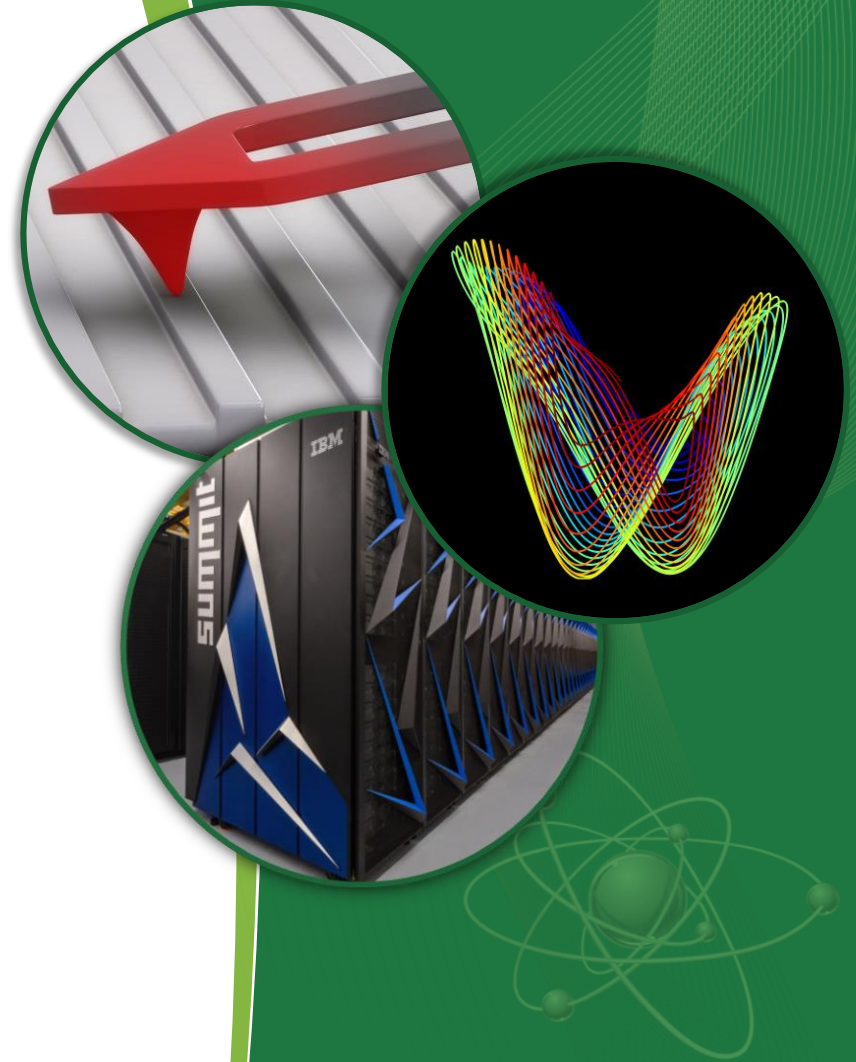


Early experiences with Machine Learning and Deep Learning on Summit/Summit- Dev

Junqi Yin

Advanced Data and Workflows Group



Outline

- ML/DL software stack on Summit
- CORAL2 benchmark
 - Data Science benchmark
 - Big Data Analytics Suite
 - Deep Learning Suite
- ML/DL performance model: Summit-Dev to Summit
- Scaling DL
 - Resnet50 on ImageNet
 - Lessons learned from exa-scale DL on Summit
- Discussion: ML vs DL use cases

ML/DL software stack on Summit (current plan and subject to change)

- Native installation
- IBM PowerAI container

</gpfs/wolf/stf011/world-shared>

- Custom container with Singularity (in planning)

Framework Version	Native	PowerAI Container	Custom Container	Python Wheels
Tensorflow	1.12	1.10, 1.8	1.9	tensorflow-1.12.0-cp36- cp36m-linux_ppc64le.whl
Pytorch	1.0rc1	0.4.1	0.4.1	torch-1.0.0a0+ff608a9- cp36-cp36m- linux_ppc64le.whl
R/PbdR	1.1		1.1	
SnapML		1.0.0		

Yin, Junqi / mldl-hpc · GitLab

https://code.ornl.gov/jqyin/mldl-hpc

GitLab Projects Groups Activity Milestones Snippets

mldl-hpc

Project Details Activity Cycle Analytics Repository Issues Merge Requests CI / CD Registry Wiki Snippets Settings

documentation Merge branch 'patch-1' into 'master'

tutorial add native support

utils add native support

wheels update tf wheel to v1.12.0

README.MD fix typo

README.MD

ML/DL software stack

Framework\Version	Native	PowerAI Container	Custom Container
Tensorflow	1.12.0	1.10.0	1.8
Pytorch	1.0	0.4.1	0.4.1
PbdR			
SnapML		1.0.0	

Wheels	CUDA:9.2.148 CUDNN:7.4.1 NCCL:2.3.7		
Tensorflow	tensorflow-1.12.0-cp36m-linux_ppc64le.whl		
Pytorch	torch-1.0.0a0+ff608a9-cp36m-linux_ppc64le.whl		

Documentation

[PowerAI on Summit](#)

Tutorial

[Keras Pytorch Tensorflow on Summit](#)

<< Collapse sidebar

<https://code.ornl.gov/summit/mldl-stack>

mldl-stack · GitLab

https://code.ornl.gov/summit/mldl-stack

GitLab Projects Groups Activity Milestones Snippets

mldl-stack

Overview Details Activity Contribution Analytics Issues Merge Requests Members

summit > mldl-stack > Details

mldl-stack

Install scripts, benchmarks, dependencies, utils, documentation for various deep learning software libraries.
Currently available libraries: Tensorflow and Pytorch.

Global

Filter by name... Last created

tensorflow
Install scripts, containers, wheels, benchmarks, utils, documentation for latest supported and legacy versions of tensorflow minutes ago

pytorch
Install scripts, containers, dependencies, wheels, benchmarks, utils, documentation for latest supported and legacy versions of pytorch a minute ago

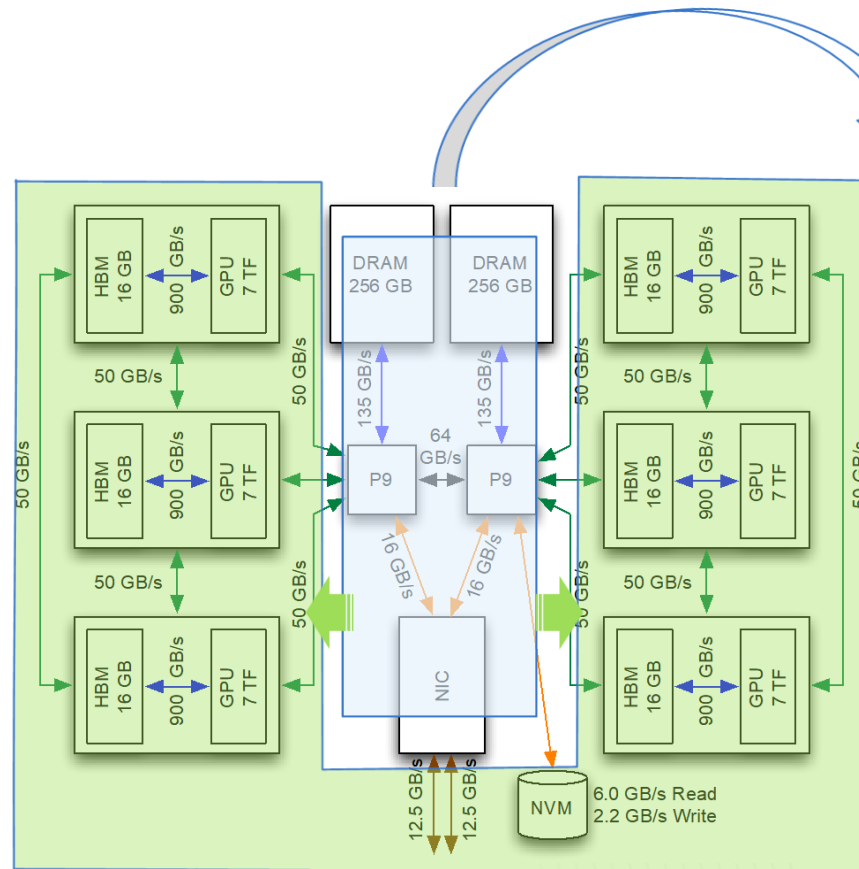
<< Collapse sidebar

<https://code.ornl.gov/jqyin/mldl-hpc>

CORAL-2 Data Sciences Benchmarks

Benchmarks	Description
Big Data Analytic Suite	PCA, K-Means, and SVM (based on pbdR)
Deep Learning Suite	CANDLE, CNN, RNN, and ResNet-50 (distributed)

Deep Learning Codes (CNN; ResNet50; ..) excel here with NVM and GPUs enabling tensor operations.

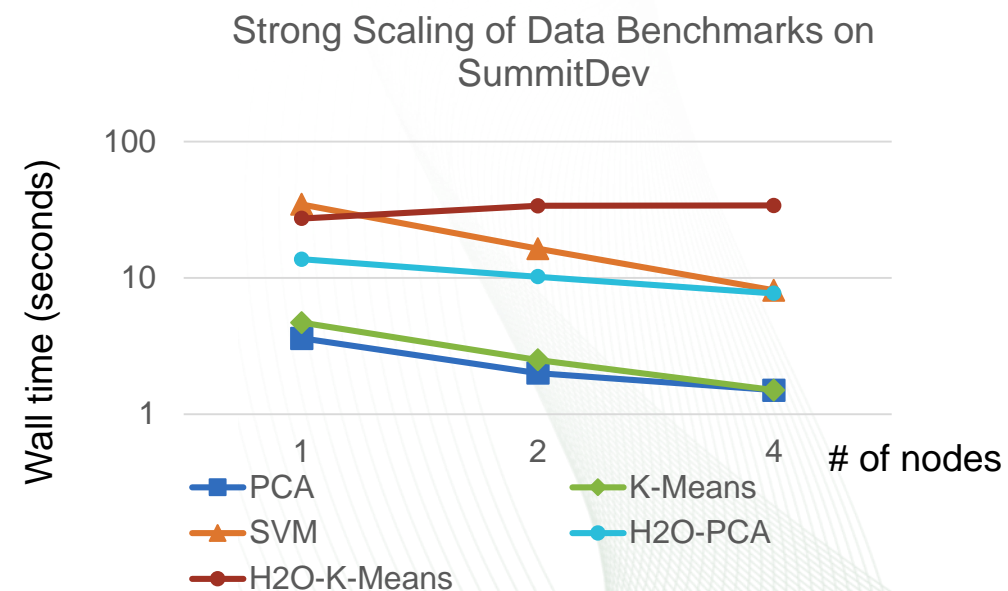
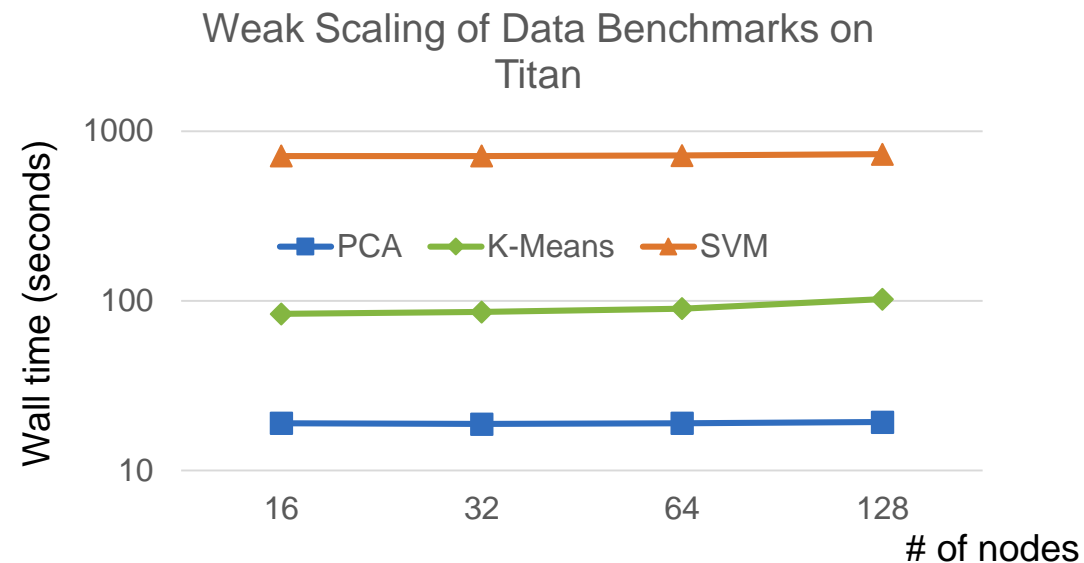
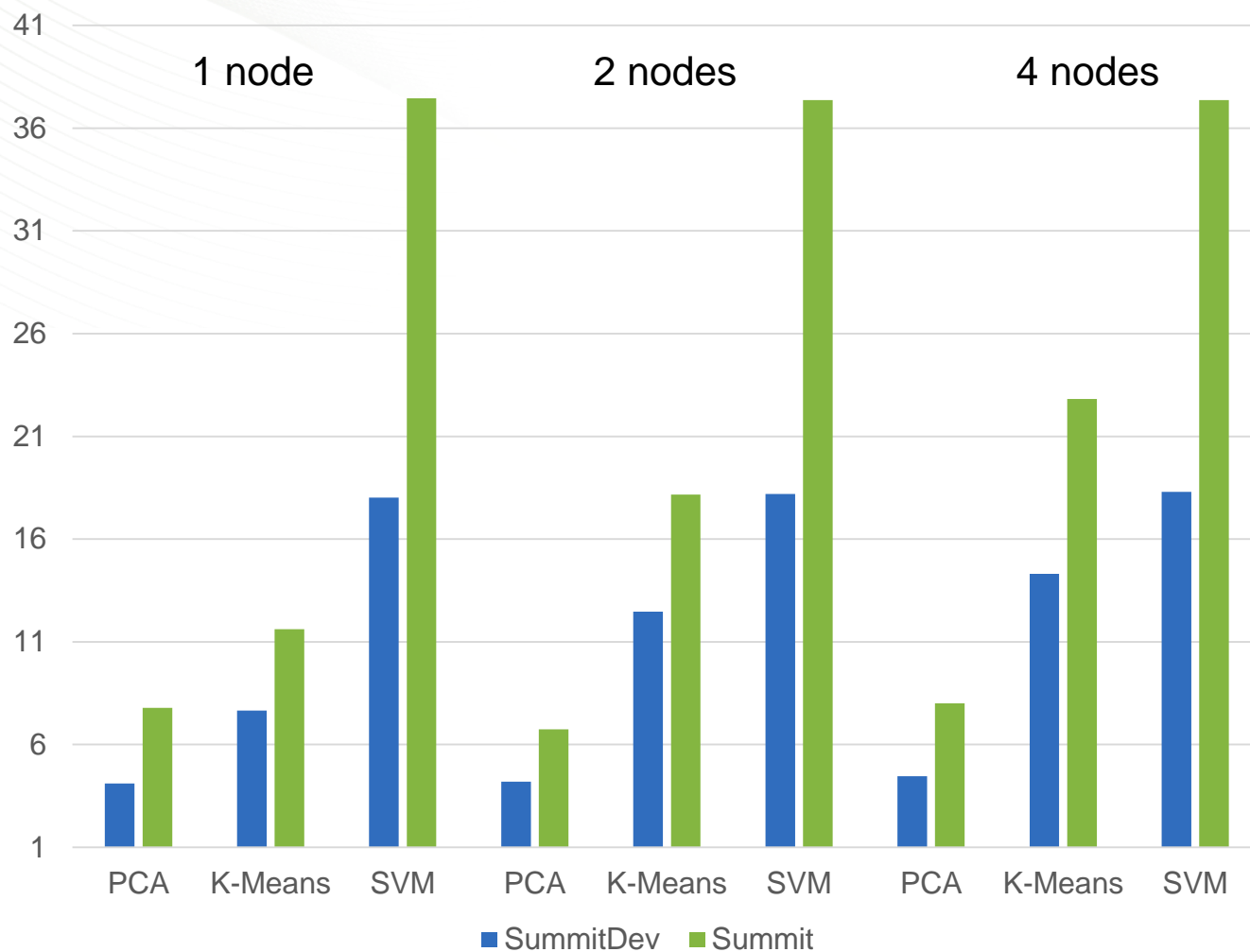


