Artificial Intelligence

# Agenda

- Why do we need to scale machine learning?

- What makes it hard to scale and how we are addressing it

- Real-world applications

- Hardware roadmap, software tools and frameworks update

IDF16
INTEL DEVELOPER FORUM

# Deep Learning: Scoring or Inferencing



Person

Hidden Layers

Deep Neural Network Model

# Deep Learning: Scoring or Inferencing



Forward Propagation

Person

Hidden Layers

Deep Neural Network Model

# Deep Learning: Training
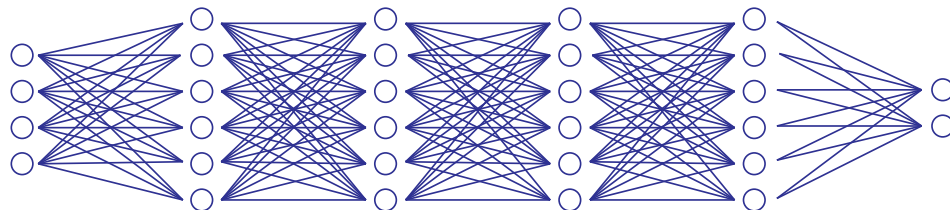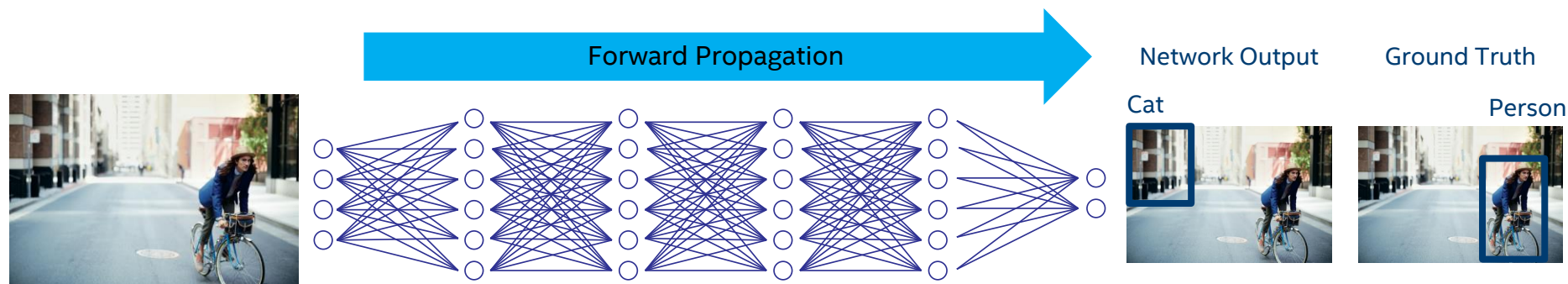
* Shihao Ji, S. V. N. Viswanathan, Nadathur Satish, Michael Anderson, and Pradeep Dubey. Blackout: Speeding up Recurrent Neural Network Language Models with very large vocabularies. http://arxiv.org/pdf/1511.06909v5.pdf. ICLR 2016

# Deep Learning: Training



Forward Propagation

Network Output

Cat

Ground Truth

Person

* Shihao Ji, S. V. N. Viswanathan, Nadathur Satish, Michael Anderson, and Pradeep Dubey. Blackout: Speeding up Recurrent Neural Network Language Models with very large vocabularies. http://arxiv.org/pdf/1511.06909v5.pdf. ICLR 2016

# Deep Learning: Training



Forward Propagation

Backward Propagation

Network Output
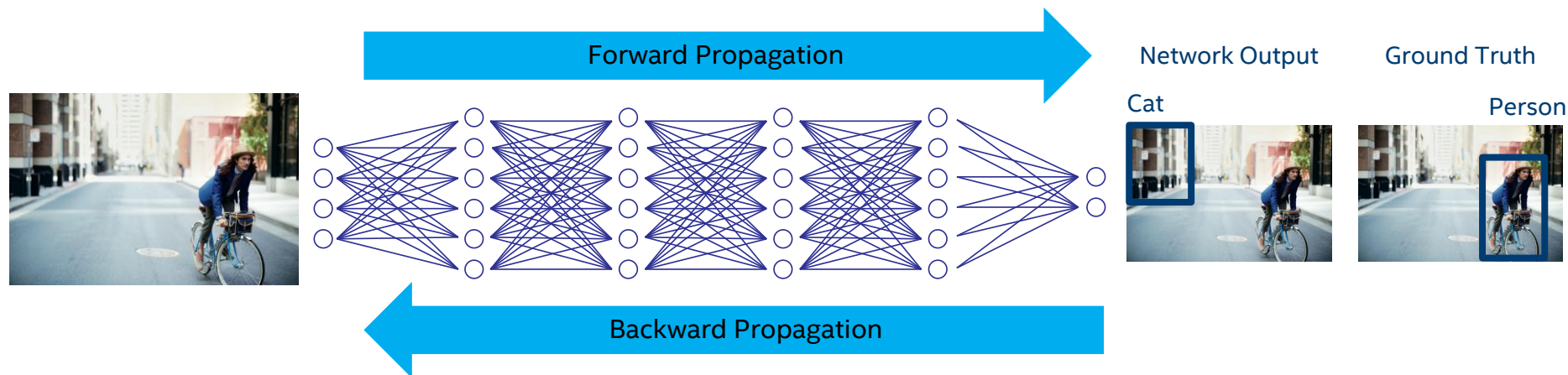
Cat

Ground Truth

Person

* Shihao Ji, S. V. N. Viswanathan, Nadathur Satish, Michael Anderson, and Pradeep Dubey. Blackout: Speeding up Recurrent Neural Network Language Models with very large vocabularies. http://arxiv.org/pdf/1511.06909v5.pdf. ICLR 2016

