

# Architectural Challenges Emerging From The Convergence of Big Data, HPC and AI

PDSW-DISC Keynote @ SC18

✉ ssukumar@cray.com

🐦 @Rangan\_Sukumar

in <https://www.linkedin.com/in/rangan/>



CRAY



# SAFE HARBOR STATEMENT

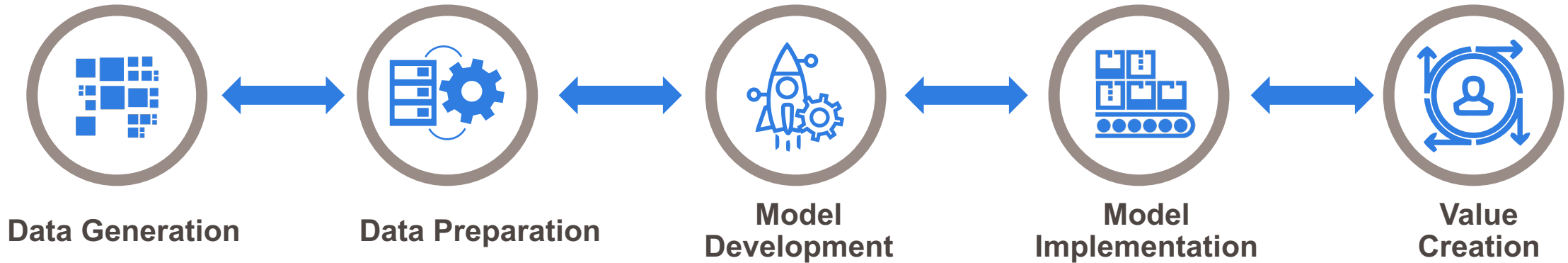
This presentation may contain forward-looking statements that are based on our current expectations. Forward looking statements may include statements about our financial guidance and expected operating results, our opportunities and future potential, our product development and new product introduction 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.





# ARCHITECTING A CONVERGENT SYSTEM



## Levels of Architectural Maturity Towards Convergence

- **Level 1:** Can a system run HPC, Big Data and AI applications/workloads?
- **Level 2:** Can a system execute a “convergent workflow” consisting of HPC, Big Data and AI tools, codes and frameworks in reasonable time?
- **Level 3:** Can a system accelerate/scale the workflow when required?
- **Level 4:** Given a workflow, is this the top “performant” system one can build?



# THE CHALLENGE OF ARCHITECTING HPC SYSTEMS



- **Design specifications**

- maximize(performance-per-\$)
- minimize(\$-to-insight)
- maximize(architected performance \* community productivity)  $\leq$  budget
- minimum(benchmark-performance)  $\geq$  scaling factor
- maximum(app-to-app performance variation)  $\leq$  epsilon
- minimize(operating costs  $\sim$  power, downtime, human resources)

- **Figures-of-merit**

- FLOPS, Programmability, Utilization, Benchmark, Scientific Innovation, ...



# THE CHALLENGE OF ARCHITECTING BIG-DATA SYSTEMS



- **Design specifications**

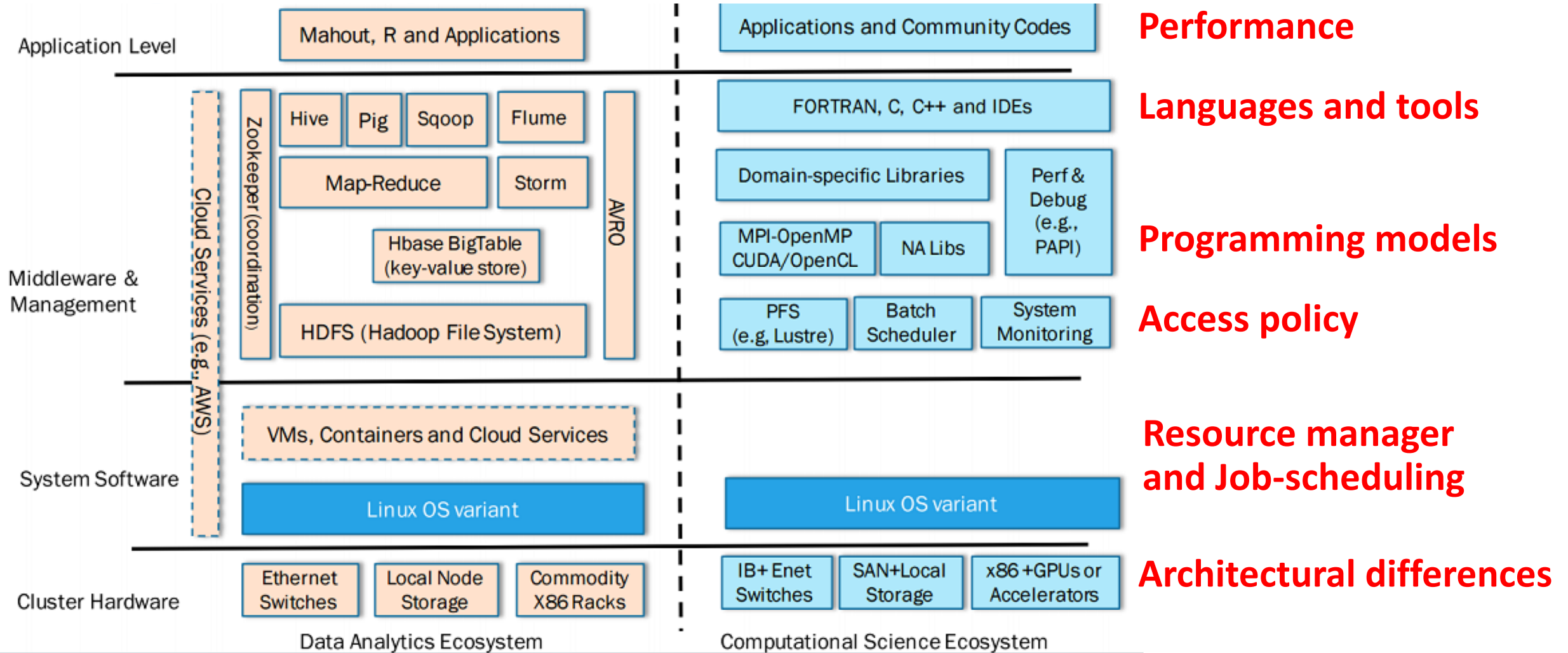
- maximize(ROI-per-byte)
- maximize(capacity \* consistency \* availability \* fault-tolerance)
- maximize(open-source tool support)
- minimize(time-to-prototype + time-to-production)
- minimize(security risk)
- minimize(operating costs ~ power, downtime, human resources)

- **Figures-of-merit**

- ROI, Elasticity, Multi-tenancy, Ease-of-use, Time-to-accuracy,...



# TODAY: IT IS THE TALE OF TWO ECOSYSTEMS



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



# REQUIREMENTS: MODE OF OPERATION



	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-like
Preferred algebra	Dense Linear	Set-theoretic / Relational
Memory access	Predictable	Random
Storage	Centralized, POSIX/RAID	Decentralized, Duplication

# REQUIREMENTS : WORKFLOWS + WORKLOAD



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



# REQUIREMENTS: PROCESS AND DEPLOYMENT



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

# REQUIREMENTS: USER EXPERIENCE



	Scientific Computing	Enterprise Computing
Programming	Vendor libraries and compilers	Open-source services and APIs
Preferred Deployment	Bare metal	Virtualized, Containers
Popular Architecture	Homogenous	Heterogenous
“Systems” Literacy	High	Low
Scheduling	Batch	Interactive, Persistent
Resource Managers	SLURM, SGE, etc.	Mesos, Kubernetes, etc.
Data in	Files	Databases (in-memory, schema)
Software	Write-once Run-many	Write-Many Run-Many
Access Interface	Terminal, SSH	Jupyter, Web-based IDEs



# IS CONVERGENCE NECESSARY?

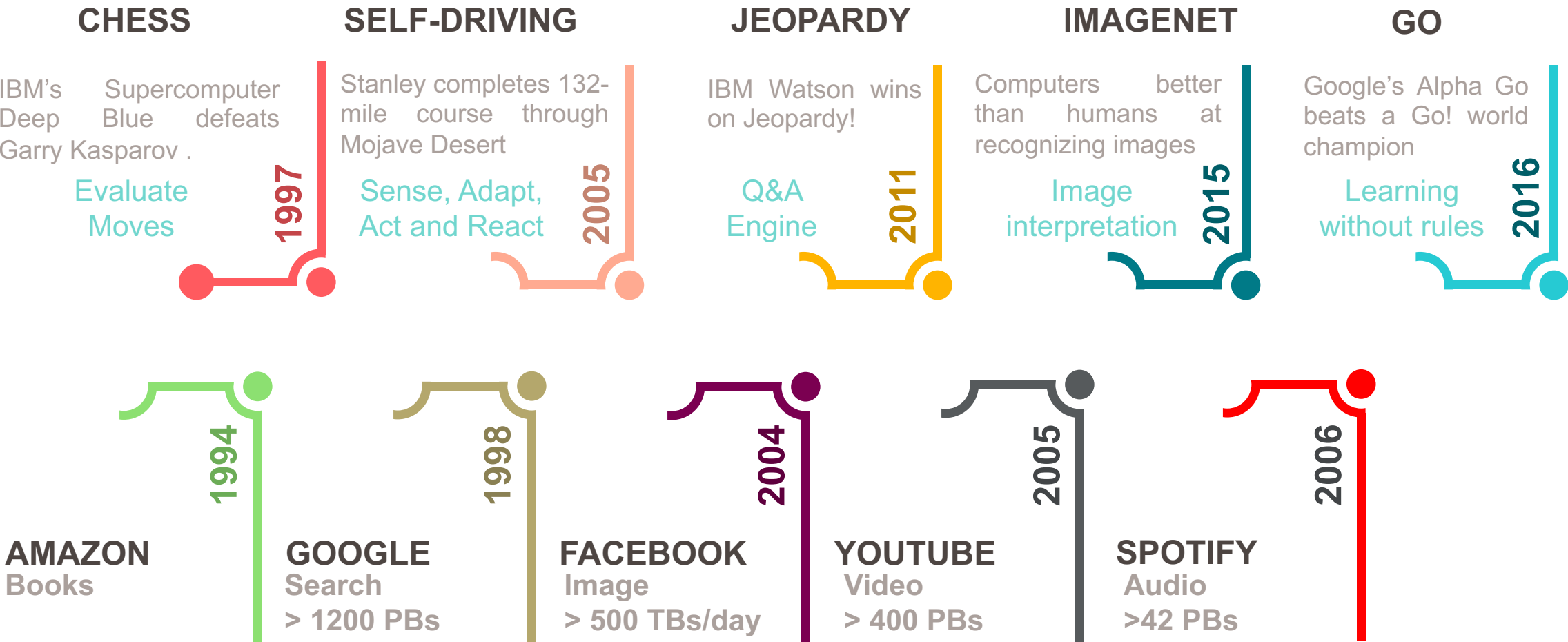
AI methods benefit from HPC best practices



