# Neural Networks for Modeling and Control of Particle Accelerators

Auralee Edelen

# We rely heavily on operators for day-to-day control tasks …



*Fermilab Control Room Photo: Reidar Hahn, FNAL*

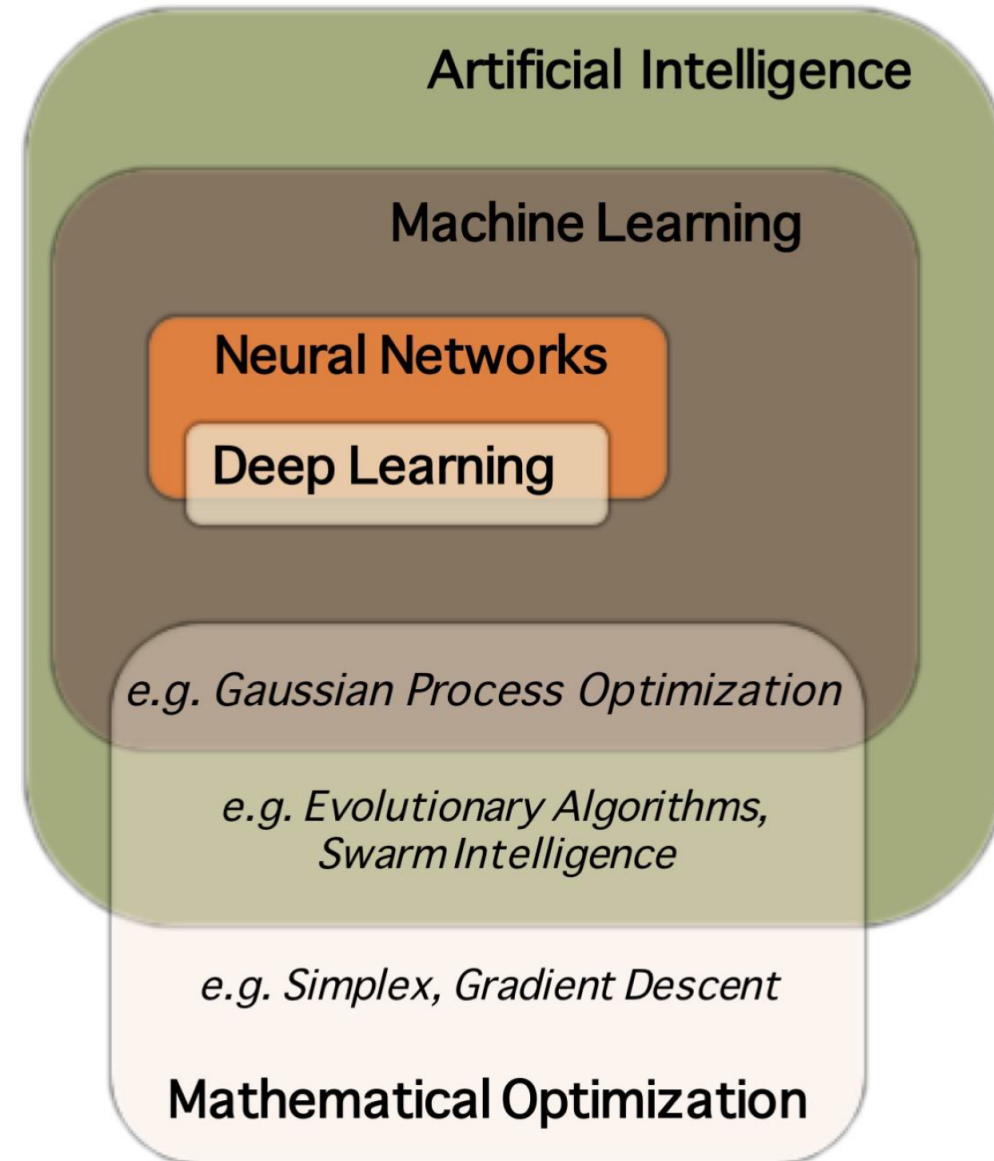*…so what can we learn from them, and what analogous techniques can we use?*

# Inspiration from Operators



Model Learning

Prediction

Planning

Diagnostic Analysis

Optimization

Learning Control

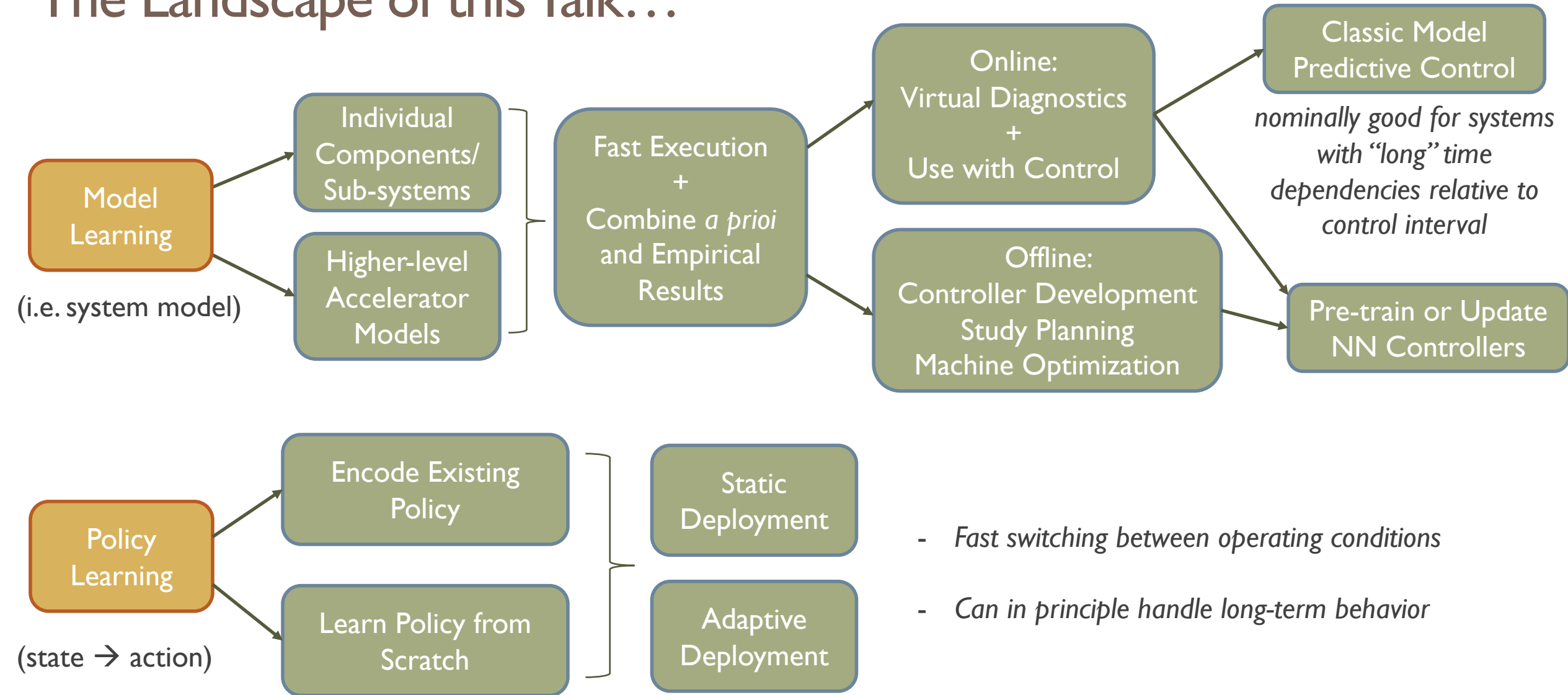*Fermilab Control Room Photo: Reidar Hahn, FNAL*

# Field Taxonomy (as of now...)

- Artificial Intelligence (AI)
  - *Concerned with enabling machines to exhibit aspects of human intelligence: knowledge, learning, planning, reasoning, perception*
  - Narrow AI: focused on a task or similar set of tasks
  - General AI: human-equivalent or greater performance on any task

- Machine Learning (ML)
  - *Enabling machines to complete tasks without being explicitly programmed*
  - Common tasks: Regression, Classification, Clustering, Dimensionality Reduction

- Neural Networks (NNs)
  - *An approach within ML that uses many connected processing units*
  - Many different architectures and training techniques

- Deep Learning (DL)
  - *Learning hierarchical representations*
  - Right now, largely synonymous with deep (many-layered) NN approaches

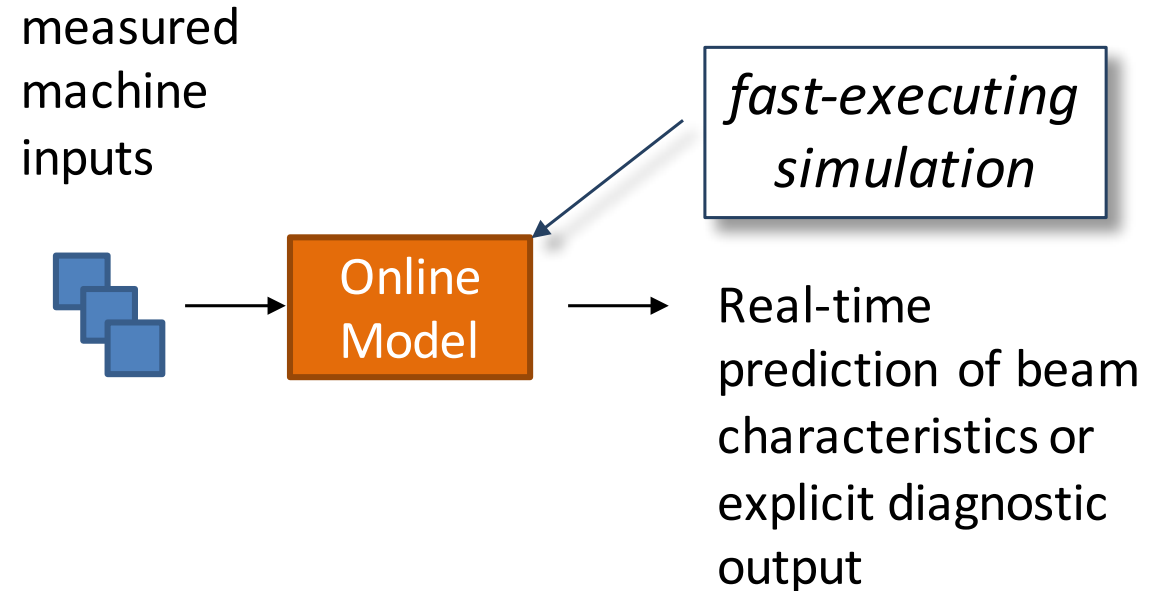***Note that these definitions are not rigid: there is a lot of fluidity in the field***

**Artificial Intelligence**

**Machine Learning**

**Neural Networks**

**Deep Learning**

*e.g. Gaussian Process Optimization*

*e.g. Evolutionary Algorithms, Swarm Intelligence*

*e.g. Simplex, Gradient Descent*

**Mathematical Optimization**

# The Landscape of this Talk…

**Model Learning**

(i.e. system model)

- Individual Components/ Sub-systems
- Higher-level Accelerator Models

Fast Execution + Combine *a prioi* and Empirical Results

Online: Virtual Diagnostics + Use with Control

Offline: Controller Development Study Planning Machine Optimization

Classic Model Predictive Control

*nominally good for systems with "long" time dependencies relative to control interval*

Pre-train or Update NN Controllers

**Policy Learning**

(state → action)

- Encode Existing Policy
- Learn Policy from Scratch

Static Deployment

Adaptive Deployment

- *Fast switching between operating conditions*
- *Can in principle handle long-term behavior*

*For all of the above, can in principle include image-based diagnostics directly*

# Online Modeling

- Use a machine model during operation
-

- Ideally:
  - Fast-executing, but accurate enough to be useful
  - Use measured inputs directly from machine
  - Combine *a priori* knowledge + learned parameters

