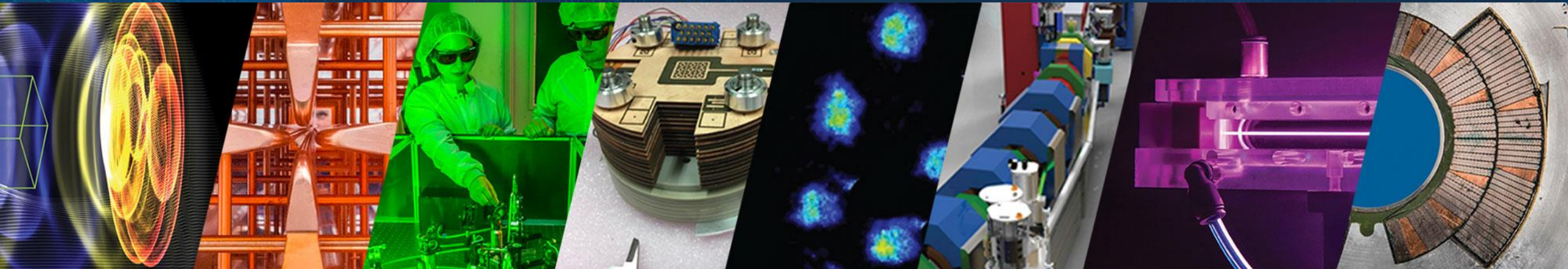


AI/ML for Accelerator Design/Control at LBNL

Remi Lehe
Lawrence Berkeley National Laboratory



June 13, 2024



ACCELERATOR TECHNOLOGY &
APPLIED PHYSICS DIVISION



U.S. DEPARTMENT OF
ENERGY

Office of
Science

AI/ML for accelerators is used across several groups in the ATAP division



ACCELERATOR TECHNOLOGY &
APPLIED PHYSICS DIVISION

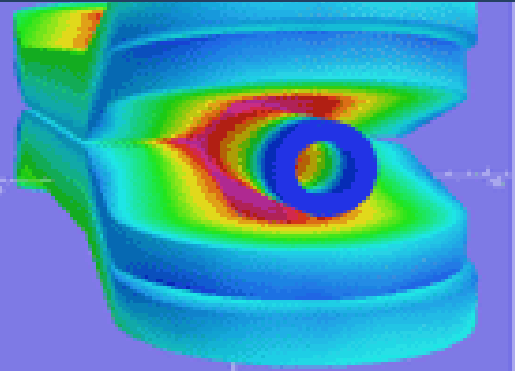
Advanced Light Source Accelerator Physics



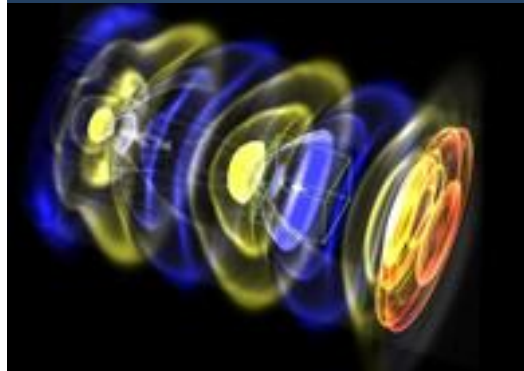
Berkeley Lab Laser Accelerator Center



Superconducting Magnet Program



Advanced Modeling Program



Fusion Science and Ion Beam Technology Program



Berkeley Accelerator Controls & Instrumentation (BACI)



Outline

- AI/ML for design optimization of accelerators
- AI/ML for accelerator operation
- Adapting simulation tools for integration with AI/ML

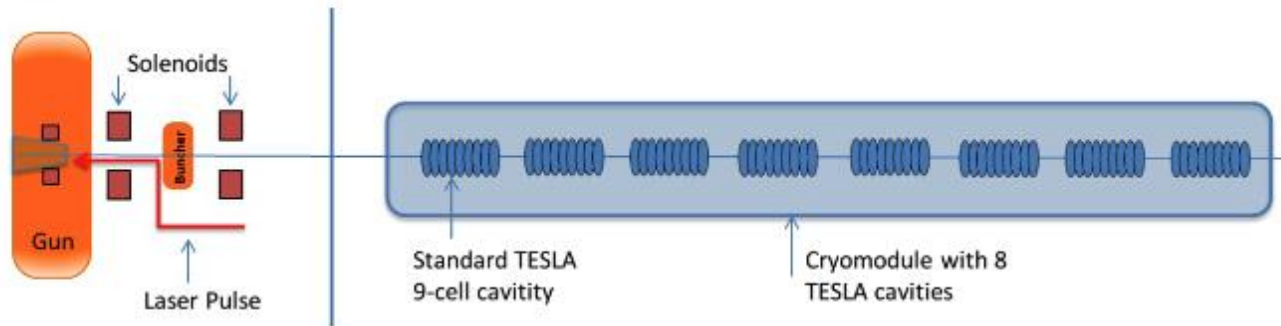
We use ML to accelerate simulation-based design of particle accelerators.

- Designing accelerators often involve **simultaneously tuning** many parameters (focusing, accelerating cavities, etc.) to reach the design with **optimal performance**,

i.e. maximize $f(\mathbf{x})$

\mathbf{x} : vector of accelerator parameters

f : function to maximize (“objective function”)



$$f = \epsilon_{\perp}$$

$$\mathbf{x} = \begin{pmatrix} B_{solenoid} \\ E_{cavity} \end{pmatrix}$$

- Each combination of parameters needs to be evaluated with **expensive simulations**.

We use ML to accelerate simulation-based design of particle accelerators.

“Conventional” optimization algorithms:

e.g.

- Gradient descent
- Genetic algorithms
- Nelder-Mead algorithm (a.k.a. simplex)
- ...

