



# Online Autonomous Tuning of the FRIB Accelerator Using Machine Learning

Peter N. Ostroumov, [ostroumov@frib.msu.edu](mailto:ostroumov@frib.msu.edu)

Co-PIs: Kilean Hwang, Dean Lee, Alexander Scheinker

Contributors: Kei Fukushima, Tomofumi Maruta, Alexander Plastun, Jinyu Wan, Tong Zhang

2024 NP AI-ML PI Exchange Meeting

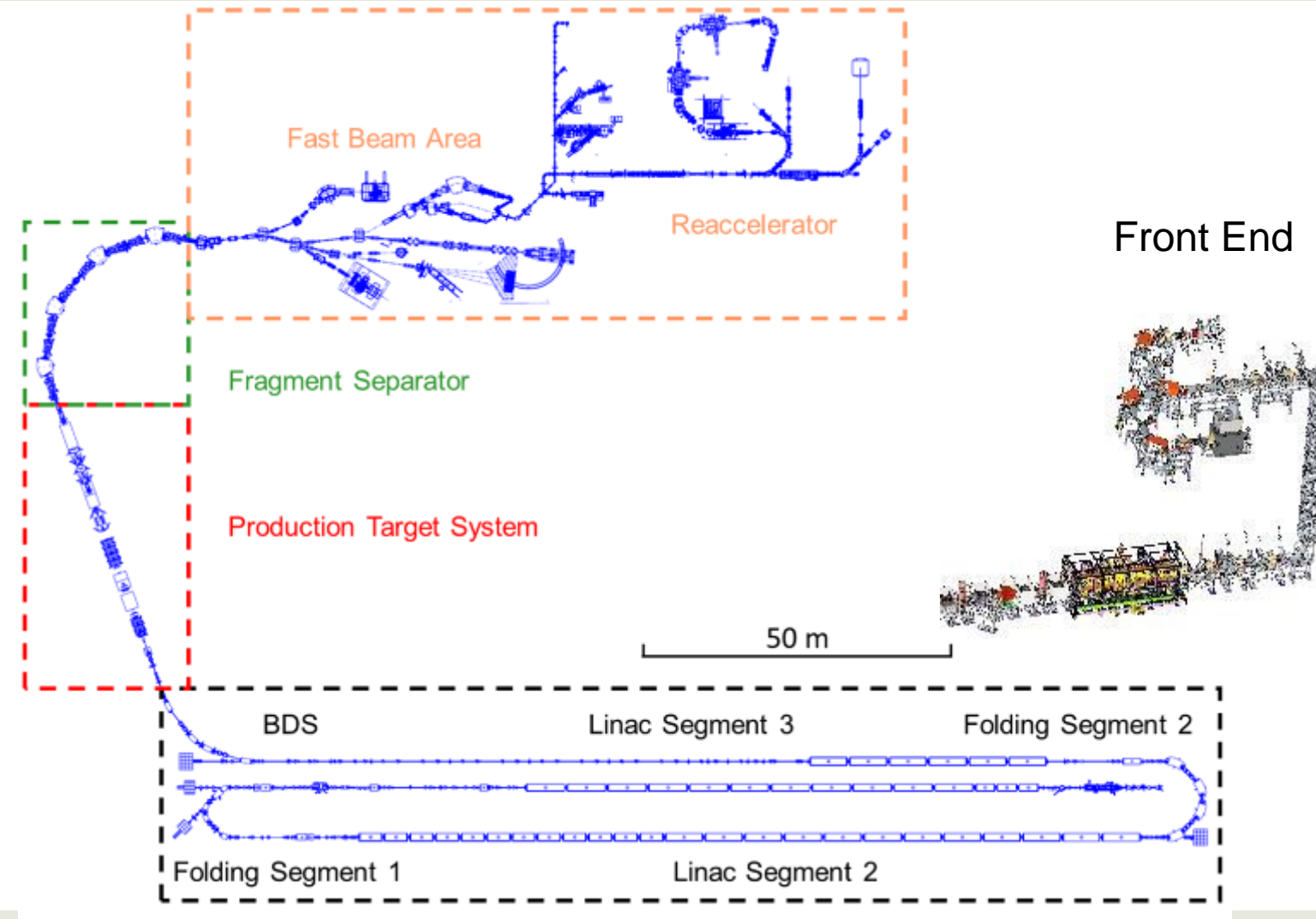


# Content

- FRIB Accelerators
- Accelerator tuning approach
- Examples of ML tasks
- Bayesian optimization (BO)
- Customized BO
- Virtual 4D Phase Space Diagnostics
- Generic BO application for accelerator tuning
- Fast NN model trained on physics equations for online application
- Virtual diagnostics for bunch length measurements
- Virtual diagnostics for beam quadrupolar moments and calculation of Courant Snyder parameters



# FRIB Layout

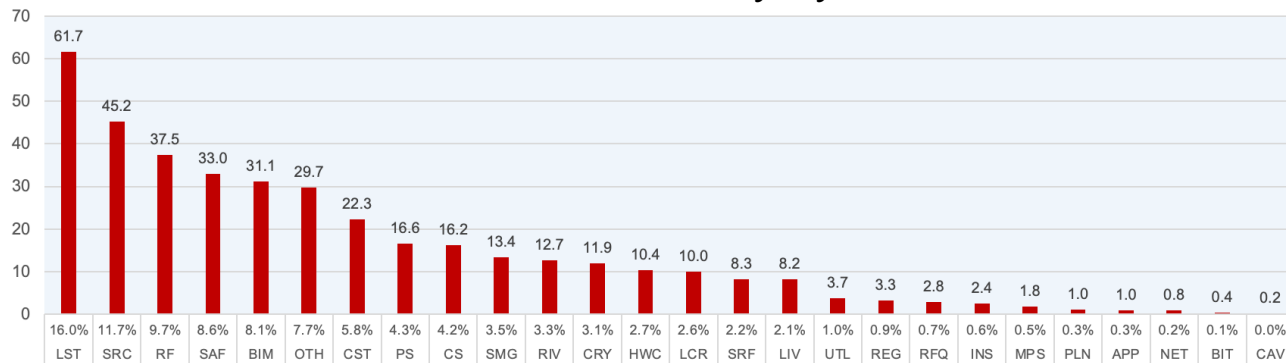


- CW Linear accelerator
  - 200 MeV/u - U, 320 MeV/u - O
  - RFQ,  $q/A=1/7$
  - 324 SC cavities in 46 cryomodules
- 80-meter-long space for upgrade to 400 MeV/u of uranium
- Fragmentation target/Beam dump
- 2-stage fragment separator
- Fast beamlines
- 6 MeV/u Re-accelerator
- In operation since May 2022
- 44 nuclear science experiments conducted since the commencement of user operation

# FRIB Operation

- Commencement of User Operation – May 2022
- The first two years of operation
  - Delivered 8500 beam hours at full energy for science and 4000 hours at lower energies, up to 40 MeV/u for Single Event Experiments (SEE)
    - » Beam availability in the first year was 92% and in the second year - 94%
  - 44 science experiments were carried out; the results were reported in multiple PRL papers.
  - More than 270 unstable isotopes produced with primary beams of  $^{18}\text{O}$ ,  $^{20}\text{Ne}$ ,  $^{28}\text{Si}$ ,  $^{36}\text{Ar}$ ,  $^{40}\text{Ar}$ ,  $^{48}\text{Ca}$ ,  $^{64}\text{Zn}$ ,  $^{70}\text{Zn}$ ,  $^{82}\text{Se}$ ,  $^{86}\text{Kr}$ ,  $^{124}\text{Xe}$ ,  $^{198}\text{Pt}$  and  $^{238}\text{U}$  accelerated up to 300 MeV/u
- Both primary and secondary beams are extremely stable during the experiment

Breakdown hours by systems



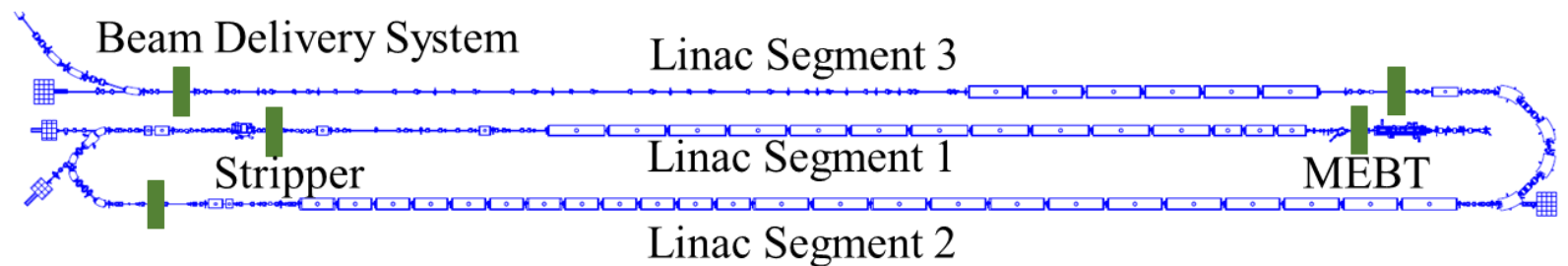
Main systems contributing to downtime

- Lithium stripper
- Ion Source
- RF power
- Safety systems
- Beam instrumentation



# Accelerator Tune Development

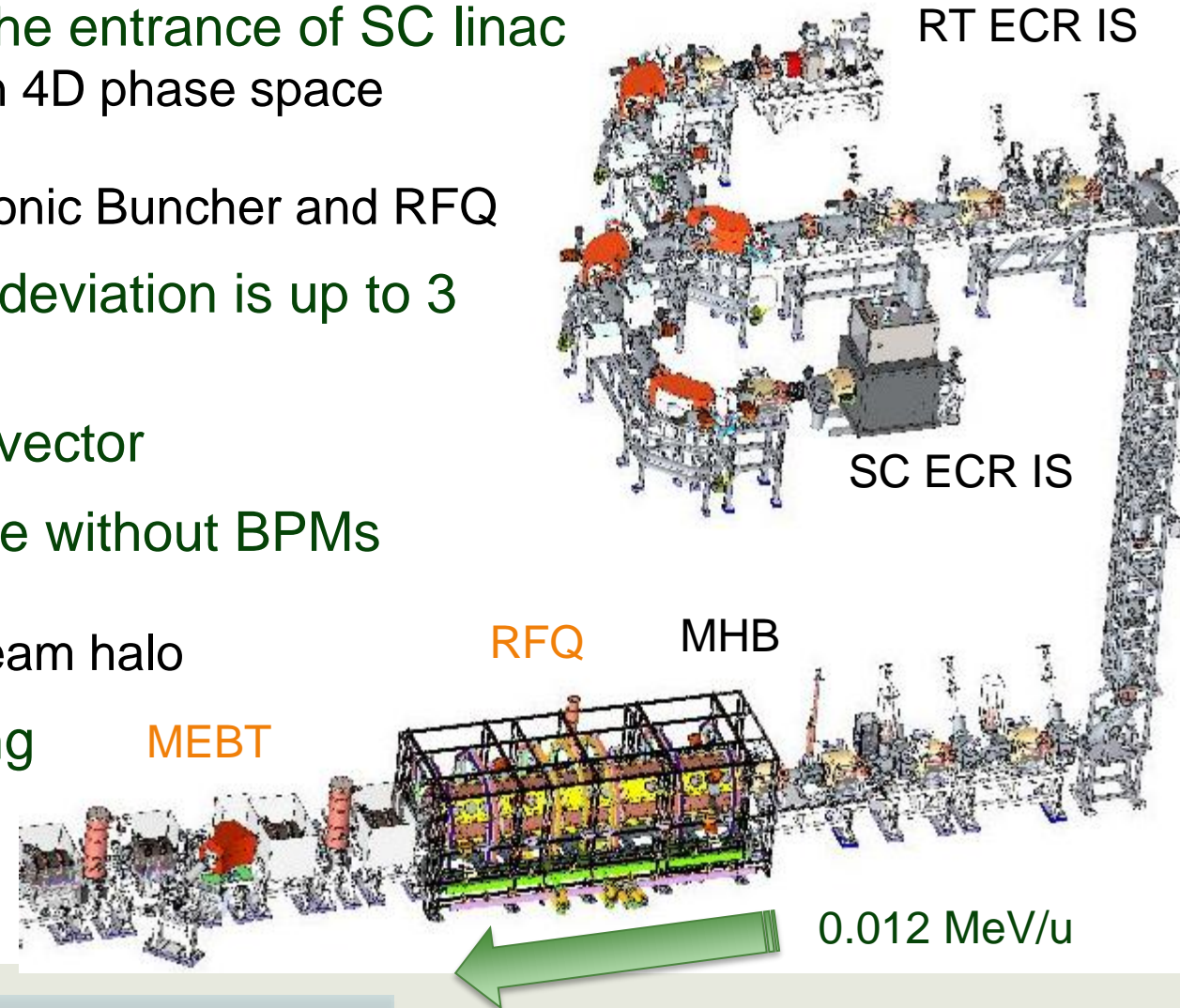
- New tunes are required for new beams, more charge states, and transition to higher power
- Virtual accelerator is based on second order envelope code and simulations of the longitudinal dynamics in realistic fields
  - Equivalent to the digital twin of the linac
  - Provides pre-setting of the entire accelerator
- Low-loss accelerator tune is achieved by transverse matching in several sections
  - The setting is verified by envelope mapping in both longitudinal and transverse planes
- Frequent ion species switching; typically occurs every week
  - Rapid beam tuning essential to provide more beam time for science
- Our Solution: Advanced beam tuning applications based on
  - Physics simulation
  - Machine Learning (ML)
  - Classical Optimization (e.g. Nelder-Mead)





# Examples of ML Task in Front End

- Problem in restoring previously tuned beam at the entrance of SC linac
  - ECR produces slightly different beam distribution in 4D phase space
  - The beam central trajectory deviates
  - Transverse-longitudinal coupling in the Multi-Harmonic Buncher and RFQ
- Result: at the MEBT, transverse beam centroid deviation is up to 3 mm; phase deviation is up to 4 degrees
- Task: match beam centroid to previously tuned vector
- Steer beam through **two apertures** in straight line without BPMs
  - Limited fast beam centroid diagnostics.
  - Maximize beam current through apertures  $\Rightarrow$  cut beam halo
- **RFQ** transmission and **MEBT** 6D beam centering
- LS1 longitudinal beam center restoration



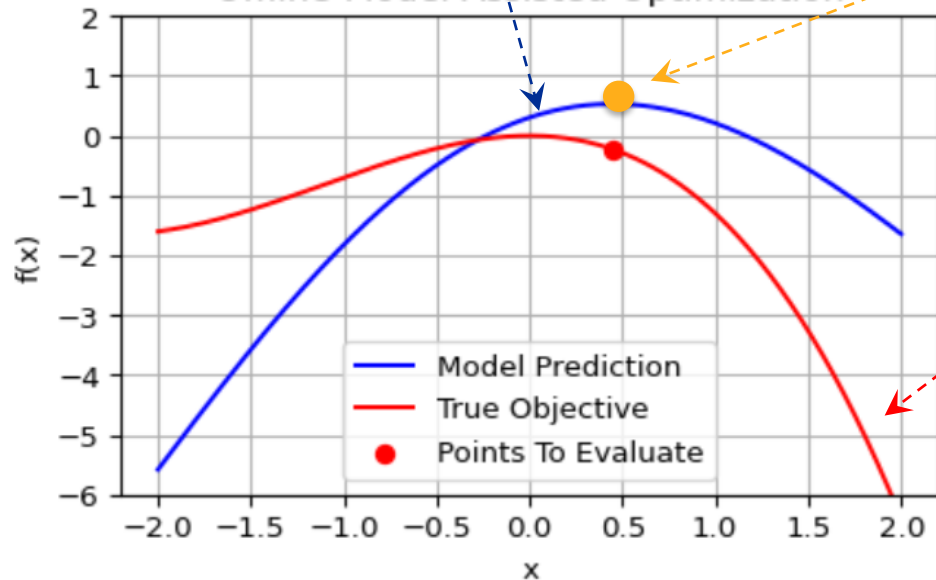
# Surrogate Model Assisted Optimization

- We consider problem of finding optimum set of control  $x^*$  that maximize an **objective**  $f(x)$

