

Efficient simulations of ML and LQCD

Denis Boyda, ALCF

Efficient Simulations on GPU hardware, October 24-27, 2022

Outline

- One slide intro to ML
- Strategy for ML for LQCD software development
- Scaling the number of ML experiments / MLOPs
- Weak Scaling / Increasing Batch Size
- Towards larger memory limits:
 - □ Checkpointing
 - Model (pipeline/tensor) parallelism or strong scaling
- Efficient Simulations on one GPU
 - □ Profiling PyTorch
 - □ Move functions to GPU / Autograd function
 - □ JIT compiler
 - □Custom kernels

One slide intro to ML



CRASH COURSE: DEEP LEARNING



"neural network"

Argonne 🛆

$$y = \sigma(A_3 \sigma(A_2 \sigma(A_1 x + b_1) + b_2) + b_3)$$

Argonne National Laboratory is a U.S. Department of Energy laboratory managed by UChicago Argonne, LLC

Slide credit: Bethany Lusch, ALCF

Forward path

- a layer get output from previous layer,
- does matmul with weights,
- applies activation function,
- sends data to the next layer.
- Computations repeated for all layers

Backward path

- Gradients are computed using chain rule
- Weights are updated using gradients

Training is done in batches. Stochastic gradient descent is used



Measurements



- Observables calculation takes similar or larger than generation resources
- · Use ML regression to compute them faster
- N configurations, N_{train} + N_{bs}measurements of 0, N_{train} measurements are used for training
- effective samples size is increased from N_{train} + N_{bs} to N
- · Correct bias



Can we decrease time for observable measurements and get higher statistics?



Normalizing flows: better uncertainty qualification at finite statistics

Lattice Schwinger model near criticality





Compute QCD phase diagram in (T, mu) with normalizing flows



Thermodynamical properties of QGP and EoS with normalizing flows



Increase signal to noise ratio via contour deformation



Can we apply it for viscosity computations in full QCD?

Reconstructing QCD Spectral Functions with Gaussian Processes



Did we improve a solution of inverse problem?

Slides credit: Denis Boyda, Machine Learning techniques in lattice QCD, 2022 RHIC/AGS Annual Users' Meeting, June 7, 2022



Machine Learning in Nuclear Physics

Amber Boehnlein, Markus Diefenthaler, Nobuo Sato, Malachi Schram, and Veronique Ziegler

Thomas Jefferson National Accelerator Facility 12000 Jefferson Avenue Newport News, Virginia, 1154

Cristiano Fanelli

Laboratory for Nuclear Science. Massachusetts Institute of Technology, Cambridge, MA 02139, USA The NSF AI Institute for Artificial Intelligence and Fundamental Interactions

Morten Hjorth-Jensen

Facility for Rare Isotope Beams and Department of Physics and Astronomy, Michigan State University, MI 48824, 1154 Department of Physics and Center for Computing in Science Education, University of Oslo, N-0316 Oslo, Norway

Tanja Horn

202

May

2

[nucl-th]

2

02309v

12.

arXiv:21

Department of Physics, The Catholic University of America, N. E. Washington, DC 20064, USA Thomas Jefferson National Accelerator Facility, 12000 Jefferson Avenue Newport News, Virginia USA

Michelle P. Kuchera

Department of Physics and Department of Mathematics & Computer Science, Davidson College. North Carolina 28035, USA

Dean Lee, Witold Nazarewicz, and Peter Ostroumov

Facility for Rare Isotope Beams and Department of Physics and Astronomy, Michigan State University, MI 48824. USA

Kostas Orginos

Department of Physics, William & Marv. Williamsburg, Virginia, 1154 Thomas Jefferson National Accelerator Facility 12000 Jefferson Avenue, Newport News, Virginia, USA

Alan Poon and Xin-Nian Wang

Nuclear Science Division, Lawrence Berkeley National Laboratory. 1 Cyclotron Road. Berkeley, California 94720, USA

Alexander Scheinker

Accelerator Operations and Technology Division Applied Electrodynamics Group, Los Alamos National Laboratory Los Alamos, New Mexico 87544, USA

Michael S. Smith

Physics Division Oak Ridge National Laboratory, Oak Ridge, Tennessee, 37831-6354, LISA

Long-Gang Pang

Key Laboratory of Quark and Lepton Physics, Institute of Particle Physics, Central China Normal University, Wuhan 430079, China

Advances in machine learning methods provide tools that have broad applicability in scientific research. These techniques are being applied across the diversity of nuclear physics research topics, leading to advances that will facilitate scientific discoveries and societal applications. This Colloquium provides a snapshot of nuclear physics research which has been transformed by machine learning techniques.