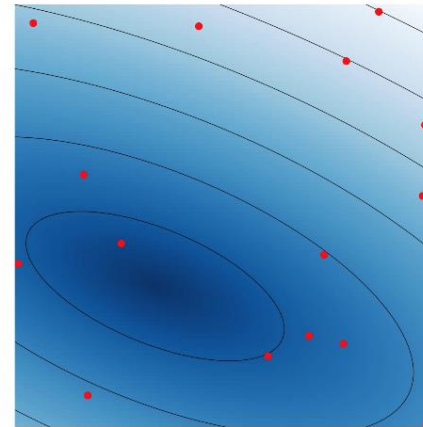
The next evaluations are based on simple rules that the depend on the **last few evaluations**.

Typically require **many** evaluations of f .

Optimization algorithms based on machine learning:

Progressively learn a **global model** of the objective function $f(x)$ over the parameter space.

Use this model to **only evaluate** the most promising x .



Model of f

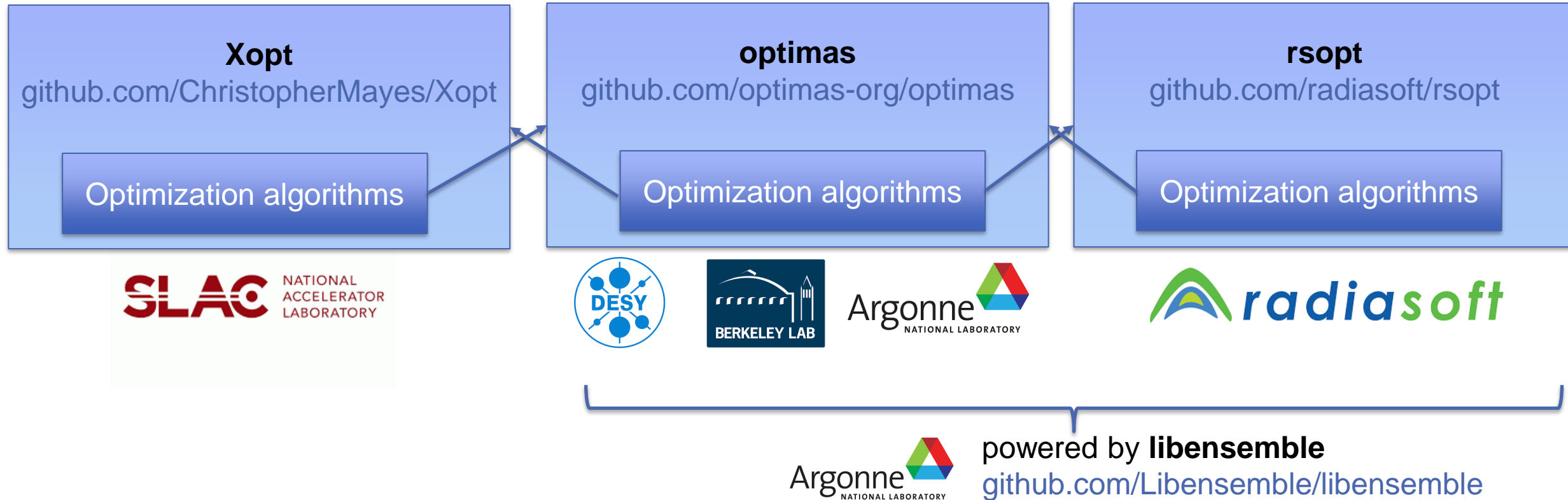
Typically require **fewer** evaluations of f .

A. Ferran Pousa et al., Bayesian optimization of laser-plasma accelerators assisted by reduced physical models, PRAB (2023)

Y. Lu et al., Demonstration of machine learning-enhanced multi-objective optimization of ultrahigh-brightness lattices for 4th-generation synchrotron light sources, NIMA (2023)

We are fostering interoperability across open-source optimization software.

- Several **open-source optimization frameworks** are being used in the accelerator community (each with their respective strengths)



- Ongoing efforts by the developers to **standardize optimizers** and **foster interoperability**.

Outline

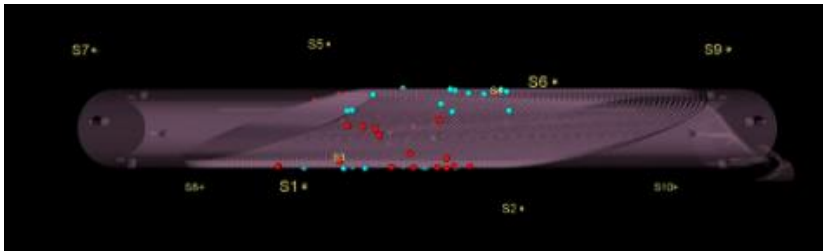
- AI/ML for design optimization of accelerators
- AI/ML for accelerator operation
- Adapting simulation tools for integration with AI/ML

AI/ML for accelerator operations covers a broad range of topics.

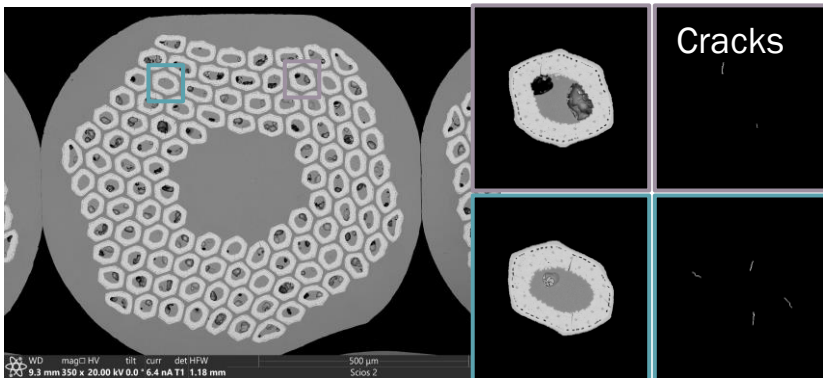
Superconducting magnets verification/protection:

- Detection/classification of **quench precursors** from acoustic emission.