IDF16
INTEL DEVELOPER FORUM

# Deep Learning: Training



Forward Propagation

Backward Propagation
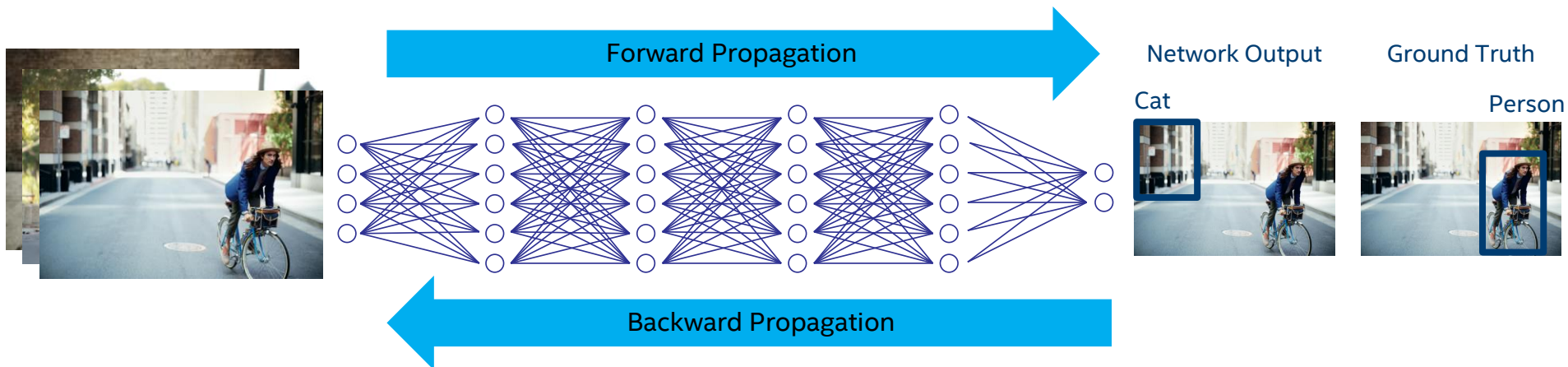
Network Output
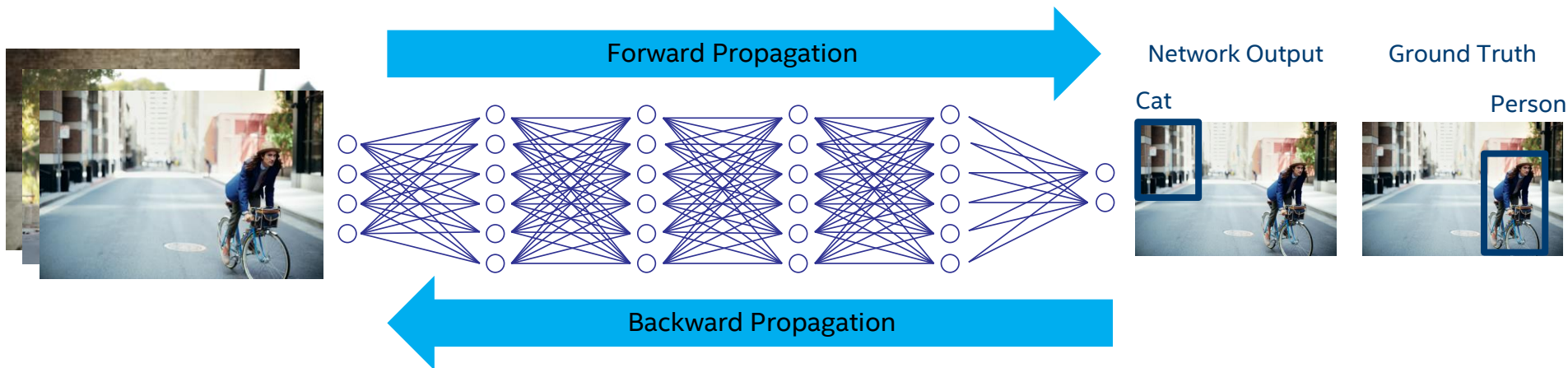
Cat

Ground Truth

Person

* Shihao Ji, S. V. N. Viswanathan, Nadathur Satish, Michael Anderson, and Pradeep Dubey. Blackout: Speeding up Recurrent Neural Network Language Models with very large vocabularies. http://arxiv.org/pdf/1511.06909v5.pdf. ICLR 2016