# CONVERGENCE HEADLINES: BIG DATA + AI



## MILESTONES IN ARTIFICIAL INTELLIGENCE



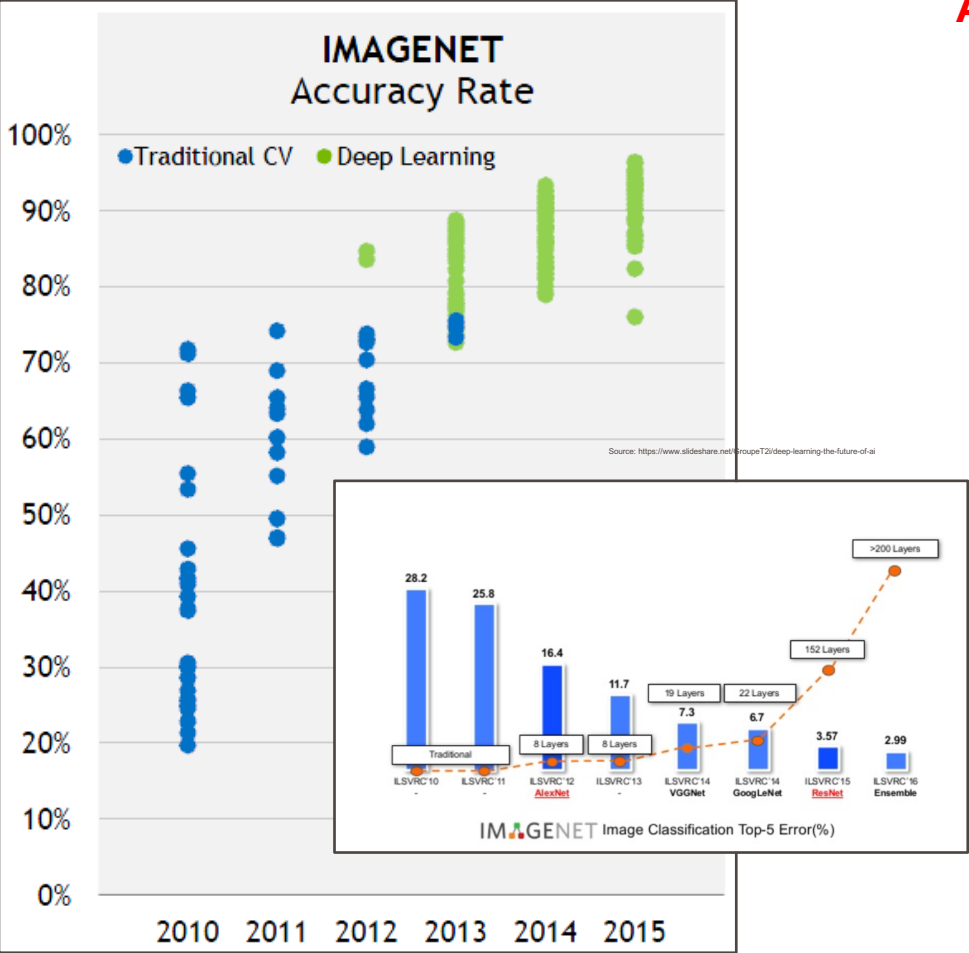
GROWTH IN UNSTRUCTURED BIG DATA



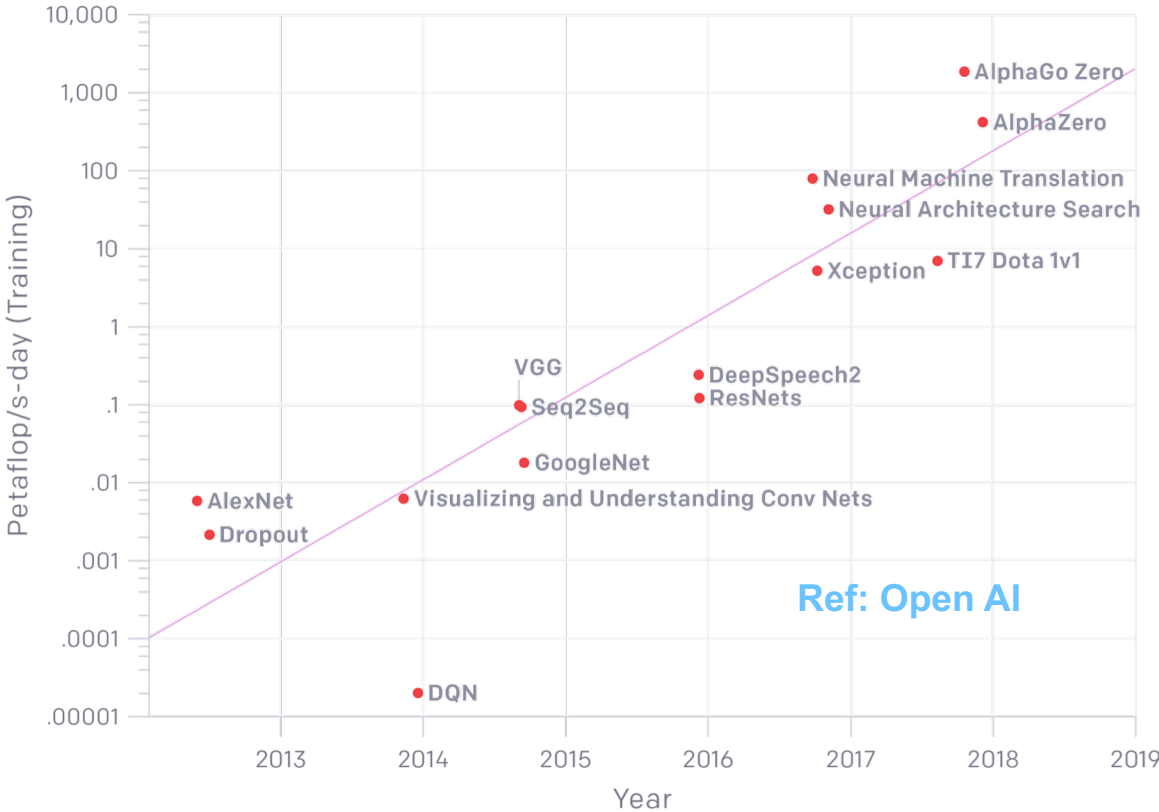
# HPC-SCALE REQUIREMENTS AT AI PRACTITIONERS



Source: NVIDIA-ces-2016-press-conference



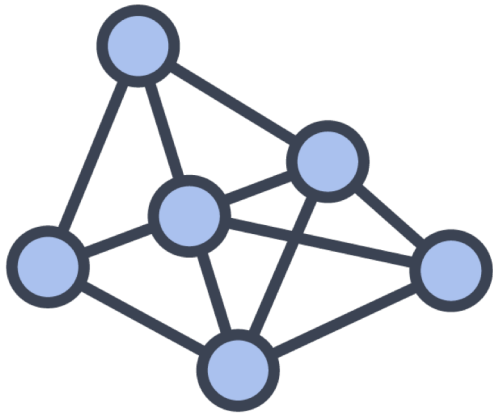
## ALEXNET TO ALPHAGO ZERO: A 300,000x INCREASE IN COMPUTE



# BENEFITS OF HPC ADOPTION



## Graph Analytics



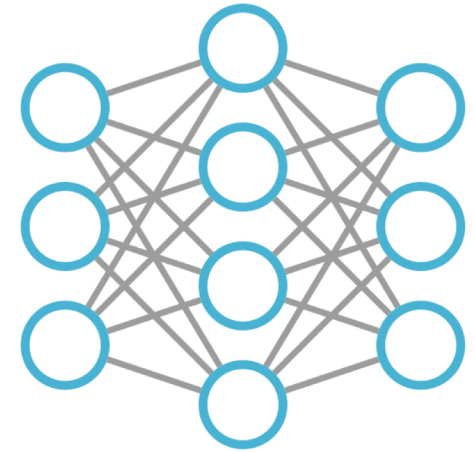
Handle 1000x bigger datasets with a 100x better speed-up with queries

## Matrix Methods

$$\begin{bmatrix} \dots & \dots \\ \vdots & \vdots \\ \dots & \dots \end{bmatrix} * \begin{bmatrix} \vdots & \dots & \vdots \\ \vdots & \dots & \vdots \\ \vdots & \dots & \vdots \end{bmatrix} \approx \begin{bmatrix} \vdots & \dots & \vdots \\ \vdots & \dots & \vdots \\ \vdots & \dots & \vdots \end{bmatrix}$$

Get 2-26x over Big Data Frameworks like Hadoop, Spark (for the same cluster-size)

## Deep Learning



95%+ scalability efficiency that can reduce training time from days to hours

### Best practices:

- Application fine-tuning / Performance optimization
- High-performance interconnect
- Algorithmic cleverness to trade compute and i/o
- Overlap compute and i/o with programming model



# HPC BUILDS THE MODELS OF TOMORROW



MORE DATA, BIGGER MODELS, NEED FOR MORE EFFICIENT AND PRODUCTIVE HARDWARE

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

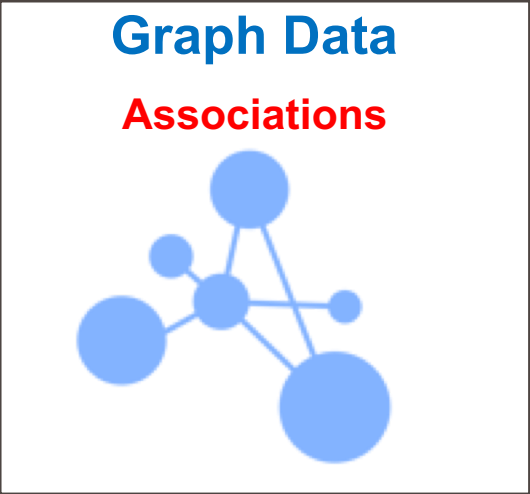
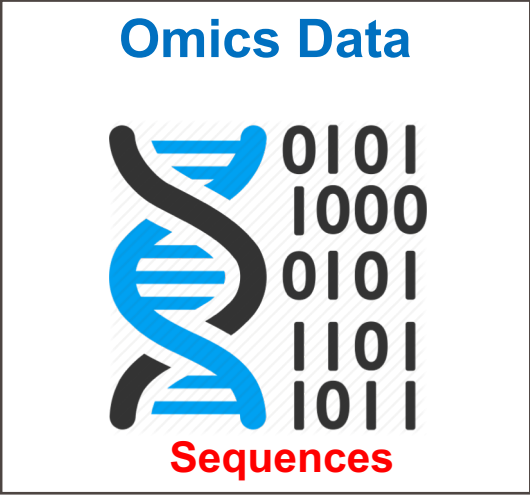
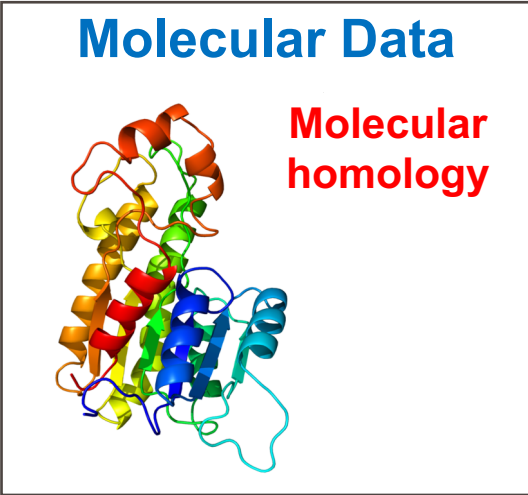
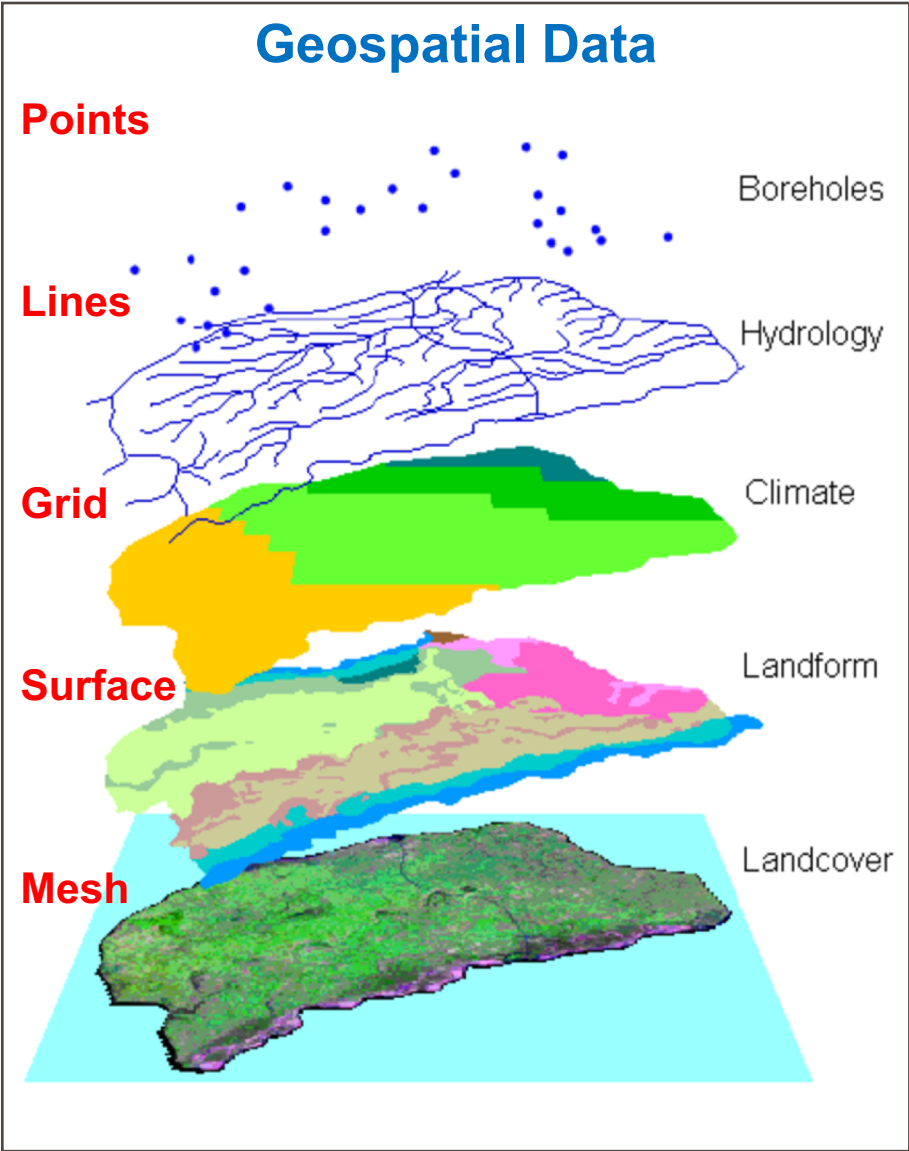


