## Malachi Schram

Head of the Data Science Department On behalf of the research from the JLab Data Science Department Newport News, Virginia June 13th, 2023



TJNAF is managed by Jefferson Science Associates for the US Department of Energy Mission:

- Provide world-class data science solutions to advance research in <u>nuclear physics</u> by working with the subject matter experts at Jefferson Lab, partnering universities and Labs, and the Department of Energy.
- Provide world-class data science solutions to scientific applications relevant to the regional scientific community

Vision:

- Expand the <u>capability</u> and <u>capacity</u> of data science at JLab
- Create a <u>collaborative</u> data science research hub to:
  - 1. Work with regional partners on challenging scientific problems
  - 2. Champion education and research opportunities with regional universities and industry
  - 3. Reduce the carbon footprint by optimizing the data science workflow and algorithms



# **Current Portfolio**

DOE Nuclear Physics:

- Quantom SciDAC (with ANL, VTech, ODU)
- Working with the experimental Halls (Tracking, etc.)
- Data Science contributing effort for AIEC (lead by EPSCI)

DOE Basic Energy Science:

- Machine Learning for Improving Accelerator and Target Performance (with ORNL)
- Collaborating with SLAC on application of ML-based controls for accelerators

DOE Advanced Scientific Computing Research:

- Data-Driven Decision Control for Complex Systems (with PNNL, ORNL, UC) Non-DOE:
  - Hampton Roads Digital Twin (with ODU)

Laboratory Directed Research & Development (LDRD):

- Multi-objective Optimization of Heat Load and Trip Rates in CEBAF (FY22)
- Adaptive Strategies for Optimal Computing Availability (FY23)



# **JLab Data Science Pillars**

# • Applications:

- Nuclear Physics
- Advanced Scientific Computing
- Health & Climate
- Focused Methods & Algorithms:
  - Uncertainty Quantification
  - Interpretability and Explainability
  - Design & Control
- Infrastructure:
  - JLab ML & Data Hub
  - JLab Data Science software



# **DOE ASCR - BRN for SciML**



Figure 1: Foundational research themes of SciML must tackle the challenges of creating domainaware, interpretable, and robust ML formulations, methods, and algorithms.



Figure 2: Opportunities for SciML impact arise in scientific inference and data analysis; in MLenhanced modeling and simulation; in intelligent automation and decision support; and in related applications.



#### **BASIC RESEARCH NEEDS FOR** Scientific Machine Learning Core Technologies for Artificial Intelligence





oplimal decisions for complex systems

# **Uncertainty Quantification for ML**

Develop methods that include uncertainty estimates in machine learning models

- <u>Applications</u>:
  - Data driven ML-based surrogate models
  - Real time controller
  - Anomaly detections
- <u>Requirements</u>:
  - Out-of-distribution uncertainties
  - Auto-calibration
  - Single inference
- Hardware considerations:
  - Memory
  - Inference time
  - Performance trade-off due to approximations





# **Uncertainty Quantification For Reliable AI/ML**

- Goal: Develop and integrate new UQ methods for reliable AI/ML
- Problem space:
  - Applications with high-dimensional continuous input features
  - Focused on large data sets for DOE applications
  - Safety constraints that should never or at least rarely be violated.
  - Inference that must happen in real-time at the control frequency of the system.
- Applications:



son Lab



I Ampl

# **Uncertainty Aware Siamese Model ("Classification")**

- We enhanced our models by adding GP approximation layer which provides the uncertainty estimate
- Results from similarity model showed a ~4x improvement in performance over previously published results, it is also much better than a vanilla Auto-encoder
- The ROC curves show true fault detection rate above 60% while keeping the false alarms below 0.5% (not optimized)
- We introduced an out-of-domain anomaly, labelled 1111 (red), the UQ-based model performed similar in classifying the anomalies and indicated high uncertainty (as expected)

W. Blokland, K. Rajput, M. Schram, T. Jeske, et al 2022 Phys. Rev. Accel. Beams 25, 122802



# Data Driven UQ ML-based Surrogate Models (Regression)

- Compare different techniques: DQR, BNN, DNGPA
- DQR models have great performance for training distribution but not for OOD
- BNN models do a better job to estimate OOD
- DGPA models are distance aware by design resulting in better OOD estimation