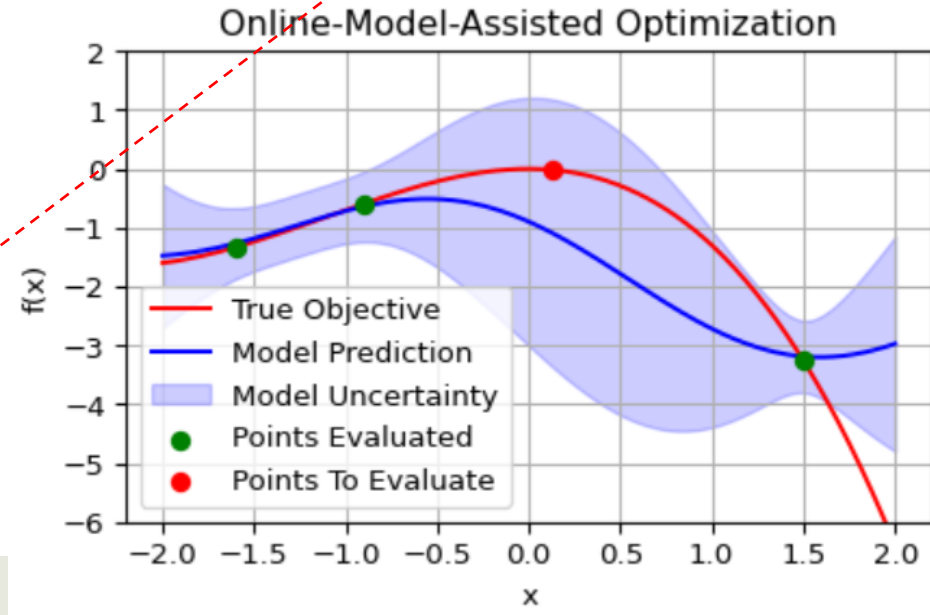
$$x^* = \operatorname{argmax} f(x)$$

- If we have a **model**  $m(x) \simeq f(x)$ , can find solution candidate  $x_0 = \operatorname{argmax} m(x)$

Optimization based on the physics or data-based deterministic surrogate model



Online Bayesian optimization

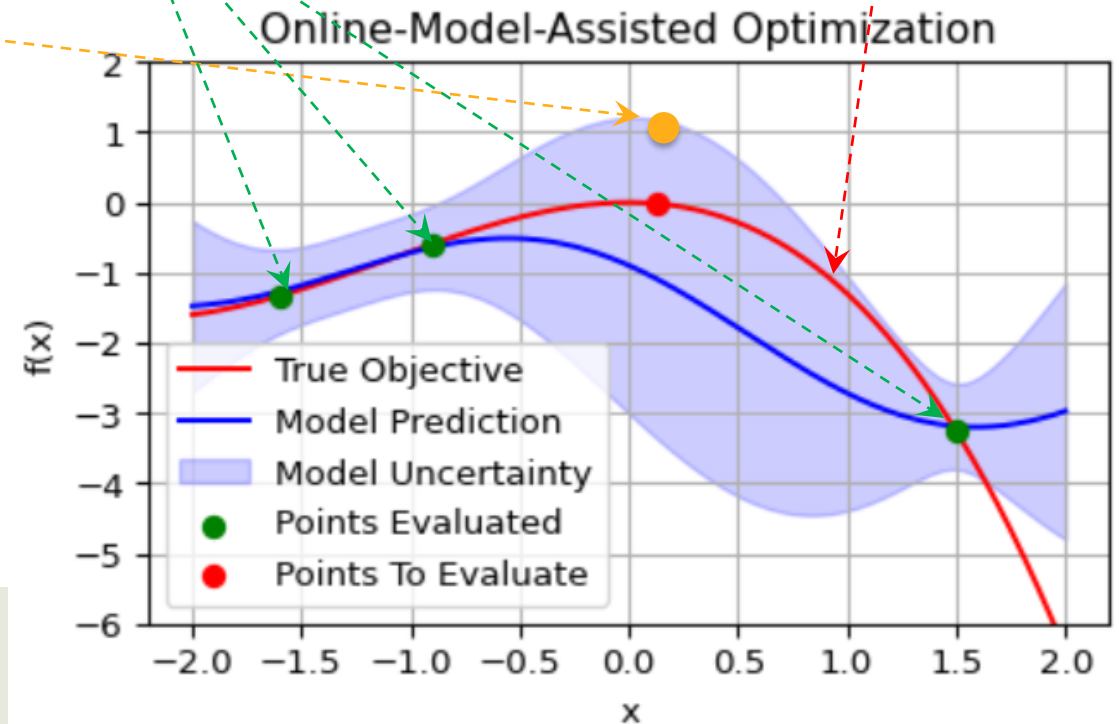
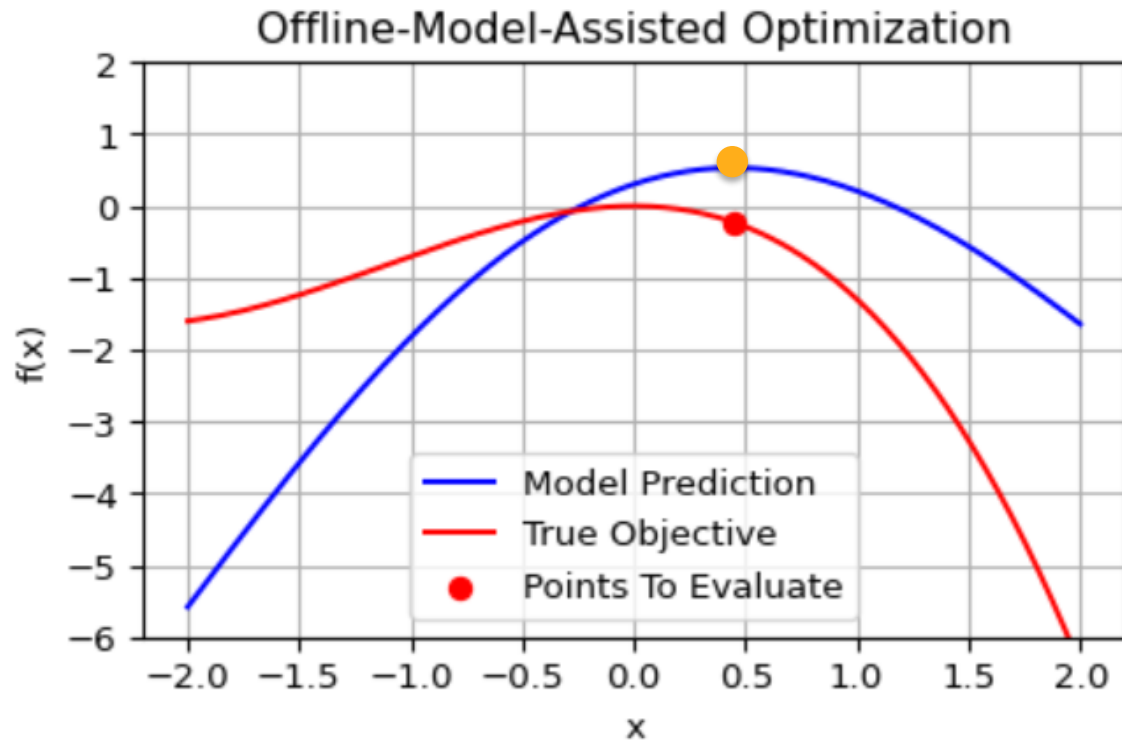


# Surrogate Model Assisted Optimization

- We consider problem of finding optimum set of control  $x^*$  that maximize an **objective**  $f(x)$

$$x^* = \operatorname{argmax} f(x)$$

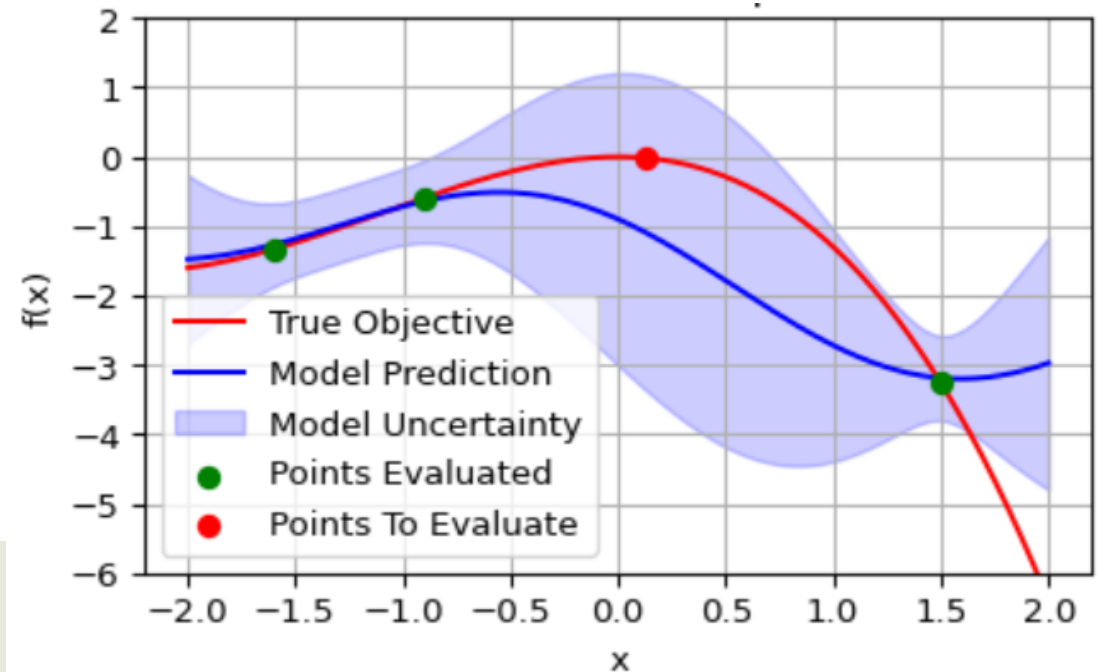
- If we have a **model**  $m(x) \simeq f(x)$ , can find solution candidate  $x_0 = \operatorname{argmax} m(x)$
- If we have a data-based probabilistic model  $\mathcal{P}(y|x, \mathcal{D})$ , we can query next candidate  $x_n$  that can account for model uncertainty.





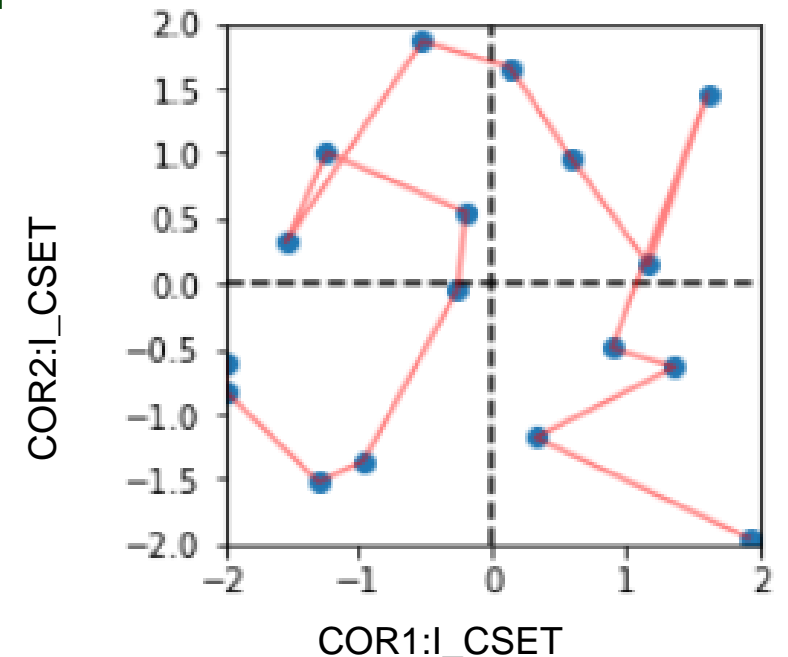
# Bayesian Optimization

- We consider the objectives that takes few to several seconds to evaluate
  - » Power supply ramping from old control  $x_i$  to new control  $x_{i+1}$  : generally a few sec but can take up to 30 sec
  - » Average out noisy measurements for a few sec.
    - Due the ramping time, a few sec of averaging cost will benefit for training surrogate model than noisier data.
  - Depending on the problem size, the overall time budget ranges from 2 to 20 minutes.
- Given the few to several seconds of evaluation time, BO is an appropriate choice
  - Due to its sample efficiency.
  - A few seconds of numerical cost ( for less than 10 dim, and less than 200 data points )
  - The scalability issue of computational complexity of BO won't be a limiting factor as the maximum number of function evaluation will not exceed a few hundreds due to tight time budget.



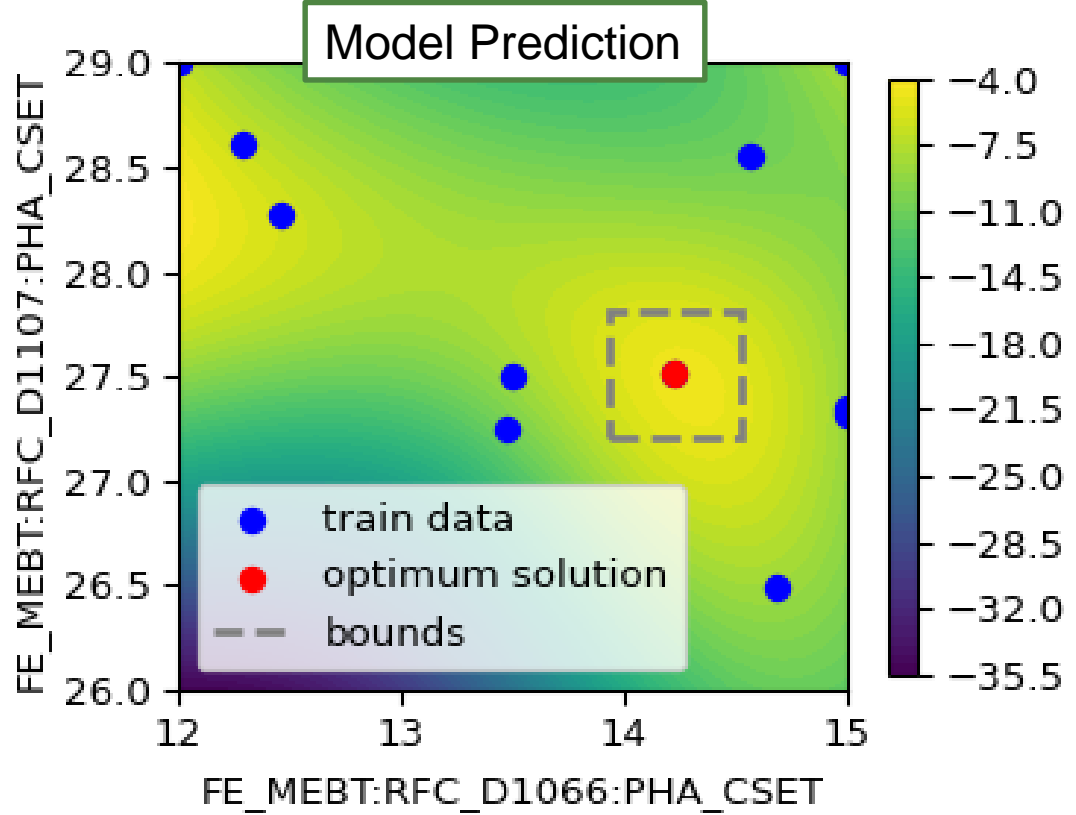
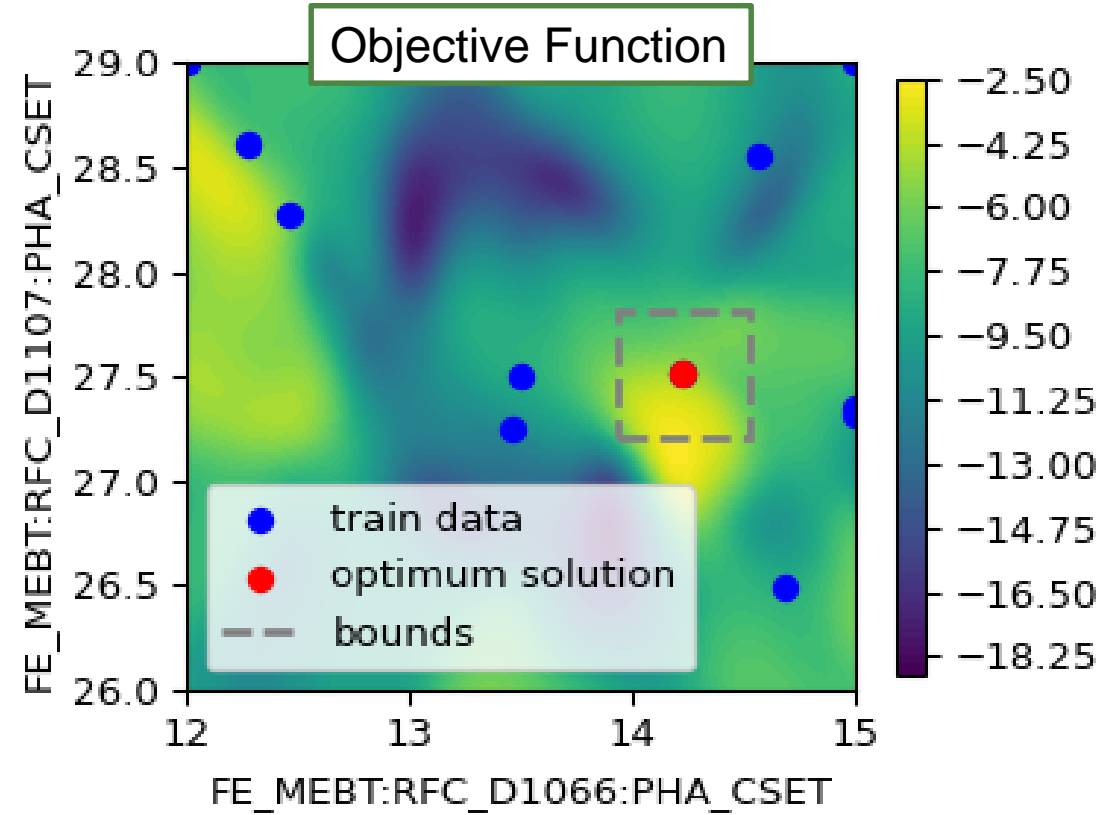
# Customized BO