# IS CONVERGENCE NECESSARY?

AI methods offer tremendous capabilities for scientific data



# DIFFERENT SHAPES OF SCIENTIFIC DATA

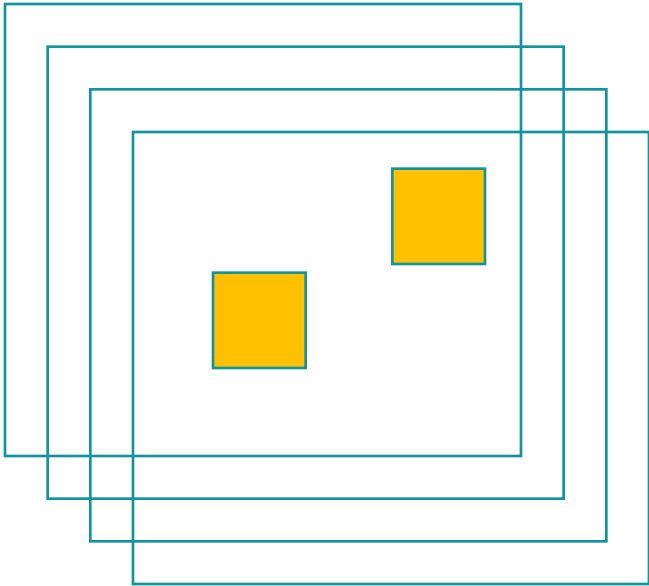


2D, 3D, 4D volumes, Higher precision (32, 64 bit), Higher # channels (3, 16, 1024), Sparse + Dense, Resolution



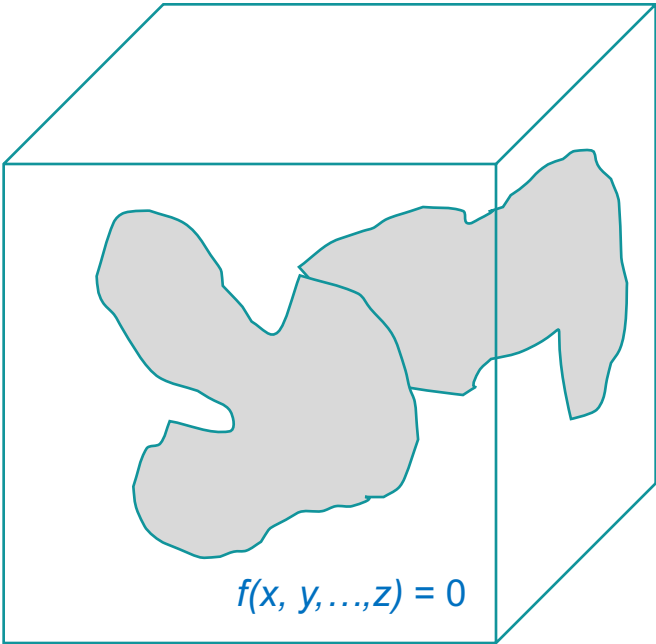
# DIFFERENT SEARCH SPACES OF SCIENTIFIC DATA

## Feature-based



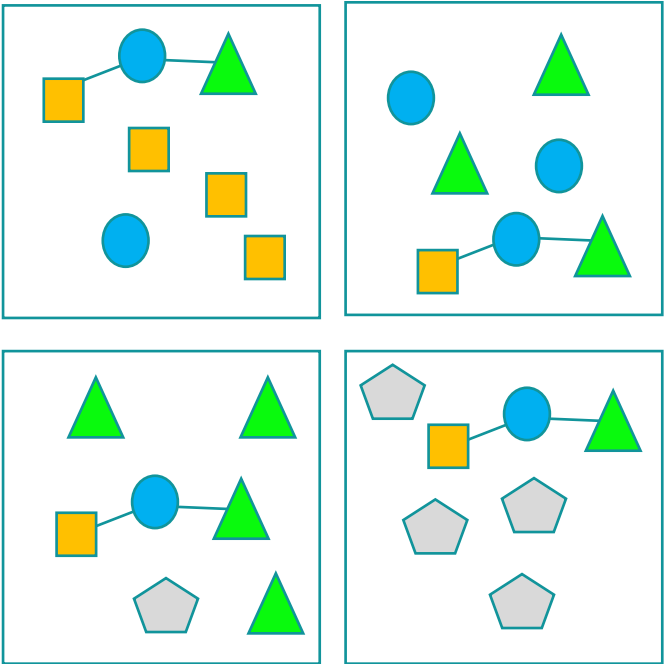
Structured Data  
Feature-space is ill-posed  
Search is well-defined

## Function-based



Unstructured Data  
Feature-space is theoretical  
Search is empirical

## Pattern-based



Semi-Structured Data  
Feature-space to be discovered  
Search is P or NP-hard

**Opportunity to create the "models of tomorrow"**

# MODELS EXPLORE AUTO-ENGINEERED FEATURE SPACES

Path towards explain-ability



Capsule Networks  
Mixture of Experts  
Neural Collaborative Filtering  
Block-Sparse LSTM

ReLu  
Dropout  
BatchNorm  
Decoder/Encoder  
Pooling  
LSTM  
GRU

WaveNet  
Attention  
Beam Search

Domain-specific Models  
Reinforcement Learning  
Generative Models

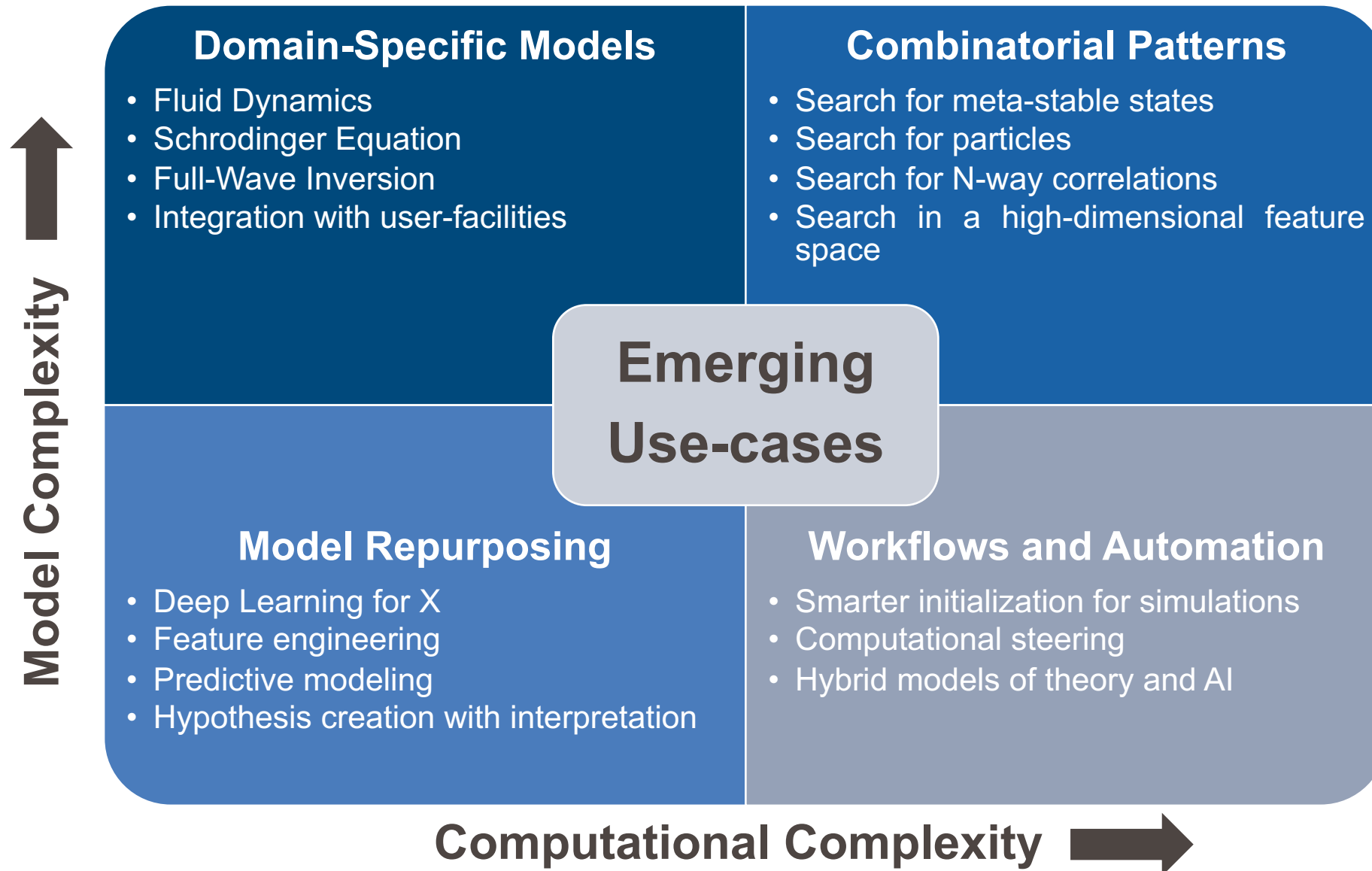
DQN  
Simulation  
DDPG

Recurrent Neural Networks  
Convolutional Neural Models

ConditionalGAN  
3D-GAN  
CoupledGAN  
SEGAN  
MedGAN



# AI ADOPTION IS INCREASING IN THE SCIENCES





Deep Learning at 15PF: Supervised and Semi-Supervised Classification for Scientific Data

Thorsten Kurth\*, Jian Zhang†, Nadathur Satish†, Ioannis Mitliagkas†, Evan Racah\*, Mostofa Ali Patwary†, Tareq Malas‡, Narayanan Sundaram†, Wahid Bhimji\*, Mikhail Smolenskiy\*, Josh DeSa\*, Mikhail Shiryayev¶, Srinivas Sridharan||, Prabhat\*, P

**Abstract**—This paper presents the first, 15-PetaFLOP Deep Learning system for solving scientific pattern classification problems on contemporary HPC architectures. We develop supervised convolutional architectures for discriminating signals in high-energy physics data as well as semi-supervised architectures for localizing and classifying extreme weather in climate data. Our Intelcapped implementation obtains ~2TFLOP/s on a single Cori Phase-II Xeon-Phi node. We use a hybrid strategy employing synchronous node-groups, while using asynchronous

poised to have a m there are unique c first. The primary cha tities of complex, Deep Learning im verge on O(10) C datasets are TBs-R contain dozens of c

Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteva1\*, Brett Kuper1\*, Roberto A. Novoa2,3, Justin Ko2, Susan M. Swetter2,4, Helen M. Blau5 & Sebastian Thrun6

Skin cancer, the most common human malignancy1-3, is primarily diagnosed visually, beginning with an initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Automated classification of skin lesions using images is a challenging task owing to the fine-grained variability in the appearance of skin lesions. Deep convolutional neural networks (CNNs)4,5 show potential for general and highly variable tasks across many fine-grained object categories6-11. Here we demonstrate classification of skin lesions using a single CNN, trained end-to-end from images directly, using only pixels and disease labels as inputs. We train a CNN using a dataset of 129,450 clinical images—two orders of magnitude larger than previous datasets12—consisting of 2,032 different diseases. We test its performance against 21 board-certified dermatologists on biopsy-proven clinical images with two critical binary classification use cases: keratinocyte carcinomas versus benign seborrheic keratoses; and malignant melanomas versus benign nevi. The first case represents the identification of the most common cancers, the second represents the identification of the deadliest skin cancer. The CNN achieves performance on par with all tested experts across both tasks, demonstrating an artificial intelligence capable of classifying skin cancer with a level of competence comparable to dermatologists. Outfitted with deep neural networks, mobile devices can potentially extend the reach of dermatologists outside of the clinic. It is projected that 6.3 billion smartphone subscriptions will

images (for example, smartphone images) exhibit variability in factors such as zoom, angle and lighting, making classification substantially more challenging13,24. We overcome this challenge by using a data-driven approach—1.41 million pre-training and training images make classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and extraction of domain-specific visual features before classification. By contrast, our system requires no hand-crafted features; it is trained end-to-end directly from image labels and raw pixels, with a single network for both photographic and dermoscopic images. The existing body of work uses small datasets of typically less than a thousand images of skin lesions16,18,19, which, as a result, do not generalize well to new images. We demonstrate generalizable classification with a new dermatologist-labelled dataset of 129,450 clinical images, including 3,374 dermoscopy images. Deep learning algorithms, powered by advances in computation and very large datasets25, have recently been shown to exceed human performance in visual tasks such as playing Atari games26, strategic board games like Go27 and object recognition6. In this paper we outline the development of a CNN that matches the performance of dermatologists at three key diagnostic tasks: melanoma classification, melanoma classification using dermoscopy and carcinoma classification. We restrict the comparisons to image-based classification. We utilize a GoogleNet Inception v3 CNN architecture9 that was pre-trained on approximately 1.28 million images (1,000 object categories)

