# Bright Computing

## Advanced cluster & cloud management made easy

Panos Labropoulos
Bright Computing, Inc.

1 ABOUT

# About Bright

## Who we are

- Enterprise infrastructure software company
- Incorporated in 2009
- Headquarters
  - San Jose, California
  - Amsterdam, Netherlands

## What we do

- Bare-metal & cloud provisioning, monitoring, management of clustered systems in the data center: Hadoop, private cloud, HPC, storage, database and other "clustered" systems
- From a number of servers to →
  - HPC
  - Hadoop cluster/ Machine Learning
  - OpenStack private cloud

## Our success

- Customers around the world: Boeing, ING Bank, Lincoln Financial, Novartis, NASA, DoD, DoE, Stanford, Oracle, HP, Intel, Sinopec and >600 more
- Resellers around the world:
  - OEMs: Atos/Bull, Cray, DDN, Dell, Huawei, Lenovo, SGI, HPE
  - Integrators: >75 around the world
- Award winning ...

**2** CUSTOMERS

Bright Computing

**Enterprise**

**Government**

**Education**

# Embedded in other solutions

In the open seas

InfiniCortex 2015

10

# Dark sites

# Partners

## Resellers / OEMs / ISVs

DELL

Bull
atos technologies

HUAWEI

CRAY

Hewlett Packard
Enterprise

lenovo

Bright Cluster Manager users

Created with mapchart.net ©

- From  4 to 8,000+ nodes
- Single user to few million jobs/day
- Multiple architectures: x86, Power 8/9 and soon ARMv8
- Multiple Linux distributions: RHEL/Centos/SL 6 and 7 , SLES 11 & 12, Ubuntu LTS
- Multiple interconnects: Infiniband, OmniPath, HSE, ROCE, usNIC, PCI-E switches, iWARP, NVlink
- Multiple types of accelerators, NVIDIA GPUs, AMD GPUs, Intel Xeon Phis, FPGAs and other ASICs
- Bare-metal vs virtualized
- Bare-metal vs containerized
- On-premise vs on the cloud

**Heterogenous clusters are the rule, not the exception!**

**Bright** Computing

# 3 THE PROBLEM

- Clusters are still hard to setup, operate, use and manage

- New buyers require "ease-of-everything"

- The shortage of skilled people is a major hindrance

- Software is still the #1 roadblock

- Better management software is needed

# A Better Solution

- Bright Cluster Manager takes a much more fundamental & integrated approach
  - Designed and written from the ground up
  - Single cluster management agent provides all functionality
  - Single, central database for configuration and monitoring data
  - Single UI for ALL cluster management functionality
  - Enterprise-level support

- Which makes Bright Cluster Manager …
  - Extremely easy to use
  - Extremely scalable
  - Secure & reliable
  - Complete
  - Flexible
  - Maintainable

# Bright Architecture

Bright Cluster Manager and Deep Learning

- BCM DL packages is the default AI stack for a number of OEMs including:

**Dell EMC Machine Learning and Deep Learning technical specifications**

| Configuration | Single node 'D' | Medium node 'G' | Medium node 'K' | Large multi-node 'G' | Multi-node NVLink GPU 'K' |
|---|---|---|---|---|---|
| | Well-suited for use as a dedicated server to train models before deploying to production | Suited for larger training jobs, with support for 4 GPUs | Peer-to-peer interconnect between GPUs provides better training throughput | Suited for scaling out deep learning frameworks at a cluster level | Peer-to-peer interconnect between GPUs provides better training throughput |
| Compute | PowerEdge C4130 Server | PowerEdge C4130 Server | PowerEdge C4130 Server | PowerEdge C4130 Server | PowerEdge C4130 Server |
| Processor | 2 x Intel Xeon E5-2690 v4 | 2 x Intel Xeon E5-2690 v4 | 2 x Intel Xeon E5-2690 v4 | 2 x Intel Xeon E5-2690 v4 | 2 x Intel Xeon E5-2690 v4 |
| Memory | 128GB DDR4 @ 2,400 MHz | 256GB DDR4 @ 2,400 MHz | 256GB DDR4 @ 2,400 MHz | 256GB DDR4 @ 2,400 MHz | 256GB DDR4 @ 2,400 MHz |
| Solid state drives (SSDs) | 2 x 200GB 1.8 inch SSDs | 2 x 200GB 1.8 inch SSDs | 2 x 200GB 1.8 inch SSDs | 2 x 200GB 1.8 inch SSDs | 2 x 200GB 1.8 inch SSDs |
| Networking | NA | NA | NA | Mellanox® InfiniBand® EDR 100Gb/s with ConnectX-4 VPI adapter card | Mellanox InfiniBand EDR 100Gb/s with ConnectX-4 VPI adapter card |
| Accelerator | 2 x NVIDIA Tesla P100 16GB PCIe | 4 x NVIDIA Tesla P100 16GB PCIe | 4 x NVIDIA Tesla P100 SXM2 16GB | 4 x NVIDIA Tesla P100 16GB PCIe | 4 x NVIDIA Tesla P100 SXM2 16GB |
| Software | Bright Computing Solution for Deep Learning | Bright Computing Solution for Deep Learning | Bright Computing Solution for Deep Learning | Bright Computing Solution for Deep Learning | Bright Computing Solution for Deep Learning |