- **Tightly avoid accelerator idle time, threaded computation**
  - Simultaneously queries candidate solutions while the beam measurements are ongoing
  - Candidate search terminates upon objective evaluation completion by machine
- **Resetting time awareness for max beam time utilization**
  - Particularly, power supply of corrector **polarity change**.
- **Order evaluation candidates (initialization points or multi-batch query) to minimize resetting cost**



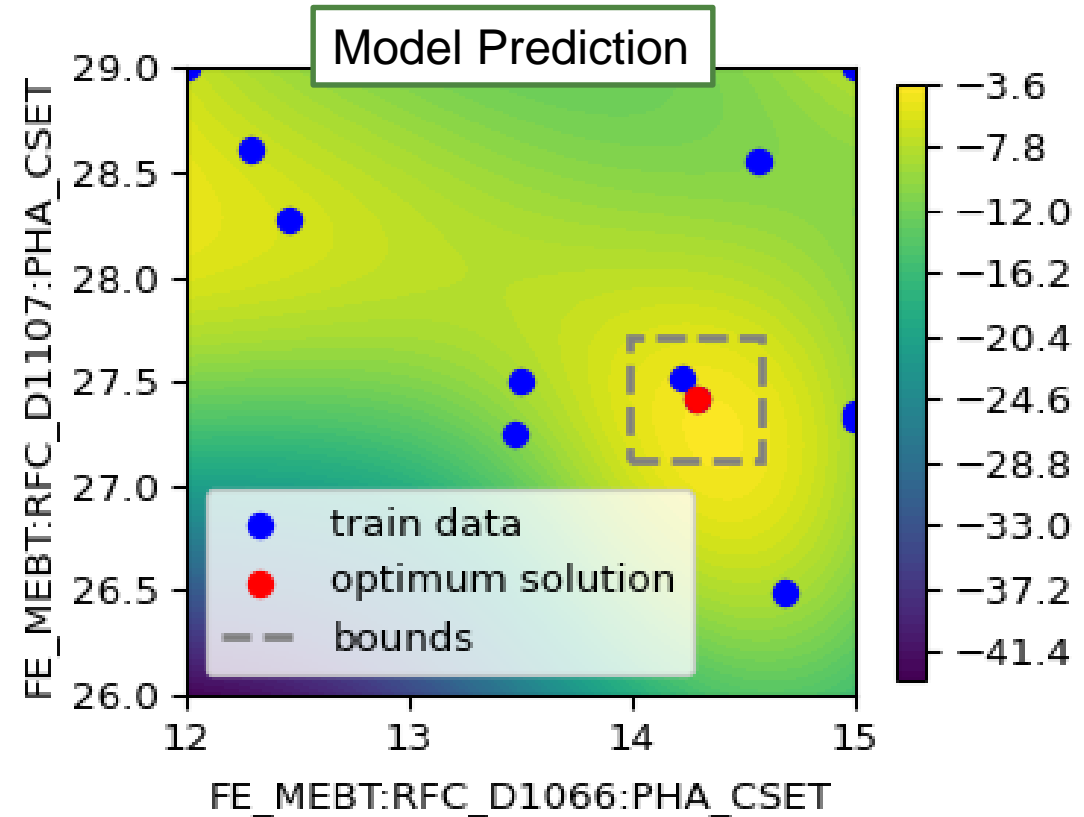
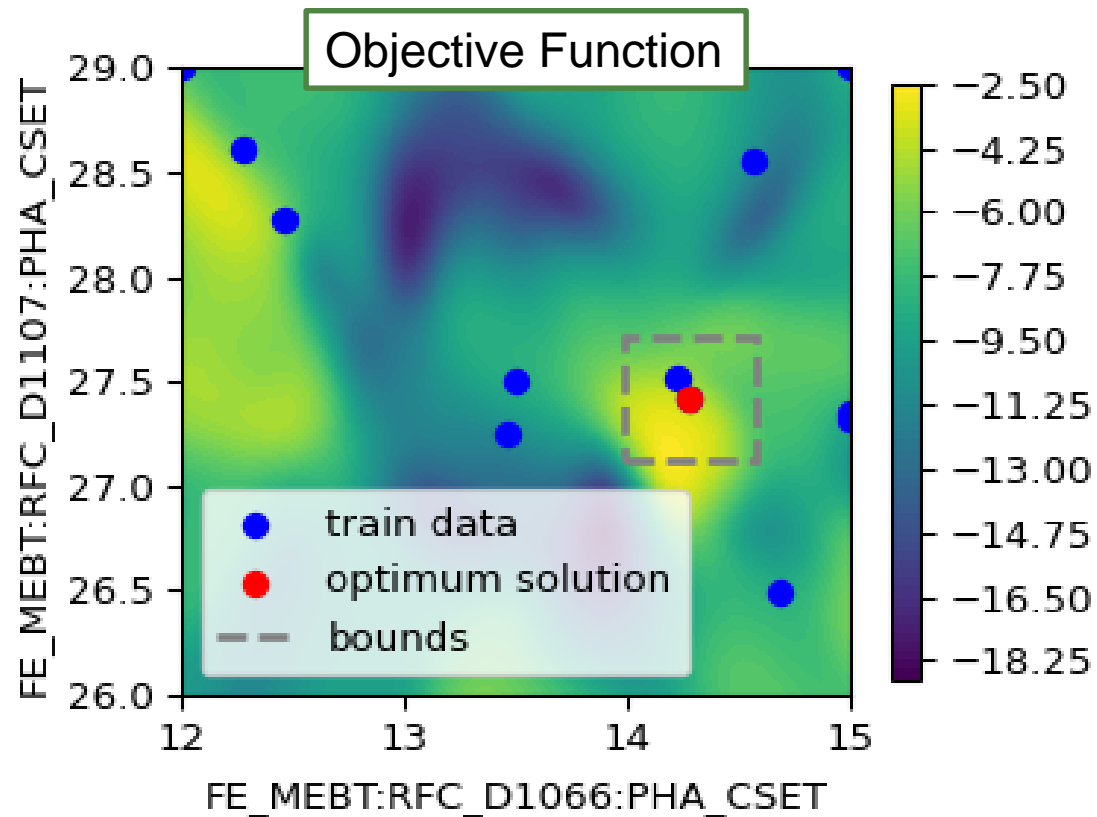
# Customized BO: Global

- Transit from Global (over whole control domain) to Local (over narrow domain) optimization
  - GlobalBO in a limited beam-time remains challenging  $\Rightarrow$  Completion through LocalBO
- Incorporate archived data or simulation model
  - Through prior mean assisted BO (pmBO)



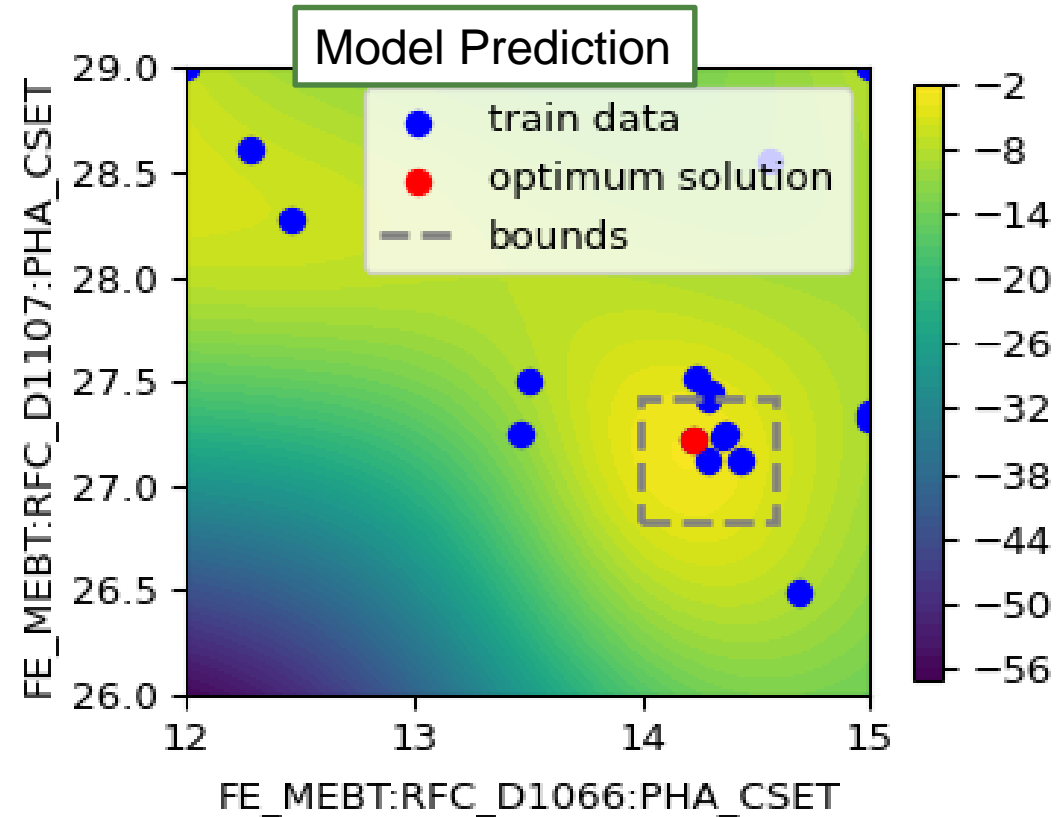
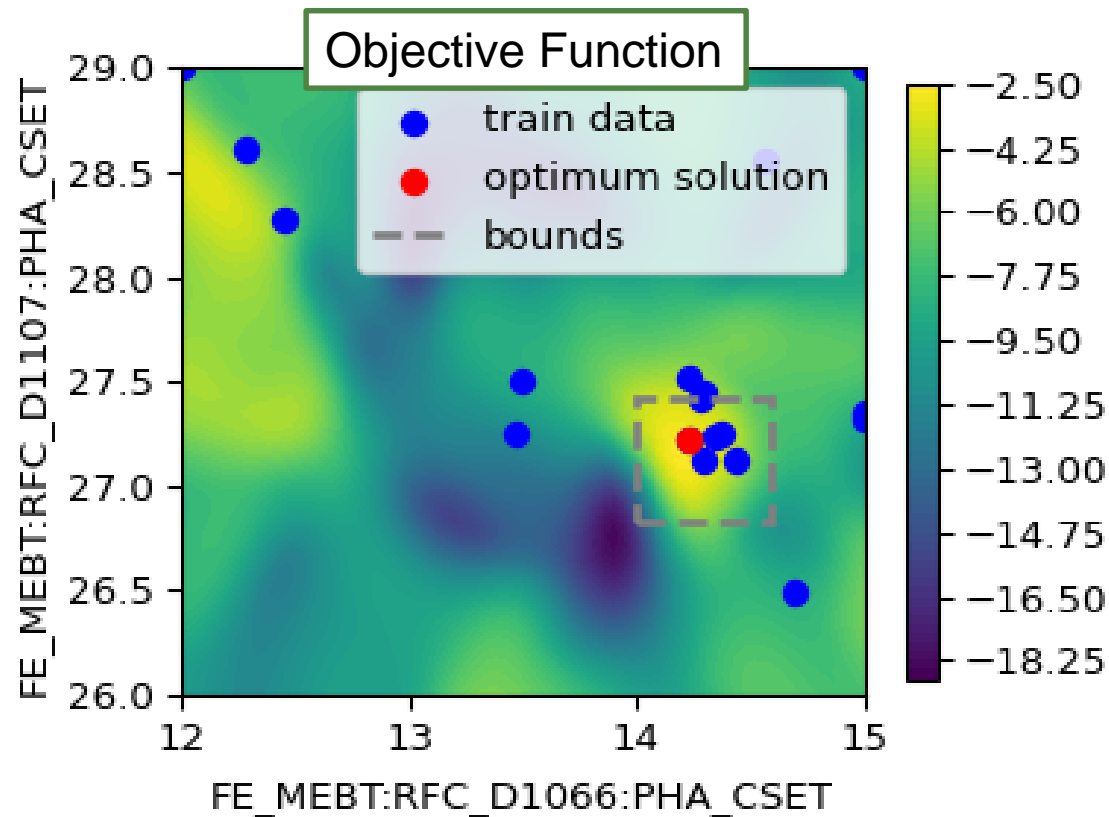
# Customized BO: Global to Local

- Transit from Global (over whole control domain) to Local (over narrow domain) optimization
  - GlobalBO in a limited beam-time remains challenging  $\Rightarrow$  Completion through LocalBO
- Incorporate archived data or simulation model
  - Through prior mean assisted BO (pmBO)

