



DISTRIBUTED DL/ML SOLUTIONS FOR HPC SYSTEMS

Dr. Fabio Baruffa

Sr. Technical Consulting Engineer, Intel IAGS

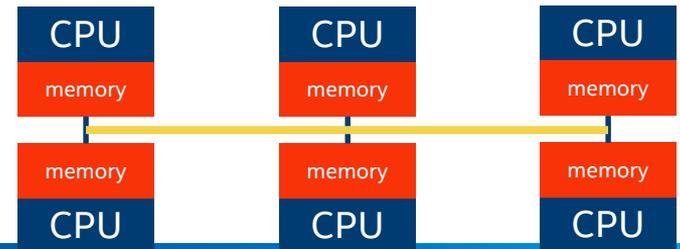
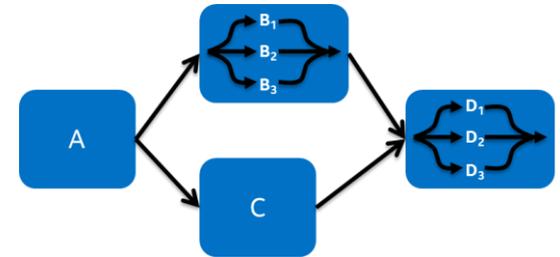
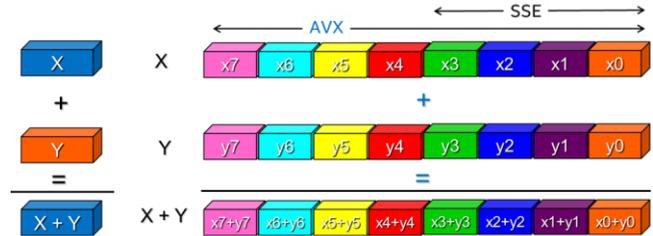




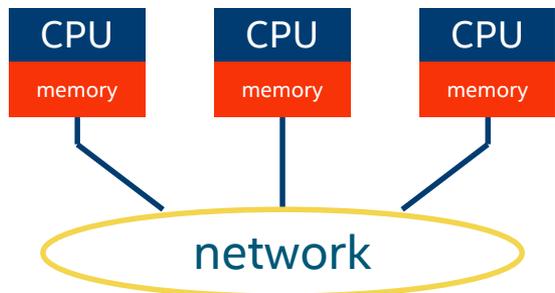
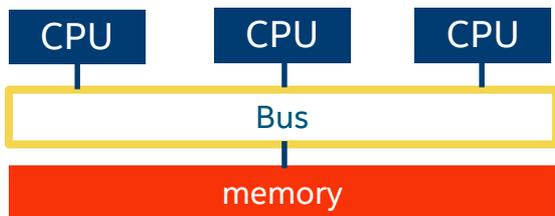
BASIC CONCEPTS ON DISTRIBUTED COMPUTING

TYPES OF PARALLELISM

- **SIMD**: Single instruction multiple data (Data Parallel)
 - The same instruction is simultaneously applied on multiple data items
- **MIMD**: Multiple instructions multiple data (Task Parallel)
 - Different instructions on different data
- **SPMD**: Single program multiple data (MPI Parallel)
 - This is the message passing programming on distributed systems



SHARED VS DISTRIBUTED MEMORY SYSTEM



- **Shared memory**

- There is a unique address space shared between the processors
- All the processors can access the same memory

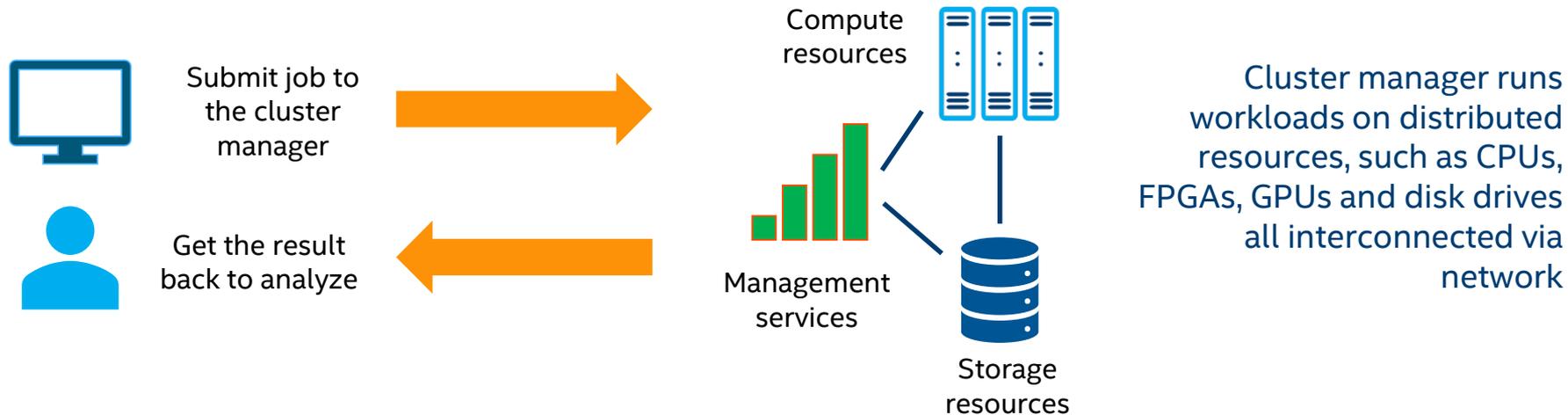
- **Distributed memory**

- Each processor has its own local memory
- Messages are exchanged between the processors to communicate the data

WHAT IS HIGH-PERFORMANCE COMPUTING (HPC)?

Leveraging distributed compute resources to solve complex problems with large datasets

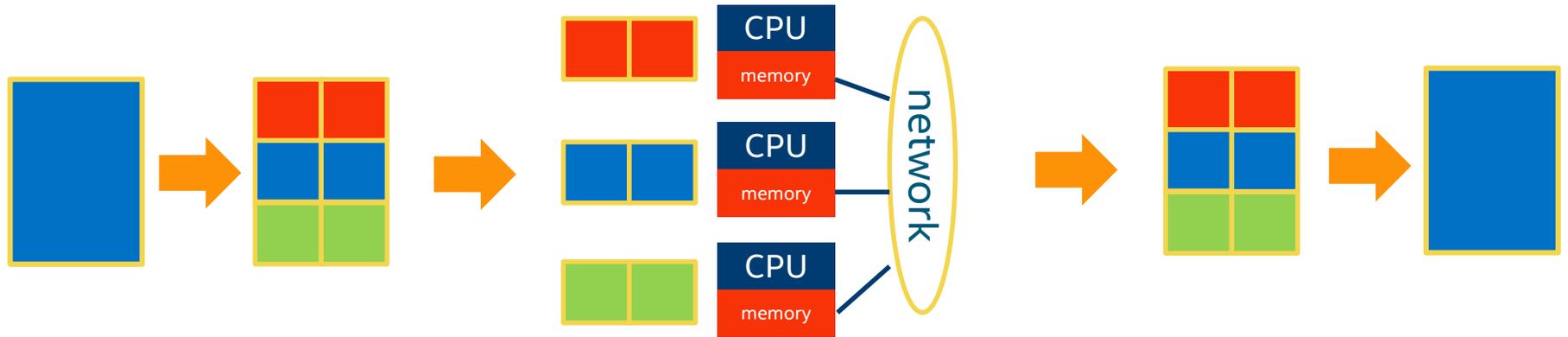
- Terabytes to petabytes to zettabytes of data
- Results in minutes to hours instead of days or weeks



