# AI for HPC and HPC for AI Workflows: The Differences, Gaps and Opportunities with Data Management

## @SC Asia 2018

**Rangan Sukumar, PhD**
**Office of the CTO, Cray Inc.**

# Safe Harbor Statement

This presentation may contain forward-looking statements that are based on our current expectations. Forward looking statements may include statements about our financial guidance and expected operating results, our opportunities and future potential, our product development and new product introduction plans, our ability to expand and penetrate our addressable markets and other statements that are not historical facts. These statements are only predictions and actual results may materially vary from those projected. Please refer to Cray's documents filed with the SEC from time to time concerning factors that could affect the Company and these forward-looking statements.

# AI for HPC and HPC for AI @ Cray



**CRAY URIKA-GX** → **CS-STORM 500NX** → **CS-STORM 500GT** → **CRAY XC SERIES** **CRAY URIKA-XC**

**Integrated Analytics and AI platform for Data Preparation and Machine Learning**

**Dense GPU systems with broad support for NVIDIA® Tesla® Accelerators and FPGAs**

**Scalable high performance supercomputers with Analytics and AI/DL**

# AI for HPC and HPC for AI: Today's talk

**CRAY**

**Deloitte.**

**Big Four Consulting Firm**

**Urika-GX**

"Teaming with Cray was a clear choice and allows for versatility combined with speed to tackle big data problems"

-- Deloitte Advisory Cyber Risk Services

**Top 5 Global Pharma**

**CS-Storm Dense GPU Cluster**

Supporting core research and development in areas including chem-informatics and large machine and deep learning workloads

**Fortune 20 Global Technology Company**

**Cray CS-Storm Dense GPU Cluster**

Recent win to support machine and deep learning workloads including autonomous vehicles

**Stanford University**

**The Stanford Research Computing Facility**

**Cray CS-Storm "XStream" Dense GPU Cluster**

Research in astrophysics, structural biology and bioinformatics, materials modeling, and climate
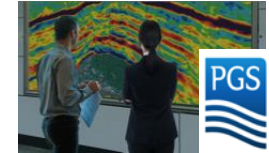
Recent breakthrough in 3D deep convolutional neural networks for amino acid environment similarity analysis

**Argonne NATIONAL LABORATORY**

**Argonne National Laboratory**

**Cray XC40**

Predicting how specific patients and tumors respond to different types of drugs using the scalable deep neural network code called CANDLE — or CANcer Distributed Learning Environment

**PGS**

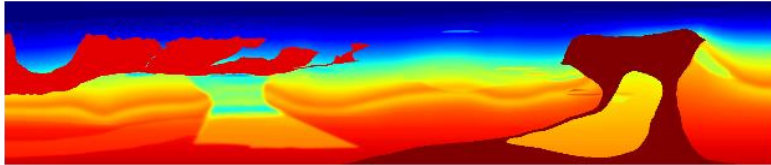**Leading Seismic Processing Services Company**

**Cray XC40**

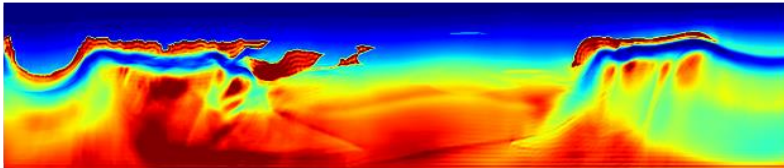Machine Learning at Scale for Full Waveform Inversion

PGS applied machine learning optimization techniques such as regularization and steering to determine the velocity model in a Full Waveform Inversion seismic imaging workload.
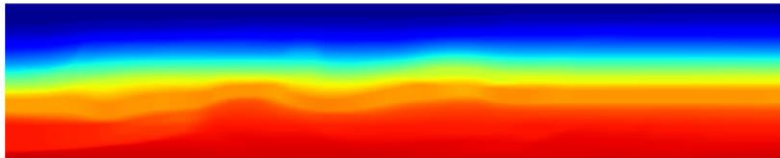
# Use-case: AI for HPC Application



**"Ground truth" Benchmark**



Synthetic benchmark with known subsurface geometry and seismic reflection data.