Exascale Deep Learning for Climate Analytics

Thorsten Kurth  
Lawrence Berkeley National Laboratory  
Berkeley, CA 94720, USA

Sean Treichler  
NVIDIA  
Santa Clara, CA 95051, USA

Data-driven Fluid Simulations using Regression Forests

Lubor Ladický\*† SoHyoon Jeong\*† Barbara Solenthaler† Marc Pollefeys† Markus Gross†  
ETH Zurich ETH Zurich ETH Zurich ETH Zurich ETH Zurich  
Disney Research Zurich



doi:10.1038/nature21056

capable of simulating millions of particles in realtime. Our promising physics-based simulations in time-critical settings, where running

1 Introduction

Computing high-resolution fluid simulations with tradition of-the-art approaches is very challenging, as they require dous computational resources to compute a scene with m particles. The main bottleneck is the severe restriction on step size needed to guarantee stability, and thus simulations are typically in the range of hours to days on high-end co making it impossible to achieve high-resolution fluids in re

The standard Smoothed Particle Hydrodynamics (LUC (SPH) method approximates continuous quantities in the Stokes differential equations using discrete particles with a ate smoothing kernel and replaces a continuous advective advection of particles. The approach does not deal with the pressibility constraint directly, which causes significant vis pleasant artifacts.

Recent work has addressed this problem and either enforce sity invariance condition or a divergence-free velocity f predictive-corrective scheme has been introduced where values are iteratively updated to satisfy the zero compress strain [Solenthaler and Pajarola 2009]. The performance i further improved by discretizing and iteratively solving ture Poisson equation [Hansen et al. 2013]. Most recently tion based fluid (PBF) approach has been presented [Mac Mueller 2013], where first all particles are advected, and tected to the manifold of feasible solutions by iteratively ing positions of particles to satisfy the incompressibility co PBF allows to use a larger time step compared to its counter density invariance condition has also been combined with scale scheme [Horvath and Solenthaler 2013] to further s the computation. Despite all these improvements, high-re fluids are still computed offline.

An alternative to particle-based approaches are grid-based ods that approximate continuous quantities on a discrete grid [Enright et al. 2002]. Incompressibility is enforced grid by solving the Poisson's equation, making the veloc divergence-free. To speed up grid-based simulations, the space can be restricted to simpler topology [Chentanez an 2010; Chentanez and Müller 2011]. To avoid discretization of grid-based methods, the hybrid FLIP model [Zhu and

Fast and Accurate Modeling of Molecular Atomization Energies with Machine Learning

Matthias Rupp,1,2 Alexandre Tkatchenko,3,2 Klaus-Robert Müller,1,2 and O. Anatole von Lilienfeld4,2,¶  
1Machine Learning Group, Technical University of Berlin, Franklinstr 28/29, 10587 Berlin, Germany  
2Institute of Pure and Applied Mathematics, University of California Los Angeles, Los Angeles, CA 90095, USA  
3Fritz-Haber-Institut der Max-Planck-Gesellschaft, 14195 Berlin, Germany  
4Argonne Leadership Computing Facility, Argonne National Laboratory, Argonne, Illinois 60439, USA  
(Dated: September 14, 2011)

We introduce a machine learning model to predict atomization energies of a diverse set of organic molecules, based on nuclear charges and atomic positions only. The problem of solving the molecular Schrödinger equation is mapped onto a non-linear statistical regression problem of reduced complexity. Regression models are trained on and compared to atomization energies computed with hybrid density-functional theory. Cross-validation over more than seven thousand small organic molecules yields a mean absolute error of ~10 kcal/mol. Applicability is demonstrated for the prediction of molecular atomization potential energy curves.

Solving the Schrödinger equation (SE),  $H\Psi = E\Psi$ , for assemblies of atoms is a fundamental problem in quantum mechanics. Alas, solutions that are exact up to numerical ory (DFT) level [2, 13, 14], any other training set or level of theory could be used as a starting point for subsequent ML training. Cross-validation on 7165 molecules yields which is an or-ting bonds or

Hierarchical attention networks for information extraction from cancer pathology reports

Shang Gao,1 Michael T Young,1 John X Qiu,1 Hong-Jun Yoon,1 James B Christian,1 Paul A Fearn,2 Georgia D Tourassi,1\* and Arvind Ramanathan1,\*

1Computational Science and Engineering Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA and 2Surveillance Informatics Branch, Division of Cancer Control and Population Sciences, National Cancer Institute, Bethesda, MD, USA

\*Corresponding Author: Arvind Ramanathan, Computational Sciences and Engineering Division, Health Data Sciences Institute, Oak Ridge National Laboratory, MS-6085, One Bethel Valley Road, Oak Ridge, TN 37831-6085, USA. Email: ramana-thana@ornl.gov. Phone: 865-576-7266. Fax: 865-241-0337

Received 12 June 2017; Revised 10 October 2017; Editorial Decision 15 October 2017; Accepted 26 October 2017

ABSTRACT

**Objective:** We explored how a deep learning (DL) approach based on hierarchical attention networks (HANs) can improve model performance for multiple information extraction tasks from unstructured cancer pathology reports compared to conventional methods that do not sufficiently capture syntactic and semantic contexts from free-text documents.

**Materials and Methods:** Data for our analyses were obtained from 942 deidentified pathology reports collected by the National Cancer Institute Surveillance, Epidemiology, and End Results program. The HAN was implemented for 2 information extraction tasks: (1) primary site, matched to 12 International Classification of Diseases for Oncology topography codes (7 breast, 5 lung primary sites), and (2) histological grade classification, matched to G1-G4. Model performance metrics were compared to conventional machine learning (ML) approaches including naive Bayes, logistic regression, support vector machine, random forest, and extreme gradient boosting, and other DL models, including a recurrent neural network (RNN), a recurrent neural network with attention (RNN w/A), and a convolutional neural network.

**Results:** Our results demonstrate that for both information tasks, HAN performed significantly better compared to the conventional ML and DL techniques. In particular, across the 2 tasks, the mean micro and macro F-scores for the HAN with pretraining were (0.852, 0.708), compared to naive Bayes (0.518, 0.213), logistic regression (0.682, 0.453), support vector machine (0.634, 0.434), random forest (0.698, 0.508), extreme gradient boosting (0.696, 0.522), RNN (0.505, 0.301), RNN w/A (0.637, 0.471), and convolutional neural network (0.714, 0.460).

**Conclusions:** HAN-based DL models show promise in information abstraction tasks within unstructured clinical pathology reports.

**Key words:** clinical pathology reports, information retrieval, recurrent neural nets, attention networks, classification

arXiv:1708.05256v1 [cs.LG] 17 Aug 2017



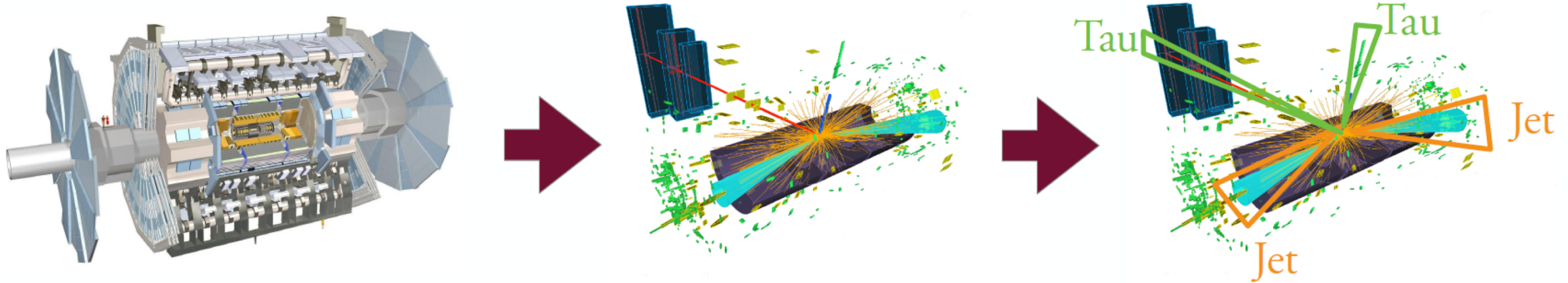
# EXAMPLES OF CONVERGENT WORKFLOWS

Big Data and AI @ HPC Centers



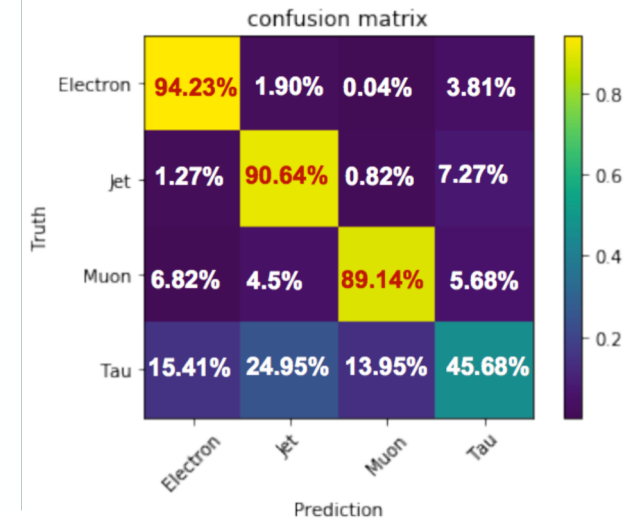
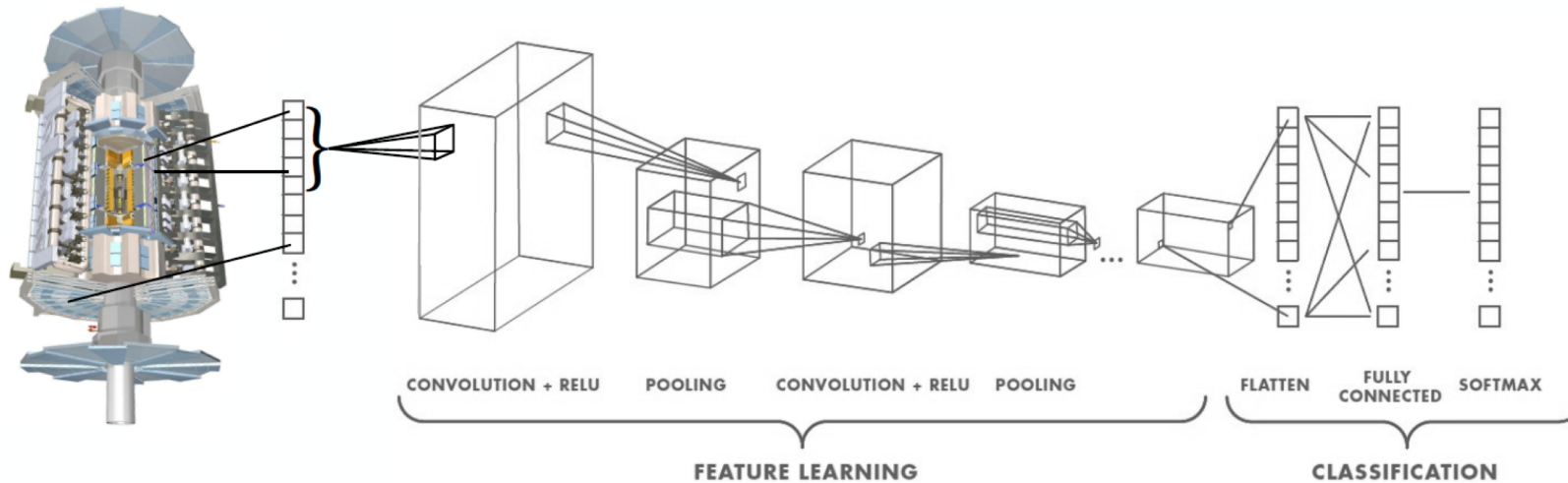
# #1: COMPUTING AT SCIENTIFIC FACILITIES