DOMAIN DECOMPOSITION METHOD FOR HPC

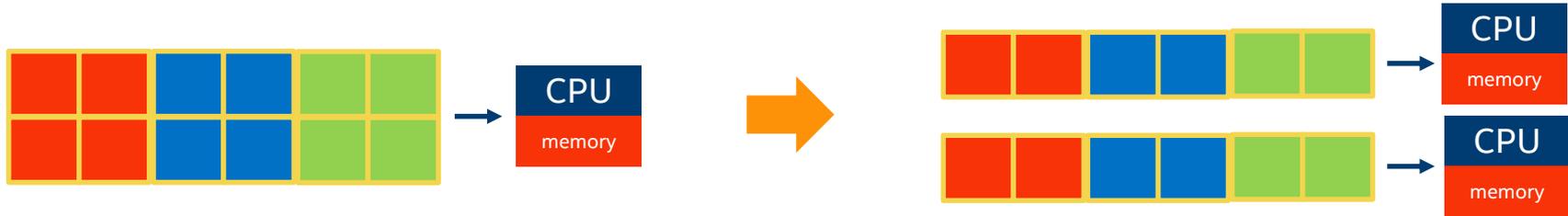
The domain decomposition is a technique for dividing a computational problem in several parts (domains) allowing to solve a large problem on the available resources

- *Partition* the data, assign them to each resource and associate the computation
- *Communication* happens to eventually exchange intermediate results
- *Aggregate* the results from the different resources

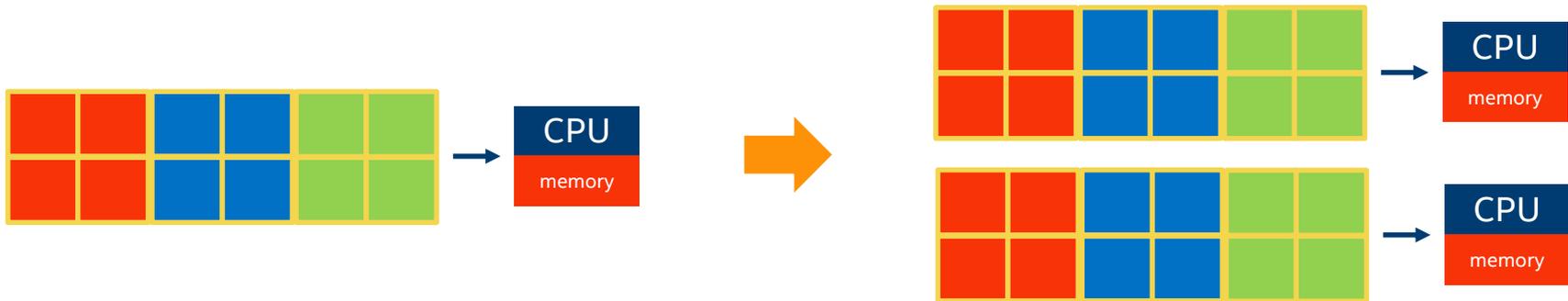


SCALING ASPECTS OF DISTRIBUTED COMPUTING

- Strong scaling: how the time to solution changes by increasing the compute resources for a fixed *total* problem size

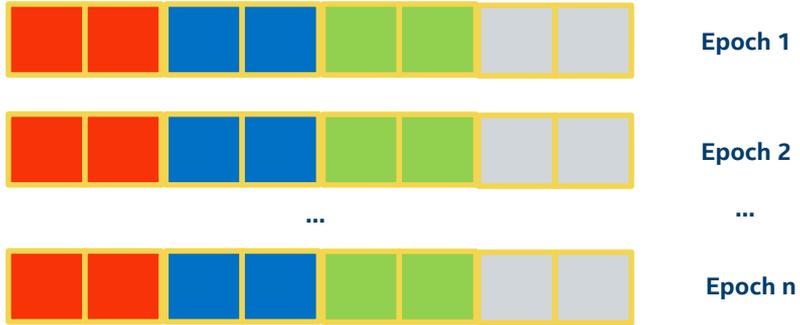


- Weak scaling: how the time to solution changes by increasing the compute resource for a constant problem size *per process*



HOW DO WE REDUCE THE COMPUTATIONAL TIME?

Number of training data set = 8



We could use a strong scaling approach to reduce the time for all the epochs

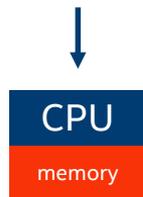
STRONG SCALING ON TRAINING SET

Number of training data set = 8



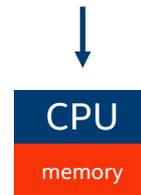
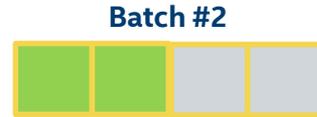
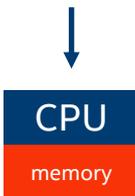
We divide the dataset of 8 training samples into 2 batches of Batch size 4

Model will be updated after each batch of 4 samples



Batch size = 4

Going beyond reducing the Batch size and increasing the #CPUS (strong scaling) can cause a loose of performance in the model

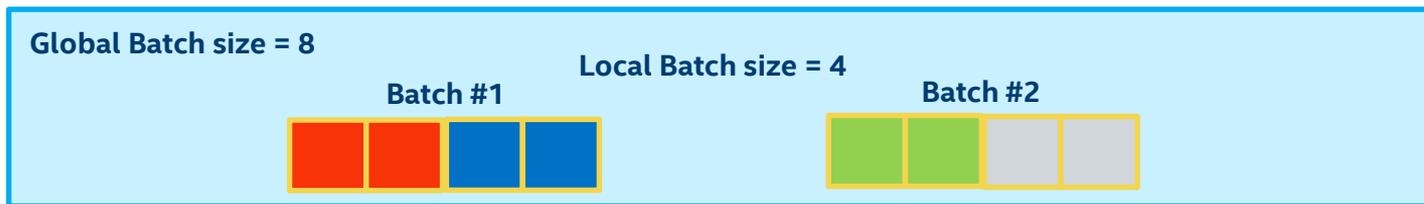


STRONG AND WEAK SCALING ON TRAINING SET

Number of training data set = 8



We keep the same Batch size/CPU, increasing the overall Batch size



We update the model after the Global batch size is reached, reducing the number of iterations per epoch

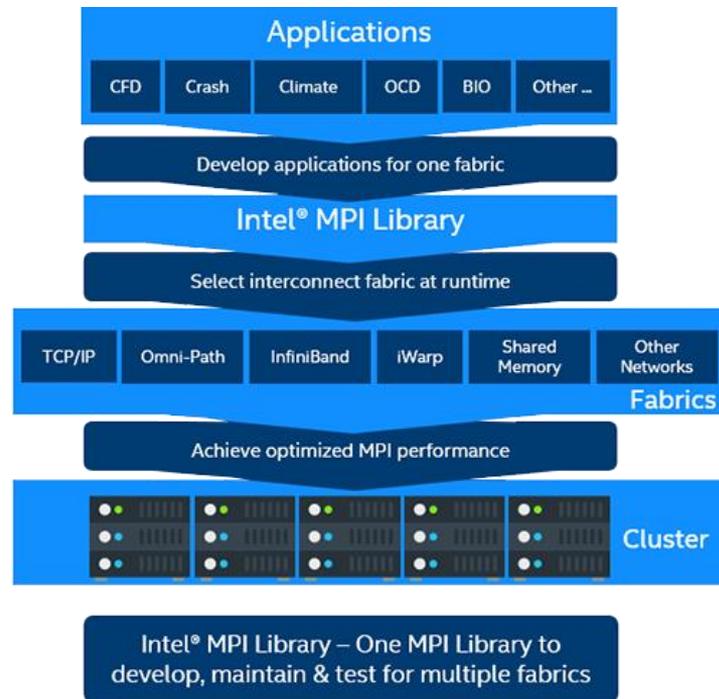


MESSAGE PASSING INTERFACE (MPI)

MPI is a standard which gets implemented in form of libraries for inter-process communication and data exchange.