Traditional Node:
PCA, K-Means, etc.
excel due to the
node's memory,
CPU, and on-chip
bandwidth



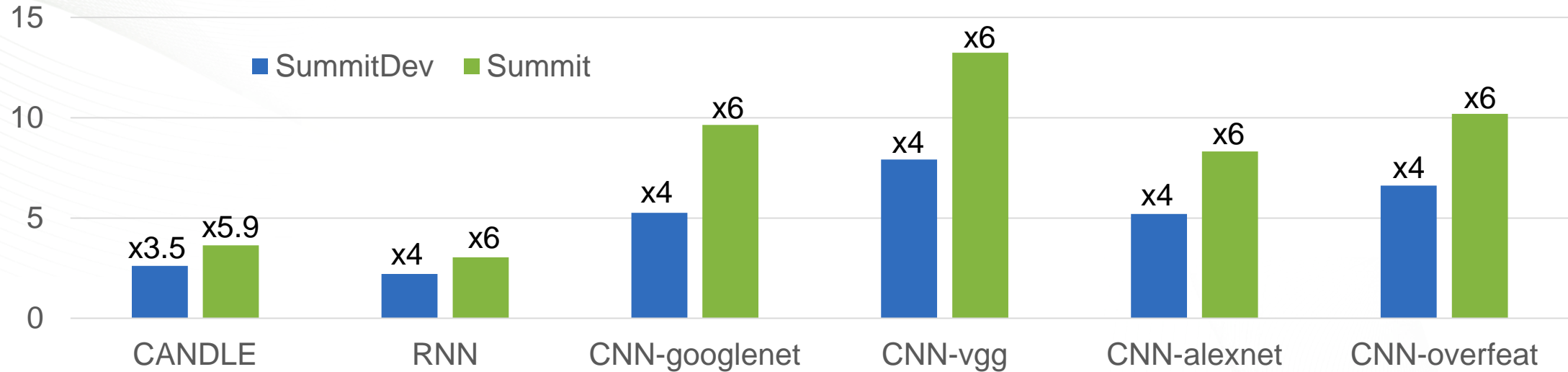
Big Data Analytic Suite

Speedup Over Titan Baseline for CORAL-2 Big Data Benchmarks (based on pbdR)

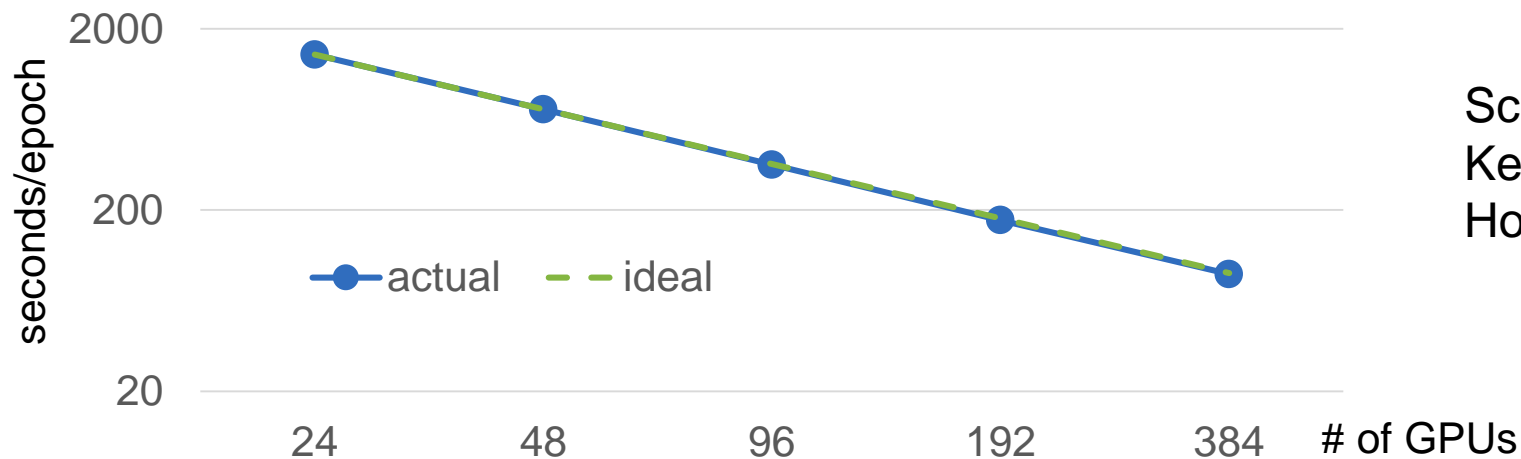


Deep Learning Suite

Speedup Over Titan Baseline for CORAL-2 Deep Learning Benchmarks



Strong Scaling of ResNet-50 on Summit



Scaling of Resnet-50 based on Keras (Tensorflow backend) and Horovod on ImageNet data

Performance model for BDAS

Architecture			Power9														
Workload(4aff10a)			Kmeans														
Input size																	
SMT (thread rank)																	
walltime(s)																	
			PCA			PCA			PCA			PCA			PCA		
			8GB			64GB			8GB			64GB			64GB		
			1	2(2 42)	4(4 42)	1	2(1 84)	4(2 84)	1	2(2 42)	4(4 42)	1	2(2 42)	4(4 42)	1	2(2 42)	4(4 42)
			1.9	2.1	2.4	13	12.2	13.5	3.1	3	3	26.2	26.5	26.6	1	16.7	16.7
Power8	PCA	1	6.8	3.58x	3.24x	2.83x											
		8GB	3.6	1.89x	1.71x	1.50x											
		4(2 40)	3.7	1.95x	1.76x	1.54x											
		8(4 40)	4	2.11x	1.90x	1.67x											
		1	53.1				4.08x	4.35x	3.93x								
		64GB	28.2				2.17x	2.31x	2.09x								
		4(1 80)	24.8				1.91x	2.03x	1.84x								
		8(2 80)	25.9				1.99x	2.12x	1.92x								
		1	4.7							1.52x	1.57x	1.57x					
		8GB	4.8							1.55x	1.60x	1.60x					
		4(4 20)	4.8							1.55x	1.60x	1.60x					
		8(8 20)	4.9							1.58x	1.63x	1.63x					
	Kmeans	1	72.8										2.78x	2.75x	2.74x		
		64GB	49.2										1.88x	1.86x	1.85x		
		4(1 80)	35.3										1.35x	1.33x	1.33x		
		8(2 80)	35.8										1.37x	1.35x	1.35x		
		1	34.6													2.07x	
	SVM	8GB	33.1													1.98x	
		4(2 40)	34.7													2.08x	
		8(2 80)	37.5													2.25x	
		1	286.6														
		64GB	285.7														
		4(1 80)	283.1														
		8(1 160)	291.5														

$$\log(Perf_{Power}) = Architecture + Size + Workload + Threads$$

Performance model for DL workloads

Architecture					Volta											
	Workload				CNN				RNN		Comm					
		Implementation			WINOGRAD_NONFUSED	IMPLICIT_PRECOMP_GEMM			LSTM		GRU		NCCL		MPI	
			Precision		fp32	fp16	fp32	fp16	fp32	fp16	fp32	fp16	fp32	fp16	fp32	fp16
				walltime(s)		2.99		1.98		5.1		227.98		0.33		0.46
Pascal	CNN	WINOGRAD_NONFUSED	fp32													
			fp16	5.1		1.71x										
		IMPLICIT_PRECOMP_GEMM	fp32													
			fp16	3.1				1.57x								
	RNN	LSTM	fp32													
			fp16	7.4					1.45x							
		GRU	fp32													
			fp16	359.5						1.58x						
	Comm	NCCL	fp32													
			fp16	0.3									1.27x			
		MPI	fp32													
			fp16	0.9												1.91x
		Problem size :														
		WINOGRAD_NONFUSED: input: 112x112xx64x16 filter: 3x3x128														
		IMPLICIT_PRECOMP_GEMM: input: 112x112xx64x8 filter: 3x3x128														
		lstm:1024-64-25 gru: 1024-64-1500 (RNN)														
		100000 4nodes (Comm)														

Takeaway - ML

- Per node, expect ~2x over SummitDev, up to ~35x over Titan.
- OpenBLAS provides close performance as IBM ESSL, although ESSL seems to handle SMT better.
- Use SMT=1/2 on Summit SMT=2/4 on SummitDev for pbdR and oversubscribe threads.
- Use RAPIDS, H2O4GPU, SnapML (close source), etc to take advantage of GPUs

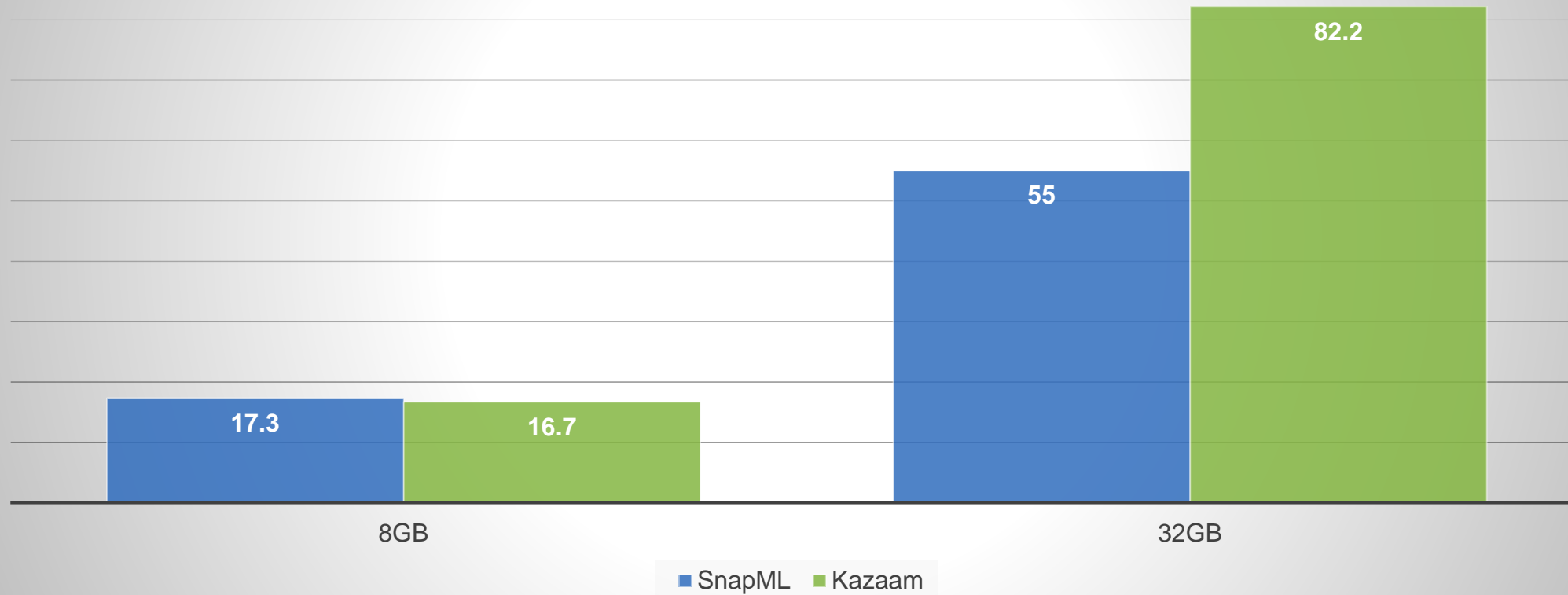
For more details, please refer to [arXiv:1811.02287](https://arxiv.org/abs/1811.02287)

Takeaway - DL

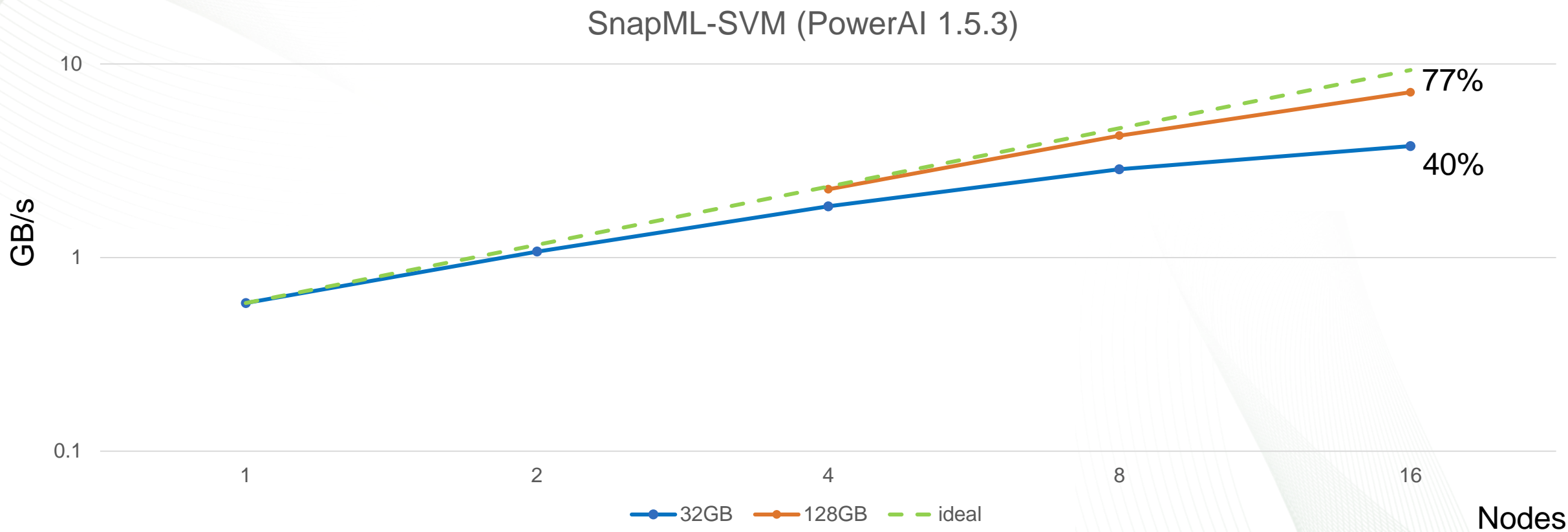
- Per node, expect $\sim 2.5x$ over Summit-Dev, up to $\sim 80x$ over Titan.
- Average $\sim 60x$ for CNN workloads, $\sim 20x$ for RNN workloads, over Titan
- $\sim 1.5x$ in communication over Summit-Dev
- Near ideal scaling for Keras (Tensorflow backend) + Horovod up to 64 nodes for Resnet50 on ImageNet