State-of-the-art: Took 3000 collaborators nearly 10 years to build



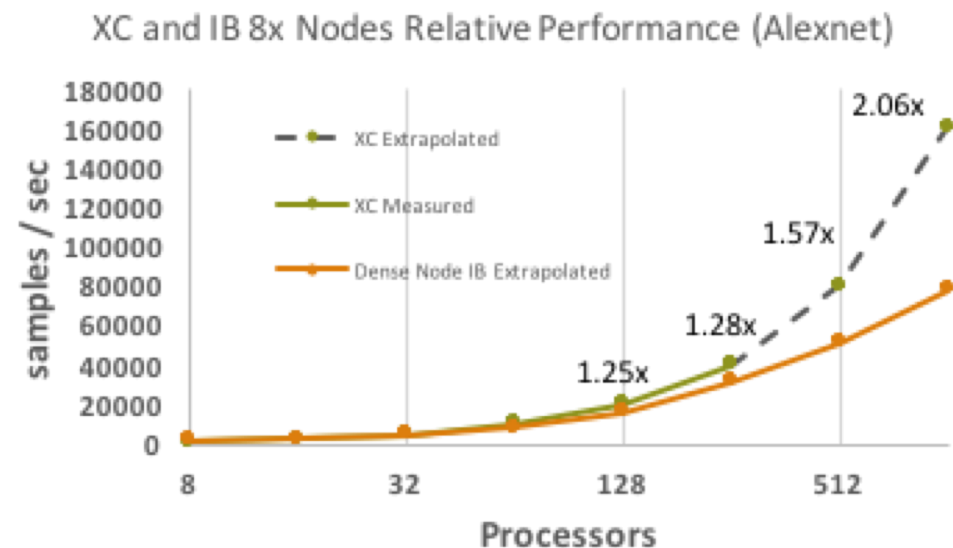
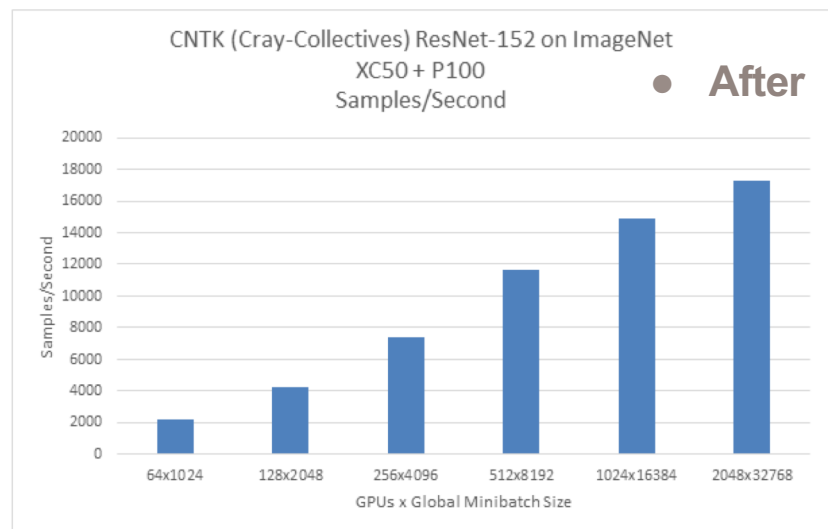
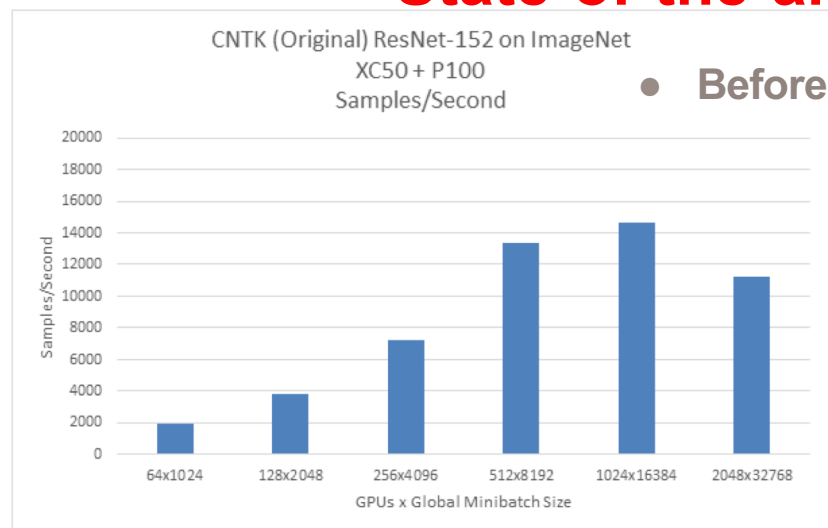
Using the labels from previous data analysis efforts....

Saves compute cycles....

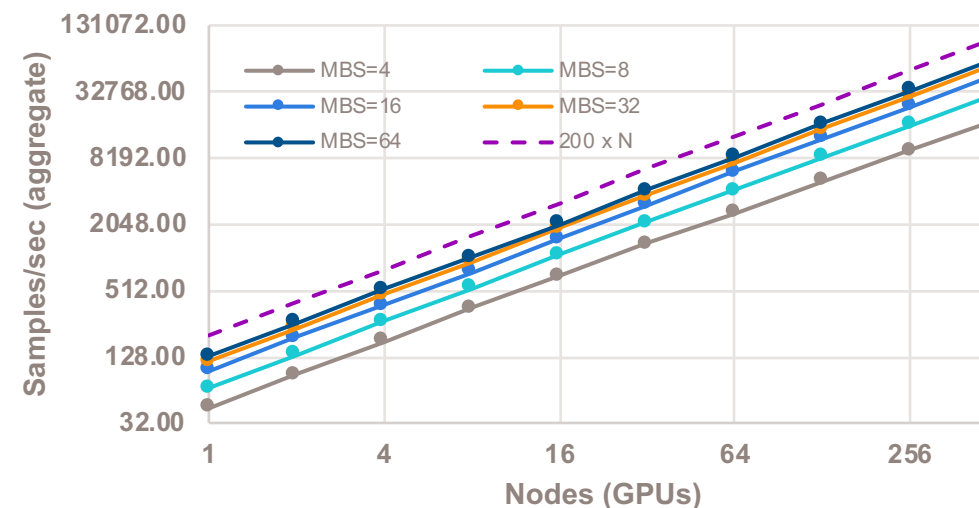


# #2: DISTRIBUTED TRAINING OF AI/ML CODES

State-of-the-art: Takes 6+days to train a model



Inception v3 Performance on XC50



Convergence “trains” in minutes



# #3: INTERACTIVE ANALYSIS OF BIG DATA



State-of-the-art: Exploratory data analysis is not interactive

- Cori @ NERSC
- 1630 compute nodes
- Memory: 128 GB/node,
- 32 2.3GHz Haswell cores/node

Gittens, Alex, et al. "Matrix factorizations at scale: A comparison of scientific data analytics in spark and C+ MPI using three case studies." , *IEEE International Conference on Big Data*.2016.

Science Area	Format/Files	Dimensions	Size
MSI	Parquet/2880	8,258,911 × 131,048	1.1TB
Daya Bay	HDF5/1	1,099,413,914 × 192	1.6TB
Ocean	HDF5/1	6,349,676 × 46,715	2.2TB
Atmosphere	HDF5/1	26,542,080 × 81,600	16TB

Convergence enables  
“iterative discovery”

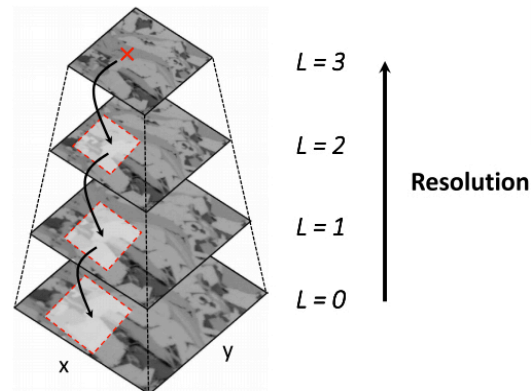
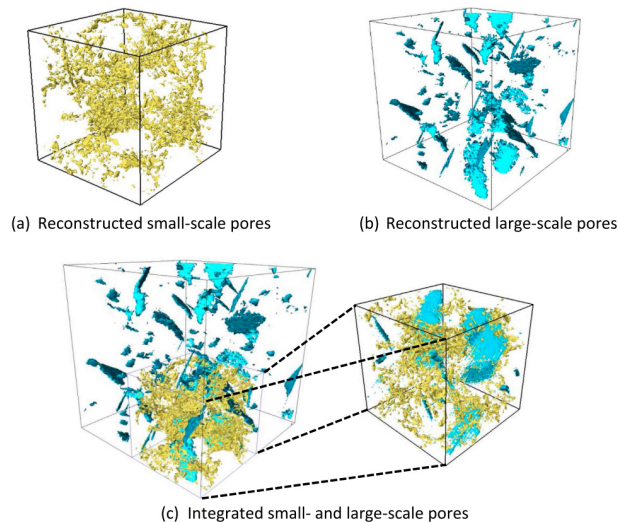
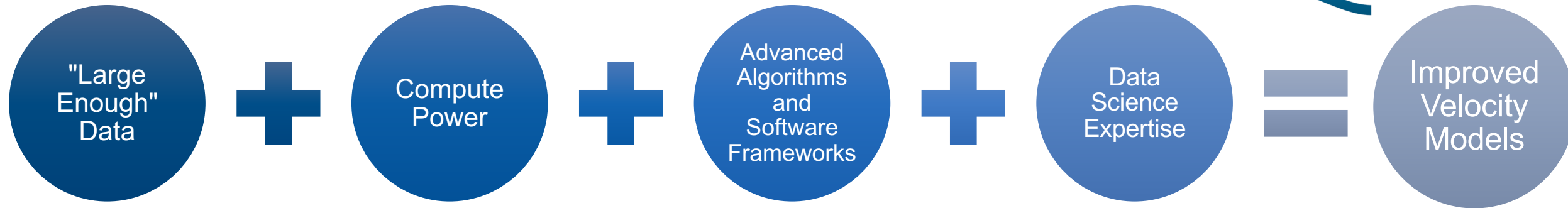
	Nodes / cores	MPI Time	Spark Time	Gap
NMF	50 / 1,600	1 min 6 s	4 min 38 s	4.2x
	100 / 3,200	45 s	3 min 27 s	4.6x
	300 / 9,600	30 s	70 s	2.3x
PCA (2.2TB)	100 / 3,200	1 min 34 s	15 min 34 s	9.9x
	300 / 9,600	1 min	13 min 47 s	13.8x
	500 / 16,000	56 s	19 min 20 s	20.7x
PCA (16TB)	MPI: 1,600 / 51,200 Spark: 1,522 / 48,704	2 min 40 s	69 min 35 s	26x

# #4: HYBRID THEORY+AI MODELS



CRAY

State-of-the-art: Works when manually-tuned

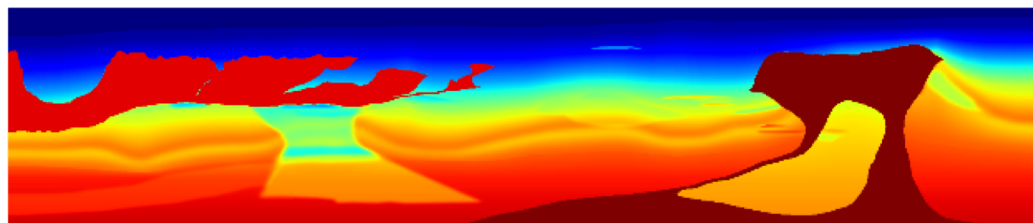


$$\min_m \underbrace{\|F(m) - d\|^2}_{\text{data misfit}} + \lambda \left( \underbrace{\|\Theta \nabla(m - m_r)\|_{L^1}}_{\text{weighted TV}} + \underbrace{\int_{\Omega} P \cdot \nabla(m - m_r)}_{\text{Steerable Variation}} \right)$$