IDF16
INTEL DEVELOPER FORUM

# Deep Learning: Training



Forward Propagation

Network Output

Cat

Ground Truth
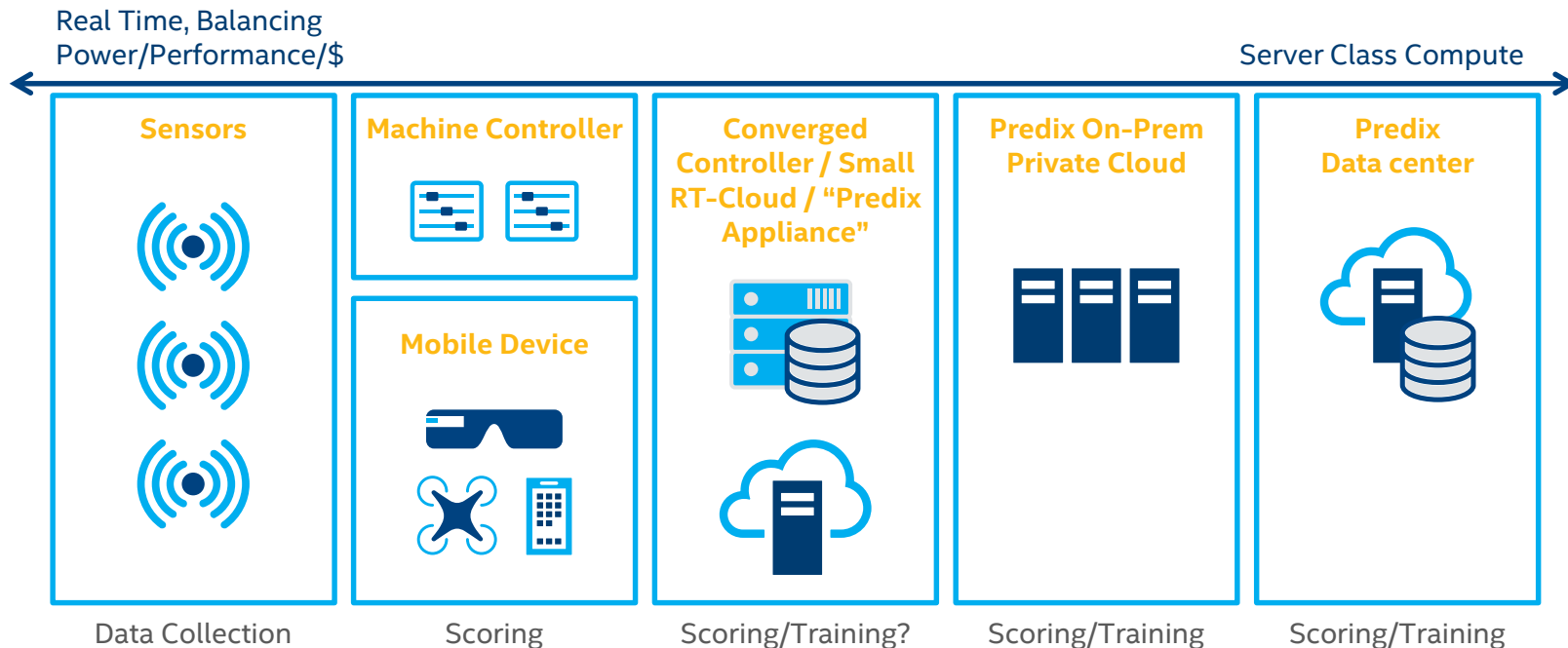
Person

Backward Propagation

Complex Networks with billions of parameters can take days to train on a modern processor*
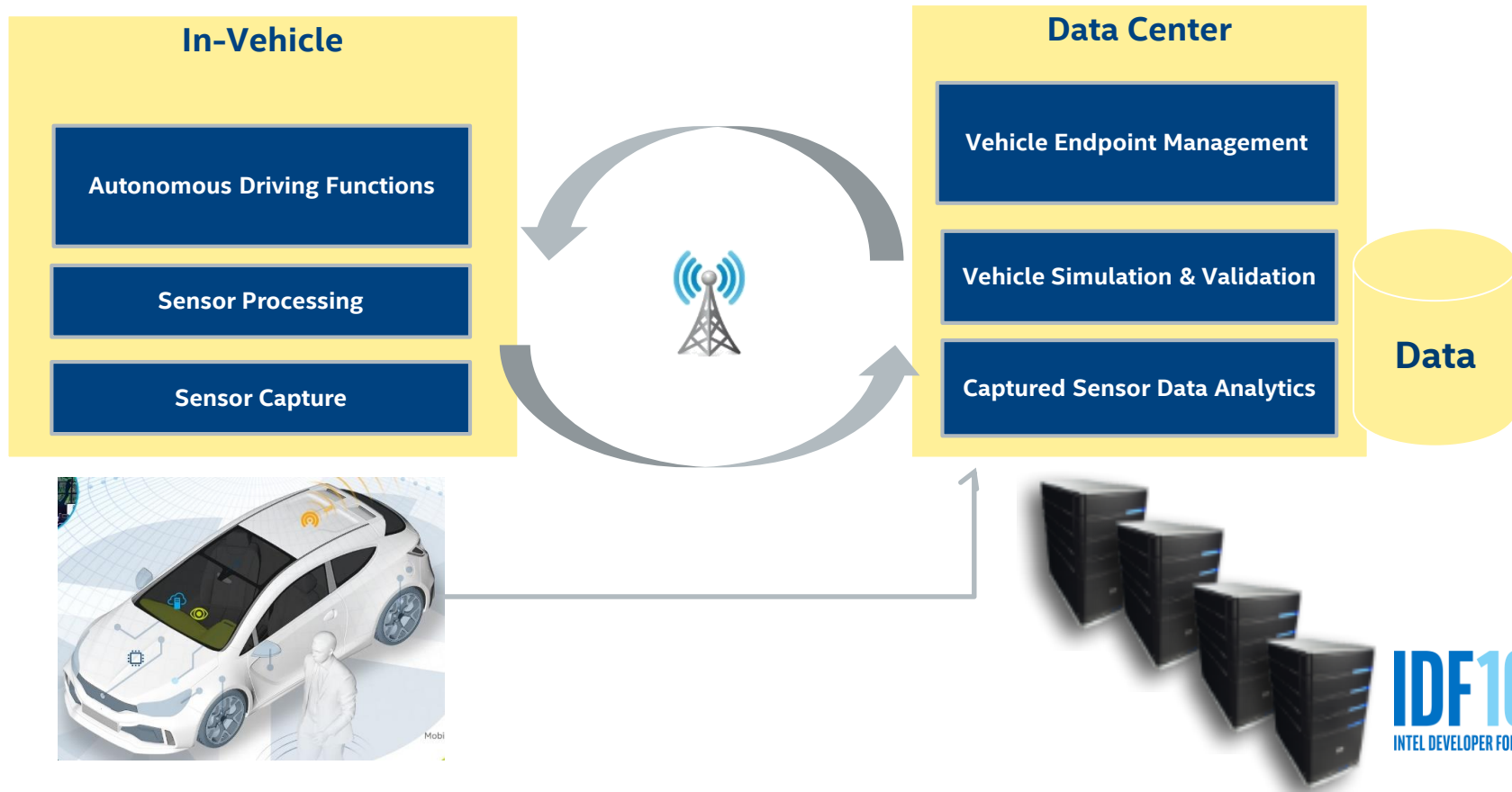
Hence, the need to reduce time-to-train using a cluster of processing nodes

* Shihao Ji, S. V. N. Viswanathan, Nadathur Satish, Michael Anderson, and Pradeep Dubey. Blackout: Speeding up Recurrent Neural Network Language Models with very large vocabularies. http://arxiv.org/pdf/1511.06909v5.pdf. ICLR 2016

IDF16
INTEL DEVELOPER FORUM

# Machine Learning Continuum: Connected Factory



Real Time, Balancing Power/Performance/$ ← → Server Class Compute

| Sensors | Machine Controller / Mobile Device | Converged Controller / Small RT-Cloud / "Predix Appliance" | Predix On-Prem Private Cloud | Predix Data center |
| --- | --- | --- | --- | --- |
| Data Collection | Scoring | Scoring/Training? | Scoring/Training | Scoring/Training |

IDF16 INTEL DEVELOPER FORUM

# Machine Learning Continuum: Self-Driving Cars



**In-Vehicle**

- Autonomous Driving Functions
- Sensor Processing
- Sensor Capture

**Data Center**

- Vehicle Endpoint Management
- Vehicle Simulation & Validation
- Captured Sensor Data Analytics

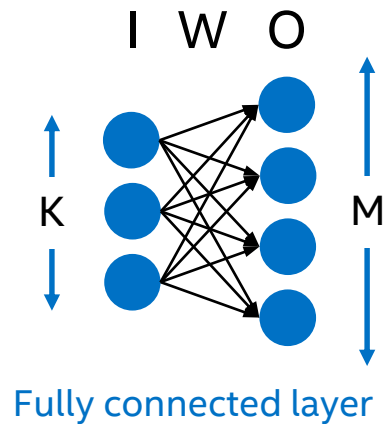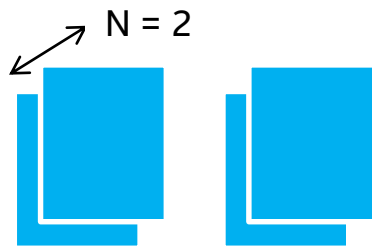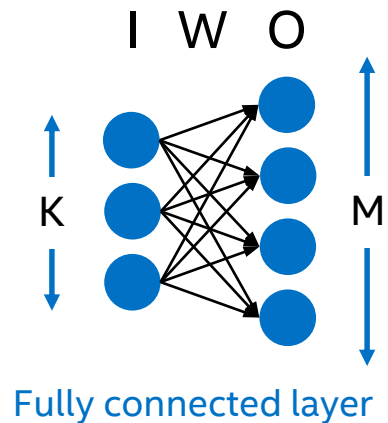**Data**

IDF16
INTEL DEVELOPER FORUM

# Agenda

- Why do we need to scale machine learning

- What makes it hard to scale and how we are addressing it

- Real-world experience of an industry leader

- Hardware roadmap, software tools and frameworks update

- Summary

IDF16
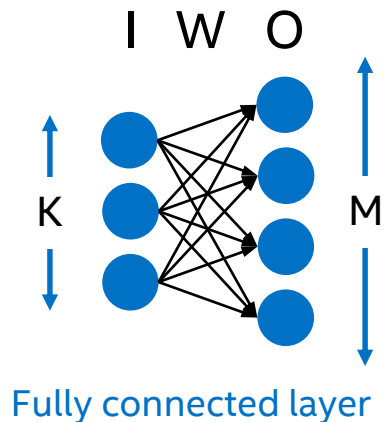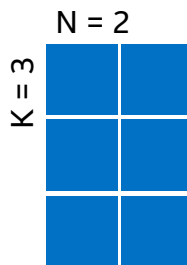INTEL DEVELOPER FORUM

# Compute Kernels

# Compute Kernels

I  W  O



Fully connected layer

# Compute Kernels



I   W   O

K

M

Fully connected layer

N = 2

# Compute Kernels

$$I \in R^{KxN}$$
*Input*

$$W \in R^{MxK}$$
*Weights or model*

$$O \in R^{MxN}$$
*Output or activations*

Fully connected layer

N = 2

# Compute Kernels



Fully connected layer

$$I \in R^{KxN}$$
*Input*

$$W \in R^{MxK}$$
*Weights or model*

$$O \in R^{MxN}$$
*Output or activations*

Forward propagation: (M x K) * (K x N)

Backward propagation: (M x K)$^T$ * (M x N)