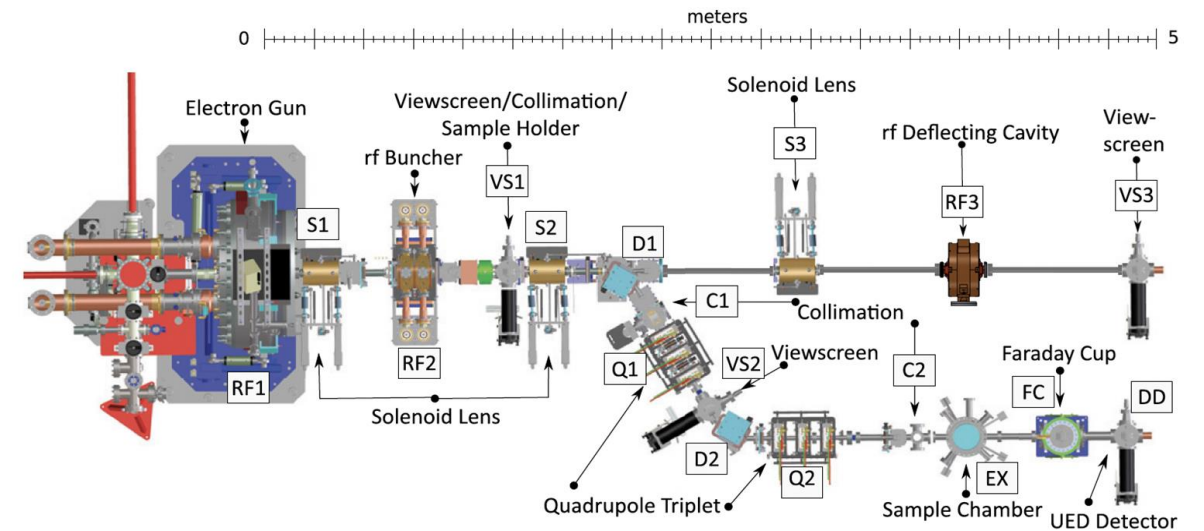
M. Marchevsky, [arXiv:2203.08871](https://arxiv.org/abs/2203.08871) (2022)



- Check for **cracks/defects** in superconducting cables with automated image analysis



Control and autotuning at compact LBNL accelerators (e.g., HiRES & NDCX-II):



A. Scheinker, et. al., [Scientific Reports](https://doi.org/10.1038/s41598-021-00000-0) (2021)

A. Scheinker, et. al, [Phys. Rev. E](https://doi.org/10.1103/PhysRevE.107.014801) 107 (2023)

F. Cropp, et. al., [Phys. Rev. Accel. Beams](https://doi.org/10.1103/PhysRevAccelBeams.26.014801) 26 (2023)

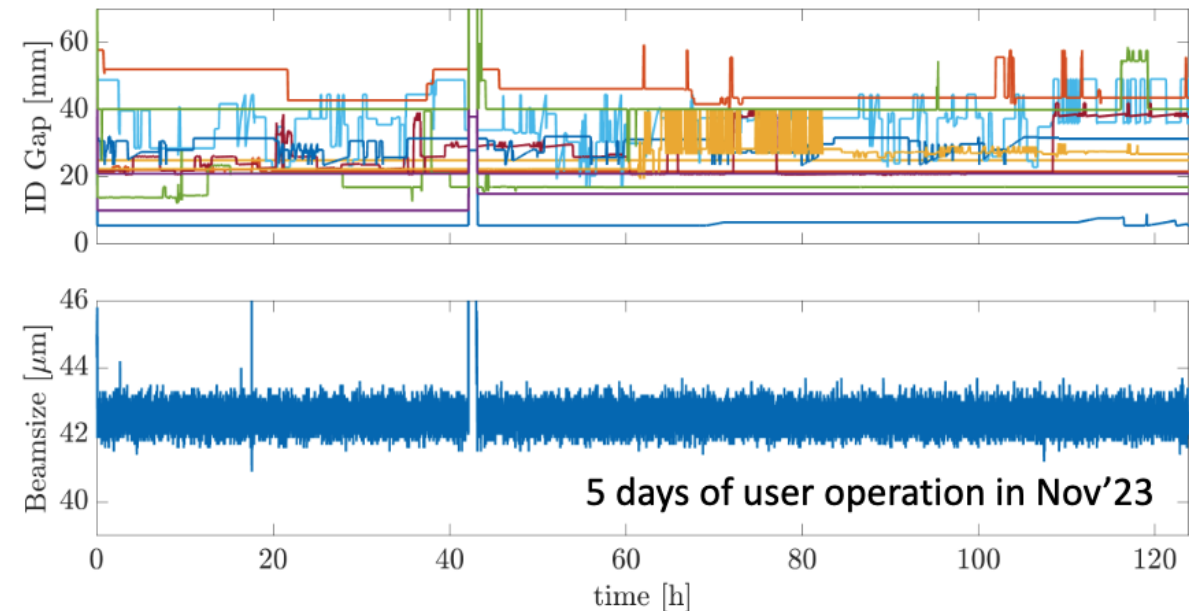
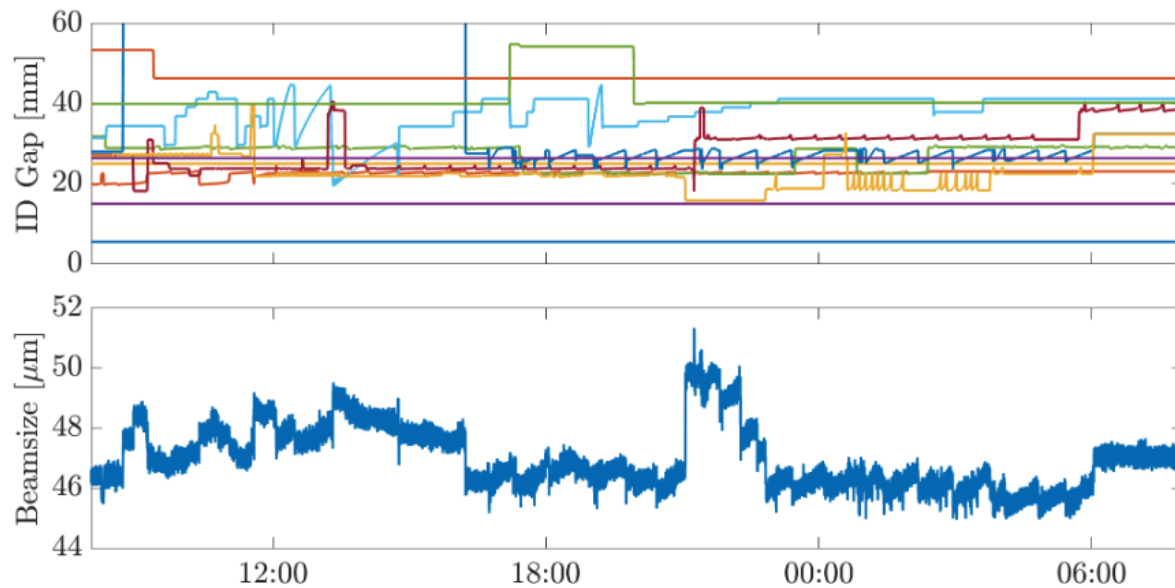
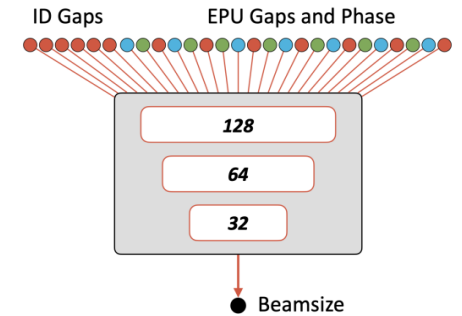
ML is used in operation at the Advanced Light Source, to stabilize beam size.