- Applications:
  - A tool for operators + virtual diagnostic
  - Predictive control
  - Help flag aberrant behavior
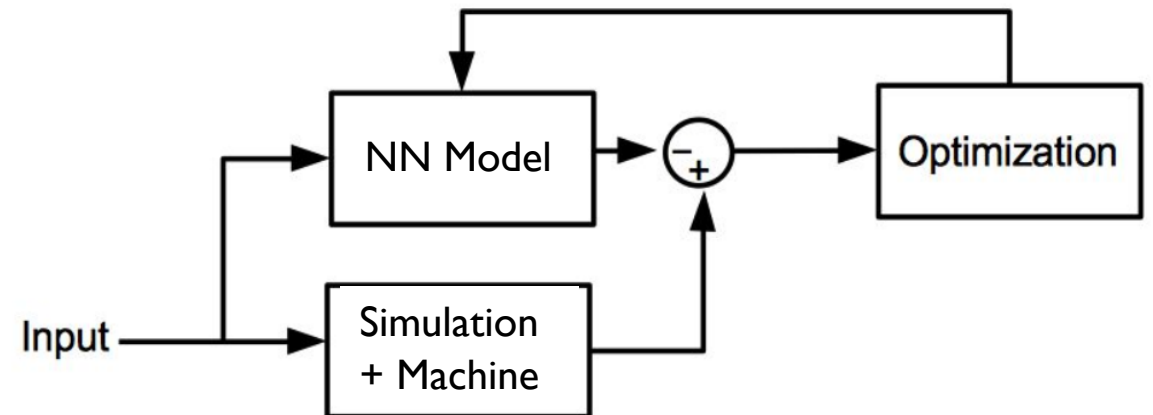  - *Bonus: control system development*

measured
machine
inputs

Online
Model

*fast-executing
simulation*

Real-time
prediction of beam
characteristics or
explicit diagnostic
output

# Online Modeling

- Use a machine model during operation
  -
- Ideally:
  - Fast-executing, but accurate enough to be useful
  - Use measured inputs directly from machine
  - Combine *a priori* knowledge + learned parameters

- Applications:
  - A tool for operators + virtual diagnostic
  - Predictive control
  - Help flag aberrant behavior
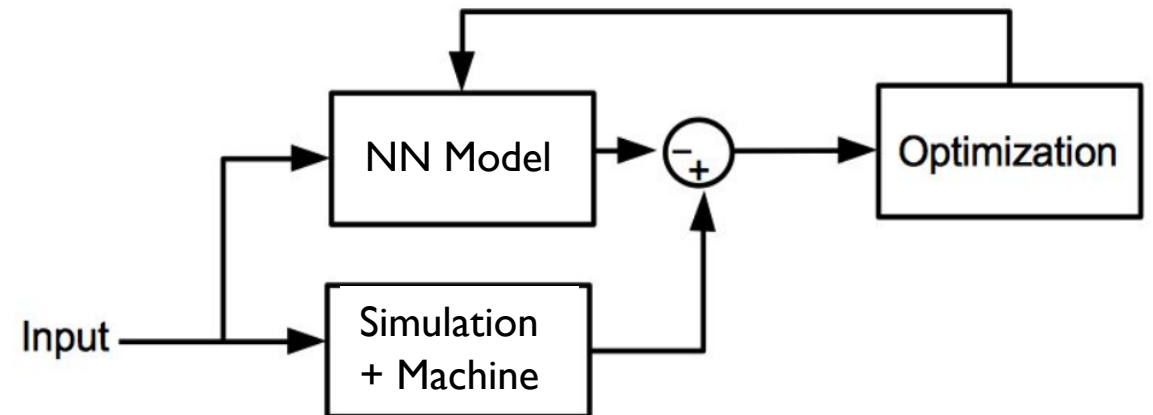  - *Bonus: control system development*

One approach: faster modeling codes

Simpler models (tradeoff with accuracy)

analytic calculations         *e. g.  J. Galambos, et al., HPPA5, 2007*

Parallelization and GPU-acceleration of existing codes

HPSim/PARMILA              *X. Pang, PAC13, MOPMA13*

*elegant*                          *I. V. Pogorelov, et al., IPAC15, MOPMA035*

Improvements to modeling algorithms

*Lorentz-boosted frame*        *J.-L. Vay, Phys. Rev. Lett. 98 (2007) 130405*

# Online Modeling

- Use a machine model during operation
.

- Ideally:
  - Fast-executing, but accurate enough to be useful
  - Use measured inputs directly from machine
  - Combine *a priori* knowledge + learned parameters

- Applications:
  - A tool for operators + virtual diagnostic
  - Predictive control
  - Help flag aberrant behavior
  - *Bonus: control system development*

Another approach: machine learning model

Once trained, neural networks can execute quickly

Train on data from slow, high-fidelity simulations
+
Train on measured data

# Online Modeling

- Use a machine model during operation

- Ideally:
  - Fast-executing, but accurate enough to be useful
  - Use measured inputs directly from machine
  - Combine *a priori* knowledge + learned parameters

- Applications:
  - A tool for operators + virtual diagnostic
  - Predictive control
  - Help flag aberrant behavior
  - *Bonus: control system development*

Another approach: machine learning model

Once trained, neural networks can execute quickly

Train on data from slow, high-fidelity simulations

\+

Train on measured data



*An initial study at Fermilab:*
*A. L. Edelen, et al. NAPAC16, TUPOA51*
*One PARMELA run with 2-D space charge: ~ 20 minutes*
*Neural network model: ~ a millisecond*

# Virtual Diagnostics

*Predict what diagnostics might look like when they are unavailable or don't exist*

fast-executing simulation

measured machine inputs → Online Model → Real-time prediction of beam characteristics or explicit diagnostic output

# Virtual Diagnostics

*Predict what diagnostics might look like when they are unavailable or don't exist*

*fast-executing simulation*

measured machine inputs → **Online Model** → Real-time prediction of beam characteristics or explicit diagnostic output



e.g. GPU-accelerated HPSim at LANSCE (based on PARMILA)

*X. Pang, et al., PAC13, MOPMA13*
*X. Pang, IPAC15, WEXC2*
*X. Pang and L. Rybarcyk, CPC185, is. 3 (2014)*
*L. Rybarcyk, et al., IPAC15, MOPWI033*
*L. Rybarcyk, HB2016, WEPM4Y01*

# Virtual Diagnostics

*Predict what diagnostics might look like when they are unavailable or don't exist*

fast-executing simulation

measured machine inputs → **Online Model** → Real-time prediction of beam characteristics or explicit diagnostic output

(ML model)

measured machine inputs → **Online Model** → diagnostic measurements

# Virtual Diagnostics

*Predict what diagnostics might look like when they are unavailable or don't exist*

# Virtual Diagnostics

*Predict what diagnostics might look like when they are unavailable or don't exist*

fast-executing simulation

measured machine inputs → Online Model → Real-time prediction of beam characteristics or explicit diagnostic output

(ML model)

measured machine inputs → Online Model → diagnostic measurements ✗

diagnostic prediction

- *moved to another location*
- *destructive, cannot always use*
- *blocked for update time*

# Virtual Diagnostics at Fermilab's FAST Facility



Multi-slit emittance measurement after the second capture cavity (X107 to X111) takes 10-15 seconds
→ *can we get an online prediction of what this intercepting diagnostic would show?*

# Initially limit the scope…



Solenoid Current

Phases (Gun, CC1, CC2)

Initial Bunch Properties
(charge, length, $\varepsilon_{x,y}$, *x-y corr.*)

Neural
Network

Transverse Sigma Matrix

Average Beam Energy

Transmission

$\varepsilon_{x,y} \quad \alpha_{x,y} \quad \beta_{x,y}$

# Could in principle use measured data alone, but want to be efficient with machine time

→ *use simulation data to fill out the training set*



gun phase scans
solenoid current scans
*(with two different laser intensities)*

beam   mask   screen

fit to obtain
subset of phase
space parameters

$\langle y^2 \rangle$
$\langle yy' \rangle$
$\langle y'^2 \rangle$

other setting
combinations

OPAL   elegant

cathode → CC2
*with 3-D space charge routine*

$+$   full sigma
matrix

# Training on imperfect simulations … NN only as good as the simulation

Poor agreement between simulation and measured data for some input/output relationships

→ *can we update the NN model with measured data without disrupting the other predictions?*
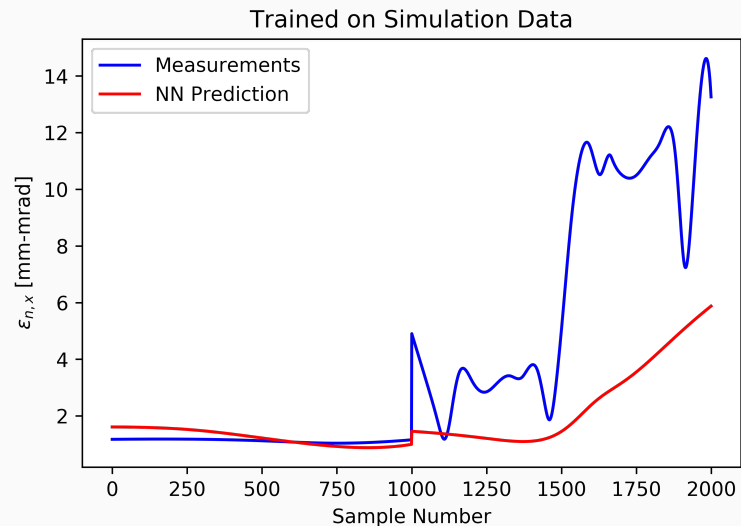
## Solenoid Scan