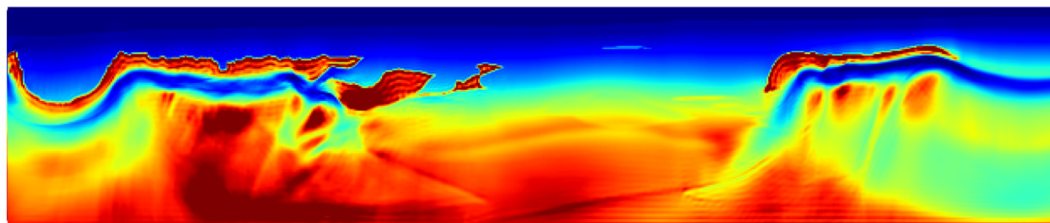
$$d_{js}(\mathbf{DI}||\mathbf{M}) = \frac{1}{2} \int_{-\infty}^{+\infty} \mathbf{DI}_i \log \left( \frac{\mathbf{DI}_i}{\mathbf{M}_i} \right) dx + \frac{1}{2} \int_{-\infty}^{+\infty} \mathbf{M}_i \log \left( \frac{\mathbf{M}_i}{\mathbf{DI}_i} \right) dx$$

# #4: HYBRID THEORY+AI MODELS

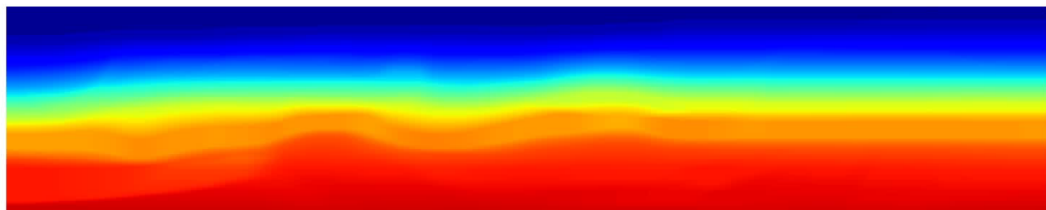
“Ground truth” Benchmark



Conventional FWI



FWI with Machine Learning



Synthetic benchmark with known subsurface geometry and seismic reflection data.

Conventional FWI attempts to derive a more accurate velocity model.

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

**Convergence enables “higher fidelity” to reality**



# TRENDS AND CHALLENGES



# TREND: MULTI-MODEL DATA

CRAY

INCREASING SPATIAL, TEMPORAL AND SPECTRAL RESOLUTION

Transactional



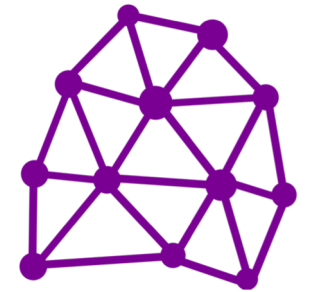
Conversational



Genomics



Graphs



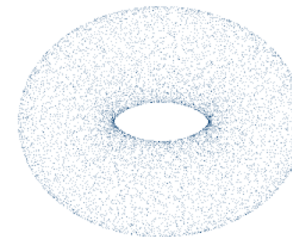
Documents



Images



3D Point clouds



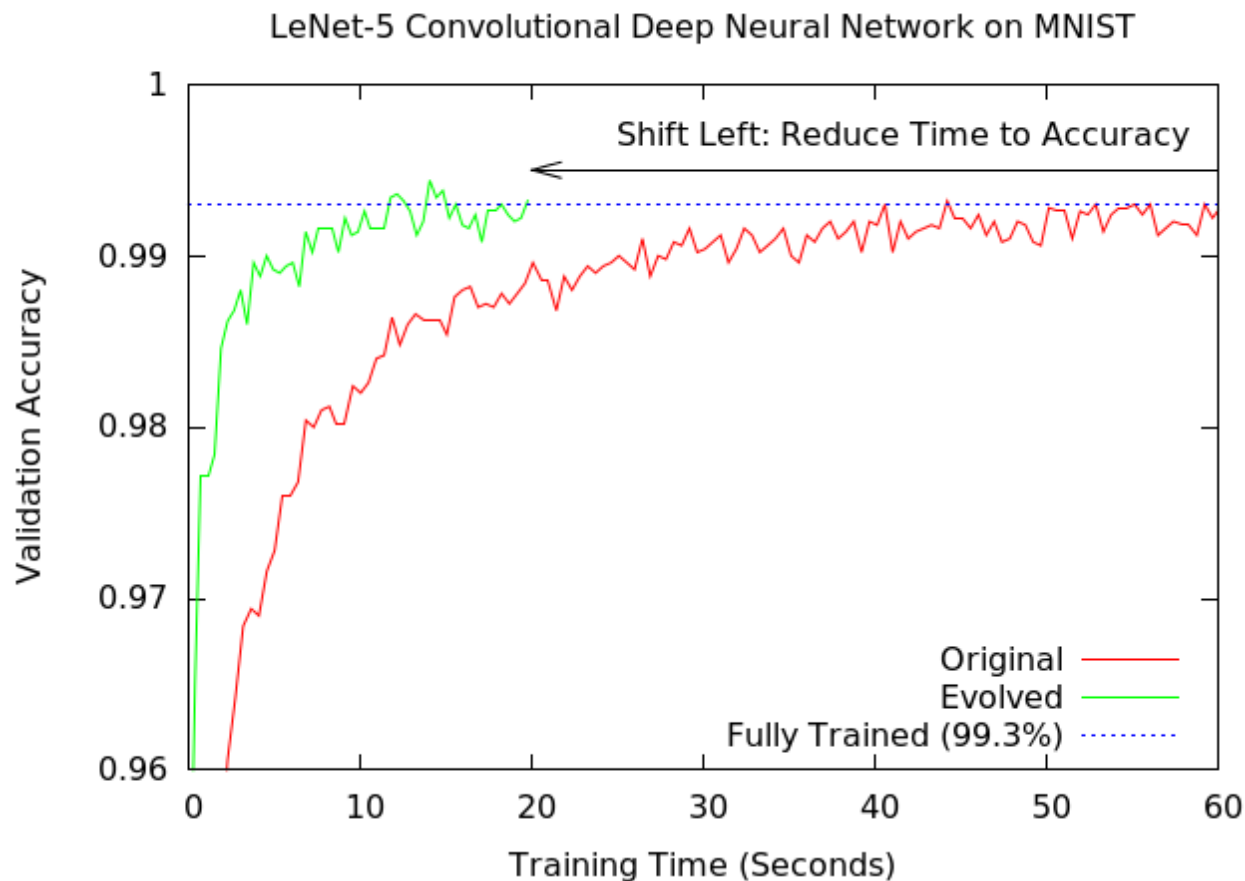
Sensors



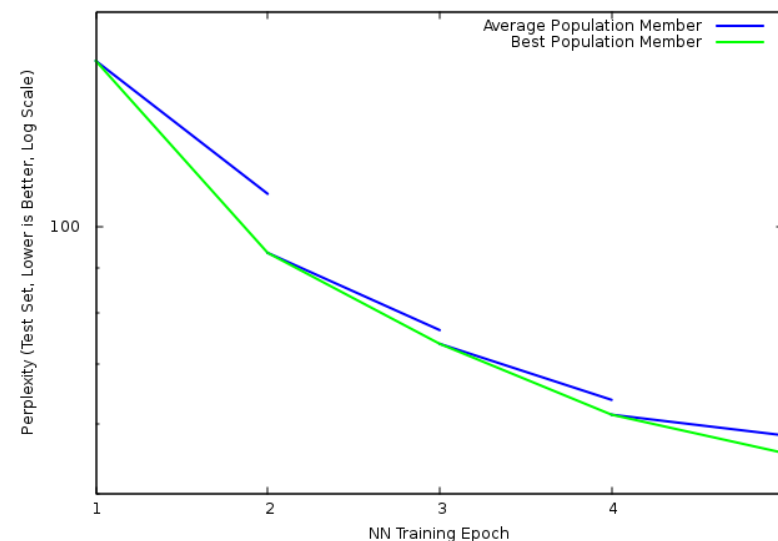
# CHALLENGE: I/O PATTERNS OF HPO LESS STUDIED



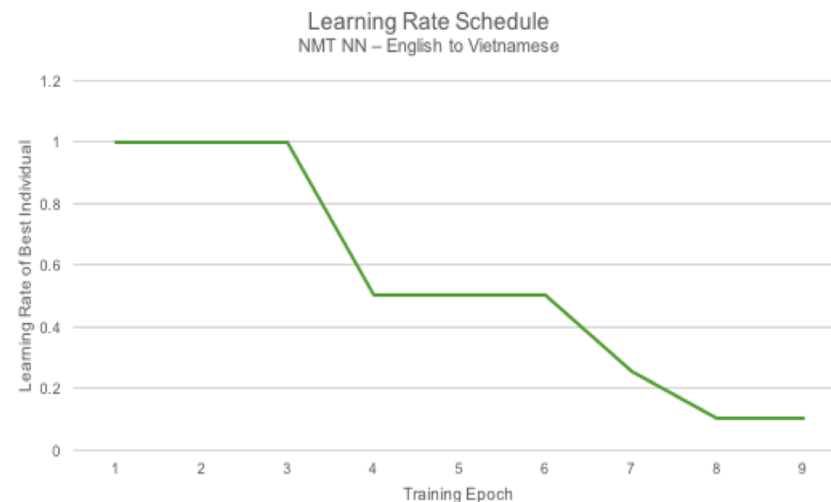
## HYPERPARAMETER OPTIMIZATION CRITICAL FOR MULTI-MODEL DATA



Source: Aaron Vose, Cray Performance Team



### Learning the optimal topology



### Learning a “learning-rate” schedule

# TREND : PROLIFERATION OF APPLIANCES

CRAY



Massively parallel processing databases



“Analytics is retrieval”

Distributed Analytics on Storage



“Take compute to cheap storage”

Distributed-memory Analytics



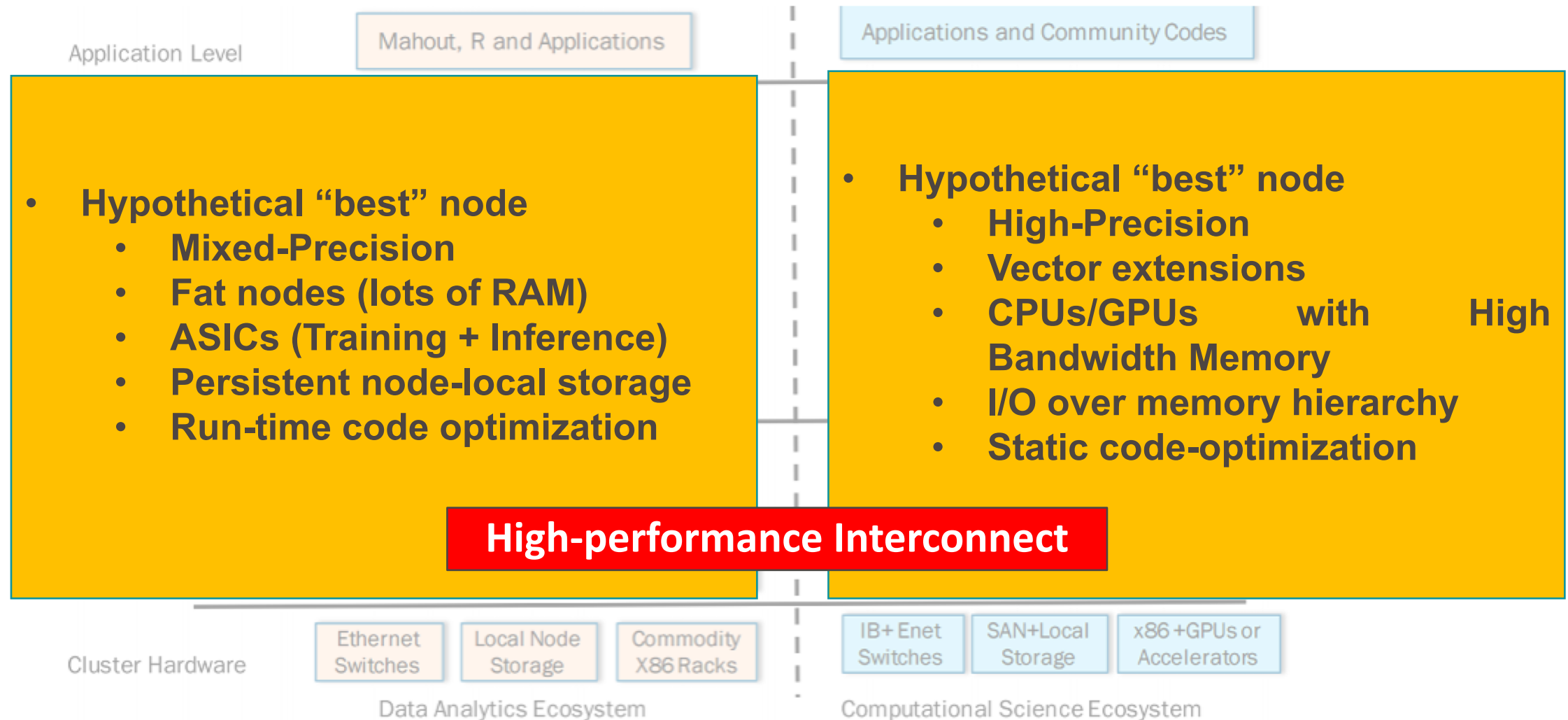
“Algorithm is made to work on distributed memory-chunks”

... SOLVES THE COMPUTE PROBLEM BUT CREATES NEW I/O PROBLEMS.