CONTENTS

2. Ensemble generation

D. High-Energy Nuclear Theory

A. Streaming Detector Readout

B. Reconstruction and Analysis

3. Particle identification

5. Event reconstruction

1. Charged particle tracking

4. Event and signal classification

1. Bayesian inference

4. Miscellaneous

IV. Experimental Methods

2. Calorimetry

6. Spectroscopy

3. Correlation function estimators

2. Inversion problems with ML

3. Other applications of ML methods

CONTENTS			C. Experimental Design
I.	Introduction	2	1. Design for detector systems 2. Interface with Theory
Π.	Machine learning for nuclear physics in broad strokes	3	V. Accelerator Science and Operations A. ML-based surrogate models for accelerator models
ш.	Nuclear Theory	5	B. Anomaly detection and classification
	A. Low-Energy Nuclear Theory	6	C. Control optimization
	 Early applications of machine learning Machine learning for data mining 	6 6	D. Adaptive ML for non-stationary systems
	3. Properties of heavy nuclei and nuclear density		VI. Nuclear Data
	functional theory	6	A. Overhauling the Nuclear Data Pipeline
	4. Nuclear properties with ML	7	B. Improving Compilations and Evaluations
	5. Nuclear shell model applications	7	C. Building Emulators and Surrogate Models
	6. Effective field theory and A-body systems	7	
	7. Nuclear reactions	9	VII. Summary and perspectives
	8. Neutron star properties and nuclear matter		
	equation of state	9	Acknowledgments
	B. Medium-Energy Nuclear Theory	10	
	1. Bayesian inference	10	References
	2. Simultaneus extraction paradigm	10	
	3. LQCD and experimental global analysis	11	
	C. Lattice QCD	11	
	1. The sign problem at non-zero density	11	I. INTRODUCTION

12

13

13

13

15

15

15

15

16

I. INTRODUCTION

This Colloquium represents an up-to-date summary of work in the application of machine learning (ML) in nuclear science, covering topics in nuclear theory, experimental methods, accelerator technology, and nuclear 13 data. An overview of the use of artificial intelligence (AI) 14 techniques in nuclear physics which aimed at identifying commonalities and needs has been provided in Bedaque et al. (2021).

2

18

18

18

19

19

20

21

21

22

22

23

23

23

Nuclear physics is a well-established field, with more than a century of fundamental discoveries covering a huge span of degrees of freedom, energy scales and length scales, ranging from our basic understanding of funda-

17 mental constituents of matter to the structure of stars



Machine Learning in Nuclear Physics

Amber Boehnlein, Markus Diefenthaler, Nobuo Sato, Malachi Schram, and Veronique Ziegler

Thomas Jefferson National Accelerator Facility, 12000 Jefferson Avenue, Newport News, Virginia, USA

Cristiano Fanelli

USA

2022

2 May

[nucl-th]

arXiv:2112.02309v2

Laboratory for Nuclear Science, Massachusetts Institute of Technology, Cambridge

Alexander Scheinker

Accelerator Operations and Technology Division Applied Electrodynamics Group, Los Alamos National Laboratory, Los Alamos, New Mexico 87544, USA

Michael S. Smith

Physics Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, 37831-6354, USA

Reduce time for one cycle

Nuclear Science Division, Lawrence Berkeley National Laboratory 1 Cyclotron Road, Berkeley, California 94720, USA

Alan Poo

B. Reconstruction and Analysis

- 1. Charged particle tracking
- 2. Calorimetry 3. Particle identification
- 3. Particle identification 4. Event and signal classification
- 5. Event reconstruction
- 6. Spectroscopy

et al. (2021).

15

- Nuclear physics is a well-established field, with more than a century of fundamental discoveries covering a huge
- span of degrees of freedom, energy scales and length
- scales, ranging from our basic understanding of funda-

2

17 mental constituents of matter to the structure of stars

7





Number of ML experiments in our data bases. Development of flow-base models for





MLOps (Machine Learning Model Operationalization Management) provides

- Hyperparameters/configuration tracking
- Live information (stdout, stderr, results)
- Artifacts (models, datasets) control and versioning
- Code control and versioning
- Environment configuration
- Fail trace