### Simulation Data Only

**Trained on Simulation Data**

**Trained on Simulation Data**

### Updated with Measured Data

**Trained on Simulation + Measured Data**

**Trained on Simulation + Measured Data**

## Phase Scan

**Trained on Simulation**

**Trained on Simulation + Measured Data**

Solenoid Scan

Phase Scan

**Simulation Data Only**

Trained on Simulation Data

Trained on Simulation Data

Trained on Simulation

**Updated with Measured Data**

Trained on Simulation + Measured Data
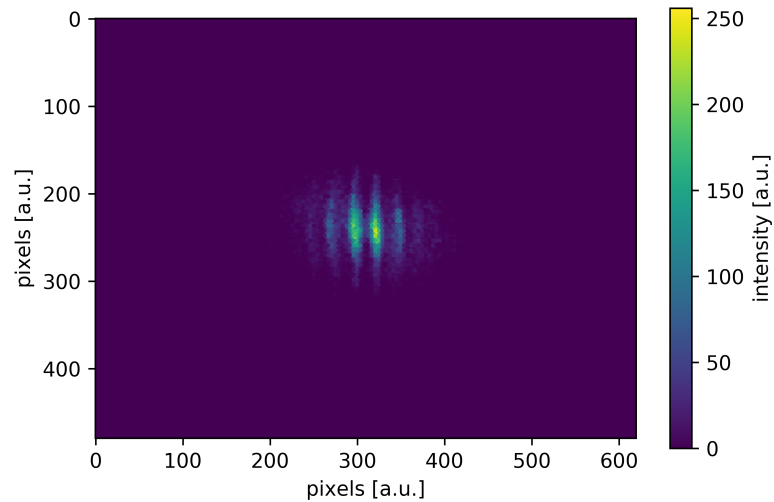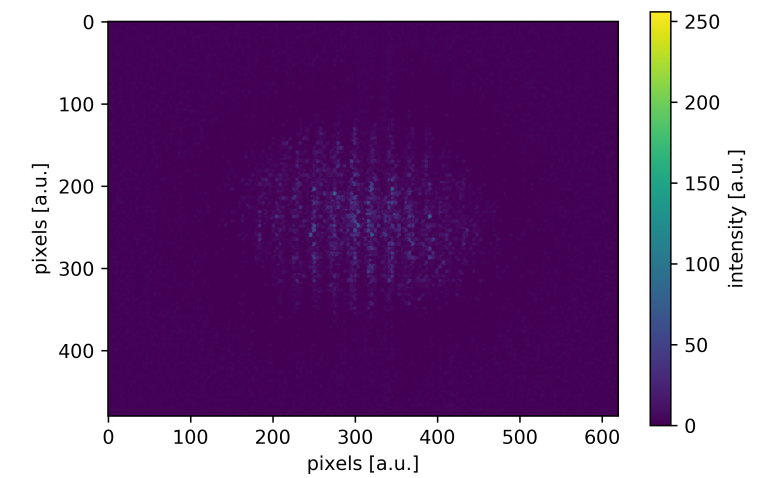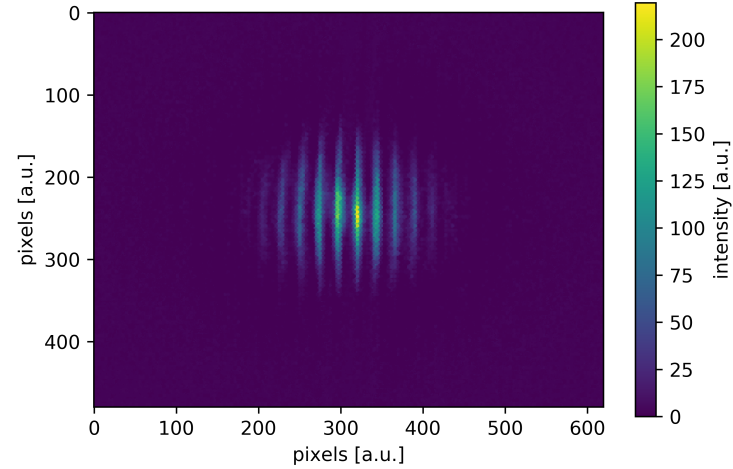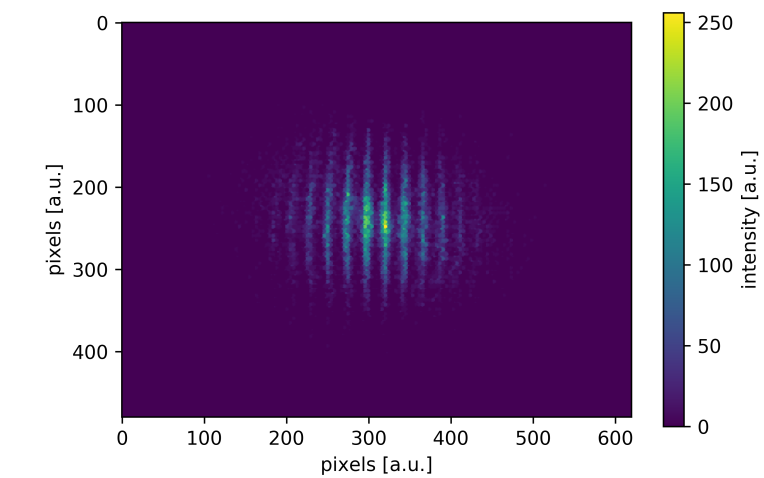
Trained on Simulation + Measured Data

Trained on Simulation + Measured Data

*Why bother with simulation at all? → Rough initial solution facilitates training with small amount of measured data*

# Predicting Image Output Directly

Simulated

NN Predictions

Difference

*Work with J. Edelen, D. Edstrom, A. Halavanau, J. Ruan, P. Piot, A. Romanov, S.G. Biedron, S.V. Milton*

# Bigger Picture

Fast-executing, accurate machine model

*Online:*     *facilitate studies*

*Offline:*     *study planning*
*downstream component design*
*controller training*

*Earlier work: account for changes in laser spot*
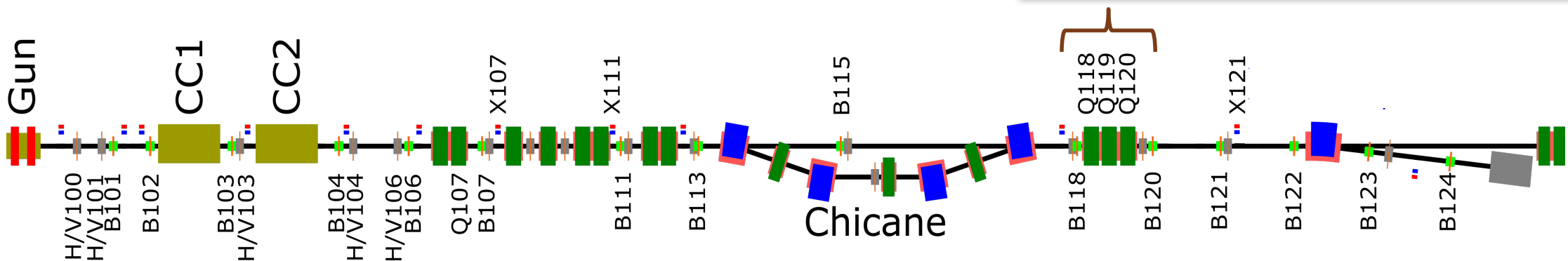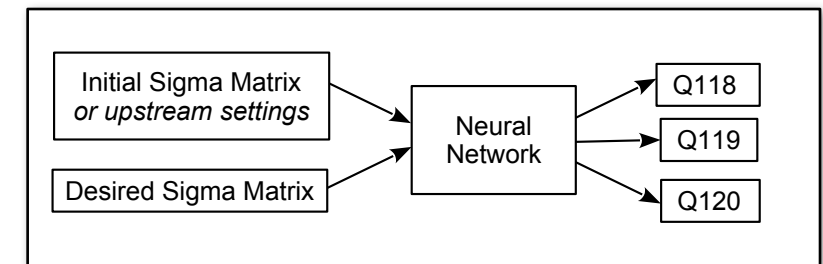*A. L. Edelen, et al. NAPAC16, TUPOA51*



One piece of a larger set of studies:

- *Accounting for laser spot changes*
- *NN controller (starting with round-to-flat beam transform)*
- ***The vision is to combine these***

*Ongoing work: NN-based round-to-flat beam transform*

# Fast Switching Between Trajectories

*Work with C. Tennant and D. Douglas, JLab*

- 76 BPMs, 57 dipoles, 53 quadrupoles
- Traditional approach has never worked (linear response matrix)
- Rely on a few experts for steering tune-up
- Want to specify small offsets in trajectory at some locations
- Didn't initially have an up-to-date machine model available

Learn responses (NN model) from tune-up data and dedicated study time:
dipole + quadrupole settings → predict BPMs + transmission

Train controller (NN policy) offline using NN model:
desired trajectory → dipole settings
(and penalize losses + large magnet settings)



JLab

# Fast Switching Between Trajectories

Main anticipated advantage of NN over standard approach:

*Adaptive control policy → adjust without interfering with operation for response measurements as often?*

*Handling of trajectories away from BPM center (nonlinear)*
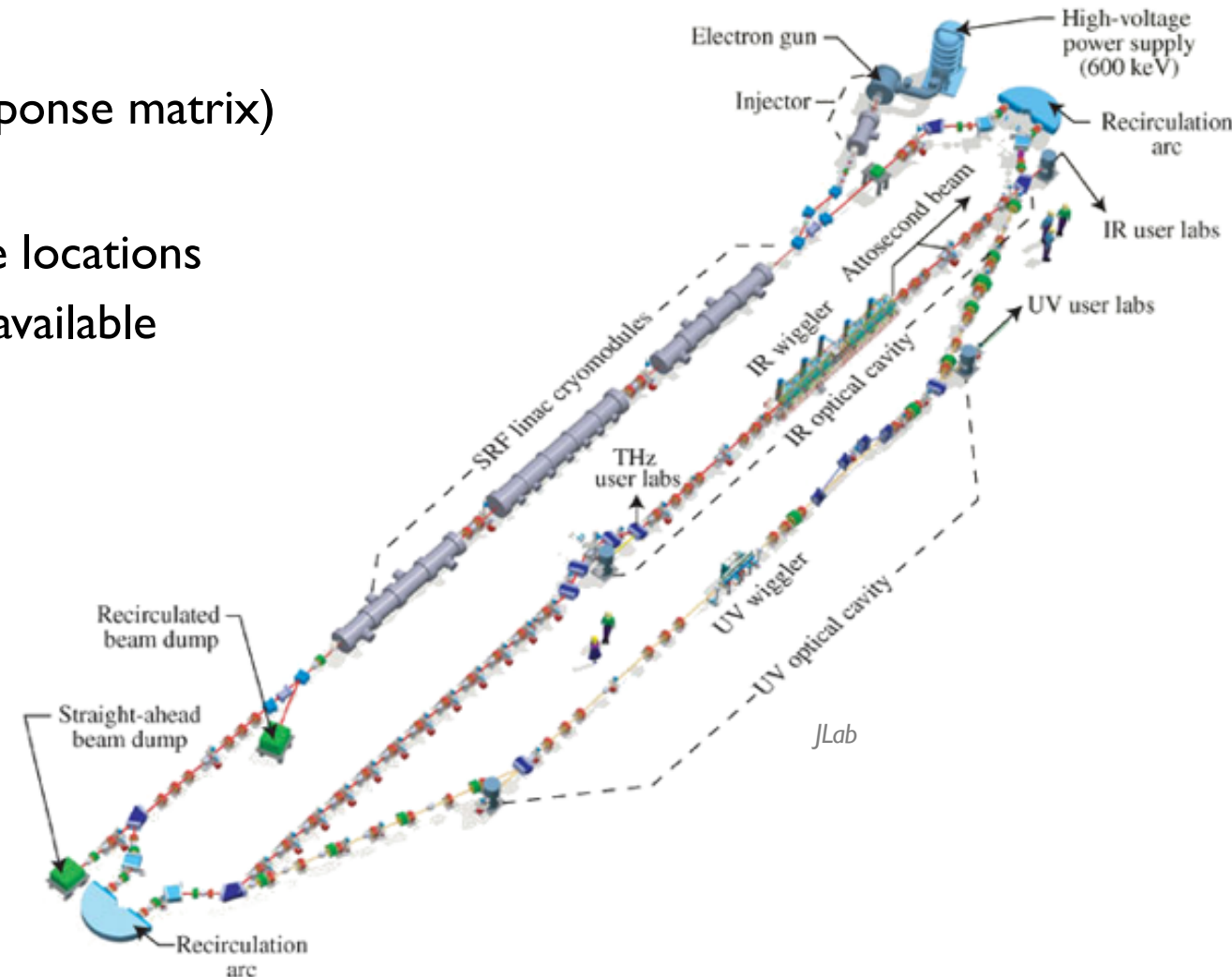
*But, need to quantify this …*