Function categories:

- Point-to-point communication
- Collective communication
- Communicator topologies
- User-defined data types
- Utilities (for example, timing and initialization)





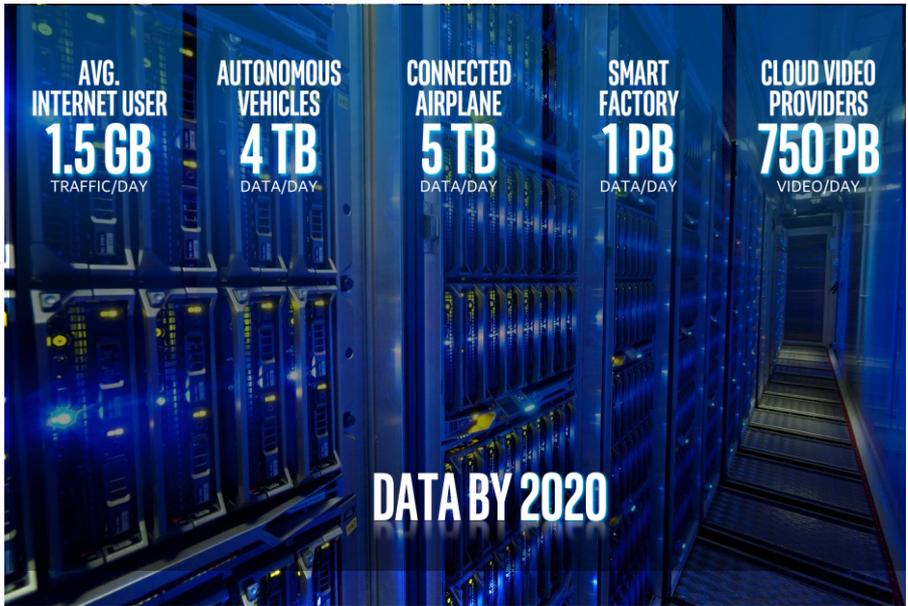
DISTRIBUTING STRATEGY FOR MACHINE LEARNING

FROM PROTOTYPE TO PRODUCTION

```
In [11]:
sns.set()
sns.pairplot(df_train[cols], size = 2.5)
plt.show();
```



PERFORMANCE →



<https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python>

Optimization Notice

Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.



WHY DISTRIBUTED ML/DL

- Most Machine Learning tasks assume the data can be easily accessible, but:
 - Data loading on a single machine can be a bottleneck in case of large amount of data
 - To run production applications large memory systems is required (data not fitting in the local computer RAM)
 - Traditional sequential algorithms are not suitable in case of distributed memory system
- Time to solution is critical on highly competitive market.

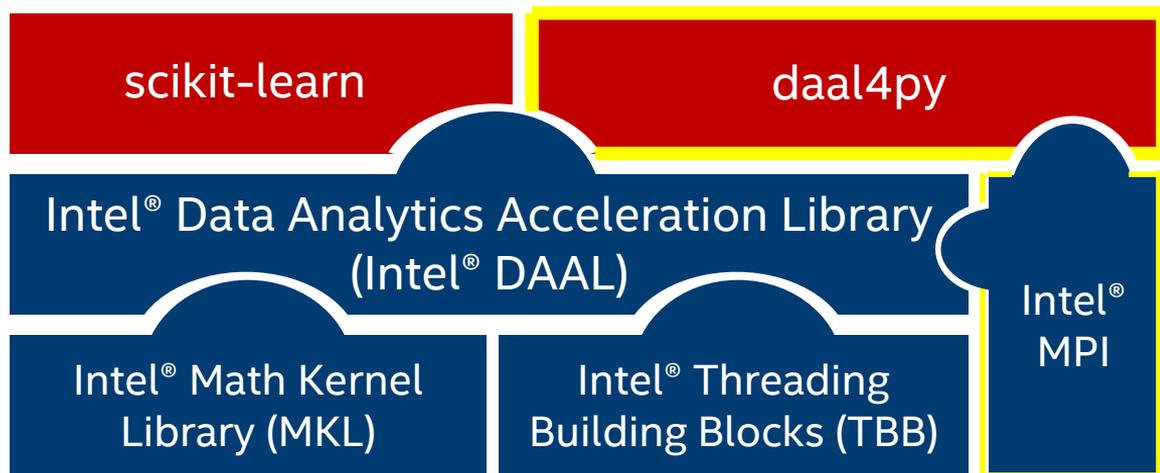
WHY DISTRIBUTED ML/DL

- Deep Learning training takes time:
 - Computational complexity of DL training can be up to 100+ ExaFLOP (1 ExaFLOP = 10^{18} op);
 - Typical single node performance is up-to tens of TeraFLOPS (1 TF = 10^{12} op/sec);
 - Peak performance of most powerful HPC clusters is up-to tens of PetaFLOPS (1 PF = 10^{15} op/sec).
- **Time to solution is critical on highly competitive market.**

DAAL4PY: ACCELERATED ANALYTICS TOOLS

- Package created to address the needs of **Data Scientists and Framework Designers** to harness the **Intel® Data Analytics Acceleration Library (DAAL)** with a **Pythonic API**
- For scaling capabilities, **daal4py** also provides the ability to do distributed machine learning using **Intel® MPI library**
- **daal4py** operates in SPMD style (Single Program Multiple Data), which means your program is executed on several processes (e.g. similar to MPI)
- The use of MPI is not required for **daal4py**'s SPMD-mode to work, all necessary communication and synchronization happens under the hood of daal4py
- It is possible to use **daal4py** and **mpi4py** in the same program

SCALING MACHINE LEARNING BEYOND A SINGLE NODE



Simple Python API
Powers scikit-learn

Powered by DAAL

Scalable to multiple nodes

Try it out! `conda install -c intel daal4py`

K-MEANS USING DAAL4PY

```
import daal4py as d4p

# daal4py accepts data as CSV files, numpy arrays or pandas dataframes
# here we let daal4py load process-local data from csv files
data = "kmeans_dense.csv"

# Create algob object to compute initial centers
init = d4p.kmeans_init(10, method="plusPlusDense")
# compute initial centers
ires = init.compute(data)
# results can have multiple attributes, we need centroids
Centroids = ires.centroids
# compute initial centroids & kmeans clustering
result = d4p.kmeans(10).compute(data, centroids)
```

DISTRIBUTED K-MEANS USING DAAL4PY

```
import daal4py as d4p

# initialize distributed execution environment
d4p.daalinit()

# daal4py accepts data as CSV files, numpy arrays or pandas dataframes
# here we let daal4py load process-local data from csv files
data = "kmeans_dense_{}.csv".format(d4p.my_procid())

# compute initial centroids & kmeans clustering
init = d4p.kmeans_init(10, method="plusPlusDense", distributed=True)
centroids = init.compute(data).centroids
result = d4p.kmeans(10, distributed=True).compute(data, centroids)
```

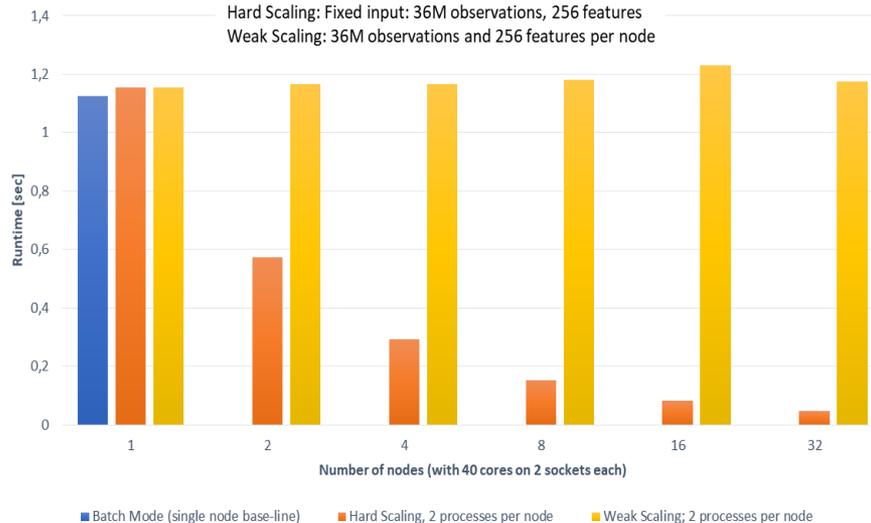
```
mpirun -n 4 python ./kmeans.py
```

STRONG & WEAK SCALING VIA DAAL4PY

Hardware	Intel(R) Xeon(R) Gold 6148 CPU @ 2.40GHz, EIST/Turbo on
	2 sockets, 20 Cores per socket
	192 GB RAM
	16 nodes connected with Infiniband
Operating System	Oracle Linux Server release 7.4
Data Type	double