- The ALS is a storage ring light source.
- The parameters of **insertion devices** (e.g. undulators) are frequently changing, due to changing modes of operation.
- These changes affect the **beam size**.



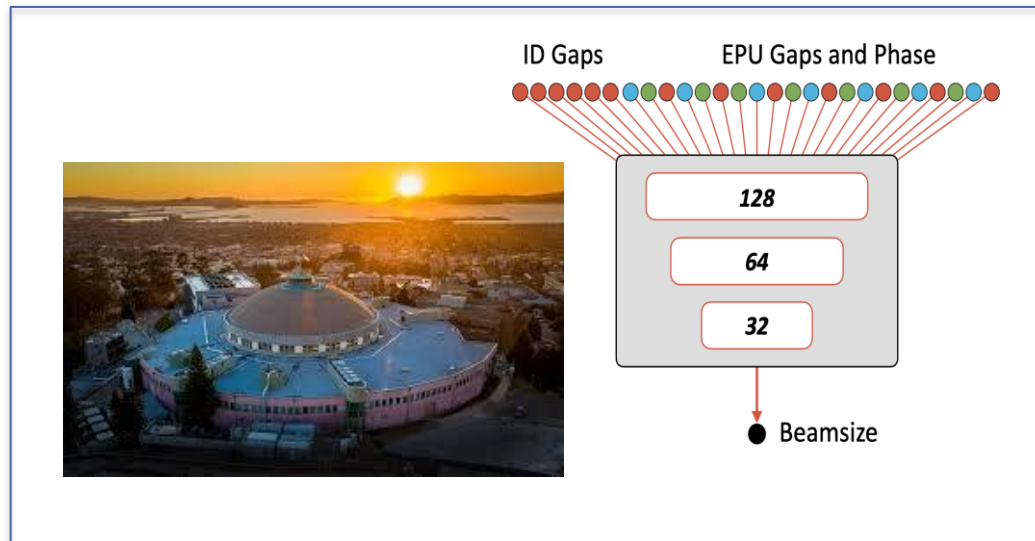
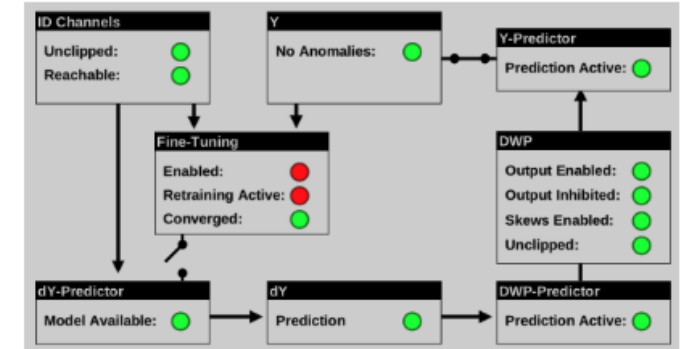
- The ALS now uses a neural network to **predict** changes in beam size and **correct** them.
Leeman et al., PRL 123, 194801 (2019)
Hellert et al., accepted in PRAB (2024)