# CRAY® CS-STORM™ ACCELERATED GPU SYSTEM

Cray® CS-Storm™ cluster supercomputers tackle the toughest extreme HPC and artificial intelligence (AI) workloads. Designed for speed, architected for scale and integrated for production use, the Cray CS-Storm GPU-accelerated supercomputer is your path to exploiting the performance available from the latest NVIDIA® Tesla® GPUs.

**AI Software**

Supports Bright Cluster Manager™ for HPC and an optional Bright Deep Learning extension

Machine learning frameworks including Caffe, Chainer, CNTK, Keras, MXet, Pytorch, Scikit-learn, Tensorflow, Theano, Torch, Dynet, Intel Neon, TensorFlowOnSpark, BigDL, Deeplearn4j

# DWell benchmarks

MXNet Inception-BN ILSVRC12

http://en.community.dell.com/techcenter/high-performance-computing/b/general_hpc/archive/2016/11/11/deep-learning-performance-with-p100-gpus

# DS Ecosystem

- Easily call deep learning libraries within existing Big Data workflows, making it immediately available to Big Data application developers.

- Seamlessly perform transfer learning of deep learning models

- Leverage existing Big Data installations to productionize high quality models at scale

- Exploit the already excellent support of Big Data tools when it come to SQL functions

- Process more easily  complex data such as images, audio and video from Big Data pipelines

- Support all/most popular frameworks from Python, Java, Scala and R

# A Deep Learning pipeline



Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

Zhang et al.

31

**Glue code**: DL package are developed as stand-alone solutions. Incorporating those in a data reduction pipeline results a glue code paradigm (5% DL code, 95% glue code). Makes it hard to take advantage of domain-specific properties

**Pipeline jungles**: Preparing data in an ML-friendly format may become a jungle of scrapes, joins, and sampling steps, often with intermediate files output. Often the result of separation between "Research" and "Engineering"

**Dead experimental code paths**:  It becomes increasingly attractive in the short term to perform experiments with alternative frameworks.  Increasing difficulties of maintaining backward compatibility and an  complexity.

# Goals (contd.)

- Turn-key solution: keys included (NCCL, CuDNN, HPC-X, MKL, CUDA)

- Optimized out of the box

- Native applications for the underlying distro

- Not enforcing the use of docker

- Use native package manager to install frameworks + dependencies

- Install on shared storage

- Much more aggressive optimizations on CPU part

- Integration with HPX/Big Data enabled (e.g. Tensorflow understands AmazonS3 and HDFS)

- QA testing and release engineering & leveraging expertise and feedback from more than 500 HPC sites

- -dev packages

- Common environment across all supported distros and archs (CUDA 9.1, CuDNN 7.0, MKL/OpenBLAS 0.2.20, GCC_cuda, OMPI 3.x+Mellanox extensions, NCCL2)

- Native Python + Python 3.6 support

- Examples should work out of the box

# With Bright:  Two simple commands

```
# yum install tensorflow cntk mxnet digits

# yum --installroot=/cm/images/default-image \
install cm-ml-distdeps
```

- 1st command installs Tensorflow etc in a shared directory on the head node; It is immediately available on every node in cluster
- Yum installs all dependencies for DIGITS including Caffe, Torch, and all Python dependencies
- 2nd command installs all library dependencies into default-image

# Running an application is really simple

```
[root@bright80-r720 ~]# module load tensorflow/1.5.0rc1

[root@bright80-r720 ~]# cd  models/image/imagenet

[root@bright80-r720 imagenet]# python classify_image.py --image_file=Indochinese-Tiger-Zoo.jpg
I tensorflow/core/common_runtime/gpu/gpu_init.cc:102] Found device 0 with properties:
name: Tesla K40c
major: 3 minor: 5 memoryClockRate (GHz) 0.745
pciBusID 0000:05:00.0
Total memory: 11.92GiB
Free memory: 11.78GiB


tiger, Panthera tigris (score = 0.71628)
tiger cat (score = 0.11725)
lynx, catamount (score = 0.00376)
jaguar, panther, Panthera onca, Felis onca (score = 0.00371)
cougar, puma, catamount, mountain lion, painter, panther, Felis concolor (score = 0.00218)
[root@bright73-r720 imagenet]#
```