daal4py Linear Regression Distributed Scalability

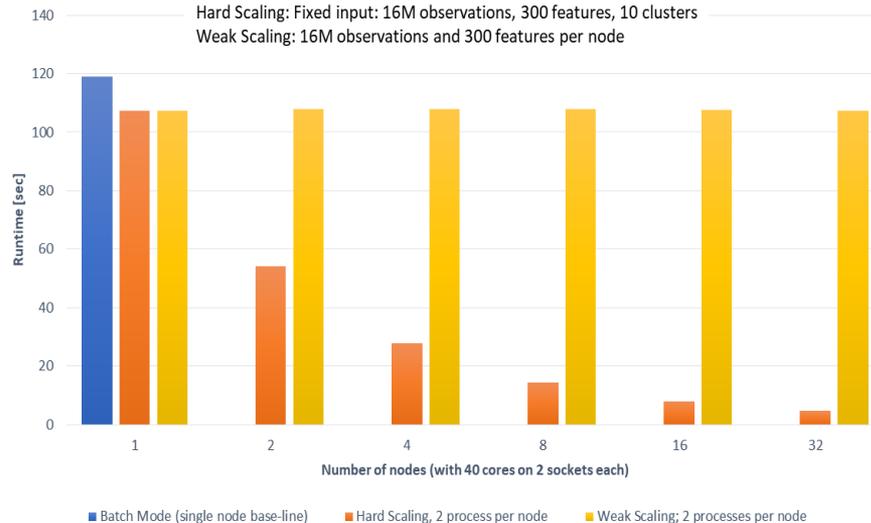
Hard Scaling: Fixed input: 36M observations, 256 features
Weak Scaling: 36M observations and 256 features per node



On a 32-node cluster (1280 cores) daal4py computed linear regression of 2.15 TB of data in 1.18 seconds and 68.66 GB of data in less than 48 milliseconds.

daal4py K-Means Distributed Scalability

Hard Scaling: Fixed input: 16M observations, 300 features, 10 clusters
Weak Scaling: 16M observations and 300 features per node



On a 32-node cluster (1280 cores) daal4py computed K-Means (10 clusters) of 1.12 TB of data in 107.4 seconds and 35.76 GB of data in 4.8 seconds.

Optimization Notice

Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.



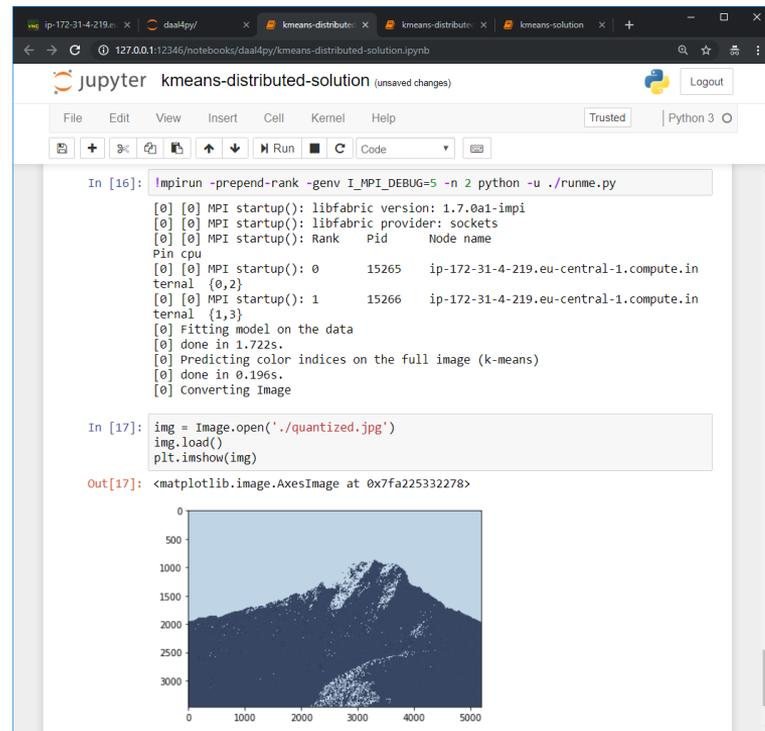
HANDS-ON

DISTRIBUTED K-MEANS USING DAAL4PY

- 1) Performs a pixel-wise Vector Quantization (VQ) using K-Means
- 2) Implemented the domain decomposition according to:
 - `d4p.num_procs()`
 - `d4p.my_procid()`
- 3) Using the distributed algorithm from Daal4Py
 - `d4p.kmeans_init(n_colors, method="plusPlusDense", distributed=True)`
- 4) What is the meaning of `d4p.daalinit()` & `d4p.daalfini()`?
- 5) How does threading compare to multiprocessing in terms of performance?

DISTRIBUTED K-MEANS SUMMARY

- Each process (MPI rank) get's a different chunk of data
- Only process #0 reports results
- Inference is using the same routines as training with 0 maximum iterations and centroid assignment
- There is no oversubscription since DAAL only sees the cores “owned” by the corresponding MPI rank

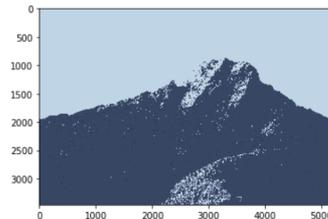


```
In [16]: mpirun -prepend-rank -genv I_MPI_DEBUG=5 -n 2 python -u ./runme.py

[0] [0] MPI startup(): libfabric version: 1.7.0a1-impi
[0] [0] MPI startup(): libfabric provider: sockets
[0] [0] MPI startup(): Rank   Pid   Node name
Pin  cpu
[0] [0] MPI startup(): 0     15265 ip-172-31-4-219.eu-central-1.compute.in
terminal {0,2}
[0] [0] MPI startup(): 1     15266 ip-172-31-4-219.eu-central-1.compute.in
terminal {1,3}
[0] Fitting model on the data
[0] done in 1.7225s.
[0] Predicting color indices on the full image (k-means)
[0] done in 0.1965s.
[0] Converting Image

In [17]: img = Image.open('./quantized.jpg')
img.load()
plt.imshow(img)

Out[17]: <matplotlib.image.AxesImage at 0x7fa225332278>
```

