# Deploying AI Frameworks on Secure HPC Systems with Containers.

26.09.2019| Atanas Atanasov, Fabio Baruffa, David Brayford, Sofia Vallecorsa, Walter Riviera

# Legal Notices & Disclaimers

This document contains information on products, services and/or processes in development.  All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at intel.com, or from the OEM or retailer. No computer system can be absolutely secure.

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 http://www.intel.com/performance.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings.  Circumstances will vary.  Intel does not guarantee any costs or cost reduction.

Statements in this document that refer to Intel's plans and expectations for the quarter, the year, and the future, are forward-looking statements that involve a number of risks and uncertainties. A detailed discussion of the factors that could affect Intel's results and plans is included in Intel's SEC filings, including the annual report on Form 10-K.

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.

No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.
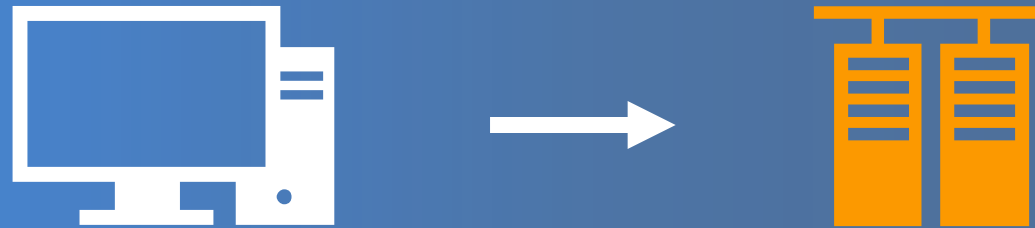
Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

Intel, the Intel logo, Pentium, Celeron, Atom, Core, Xeon, Movidius and others 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.

© 2019 Intel Corporation.

# High Performance AI (HPAI) in a **Container**

Transition AI algorithms from the
**laptop to supercomputer**
with minimal effort

**"It just works"**

# HPAI =

**M&S**

- Equation based on model
- Computing driven
- Numerically intensive
- Creates simulations
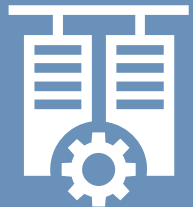- Monte Carlo
- Larger problems
- Iterative methods
- PDE

**+**

- Linear algebra
- Matrix operations
- Iterative methods
- Compute intensive
- Data transfer
- Predictive
- Probabilities
- Stencil codes
- Calculus
- Pattern recognition
- Graphs

**Analytics**

- Finds patterns
- Correlations in data
- Logic driven
- Creates inferences
- Knowledge discovery
- Graphs
- Data-driven science
- Predictions
- CNN
- RNN