- Used in routine operation

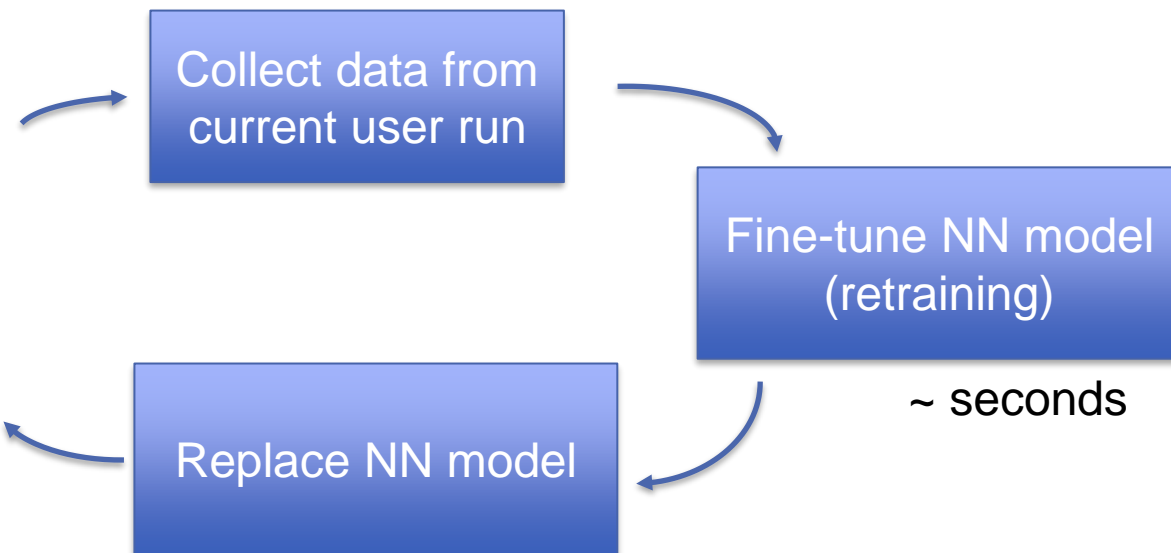


ML is used in operation at the Advanced Light Source, to stabilize beam size.

- The machine configuration **drifts** over time. To always **remain accurate**, the neural network needs to be periodically **retrained**.
- The ALS team developed an extensive framework to automatically **monitor, retrain, archive and deploy** the neural network.

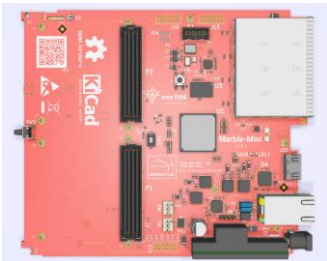


Online model



The BACI group is developing tools to run AI/ML on FPGA.

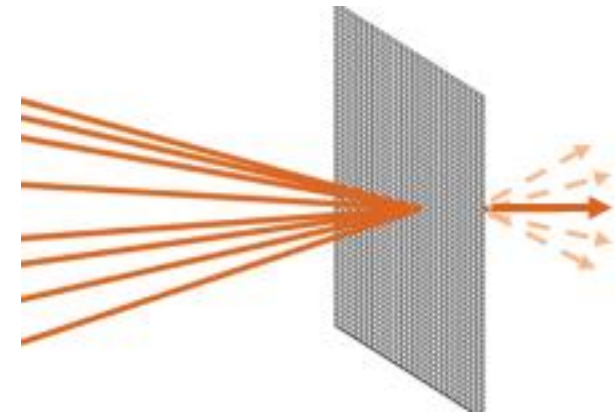
- Some accelerator control applications require **very low latency**, necessitating FPGA hardware.
- The **Marble-Mini** is an FPGA-based board that is widely used in low-level RF control for accelerators.



github.com/BerkeleyLab/Marble-Mini

- The BACI group at LBNL is developing a framework to run neural networks on Marble-Mini.
L.Doolittle, Q. Du and D. Wang, Software Disclosure 2024-008: Generic Multi-Layer Perceptron Inference Accelerator on FPGA (vneuron) v1.0
+ LBNL LDRD led by Dan Wang

- **Example use case:** neural network that controls the phases of 9 coherently combined laser beams

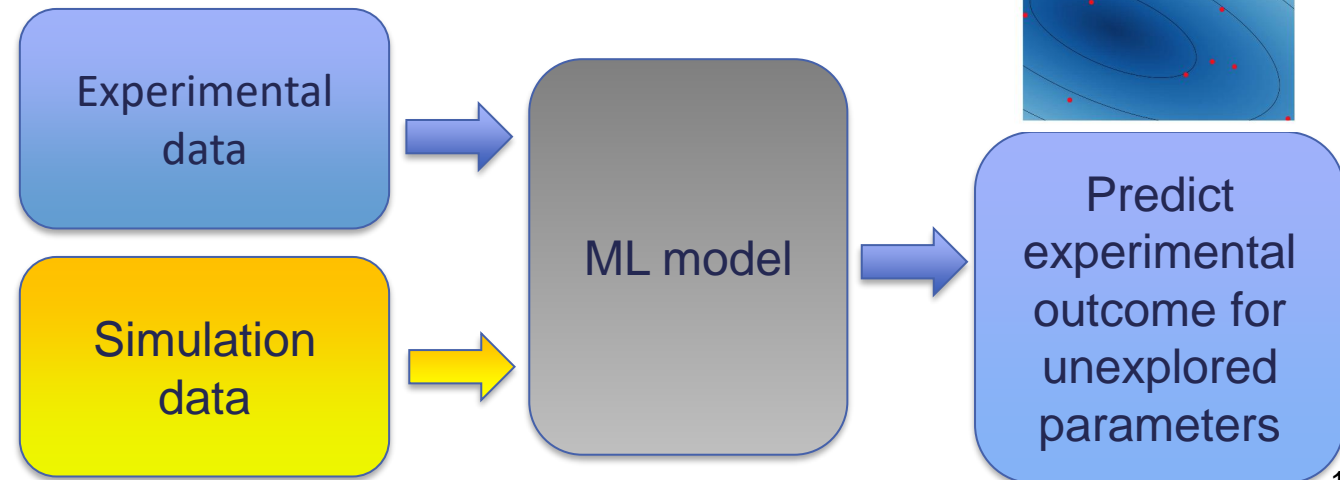
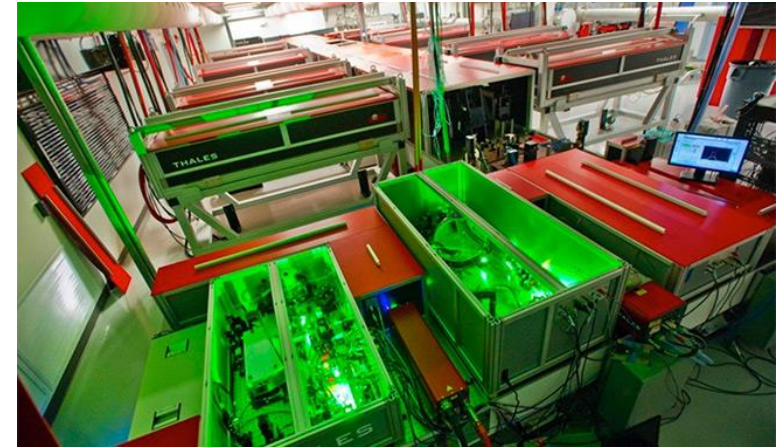


- 3-layer neural network with ~1600 weights
Inference time: 131 cycles (1046 ns)
Uses only 16% of the chip resources

We are working towards digital twins to guide accelerator tuning in real time.