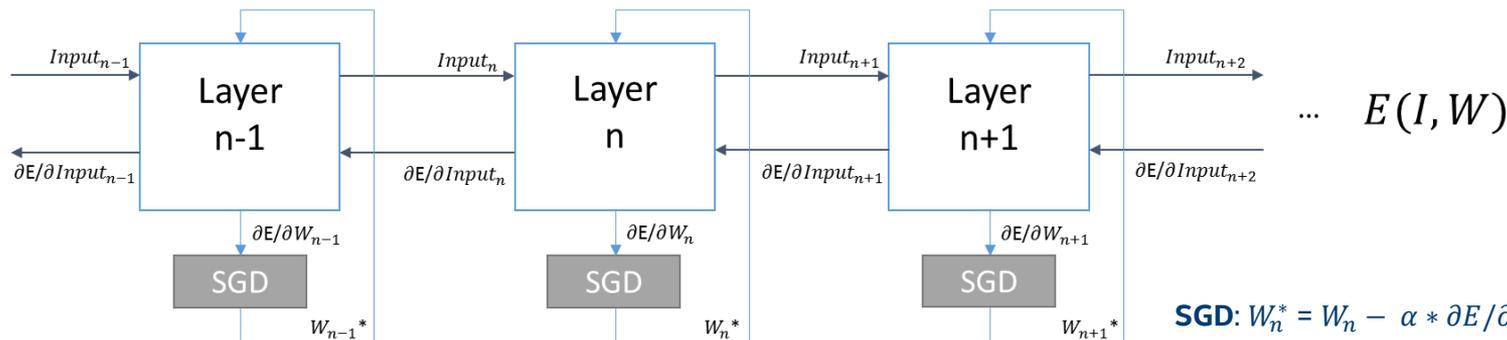




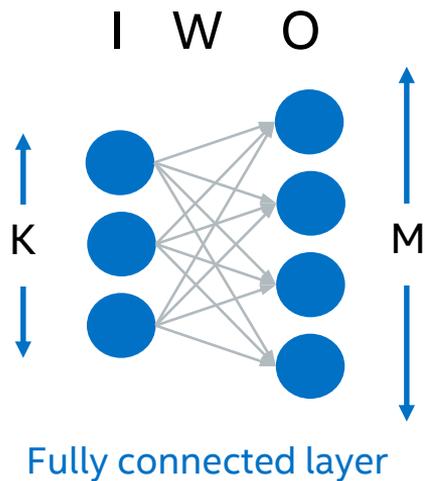
DISTRIBUTING STRATEGY FOR DEEP LEARNING

DEEP LEARNING TRAINING PROCEDURE

- **Forward propagation:** calculate loss function based on the input batch and current weights;
- **Backward propagation:** calculate error gradients w.r.t. weights for all layers (using chain rule);
- **Weights update:** use gradients to update weights; there are different algorithms exist - vanilla SGD, Momentum, Adam, etc.



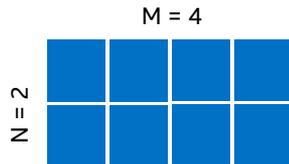
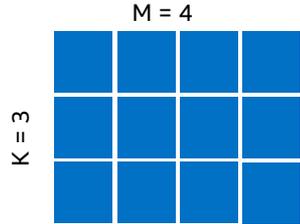
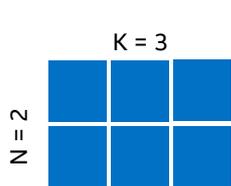
PARALLELISM OPTIONS



$I \in \mathbb{R}^{N \times K}$
Input

$W \in \mathbb{R}^{K \times M}$
Weights
or model

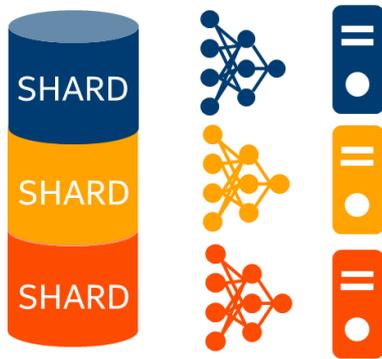
$O \in \mathbb{R}^{N \times M}$
Output
or activations



Several options for parallelization

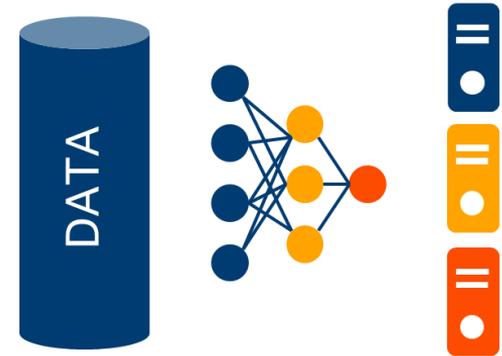
NEURAL NETWORK PARALLELISM

DATA PARALLELISM



Data is processed in increments of N .
Work on minibatch samples and distributed among the available resources.

MODEL PARALLELISM



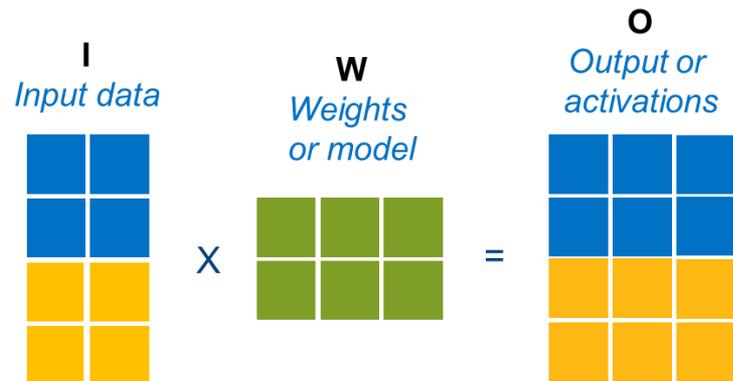
The work is divided according to the neurons in each layer. The sample minibatch is copied to all processors which compute part of the DNN.

source: <https://arxiv.org/pdf/1802.09941.pdf>

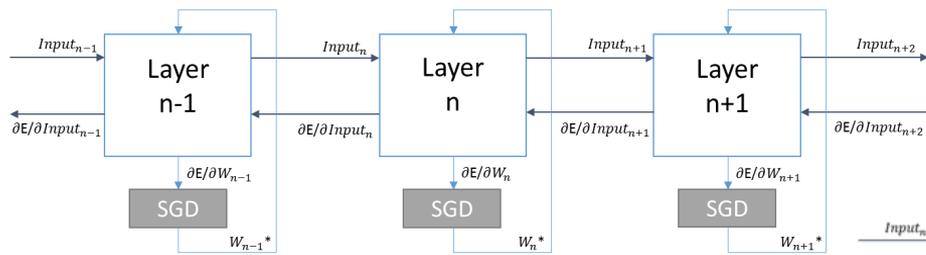
MULTI-NODE PARALLELIZATION

- Data parallelism:

- Replicate the model across nodes;
- Feed each node with its own batch of input data;
- Communication for gradients is required to get their average across nodes;
- Can be either
 - *AllReduce* pattern
 - *ReduceScatter* + *AllGather* patterns