**Conventional FWI**



Conventional FWI attempts to derive a more accurate velocity model.
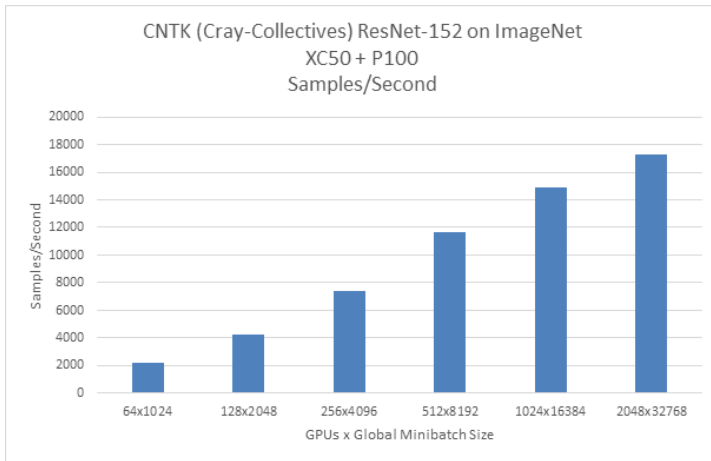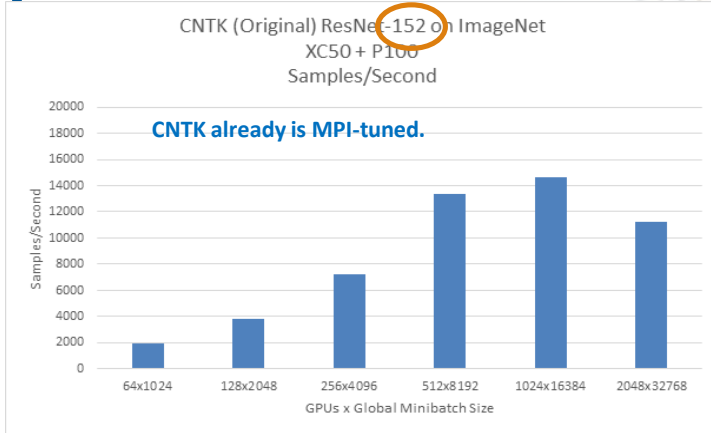
**FWI with Machine Learning**



Using machine learning (regularization and steering) to guide the convergence process.

# Use-case: HPC for AI Application



"Piz Daint is a supercomputer with Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , 4888 NVIDIA Tesla P100."