M. Schram, K. Rajput, et al 2023 Phys. Rev. Accel. Beams 26, 044602



#### • Goal:

- Develop a UQ model to predict the capacitance of 3 capacitors in the SNS HVCM system using currently available sensor data
- Techniques:
  - Singular Value Decomposition, Residual Networks, Distance Preservation and bi-Lipschitz Constraints, Gaussian Process Approximation
- Results:
  - Achieves <1% error with accurate UQ for in-distribution data
  - Provides significantly better OOD performance than other methods
  - Submitted to Machine Learning with Applications

Distance Preserving Machine Learning for Uncertainty Aware Accelerator Capacitance Predictions

Steven Goldenberg<sup>a,\*</sup>, Malachi Schram<sup>a</sup>, Kishansingh Rajput<sup>a</sup>, Thomas Britton<sup>a</sup>, Chris Pappas<sup>b</sup>, Dan Lu<sup>b</sup>, Jared Walden<sup>b</sup>, Majdi I. Radaideh<sup>c</sup>, Sarah Cousineau<sup>b</sup>, Sudarshan Harave<sup>b</sup>

<sup>a</sup> Thomas Jefferson National Accelerator Facility Newport News VA 23606 USA <sup>b</sup>Oak Ridge National Laboratory Oak Ridge TN 37830 USA <sup>c</sup>Department of Nuclear Engineering and Radiological Sciences, The University of Michigan, Ann Arbor, Michigan 48109, USA





# MULTI-MODULE CVAE TO PREDICT HVCM FAULTS IN THE SNS

- <u>Goal</u>: Predict an upcoming machine failure before it occurs to improve the reliability of the HVCMs and reduce the down time for the SNS facility
- **How:** We use pulses leading to failure to allow for future forecasting
- <u>Method</u>: Condition the Variational Autoencoder (VAE) model on the module unique identifier to learn the association between waveforms and their modules
  - By using all 15 modules, we eliminate the need to train a Singlemodule and increase the number of statistics
- <u>Paper</u>: Accepted for publications at Machine Learning with Applications journal: <u>https://arxiv.org/pdf/2304.10639.pdf</u>



Figure 2: Compare the AUC values between single-module and multi-module using six types of faults across several modules. The error bar is ± 1 Standard Deviation (SD) error generated by sampling the latent Z of each Method.



Figure 1: KDE distributions of the MSE from reconstructing normal (green) and faulty waveforms (red) for three fault types, with the corresponding ROC curve for each fault.



Figure 3: Loss surface, where x- and y-axis are two random directions in weights space generated using filter normalization method. Multi-Module shows convex-like surface, while Single-module hs chaotic behaviour.



# Improving System Controls Through AI/ML

- Goal: Apply UQ based AI/ML methods to improve system controls
- Applications: Detector and Accelerator controls
- Techniques: Bayesian Optimization, Genetic Algorithms, Model Predictive Control, Reinforcement Learning



T. Jeske *et al* 2022 *JINST* **17** C03043 T. Jeske *et al* 2023 *J. Phys.: Conf. Ser.* **2438** 012132 June 15, 2023



M. Schram *et al* 2022 Phys. Rev. Accel. Beams 24, 104601 Jefferson Lab

# **GRAPH NEURAL NETWORKS FOR TRACKING**

#### **Initial Results**

- **Goal:** Improve charged particle tracking with the help of Graph Neural Networks (GNNs) and potentially run in near real time on edge computing
- Application: Charged Particle Tracking at GlueX, CLAS, EIC, and many more...
- Techniques: Graph Analytics, GNNs, GNNs on FPGA, Physics Informed ML



Predicting True edges belonging to tracks with GNNs

Train AUC:0.9927

1.0

0.8

9.0 Rate

0.4

0.2

Positive

# FEMTOSCALE IMAGING OF NUCLEI USING EXASCALE PLATFORMS

- The goal is to extract a quark and gluon tomography of nuclei and answer important questions on the nature of visible matter at the femtoscale
- Develop modular components to dynamically compose workflows
- Multiple AI/ML components that need to scale to LCFs