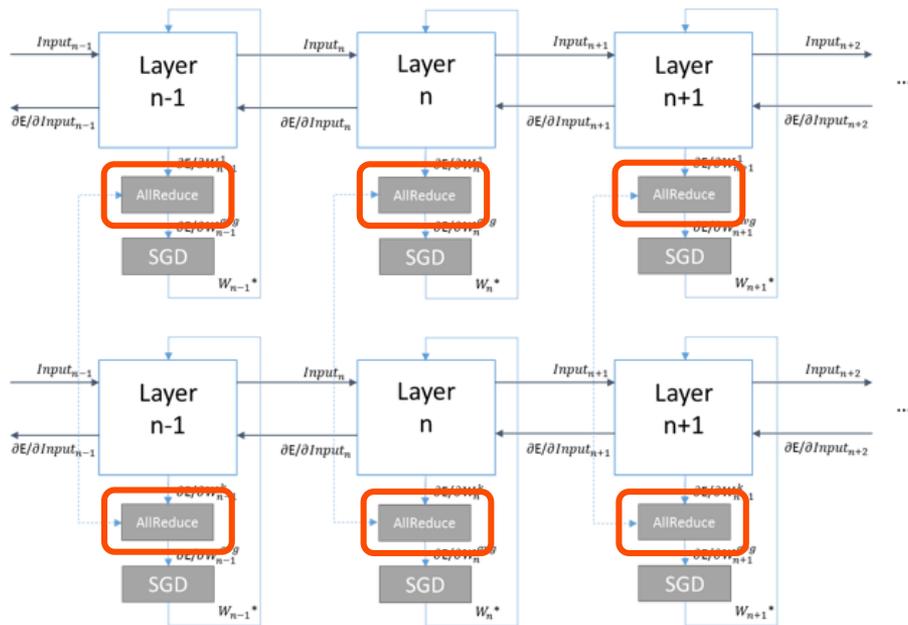


DATA PARALLELISM



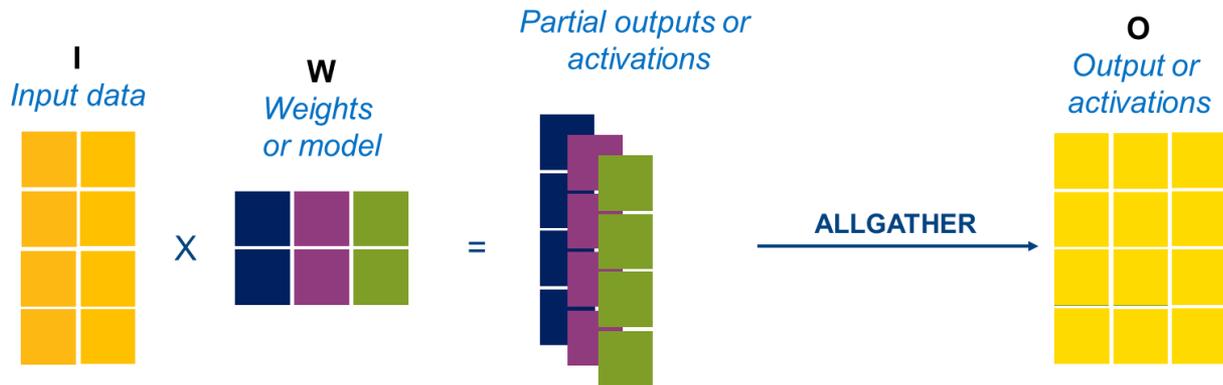
... $E(I, W)$ ← Single node

Multi node
Data Parallelism



MULTI-NODE PARALLELIZATION

- Model parallelism:
 - Model is split across nodes;
 - Feed each node with the same batch of input data;
 - Communication for partial activations is required to gather the result;



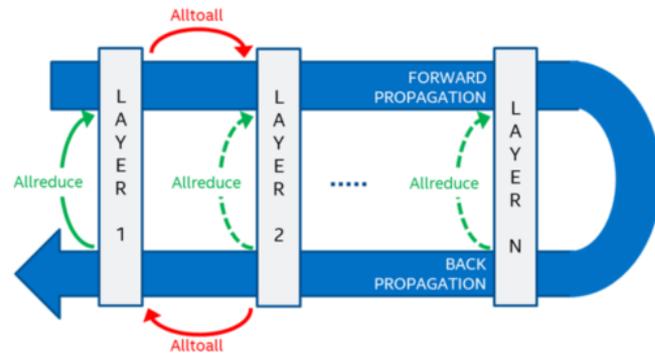
MULTI-NODE PARALLELIZATION

- **What parallelism flavor to use?**
 - Use model parallelism when volume of gradients is much higher than volume of activations or when model doesn't fit memory;
 - Use data parallelism otherwise;
 - Parallelism choice affects activations/gradients ratio
 - Data parallelism at scale makes activations \ll weights
 - Model parallelism at scale makes weights \ll activations
 - There're also other parallelism flavors – pipelined, spatial, etc.

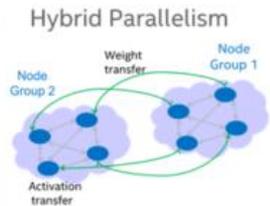
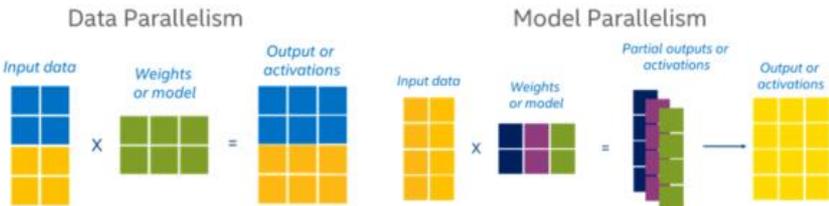
INTEL[®] MACHINE LEARNING SCALING LIBRARY (MLSL)

Distributed Deep Learning Requirements:

- ✓ Compute/communication overlap
- ✓ Choosing optimal communication algorithm
- ✓ Prioritizing latency-bound communication
- ✓ Portable / efficient implementation
- ✓ Ease of integration with quantization algorithms
- ✓ Integration with Deep Learning Frameworks



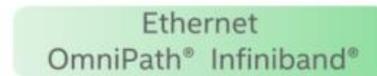
Communication dependent on work partitioning strategy
 Data parallelism = Allreduce (or) Reduce_Scatter + Allgather
 Model parallelism = AlltoAll



Numerous DL Frameworks



Multiple NW Fabrics



INTEL[®] MACHINE LEARNING SCALING LIBRARY (MLSL)

<https://github.com/01org/MLSL/releases>

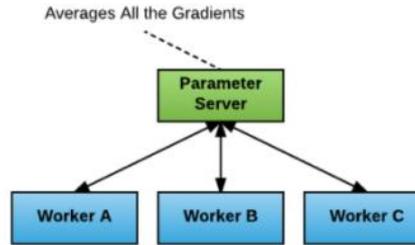
Some of the Intel MLSL features include:

- Built on top of MPI, transparently supports various interconnects: Intel[®] Omni-Path Architecture, InfiniBand*, and Ethernet;
- Optimized to drive scalability of DL communication patterns
- Ability to trade off compute for communication performance – beneficial for communication-bound scenarios
- New domain-specific features are coming soon

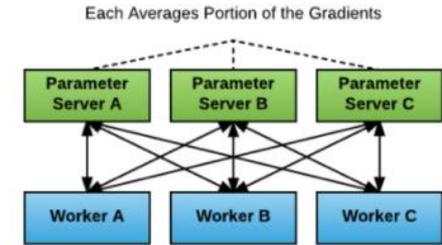
DISTRIBUTED TENSORFLOW*

Distributed
Tensorflow with
Parameter Server

With
Parameter
Server →



or



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

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

Optimization Notice

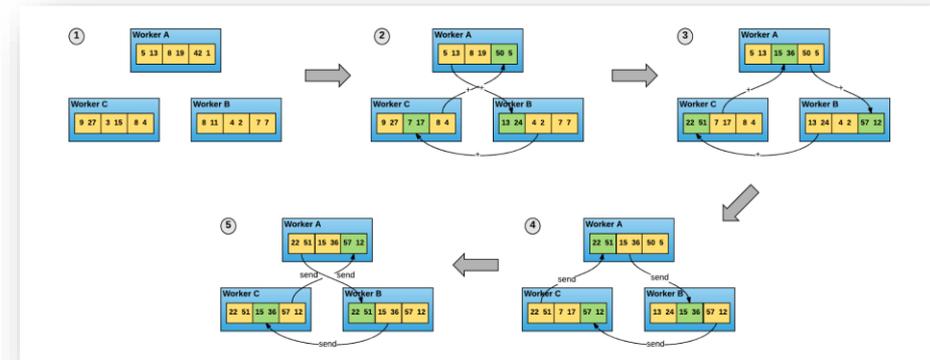
Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.

Intel and the Intel logo 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. Copyright © 2018, Intel Corporation.



DISTRIBUTED TENSORFLOW* WITH HOROVOD

ring-allreduce



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



Uber's open source
Distributed training
framework for TensorFlow

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

- Horovod is a distributed training framework for TensorFlow, Keras, PyTorch, and MXNet.
- The goal of Horovod is to make distributed Deep Learning fast and easy to use
- Horovod core principles are based on MPI concepts such as **size**, **rank**, **local rank**, **allreduce**, **allgather** and **broadcast**.
- Separate infrastructure with model development
- Advantages
 - Minimal code changes to run distributed TensorFlow
 - Network-optimal
 - No parameter server

More info: <https://github.com/horovod/horovod/>

Optimization Notice

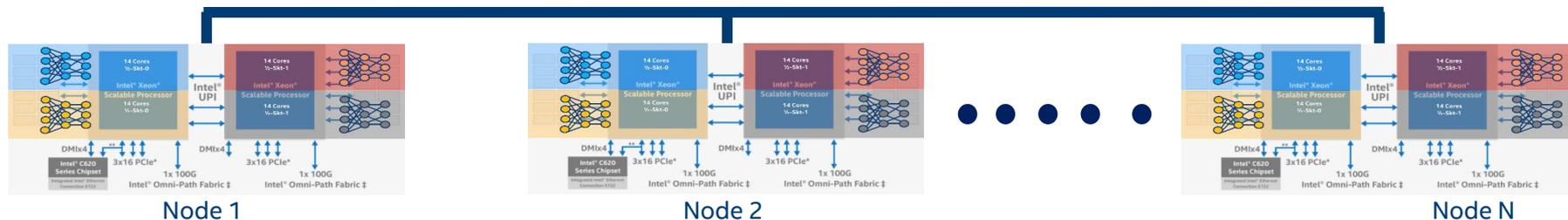
Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.