## CNTK (Original) ResNet-152 on ImageNet
### XC50 + P100
### Samples/Second

**CNTK already is MPI-tuned.**

## CNTK (Cray-Collectives) ResNet-152 on ImageNet
### XC50 + P100
### Samples/Second

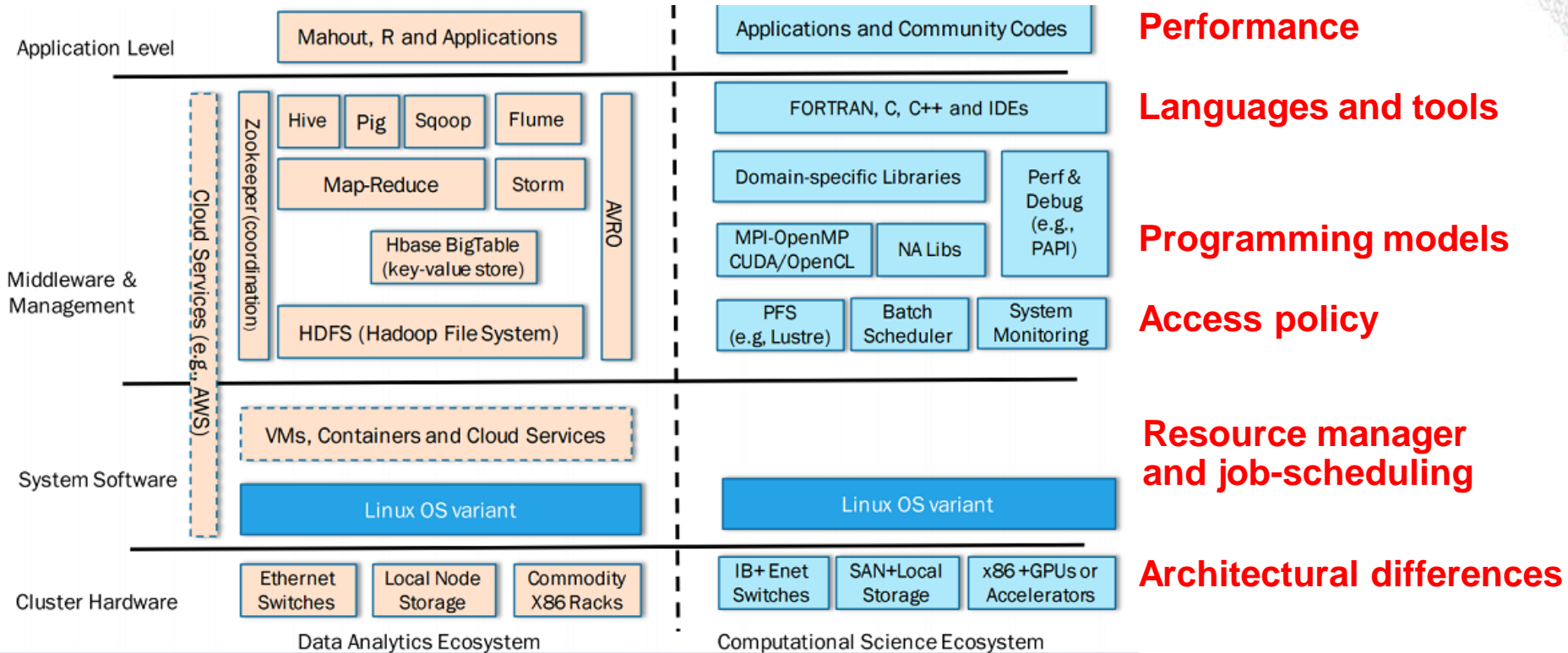## Microsoft, Cray claim deep learning breakthrough on supercomputers

Project could allow larger and more complex deep learning workloads on supercomputers.

By Steve Ranger | December 7, 2016 -- 12:09 GMT (04:09 PST) | Topic: Data Centers