Learn responses (NN model) from tune-up data and dedicated study time:
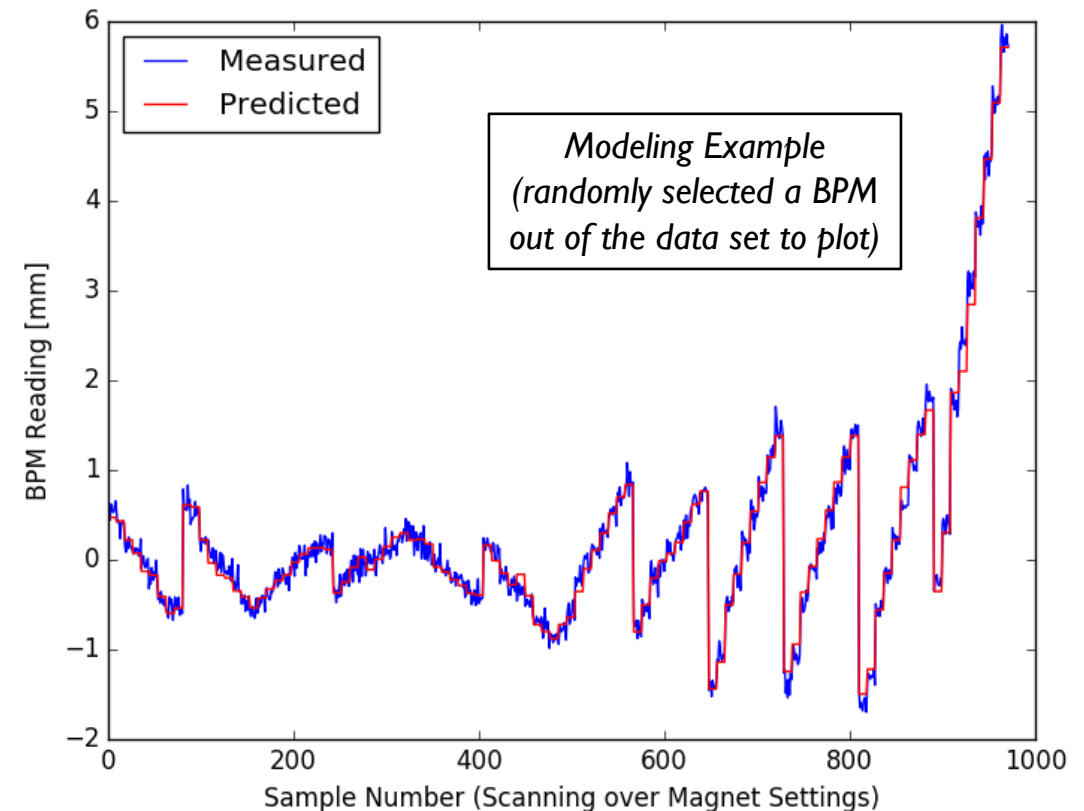dipole + quadrupole settings → predict BPMs + transmission

Train controller (NN policy) offline using NN model:
desired trajectory → dipole settings
(and penalize losses + large magnet settings)

Preliminary Results:

*Model Errors for BPMs:*

| | | |
|---|---|---|
| Training Set: | 0.07 mm MAE | 0.09 mm STD |
| Validation Set: | 0.08 mm MAE | 0.07 mm STD |
| Test Set: | 0.08 mm MAE | 0.03 mm STD |

*Controller*: random initial states → on average within 0.2 mm of center immediately



Modeling Example
(randomly selected a BPM out of the data set to plot)

# Switching Between User Requests

- FEL facilities support a wide variety of scientific endeavors (e.g. imaging protein structures[1], understanding processes like photosynthesis[2], origin of material properties[3])

- Need to accommodate requests for a wide variety of photon beam characteristics

- May switch as often as every few days

- Have save/restore settings, but these are discrete, and there can be some drift in the machine

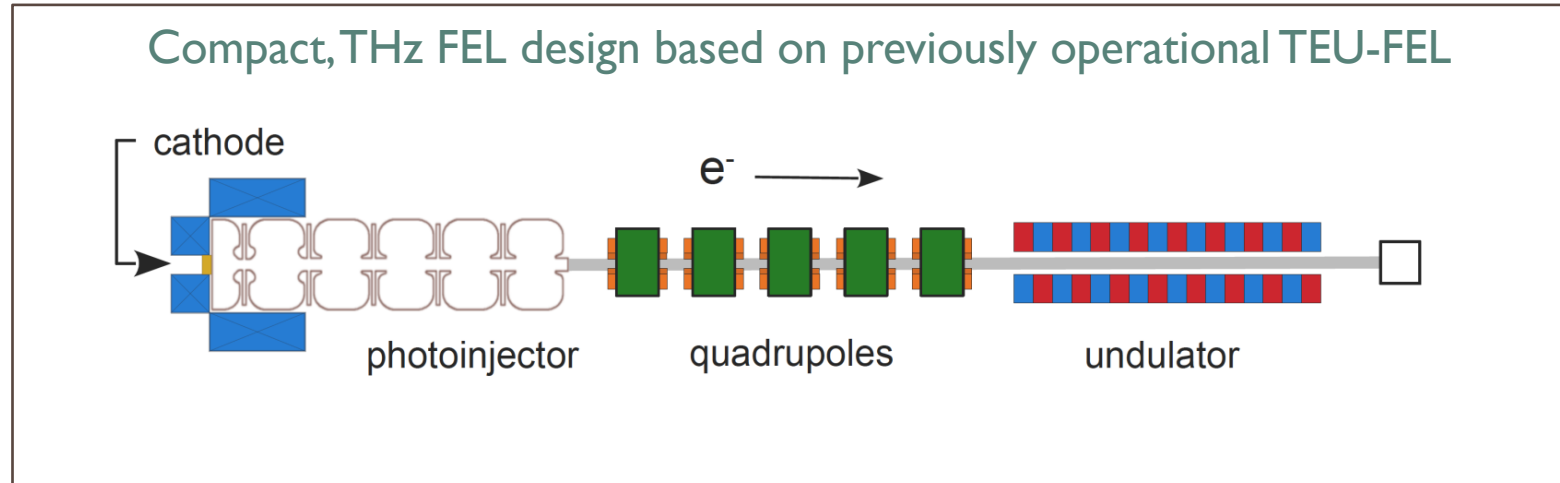- Time spent tuning = reduced scientific output for a given operational budget

*Would be nice to have a tool that can quickly give suggested settings for a given photon beam request, is valid globally, and can adapt to changes over time*



*e.g. the Linac Coherent Light Source*
(image: lcls.slac.standford.edu)

[1] J.-P. Colletier, et al.,"De novo phasing with X-ray laser reveals mosquito larvicide BinAB structure," Nature , vol. 539, pp. 43–47, Sep. 2016.
[2] I. D. Young, et al., "Structure of photosystem II and substrate binding at room temperature," Nature , vol. 540, pp. 453–457, Nov. 2016.
[3] M. P. Jiang, et al., "The origin of incipient ferroelectricity in lead telluride," Nature Communications, vol. 7, no. 12291, Jul. 2016.
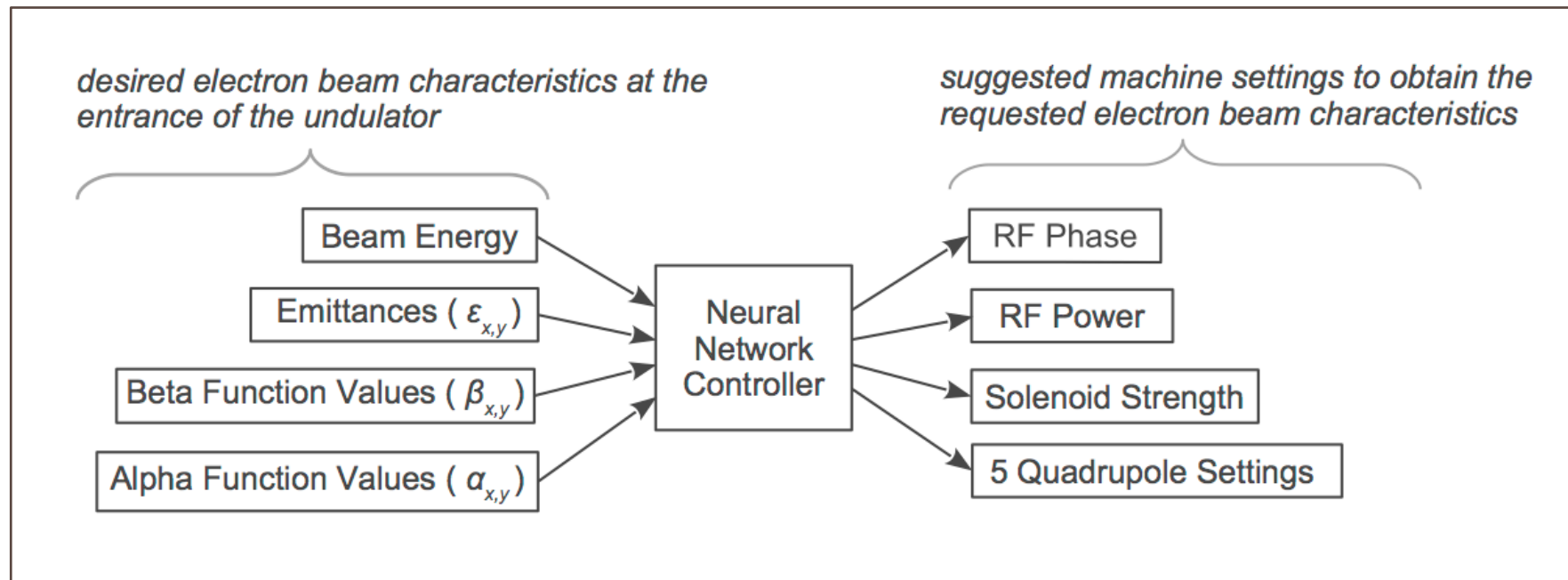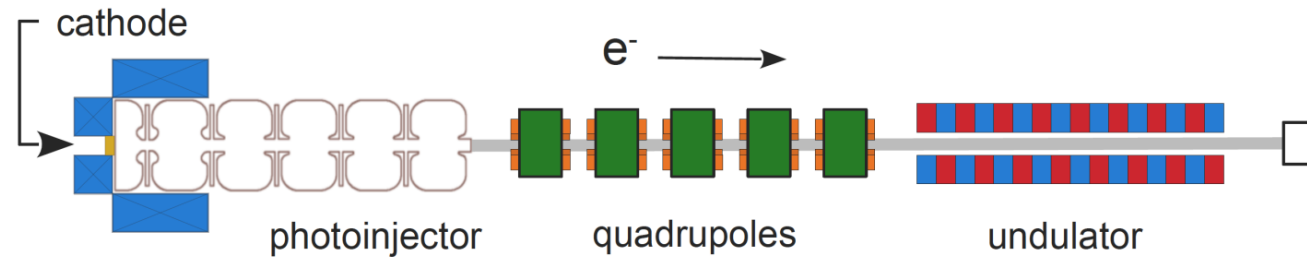
# Starting Smaller: A Case Study



Compact, THz FEL design based on previously operational TEU-FEL

cathode

e⁻ →

photoinjector    quadrupoles    undulator

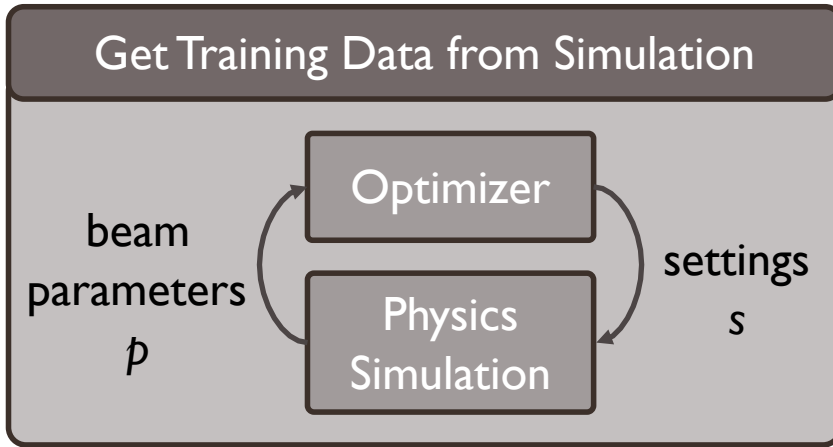3 – 6 MeV electron beam
200 – 800 $\mu$m photon beam

Previously operated at University of Twente in the Netherlands

Was going to be re-built at CSU

This is an appealing system for an initial study because it has a small number of machine components, yet it exhibits non-trivial beam dynamics.