and an efficient way of analyzing experiments though

- Dashboard
- API to database











oainmo	ard ut	(ml_generation_u1)	su2 (mi_generation_s	u2) su3 (ml_generatio	on_su3) phi4 (mi_ge	neration_phi4) ferm	(ml_generation_ferm)	QCD (ml_generation_qcd)	
Status:	7 selected	•		F	ilters: Column Name	<pre></pre>	== V Enter	Value Value	Filter
	:: Id	:: Experiment Name	:: Hostname	:: Format	:: Command	:: Start Time	:: Status	:: Tags	:: Notes
	1386	phi4_rev	p54n1	MongoObserver-0.7.0	run_train	2022-08-31T16:33:37	COMPLETED	Add Tags 🗸	cf 325, AdamW
	1385	phi4_squeeze	p56n1	MongoObserver-0.7.0	run_train	2022-08-30T05:05:11	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 💉
	1384	phi4_squeeze	p56n4	MongoObserver-0.7.0	run_train	2022-08-30T02:29:37	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
	1383	phi4_squeeze	p56n2	MongoObserver-0.7.0	run_train	2022-08-29T11:53:13	COMPLETED	Add Tags 🗸 🗸	Enter Notes 🖍
	1382	phi4_squeeze	p55n4	MongoObserver-0.7.0	run_train	2022-08-29T11:53:02	COMPLETED	Add Tags 🗸 🗸	Enter Notes 🖍
	1381	phi4_squeeze	p55n1	MongoObserver-0.7.0	run_train	2022-08-29T11:53:02	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
	1380	phi4_squeeze	p54n3	MongoObserver-0.7.0	run_train	2022-08-29T11:52:53	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
	1379	phi4_squeeze	p56n3	MongoObserver-0.7.0	run_train	2022-08-29T11:51:46	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
	1378	phi4_squeeze	p51n4	MongoObserver-0.7.0	run_train	2022-08-29T11:51:46	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
•	1377	phi4_squeeze	p54n3	MongoObserver-0.7.0	run_train	2022-08-29T11:51:25	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍
	1376	phi4_squeeze	p55n4	MongoObserver-0.7.0	run_train	2022-08-29T11:51:24	PROBABLY_DEAD	Add Tags 🗸 🗸	Enter Notes 🖍



Compare 2 Runs [1386, 1385]





Compare 2 Runs [1386, 1385]

Captured Out Runs: 2 selected • Metric Label: Id • Run 1: 1386 • •	Metrics Plot						
Run 1: 1386 V { B: null, F =	Captured Out	Runs: 2 sele	cted - Metric Label: _id	•			
{ B: null,		Run 1:	1386		Run 2:	1385	~
<pre>k: nul, L: mui, M2: mui,</pre>	<pre>{ B: null, K: null, L:</pre>	<pre>>, Coupling", 111, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,</pre>					





Weak Scaling Increasing Batch Size

Introducing of Batch Size

In LQCD one deals with one configuration U_{μ} at every moment



Image credit to EPJ Web of Conferences 245, 09008 (2020), CC

Increasing Batch Size

- Speeds up training
- Increases GPU utilization
- Does not increase communications between GPUs

Increasing batch size is the most common way to increase GPU utilization and do distributed training in ML

In ML one work with **batch** of configurations at every moment due to Stochastic Gradient Descent



Image credit to EPJ Web of Conferences 245, 09008 (2020), CC



Scaling Batch Size Has Limits



C. J. Shallue et.al. Measuring the Effects of Data Parallelism on Neural Network Training. Journal of Machine Learning Research 20 (2019) 1-49, CC

Three scaling regimes:

- perfect scaling
- diminishing returns
- maximal data parallelism



Scaling Batch Size in Flow-based ML model for U(1) in 2D



Distributed Training: Data parallelism





Scaling Data Parallelism



Horovod was scaled on Summit up to 4600 nodes (1 Summit node = 6 Nvidia V100) with 93% efficiency



Figure 6. Scaling efficiency and Sustained Performance of distributed Deep Learning using the improved gradient reduction strategies up to 27,600 V100 GPUs.

Horovod has knobs to tune for reaching high efficiency!

O PyTorch

- Scales out of the box without proper knobs tuning
- Up to my best knowledge there are no numerical demonstrations of scaling up to 4-5k nodes, but there is no reason not to scale

PyTorch DDP should scale out of the box





Towards larger memory limits:

Memory Allocation sources in ML

There are two main sources of memory allocations in ML

Storing weights and gradients for model as well as optimizer states

bottleneck for large models

- Storing activations
- bottleneck for large configurations, LQCD!

Storing activations allows to save compute.

```
Without storing activations compute \sim O(n^2)
```

With storing activations compute $\sim O(n)$



Gradients checkpointing - 1

Through recomputation of activations during backward path memory cost can be reduced to $O(\sqrt{n})$, where n is number of layers of NN with only **extra forward path** (30% of computation time)

Tianqi Chen, Bing Xu, Chiyuan Zhang, Carlos Guestrin, Training Deep Nets with Sublinear Memory Cost, arXiv:1604.06174

Forward Path

- Input to every section is saved in memory
- Activations inside sections are not saved Backward path
- Recompute forward path inside a section with storing intermediate activations
- Do backward path as usual inside the section





Gradients checkpointing - 2



Forward Path

- Input to every section is saved in memory
- Activations inside sections are not saved Backward path
- Recompute forward path inside a section with storing intermediate activations
- Do backward path as usual inside the section