# Developer-friendly

- **Some users are Dl framework developers**

```
[root@deeplearning ~]# module load tensorflow bazel
[root@deeplearning ~]# cd SRC/tensorflow
[root@deeplearning tensorflow]# bazel build --config=opt --config=cuda //tensorflow/tools/pip_package:build_pip_package
WARNING: /root/SRC/tensorflow/tensorflow/contrib/learn/BUILD:15:1: in py_library rule //tensorflow/contrib/learn:learn:
ndle:exporter': Use SavedModel Builder instead.
WARNING: /root/SRC/tensorflow/tensorflow/contrib/learn/BUILD:15:1: in py_library rule //tensorflow/contrib/learn:learn:
ndle:gc': Use SavedModel instead.
INFO: Found 1 target...
Target //tensorflow/tools/pip_package:build_pip_package up-to-date:
  bazel-bin/tensorflow/tools/pip_package/build_pip_package
INFO: Elapsed time: 0.749s, Critical Path: 0.00s
[root@deeplearning tensorflow]#
```

All the necessary tools to get started with Machine Learning development are installed by default.

# Available Machine Learning Frameworks (January 2018)

| Framework | Graph | Version | Notes |
|---|---|---|---|
| caffe)(-nv) | static | 0.16.5 | half-precision support, nccl2 |
| caffe2 | static | 0.8.1 | OpenMPI, nccl2 |
| caffe-mpi | static | 2.0 | OpenMPI+HPCX acceleration, nccl2 |
| chainer | dynamic | 3.2.0 | |
| CafffeOnSpark | static | 0.1 | Spark 2.x |
| CNTK | static | 2.3.1 | OpenMPI, MKL |
| Keras | N/A | 2.1.2 | Theano, CNTK and Tensorflow backends |
| MXNet | static | 1.0 post4 | Scala bindings |
| Pytorch | dynamic | 0.4.0a0 | +Patches for CUDA 9.1 support and MPI Allreduce, nccl2, OpenmPI 3 + HPC-X |
| Scikit-learn | static | 0.19.1 | |
| Tensorflow | static* | 1.5.0 | Infiniband, HDFS, SPARK, Mesos, Docker, Kubernetes, Fold, Serving, nccl2 |
| Theano | static | 1.0.1 | |
| Torch | static | 7.0 | + patches for CUDA 9.1 support, nccl2, FP16 fixes |
| *dynet* | *dynamic* | *2.0.2* | |
| *Intel Neon* | *static* | *2.6.0* | |
| *TensorFlowOnSpark* | *static* | *0.12* | |
| *BigDL* | *N/A* | *0.1.1* | |
| *deeplearning4j* | | *0.5.0* | |

# Available supporting libraries/tools (January 2018)

| Package | Version | Notes |
|---|---|---|
| bazel | 0.9.0 | 0.6.7 on Power |
| cub | 1.7.4 | |
| CUDA Driver | 384.111 | Tesla-certified version |
| CUDA | 9.1 | Only 9.1 is supported for DL |
| CUDNN | 6.0,7.0 | |
| digits | 6.0.1 | Plus bug fixes from 'master' branch |
| gflags | 2.2.1 | |
| glog | 0.3.5 | |
| HDF5 | 1.8.18 | |
| leveldb | 1.20 | |
| lmdb | 0.9.21 | Scala bindings |
| NCCL | 1.3.4-1 and 2.0 | Patches for CUDA 9.1 support to NCCL 1.3.4 |
| OpenCV | 3.1.0 | 3.2.0 does not work with CNTK |
| Protobuf | 3.2.0 and 2.5.4 | |
| eigen3 | | |
| python | 3.6.3 + 2.7.5 | |
| Mellanox HPC-X | 2.0 | |
| EasyBuild | 3.5.1 | |
| OpenMPI | 3.0.0 | CUDA 9.1 and Mellanox MXM, Knem, HCOLL, FCA support |
| TensorRT | 3.1.2 | |