Intel and the Intel logo 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. Copyright © 2018, Intel Corporation.



DISTRIBUTED TRAINING WITH HOROVOD* MPI LIB

Interconnect Fabric (Intel® OPA or Ethernet)



Distributed Deep Learning Training Across Multiple nodes

Each node running multiple workers/node

Uses optimized MPI Library for gradient updates over network fabric

Caffe – Use Optimized Intel® MPI ML Scaling Library (Intel® MLSL)

TensorFlow* – Uber horovod MPI Library

Intel Best Known Methods: <https://ai.intel.com/accelerating-deep-learning-training-inference-system-level-optimizations/>

Optimization Notice

Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.

Intel and the Intel logo 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. Copyright © 2018, Intel Corporation.



HOROVOD: HOW TO CHANGE THE CODE

- Add `import horovod.tensorflow as hvd` and run `hvd.init()` in the beginning of the program
- Scale the learning rate by number of workers. Effective batch size in synchronous distributed training is scaled by the number of workers. An increase in learning rate compensates for the increased batch size.
- Wrap optimizer in `hvd.DistributedOptimizer`. The distributed optimizer delegates gradient computation to the original optimizer, averages gradients using `allreduce` or `allgather`, and then applies those averaged gradients.
- Add `hvd.BroadcastGlobalVariablesHook(0)` to broadcast initial variable states from rank 0 to all other processes. This is necessary to ensure consistent initialization of all workers when training is started with random weights or restored from a checkpoint. Alternatively, if you're not using `MonitoredTrainingSession`, you can simply execute the `hvd.broadcast_global_variables` op after global variables have been initialized.
- Modify your code to save checkpoints only on worker 0 to prevent other workers from corrupting them. This can be accomplished by passing `checkpoint_dir=None` to `tf.train.MonitoredTrainingSession`, if `hvd.rank() != 0`.

<https://github.com/horovod/horovod#usage>

HOROVOD 101 QUICK START



```
import horovod.tensorflow as hvd
hvd.init()
```

```
#Scale the optimizer
```

```
opt = tf.train.AdagradOptimizer(0.01 * hvd.size())
```

```
# Add Horovod Distributed Optimizer
```

```
opt = hvd.DistributedOptimizer(opt)
```

```
hooks = [hvd.BroadcastGlobalVariablesHook(0)]
```

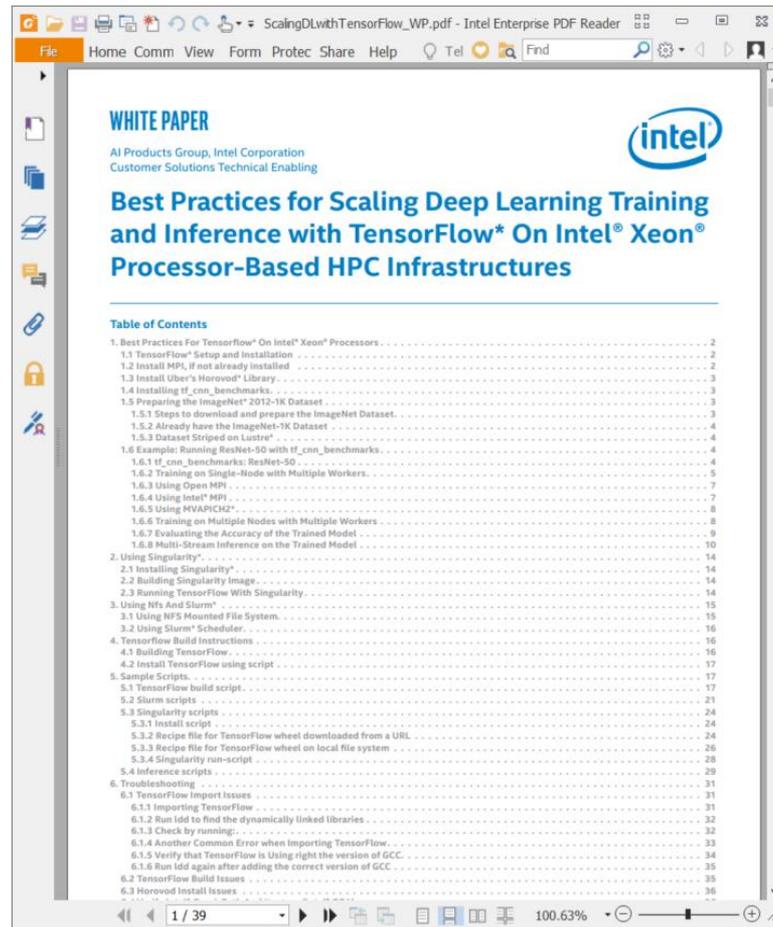
```
# Save checkpoints only on worker 0 to prevent other workers from
corrupting them.
```

```
checkpoint_dir = '/tmp/train_logs' if hvd.rank() == 0 else None
```

SCALING TENSORFLOW*

There is way more to consider when striking for peak performance on distributed deep learning training:

<https://ai.intel.com/white-papers/best-known-methods-for-scaling-deep-learning-with-tensorflow-on-intel-xeon-processor-based-clusters/>



The screenshot shows a PDF document viewer displaying a white paper. The document is titled "WHITE PAPER" and is from the "AI Products Group, Intel Corporation Customer Solutions Technical Enabling". The main title is "Best Practices for Scaling Deep Learning Training and Inference with TensorFlow* On Intel® Xeon® Processor-Based HPC Infrastructures". The document includes a "Table of Contents" with the following sections and page numbers:

Section	Page
1. Best Practices For TensorFlow* On Intel® Xeon® Processors	2
1.1 TensorFlow* Setup and Installation	2
1.2 Install MPI, if not already installed	2
1.3 Install Uber's Horovod* Library	3
1.4 Installing tf_cnn_benchmarks	3
1.5 Preparing the ImageNet* 2012-1K Dataset	3
1.5.1 Steps to download and prepare the ImageNet Dataset	3
1.5.2 Already have the ImageNet-1K Dataset	4
1.5.3 Dataset Striped on Lustre*	4
1.6 Example: Running ResNet-50 with tf_cnn_benchmarks	4
1.6.1 tf_cnn_benchmarks: ResNet-50	4
1.6.2 Training on Single-Node with Multiple Workers	5
1.6.3 Using Open MPI	7
1.6.4 Using Intel* MPI	7
1.6.5 Using MVAPICH2*	8
1.6.6 Training on Multiple Nodes with Multiple Workers	8
1.6.7 Evaluating the Accuracy of the Trained Model	9
1.6.8 Multi-Stream Inference on the Trained Model	10
2. Using Singularity*	14
2.1 Installing Singularity*	14
2.2 Building Singularity Image	14
2.3 Running TensorFlow With Singularity	14
3. Using NFs And Slurm*	15
3.1 Using NF's Mounted File System	15
3.2 Using Slurm* Scheduler	16
4. Tensorflow Build Instructions	16
4.1 Building TensorFlow	16
4.2 Install TensorFlow using script	17
5. Sample Scripts	17
5.1 TensorFlow build script	17
5.2 Slurm scripts	21
5.3 Singularity scripts	24
5.3.1 Install script	24
5.3.2 Recipe file for TensorFlow wheel downloaded from a URL	24
5.3.3 Recipe file for TensorFlow wheel on local file system	26
5.3.4 Singularity transcript	28
5.4 Inference scripts	29
6. Troubleshooting	31
6.1 TensorFlow Import Issues	31
6.1.1 Importing TensorFlow	31
6.1.2 Run ldd to find the dynamically linked libraries	32
6.1.3 Check by running	32
6.1.4 Another Common Error when Importing TensorFlow	33
6.1.5 Verify that TensorFlow is Using right the version of GCC	34
6.1.6 Run ldd again after adding the correct version of GCC	35
6.2 TensorFlow Build Issues	35
6.3 Horovod Install Issues	36