# The State of Practice: Tale of Two Ecosystems

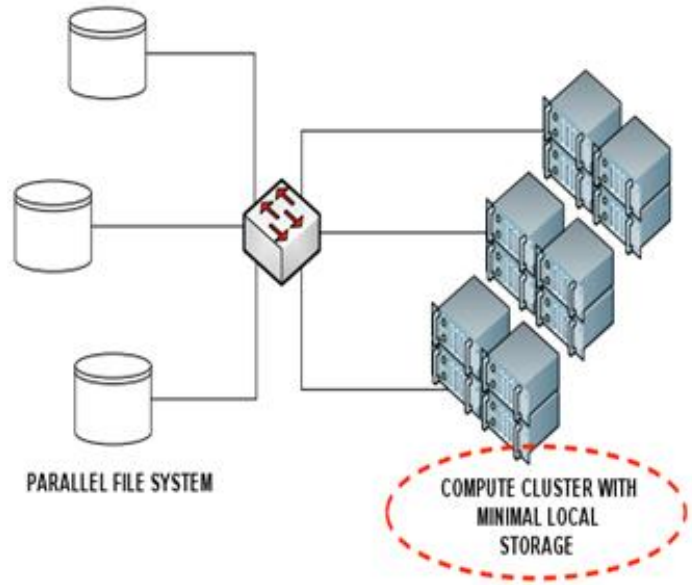

J. Dongarra et al., Exascale computing and Big Data: The next frontier, ACM Communications 2015

# The State of Practice: AI for HPC and HPC for AI



## Scientific Computing
### HPC Architectures

PARALLEL FILE SYSTEM

COMPUTE CLUSTER WITH MINIMAL LOCAL STORAGE

## Enterprise Computing
### Commodity Architectures

ETHERNET

COMPUTE/DATA CLUSTER

**Make assumptions about direction of data flow and requirements of data sizes to be moved**

# The State of Practice: Tale of Two Ecosystems

| | Scientific Computing | Enterprise Computing |
|---|---|---|
| **Primarily used for** | Solving equations | Search/Query, Machine learning |
| **Philosophy** | Send data to compute | Send compute to data |
| **Efficiency via** | Parallelism | Distribution |
| **Scaling expectation** | Strong (scale-up) | Weak (scale-out) |
| **Programming model** | MPI, OpenMP, etc. | Map-reduce, SPMD, etc. |
| **Popular languages** | FORTRAN, C++, Python | Java, Scala, Python, R |
| **Design strength** | Multi-node communication using an interconnect | Built-in job fault tolerance over Ethernet |
| **Access model** | On-premise | Cloud |
| **Preferred algebra** | Dense Linear | Set-theoretic / Relational |
| **Memory access** | Predictable | Random |
| **Storage** | Centralized, POSIX/RAID | Decentralized, Duplication |

# Terminologies: Tale of Two Ecosystems

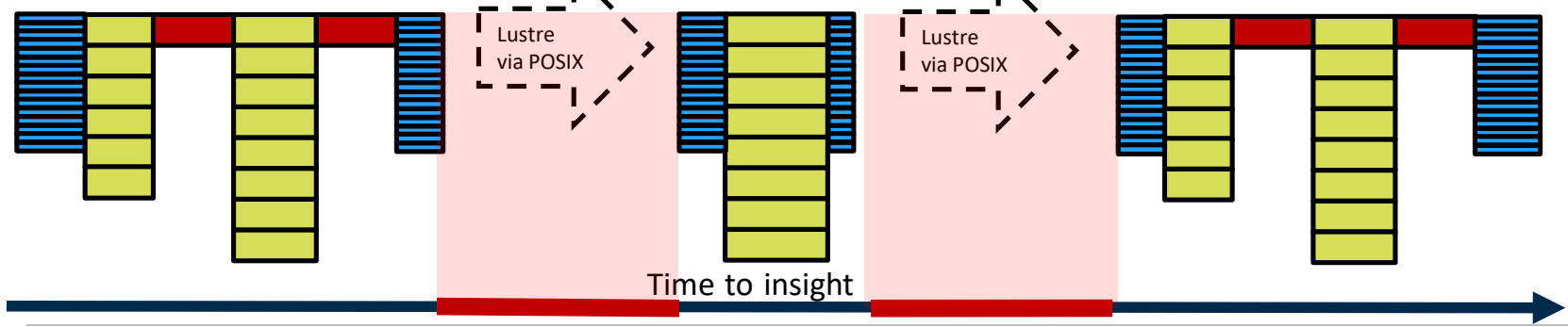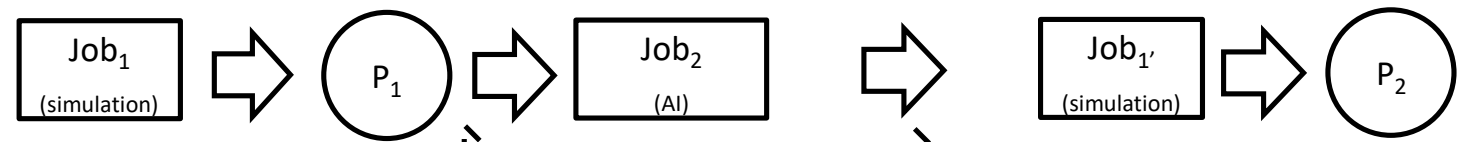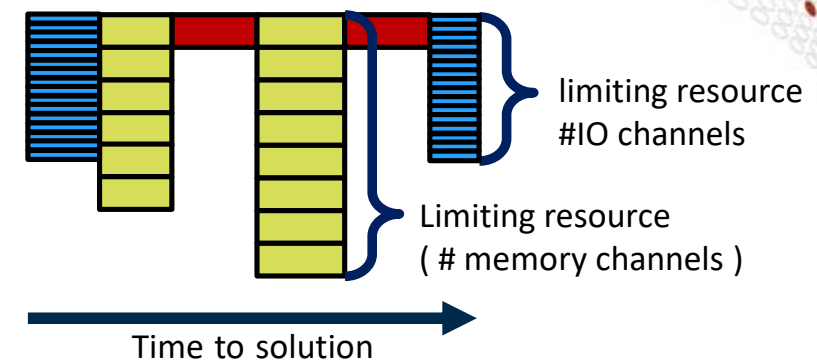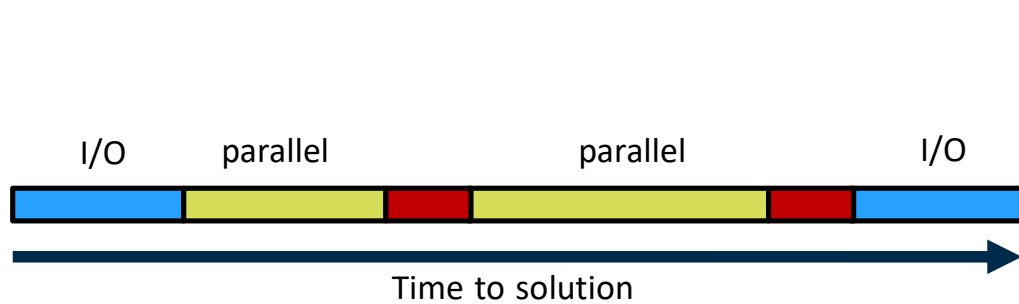| | Scientific Computing | Enterprise Computing |
|---|---|---|
| **Data (Structured)** | Vector, Matrix, Tensor | Table, Key-Values, Objects |
| **Data (Unstructured)** | Mesh, Images (Physics-based) | Documents, Images (Camera) |
| **Visualization** | Voxel, Surface, Point Clouds | Word Cloud, Parallel Coordinates, BI Tools |
| **Validation** | Cross-validation (ROC curves, statistical significance) | Manual / Subject matter expert, A/B testing |
| **Extract, Transform, Load** | Fourier, Wavelet, Laplace, etc. Cartesian, Radial, Toroidal, etc. | File-format transformations e.g. CSV to VRML |
| **Search (Query)** | Properties such as periodicity, self-similarity, anomaly, etc. | SQL, SPARQL, etc. (Sum, Average, Group by) |
| **Funding Model** | Non-profit grand challenge (Answer matters) | Value-driven (Cost matters) |