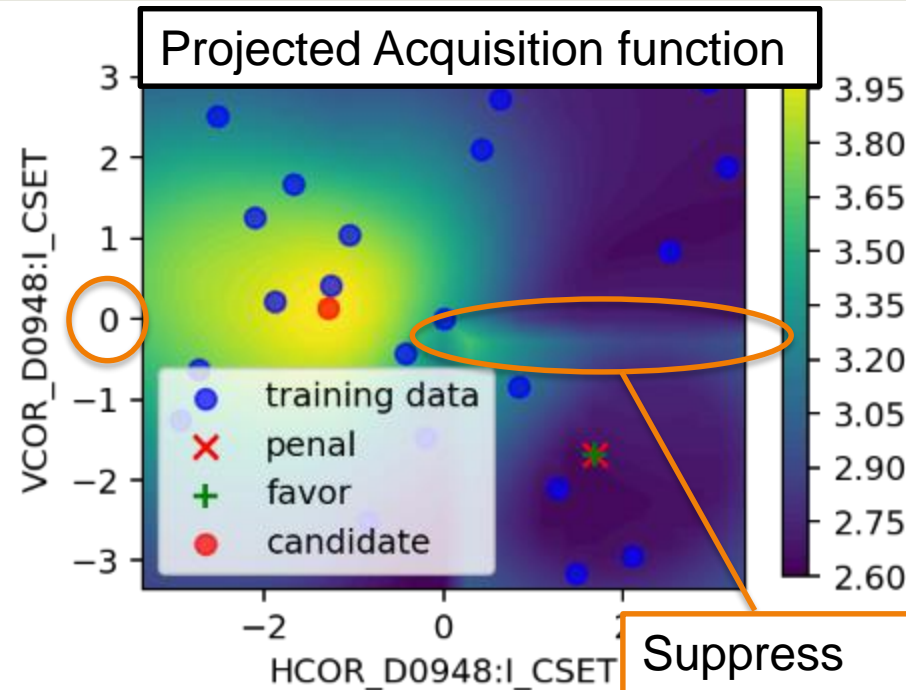


# Customized BO: Global to Local

- Transit from Global (over whole control domain) to Local (over narrow domain) optimization
  - GlobalBO in a limited beam-time remains challenging  $\Rightarrow$  Completion through LocalBO
- Incorporate archived data or simulation model
  - Through prior mean assisted BO (pmBO)



# Customized BO: Maximize Transmission through Two Apertures



Suppress polarity change by favoring current quadrant

- **Tightly avoid machine idle time**
  - Simultaneously queries candidate solutions while the machine evaluates objectives.
  - Candidate search terminates upon objective evaluation completion by machine
- **Ramping time awareness for max beam time utilization**
  - Particularly, power supply of corrector **polarity change**.
- **Order evaluation candidates (initialization points or multi-batch query) to minimize ramping cost**

$$f_{penal} = -C_{penal} e^{-(x-x_{penal})^2 / L_{penal}^2}$$

$$f_{favor} = +C_{favor} e^{-(x-x_{favor})^2 / L_{favor}^2}$$

$$f_{polarity} = \begin{cases} +C_{polarity} & \text{if } \text{sign}(x) = \text{sign}(x_{current}) \\ 0 & \text{else} \end{cases}$$

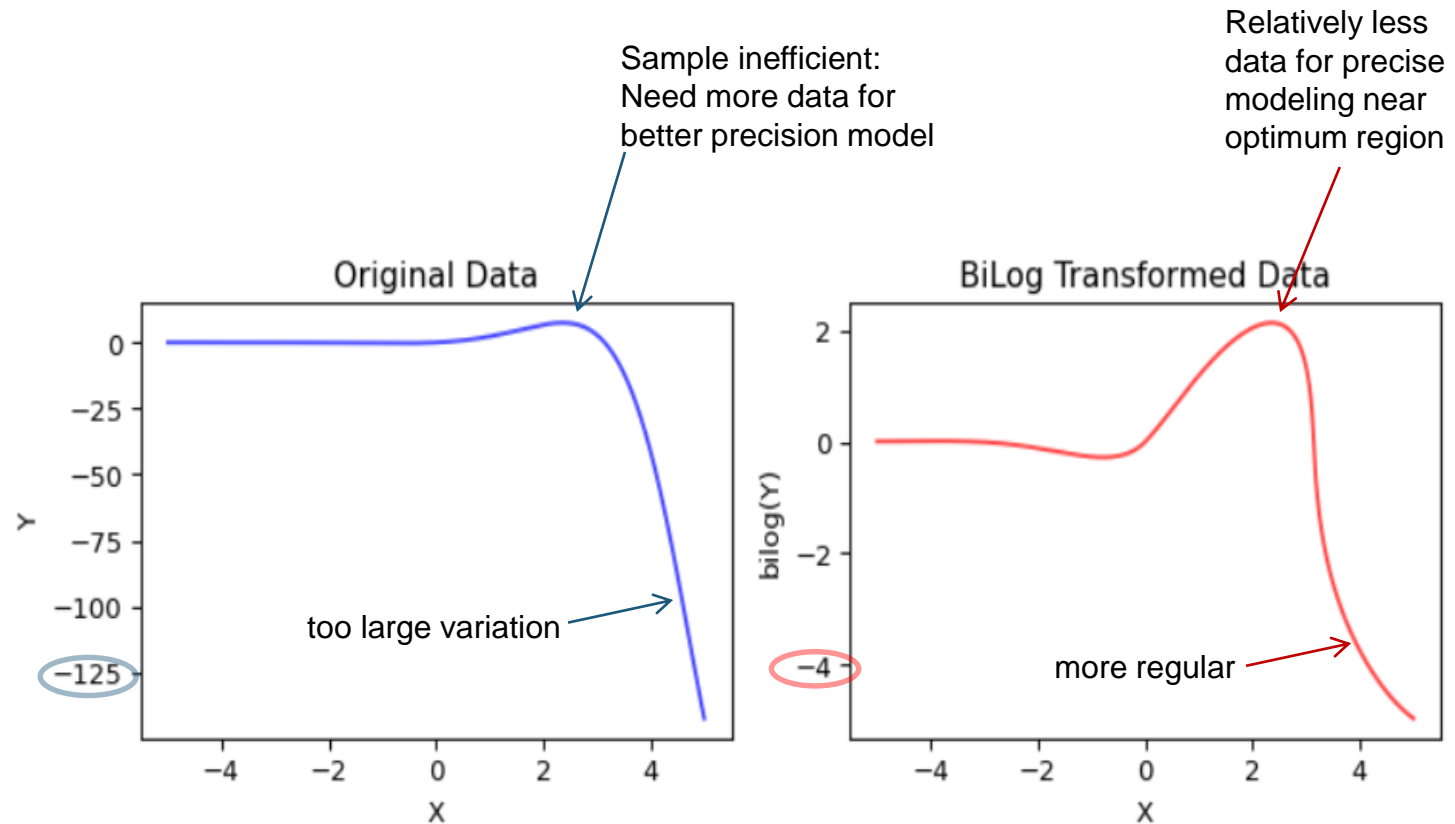


# Objective Function Construction

- Regularity of the objective function strongly effect optimization performance
- Objective function construction template saves beam time!
- We implemented intuitive scalarization of multiple objectives.

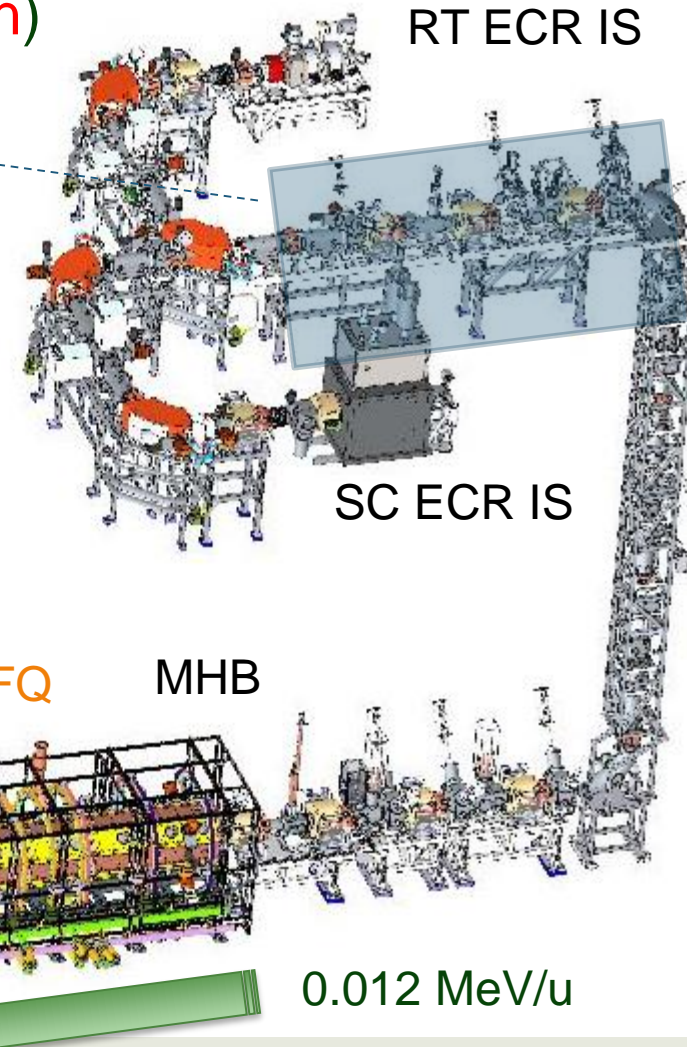
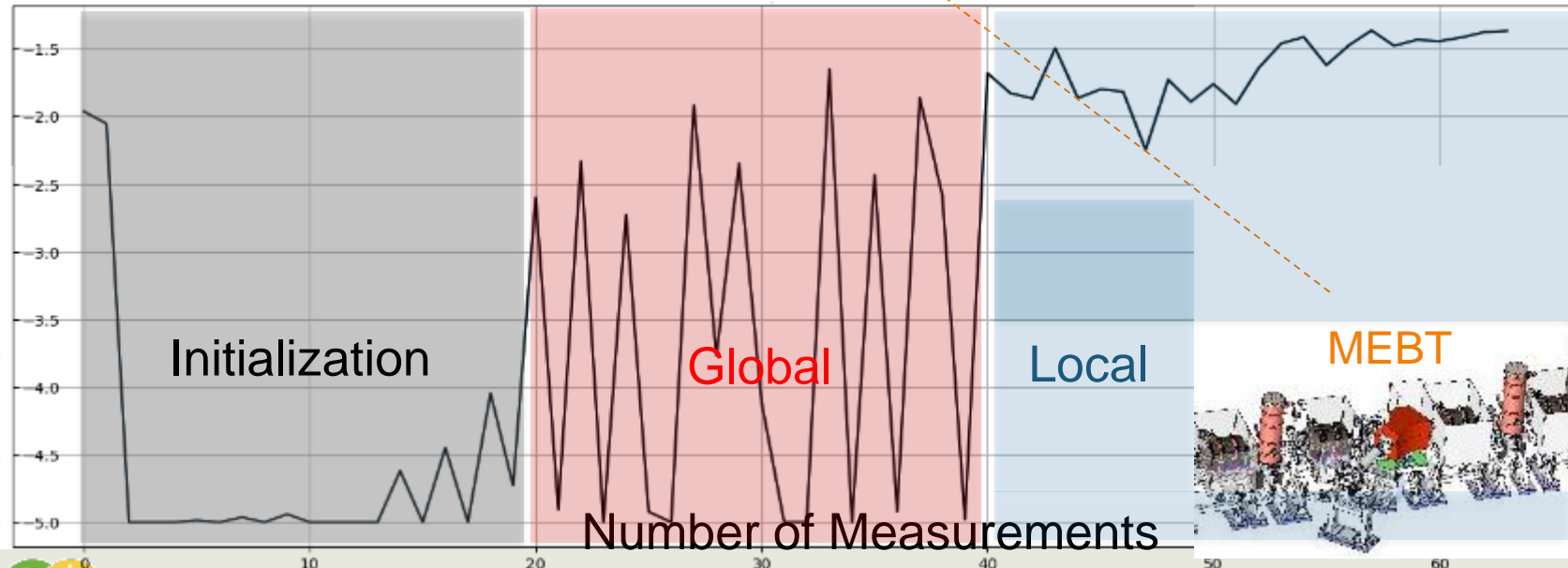
$$Obj = \sum_i w_i \varphi \left( \frac{v_{i,meas} - v_{i,target}}{tolerance_i} \right)$$

$$Obj_{bilog} = \text{sign}(Obj) \log(1 + |Obj|)$$



# Front End Tuning

- Steer beam through two apertures in straight line without BPMs (4 min)
  - Limited fast beam centroid diagnostics.
  - Maximize beam current through apertures  $\Rightarrow$  cut beam halo
- RFQ transmission and MEBT 6D beam centering (18 min)
- LS1 longitudinal beam center restoration (2 min)
  - Objective of FC reading after two apertures



# Motivation: Accelerators are complex time-varying systems (especially FRIB!)

- Collective effects: space charge forces.

## Dataset for AI training summary

13 beam species

Charge states:

26, 22, 23, 24, 25, 27, 28, 1, 2, 3, 4, 5, 6

Ion mass numbers:

124, 124, 124, 124, 124, 124, 124, 16, 16, 16, 16, 16, 16

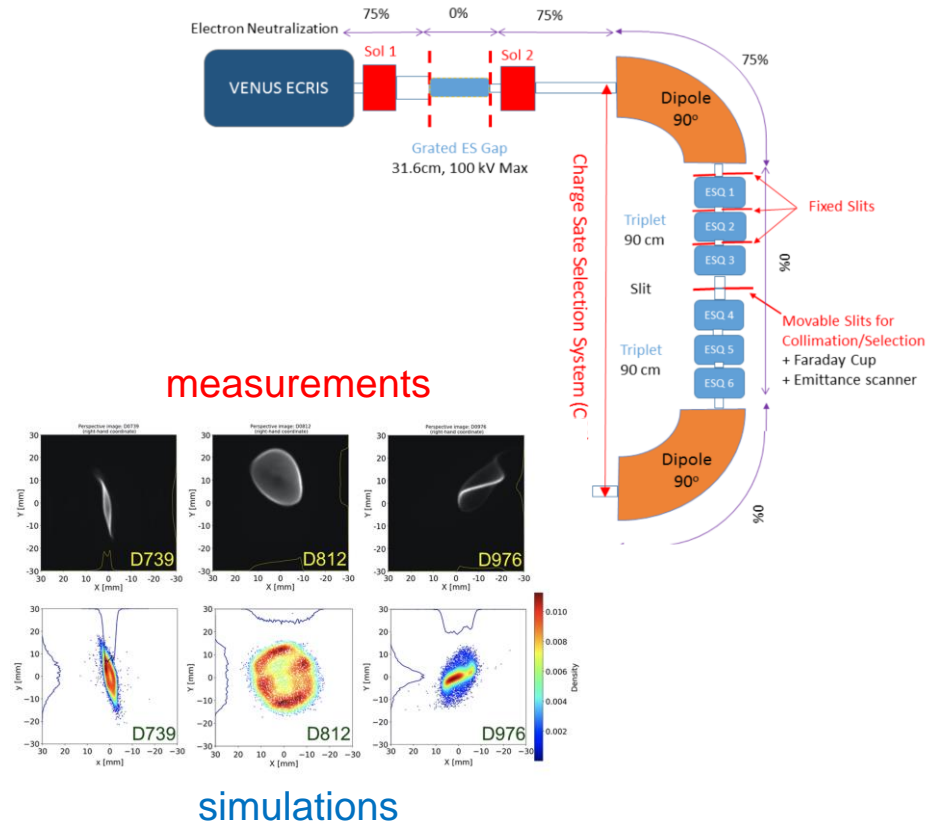
Macroparticles per charge state: 300 K

Total macroparticles per simulation: 3.9 M

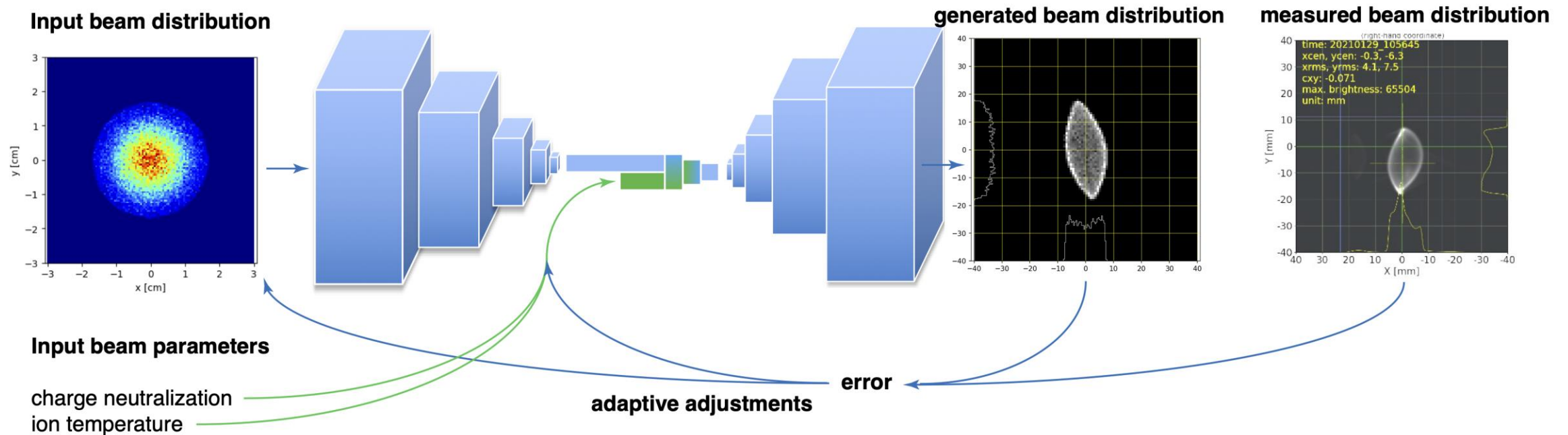
Time for 1 simulation, using 96 cores: ~6 hours

