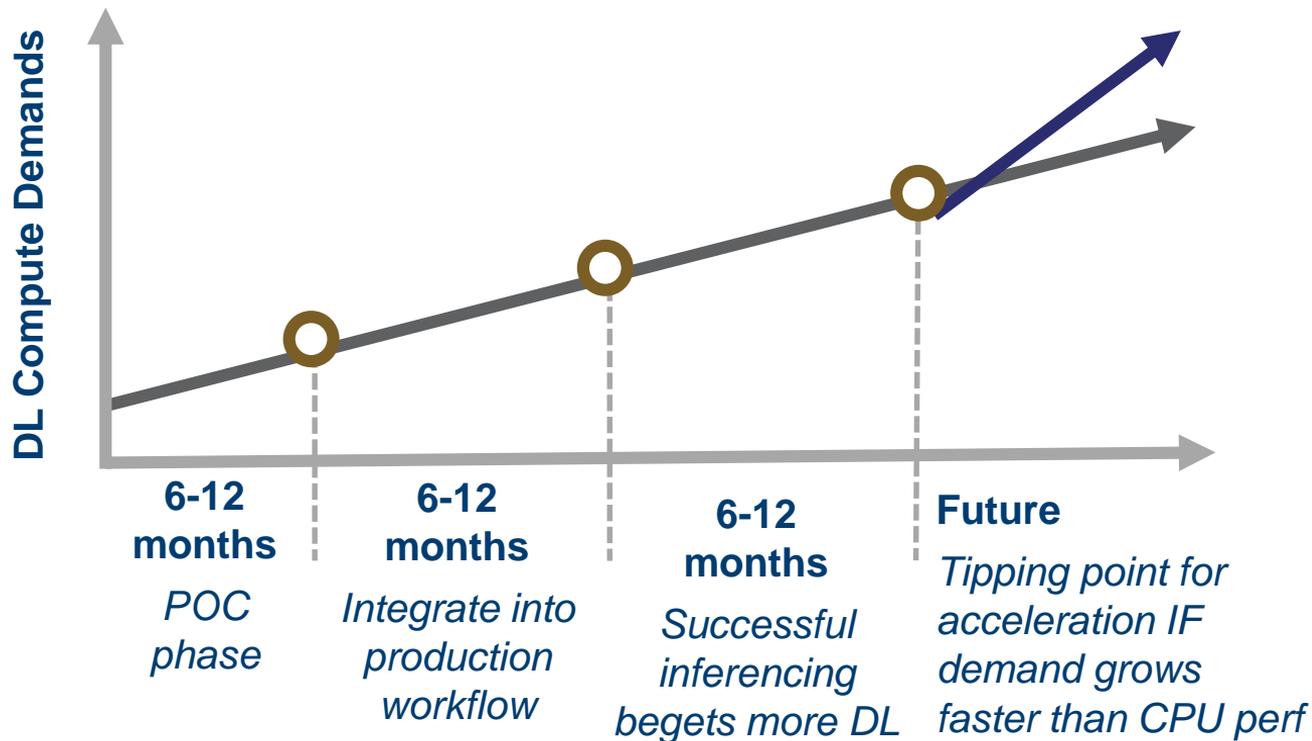
Custom Inference

**DEEP  
LEARNING**

**OPEN. FOR BUSINESS.**



# Deep Learning Journey



OPEN. FOR BUSINESS.



OCP  
SUMMIT

# Deep Learning Acceleration Options

	Utilization = TCO	Intermittent training Intermittent or low level of inferencing	Year Round intense training High inferencing compute utilization/node
	Time-to-train	Time-to-train in ~1 day is sufficient with access to additional compute for sale out	Time-to-train in hours is critical for scientist time year round
	Time-to-solution	Fastest by leveraging existing analytics, storage & management infrastructure	Due to high utilization at scale integrating dedicated infrastructure is preferred
	Flexibility & scalability	Resource allocation of shared compute & extend storage pool for unstructured data	No requirement to reallocate compute cycles throughout the year
	Locality	Distributed training and inferencing on edge and premise (eg; factory)	High training and batch inferencing utilization in core datacenter
	Workload	Hybrid of general compute & NN acceleration High definition images exceeding 32GHBM	Mixed precision linear algebra Model layers, weights and data can fit in HBM