Weight update: (M x N) * (K x N)$^T$

IDF16
INTEL DEVELOPER FORUM

# Parallelism Options

**I**

*Input data*

**W**

*Weights or model*

**O**

*Output or activations*

# Parallelism Options



**I**
*Input data*

**W**
*Weights or model*

**O**
*Output or activations*

Data

# Parallelism Options

**I**
*Input data*

**W**
*Weights or model*

**O**
*Output or activations*

Model

# Parallelism Options

**I**
*Input data*

**W**
*Weights
or model*

**O**
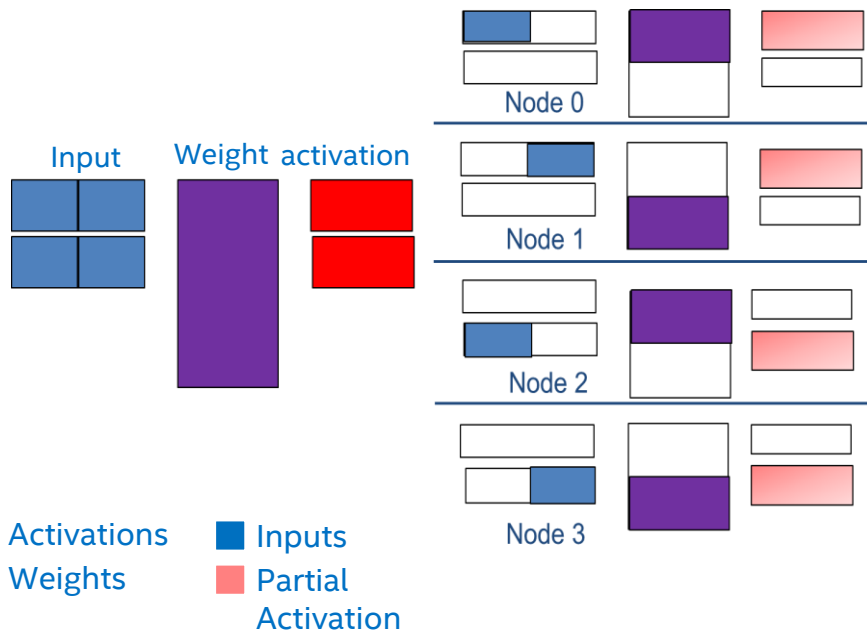*Output or
activations*

Hybrid

# Data/Model Parallelism at Scale

- General rule of thumb
  - Use data parallelism when activations > weights
  - Use model parallelism when weights > activations

- Implications of data and model parallelism
  - Data parallelism at scale makes activations << weights
  - Model parallelism at scale makes weights << activations
  - Compute efficiency goes down due to skewed matrices
  - Communication time dominates at scale
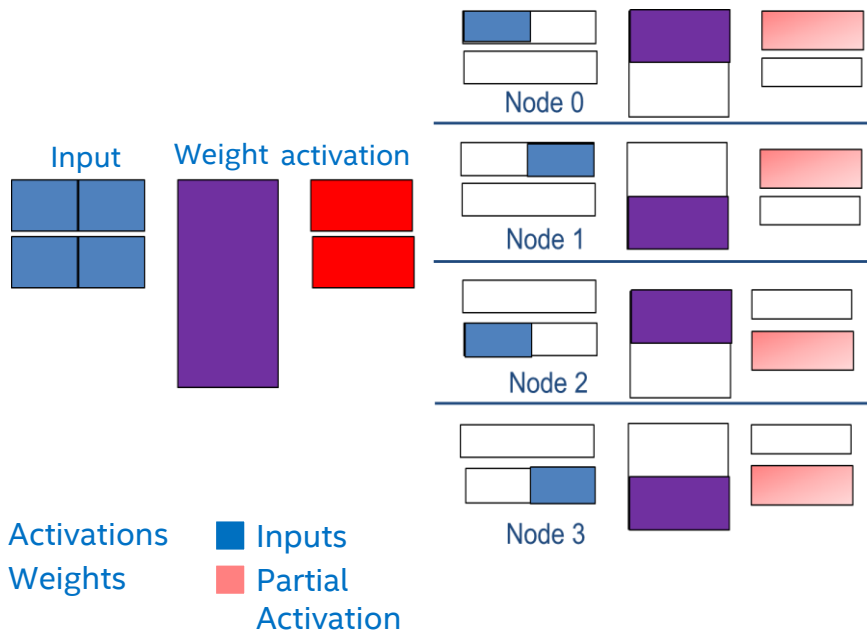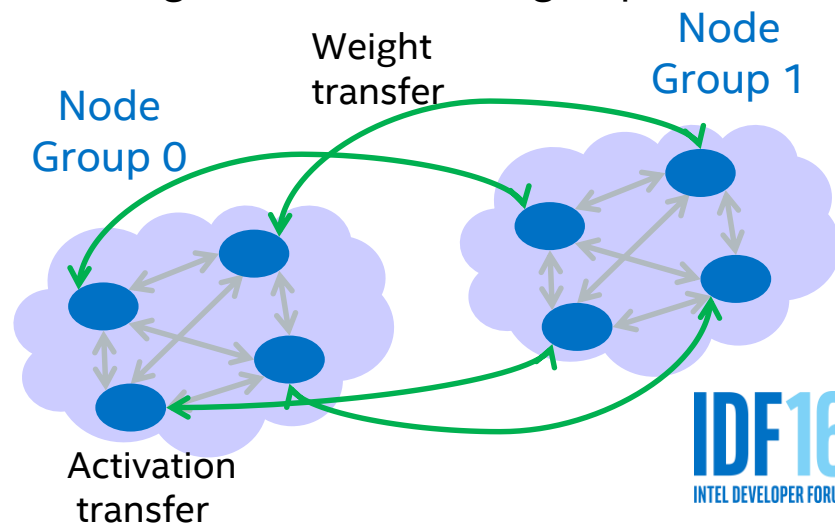
IDF16
INTEL DEVELOPER FORUM

# Addressing the scaling challenge

- Hybrid parallelism to improve compute efficiency

  - Partition across activations and weights to minimize skewed matrices
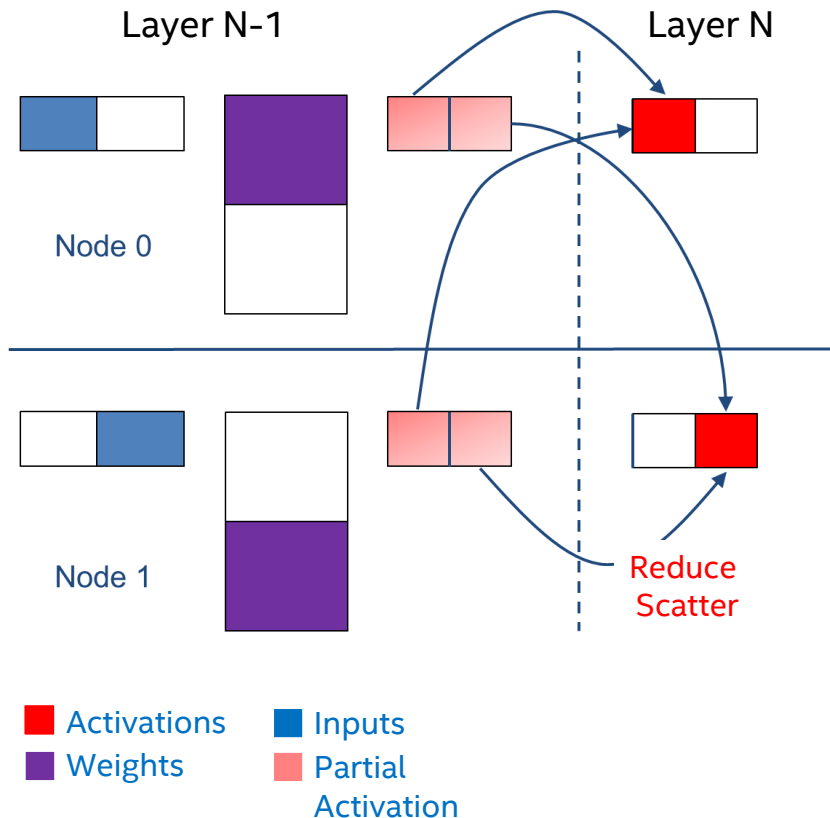
# Addressing the scaling challenge

- **Hybrid parallelism to improve compute efficiency**
  - Partition across activations and weights to minimize skewed matrices