IBM's SnapML

SVM Benchmark (seconds)

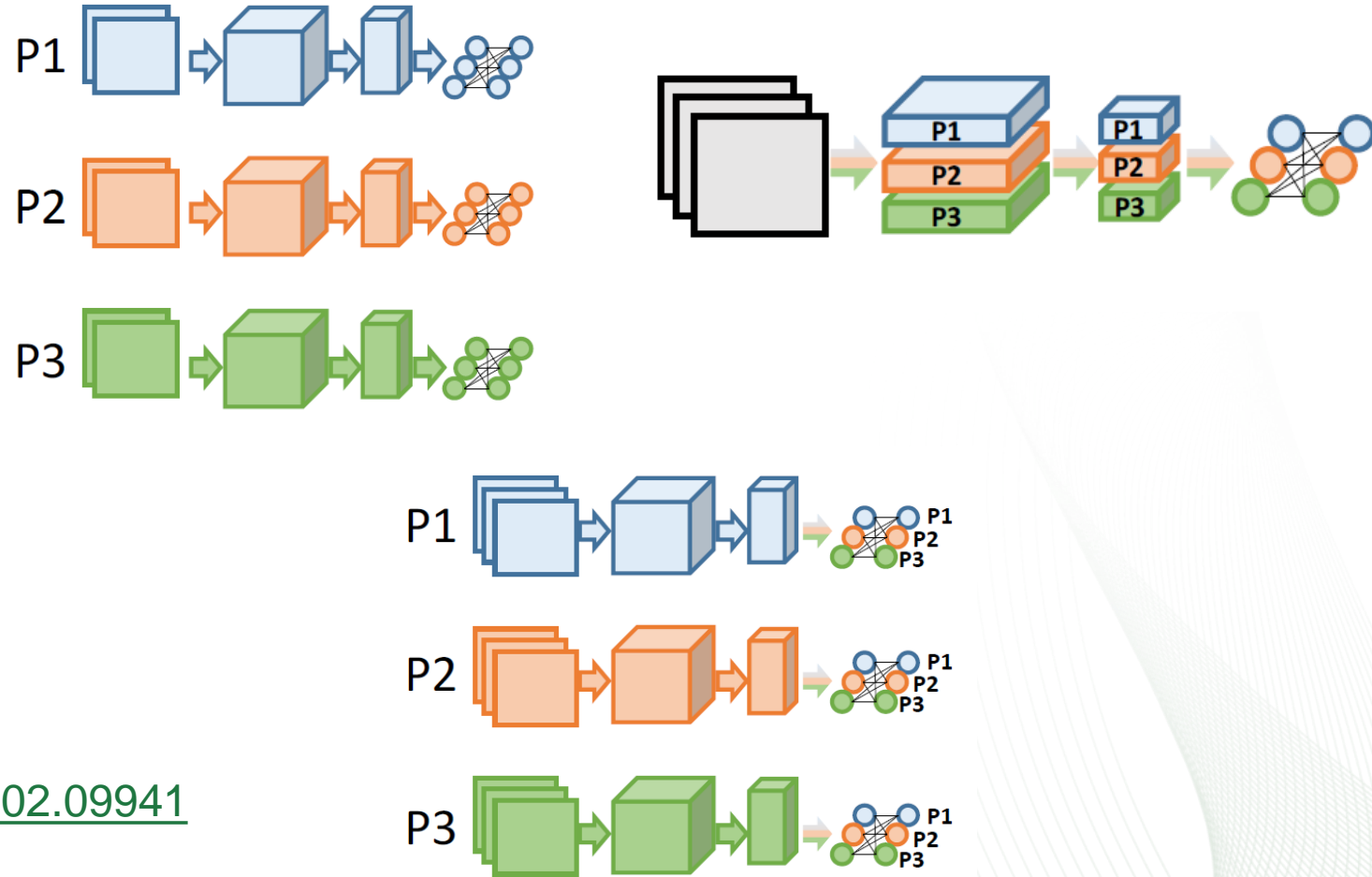


IBM's SnapML



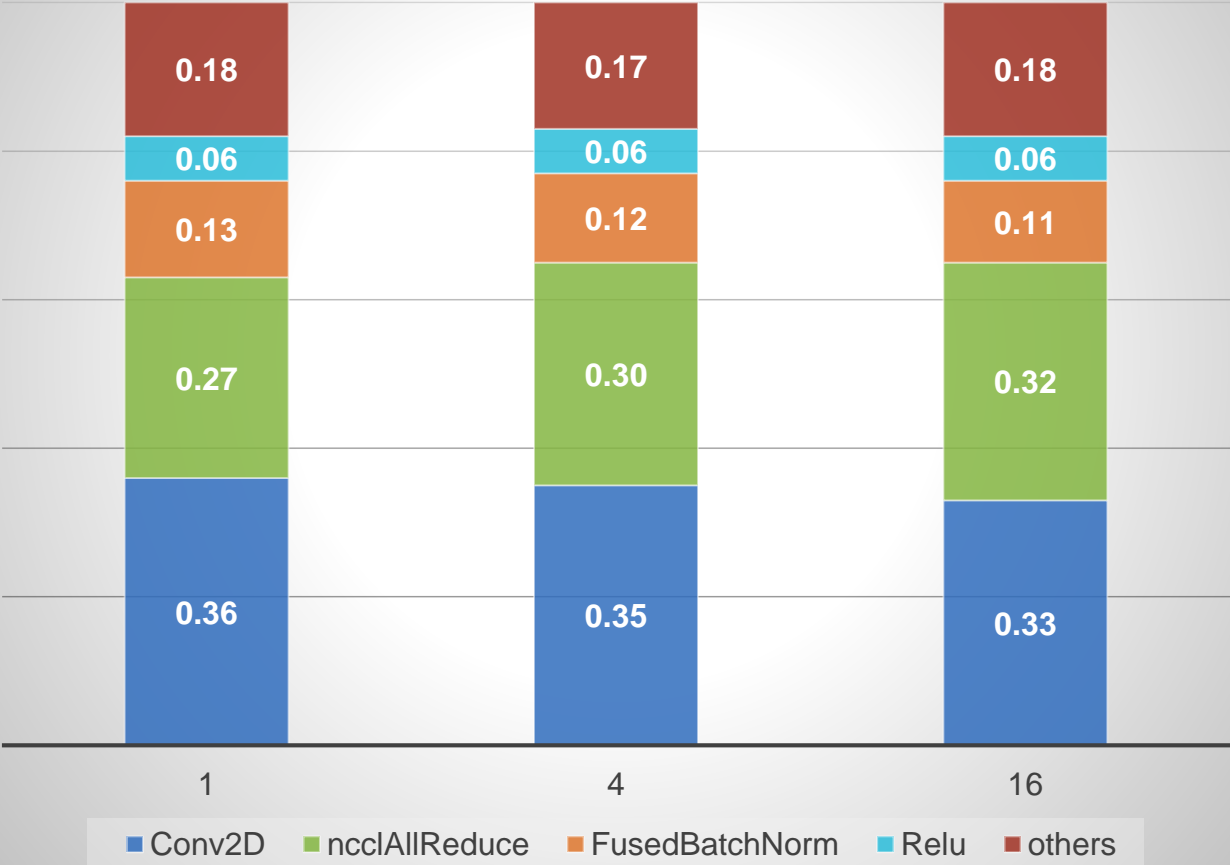
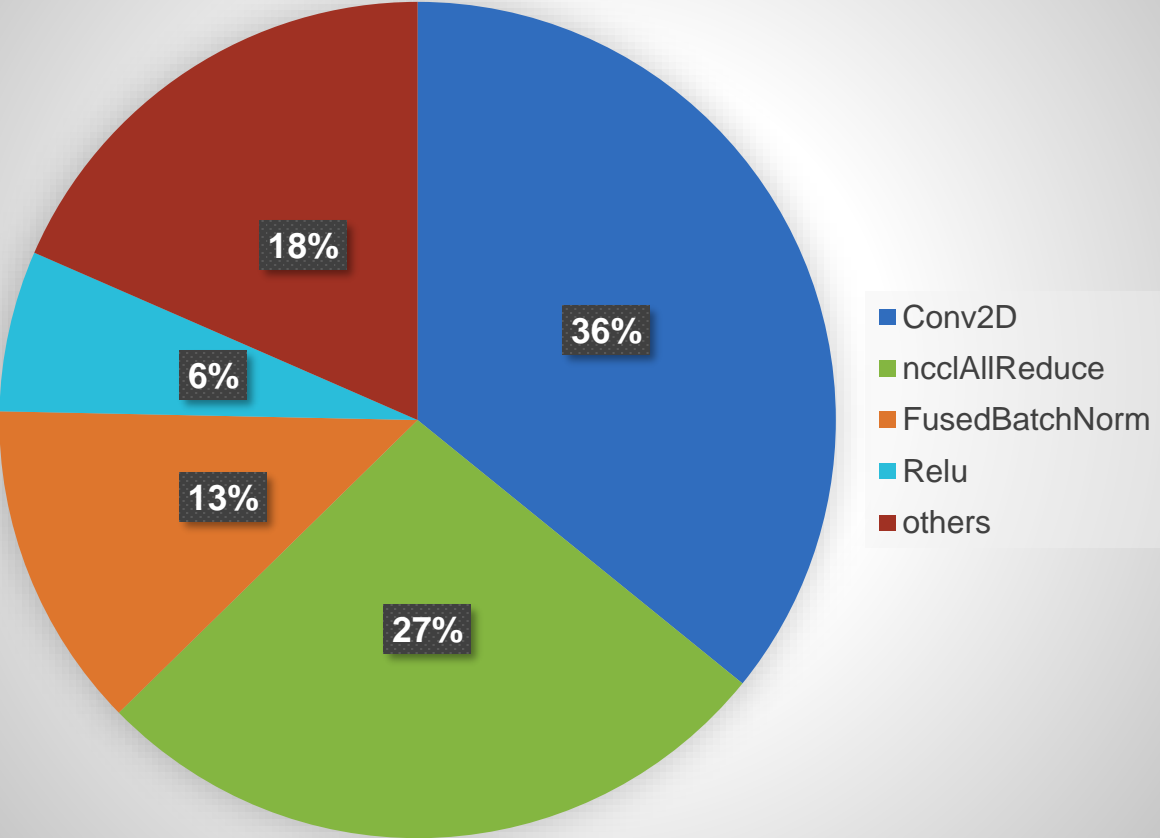
Distributed deep learning

- Data parallel
 - **Synchronized**
 - Stale
 - Asynchronized
- Model parallel
- Hybrid



Review: [arXiv:1802.09941](https://arxiv.org/abs/1802.09941)

TensorFlow Resnet50 profiling on Summit

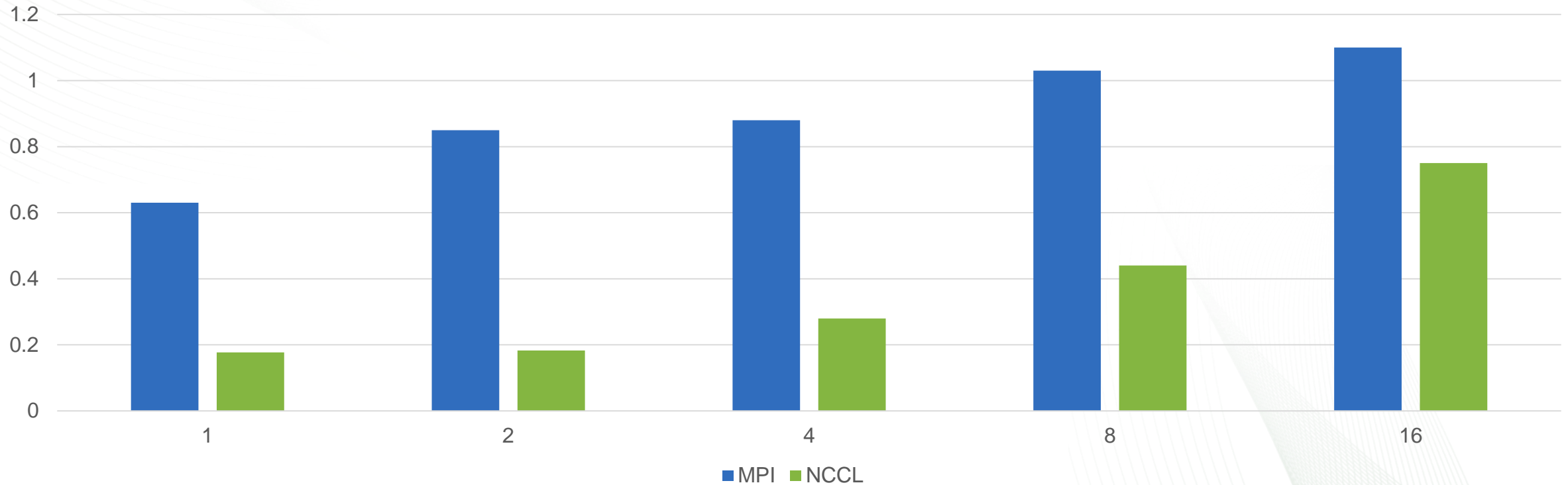


“mini-MPI” for distributed deep learning

- NCCL (Nvidia): collective multi-GPU communication
- Horovod (Uber): Tensorflow and Pytorch support
 - NCCLReduceScatter - MPIAllreduce – NCCLAllgather for data divisible by `local_rank()`
 - NCCLReduce - MPIAllreduce – NCCLBcast for the remainder
 - Tensor Fusion: fuse small allreduce tensor operations into larger ones for performance gain
 - Compression (cast vars to fp16) before allreduce
- GLOO (Facebook): Pytorch support
- DDL (IBM): Tensorflow, Pytorch, Caffe support. Close source.