- **Real-time accelerator tuning** requires to simultaneously adjust multiple parameters (e.g. focusing elements, accelerator cavities) to achieve **optimal operation**.
- Corresponds to optimization in high-dimensional space, **time-consuming** if the number of parameters is large.
- Leveraging information from **numerical simulations** can significantly speed-up real-time tuning.

e.g. Hanuka et al., Physics model-informed Gaussian process for online optimization of particle accelerators (2021)



We are working towards digital twins to guide accelerator tuning in real time.

Framework for deployment ("Superfacility")



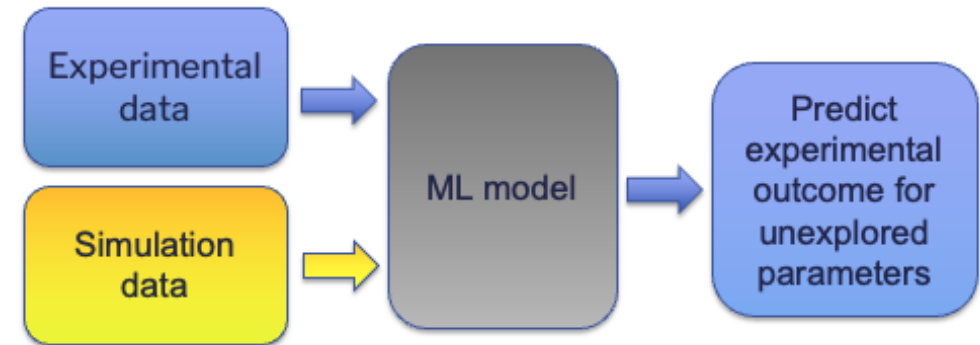
Experimental facility:

- Real-time tuning
- Data collection
- Ongoing collaboration with SLAC to develop corresponding software framework. (funding: LDRD + NESAP program)
- Will leverage existing software for parts of this workflow e.g., github.com/slaclab/lume-services

Supercomputer:

- Updates ML model
- Launch new simulations

ML challenges



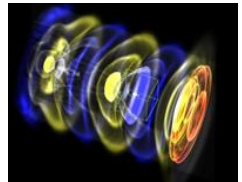
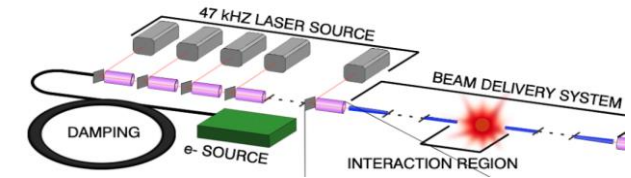
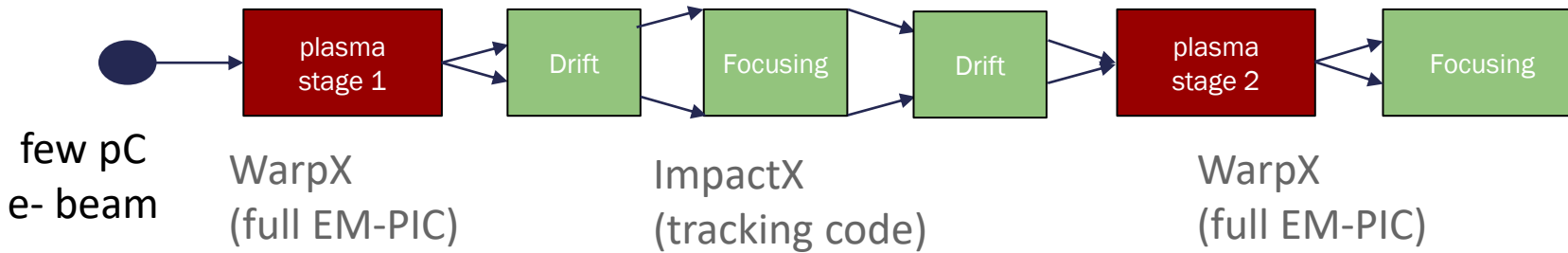
- Experimental and simulation data **oftentimes don't match exactly**. The ML model needs to **handle them differently**.
- Several possible techniques:
 - Multi-fidelity Gaussian process *A. Ferran Pousa, PRAB (2023)*
 - Calibration of neural networks *T. Boltz, arxiv 2403.0322 (2024)*

Outline

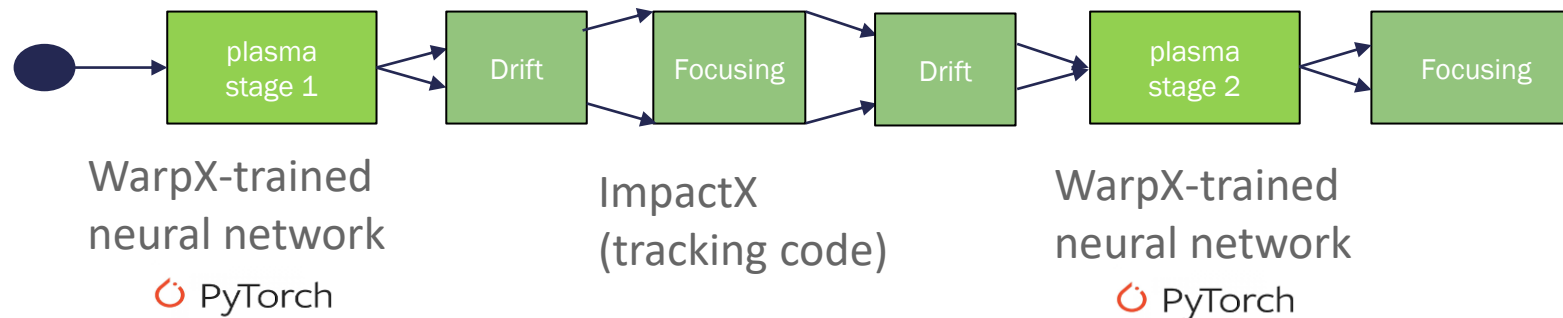
- AI/ML for design optimization of accelerators
- AI/ML for accelerator operation
- Adapting simulation tools for integration with AI/ML

ML surrogate models can speed-up simulations.