- **Node groups to improve communication efficiency**
  - Avoid global transfer of activations and weights via node groups
    - Activations transfer within a group
    - Weight transfer across groups
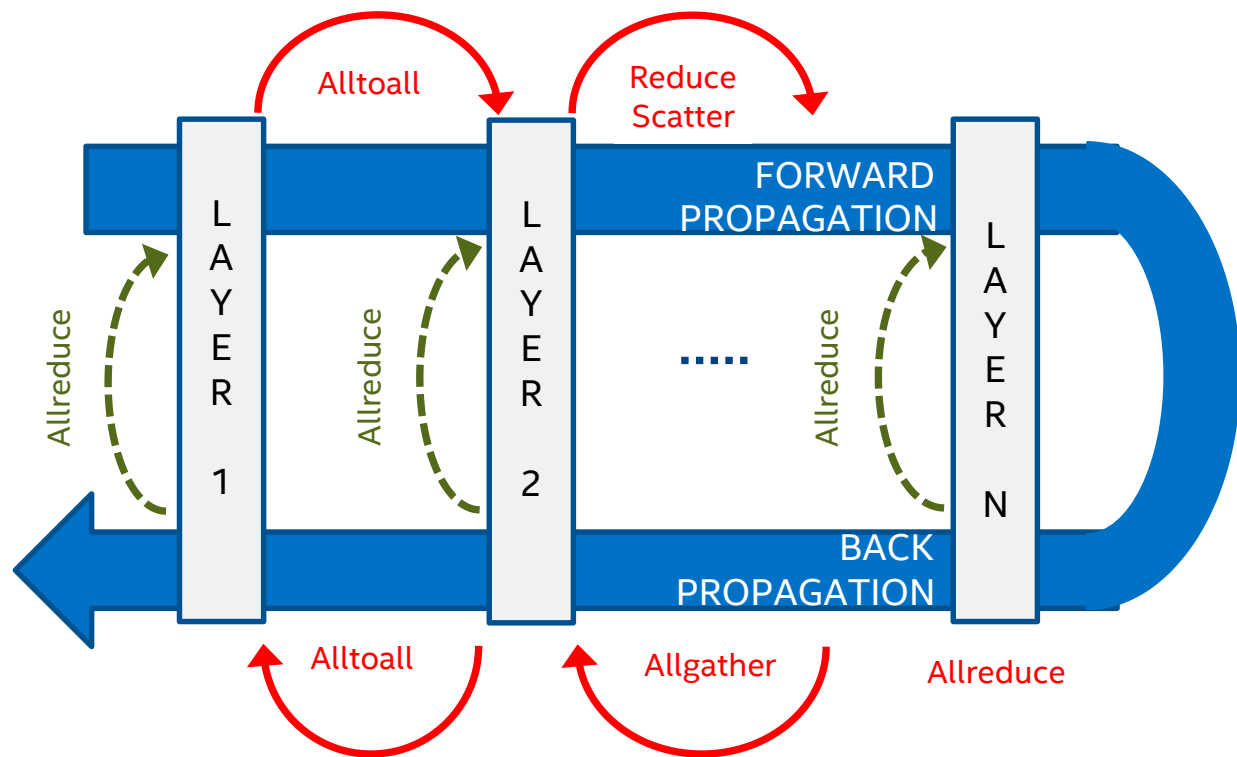
# Communication Patterns in Deep Learning



Layer N-1                    Layer N

Node 0

Node 1

Reduce Scatter

Reduce the activations from layer N-1 and scatter at layer N

Common MPI collectives in DL
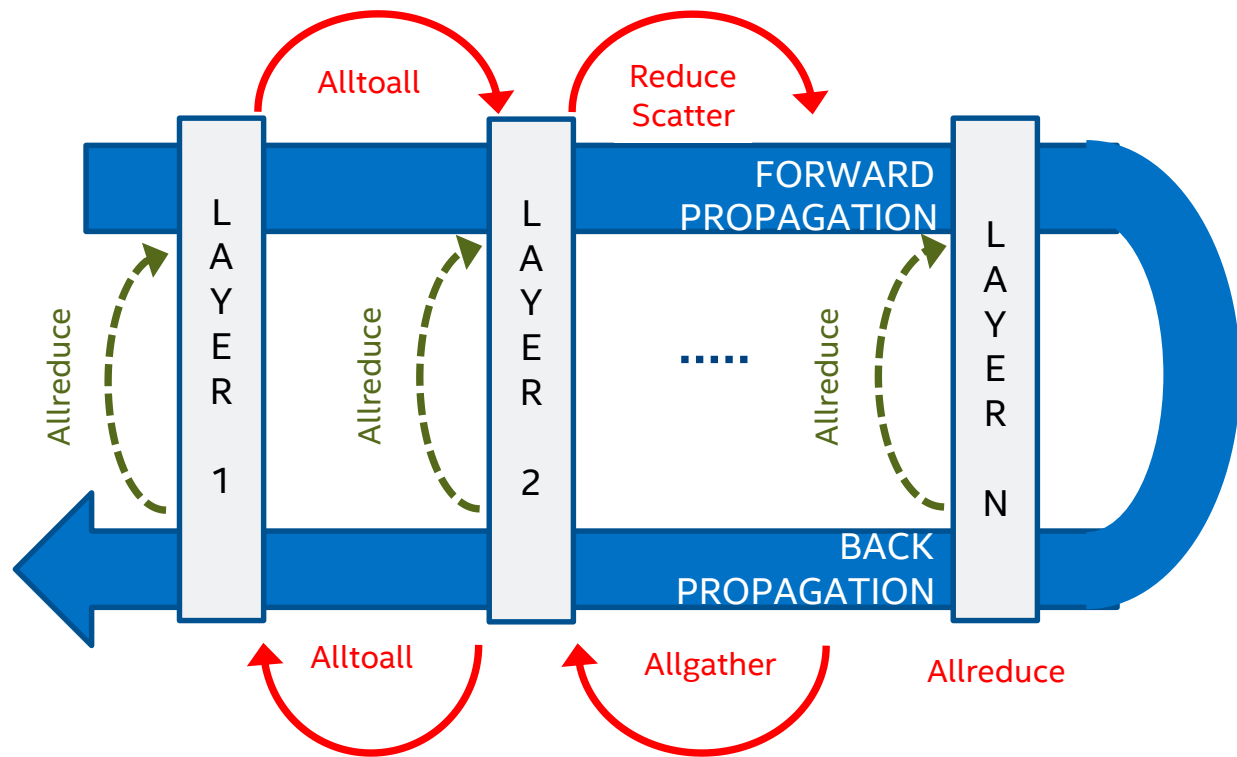• Reduce Scatter
• AllGather
• AllReduce
• AlltoAll

■ Activations   ■ Inputs
■ Weights       ■ Partial Activation

IDF16
INTEL DEVELOPER FORUM

# Communication Patterns in Deep Learning ... contd.



Activations (required immediately in next layer)

Updated weights (required during forward propagation of the corresponding layer)