NCCL vs MPI allreduce

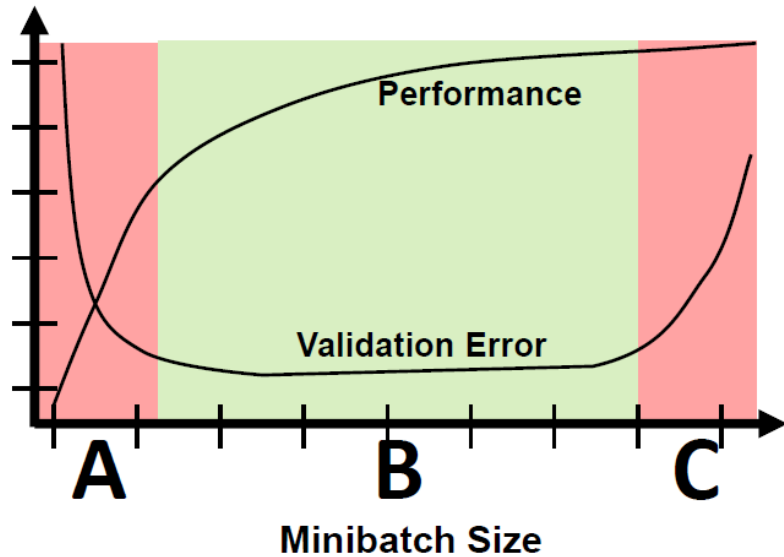
Time to transfer 10^5 floats on SummitDev



Differences in scaling up: DL VS simulation

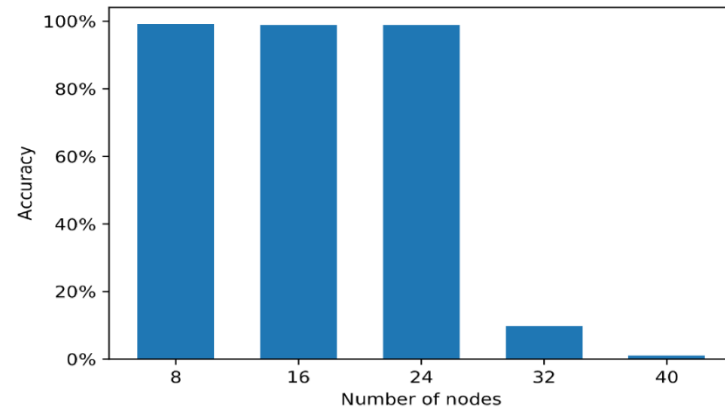
- DL is a global optimization, changing scale -> changing solution space.
 - DL usually requires changing network architecture, update scheme, etc
- Scale in OPS \neq Scale in time-to-solution (accuracy)
 - Tradeoff between more epochs and faster convergence
- High per-node OPS makes DL comm- and/or IO- bound at relatively small node count.
 - DL requires special designed comm (mainly all-reduce) and IO pipeline

Synchronized data parallel: scaling vs convergence

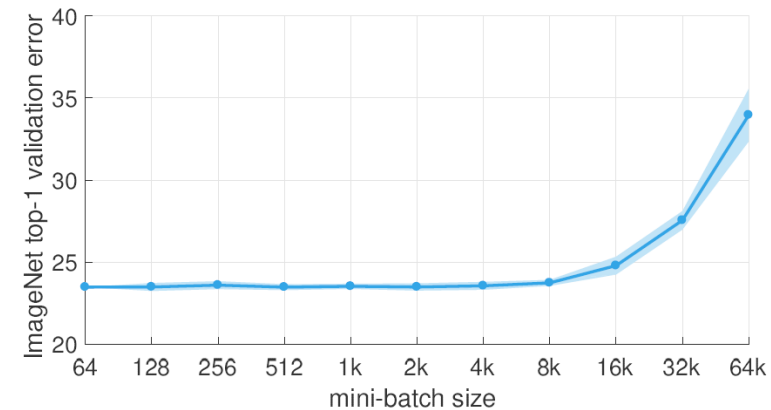


(a) Minibatch Effect on Accuracy and Performance (Illustration)

- Possible causes: “**generalization gap**” (Keskar et al. 2017)
 - loss of the explorative properties
 - tend to converge to sharp minimizers
 - model overfits the training data



Convergence of MNIST with increasing mini-batch size



(b) Empirical Accuracy (ResNet-50, figure adapted from [Goyal et al. 2017], lower is better)

Large mini-batch size training

- mini-batch size 8K (arXiv:1706.02677)
 - Warmup with default learning rate for optimizer
 - Start with learning rate multiplying # of workers
 - Decay learning rate periodically
- mini-batch size 32K
 - Layer-wise adaptive rate scaling (LARS) (arXiv:1711.04325)

State-of-the-art Imagenet training

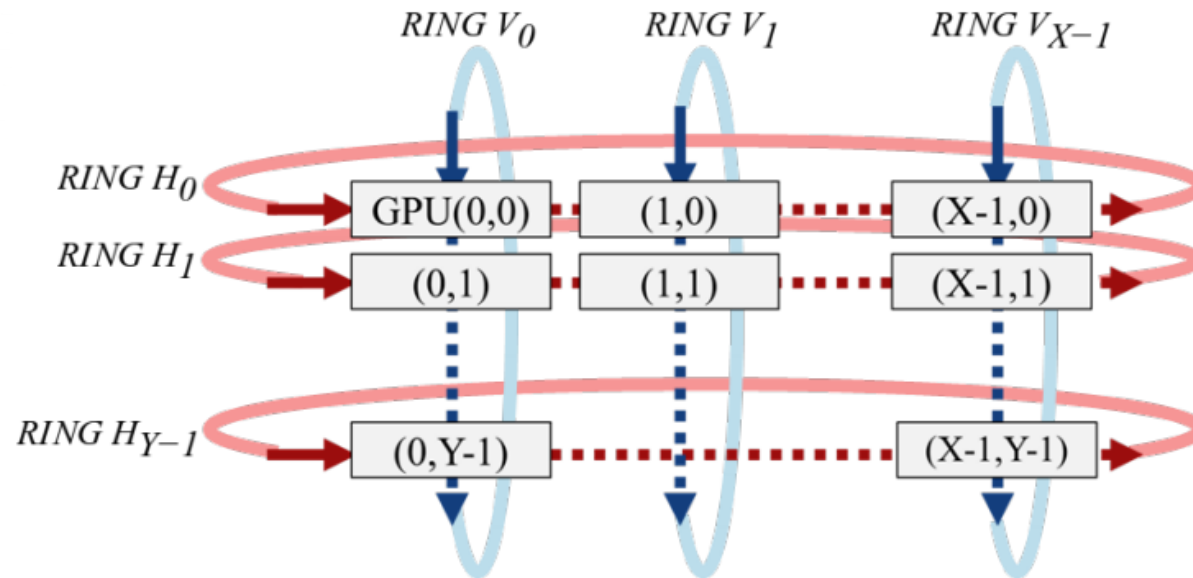
Chronology of Distributed Deep Learning Records

Table 1 : Training time and top-1 1-crop validation accuracy with ImageNet/ResNet-50

		Batch Size	Processor	DL Library	Time	Accuracy
2016	He et al.	256	Tesla P100 x8	Caffe	29 hours	75.3%
2017	Goyal et al.	8K	Tesla P100 x256	Caffe2	1 hour	76.3%
2017	Smith et al.	8K→16K	full TPU Pod	TensorFlow	30 mins	76.1%
2017	Akiba et al.	32K	Tesla P100 x1024	Chainer	15 mins	74.9%
2018	Jia et al.	64K	Tesla P40 x2048	TensorFlow	6.6 mins	75.8%
2018	Mikami et al.	34K→68K	Tesla V100 x2176	NNL	224 secs	75.03%

State-of-the-art Imagenet training (arXiv:1811.05233)

- Batch size control + LARS -> 68K mini-batch size
- 2D-Torus All-reduce communication

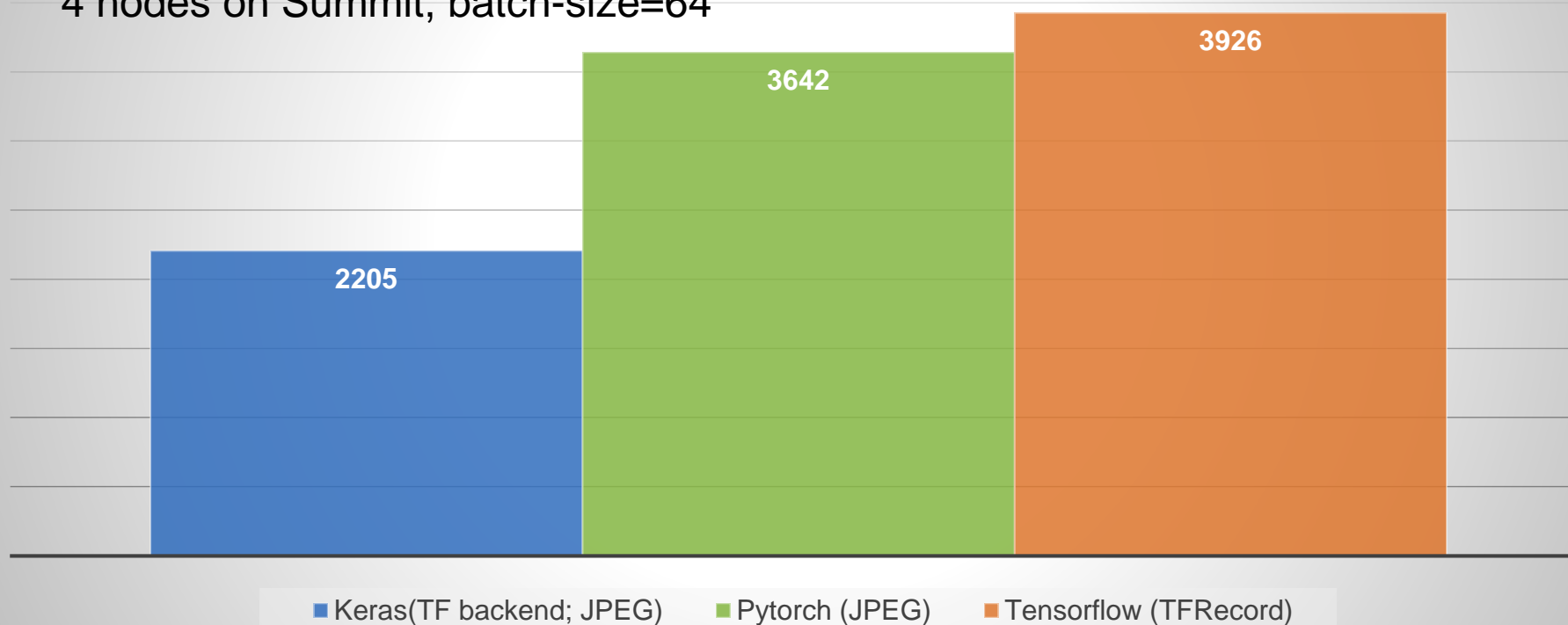


- 224s training -> 75.03% top1 accuracy and 66% scaling efficiency on 2176 V100.

Without tuning

Resnet50 on Imagenet (Images/s)

4 nodes on Summit, batch-size=64



Tuning of Tensorflow and Keras

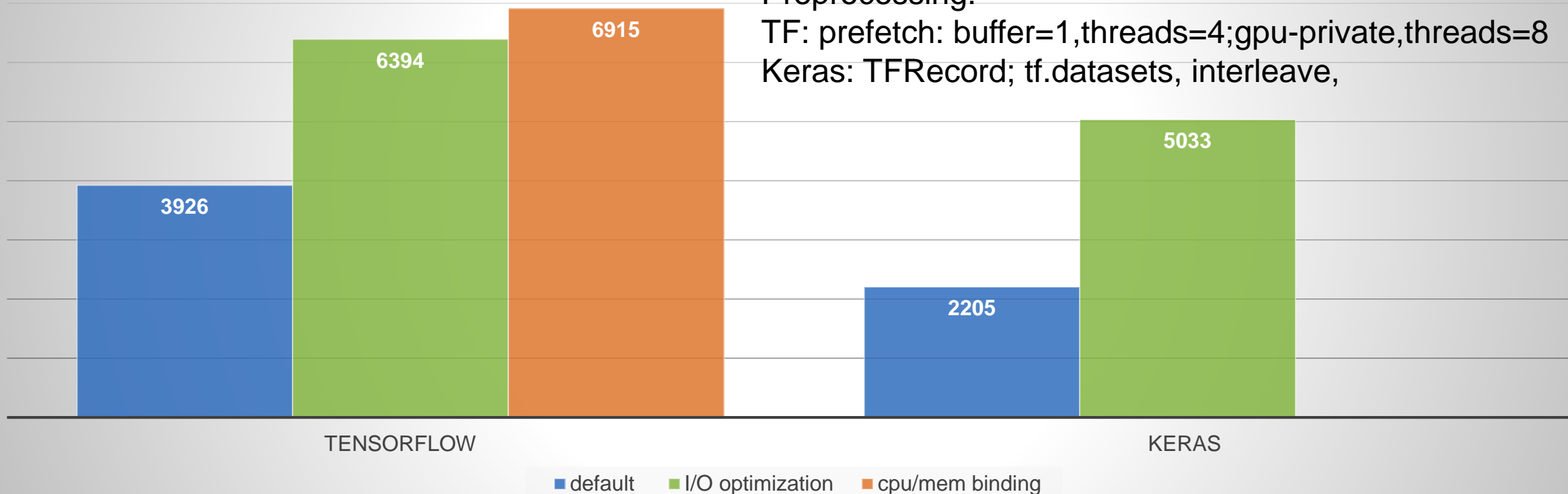
Performance Tuning

4 nodes on Summit, batch-size=64

Preprocessing:

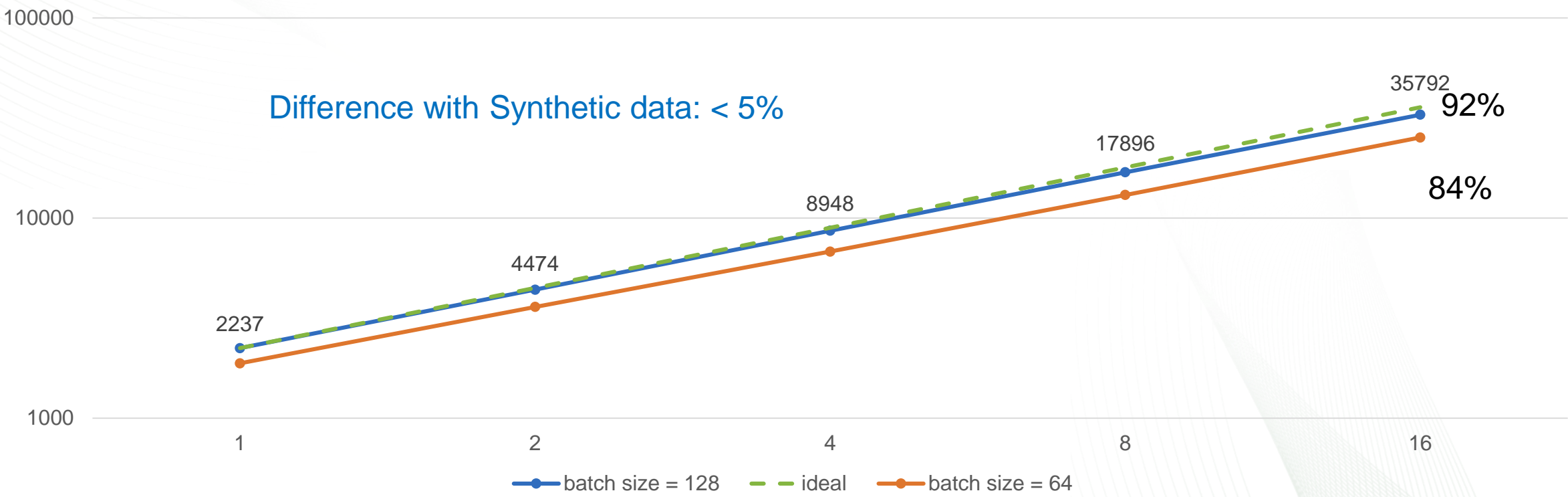
TF: prefetch: buffer=1, threads=4; gpu-private, threads=8