Sukumar, S. R., et al., (2016, December). Kernels for scalable data analysis in science: Towards an architecture-portable future. *In the Proc. Of the 2016 IEEE International Conference on Big Data,* pp. 1026-1031.
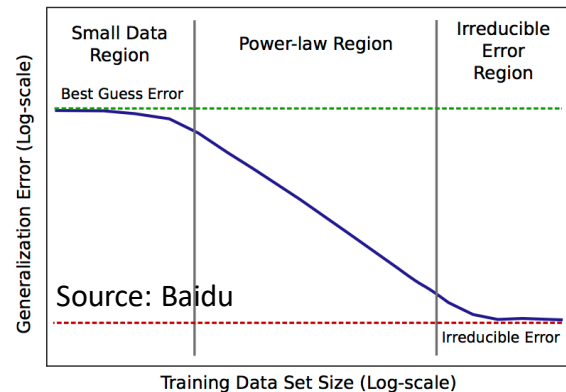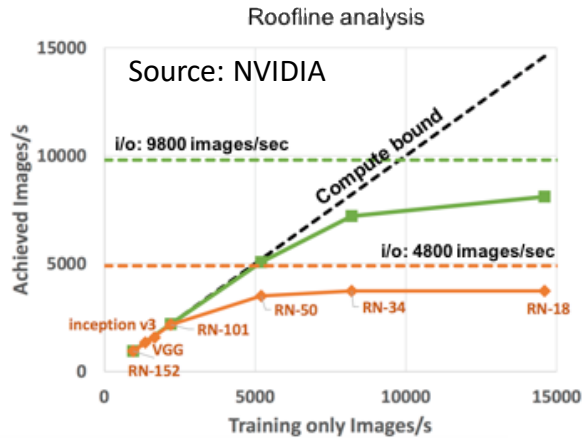
# Deep Learning: Tale of Two Ecosystems