lrz

# Requirements for AI on HPC

**Compute intensive hardware**

**Optimized AI frameworks**
TensorFlow, PyTorch, Caffe

**Optimized software**
numerical libraries, Python

**HPC specific software**
distributed computing, workload manager

**Method of deploying the AI software**
in a simple, straight-forward and flexible way

## Need to get to: "It just works"

# Key Challenges

## Package Management

### Frameworks have conflicting dependencies

The frameworks & their dependencies need to be combined in a single module

### Rapid update cycles

Provide a mechanism for users to build there own frameworks

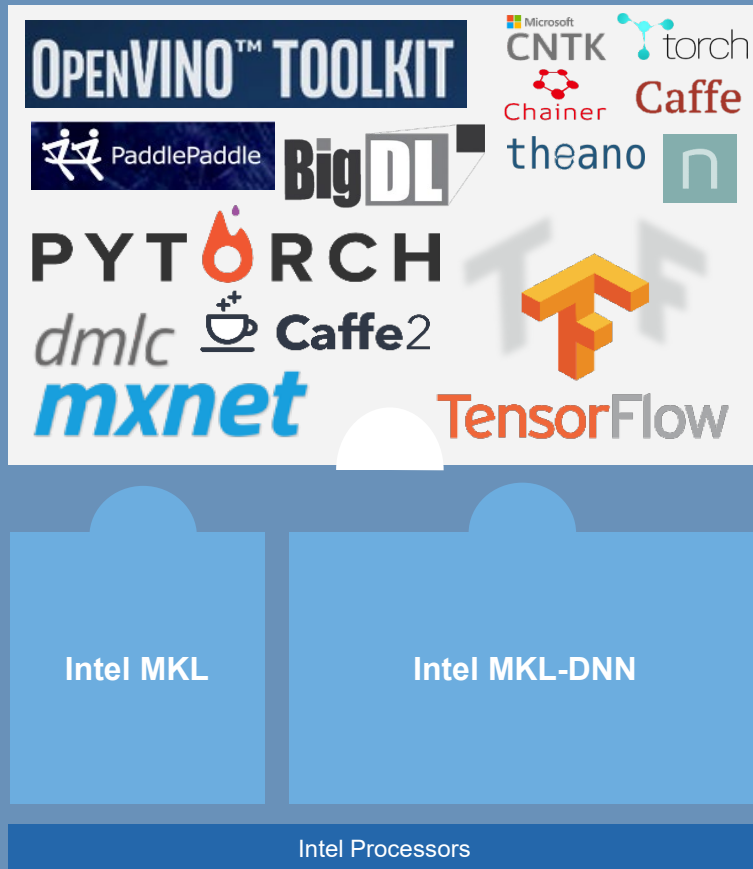## Dynamic Programming Environment

### Python dependencies

Each unique framework needs its own Python instance

### Connecting to external servers

Build frameworks on systems without internet access

# Distributed Mechanisms



## System-level Optimizations

# Detecting and Identifying High Energy Physic Particles

- CLIC Electromagnetic calorimeter
  - Sparse images
  - Highly segmented (pixelized)
  - Large dynamic range
- Segmentation is critical for particle identification and energy determination

# 3D Convolutional GAN



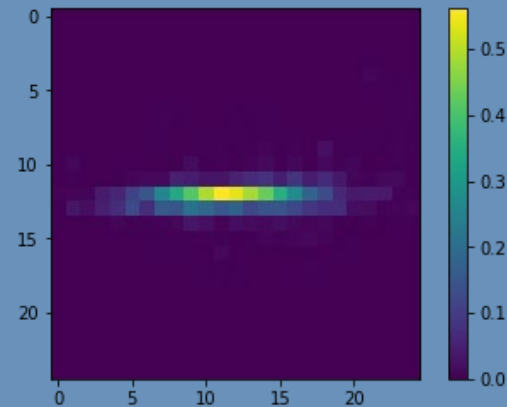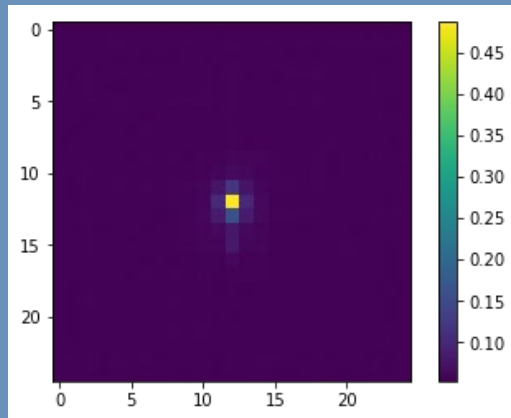**Generator**

**Discriminator**

~1M parameters

Total model Size: 3.8MB

# Charliecloud Containers in HPC



- Easy to install

- Charliecloud was developed to be run on highly secure HPC systems at US government labs

- Charliecloud runs entirely under the User ID

- Ability to run legacy design flows in containers

- Low overhead and ~ 800 lines of code

- LRZ deploys Charliecloud via Spack

- Charliecloud is available in the module system at LRZ

# Achieving High Performance AI on Secure HPC Systems

## Mechanism for deploying AI at LRZ

- Download the Intel optimized TensorFlow Docker Image (intelaipg Dockerhub)