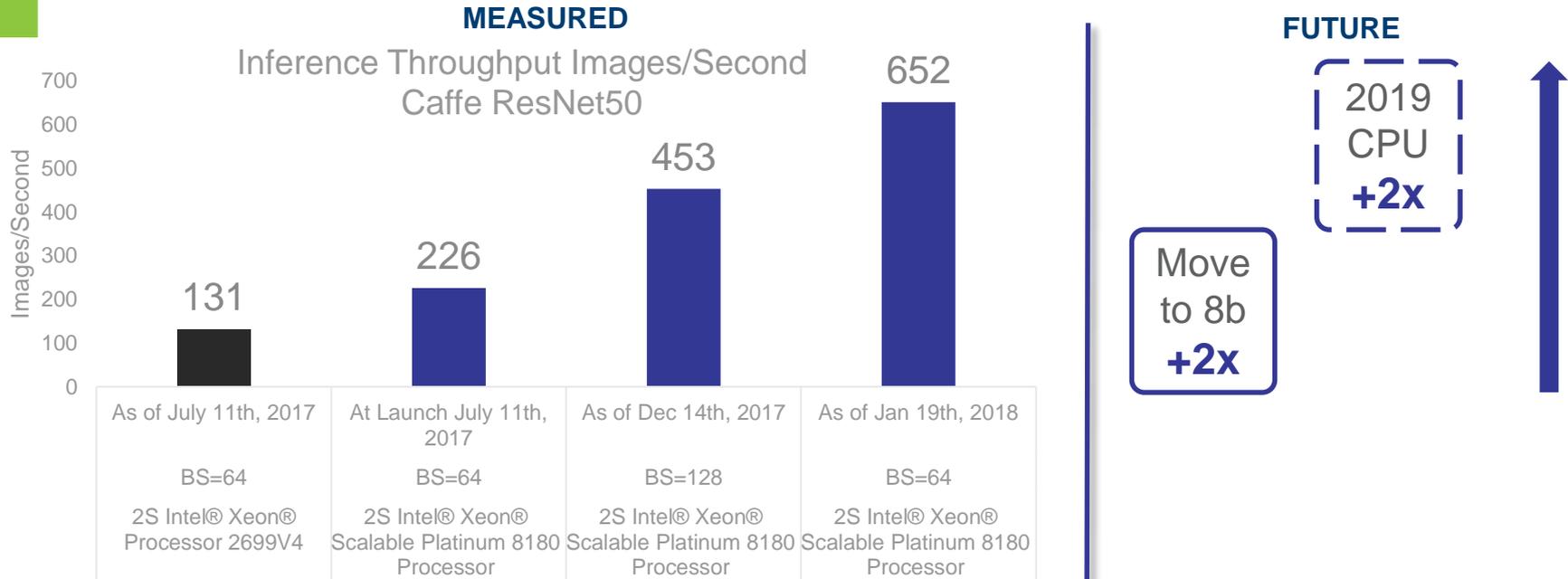
OPEN. FOR BUSINESS.





# Xeon SP Performance:

Up to 2.2X improvement with SW optimizations



Configuration Details 1, 14,

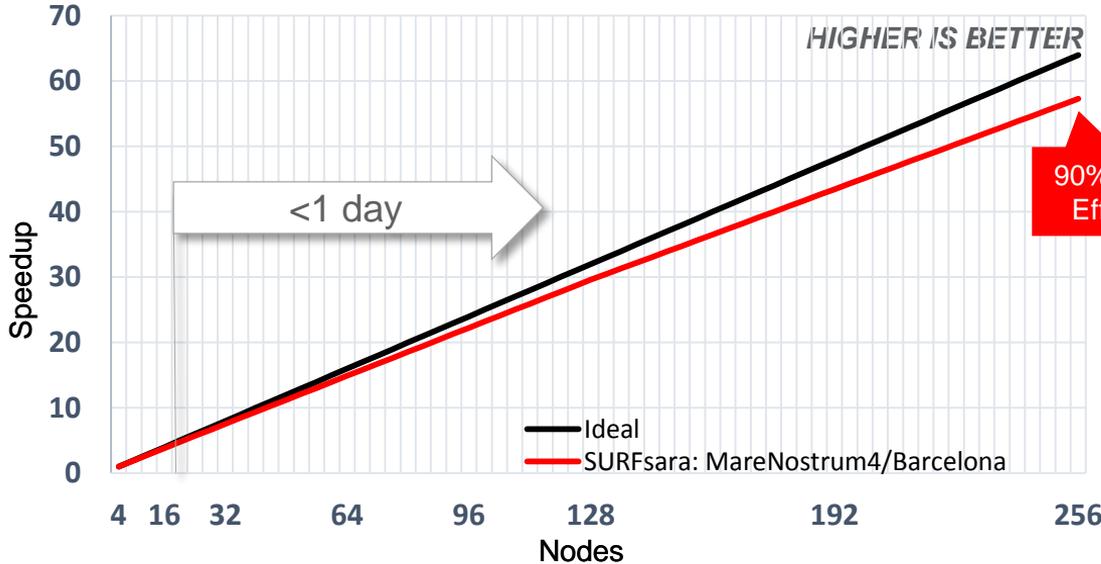
Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your system. Optimization Notice: Intel's compilers 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. 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 or estimated as of November 2017.

OPEN. FOR BUSINESS.



# Xeon SP Scaling on CPU

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



- 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)**

90% scaling efficiency with up to 74% Top-1 accuracy on 256 nodes

Configuration Details 2: Slide 127

Performance estimates were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown." Implementation of these updates may make these results inapplicable to your device or system. 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.

Optimization Notice: Intel's compiler may or may not optimize for the same degree of 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.

OCF  
SUMMIT



# Intel® Nervana™ neural network processor (NNP)

## 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<sup>1</sup>



¥ Formerly codenamed as the Crest Family

<sup>1</sup>Source: [https://newsroom.intel.com/news-releases/intel-ai-day-news-release/?\\_ga=2.26542141.1088441208.1508441324-198894050.1498491572](https://newsroom.intel.com/news-releases/intel-ai-day-news-release/?_ga=2.26542141.1088441208.1508441324-198894050.1498491572).

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

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 or estimated as of November 2017

# OPEN. FOR BUSINESS.





# Summary

## FINAL THOUGHTS

- Most AI is **NOT** deep learning
- Most DL is **NOT** training
- AI **value** is in convergence

## NEXT STEPS

Start **your journey with Xeon**

- Learn: **Focused POCs**
- Plan: **Develop IT DL technical services and platforms**
- Deploy: **Integrate with existing IT workflow**

**OPEN. FOR BUSINESS.**



# Find Out More

**LEARN**

More information at [ai.intel.com](https://ai.intel.com)

**EXPLORE**

Use Intel's performance-optimized libraries & frameworks

**ENGAGE**

Contact your Intel representative for help and POC opportunities



**OPEN. FOR BUSINESS.**





# OCP SUMMIT