# CHALLENGE: NEED A SMARTER INTERCONNECT

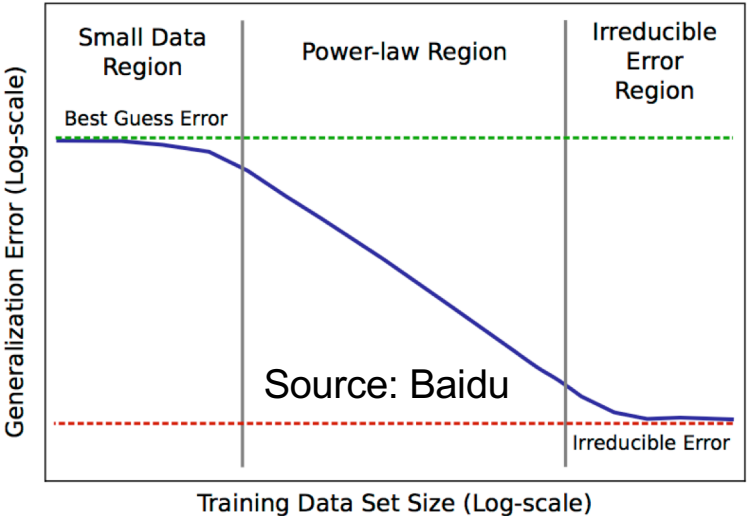
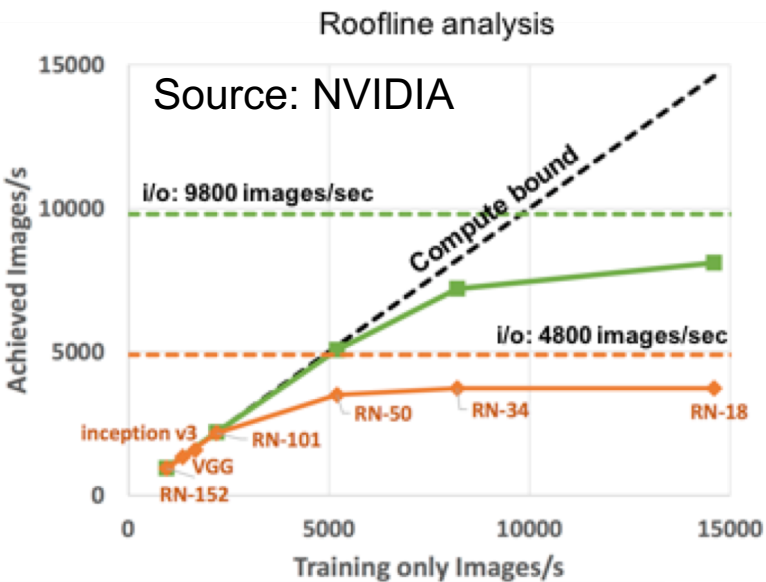
THE DATA MOVEMENT, MANAGEMENT AND I/O PATTERNS...



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

...WILL NEED A SMARTER INTERCONNECT.

# TREND: SOFTWARE TRICKS OUTPERFORMING HARDWARE



ResNet-50 Success	Time-to-accuracy	How many GPUs?	Scalability Efficiency
Facebook (Caffe2)	2 days 1 hour	352 GPUs 256	90% (large-batch)
IBM PowerAI (Caffe)	50 minutes	256 GPUs	95% (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%
Tencent	< 7 minutes	2048 GPUs	Large batch @ 64K
Fast.ai on AWS (Cost: \$40)	~18 minutes	128 GPUs	Not available (large batch)

# CHALLENGE: SOFTWARE FOR NEW HARDWARE



- Software : 7-10x improvement in time-to-accuracy in 1 year on CNNs

Method	Who?
LARS (MBS – 32K)	NVIDIA
Learning Rate schedule (~64K)	Facebook
Gradient Clipping	Microsoft
Mixed Precision Training	Baidu
Optimizer Tuning (~32K) <ul style="list-style-type: none"><li>- K-FAC</li><li>- Neumann</li></ul>	Google Research (now part of TensorFlow)
Batch Normalization	Google

- Hardware: 10-1000x in 2 years\*
  - Training
    - Intel, AMD, ARM, NVIDIA
    - Google TPU v2
    - Cerebras
    - Graphcore
    - Habana
    - **30+ startups....**
  - Inferencing
    - Intel Nervana
    - Wave Computing
    - Groq



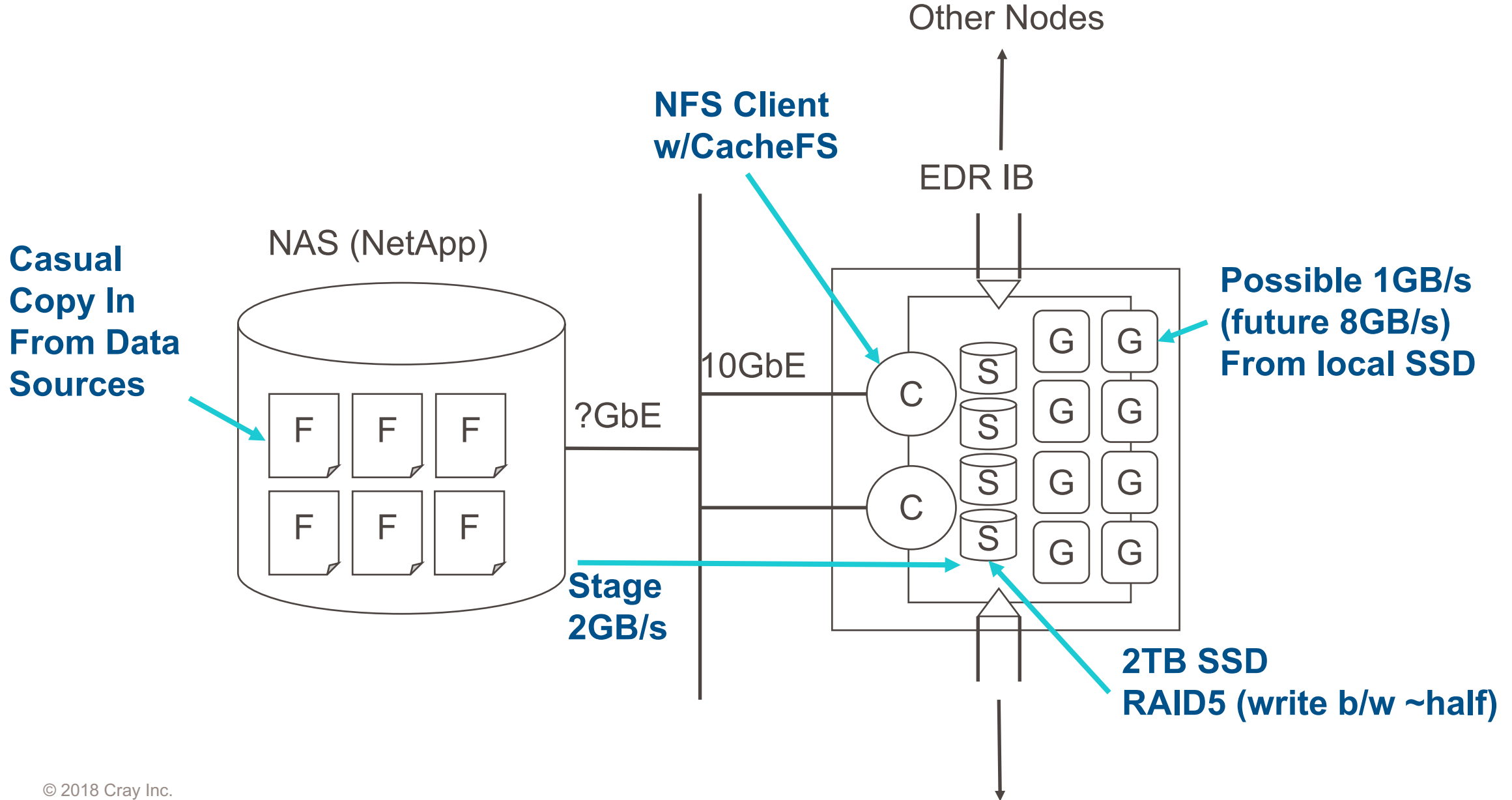
# TREND: TRIGGERED TRAINING



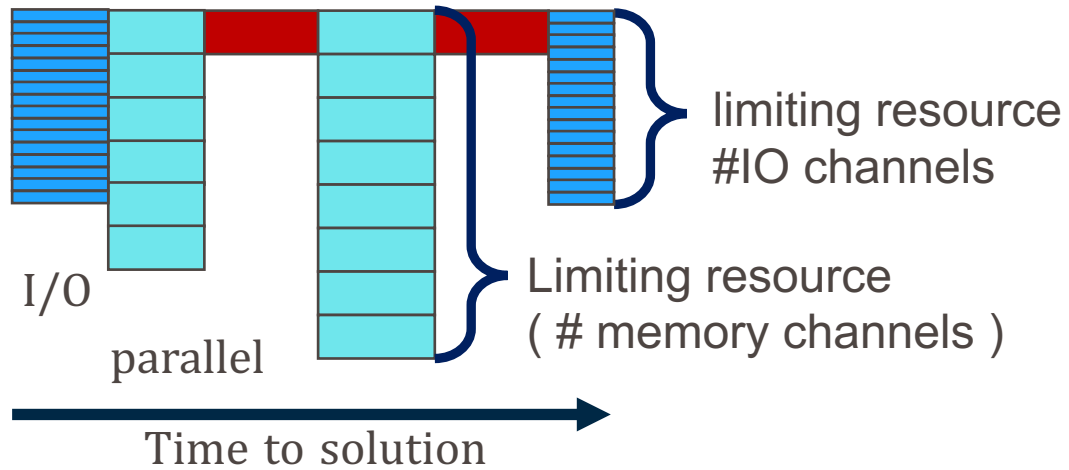
TRAINING PATTERNS DETERMINE SUPPORTING INFRASTRUCTURE FOR STORAGE AND I/O

Training	Use-case	Data size growth in unit time	Time to quality metric today	# of xPUs
Continuous	Internet-of-things	1:1	O(minutes)	O(10)
Cadence	Uber Eats prediction	n:1 ( $n \gg 1$ )	O(days)	O(10)
Delta	Speech (rare words)	n:1 ( $n \sim 1$ )	O(days)	O(1)
One-time	Lower-order physics approximations	10-100n:1	O(weeks)	O(100+)
Throughput	Speech and speaker detection	1:# of users	O(days)	O(100+)

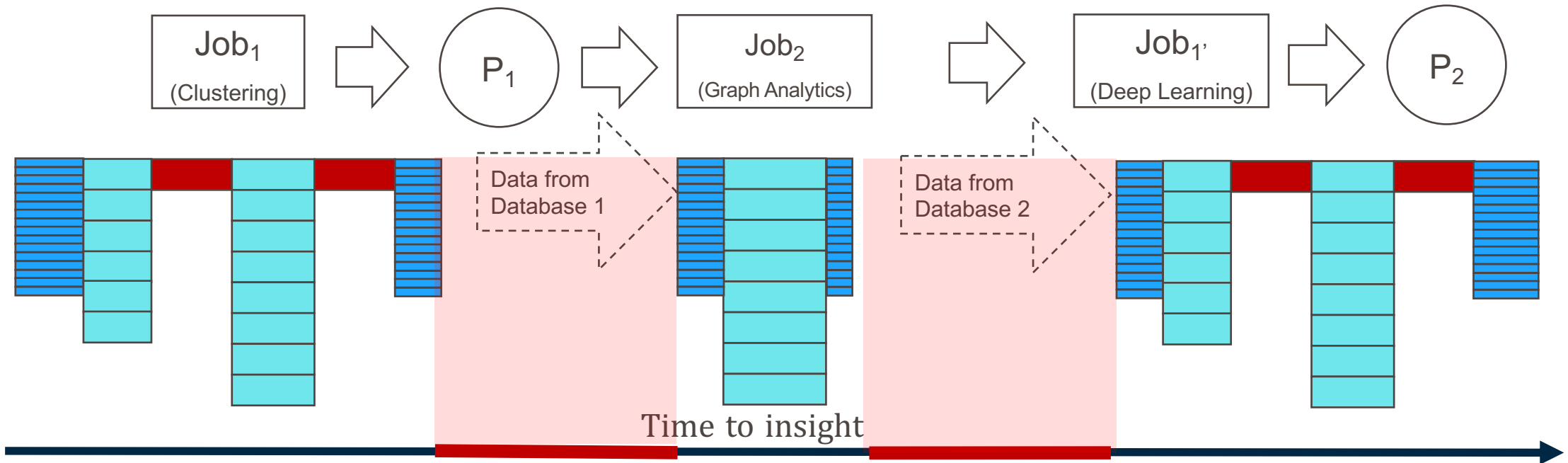
# CHALLENGE: POTENTIAL OFF-NODE I/O REQUIREMENT