- In a given lattice, some elements can be more **computationally-expensive** to model. Extreme example: laser-plasma acceleration stages



- Under certain conditions (here: negligible collective effects, specific range of parameters), computationally-expensive elements can be replaced by **ML models**, trained over **past simulations**.



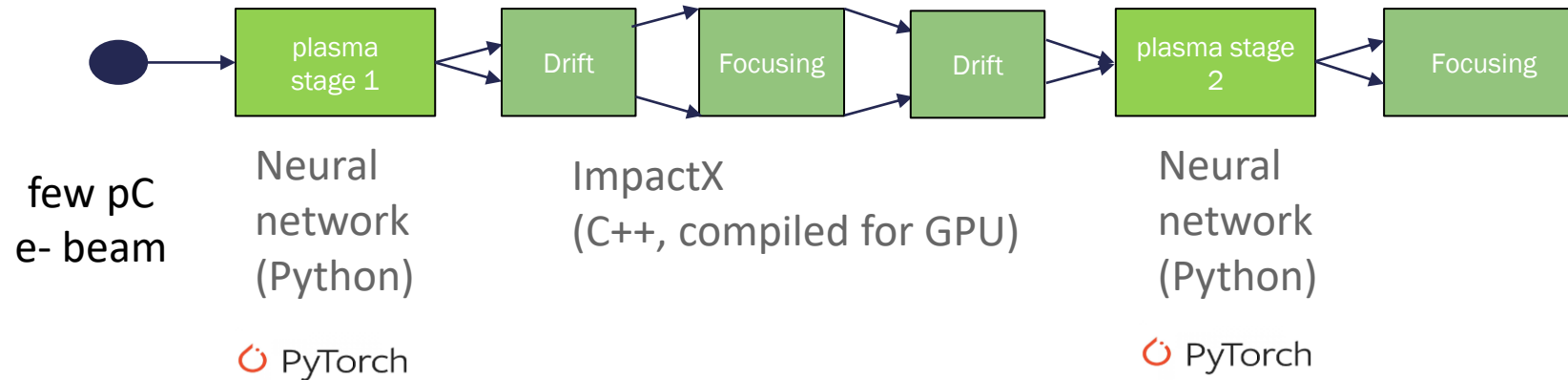
Simulation time: (with full geometry/physics)

hrs
on several GPUs

<sec
on 1 GPU

We adapt our accelerator simulation codes for seamless integration with ML surrogates.

- This workflow requires to pass **beam data** across elements (position/momentum for each particle)



- Our C++ simulation codes expose beam data (in GPU memory) through a Python interface.
→ Seamless, **GPU-Accelerated Coupling** of **accelerator simulation & ML Frameworks**.



github.com/AMReX-Codes/pyamrex
github.com/ECP-WarpX/WarpX
github.com/ECP-WarpX/ImpactX

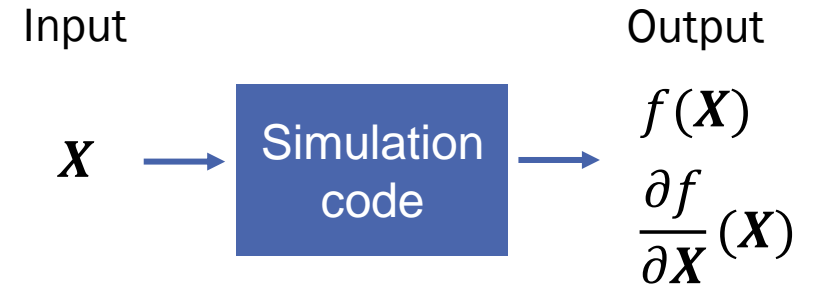
(based on pybind11
and python standards
for GPU arrays)

Even tighter ML integration can be achieved with differentiable simulations

Regular simulation code



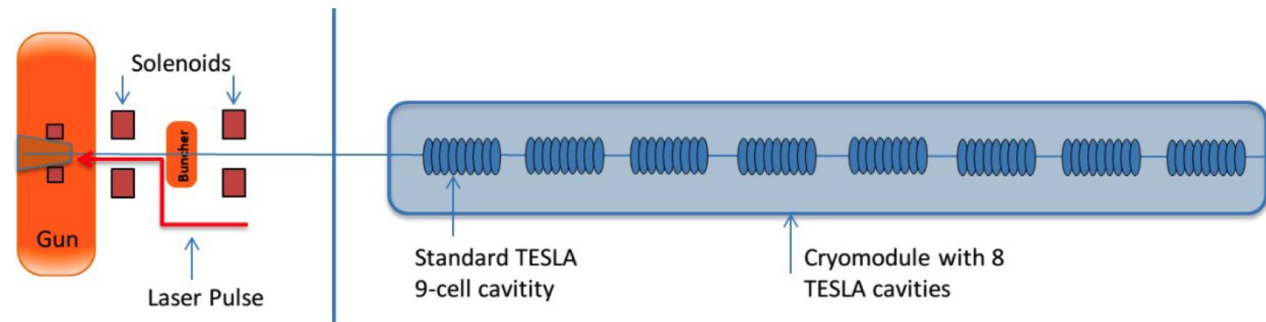
Differentiable simulation code



Example:

Input:
accelerator parameters

$$\mathbf{X} = \begin{pmatrix} B_{solenoid} \\ \varphi_{RF\ cavity} \\ E_{RF\ cavity} \\ \sigma_{beam,i} \end{pmatrix}$$



$$f = \epsilon_{\perp}$$

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial \epsilon_{\perp}}{\partial B_{solenoid}} \\ \frac{\partial \epsilon_{\perp}}{\partial \varphi_{RF\ cavity}} \\ \frac{\partial \epsilon_{\perp}}{\partial E_{RF\ cavity}} \\ \frac{\partial \epsilon_{\perp}}{\partial \sigma_{beam,i}} \end{pmatrix}$$

Differentiable codes have several advantages.