- Modify the Linux Docker image for HPC

- Modify Python to enable distributed TensorFlow execution

- Copy the training data and execution scripts to the modified Docker image

- Convert to a Charliecloud UDSS and copy the file to the HPC system

- Load the Charlicloud module

- Execute on SuperMUC-NG via Slurm

# Distributed TensorFlow Results LRZ SNG 1 MPI Rank

## 1 MPI rank & 48 OpenMP threads per node

## Intel Skylake Platinum Xeon 8174
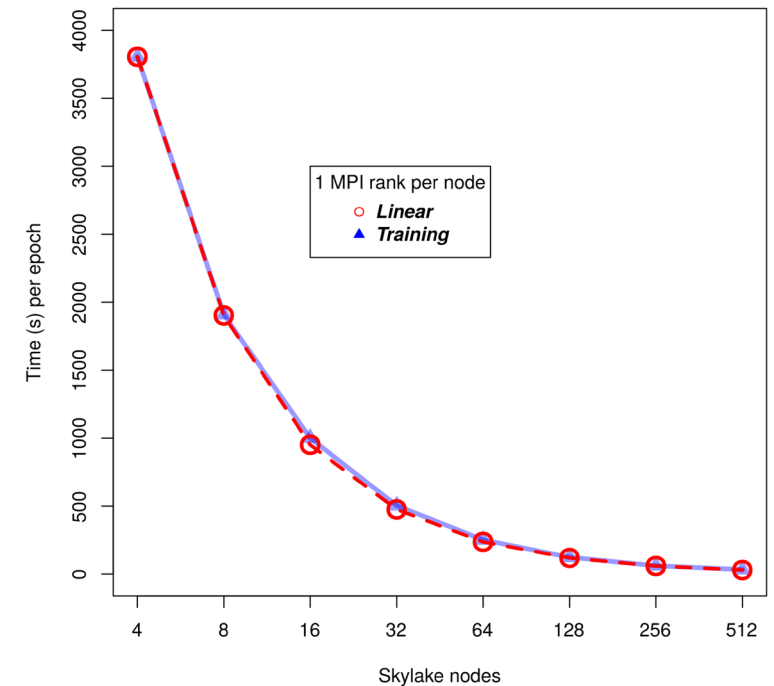
| Nodes | Training Time(S) per Epoch | Linear Time(S) per Epoch | Scaling Efficiency |
|-------|---------------------------|--------------------------|---------------------|
| 4 | 3806 | 3806 | - |
| 8 | 1910 | 1903 | 99.6% |
| 16 | 1001 | 951.5 | 95.1% |
| 32 | 504 | 475.75 | 94.4% |
| 64 | 253 | 237.87 | 94% |
| 128 | 124 | 118.93 | 95.9% |
| 256 | 61 | 59.46 | 97.5% |
| 512 | 33 | 29.73 | 90.1% |



## Throughput Overheads

| Benchmark | Free System Memory with Charliecloud (GB) | Free System Memory without Charliecloud (GB) |
|-----------|-------------------------------------------|----------------------------------------------|
| AlexNet with cifar | 331.29 | 331.33 |
| ResNet50 with imagenet | 324.47 | 324.89 |

## Memory Overheads

| Benchmark | Free System Memory with Charliecloud (GB) | Free System Memory without Charliecloud (GB) |
|-----------|-------------------------------------------|----------------------------------------------|
| AlexNet with cifar | 331.29 | 331.33 |
| ResNet50 with imagenet | 324.47 | 324.89 |

# 3DGAN Execution with 4 MPI Ranks per Node

Stampede2 @ TACC 11 OpenMP threads per MPI task Intel Skylake Platinum Xeon 8160, Standard horovod + MPI, without Charliecloud

| Nodes | Training Time(S) per Epoch | Linear Time(S) per Epoch | Scaling Efficiency |
|---|---|---|---|
| 1 | 17831 | 17831 | - |
| 2 | 8998 | 8915.5 | 99.1% |
| 4 | 4545 | 4457.75 | 98.08% |
| 8 | 2288 | 2228.87 | 97.4% |
| 16 | 1151 | 1114.44 | 96.8% |
| 32 | 581 | 557.22 | 95.9% |
| 64 | 293 | 278.61 | 95.1% |
| 128 | 148 | 139.60 | 94.1% |

SuperMUC-NG @ LRZ 12 OpenMP threads per MPI task Intel Skylake Platinum Xeon 8174, Standard horovod + MPI, with Charliecloud

| Nodes | Training Time(S) per Epoch | Linear Time(S) per Epoch | Scaling Efficiency |
|---|---|---|---|
| 4 | 959 | 959 | - |
| 8 | 507 | 479.5 | 94.6% |
| 16 | 264 | 239.75 | 90.8% |
| 32 | 137 | 119.87 | 87.5% |
| 64 | 72 | 59.93 | 83.3% |
| 128 | 39 | 29.96 | 76.8% |
| 256 | 21 | 14.98 | 71.4% |
| 512 | 12 | 7.49 | 62.5% |

# Third Quarter 2019

## Release SC'19 Denver

HPC suitable Intel optimized TensorFlow Docker image
Verified recipes to enable the deployment of AI on HPC systems using secure containers
Github repository https://github.com/DavidBrayford/HPAI

## Current Users

DLR German Aerospace Center, PyTorch, inferencing of high resolution satellite images on SuperMUC-NG