# Communication Patterns in Deep Learning ... contd.

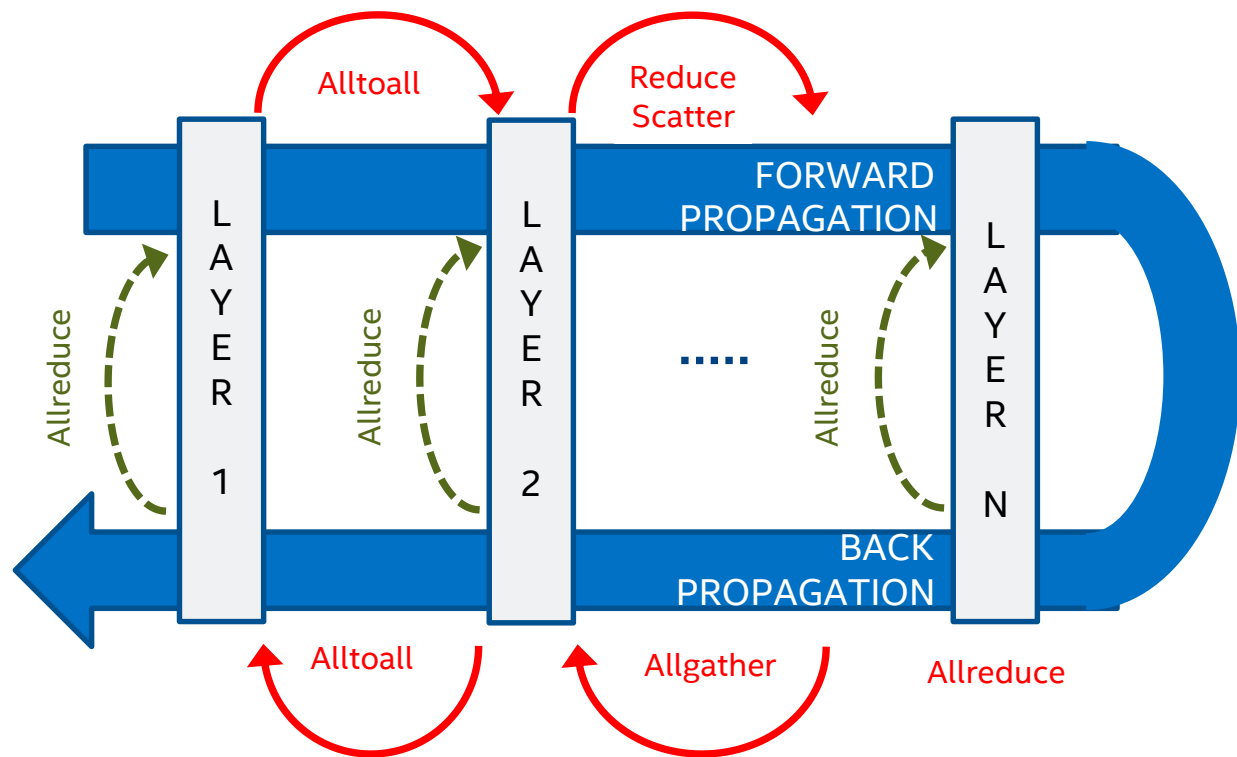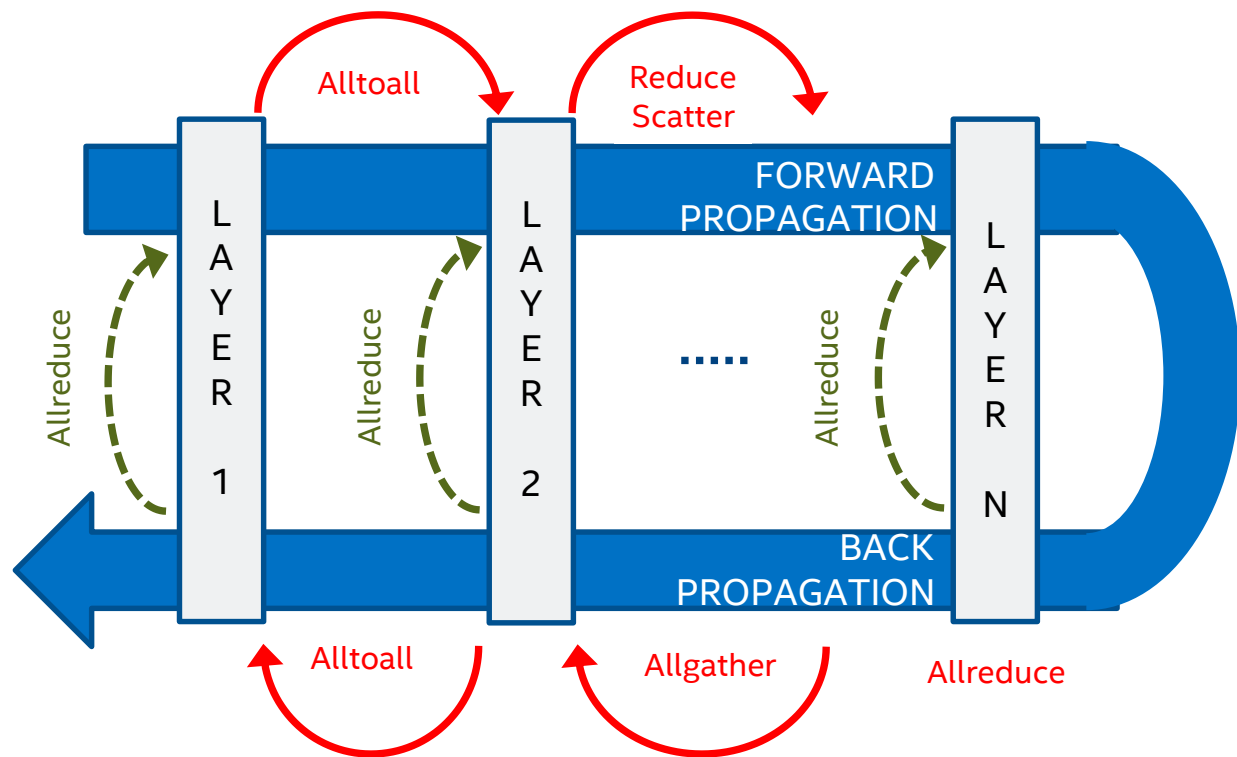

1. Optimized Collectives

FORWARD PROPAGATION

BACK PROPAGATION

Alltoall    Reduce Scatter

Alltoall    Allgather    Allreduce

Allreduce

LAYER 1   LAYER 2   .....   LAYER N

→ Activations (required immediately in next layer)

--→ Updated weights (required during forward propagation of the corresponding layer)

IDF16 INTEL DEVELOPER FORUM

# Communication Patterns in Deep Learning ... contd.



1. Optimized Collectives
2. Compute Communication Overlap

→ Activations (required immediately in next layer)

--→ Updated weights (required during forward propagation of the corresponding layer)
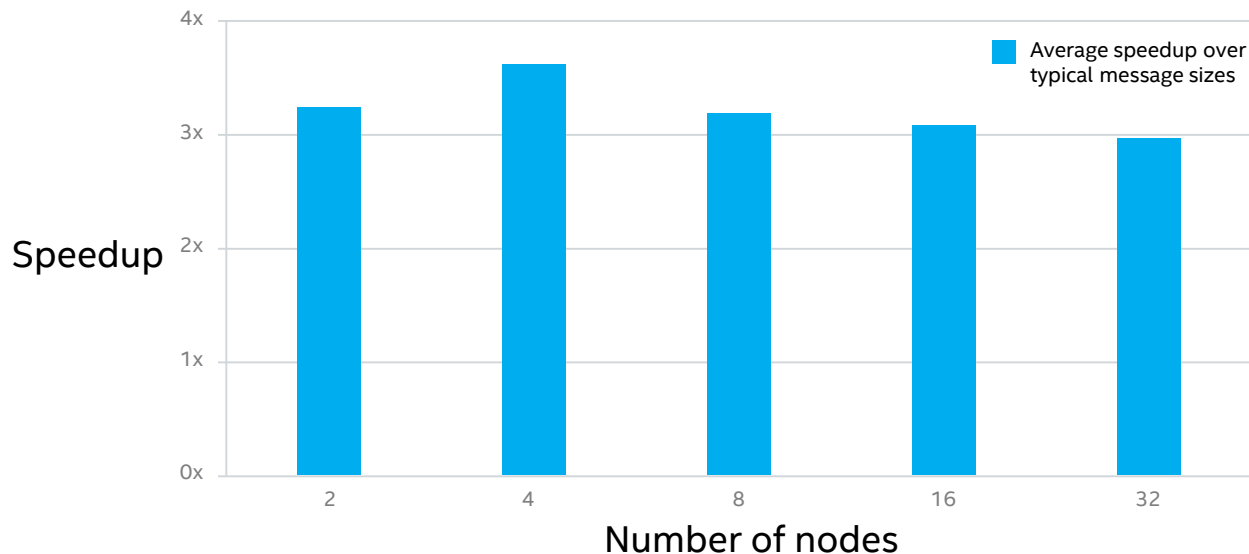
# Communication Patterns in Deep Learning ... contd.