- Python package versions are the common denominator

| | | | |
|---|---|---|---|
| chardet | 3.0.2 | pandas | 0.19.2 |
| configobj | 4.7.2 | pathlib2 | 2.2.0 |
| Cython | 0.25.2 | pip | |
| Flask | 0.10.1 | protobuf | 3.1.0.post1 |
| Flask-SocketIO | 2.6 | psutil | 3.4.2 |
| Flask-WTF | 0.12 | py4j | 0.10.4 |
| gevent | 1.2.1 | pycuda | 2015.1 |
| gunicorn | 17.5 | pycurl | 7.19.0 |
| h5py | 2.6.0 | pydot | 1.1.0 |
| hypothesis | | pydot2 | 1.0.33 |
| iniparse | 0.4 | pyliblzma | 0.5.3 |
| ipython | 5.1.0 | python-gflags | 3.1.0 |
| Jupyter | 4.3.0 | pytools | |
| JupyterHub | ` | PyYAML | 3.12 |
| kitchen | 1.1.1 | scikit-cuda | 0.5.1 |
| leveldb | 0.194 | scikit-fmm | 0.0.9 |
| lmdb | 0.87 | scikit-image | 0.12.3 |
| matplotlib | 1.5.1 | scikit-learn | |
| mercurial | 2.6.2 | scipy | |
| mock | 2.0.0 | selenium | 3.0.2 |
| netaddr | 0.7.18 | setuptools | |
| nltk | | Sphinx | 1.5.1 |
| nose | 1.3.7 | wheel | 0.29.0 |
| nose-parameterized | 0.5.0 | | |

# Always the latest upstream version

- Two week update cycle – DL frameworks change **FAST**
- Upgradable via package manager e.g. apt-get dist-upgrade

```
[root@deeplearning scripts]# ml-check-updates
PACKAGE_NAME   | CHANGED | REPOSITORY                                         | UPSTREAM_VERSION                           | BRIGHT_VERSION
bazel          | YES     | https://github.com/bazelbuild/bazel                | 0.5.2                                      | 0.5.1
caffe2         | NO      | https://github.com/caffe2/caffe2.git               | 0.7.0                                      | 0.7.1
caffe          | NO      | https://github.com/NVIDIA/caffe.git                | 0.16.2                                     | 1.0
caffe-mpi      | YES     | https://github.com/Caffe-MPI/Caffe-MPI.github.io.git | 6c2c34743583af695e8eb9fc6c946371142c3b39 | 6c2c347
caffeonspark   | YES     | https://github.com/yahoo/CaffeOnSpark.git          | 19df500abe3f0d09511b6434a0ea0bb52a6e8124  | 0.1
chainer        | YES     | https://github.com/chainer/chainer.git             | 3.0.0a1                                    | 2.0.0
cm-protobuf3   | YES     | https://github.com/google/protobuf.git             | 3.3.2                                      | 3.1.0
cntk           | NO      | https://github.com/Microsoft/CNTK.git              | 2.0                                        | 2.0
cub            | YES     | https://github.com/NVlabs/cub.git                  | 1.7.0                                      | 1.6.4
digits         | NO      | https://github.com/NVIDIA/DIGITS.git               | 5.0.0                                      | 5.0.0
gflags         | NO      | https://github.com/gflags/gflags.git               | 2.2.0                                      | 2.2.0
glog           | YES     | https://github.com/google/glog.git                 | 0.3.5                                      | 0.3.3
keras          | NO      | https://github.com/fchollet/keras.git              | 2.0.5                                      | 2.0.5
leveldb        | YES     | https://github.com/google/leveldb.git              | 1.20                                       | 1.18
lmdb           | YES     | https://github.com/LMDB/lmdb.git                   | 0.9.21                                     | 0.9.18
mxnet          | NO      | https://github.com/dmlc/mxnet.git                  | 0.10.0                                     | 0.10.0
nccl           | YES     | https://github.com/NVIDIA/nccl.git                 | 1.3.4-1                                    | 1.3.4
openblas       | YES     | https://github.com/xianyi/OpenBLAS.git             | 0.2.19                                     | 0.2.18
opencv3        | YES     | https://github.com/opencv/opencv.git               | 3.2.0                                      | 3.1.0
pytorch        | NO      | https://github.com/pytorch/pytorch.git             | 0.1.12                                     | 0.1.12
tensorflow     | YES     | https://github.com/tensorflow/tensorflow.git       | 1.2.1                                      | 1.2.0
theano         | NO      | https://github.com/Theano/Theano.git               | 0.9.0                                      | 0.9.0
torch7         | YES     | https://github.com/torch/torch7.git                | 11321f7bf65573849be90650bb376f05dfd79c60  | 7.0
[root@deeplearning scripts]#
```

- Risk of breaking things or making users unhappy if they are unprepared for and upgrade