# *Intermediate goal: get the right beam parameters at the undulator entrance*

## Get Training Data from Simulation

Optimizer

Physics
Simulation

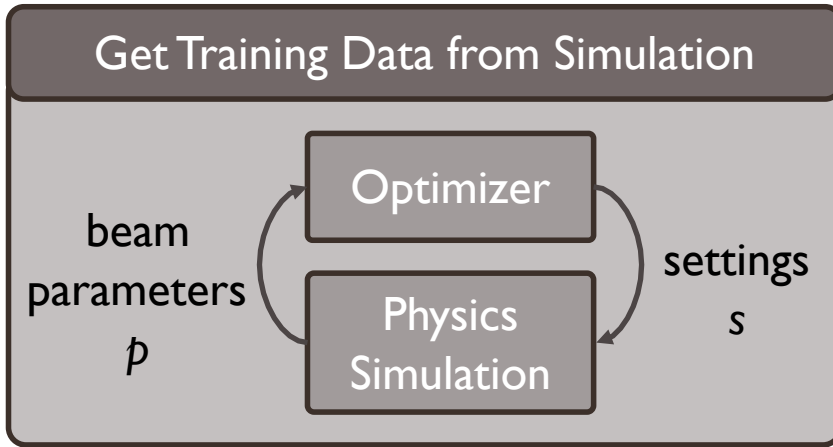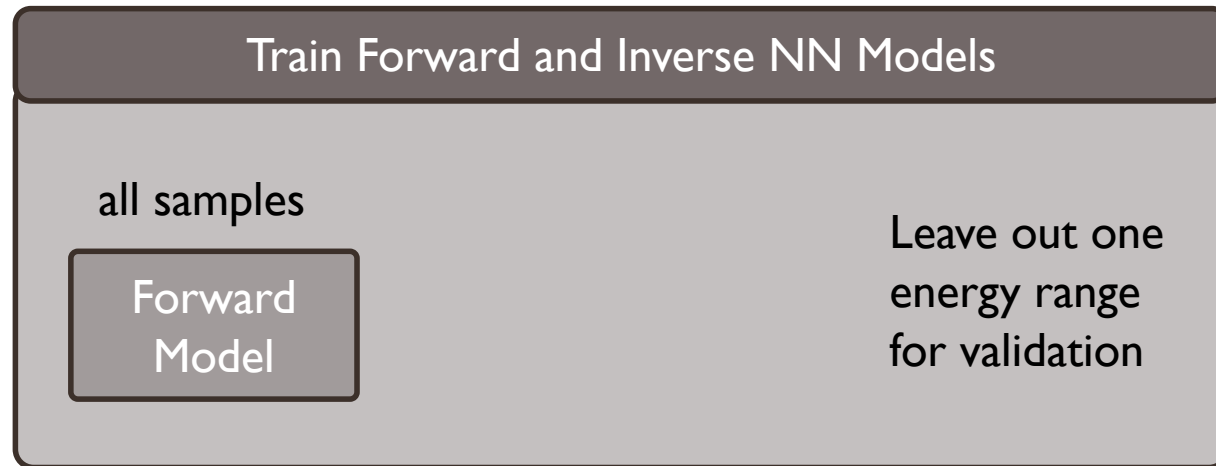beam
parameters
$p$

settings
$s$

repeat for different target energies

Don't always have a good physics-based model for particle accelerators, so what's in the data archive of a real facility?

*Noisy data + tuning around roughly optimal settings*

## Get Training Data from Simulation

beam parameters $p$

Optimizer

Physics Simulation

settings $s$

repeat for different target energies

## Train Forward and Inverse NN Models

all samples

Forward Model

Leave out one energy range for validation

Don't always have a good physics-based model for particle accelerators, so what's in the data archive of a real facility?

*Noisy data + tuning around roughly optimal settings*

## Get Training Data from Simulation

Optimizer

Physics Simulation

beam parameters $p$

settings $s$

repeat for different target energies

## Train Forward and Inverse NN Models

all samples

Forward Model

Leave out one energy range for validation

Don't always have a good physics-based model for particle accelerators, so what's in the data archive of a real facility?

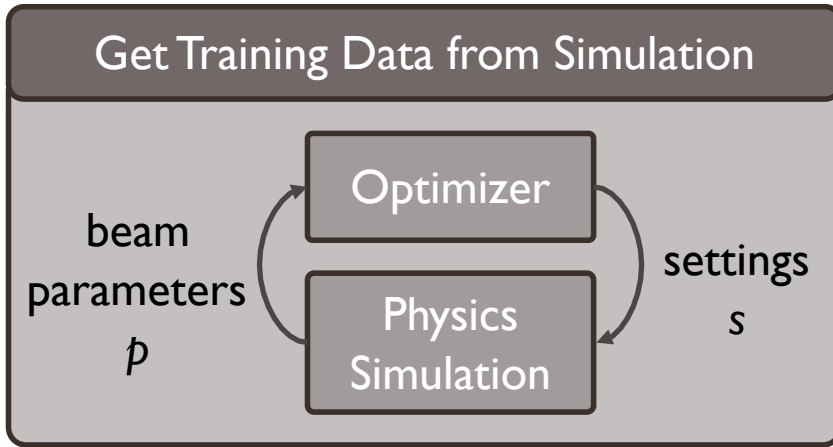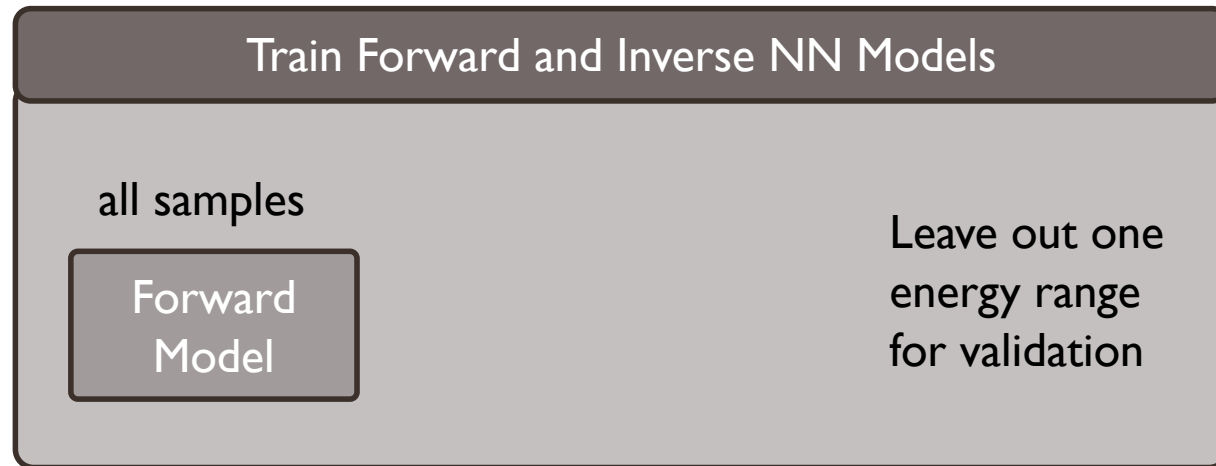*Noisy data + tuning around roughly optimal settings*

Want to use the existing data to initialize control policy

**Get Training Data from Simulation**

Optimizer

beam parameters $p$

settings $s$

Physics Simulation

repeat for different target energies

**Train Forward and Inverse NN Models**

all samples

converged samples (optimal settings)

Forward Model

Inverse Model

Leave out one energy range for validation

Initial Policy

Don't always have a good physics-based model for particle accelerators, so what's in the data archive of a real facility?
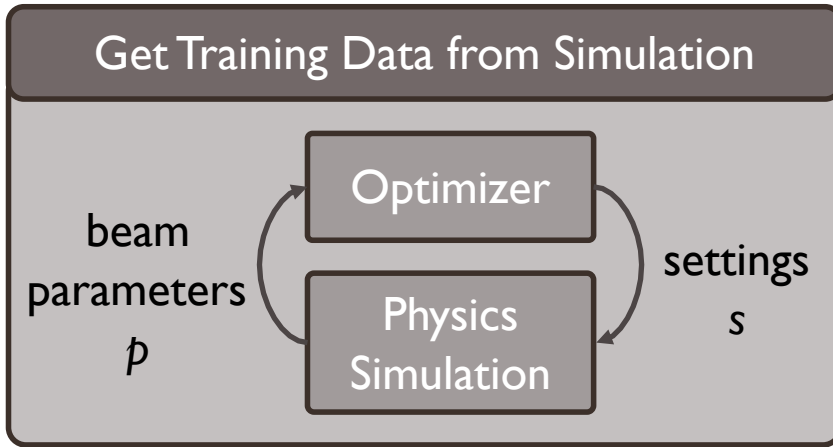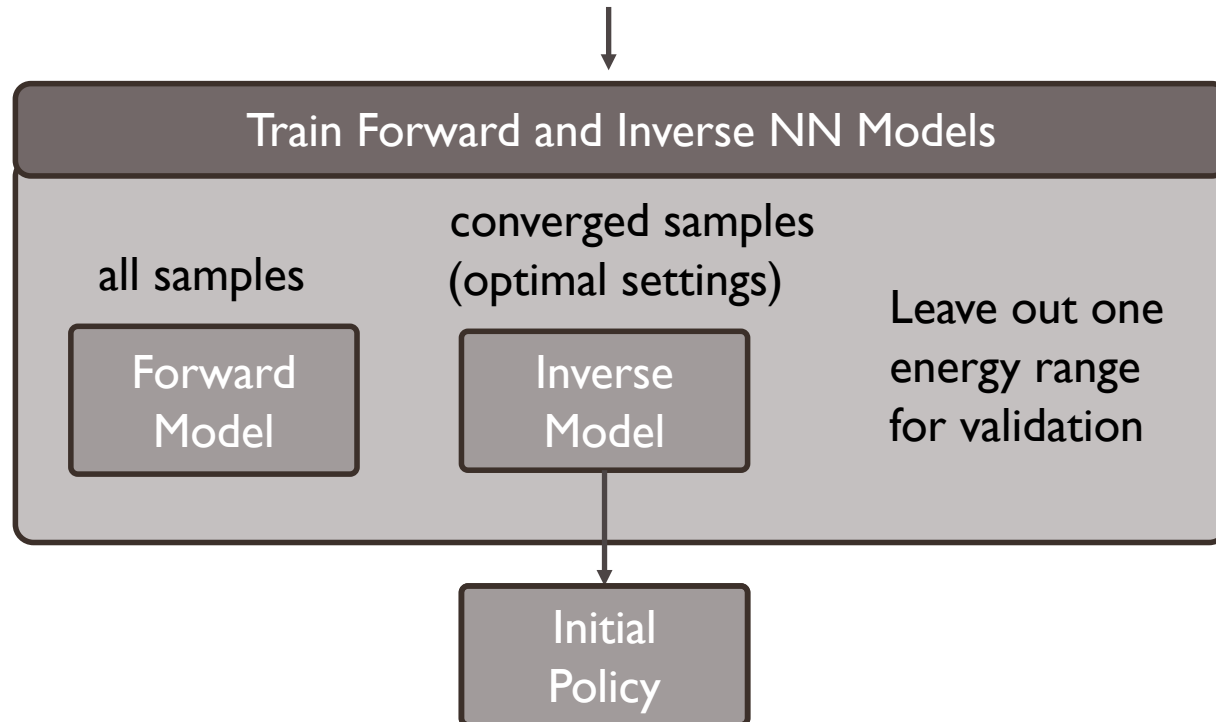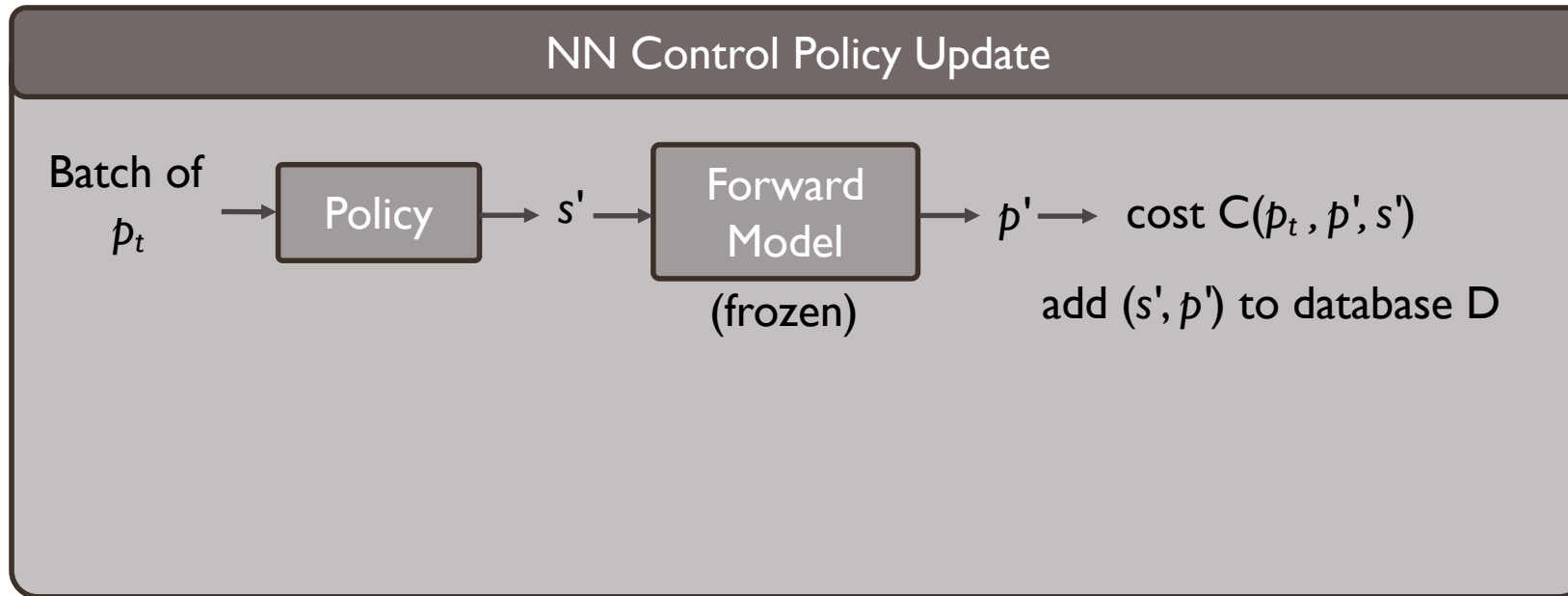
*Noisy data + tuning around roughly optimal settings*