| | Scientific Computing | Enterprise Computing |
|---|---|---|
| **Model** | Domain-specific | CNN, RNN, LSTM, GAN etc. |
| **Baseline** | Theoretic e.g. Navier Stokes | Humans, Other ML algorithms |
| **Parallelism** | Model, Ensemble | Data |
| **Use Case** | Computational Steering Proxy models | Speech, Test Image interpretation Hyper-personalization |
| **Source File System** | Lustre and GPFS | HDFS, S3, NFS etc. |
| **Figure of Merit** | Interpretability, Feasibility | Time-to-accuracy, Model-size |
| **Training Data** | $O(GBs)$ per sample, $O(10^3)$ samples, $O(10)$ categories | $O(KBs)$ per sample, $O(10^6)$ samples, $O(10^4)$ categories |
| **Data Model** | HDF5, NETCDF | Relational, Document, Key-Value |

# The Convergence: AI for HPC



I/O    parallel    parallel    I/O

Time to solution

limiting resource
#IO channels

Limiting resource
( # memory channels )

Time to solution

Job$_1$ (simulation) → P$_1$ → Job$_2$ (AI) → Job$_{1'}$ (simulation) → P$_2$

Lustre via POSIX

Lustre via POSIX

Time to insight

COMPUTE | STORE | ANALYZE

# The Convergence: HPC for AI

## Roofline analysis

Source: NVIDIA



Source: Baidu

## Opportunity for productivity with strong scaling

| ResNet-50 Success | Time-to-accuracy | How many GPUs? | Scalability Efficiency |
|---|---|---|---|
| Facebook (Caffe2) | 2 days<br>1 hour | 352 GPUs<br>256 | 90%<br>(large-batch) |
| IBM PowerAI (Caffe) | 50 minutes | 256 GPUs | 95%<br>(large-batch) |
| Google (TensorFlow) | ~24 hours | 64 TPUs | >90% |
| Preferred Networks (Chainer) | 15 minutes | 1000 GPUs | >90% |
| Cray @ CSCS (Tensorflow) | <14 minutes | 1000 GPUs | ~>95% |

**Productivity is performance and performance translates to productivity...**
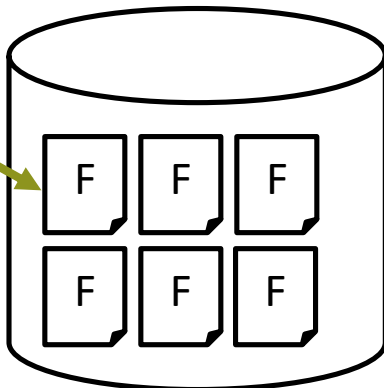
# Bottlenecks Today



**Potential off-node I/O requirement**

Other Nodes

**NFS Client w/CacheFS**

EDR IB

**Casual Copy In From Data Sources**

NAS (NetApp)

**Possible 1GB/s (future 8GB/s) From local SSD**

10GbE

?GbE

**Stage 2GB/s**

**2TB SSD RAID5 (write b/w ~half)**

COMPUTE | STORE | ANALYZE

# Bottlenecks Looking Ahead…

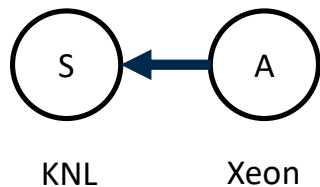| Figures-of-merit | State-of-practice | Projected 2 years ahead |
|---|---|---|
| Training-time to best accuracy | 5+ days | 2+ hours |
| Model Cost / TB (AWS GPUs) | ~$25K<br>(ResNet training on 80 GPUs for 5 days) | ~10K |
| Hardware Efficiency | O(~25 Gflops)<br>Network Depth: Flops::20x: 16x<br>(based on AlexNet-2012 and ResNet-2015) | O(Teraflops) |
| Statistical Efficiency | O(~25 Gflops)<br>Depth: Accuracy:: 20x:13+<br>(based on AlexNet-2012 and ResNet-2015) | O(Teraflops) |
| Need for compute as data grows | O(~465 Gflops)<br>Data: Flops: Error:: 2x: 5x: 3+<br>(based on DeepSpeech1 and DeepSpeech2) | O(Petaflops) |
| Training Cadence | ~ Monthly | ~ Daily |
| # of models per organization | 1x | 10-100x |