1. Optimized Collectives
2. Compute Communication Overlap
3. Smart Message and Task Scheduling

→ Activations (required immediately in next layer)
-- → Updated weights (required during forward propagation of the corresponding layer)

IDF16
INTEL DEVELOPER FORUM

# Scaling Deep Learning Communication Primitives

## MPI ALLReduce performance on Intel® Xeon Phi™ Knights Landing



**Deep learning specific optimizations result in 3X speedup for the Allreduce collective**
**(average for the message profile of 16KB – 16MB floats)**

# Benefits of Multinode Optimizations



Overfeat-FAST: Scaling Efficiency

**Higher is better**

Overfeat-FAST: Compute Communication Breakdown

**Higher % of compute (green) is better**

Scaling efficiency without multi-node optimizations drops 1.3-2.1X for large node counts

IDF16 INTEL DEVELOPER FORUM

# Scaling Training Time of a common Neural Network

Speed up –
Higher is better



Bar chart values:
- 1 node: 1.0x
- 2 nodes: 2.0x
- 4 nodes: 3.9x
- 8 nodes: 7.3x
- 16 nodes: 15.3x
- 32 nodes: 26.7x

Number of nodes

Topology: **AlexNet***          Dataset: **Large image database**

With convergence and I/O overhead included

* https://github.com/soumith/convnet-benchmarks/blob/master/caffe/imagenet_winners/alexnet.prototxt

IDF16
INTEL DEVELOPER FORUM

# Scaling efficiently popular neural network topologies



Scaling Neural Network Training on
Intel® Xeon Phi™ Knights Landing

Dataset: **Large image database**    Without convergence and I/O overhead

IDF16
INTEL DEVELOPER FORUM

# The Virtuous Cycle of Compute



Training

Scoring

# The Virtuous Cycle of Compute

# Delivering Trained Model to Edge Device



Accuracy

Compression

Throughput
(Images/sec)

# Delivering Trained Model to Target Device †



Accuracy

1x (no accuracy loss)

10x

Compression

1.1x

Throughput
(Images/sec)

Sparsifying FC layers (e.g., Deep Compression*) -> mobile

# Delivering Trained Model to Target Device †



Accuracy

1x (no accuracy loss)

5.7x

1.7x

Compression

Throughput
(Images/sec)

Balanced sparsifying of Conv and FC layers → automotives

IDF16
INTEL DEVELOPER FORUM

# Delivering Trained Model to Target Device †



Accuracy

0.95x (5% accuracy loss)

9.0x

Compression

2.6x

Throughput
(Images/sec)

IDF16
INTEL DEVELOPER FORUM

# Agenda

- Why do we need to scale machine learning

- What makes it hard to scale and how we are addressing it

- Real-world applications

  - Introducing Dr. Amir Khosrowshahi

- Hardware roadmap, software tools and frameworks update

- Summary

**IDF16**
**INTEL DEVELOPER FORUM**

**Introducing Dr. Amir Khosrowshahi**

**CTO, Nervana Systems**

# About Nervana

A platform for machine intelligence

- Enable deep learning at scale
- Optimized from algorithms to silicon

# Application Areas

Healthcare

Agriculture

Finance

Online Services

Automotive

Energy

# Deep Learning as a Core Technology

# Neon: Nervana Python* Deep Learning Library

- User-friendly, extensible, abstracts parallelism

- Support for many deep learning models

- Interface to **Nervana** cloud

- Supports multiple backends

- Integrates with large cloud datastores

- Core routines written in assembler



See github for details

# Agenda

- Why do we need to scale machine learning

- What makes it hard to scale and how we are addressing it

- Real-world experience of an industry leader

- Hardware roadmap, software tools and frameworks update

- Summary

# Intel® Xeon Phi™ Processor Family for Performance

**Enables shorter time to train**

## Breakthrough Highly-Parallel Performance

- Up to ~6 SGEMM TFLOPs[1] per socket
- 1.38x[2] better scaling efficiency resulting in lower time to train for multi-node
- Eliminates add-in card PCIe* offload bottleneck and utilization constraints

## Removes Barriers through Integration

- Integrated Intel® Omni-Path fabric (dual-port; 50 GB/s) increases price-performance and reduces communication latency for deep learning networks

## Better Programmability

- Binary-compatible with Intel® Xeon® processors
- Open standards, libraries and frameworks

# Intel® Xeon Phi™ Processor Family for Performance

**Enables shorter time to train**

## Breakthrough Highly-Parallel Performance

- Up to ~6 SGEMM TFLOPs[1] per socket
- 1.38x[2] better scaling efficiency resulting in lower time to train for multi-node
- Eliminates add-in card PCIe* offload bottleneck and utilization constraints

## Removes Barriers through Integration

- Integrated Intel® Omni-Path fabric (dual-port; 50 GB/s) increases price-performance and reduces communication latency for deep learning networks

## Better Programmability

- Binary-compatible with Intel® Xeon® processors
- Open standards, libraries and frameworks

For an exciting new Intel® Xeon Phi™ roadmap update for machine learning/AI: Please attend Intel EVP Diane Bryant's Keynote tomorrow, Aug 17, 9am

IDF16
INTEL DEVELOPER FORUM

# Better performance in Deep Neural Network workloads with Intel® Math Kernel Library (**Intel**® **MKL**)

# Better performance in Deep Neural Network workloads with Intel® Math Kernel Library (**Intel® MKL**)



**Caffe*/AlexNet <u>single node</u> training performance**

Chart showing Performance speedup:
- Intel® Xeon® E5-2699 v4, Out-of-the-box: ~1x
- 5.8x →
- Intel® Xeon® E5-2699 v4, +Intel MKL 11.3.3: ~5.8
- 2.1x →
- Intel® Xeon® E5-2699 v4, +Intel MKL 2017: ~12
- 2x →
- Intel® Xeon Phi™ 7250, +Intel MKL 2017: **24x**

IDF16
INTEL DEVELOPER FORUM

# Better performance in Deep Neural Network workloads with Intel® Math Kernel Library (Intel® MKL)



Caffe*/AlexNet single node inference performance

IDF16
INTEL DEVELOPER FORUM

# Intel Deep Learning Software Stack and Timeline

**Intel® Math Kernel Library (Intel® MKL)**

Xeon

Xeon Phi

**Intel MKL** is SW building block to extract max Intel HW performance and provide common interface to all Intel accelerators.