Gradients checkpointing – memory cost



when we identical layers and equal splitting of NN to integer number of sections

- *s* number of sections (checkpoints)
- n number of layers in NN
- A- activation size of one layer

with
$$\mathbf{s} = \sqrt{\mathbf{n}}$$
 total memory cost $\mathcal{C}_{total} \sim \sqrt{n}$

when sections contain different layers optimal number of checkpoints

$$s = \sqrt{C_{intermediate}/C_{input}}$$







when sections contain different layers optimal number of checkpoints

$$s = \sqrt{C_{intermediate}/C_{input}}$$





Towards larger memory limits: Model (pipeline/tensor) parallelism or strong scaling

Distributed training: Model vs Data parallelism

Parallelization schemes for distributed learning





Image source is Huihuo Zheng Talk at Argonne Training Program on Extreme-Scale Computing 2019 https://extremecomputingtraining.anl.gov/files/2019/08/ATPESC_2019_Track-8_6_8-9_130pm_Zheng-Data_Parallel_DL.pdf



Tensor parallelism



- For fully-connected layers and attention heads amount of communications is proportional to number of layers and layer size
- Collective required



Tensor parallelism for LQCD – halo exchange

In LQCD just split configuration and do halo exchange



- Only halo needs to be exchanged
- Point-to-point communications





Pipeline parallelism





Pipeline parallelism: efficient training

Efficient implementation is possible and requires advanced pipeline **schedulers** Yanping Huang et. al. *GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism.* 2018. arXiv:1811.06965



Memory cost per device \sim number of micro batches \sim m



Pipeline parallelism: advanced scheduling

More efficient pipeline scheduling allows to reduce bubble time ration and memory allocation

	Bubble time ration	Activations memory cost
GPipe Y. Huang et. al. GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism. 2018. arXiv:1811.06965	$\frac{p-1}{m}$	<i>O</i> (<i>m</i>)
PipeDream-Flush D. Narayanan et.al. Memory-Efficient Pipeline-Parallel DNN Training. arXiv:2006.09503	$\frac{p-1}{m}$	<i>O</i> (<i>p</i>)
Interleaved Stages D. Narayanan et.al. Efficient Large-Scale Language Model Training on GPU Clusters Using Megatron-LM. arXiv:2104.04473	$\frac{1}{v}\frac{p-1}{m}$	<i>O</i> (<i>p</i>)

- *p* number of pipeline stages (devices)
- m number of micro-batches
- v number of model chuncs





Efficient Simulations on one GPU Profiling PyTorch Move functions to GPU / Autograd function JIT compiler Custom kernels, Sycl/cupy/Numba, DLpack

Profiling PyTorch

Python cProfile





Move all functions to GPU / Autograd function

torch.linalg.eig(A, *, out=None)

• NOTE

When inputs are on a CUDA device, this function synchronizes that device with the CPU.

```
class DiagonalizeSUN(torch.autograd.Function):
    @staticmethod
    def forward(ctx, U_re, U_im, *, diagonalize=_simple_np_diagonalize):
        U = ComplexTensor(U_re, U_im)
        d, P = diagonalize(U)
        ctx.save_for_backward(d.real, d.imag, P.real, P.imag, U.real, U.imag)
        return d.real, d.imag, P.real, P.imag
    @staticmethod
    def backward(ctx, gd_re, gd_im, gP_re, gP_im):
        d_re, d_im, P_re, P_im, U_re, U_im = ctx.saved_tensors
        gU_re, gU_im = grad_diagonalize(d_re, d_im, P_re, P_im, U_re, U_im, gd_re, gd_im, gP_re, gP_im)
        return gU_re, gU_im
```

def diagonalize_suN(U): d_re, d_im, P_re, P_im = DiagonalizeSUN.apply(U.real, U.imag) return ComplexTensor(d re, d im), ComplexTensor(P re, P im)



Just in time (JIT) Pytorch compiler



Fuse a lot of small cuda kernels to one kernel PYTORCH_NVFUSER_ENABLE=complex

@torch.jit.script

def _torch_diagonalize_su3(U):





Custom Kernels

Write custom CUDA or SYCL kernels https://pytorch.org/tutorials/advanced/cpp_extension.html

or bind to existing LQCD codes



pybind11

or use NUMBA



diag	(torch): 21.35-23.43ms
diag	(jit) : 18.01-20.22ms
diag	(sycl) : 3.48-4.00ms
diag	(numba) : 2.77-3.07ms



Summary: Invest time in development only when necessary

Acknowledgment

This work is supported by the Argonne Leadership Computing Facility, which is a U.S. Department of Energy Office of Science User Facility operated under contract DE-AC02-06CH11357



Argonne Leadership Computing Facili