# Solution: System to Eco-system Thinking

CRAY

| **Hardware** | **Software** | **Ecosystem** |
|---|---|---|

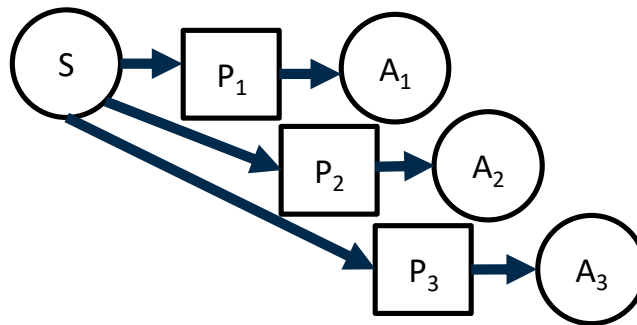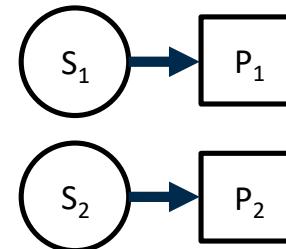| | **System** | **Function** | **Community Productivity** |
|---|---|---|---|
| **Facility Performance** | **Utilization** — Peak vs. Sustained, Performance per $ | **Application/Codes** — e.g. Deep Learning, Graph analytics | **Domain-specific Creativity** — Is there an ecosystem of sustainable community (open-source) engagement that enables vertical segments? |
| **System Performance** | **Reliability** — Faults, MTTF, Uptime    **Scalability** — Weak and strong | **Kernel/Motif** — e.g. DGEMM, SYRK, ReLU, inner product | **Code Portability** — Does a user have to rewrite code? Does vendor support code porting for novel architectures? |
| **Multi-node Performance** | **System Architecture** | **Programming Model** — e.g. MR, PGAS, GRPC | **Programmability** — Does an end-user have to learn a new language or can they launch jobs with modern tools (e.g. notebooks)? |
| **Node Performance** | **Interconnect** — eth, InfiniBand, Aries    **Provisioning** — Mesos, Moab, SLURM | **Libraries** — e.g. MKL, CUDA, libSci    **Collectives** — e.g. NCCL, MPI | **Data Pre-Processing** — Does system offer tools to optimize ETL wall-time? |
| | **Node Architecture** — # of xPUs+ cache + memory + network | | |
| **Component Performance** | **Disk** — Latency    **Memory** — Capacity, Latency    **xPU** — Speed    **i/o** | **Data Structure** — e.g. matrix, sequences, unstructured grids | **Data Movement** — Does system provide ability to run multiple frameworks/applications on the same data? |

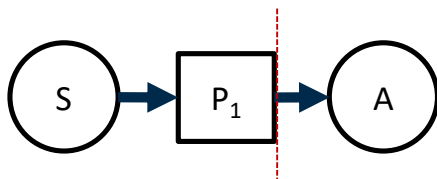# Solution: Communication-Aware Data Objects
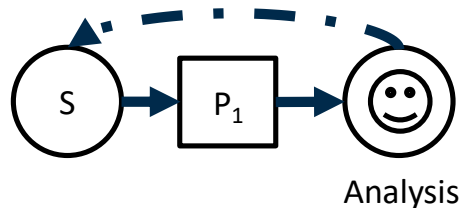