# Intel Deep Learning Software Stack and Timeline

**Intel Math Kernel Library (Intel® MKL)**

**Intel® MKL-DNN**

Xeon    Xeon Phi    FPGA

**Intel MKL-DNN** is an open source IA optimized DNN APIs, combined with Intel® MKL and build tools designed for scalable, high-velocity integration with ML/DL frameworks.

**Targeted release: Q3 2016**

Includes:
- Open Source implementations of new DNN functionality included in MKL 2017 Beta, new algorithms ahead of MKL releases
- IA optimizations contributed by community

**Intel MKL** is SW building block to extract max Intel HW performance and provide common interface to all Intel accelerators.

IDF16
INTEL DEVELOPER FORUM

# Intel Deep Learning Software Stack and Timeline

## Deep Learning Frameworks

theano · Caffe BVLC · TensorFlow · Microsoft CNTK · Google · torch

**Intel® Math Kernel Library (Intel® MKL)**

**Intel® MKL-DNN**

Xeon · Xeon Phi · FPGA

Popular Deep Learning frameworks

**Intel MKL-DNN** is an open source IA optimized DNN APIs, combined with Intel® MKL and build tools designed for scalable, high-velocity integration with ML/DL frameworks.

**Targeted release: Q3 2016**

Includes:
- Open Source implementations of new DNN functionality included in MKL 2017 Beta, new algorithms ahead of MKL releases
- IA optimizations contributed by community

**Intel MKL** is SW building block to extract max Intel HW performance and provide common interface to all Intel accelerators.

- **Multi-Node scaling for Knights Landing: Caffe\* by EoY and 1H'17 in other frameworks**

**IDF16**
**INTEL DEVELOPER FORUM**

# Intel Deep Learning Software Stack and Timeline

**Intel Deep Learning Tools**

Tools to accelerate design, training and deployment of deep learning solutions
**Targeted release: Q3'2016**

**Deep Learning Frameworks**

theano · Caffe BVLC · TensorFlow · Microsoft CNTK · Google · torch

Popular Deep Learning frameworks

**Intel® Math Kernel Library (Intel® MKL)**

**Intel® MKL-DNN**

**Intel MKL-DNN** is an open source IA optimized DNN APIs, combined with Intel® MKL and build tools designed for scalable, high-velocity integration with ML/DL frameworks.
**Targeted release: Q3 2016**

Includes:
- Open Source implementations of new DNN functionality included in MKL 2017 Beta, new algorithms ahead of MKL releases
- IA optimizations contributed by community

**Intel MKL** is SW building block to extract max Intel HW performance and provide common interface to all Intel accelerators.

Xeon · Xeon Phi · FPGA

- **Multi-Node scaling for Knights Landing: Caffe\* by EoY and 1H'17 in other frameworks**
- **Intel Deep Learning Tools with support for model compression by end of 2016**

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

IDF16
INTEL DEVELOPER FORUM

35

# Call to action

- Machine learning is the key enabler for a new virtuous cycle of compute triggered by explosion of digital data and ubiquitous connectivity

- It can vastly expand the reach of computing for applications like self-driving, agriculture, health and manufacturing

- Help machine learning unlock the true potential of AI

- Consider the full, end-to-end pipeline when you think about your AI needs.

- Try Intel MKL, Intel optimized frameworks & Intel Xeon Phi

# Summary

- Machine learning must scale out to bring down the training time of weeks/days to days/hours

- Machine learning compute infrastructure must be both performant & productive for developers, and leverage the efficiency of cloud

- Scaling distributed machine learning is challenging as it pushes the limits of available data/model parallelism and internode communication

- Intel's new deep learning tools -- with the upcoming integration of Nervana cloud stack -- are designed to hide/reduce the complexity of strong scaling time-to-train and model deployment tradeoffs on resource-constrained edge devices without compromising the performance need

IDF16
INTEL DEVELOPER FORUM

# Related Tech Sessions

ANATS01: Deep Learning Frameworks and Optimization Paths on Intel® Architecture
By Andres Rodriguez, Panchumarthy, Ravi, and Tom "Elvis" Jones, Amazon

ANATS03: Enabling an End to End Architecture for Autonomous Vehicles
By Jack Weast

ANATS05: How to Parallelize Neural Networks (xNNs) for Intel® Xeon Phi™
By Nadathur R. Satish

For more information on machine learning at Intel: intel.com/machinelearning

A PDF of this presentation is available is available from our Technical Session Catalog: www.intel.com/idfsessionsSF.  This URL is also printed on the top of Session Agenda Pages in the Pocket Guide.

IDF16
INTEL DEVELOPER FORUM

# Legal Notices and Disclaimers

IDF16
INTEL DEVELOPER FORUM

# 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 SSE3 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

# Risk Factors

The above statements and any others in this document that refer to plans and expectations for the second quarter, the year and the future are forward-looking statements that involve a number of risks and uncertainties. Words such as "anticipates," "expects," "intends," "plans," "believes," "seeks," "estimates," "may," "will," "should" and their variations identify forward-looking statements. Statements that refer to or are based on projections, uncertain events or assumptions also identify forward-looking statements. Many factors could affect Intel's actual results, and variances from Intel's current expectations regarding such factors could cause actual results to differ materially from those expressed in these forward-looking statements. Intel presently considers the following to be important factors that could cause actual results to differ materially from the company's expectations. Demand for Intel's products is highly variable and could differ from expectations due to factors including changes in business and economic conditions; consumer confidence or income levels; the introduction, availability and market acceptance of Intel's products, products used together with Intel products and competitors' products; competitive and pricing pressures, including actions taken by competitors; supply constraints and other disruptions affecting customers; changes in customer order patterns including order cancellations; and changes in the level of inventory at customers. Intel's gross margin percentage could vary significantly from expectations based on capacity utilization; variations in inventory valuation, including variations related to the timing of qualifying products for sale; changes in revenue levels; segment product mix; the timing and execution of the manufacturing ramp and associated costs; excess or obsolete inventory; changes in unit costs; defects or disruptions in the supply of materials or resources; and product manufacturing quality/yields. Variations in gross margin may also be caused by the timing of Intel product introductions and related expenses, including marketing expenses, and Intel's ability to respond quickly to technological developments and to introduce new products or incorporate new features into existing products, which may result in restructuring and asset impairment charges. Intel's results could be affected by adverse economic, social, political and physical/infrastructure conditions in countries where Intel, its customers or its suppliers operate, including military conflict and other security risks, natural disasters, infrastructure disruptions, health concerns and fluctuations in currency exchange rates. Results may also be affected by the formal or informal imposition by countries of new or revised export and/or import and doing-business regulations, which could be changed without prior notice. Intel operates in highly competitive industries and its operations have high costs that are either fixed or difficult to reduce in the short term. The amount, timing and execution of Intel's stock repurchase program could be affected by changes in Intel's priorities for the use of cash, such as operational spending, capital spending, acquisitions, and as a result of changes to Intel's cash flows or changes in tax laws. Product defects or errata (deviations from published specifications) may adversely impact our expenses, revenues and reputation. Intel's results could be affected by litigation or regulatory matters involving intellectual property, stockholder, consumer, antitrust, disclosure and other issues. An unfavorable ruling could include monetary damages or an injunction prohibiting Intel from manufacturing or selling one or more products, precluding particular business practices, impacting Intel's ability to design its products, or requiring other remedies such as compulsory licensing of intellectual property. Intel's results may be affected by the timing of closing of acquisitions, divestitures and other significant transactions. A detailed discussion of these and other factors that could affect Intel's results is included in Intel's SEC filings, including the company's most recent reports on Form 10-Q, Form 10-K and earnings release.

Rev. 4/14/15

IDF 16
INTEL DEVELOPER FORUM