Module 1







Optimize QCF parameters

# FRAMEWORK

- Developed a common and modular framework
- Includes:
  - Core base classes
  - ProxyApp, GPDs Theory
  - Experimental (filter, detector, etc.)
  - GAN workflows

cfg	Updated visualization tools
core	Implemented visulation tools and helper scripts
demos	Implement requested changes and information capture
discriminator_module	fix useBias bug in disc
eventselection_module	Implemented visulation tools and helper scripts
expdata_module	Update with code from pub_demo that appeared to be missing.
experimental_module	Add first readme version for the experimental module
generator_module	configurable-N events per parameter set
sample_data	Add sample data for test use cases.
theory_module	Updated visualization tools
utils	Updated visualization tools
workflow	Updated visualization tools

😱 schr	476 Include MacOS env builde	er (testing). 8e78elf 2 weeks a	go 🕲 111 commi
tomo	ography_toolkit_dev	Updated visualization tools	last mon
utes 📄	ts	Implement requested changes and information capture	last mon
🗋 .gitiç	nore	add scripts dir to ignore	last mon
🗅 ном	/TO_CONTRIBUTE.md	Write documentation file: HOWTO_CONTRIBUTE.md	last mon
🗋 REA	DME.md	Merge branch 'master' into 80-update-conda-requirements-file	last mon
🗅 env-	metal-arm64.yaml	Include MacOS env builder (testing).	2 weeks ag
🗋 env.y	/aml	Update the conda setup file and the requirements.txt	last mon
🗋 requ	irements.txt	Update the conda setup file and the requirements.txt	last mon
🗋 setu	р.ру	Implemented visulation tools and helper scripts	last mon
≣ REA	DME.md		4
Sci Dev re	DAC Quantom	tom project	

3. When your code is complete and updated in the branch then make a pull request

More details are found here: HOWTO\_CONTRIBUTE.md



# **PROXY APP (NOT REAL PHYSICS)**

• Only considering u and d contribution and has not physical equivalent





# **CLOSURE TEST**

- Use the Proxy App for closure tests and scaling
- We generate toy data (1M events) using fix parameters in the theory module
- We train an ensemble of 15 GAN workflows
- Results are from ideal setup



Parameter Residuals from Proxy App



son Lab

Parton Densities from Proxy App

# **ADDITIONAL PROXY APP RESULTS**

- Including some detector effects:
  - Detector with 5% / 12% resolution effect on sigma1 / sigma2
  - Added correlations between sigma1 / sigma2





# SCALING USING COMMON TOOLS

•Initial scaling studies where based on Horovod and Fairscale

•Both didn't scale very well

•The stochastic nature of the workflow doesn't work well with allreduce techniques







## **EXPLORING SCALING USING ENSEMBLES**

- Simplest approach is to have completely decoupled learning and take aggregate predictions at fixed time intervals
- This allows us to study the model stability and accelerate convergence





## **IMPACT FROM SAMPLE STATISTICS**

- Statistical impact in results: 1k samples (left) and 10k samples (right)
- Clear bias that need to be understood (difference between parent and sampled distributions)





# **BEYOND THE PROXY APP**

- The proxy app allows us to study the framework and scaling
- We need to also study potential real use case (inclusive DIS, etc.)





# LEARNING FROM GPD IMAGES

- The initial goal is to ensure that we can generate images that look like the designed "data" images
- This is the MNIST example for NP  $\ensuremath{\textcircled{}}$







# LEARNING FROM GPD IMAGES

- We now add the evolution code in the mix
- The evolution code is slow and require a lot of memory
- In fact, we cannot run on GPU due to memory requirement (on A100)





- Gradient-based approach require all components to be differentiable for backpropagation
  - This is a challenge for traditional sampling techniques
- Gradient-based optimizers store large amount of information
  - This is a problem when backpropagating through evolution code, etc.
- Elements of the current workflow have stochastic execution times which doesn't scale using allreduce methods



# Thank you