Want to use the existing data to initialize control policy → *model not invertible, but can pre-train policy with converged settings*

# Training the Control Policy

- *First: just want to switch to roughly correct settings*
- *Then, two options: efficient local tuning algorithms we already use, or online model/controller updating*



$p_t$ -- *target beam parameters*

$s'$ - *predicted optimal settings*
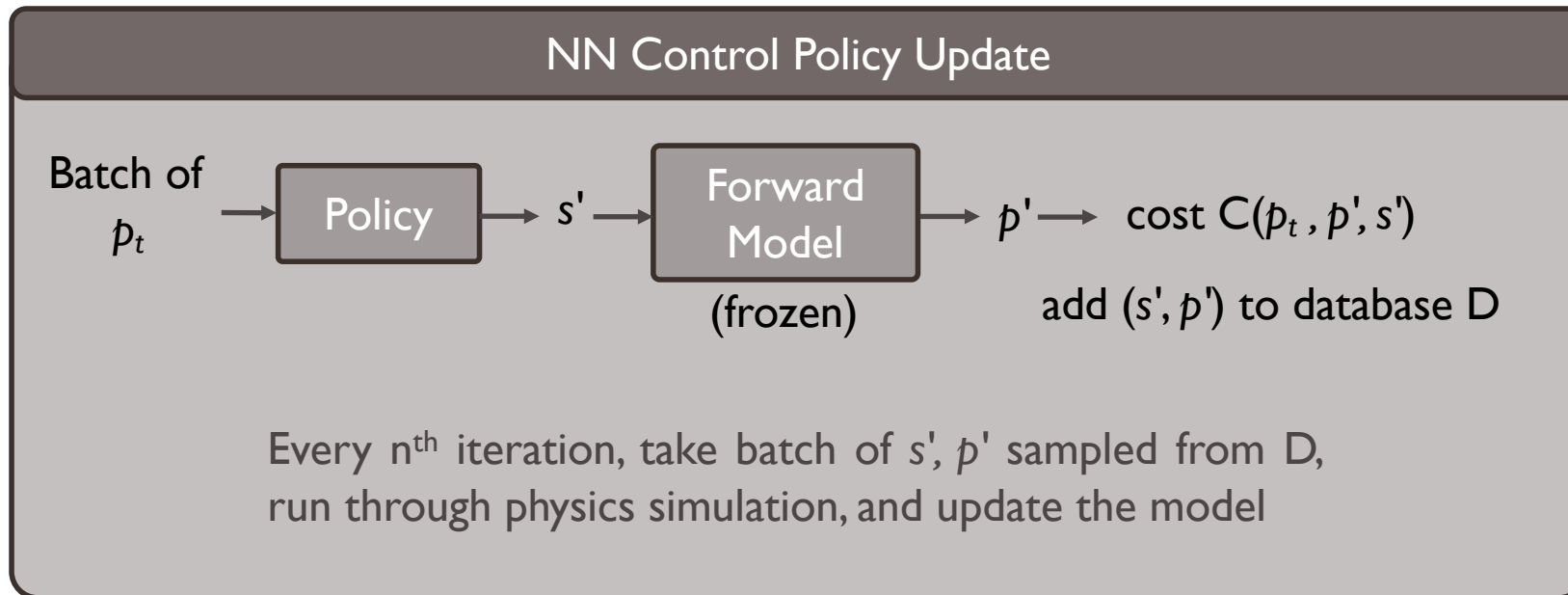
$p'$ - *predicted beam parameters*

*Cost:*
*difference between $p'$ and $p_t$*
*penalize loss of transmission*
*penalize higher magnet settings*

# Training the Control Policy

- *First: just want to switch to roughly correct settings*
- *Then, two options: efficient local tuning algorithms we already use, or online model/controller updating*

**NN Control Policy Update**

Batch of $p_t$ → **Policy** → $s'$ → **Forward Model** (frozen) → $p'$ → cost $C(p_t, p', s')$

add $(s', p')$ to database D

Every $n^{th}$ iteration, take batch of $s'$, $p'$ sampled from D, run through physics simulation, and update the model

*$p_t$ -- target beam parameters*

*$s'$ – predicted optimal settings*

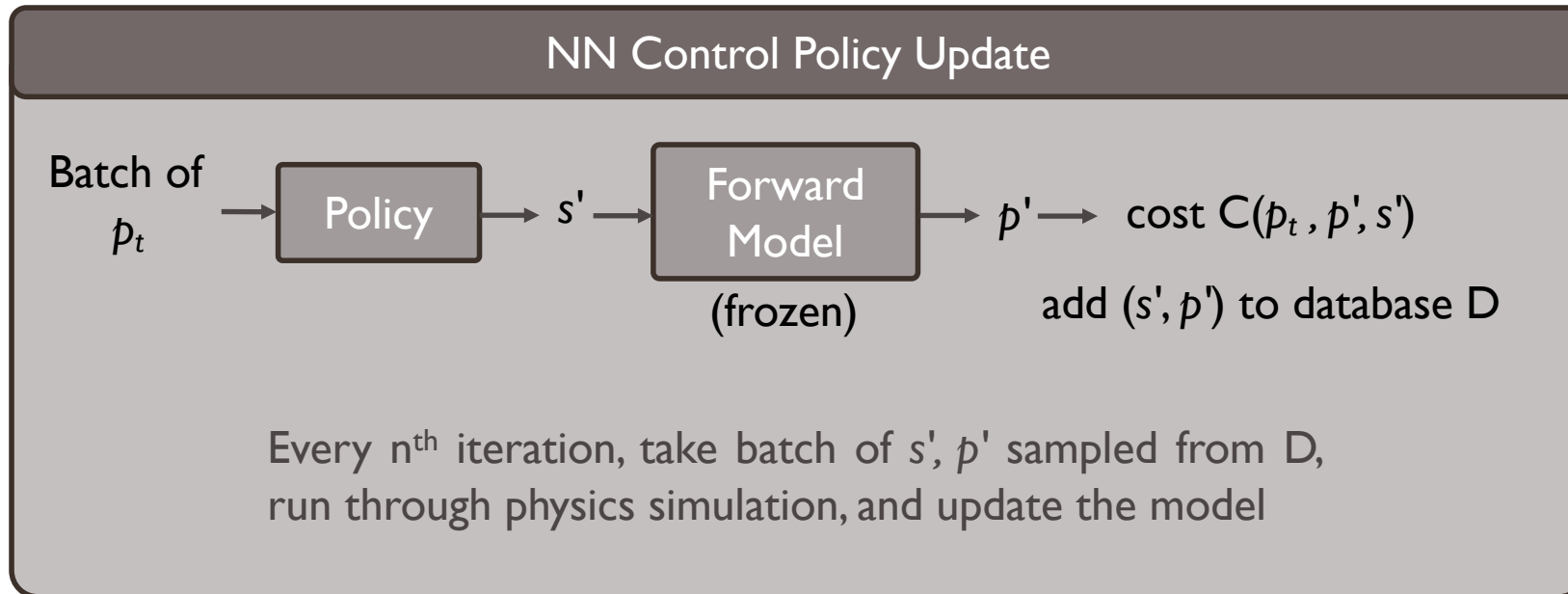*$p'$ – predicted beam parameters*

*Cost:*
*difference between $p'$ and $p_t$*
*penalize loss of transmission*
*penalize higher magnet settings*

# Training the Control Policy

- *First: just want to switch to roughly correct settings*
- *Then, two options: efficient local tuning algorithms we already use, or online model/controller updating*



NN Control Policy Update

Batch of $p_t$ → Policy → $s'$ → Forward Model (frozen) → $p'$ → cost $C(p_t, p', s')$

add $(s', p')$ to database D

Every $n^{th}$ iteration, take batch of $s'$, $p'$ sampled from D, run through physics simulation, and update the model