Async DAG-execution

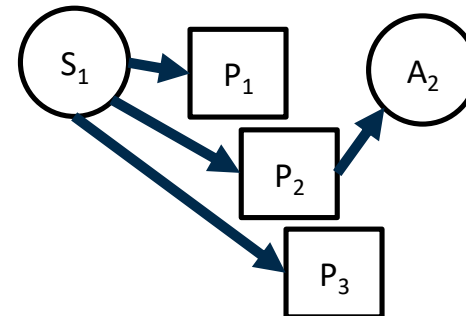Dataflow semantics

Data Management & curation

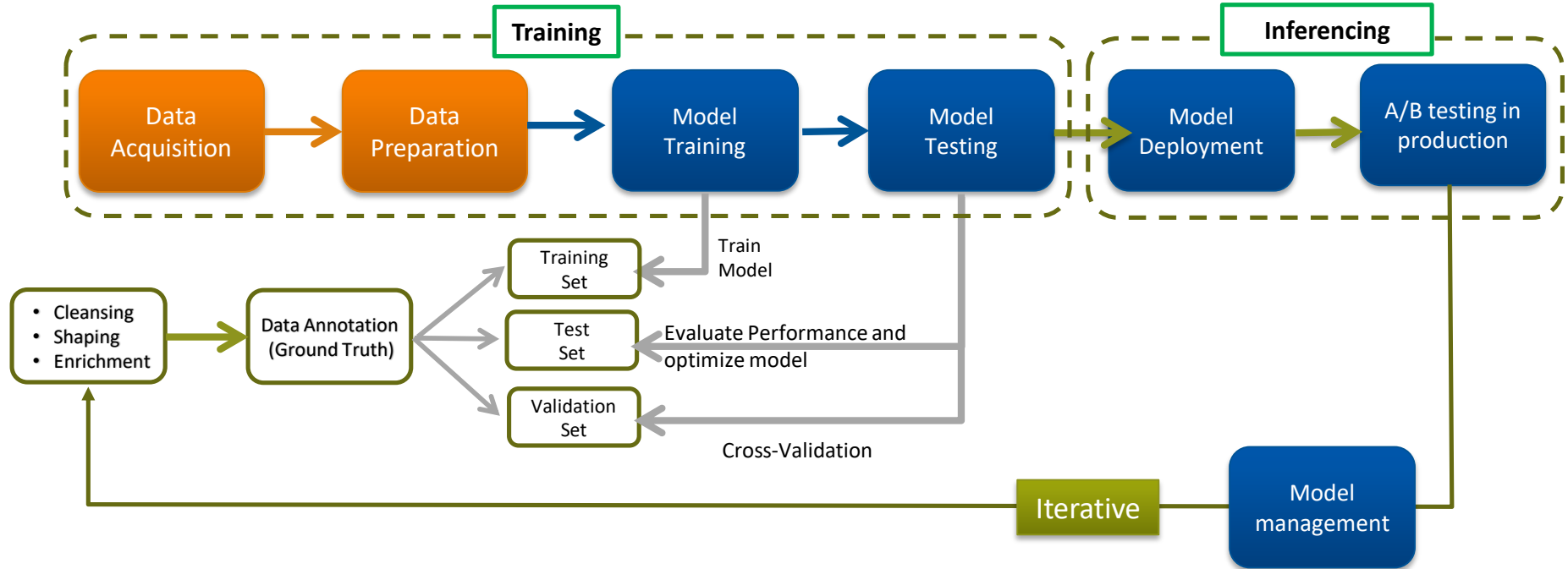Data-driven notification

Human in-the-loop

Data-dependent workflow

# Solution: HPC Best Practices for Data Management

Making data access and I/O methods available / relevant to all levels of the software stack.

| | Operating Systems | Runtimes | Systems Software | Programming Environments | Applications | Workflows |
|---|---|---|---|---|---|---|
| HBM | | memkind | | | memkind | |
| GPU MEM | | CUDA | CUDA | PTX | CUDA | |
| DRAM | C / ASM | C / ASM | C | C / ASM | C / Fortran | |
| NV-DIMM | | pmem | pmem | | pmem / pmemkind | pmem / pmemkind |
| LOCAL SSD | | | | | POSIX | POSIX |
| BURST BUFFER | | | | | DSL (e.g Datawarp) | DSL (e.g Datawarp) |
| Network SSD | | | | | POSIX | POSIX |
| DISK / PFS | POSIX / swap | | | | POSIX / MPI-IO | POSIX |
| TAPE | | | | | | TSM |
| CLOUD | | | | | | S3 |

# Solution: End-to-End Thinking with Benchmarks

# Future: Integration of Storage, Memory and Compute

- **General purpose flexibility**
  - Commodity-like configurations
- **Seamless heterogeneity**
  - CPUs, GPUs, FPGAs, ASICs
- **High-performance interconnects for data centers**
  - MPI and TCP/IP collectives, compute on the network
- **Unified software stack**
  - Programming environment for performance and productivity
- **Workflow optimization**
  - Match growth in compute and data with I/O