Total simulations run: ~420

## FRIB CSS

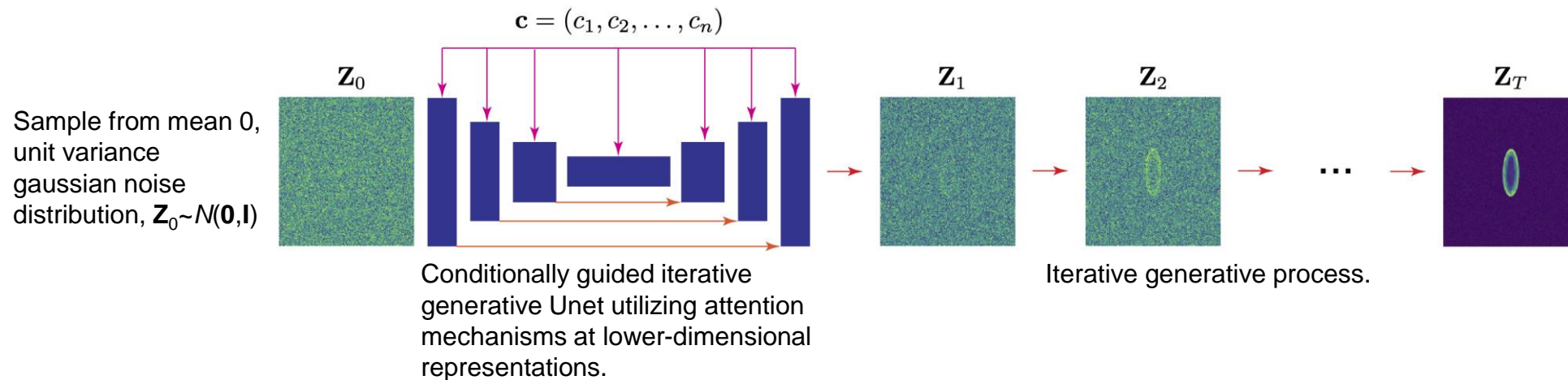


# Goal: Create Adaptive Generative Deep Learning Tools for Virtual 4D Phase Space Diagnostics



# Conditional Guided Diffusion for Creating 4D Phase Space

- Generative diffusion is the state-of-the-art AI-based method for creating high resolution representations of complex objects, such as all 6 unique projections of the FRIB beam's 4D phase space (assuming a beam uniform in  $z$  with little/no energy spread).
- Conditional vector  $c$  contains lattice parameters (such as magnet settings), which of the beam's 2D phase space projections to generate, and where along the lattice to generate it.





# Conditional Guided Diffusion for Creating 4D Phase Space

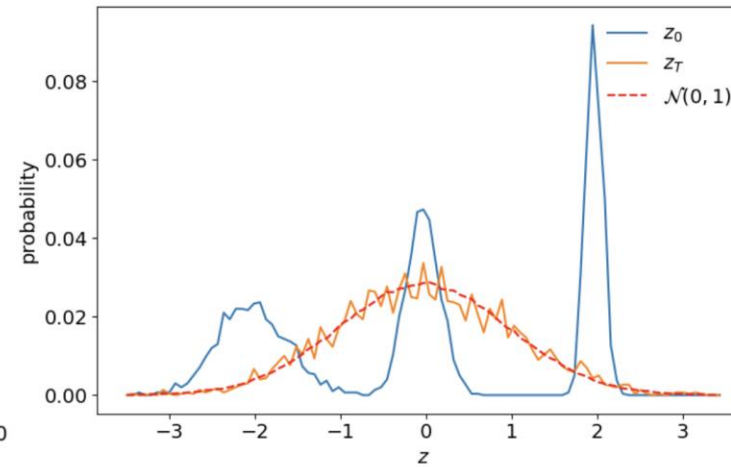
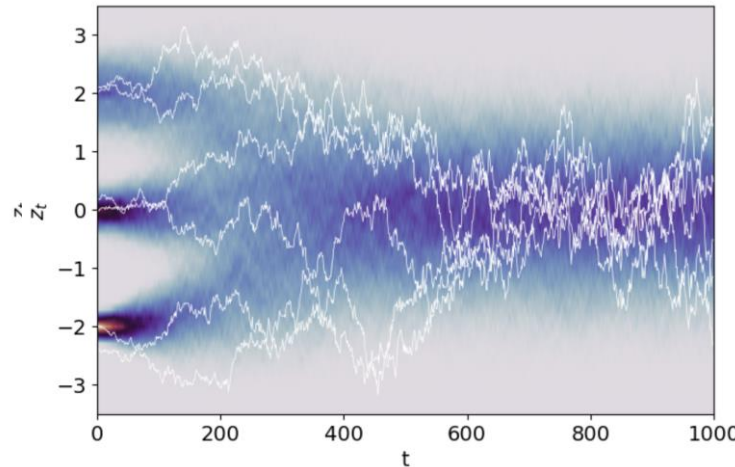
## Iterative Steps

$$\begin{aligned} \mathbf{z}_1 &= \sqrt{1 - \beta_1} \mathbf{Y} + \sqrt{\beta_1} \boldsymbol{\varepsilon}_1, & \mathbf{z}_2 &= \sqrt{1 - \beta_2} \mathbf{z}_1 + \sqrt{\beta_2} \boldsymbol{\varepsilon}_2 \\ & \vdots & & \\ \mathbf{z}_t &= \sqrt{1 - \beta_t} \mathbf{z}_{t-1} + \sqrt{\beta_t} \boldsymbol{\varepsilon}_t, & \boldsymbol{\varepsilon}_t &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \beta_t \in [0, 1], \quad t \in \{1, 2, \dots, T\} \end{aligned}$$

## Probability Distributions of a Markov Chain

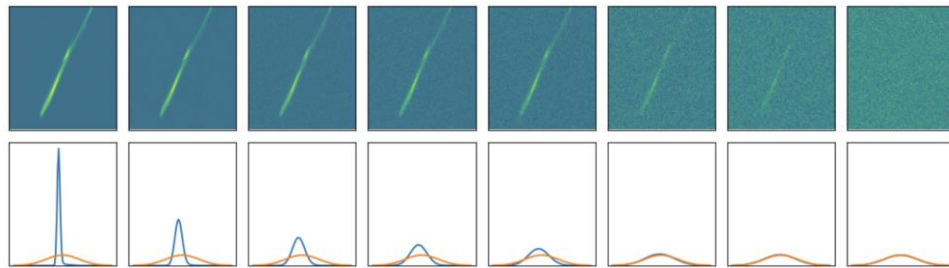
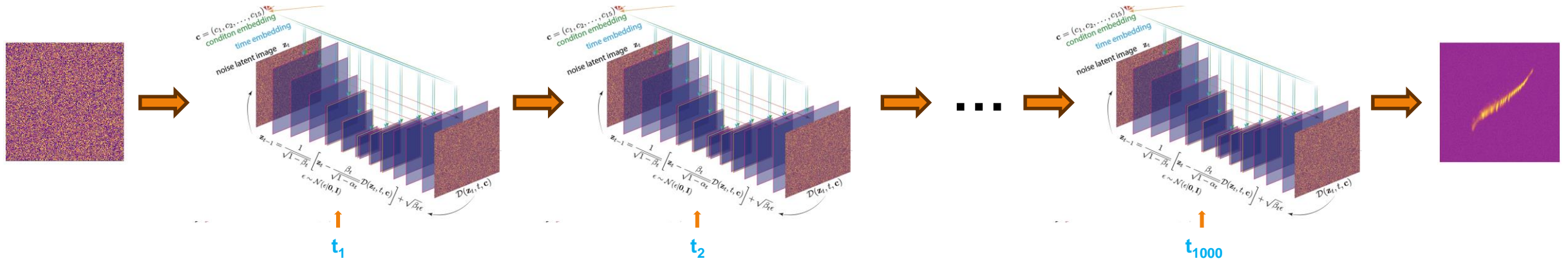
$$q(\mathbf{z}_1 | \mathbf{Y}) = \mathcal{N}_{\mathbf{z}_1} \left[ \sqrt{1 - \beta_1} \mathbf{Y}, \beta_1 \mathbf{I} \right], \quad q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}_{\mathbf{z}_t} \left[ \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I} \right] \quad q(\mathbf{z}_{1:T} | \mathbf{Y}) = q(\mathbf{z}_1 | \mathbf{Y}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$$

$$\mathbf{z}_t = \sqrt{\alpha_t} \mathbf{Y} + \sqrt{1 - \alpha_t} \boldsymbol{\varepsilon}, \quad \alpha_t = \prod_{s=1}^t (1 - \beta_s) \quad \implies \quad q(\mathbf{z}_t | \mathbf{Y}) = \mathcal{N}_{\mathbf{z}_t} \left[ \sqrt{\alpha_t} \mathbf{Y}, (1 - \alpha_t) \mathbf{I} \right] \quad \longrightarrow \quad \mathcal{N}(0, 1)$$



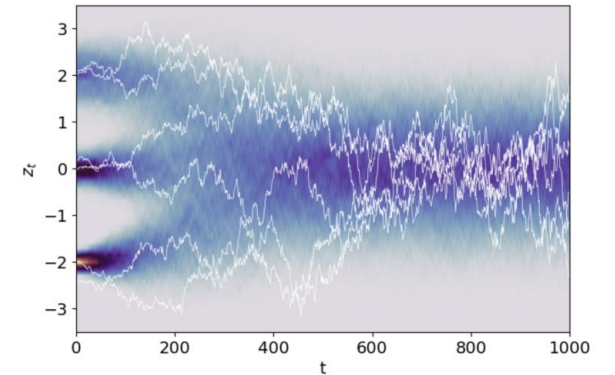


# Conditional Guided Diffusion for Creating 4D Phase Space



$$dz = \mathbf{f}(z, t)dt + g(t)d\mathbf{w}, \quad \text{SDE}$$

$$dx = [\mathbf{f}(x, t) - g(t)^2 \nabla_x \log p_t(x)]dt + g(t)d\bar{\mathbf{w}}$$



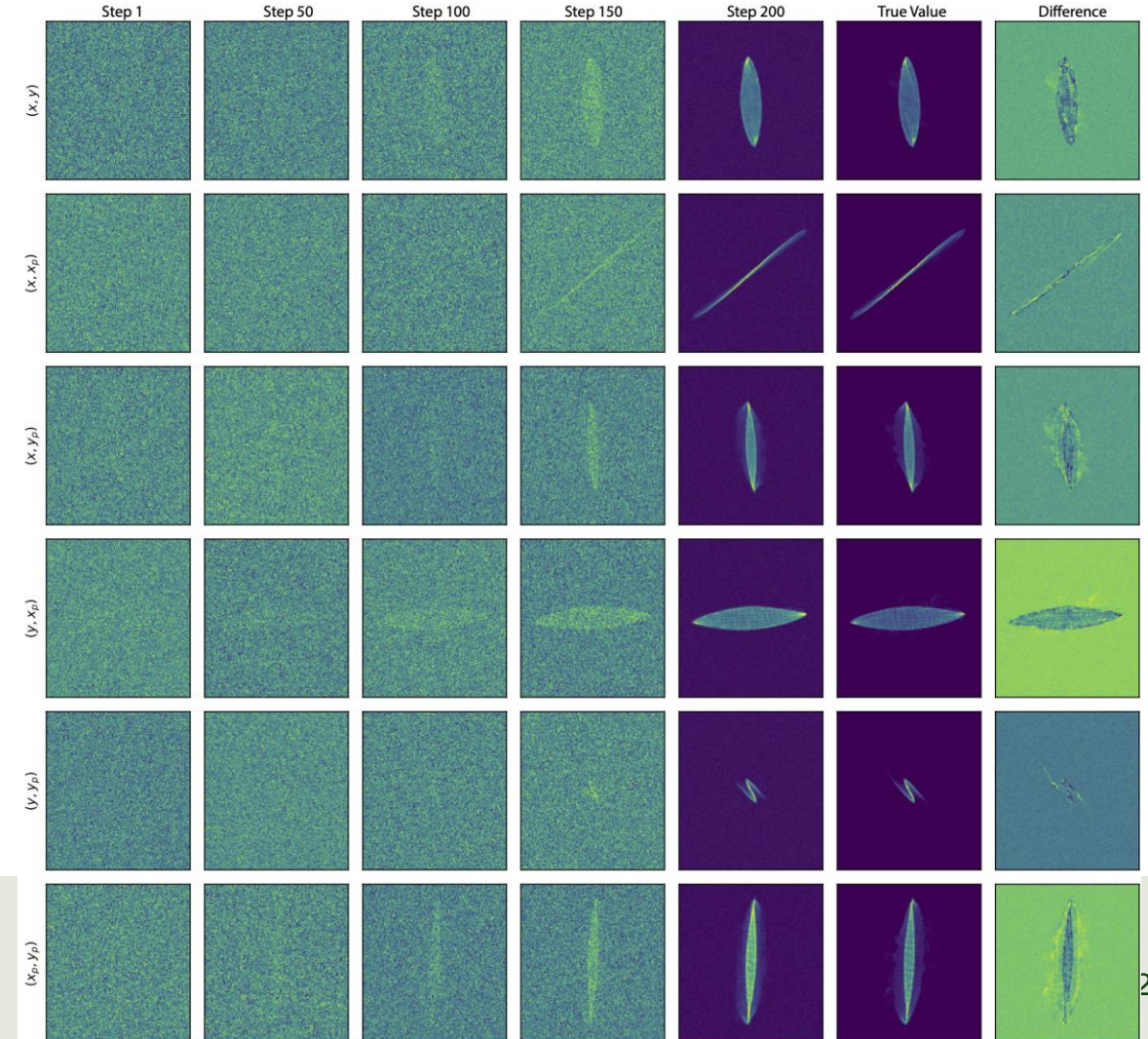
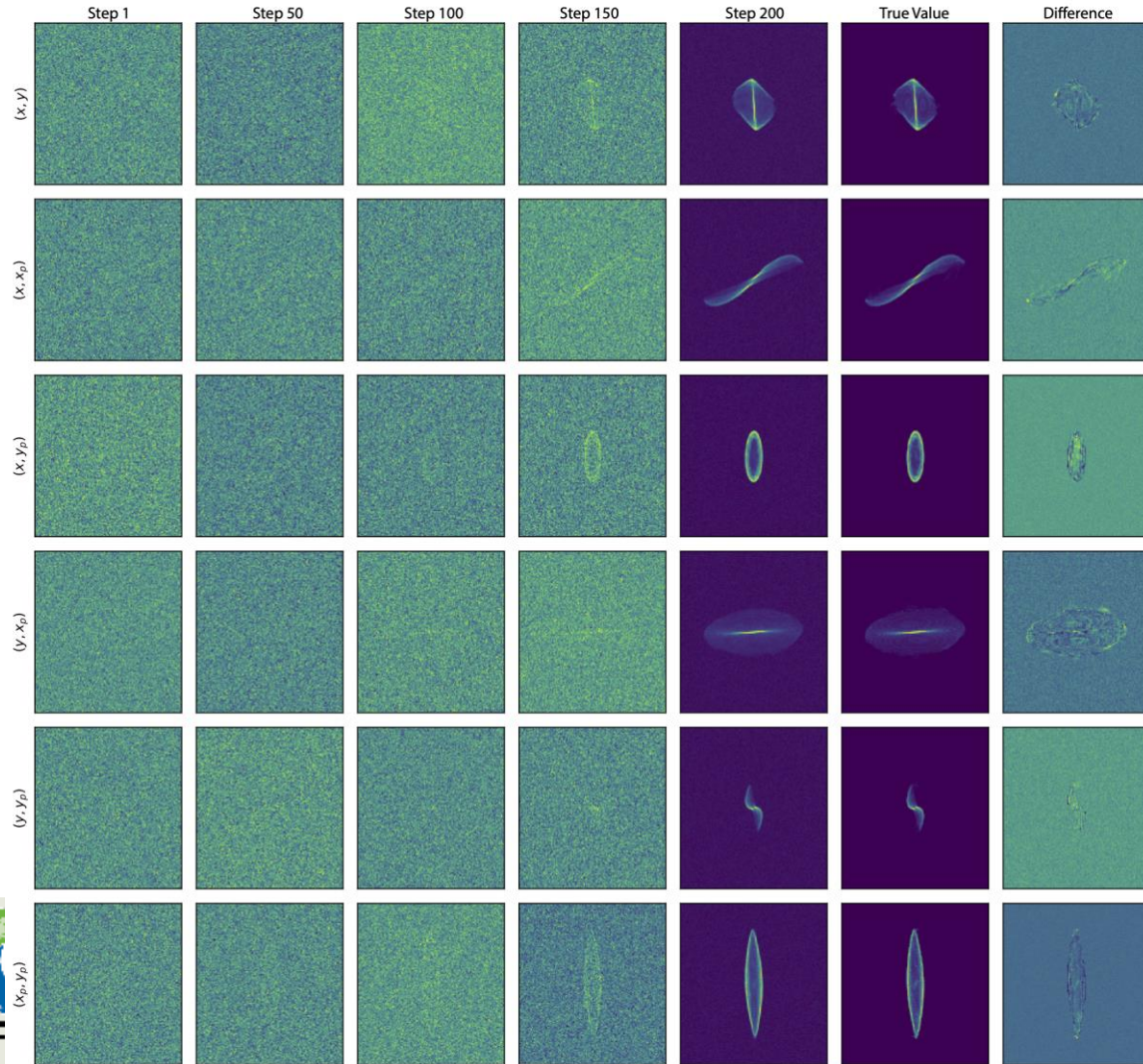
$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{r}, t) = -\frac{\hbar^2}{2m} \nabla^2 \Psi(\mathbf{r}, t) + V(\mathbf{r}) \Psi(\mathbf{r}, t)$$