$p_t$ -- *target beam parameters*

$s'$ -- *predicted optimal settings*

$p'$ -- *predicted beam parameters*

*Cost:*
*difference between $p'$ and $p_t$*
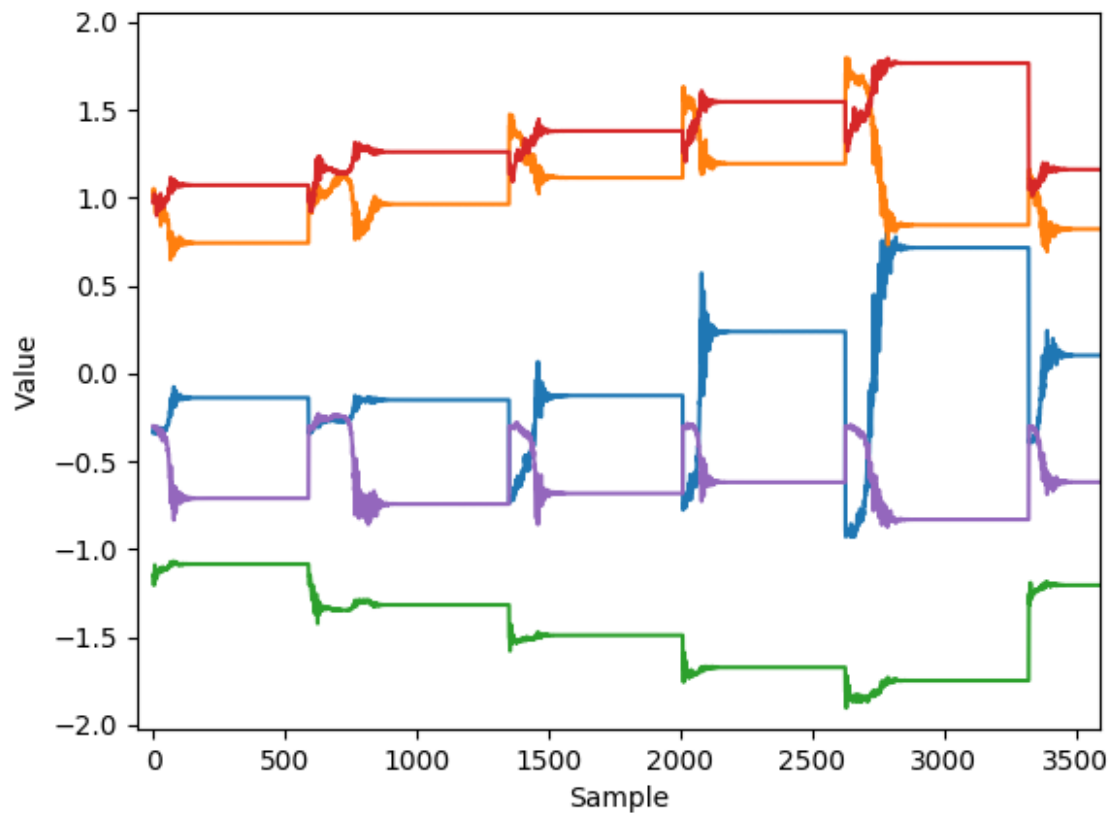*penalize loss of transmission*
*penalize higher magnet settings*

Then test policy directly on simulation

# Initial Model and Controller

Training data from simulation:
- output from each iteration of Nelder-Mead, L-BFGS
- 12 beam energies between 3.1 – 6.2 MeV (7195 samples)



*Example of what the training data looks like*
*(quadrupoles shown in this case)*

Model: 50-50-30-30 tanh nodes in hidden layers
- 8 inputs *(rf power, rf phase, sol. strength, quads)*
- 8 outputs $(\alpha_x, \alpha_y, \beta_x, \beta_y, \varepsilon_x, \varepsilon_y, E, N_p)$
- 5.7-MeV run used for validation set

First study: focus on target $\alpha, \beta$ for a given energy

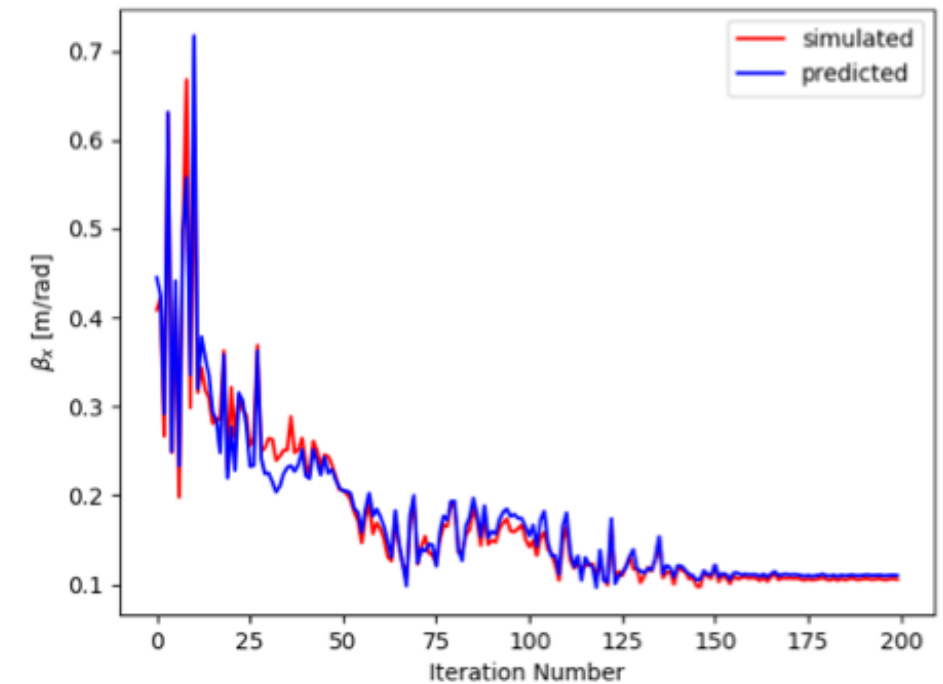Policy: 30-30-20-20 tanh nodes in hidden layers
- inputs/outputs opposite the above (except $N_p$)
- random target energies, $\alpha_{xy} = 0, \beta_{xy} = 0.106$
- exclude 4.8 – 5.2 MeV range for validation

# Initial Model and Controller Performance

## Summary of Model Performance

| Parameter | Train MAE | Train STD | Train Max | Val. MAE | Val. STD | Val. Max |
|-----------|-----------|-----------|-----------|----------|----------|----------|
| $\alpha_x$ [rad] | 0.018 | 0.042 | 0.590 | 0.067 | 0.091 | 0.482 |
| $\alpha_y$ [rad] | 0.022 | 0.037 | 0.845 | 0.070 | 0.079 | 0.345 |
| $\beta_x$ [m/rad] | 0.004 | 0.009 | 0.287 | 0.008 | 0.012 | 0.130 |
| $\beta_y$ [m/rad] | 0.005 | 0.011 | 0.357 | 0.012 | 0.017 | 0.189 |

## Example of Model Performance

# Initial Model and Controller Performance

### Summary of Model Performance

| Parameter | Train MAE | Train STD | Train Max | Val. MAE | Val. STD | Val. Max |
|---|---|---|---|---|---|---|
| $\alpha_x$ [rad] | 0.018 | 0.042 | 0.590 | 0.067 | 0.091 | 0.482 |
| $\alpha_y$ [rad] | 0.022 | 0.037 | 0.845 | 0.070 | 0.079 | 0.345 |
| $\beta_x$ [m/rad] | 0.004 | 0.009 | 0.287 | 0.008 | 0.012 | 0.130 |
| $\beta_y$ [m/rad] | 0.005 | 0.011 | 0.357 | 0.012 | 0.017 | 0.189 |

Controller ability to reach $\alpha_{x,y}$ = 0 and $\beta_{x,y}$ = 0.106 in **one iteration**

| Parameter | Train MAE | Train STD | Train Max | Val. MAE | Val. STD | Val. Max |
|---|---|---|---|---|---|---|
| $\alpha_x$ [rad] | 0.012 | 0.075 | 0.011 | 0.046 | 0.063 | 0.141 |
| $\alpha_y$ [rad] | 0.013 | 0.079 | 0.012 | 0.045 | 0.064 | 0.140 |
| $\beta_x$ [m/rad] | 0.008 | 0.004 | 0.006 | 0.006 | 0.023 | 0.008 |
| $\beta_y$ [m/rad] | 0.014 | 0.011 | 0.011 | 0.011 | 0.069 | 0.038 |

### Example of Model Performance



*What this means: for a given energy, the controller will immediately reach the desired beam size to within about 10% and the beam will be close to a waist, requiring minimal further tuning (assuming no substantial drift…)*

# Dealing with "Long-Term" Time Dependencies:
# Resonant Frequency Control in Normal Conducting Cavities

*RF electron gun at the Fermilab Accelerator Science and Technology (FAST) facility*
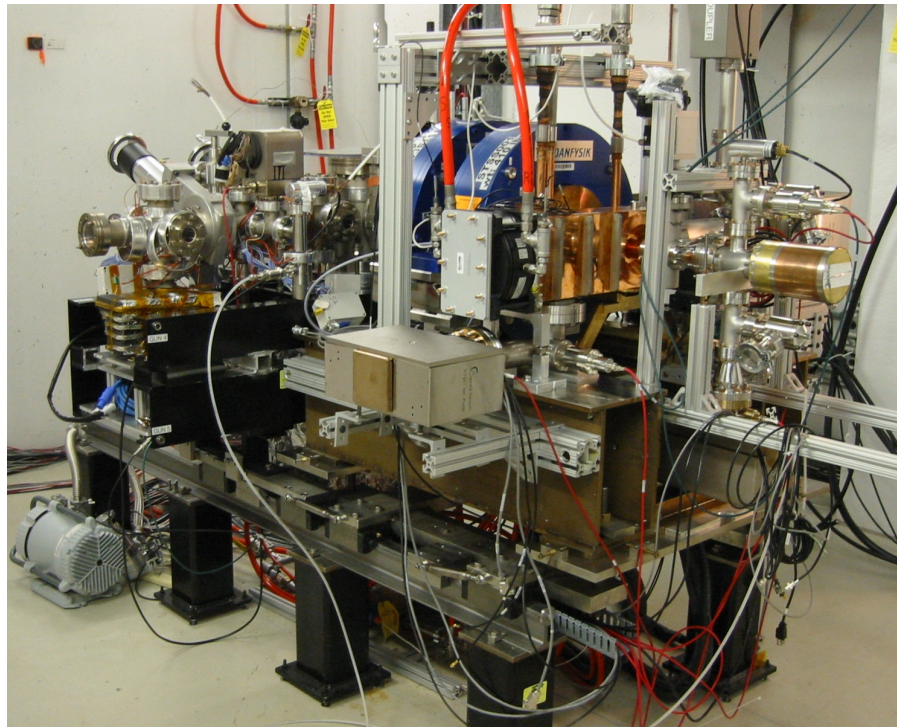
*Radio frequency quadrupole (RFQ) for the PIP-II Injector Test*
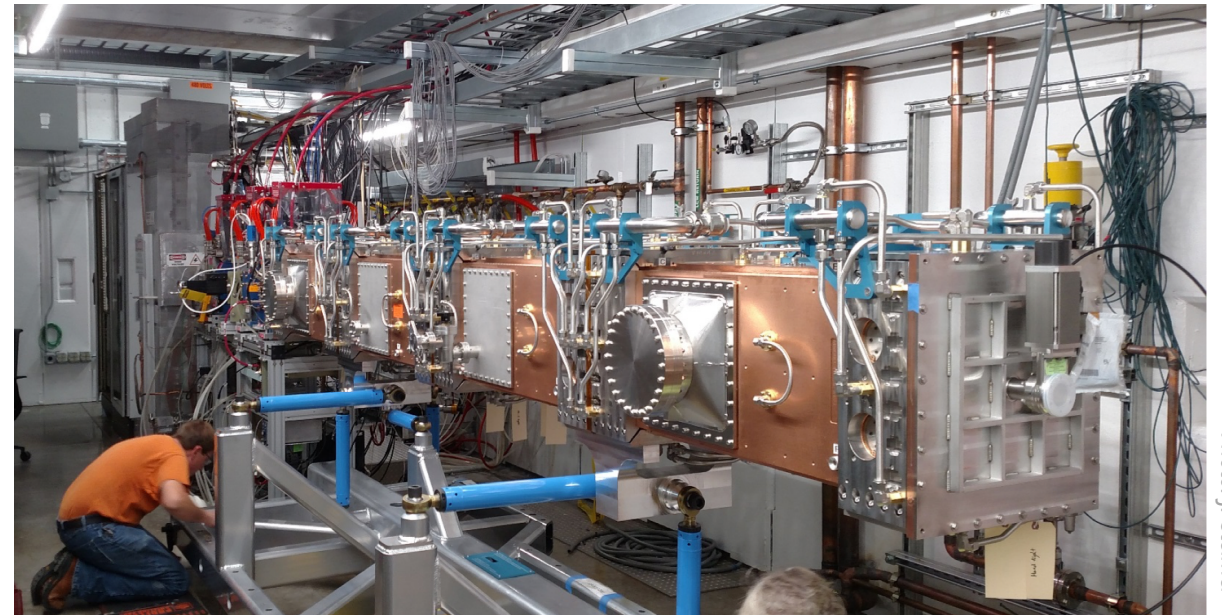