Keras: TFRecord; tf.datasets, interleave,



TF benchmark on Summit

TF CNN Benchmark Imagenet (TFRecord)



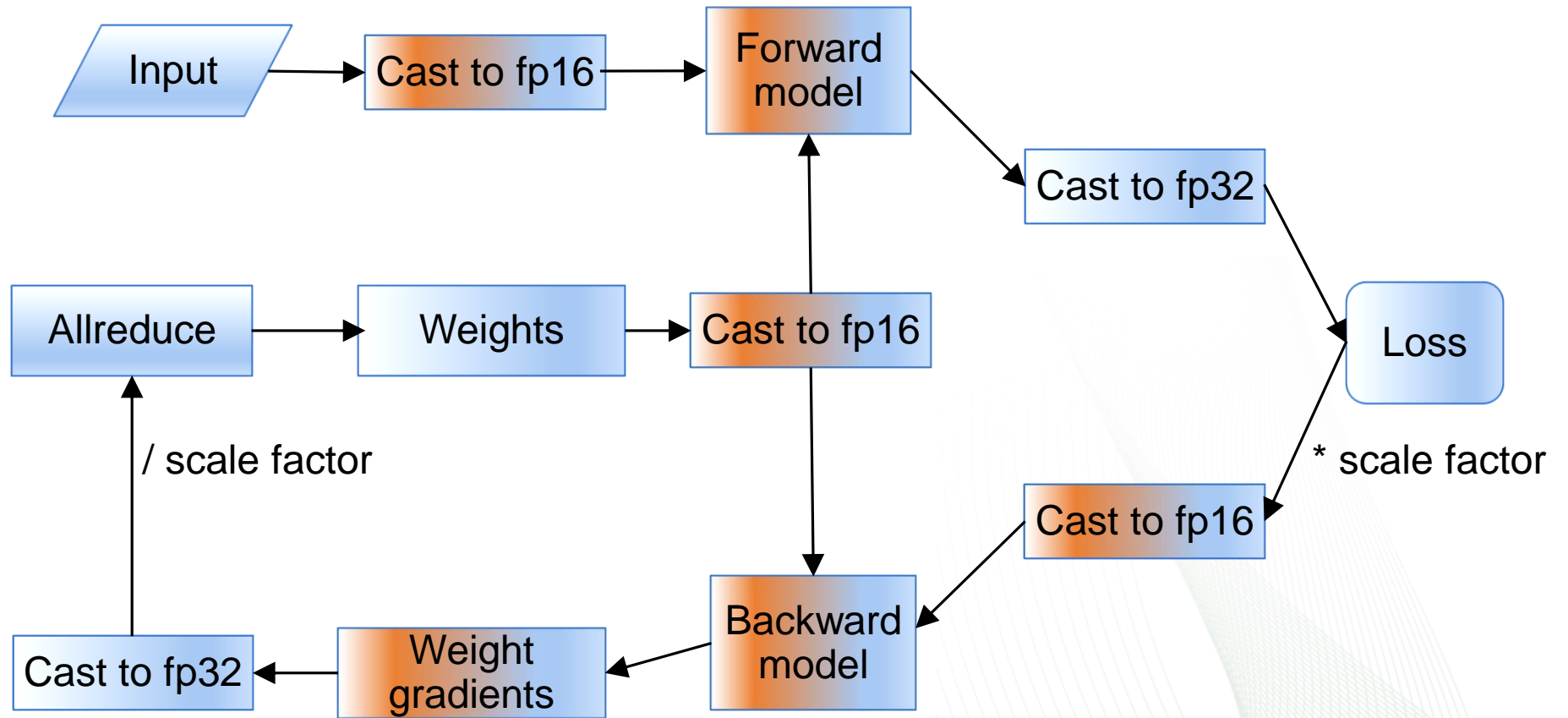
Mixed precision & Tensorcore

- Consideration

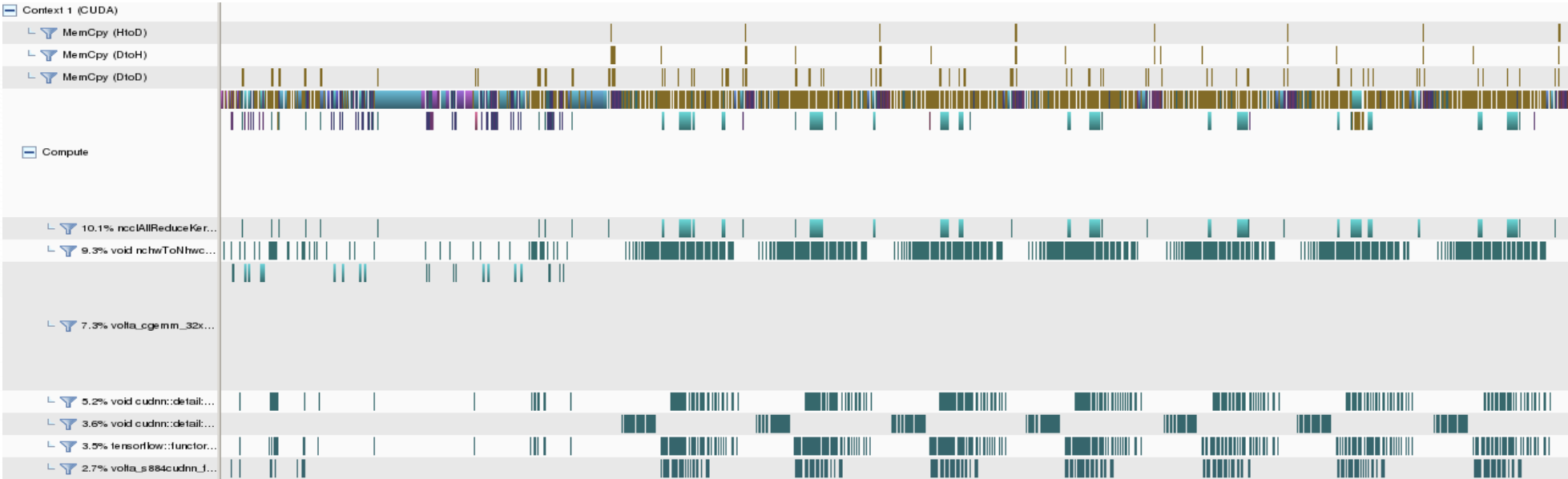
- Imprecise weights
- Gradients underflow
- Reduction overflow

- Verification

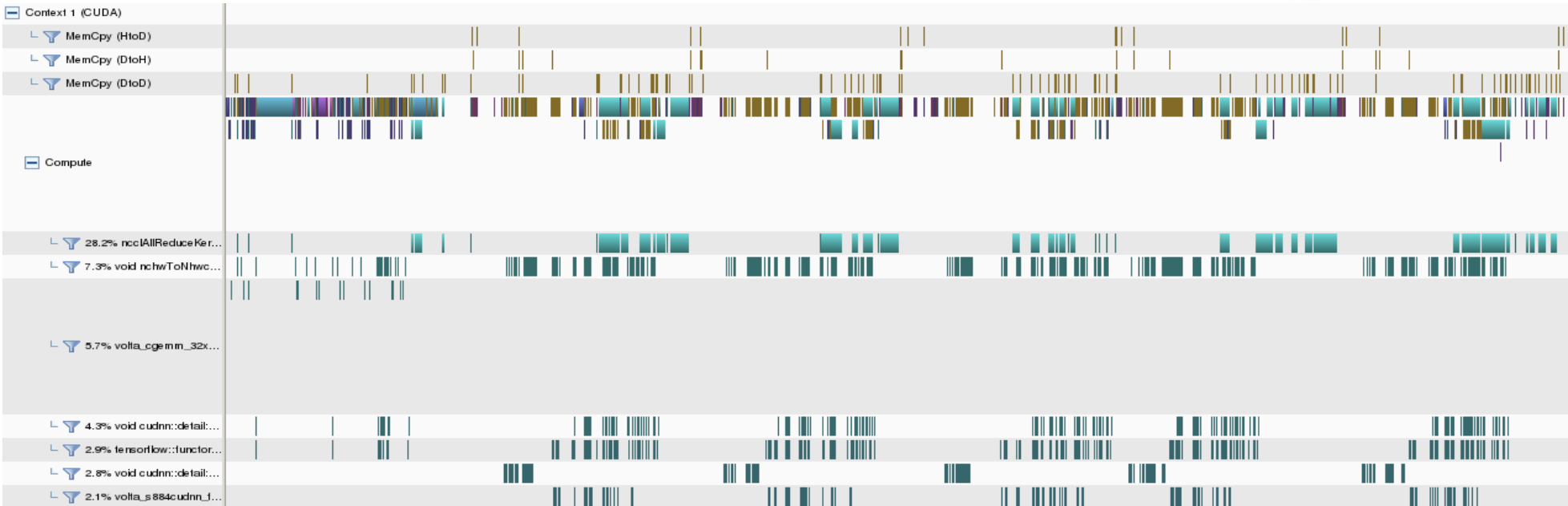
- s884cudnn



Synthetic Data



TFRecord



Lessons learned from Exa-scale DL on Summit (arXiv:1810.01993)

- Data ingestion (mostly coincide with TF performance guide)
 - Input pipeline, queueing input for compute
 - Concurrent processing with map
- Communication
 - Broadcast tree
 - Hierarchical aggregation of the control message (the order of tensors to be reduced)
 - Hybrid NCCL-MPI allreduce
 - NCCL intra node allreduce
 - 4 ranks (2 on each socket, b/c 4 IB devices) per node each MPI_Allreduce on a quarter of the data
 - NCCL intra node broadcast

Lessons learned from Exa-scale DL on Summit (arXiv:1810.01993)

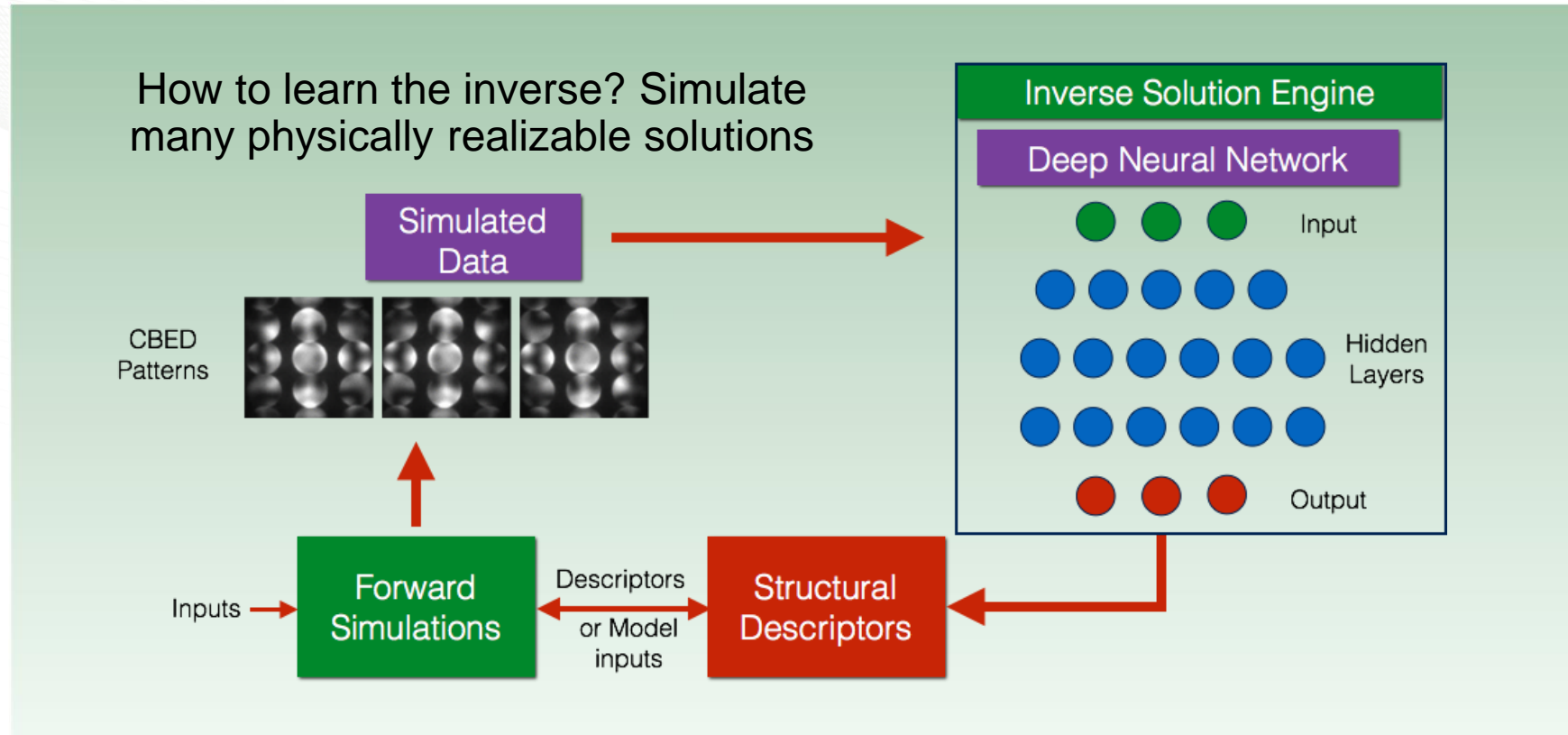
- Algorithmic considerations
 - Weighted loss, i.e. each pixel contributes differently to the loss function, specific to application (background vs area of interest)
 - LARC, a variant on LARS, for large batch sizes.
 - Multi-channel (16), more compute, more accurate
 - Gradient lag, overlap communication and computation
 - Network, larger layer, less number of layers, to improve compute intensity.