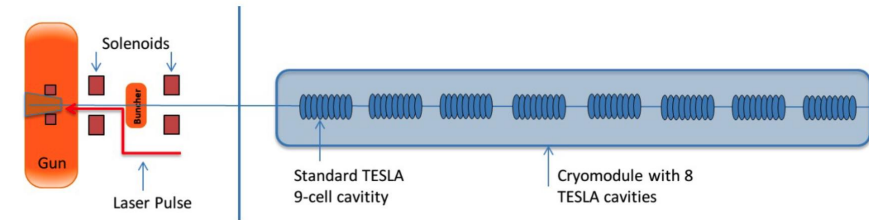
- **Sensitivity studies**

$\frac{\partial f}{\partial X}$ quantifies how sensitive the output is to the input.

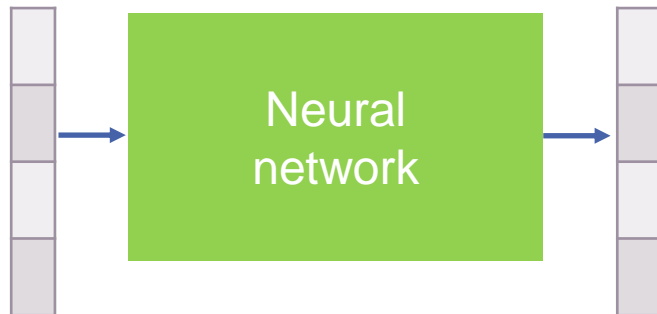
- **Optimization in high-dimensional space** (e.g. of accelerator designs)

$\frac{\partial f}{\partial X}$ can be used in **gradient-based** optimizers, which often converge faster

- **Allows training of a neural network that is combined with a differentiable code**

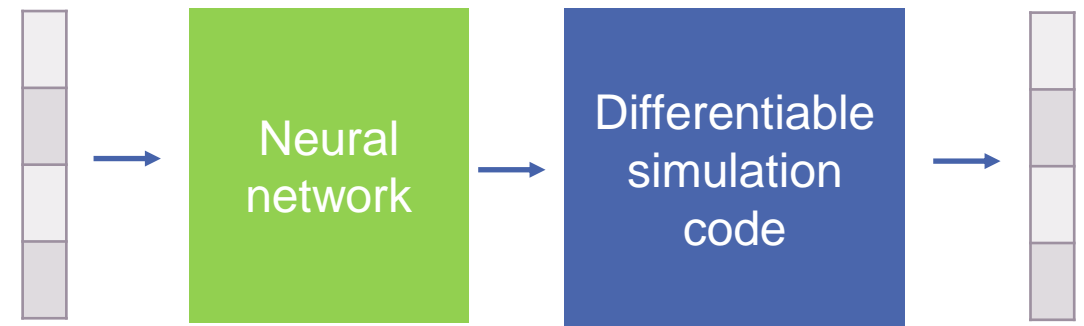


Traditional training of neural network



Input/output pairs, from a data set

Training of neural network combined with a code



Example: *R. Roussel et al., Phase Space Reconstruction from Accelerator Beam Measurements Using Neural Networks and Differentiable Simulations, PRL (2023)*

We are exploring frameworks for differentiable codes.

- Several **algorithms** are available to make a code differentiable. e.g.
J. Qiang, Differentiable self-consistent space-charge simulation for accelerator design, PRAB (2023)
- Several efforts to build differentiable **accelerator simulation codes**, based on **auto-differentiation frameworks**.

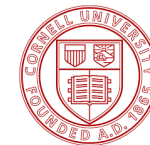
◦ pytorch  PyTorch


Cheetah: accelerator code based on pytorch
github.com/desy-ml/cheetah



◦ Julia 

Bmad-Julia: proposal to implement Bmad algorithms in Julia
github.com/bmad-sim



◦ Enzyme AD 

Takes existing C++ code and makes it auto-differentiable at compile time.
Could be leveraged to make BLAST codes (ImpactX, WarpX, ...) differentiable.



Conclusion

- AI/ML is used in **many accelerator-relevant applications**, in the ATAP division.
- In the process, we are building **software frameworks** for robust deployment of AI/ML.
- Many of these frameworks are **open-source** and can be leveraged at **other accelerator facilities**. If interested, feel free to reach out!

USPAS course on AI/ML for accelerators



Course on **Optimization and Machine Learning for Accelerators** at the **U.S. Particle Accelerator School**, since 2021

**Some of the
Instructors:**



Remi Lehe
(LBNL)



Auralee Edelen
(SLAC)



Ryan Roussel
(SLAC)

Next USPAS session:

Held in: Knoxville, Tennessee

Dates: Jan 27 - Feb 7, 2025

May include optimization and ML (not confirmed yet)

The screenshot shows a Jupyter Notebook window titled 'lab_05.ipynb'. The left sidebar displays a file explorer with a directory structure: / 2021_optimization_and_ml / labs /, containing subdirectories lab_01 through lab_05. The main content area shows a plot of a function y versus x . The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from -0.5 to 2.5. The plot features a blue curve with a shaded uncertainty region and several orange data points. Below the plot, the text 'Accelerator example' is followed by the sentence: 'Now we will use BO to optimize the strengths of quadrupoles in a triplet configuration.' Below this text is a schematic diagram of a triplet configuration of quadrupoles. It shows three black diamond-shaped quadrupoles labeled K_1 , K_2 , and K_3 arranged in a row. The quadrupoles are separated by 1 m, and the total length of the triplet is 0.1 m. The diagram is bounded by Σ_i on the left and Σ_f on the right.

Thank you