Photo: P. Stabile



Photo: J. Steimel

*"long term" in this case means responses lasting many minutes (e.g. 30), with control actions at 0.5 Hz and 1 Hz*
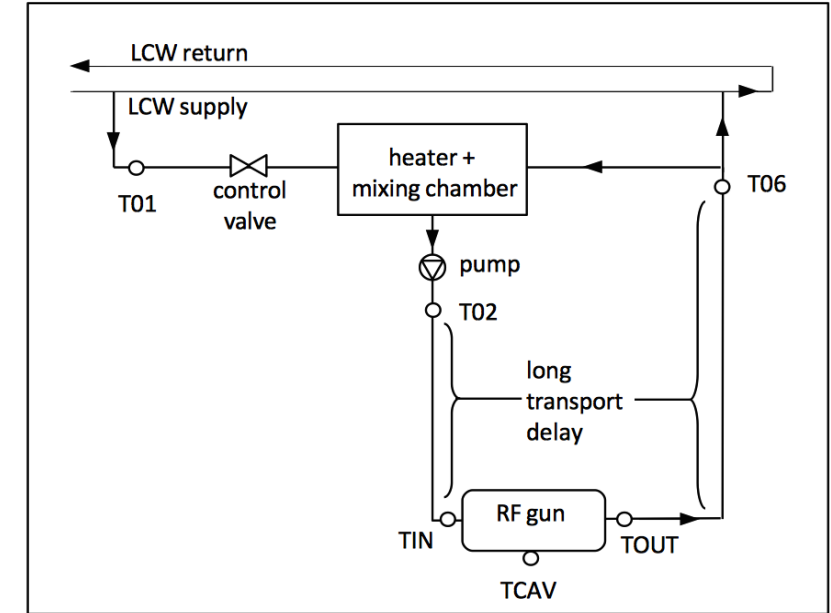
# Temperature Control for the RF Photoinjector at FAST

Resonant frequency controlled via temperature
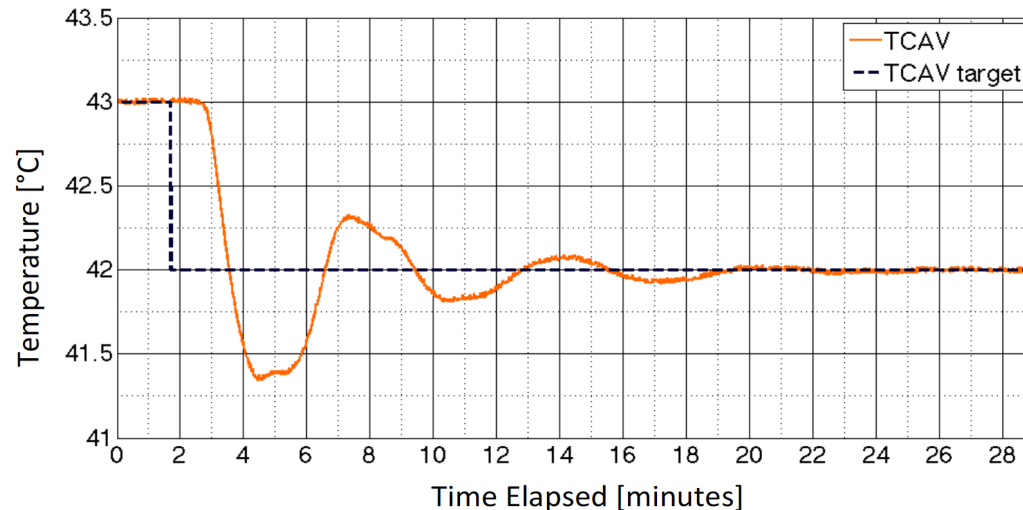
PID control is undesirable in this case:
- Long transport delays and thermal responses
- Recirculation leads to secondary impact of disturbances
- Two controllable variables: heater power + valve aperture

Applied model predictive control (MPC) with a neural network model trained on measured data:  ~ 5x faster settling time + no large overshoot
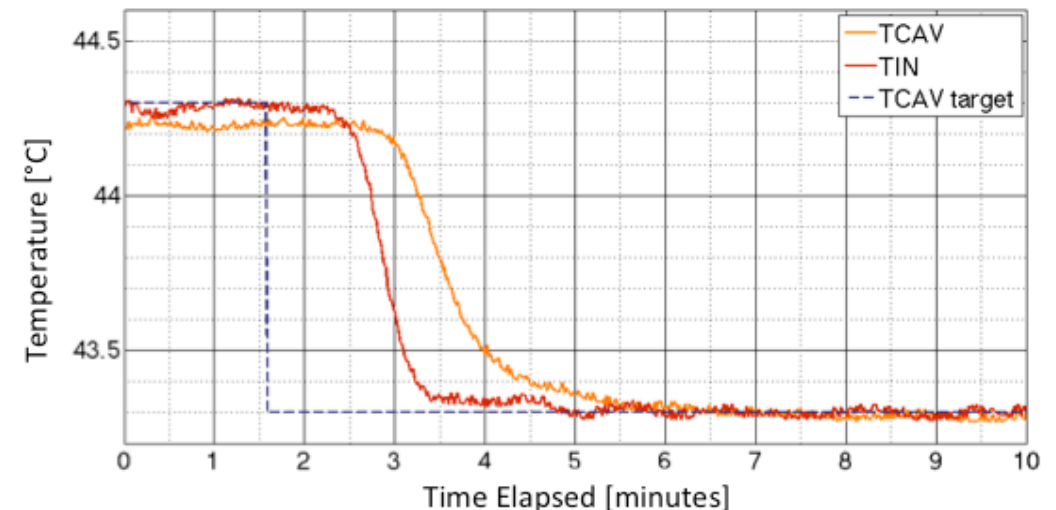
Gun Water System Layout



Existing Feedforward/PID Controller



Model Predictive Controller



*Note that the oscillations are largely due to the transport delays and water recirculation, rather than PID gains*
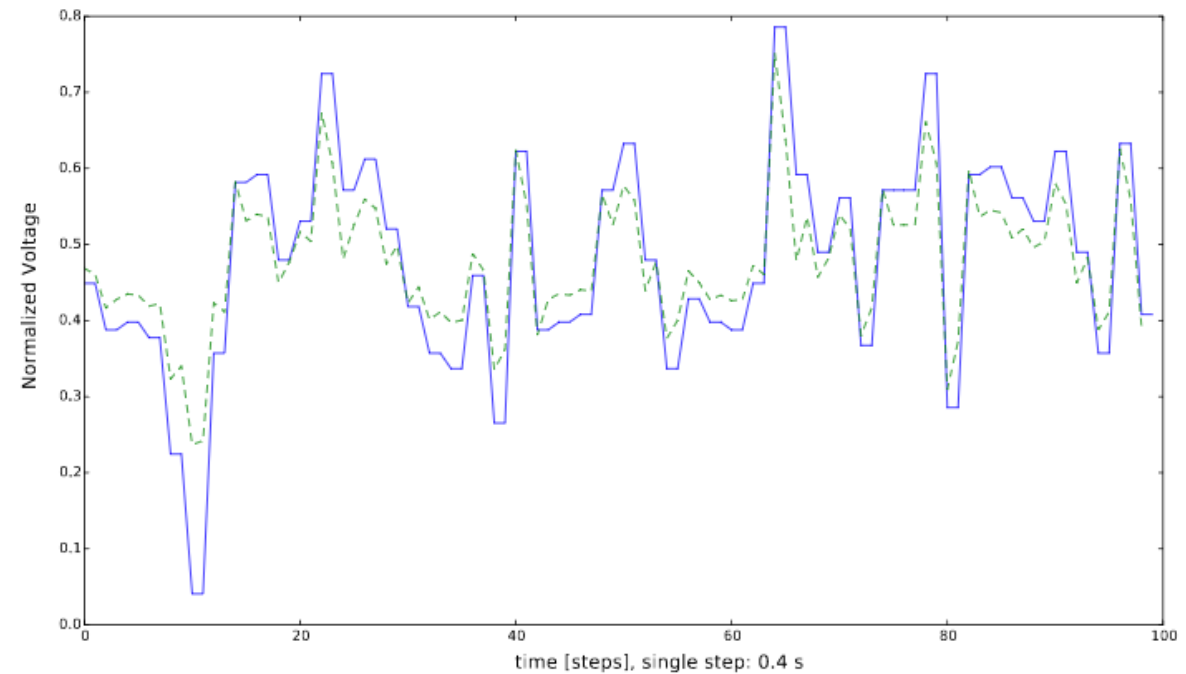
*"Some of the most dangerous malfunctions of the magnets are quenches which occur when a part of the superconducting cable becomes normally-conducting."*

**Aim:** **use a recurrent NN to identify quench precursors in voltage time series**

→ *Predict future behavior, then classify it*

Initial study with small data set:
- 425 quenches for 600 A magnets
- Used archived data from 2008 to 2016
- 16-32 previous values → predict a few time steps ahead

# Some Practical Challenges

*Need a sufficient\* amount of reliable\* data*

## Training on Measured Data

Undocumented manual changes
(e.g. rotating a BPM)

Relevant-but-unlogged variables

Availability of diagnostics

Observed parameter range in archived data

Time on machine for characterization studies
(schedule + expense)

*Ideal case:*
*- comprehensive, high-resolution data archive*
  *(e.g. including things like ambient temp./pressure)*
*- excellent log of manual changes*

## Training on Simulation Data

How representative of the real machine behavior?

Input/output parameters need to translate directly to what's on the machine (quantitatively)

High-fidelity (e.g. PIC) → time-consuming to run

Retention + availability of prior results:
(*optimize and throw the iterations away!*)

## Deployment

Initial training is on HPC systems → deployment is typically not\*
- Execution on front-end: necessary speed + memory?
- Subsequent training: on front-end or transfer to HPC?

Software compatibility for older systems:
interface with machine + make use of modern ML software libraries

I/O for large amounts of data

\* for now…

# Final Notes

- Neural networks are *very flexible tools* → *far more powerful + accessible in recent years*
- Lots of opportunities to use neural networks (and ML more broadly) to improve accelerator performance on both existing and future machines
- Transferrable between machines to some degree → *lots of potential for fruitful collaborations!*
- But, not a panacea!
  - *Simpler model-independent online optimization + simpler model-based approaches in many cases may be more appropriate*
  - *Boundaries of usefulness/reliability and tradeoff with time investment have yet to be determined rigorously*

# Final Notes

- Neural networks are very flexible tools → *far more powerful + accessible in recent years*

- Lots of opportunities to use neural networks (and ML more broadly) to improve accelerator performance on both existing and future machines

- Transferrable between machines to some degree → *lots of potential for fruitful collaborations!*

- But, not a panacea!
  - *Simpler model-independent online optimization + simpler model-based approaches in many cases may be more appropriate*
  - *Boundaries of usefulness/reliability and tradeoff with time investment have yet to be determined rigorously*

- Growing community → *two very recent workshops on ML for accelerators*

Machine Learning for Particle Accelerators
27 February – 2 March at SLAC
*Agenda/Talks: https://tinyurl.com/y988njbl*

Intelligent Controls for Particle Accelerators
30 – 31 January at Daresbury Lab
*Agenda/Talks: https://tinyurl.com/y9rg3uht*

# Final Notes

- Neural networks are very flexible tools → *far more powerful + accessible in recent years*

- Lots of opportunities to use neural networks (and ML more broadly) to improve accelerator performance on both existing and future machines

- Transferrable between machines to some degree → *lots of potential for fruitful collaborations!*

- But, not a panacea!
  - *Simpler model-independent online optimization + simpler model-based approaches in many cases may be more appropriate*
  - *Boundaries of usefulness/reliability and tradeoff with time investment have yet to be determined rigorously*

- Growing community → *two very recent workshops on ML for accelerators*

## Thanks for your attention!