DL vs conventional ML

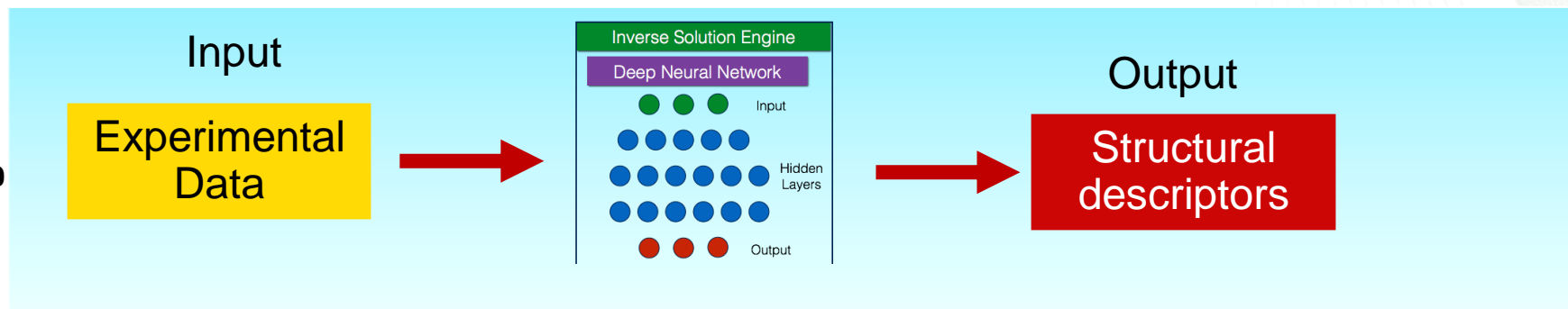
- It depends.
 - In general, DL works better for unstructured features, e.g. images, text; gradient boosting works better for data with structured ones, e.g. tabulated data; feature selection + gaussian process (equivalent infinite width neural network) works better for limited data and explainability.
- Explored in several use cases.
 - Simulation energy prediction
 - Material design (High entropy alloy)
 - Climate surrogate modelling
 - Microscopic images classification

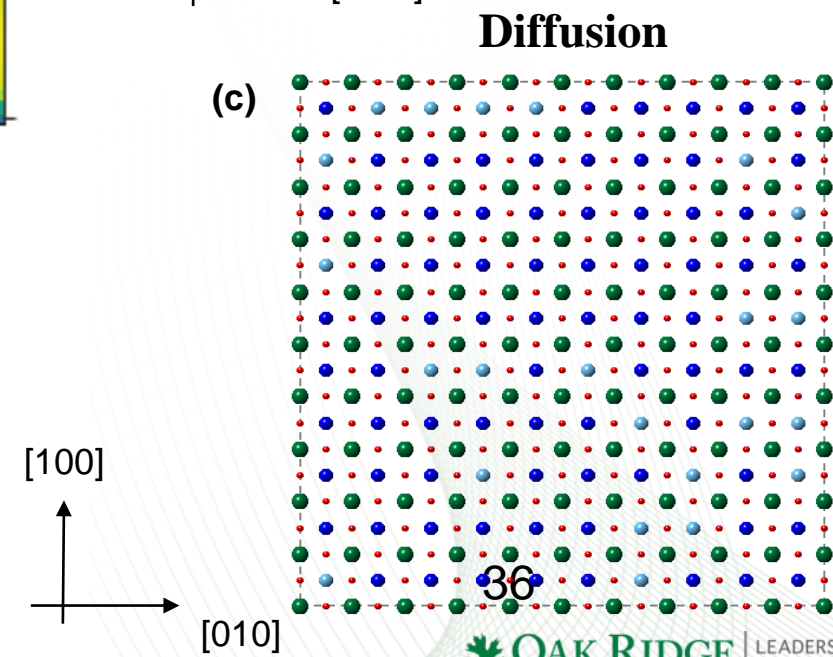
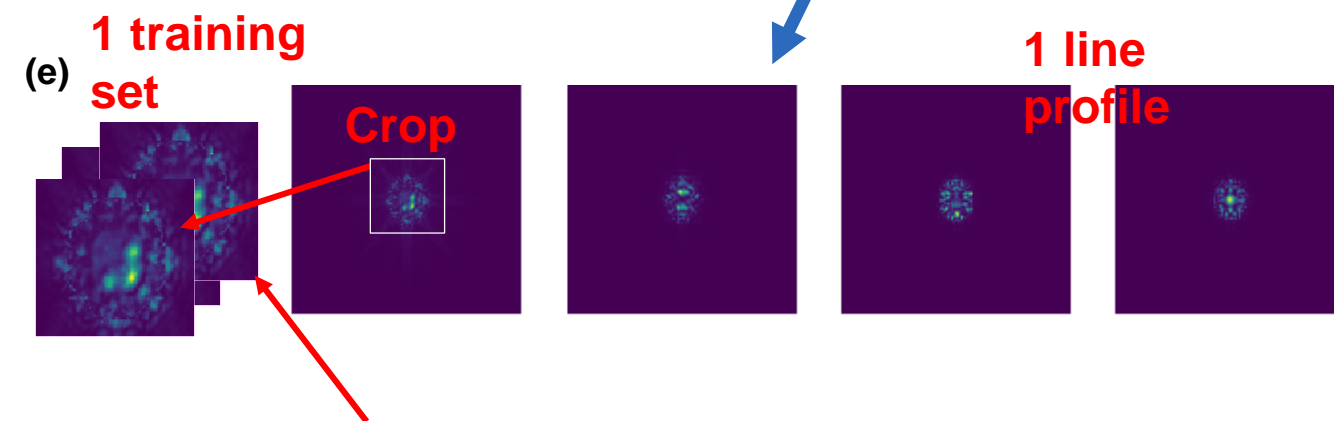
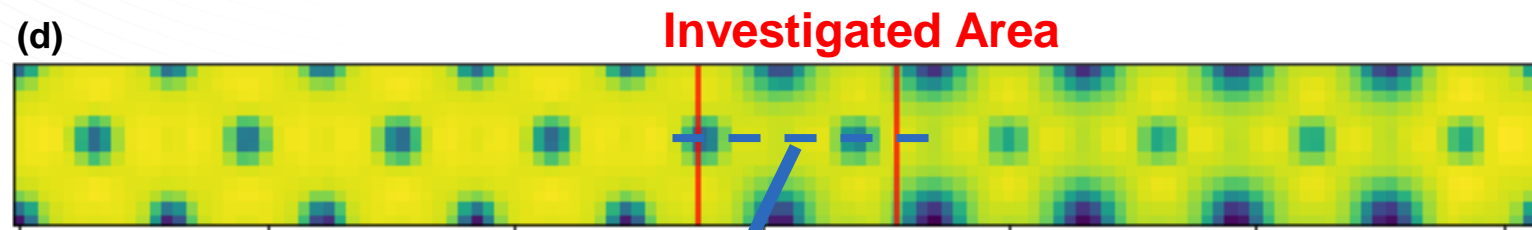
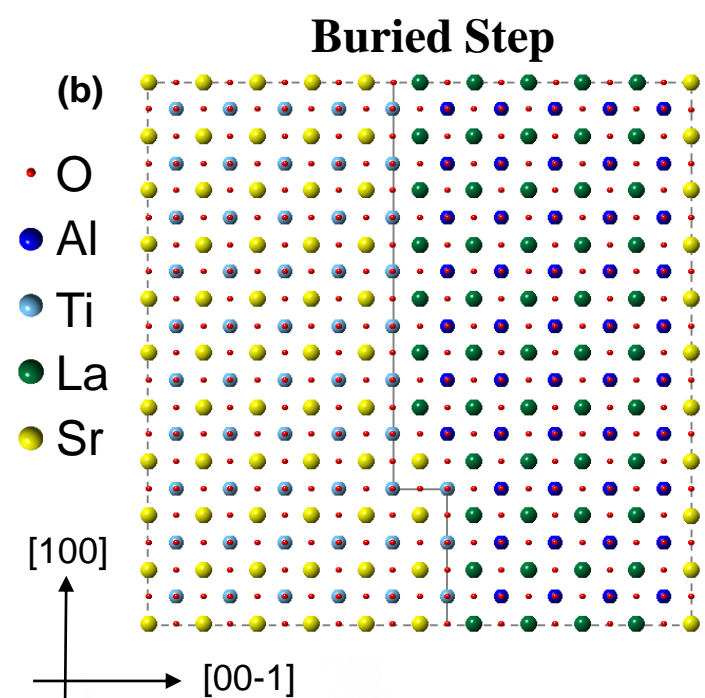
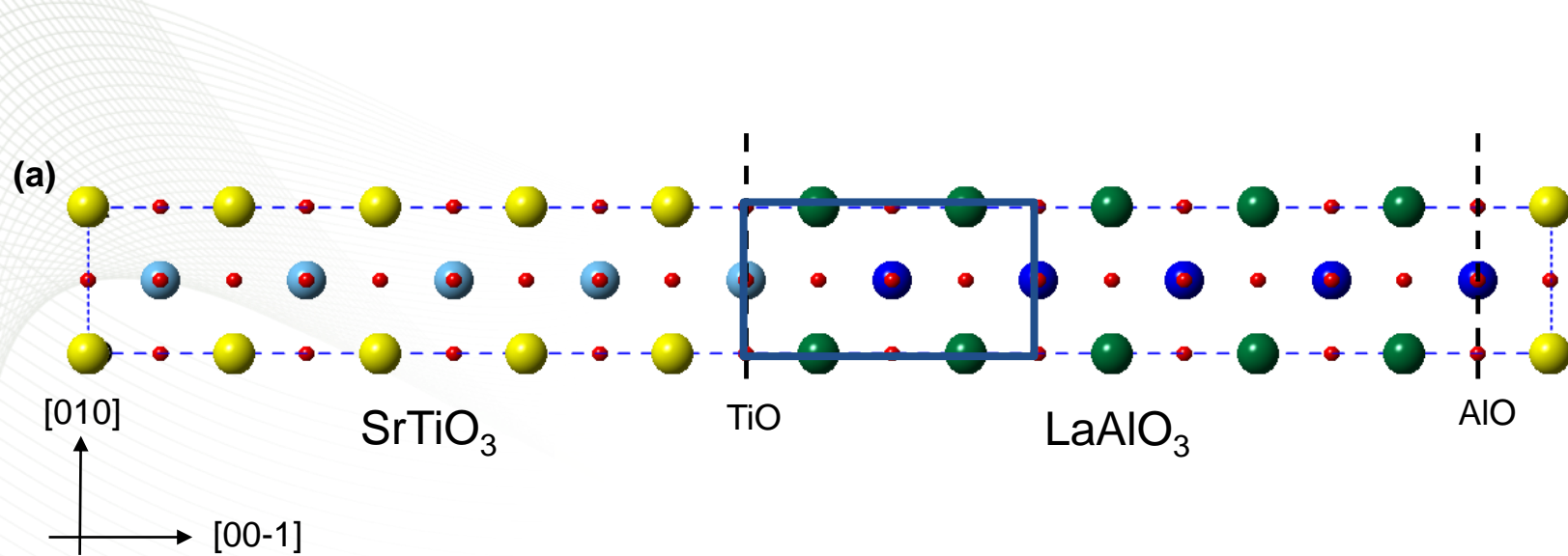
Backup slides: Use Case 1 (LDRD PI: Rama Vasudevan)

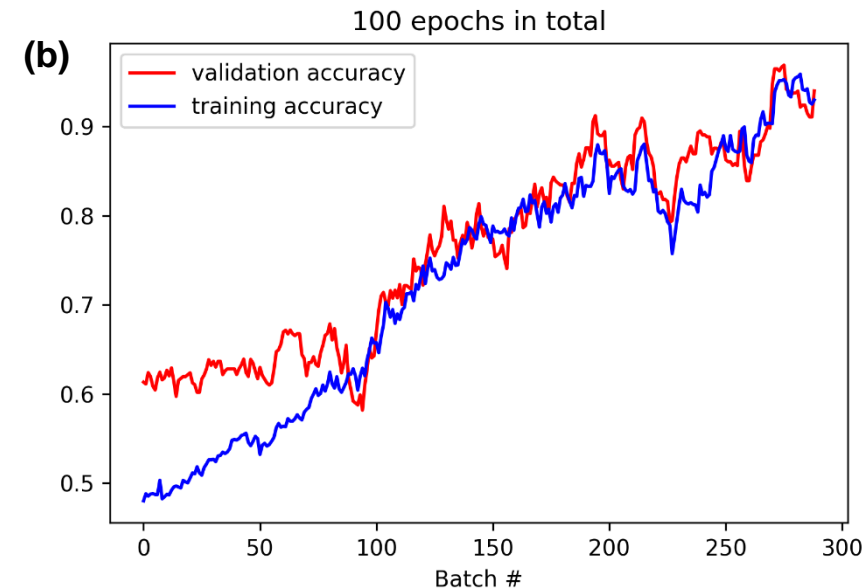
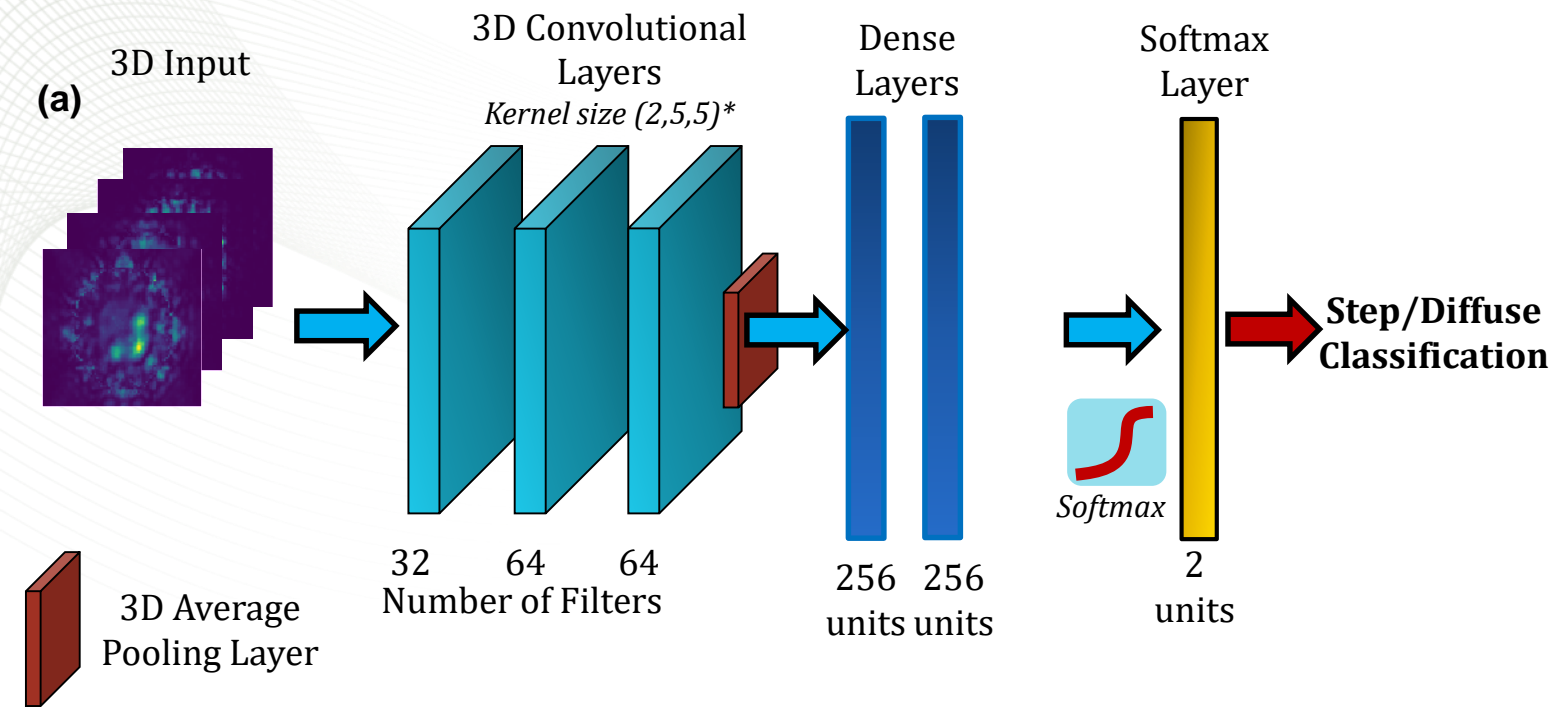
Training Phase