# Please install this package on the cluster

■ In some cases building from source is preferred as it allows to control libraries, compilers and optimizing the software for a particular architecture (AVX, network, Intel KNL, AMD GPUs,…)
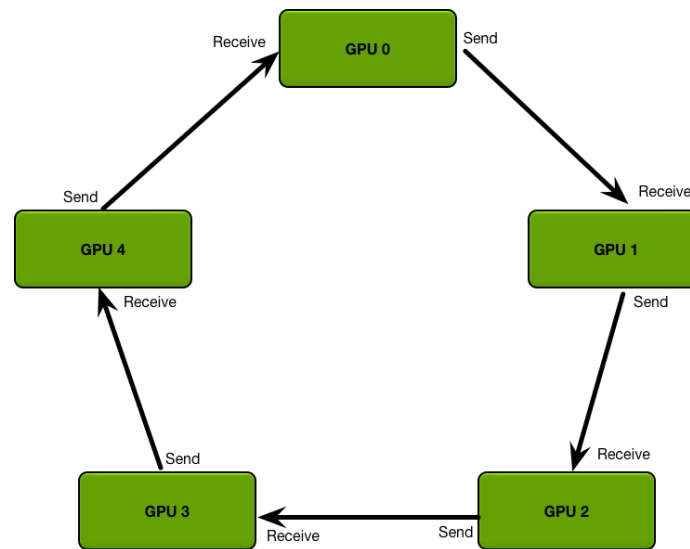
Installing (lots of) *scientific* software is typically:

- error-prone, trial-and-error
- tedious, hard to get right
- repetitive & boring (well...)
- time-consuming (hours, days, even weeks)
- frustrating ("*Pandora's box*")
- sometimes simply not worth the effort...

- ./configure --with-hcoll=/cm/shared/apps/hpcx/2.0.0/hcoll
- --with-mxm=/cm/shared/apps/hpcx/2.0.0/mxm
- --with-knem=/cm/shared/apps/hpcx/2.0.0/knem
- --with-fca
- --with-verbs --enable-mpi-thread-multiple
- --with-slurm **--with-platform=contrib/platform/mellanox/optimized**
- --with-cuda=/cm/shared/apps/cuda80/toolkit/current
- --with-cuda-libdir=/cm/local/apps/cuda-driver/libs/current/lib64  <other options>

- A ring allreduce is a bandwidth-optimal way to do an allreduce. To do the allreduce, the nodes involved are arranged in a ring:

- Each node always sends to the next clockwise node in the ring, and receives from the previous one.

- We looked at several alternatives: EasyBuild, Spack, crazy things like setting up a Jenkins CI infrastructure

Status: new
Ticket <URL: http://support.brightcomputing.com/rt/Ticket/Display.html?id=8390 >

Does Bright support or plan to support using EasyBuild to integrate HPC toolchains into the existing Bright provided modules environment?

We've done some experimenting and EB looks like an easy way for us to incorporate certain toolchains into our envionment without having to sink lots of developer time into making packages work. It would certainly help us get some toolchains into a useable state where our researchers can begin using/testing them.

http://hpcugent.github.io/easybuild/

- EasyBuild and spack made it to the shortlist but most customers chose EB

# Bundled compilers are standardize per Bright release

- Add cm- prefixed toolchains to distinguish from default ones

e.g. cmgcccuda

Need to support:

1. GCC native, 5.4, 6.4 and 7.2 (both old and new ABIs), Intel Par. Studio 2017/2018, PGI, Cray PE
2. CUDA 8.0 / 9.1
3. OpenMPI 1.10.x and 3.x, MVAPICH2

# Other requirements

- Be able to re-use the installed modules e.g. do not rebuild OpenBLAS or protobuf unless knowing what you are doing

- Make sure that easyconfigs work with the Bright toolchains (happy that it is a separate repo!)


- Packages built by the admin should be installed in /cm/shared

- Ability to generate packages (FPM)

- Ability to generate containers

# Benefits for Bright and partners

- Reduce the support load by leveraging the EasyConfigs collection

- Allow us to focus more on features that matter than building binaries to support uncommon use cases e.g. DL on very old hardware

- Provide optimized packages for each of our HW OEM partners

- Expand EasyBuild's user base

- Enterprise support

- Contribute (mostly) DL-related  EasyConfigs back to the community and expand EasyBuilds availability of supported frameworks

- Report back issues

- Bright Cluster Manager 8.1 released last Monday ships EasyBuild 3.5.1

yum/zipper/apt-get install cm-EasyBuild

module load cm-EasyBuild

Bright DL  EasyConfig repo on Github soon