# TREND: MULTI-TOOL WORKFLOWS ARE THE NORM



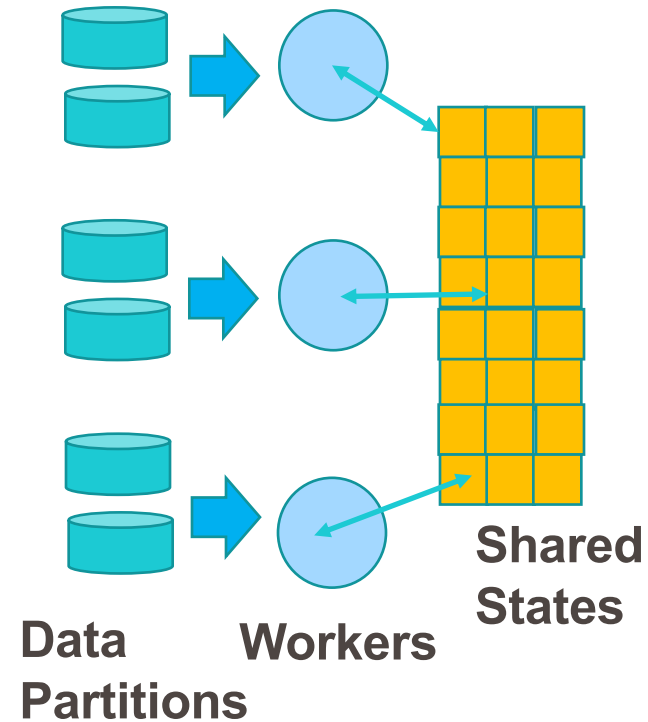
**Optimized for components but  
not the end-to-end workflow**



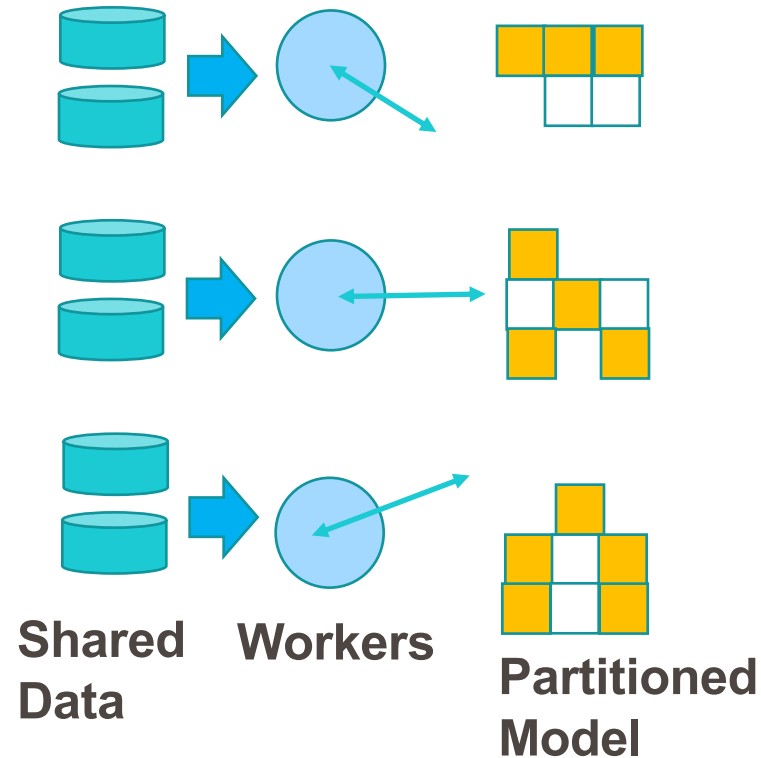


# CHALLENGE: FUTURE PARALLELISM

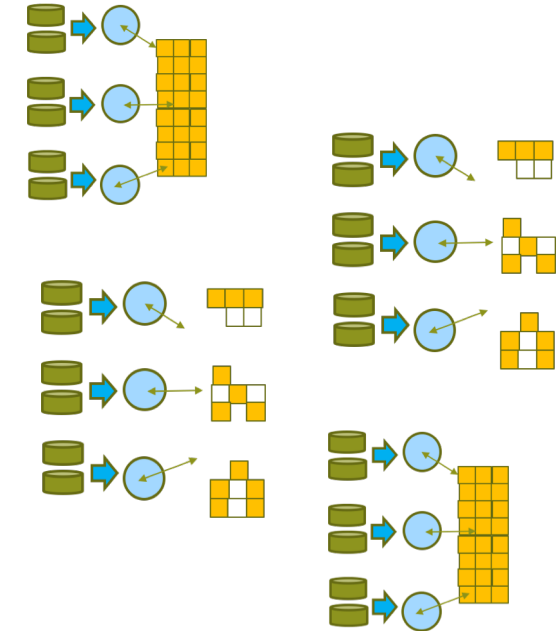
## Data-Parallelism



## Model-Parallelism

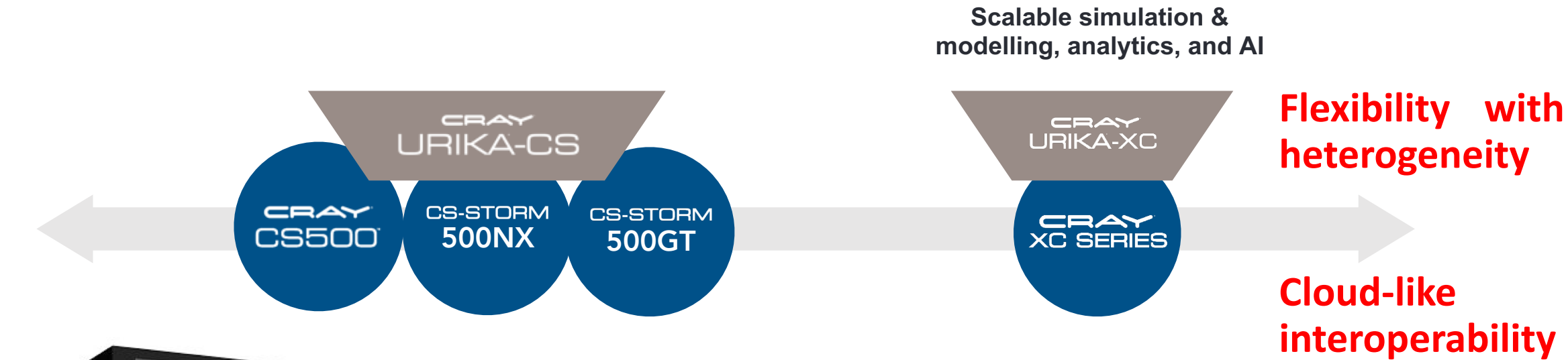


## Ensemble-Parallelism



- Model-Parallelism (Training, Inferencing)
- Higher Resolution Images
- Intra-node vs. Inter-node bandwidth

# TREND: CONVERGED HARDWARE+SOFTWARE

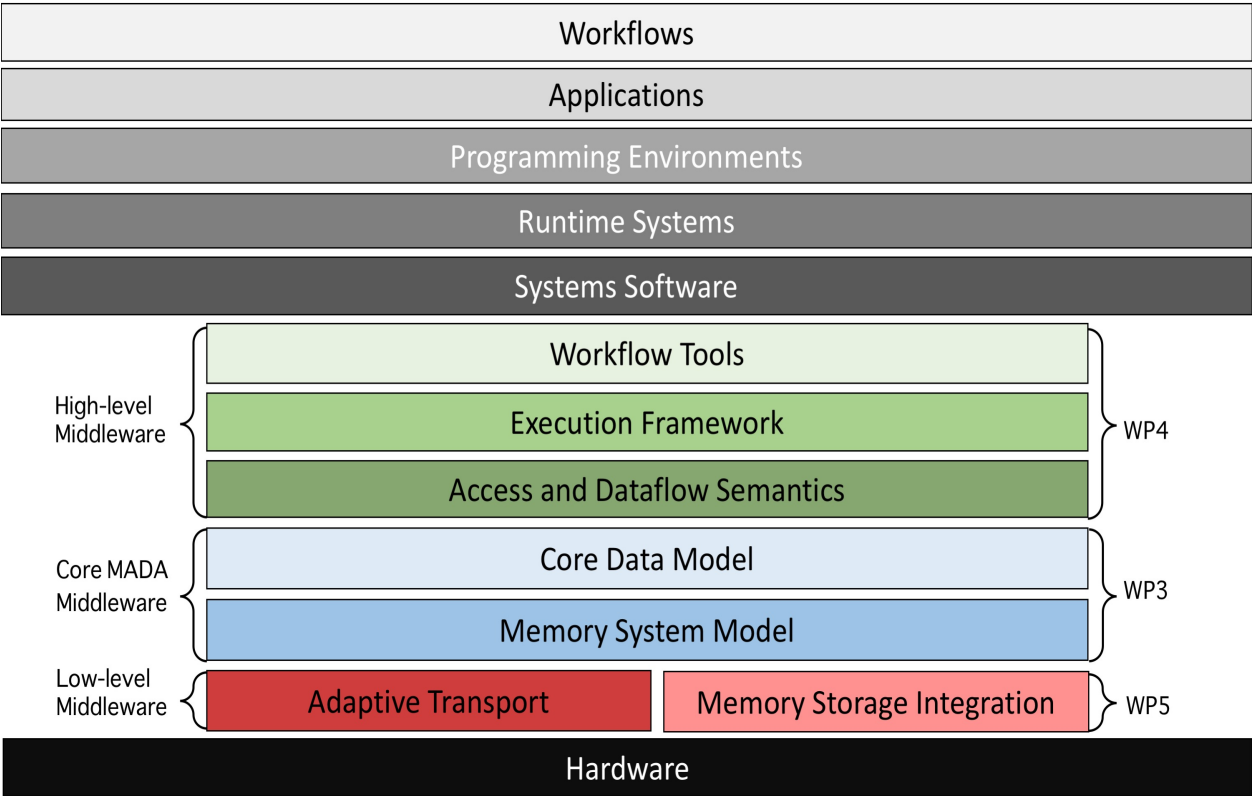


UIs: Jupyter Notebooks, TensorBoard						
MLlib, Spark SQL, Spark Streaming, GraphX	BigDL	Anaconda Python, Dask	pbdr	PyTorch	Keras, TensorFlow™	Cray Graph Engine (CGE)
Apache Spark™				Distributed Training Framework Horovod, CrayPE ML Plugin		
Intel® MKL, Intel MKL-DNN, Cray MPI						

# CHALLENGE: CONVERGENCE REQUIRES WORK



Convergence is not all hardware.....



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
	Operating Systems	Runtimes	Systems Software	Programming Environments	Applications	Workflows

Source: Adrian Tate, Cray EMEA

Lot more work before convergence can be productive....




# CHALLENGE: DELIVERING A SEAMLESS EXPERIENCE



## 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	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	Libraries e.g. MKL, CUDA, libSci	Data Pre-Processing Does system offer tools to optimize ETL wall-time?
	Provisioning Mesos, Moab, SLURM	Collectives e.g. NCCL, MPI	
Component Performance	Node Architecture # of xPUs+ cache + memory + network	Data Structure e.g. matrix, sequences, unstructured grids	Data Movement Does system provide ability to run multiple frameworks/applications on the same data?
	Disk Latency		
	Memory Capacity, Latency		
	xPU Speed		
			

# SUMMARY: SYSTEMS FOR THE FUTURE



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

# THANK YOU

QUESTIONS?