infinite?



# Preliminary Results: Conditional Guided Diffusion for Creating 4D Phase Space

- Generative diffusion is the state-of-the-art AI-based method for creating high resolution representations of complex objects, such as all 6 unique projections of the FRIB beam's 4D phase space (assuming a beam uniform in z with little/no energy spread).



# General GUI App for Bayesian Optimization

- PyBOApp is a general GUI app built upon the PHANTASY (Virtual Accelerator) framework, specifically by using its UI components and libraries.
- The App is integrated into App Launcher, an app manages the physics applications for the accelerator commissioning and operations.
- Manages Bayesian Optimization (BO) tasks with a configuration file.
- Command line interface tools are provided for testing and development.
- Debian packages were created and deployed in the Linux system





# PyBOApp: User Interface and Interactions

The screenshot displays the PyBOApp interface with several key components:

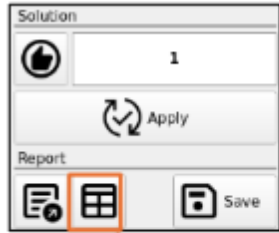
- Top Panel:** Includes "Live device power ON/OFF status" and "Live Beam Info" (showing 59 Cu 28° and 29 Cu 20.0).
- Task Configuration:** Shows "BO Task: 04-09 FC0977" with fields for Data Directory, Project Name, and Description.
- Decision Knobs Table:** A table with columns for Setpoint PV, Readback PV, Minimum, Maximum, Tolerance, Dispatch At (s), Live Setpoint, Live Readback, and Best Sotar. Row 5 is highlighted with a circled '5'.
- Actions Panel:** Contains buttons for "Run All", "Initialize", "Global", "Local", and "FineTune".
- Data Trends:** Two plots at the bottom: "Objective" (labeled with a circled '9') and "Decision Knobs" (labeled with a circled '8').



- Each task is presented as a clickable tab
- ① Defines where the data files to be saved
- ② Defines how the data is set and get
- ③ Defines the objective
- ④ Defines the budgets
- ⑤ Defines the decision knobs
- ⑥ Run the BO routines all or stage by stage
- ⑦ Get the solution, task report, save the data
- ⑧ Activate the data/figure updating
- ⑨ Data visualization
- ⑩ Undock by double-clicking

- Executing the task in a single button click
- Or, iterating each stage individually
  - 
  - 
  - 
  -
- Support Abort the running







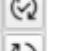





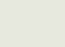
# PyBOApp: Solution Tracking and Applying




- The app keeps all the solutions during the BO run in tabular format, one can click  to read the table
- The table supports sorting, e.g. on the objective device readings
- Each row of solution could be applied to the system by clicking the  button
  - The pop-up dialog indicates how much objective change is expected
  - Press OK to apply the changes

PyBOApp: Table of Trials@ctrlm-daq1.frib.msu.edu

Report: Trials with Objective Device Readings

	time	PHY_TST:VAR_0001:X_CSET	PHY_TST:VAR_0002:X_CSET	PHY_TST:ROSENBROCK_02:Y_RD	
1	2024-11-21T09:49:27	0.000	0.000	1	
2	2024-11-21T09:51:46	3.261	5.060	3.12e+03	
3	2024-11-21T09:51:45	3.264	5.065	3.14e+03	
4	2024-11-21T09:51:19	0.000	9.790	9.54e+03	
5	2024-11-21T09:51:02	4.639	31.835	1.04e+04	
6	2024-11-21T09:49:26	-2.414	-9.725	2.41e+04	
7	2024-11-21T09:50:53	0.000	21.344	4.54e+04	
8	2024-11-21T09:51:42	5.602	10.000	4.54e+04	
9	2024-11-21T09:51:44	5.597	10.000	4.55e+04	
10	2024-11-21T09:51:39	-5.638	5.183	7.09e+04	
11	2024-11-21T09:51:37	-5.645	5.192	7.11e+04	
12	2024-11-21T09:51:35	4.193	-10.000	7.58e+04	
13	2024-11-21T09:51:36	4.202	-10.000	7.66e+04	

Since it (1) is the minimal (best), it is equivalent as click  Apply, which applies the best-so-far solution.

Apply a Solution@ctrlm-daq1.frib.msu.edu

PV	Now Set (x0)	New Set (x1)	Diff (x1-x0)
PHY_TST:VAR_0001:X_CSET	3.261	0.000	-3.261
PHY_TST:VAR_0002:X_CSET	5.060	0.000	-5.060

**Expected objective reading at the New Set with the change w.r.t. Now Set**

PHY\_TST:ROSENBROCK\_02:Y\_RD: **1 (-100%)**

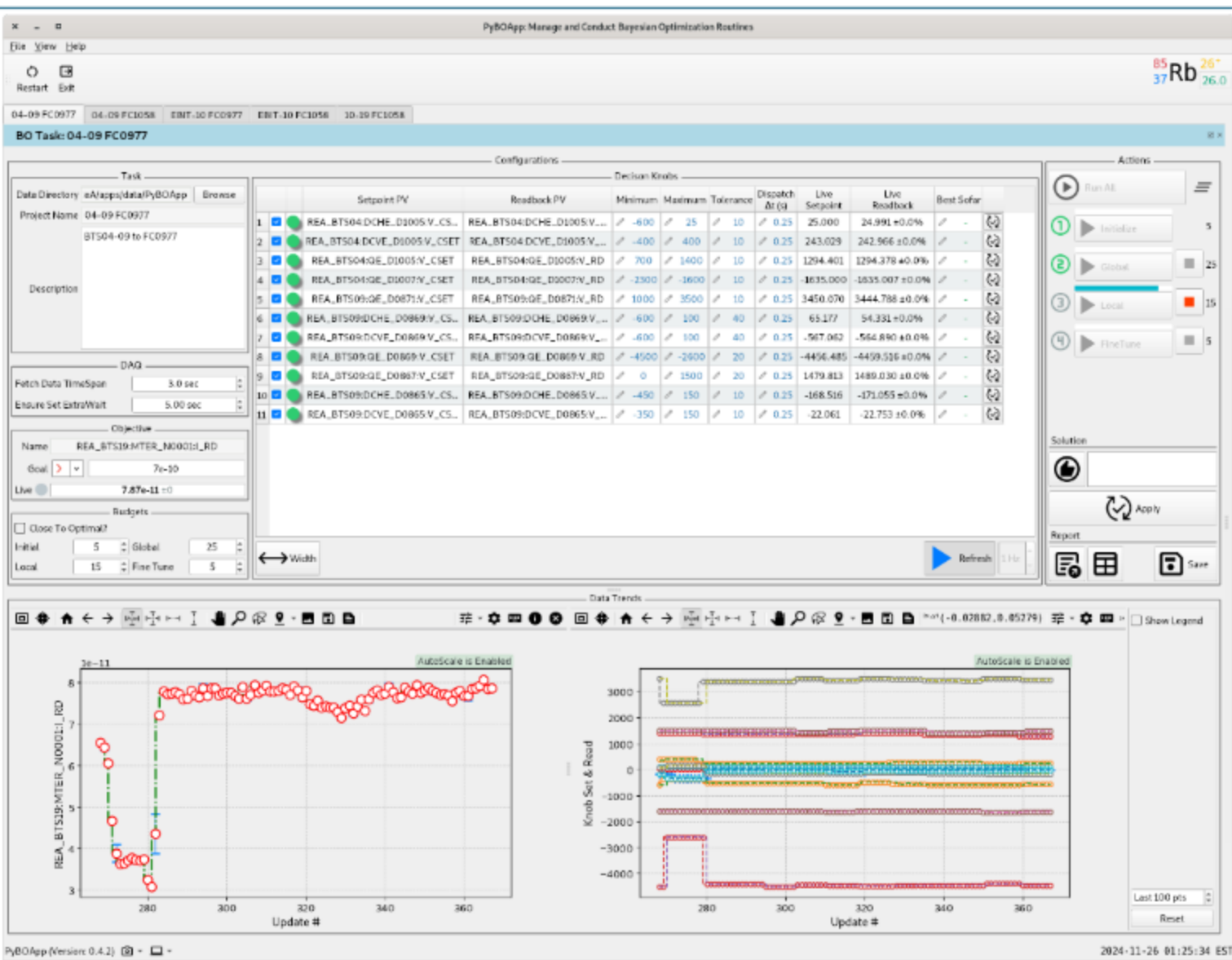
Click the **OK** button to set the **PV** with the **New Set** value for each row of the above table.

# Application on Re-Accelerator: Source Tuning

**Goal:** Maximize a Faraday Cup reading (transmission) by tuning 3 pairs of correctors (horizontal & vertical) and 5 quadrupoles, that is 11 devices in total in the upstream beamline.



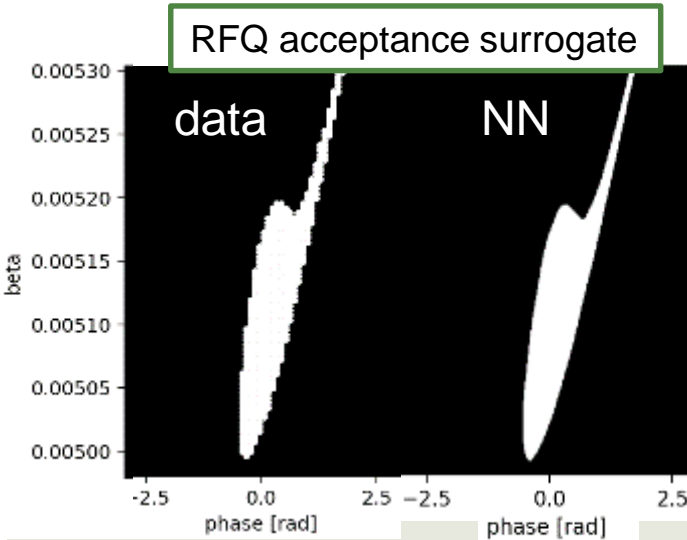
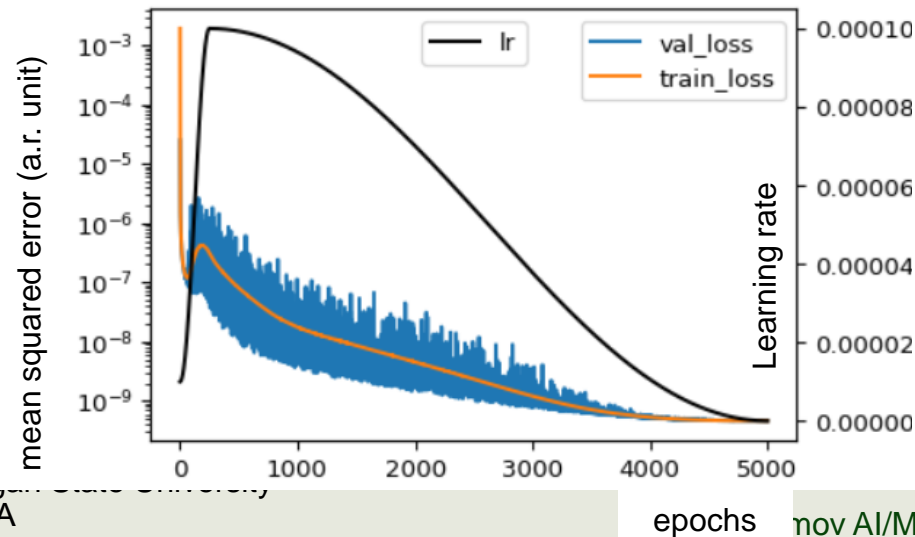
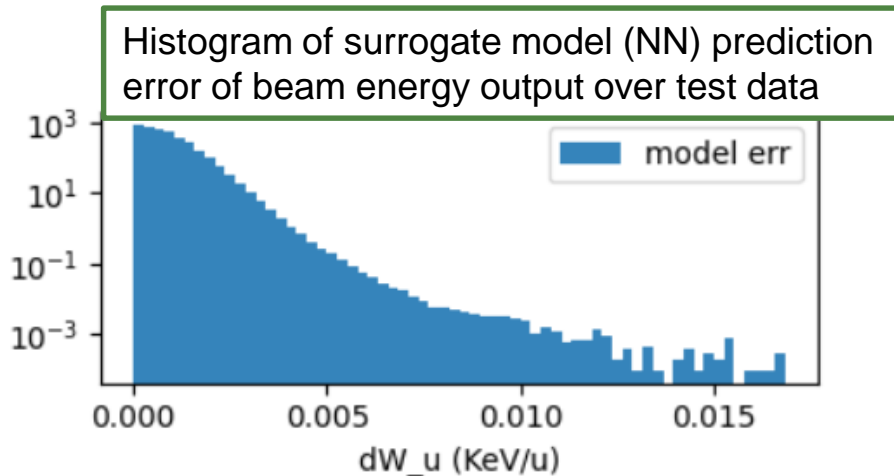
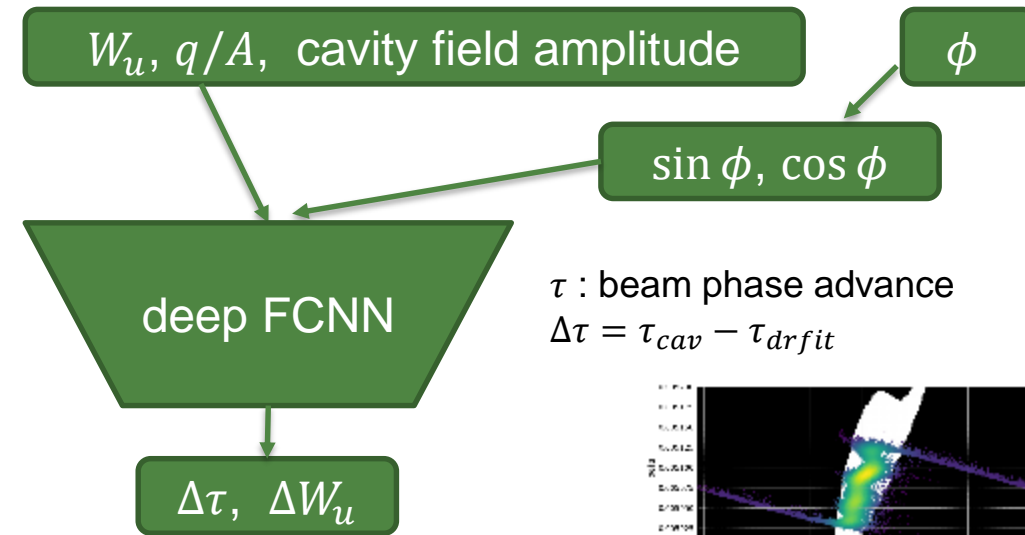
Iterating the Global, Local, FineTune run stages to find the optimal solutions