INTEL[®] MLSL BACKEND FOR HOROVOD

Install procedure:

- Install the latest versions of Intel MLSL and Intel MPI;
- `source <mlsl_install>/intel64/bin/mlslvars.sh thread`
- `source <intel_mpi_2019>/intel64/bin/mpivars.sh release_mt`
- Download Horovod and build it from source code or
 - `pip install horovod`

INTEL[®] MLSL BACKEND FOR HOROVOD

Launch procedure:

- `export MLSL_LOG_LEVEL=1`
 - output from within MLSL
- `export MLSL_NUM_SERVERS=X`
 - X is the number of cores you'd like to dedicate for driving communication
- `export MLSL_SERVER_AFFINITY=c1,c2,...,cX`
 - Core IDs dedicated to MLSL servers (uses X 'last' cores by default)
- `export HOROVOD_MLSLbackground_BGT_AFFINITY=c0`
 - Affinity for thread of Horovod
- Adjust OpenMP settings to avoid intersection with c0,c1,...,cX

Optimization Notice

Copyright © 2019, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.



HANDS-ON

TENSORFLOW+HOROVOD/CNN_MNIST-HVD.IPYNB

Delete the checkpoint if needed, otherwise TF won't train any further

```
- rm -rf checkpoints
```

Let's start changing the number of MPI tasks, what performance difference would you expect?

```
- mpirun -prepend-rank -genv OMP_NUM_THREADS=2 -genv I_MPI_DEBUG=5 -n 2 python -u cnn_mnist-hvd.py  
- mpirun -prepend-rank -genv OMP_NUM_THREADS=2 -genv I_MPI_DEBUG=5 -n 4 python -u cnn_mnist-hvd.py  
- check the size of the dataset:  
  - ls -lha ~/.keras/datasets/
```

Intel Python and Optimized Tensorflow

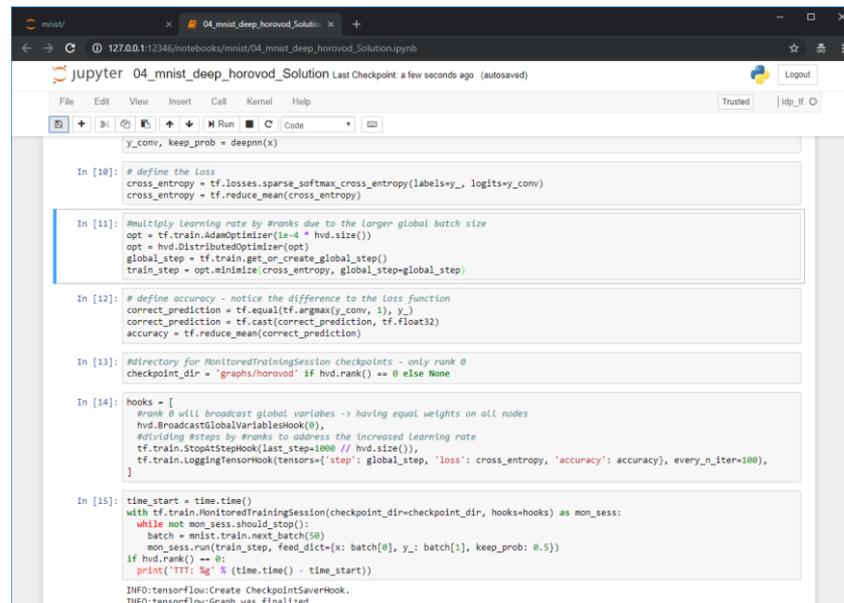
```
- source activate hvd-mpi  
- pip show tensorflow | grep Location  
  - useful to locate the TF installation for see the library linked: ldd $Location/tensorflow/libtensorflow...so  
- rm -rf /tmp/*  
- export export MKLDNN_VERBOSE=1
```

TENSORFLOW+HOROVOD/CNN_MNIST-HVD.IPYNB

- 1) How to initialize Horovod and why is it necessary?
- 2) Why is it necessary to adapt the learning rate with larger batches?
- 3) How can you dynamically adapt the learning rate?
- 4) How to identify rank #1 (0)?
- 5) Why is it necessary to adapt the number of training steps according to the number of workers / larger batches?
- 6) How can you dynamically adapt the number of training steps?
- 7) How is the single process performance vs 2 ranks vs 4 ranks?

MNIST CNN HOROVOD DEMO SUMMARY

- Horovod initializes the MPI communication underneath and therefore defines rank() and size()
- In order to reduce the Time To Train with multiple workers, therefore increasing the batch size, the learning rate needs to scale
- Same for the # of steps for training
- 4 ranks can be faster since less threading efficiency is required in small convolutions



```
mnist/ 04_mnist_deep_horovod_Solution
127.0.0.1:2346/notebooks/mnist/04_mnist_deep_horovod_Solution.ipynb
jupyter 04_mnist_deep_horovod_Solution Last checkpoint a few seconds ago (autosaved)
File Edit View Insert Cell Kernel Help Trusted | kdp_if_o

y_conv, keep_prob = deeppn(x)

In [10]: # define the loss
cross_entropy = tf.losses.sparse_softmax_cross_entropy(labels=y_, logits=y_conv)
cross_entropy = tf.reduce_mean(cross_entropy)

In [11]: # Multiply Learning rate by #ranks due to the larger global batch size
opt = tf.train.AdamOptimizer(lr * hvd.size())
opt = hvd.DistributedOptimizer(opt)
global_step = tf.train.get_or_create_global_step()
train_step = opt.minimize(cross_entropy, global_step=global_step)

In [12]: # define accuracy - notice the difference to the loss function
correct_prediction = tf.equal(tf.argmax(y_conv, 1), y_)
correct_prediction = tf.cast(correct_prediction, tf.float32)
accuracy = tf.reduce_mean(correct_prediction)

In [13]: # directory for NonStoredTrainingSession checkpoints - only rank 0
checkpoint_dir = 'graphs/horovod' if hvd.rank() == 0 else None

In [14]: hooks = [
#rank 0 will broadcast global variables -> having equal weights on all nodes
hvd.BroadcastGlobalVariablesHook(0),
#dividing #steps by #ranks to address the increased learning rate
tf.train.StopAtStepHook(last_step=1000 // hvd.size()),
tf.train.LoggingTensorHook(tensors=['step': global_step, 'loss': cross_entropy, 'accuracy': accuracy], every_n_iter=100),
]

In [15]: time_start = time.time()
with tf.train.MonitoredTrainingSession(checkpoint_dir=checkpoint_dir, hooks=hooks) as mon_sess:
while not mon_sess.should_stop():
batch = mnist.train.next_batch(50)
mon_sess.run(train_step, feed_dict={x: batch[0], y_: batch[1], keep_prob: 0.5})
if hvd.rank() == 0:
print('TTT: %g' % (time.time() - time_start))
INFO:tensorflow:Create CheckpointSaverHook.
INFO:tensorflow:Graph was finalized.
```

Legal Disclaimer & Optimization 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 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.

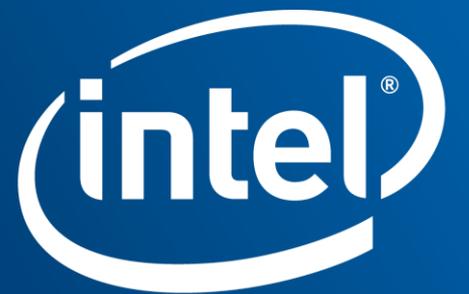
Copyright © 2019, Intel Corporation. All rights reserved. Intel, the Intel logo, Pentium, Xeon, Core, VTune, OpenVINO, Cilk, are trademarks of Intel Corporation or its subsidiaries in the U.S. and other countries.

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

BACKUP



Software