Testing Phase

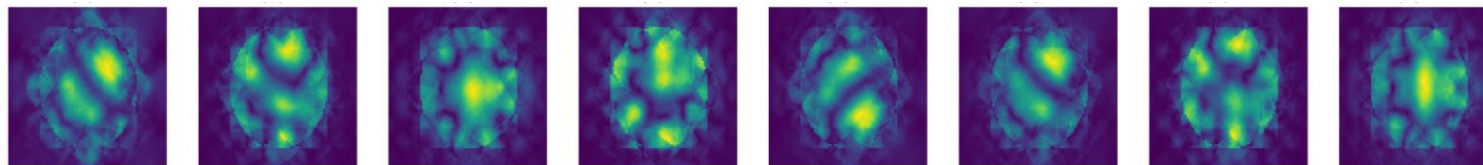




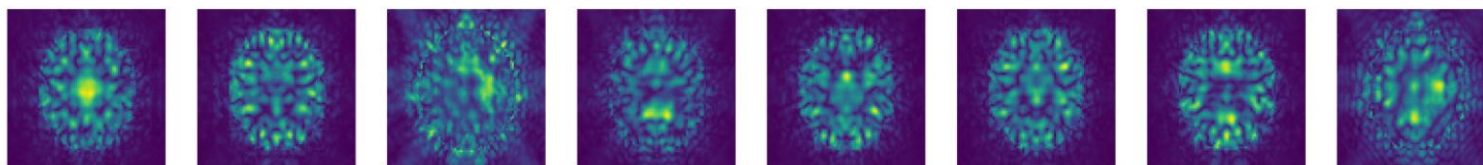


(c) Along interface →

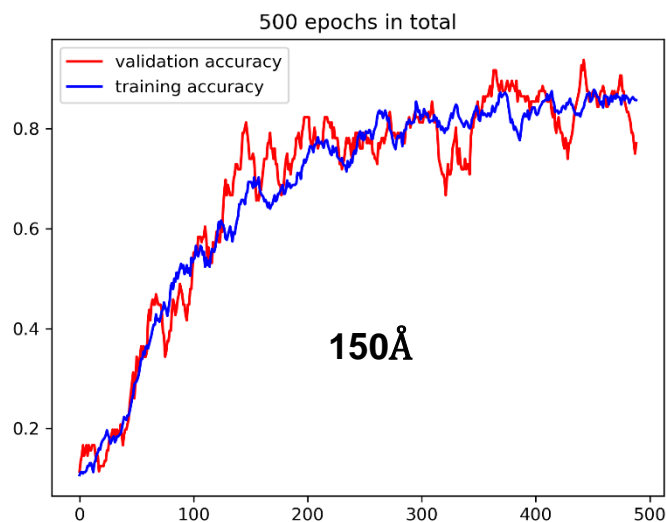
Prediction: Step ($p=0.96$). Actual: Step



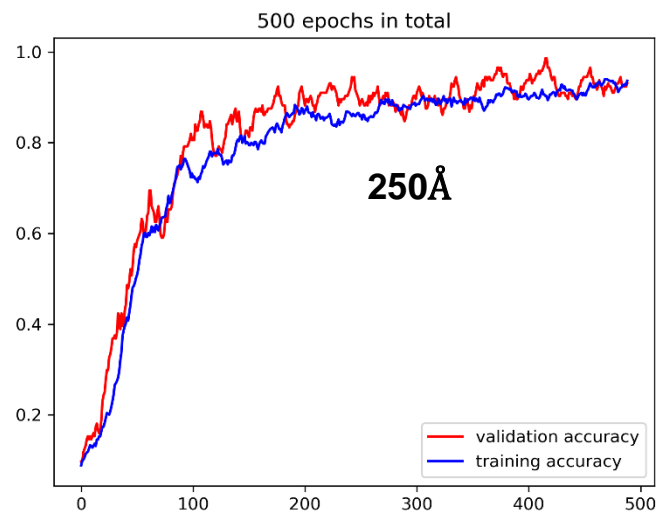
Prediction: Diffuse ($p=1.00$). Actual: Diffuse



(a)



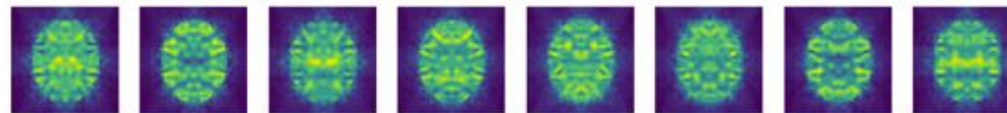
(b)



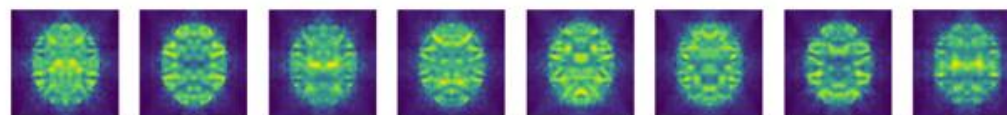
Thickness: 250Å

(c)

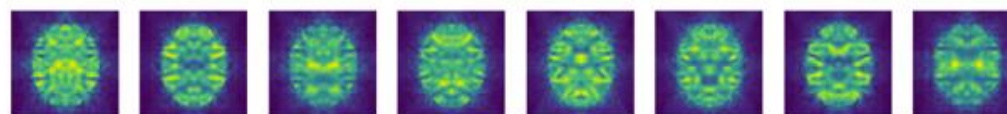
Prediction: 5 ($p=0.94$). Actual: 5



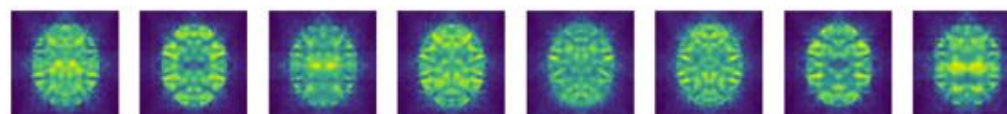
Prediction: 6 ($p=0.61$). Actual: 7



Prediction: 9 ($p=0.56$). Actual: 9



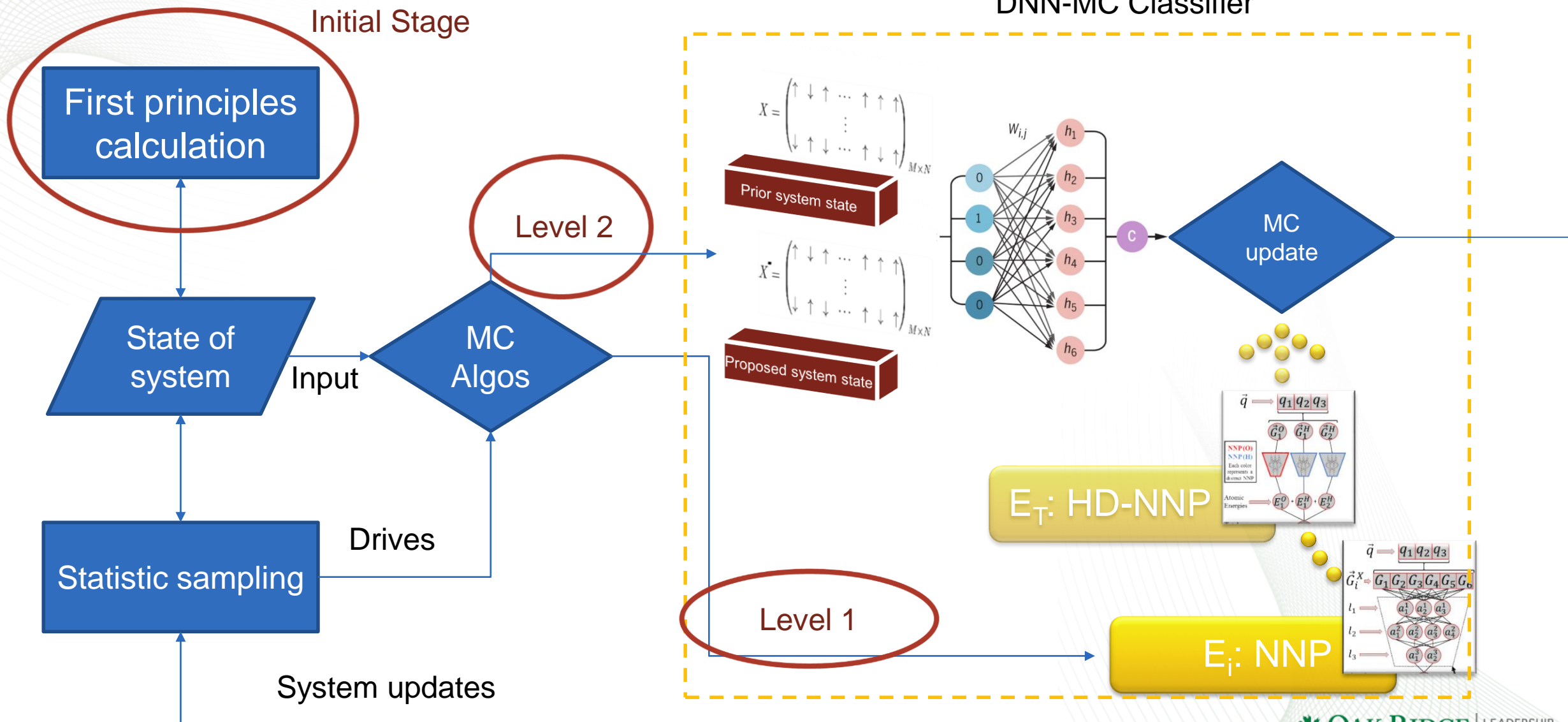
Prediction: 2 ($p=0.99$). Actual: 2



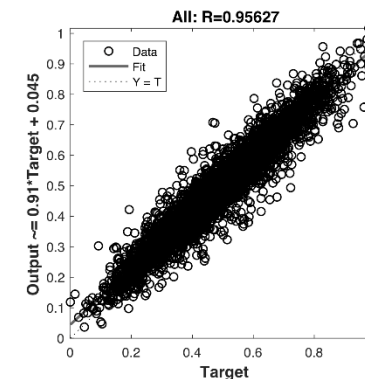
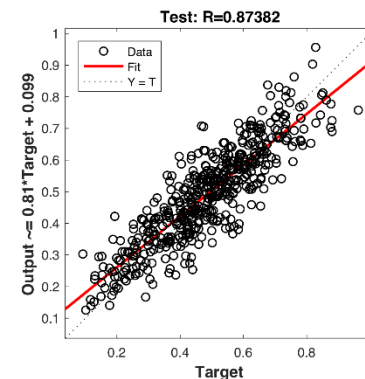
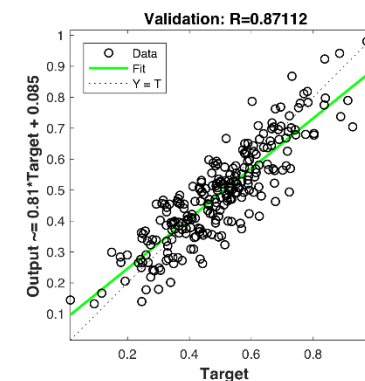
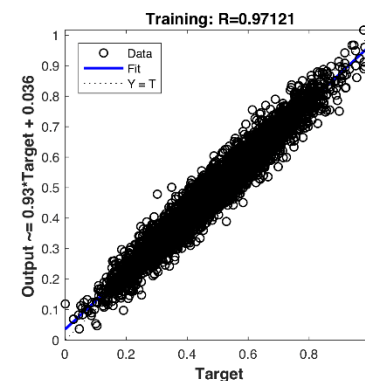
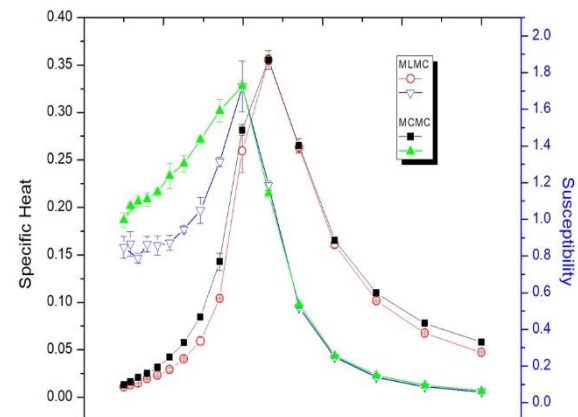
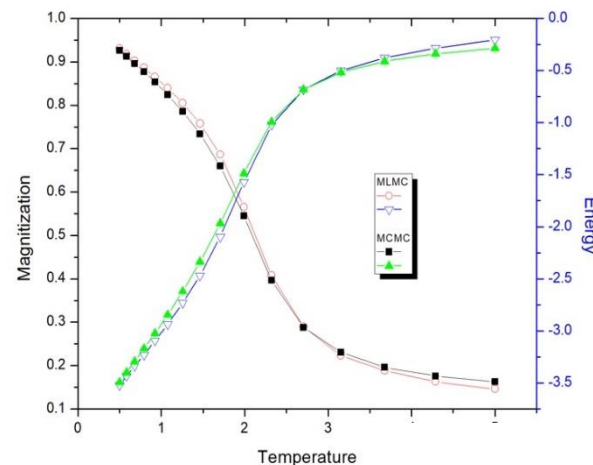
Use Case 2 (LDRD PI: Markus Eisenbach)