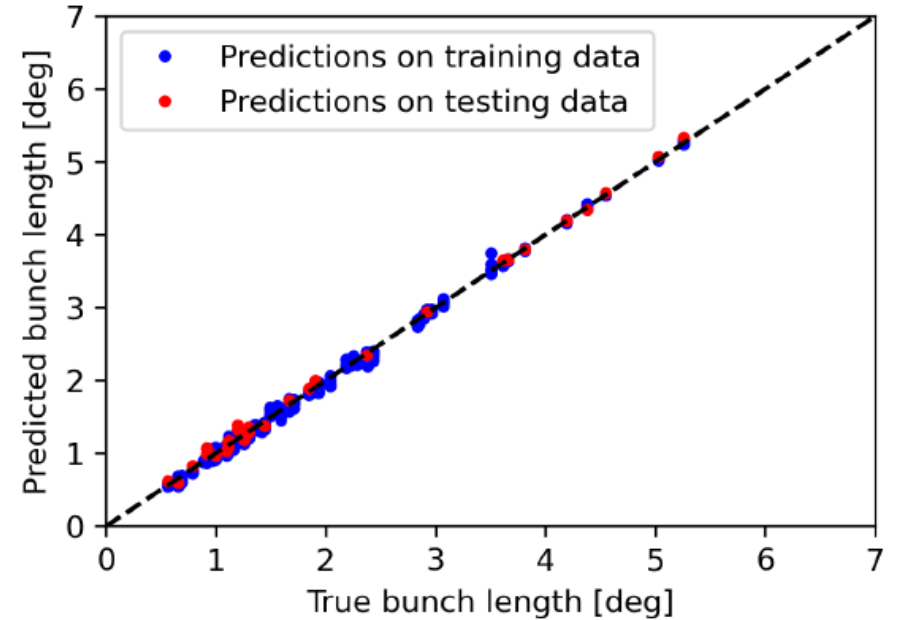
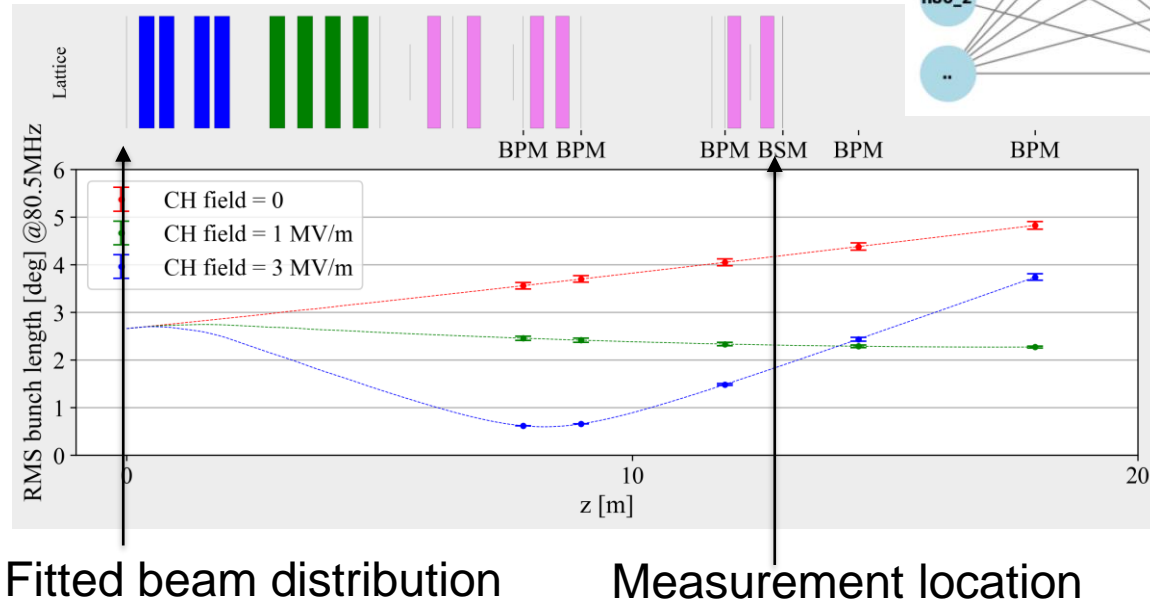
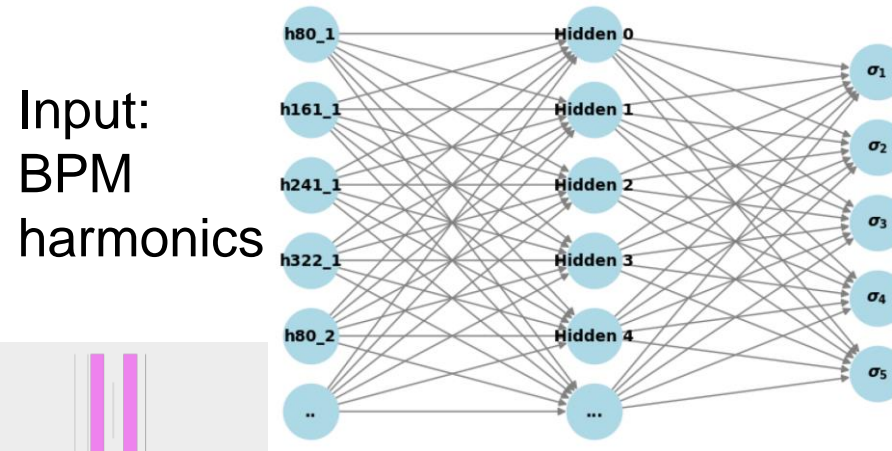
# Fast Physics Model is essential for Online Optimization

- BO can exploit physics model for online beam tuning provided the model is fast enough
  - Otherwise, constructing online model through online data is easier and faster
- Development of rapid surrogate models of physics simulators is underway.
  - surrogate modeling of 1D longitudinal RF cavity simulator achieved good accuracy and speed ( 30x faster ~ 100 us)



# Bunch Length from Beam Position Monitor (BPM)

A neural network is trained to learn the correlation of the BPM harmonics and the bunch lengths



# BPM Harmonics

For a Gaussian beam, the BPM harmonics are determined by

$$A_n = A_0 \exp\left(-\frac{(\omega_n \sigma_t)^2}{2}\right)$$

$$\frac{-\ln\left(\frac{A_n}{A_m}\right)}{(n^2 - m^2)2\pi^2 f_0^2} = \sigma_t^2$$

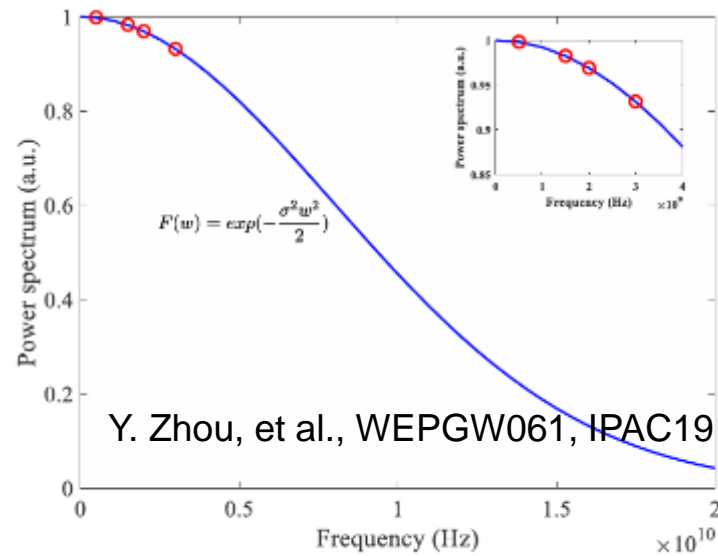
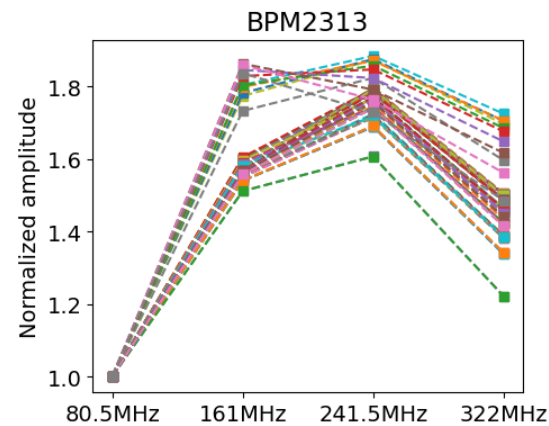
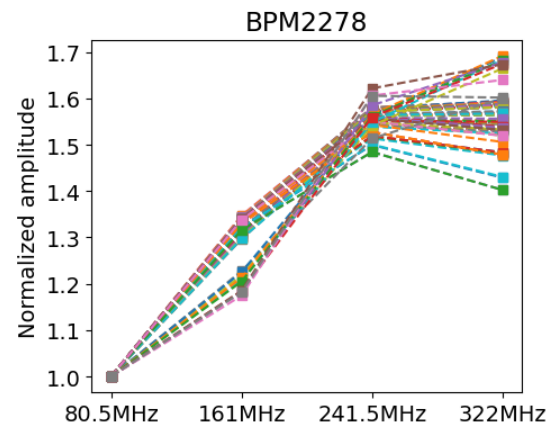
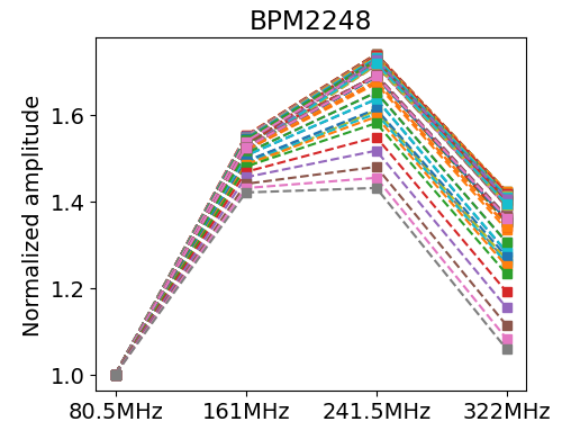
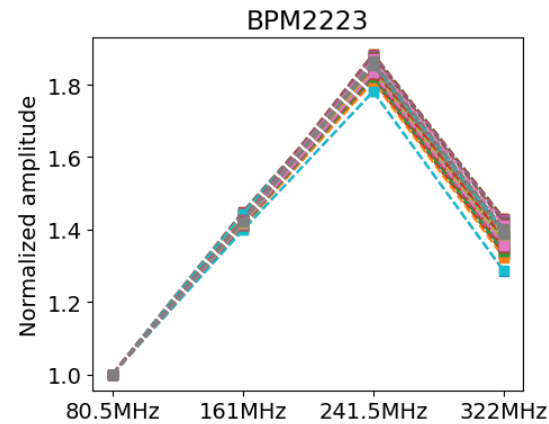
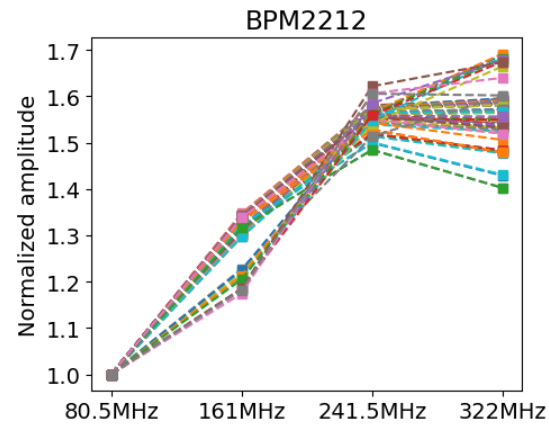


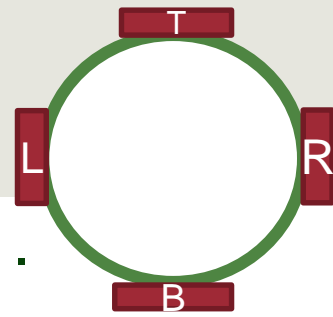
Figure 1: Power spectrum for Gaussian bunches.



However, our measurements do not show that trend.



# Beam Quadrupole Moment from BPM

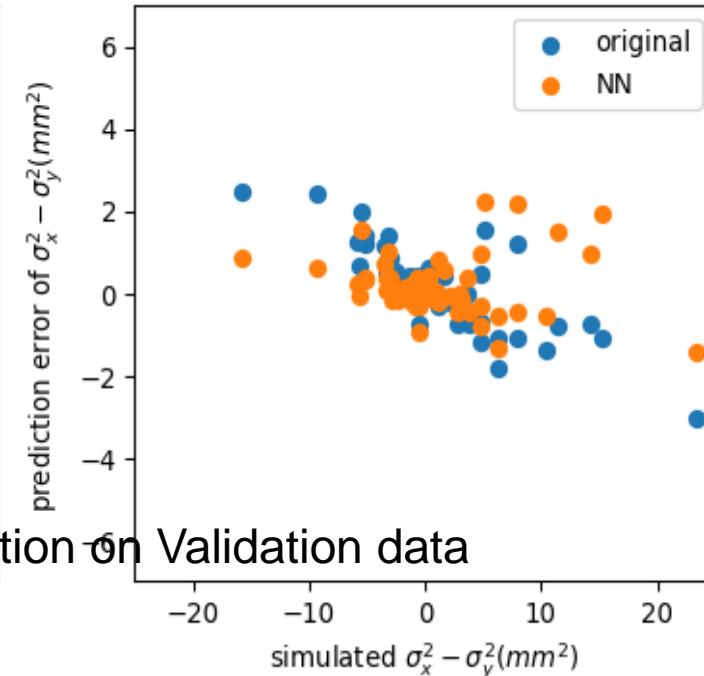
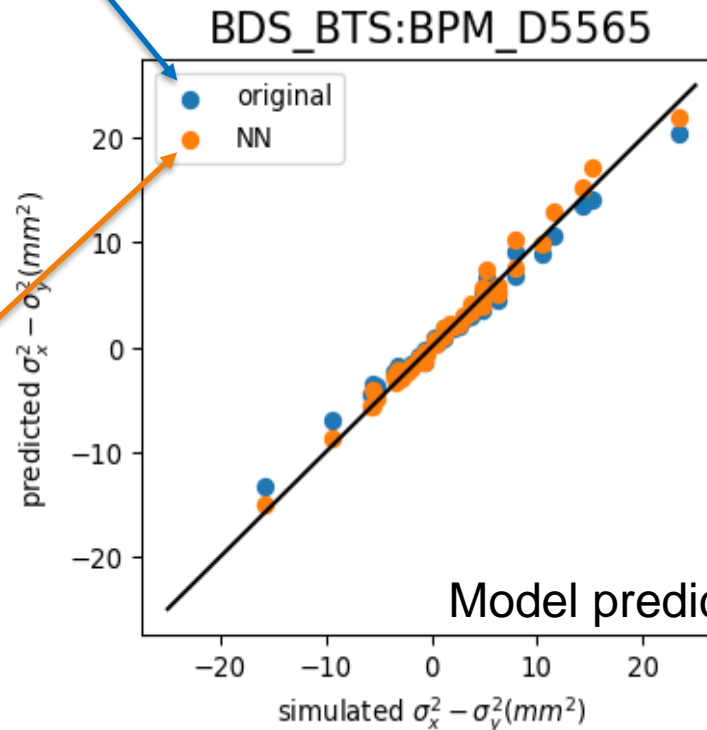


- BPM consists of 4 pick-ups to measure beam-induced signals  $U_R, U_L, U_T, U_B$ .

geometric factor

$$BPMQ = G \frac{U_R + U_L - (U_T + U_B)}{U_R + U_L + U_T + U_B} - (x^2 - y^2) \approx \sigma_x^2 - \sigma_y^2$$

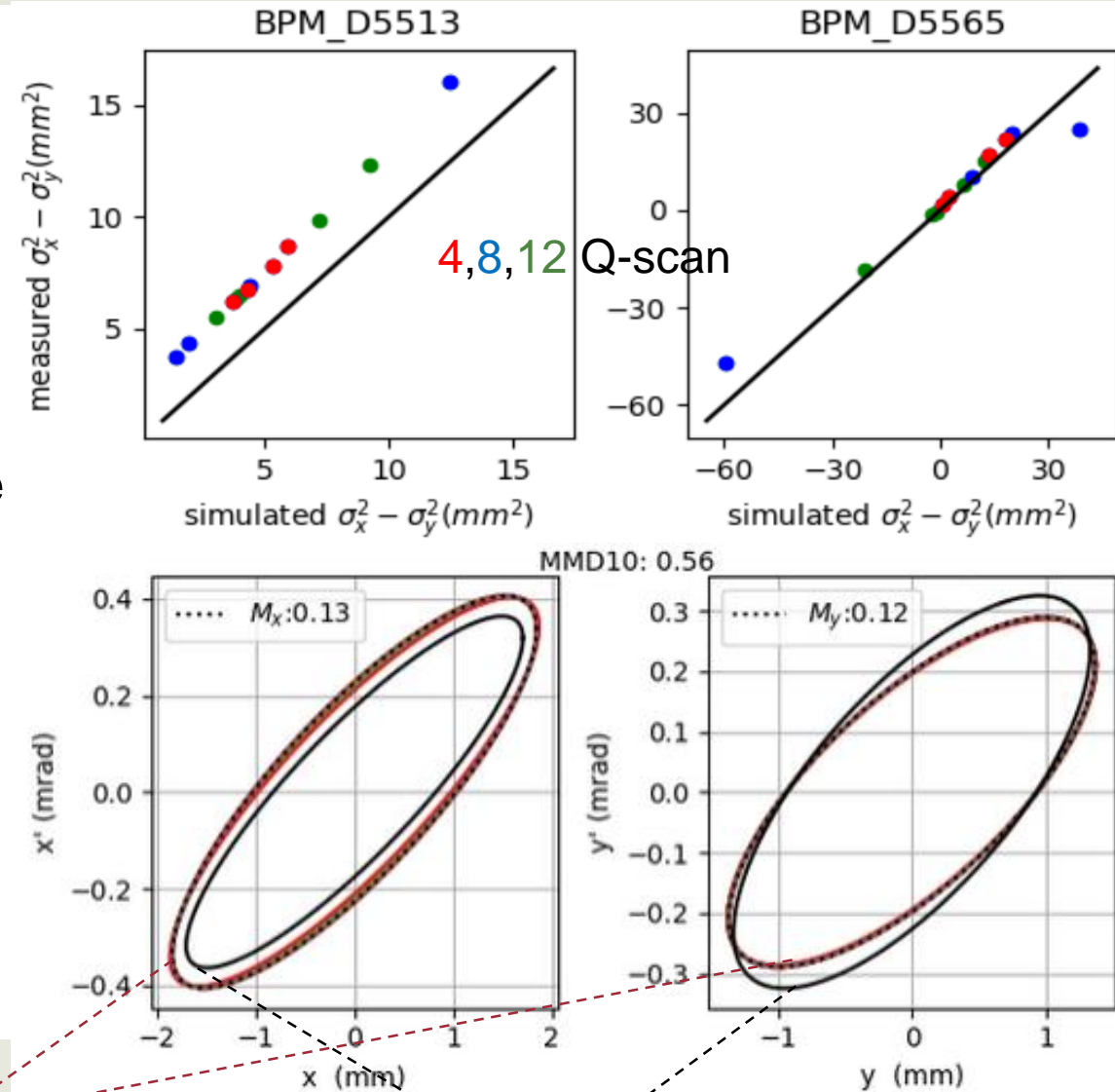
- BPM is designed to measure beam position, not the quadrupole moments
- Quadrupole moment signal strength is the 2<sup>nd</sup> order → **Weak** signal
  - Small signal error ( **calibration error**, noises ) leads to big **BPMQ error**
- Use NN to calibrate. Accuracy of data is the key for accurate NN model



$$\sigma_x^2 - \sigma_y^2 = \text{NN}(U_R, U_L, U_T, U_B)$$

# Courant-Snyder (CS) Parameters Reconstruction from BPMQ Using Bayesian Active Learning (BAL)

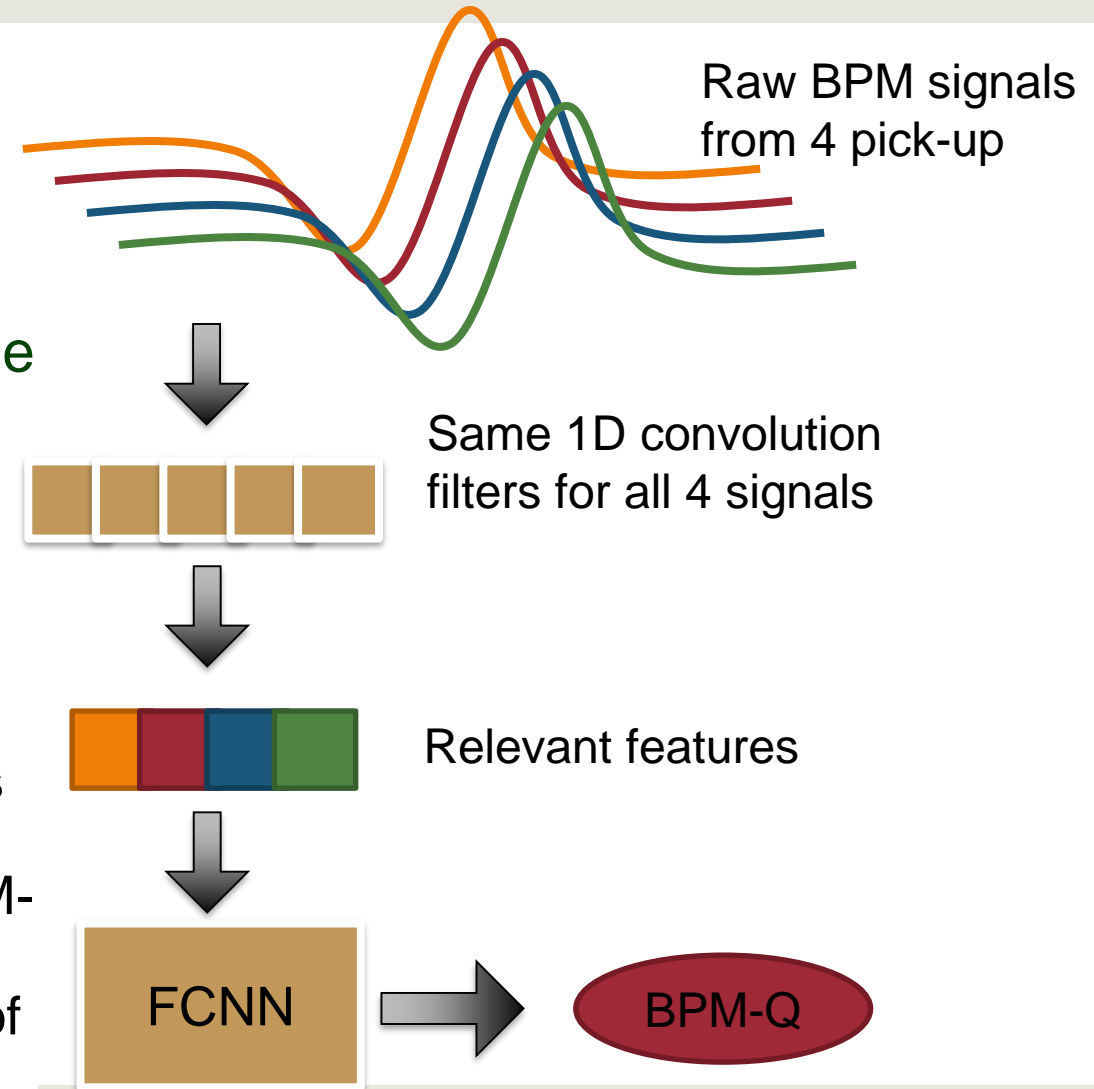
- CS parameter reconstruction using beam quadrupole moments at a few BPMs requires efficient quadrupole magnet scan.
- Ensemble of a backward-differentiable envelope simulation (with pyTorch) for surrogate model of BAL
  - Candidate quadrupole settings are queried to maximize surrogate model uncertainty of BPMQs.
- Beam test:  $^{16}\text{O}^{6+}$  ion beam (with 12 epoch of q-Scan).
  - Despite the model prediction shows some error
    - » NN model for BPM\_D5565 shows better accuracy
      - training data is more reliable due to PM5567
    - Good enough agreement between reconstructed CS using BPMQ with BAL vs PM.
      - » with BPM5565 and emittance prior





# Recent Development: Improved BPM-Q Modeling

- Amplitude of the 2<sup>nd</sup> harmonic (161MHz) of beam frequency of induced signal at BPM 4-pickups are utilized for beam position and quadrupole moment to avoid cross-talk from many other RF devices.
- Instead of the manually engineered input feature – the 161MHz component, we use NN to extract the features relevant to the BPM-Q from the full raw signal in time domain:
  - To mitigate overfitting caused by the high-dimensional input data and a relatively small training dataset, we use an encoder architecture composed of 1D convolutional layers. These layers serve to extract meaningful features by reducing input dimensionality, followed by a Fully Connected Neural Network (FCNN) that outputs the BPM-Q.
  - This resulted in about 25% improved accuracy in terms of BPM-Q error on validation data



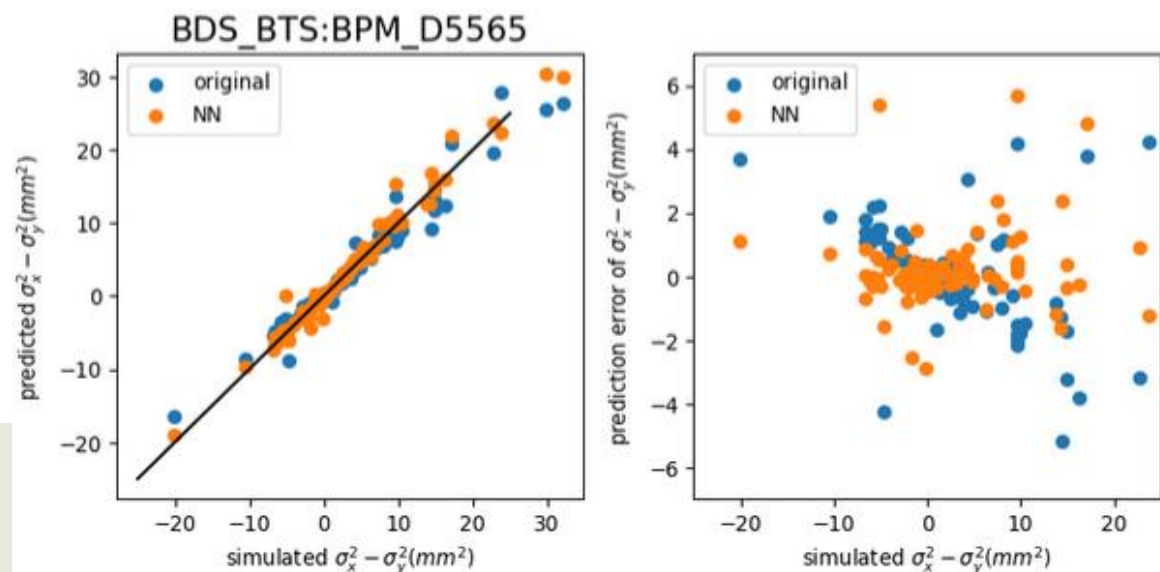
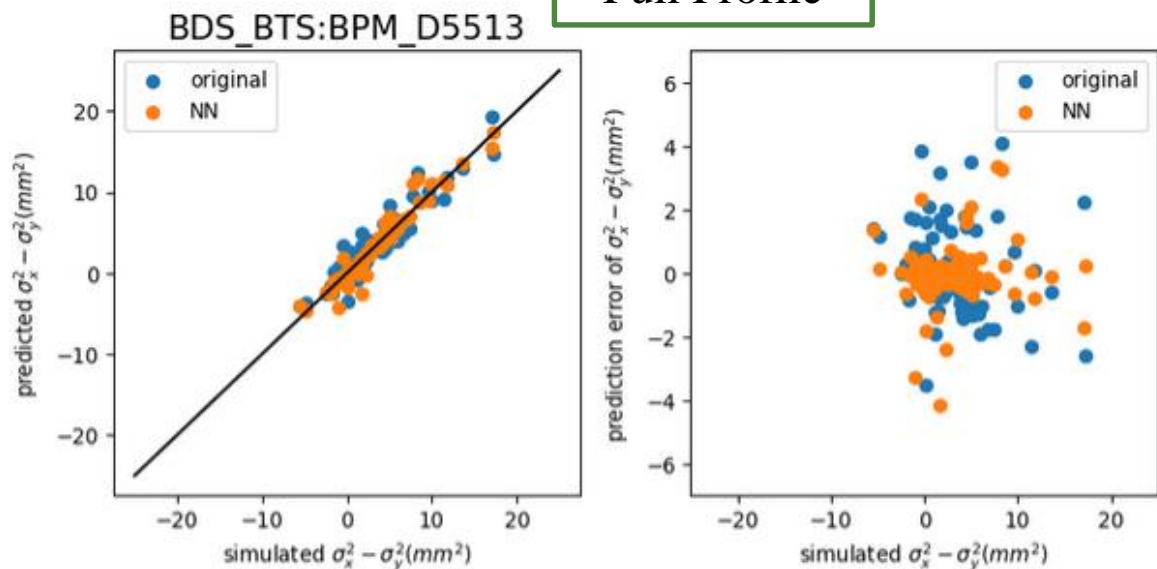


# Conclusion

- Fast beam tuning is critical for the FRIB mission
- Customized Bayesian Optimization (BO) for FRIB
  - Maximum utilization of beam time
  - Routine tasks are being established
  - Enhanced Automation, Visualization, and UI based on user feedback
- Surrogate modeling of physics simulators
  - Surrogate modeling of 1D longitudinal RF cavity simulator achieved good accuracy and speed
- Virtual Diagnostics
  - Bunch length measurement using BPM
  - Beam quadrupole moment measurement using BPM
  - Model accuracy depends on the data accuracy and quantity.
  - Strategies of CS reconstruction accuracy improvement with measurement error from virtual diagnostics is under development
    - » Even if BPMQ error is large, CS could be reconstructed within tolerable accuracy by
      - Emittance prior
      - A large number of Q-scan to avoid over-fit CS to BPM-Q error
      - Convolution filters applied for 1D bunch signals from BPM

# BPM-Q Model with Full Profile vs 161MHz Amplitude Input

Full Profile



161 MHz amplitude