Initial Stage

DNN-MC Classifier

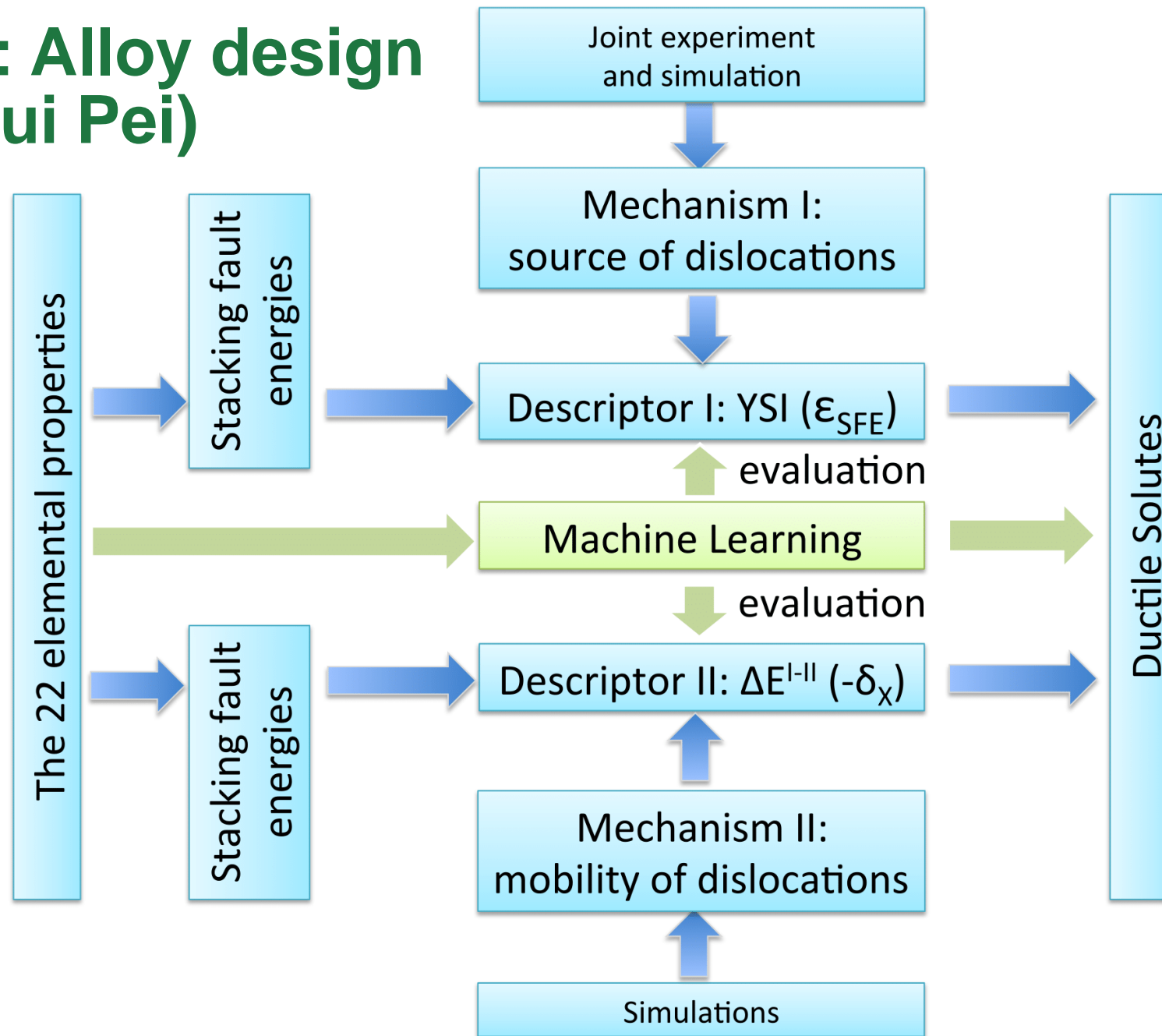


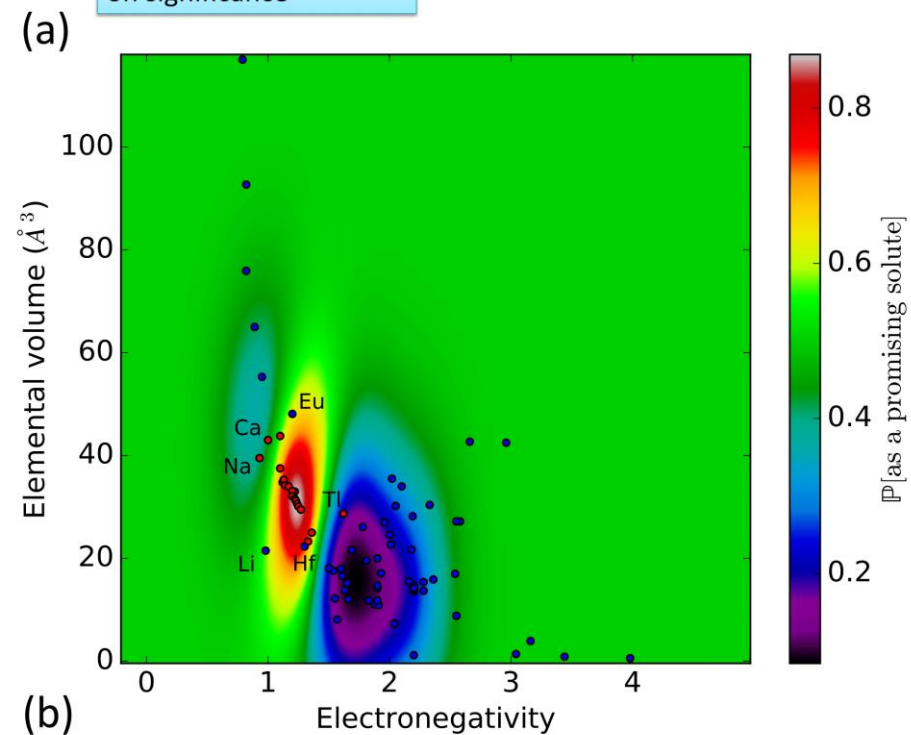
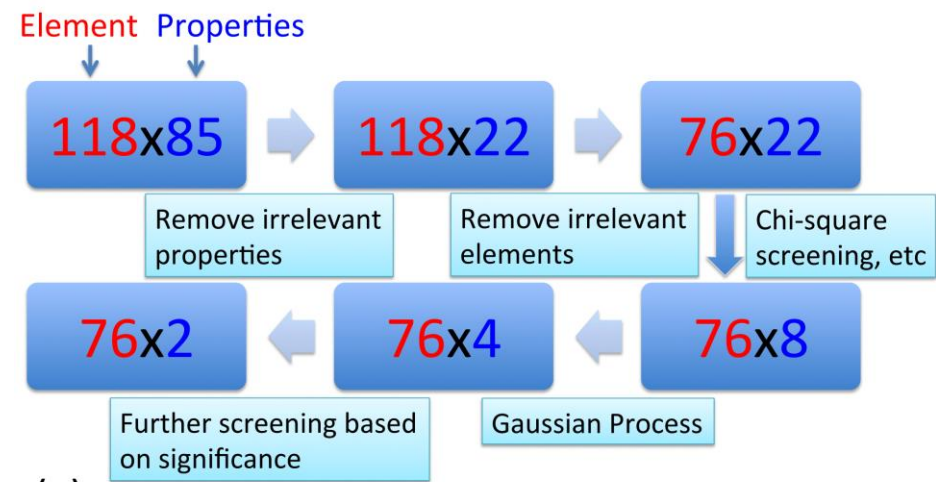
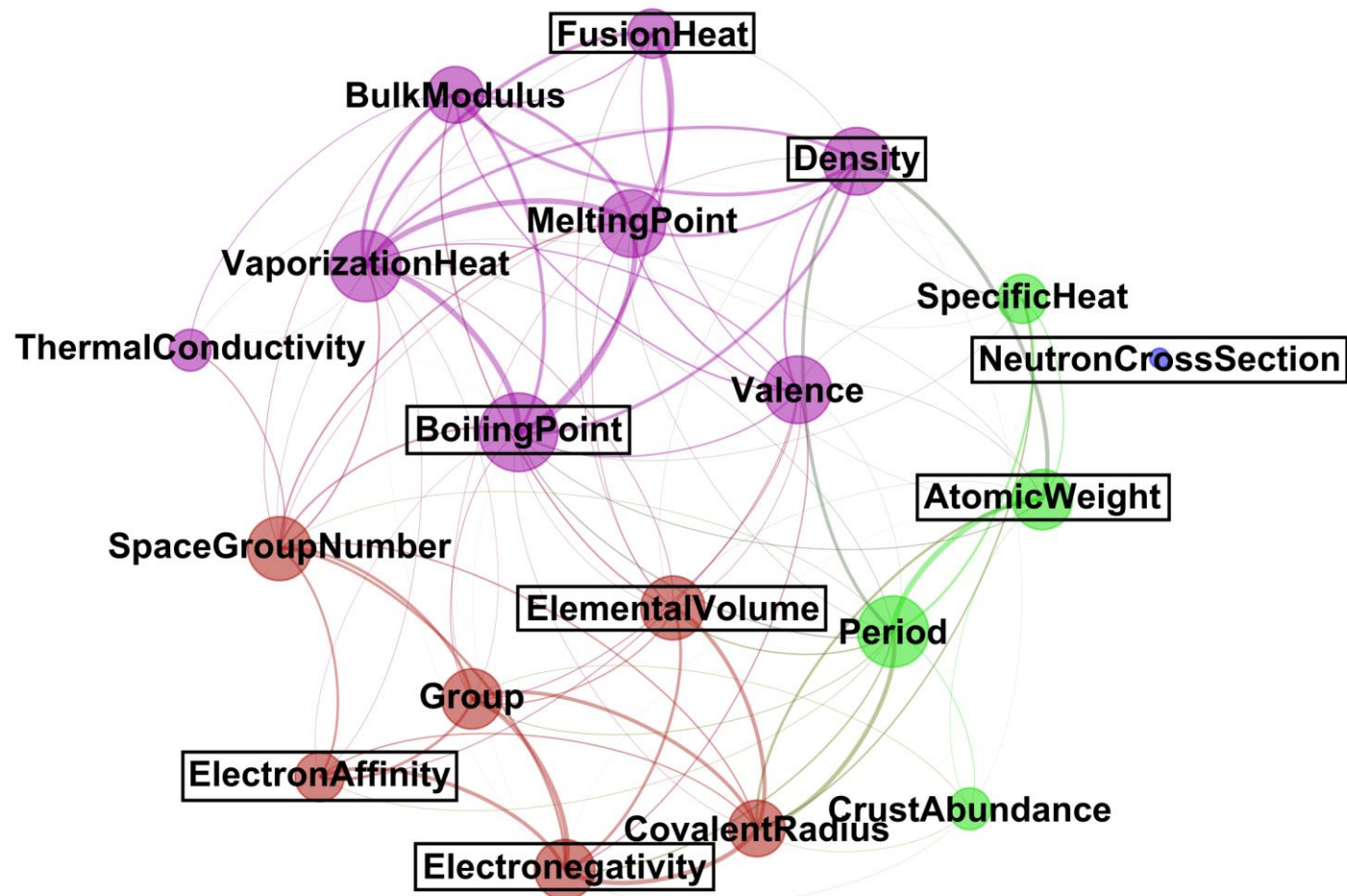
- Proof of concept for an online model (Heisenberg)
- Offline models for complex systems (Water cluster, FeCo alloy)
- Exploration of sampling algorithms (Metropolis, Wang-Landau, Nested Sampling)

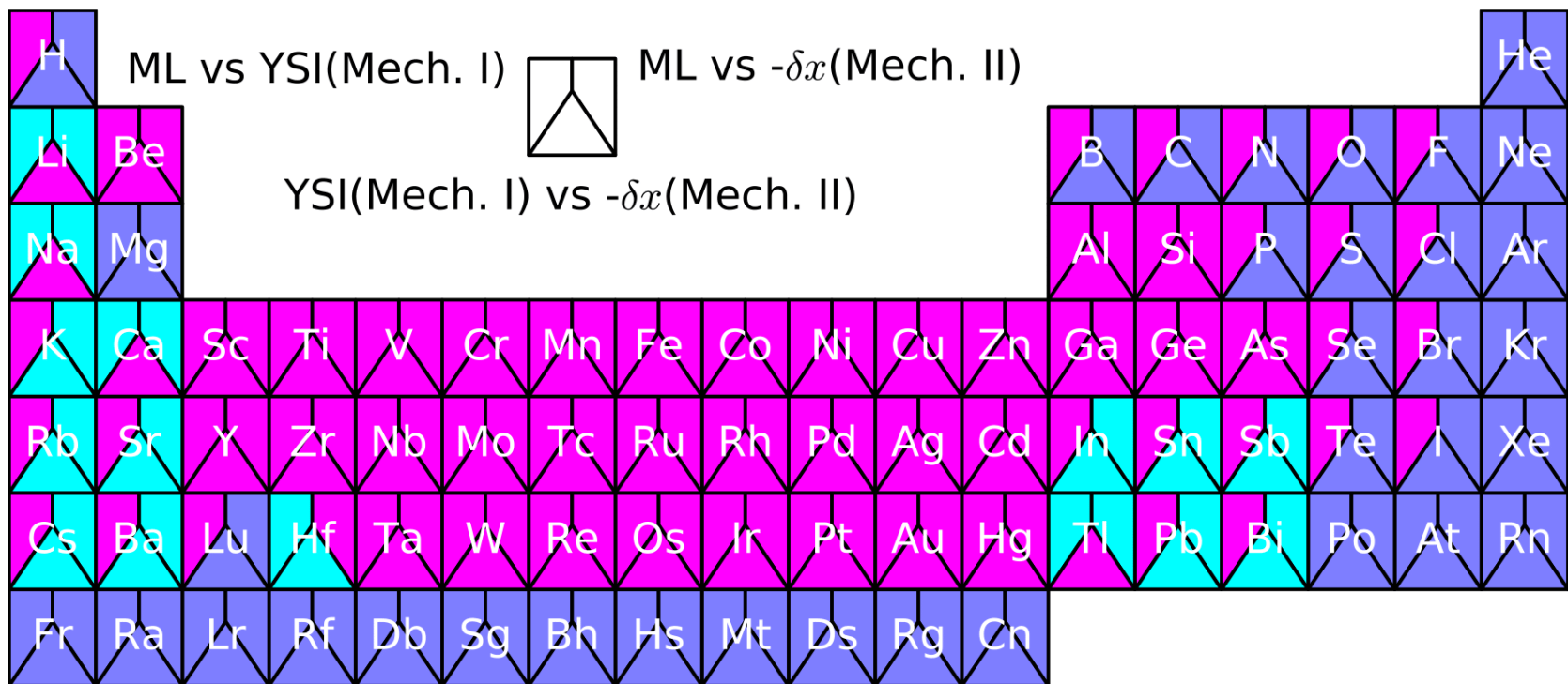


System	Sampling algorithms	Model	Accuracy
Heisenberg	Metropolis Nested Sampling	XGBoost DNN	87%
Water cluster	Wang-Landau	XGBoost DNN	91%
FeCo alloy	Metropolis	DNN	87%

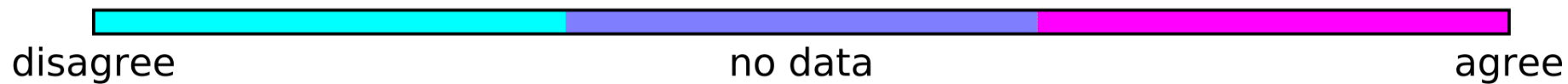
Use Case 3: Alloy design (with Zongrui Pei)







La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb
Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No



Questions?