# Interpretable Machine Learning at the European XFEL

Danilo Ferreira de Lima and *many many others* from the European XFEL

Data Analysis group, European XFEL

March 2024

# How do we see Machine Learning at the EuXFEL?

- ■ Goal: maximize scientific outcome.
- ■ But ... not all approaches are equal.
- ■ Users have the last word on how to do their experiments.
- ■ Let's manufacture consensus.
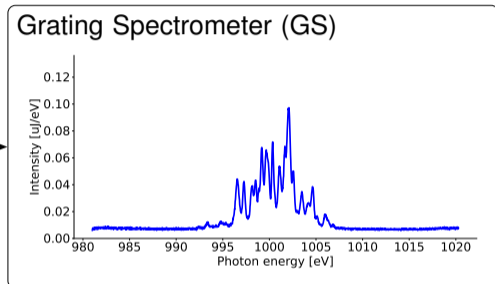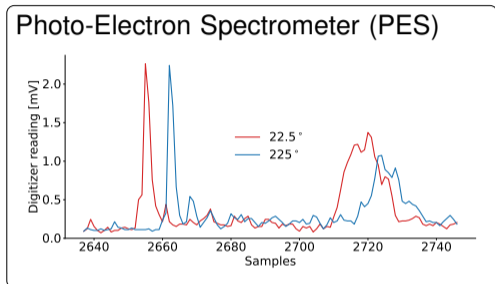
Characteristics:

- ■ *interpretability* → what do the results mean?
- ■ *context-aware* → connects to the science?
- ■ *quality control* → conditions for operation?

How to achieve it:

- ■ Clarify inner workings.
- ■ Shape it based on science.
- ■ Estimate region of validity.

# Virtual Spectrometer
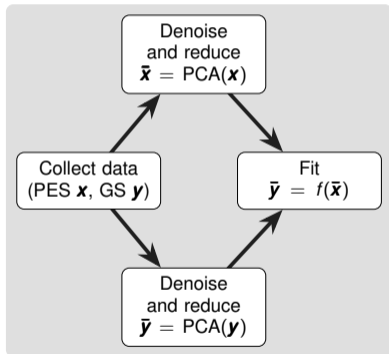
# Enhancing non-invasive X-ray diagnostics



**Photo-Electron Spectrometer (PES)**

**Grating Spectrometer (GS)**

- 🟩 Non-invasive.
- 🟩 Pulse-resolved.
- 🟥 Complex calibration.
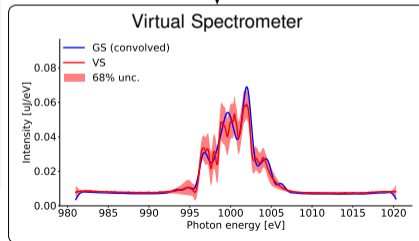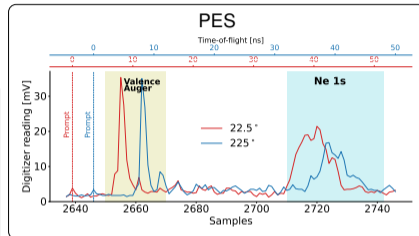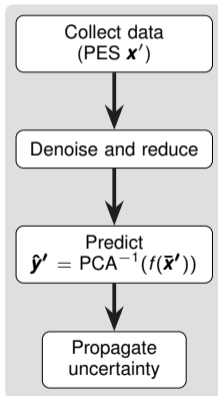- 🟥 Low resolution.

- 🟥 Invasive.
- 🟥 Train-resolved.
- 🟩 Simple calibration.
- 🟩 High resolution.

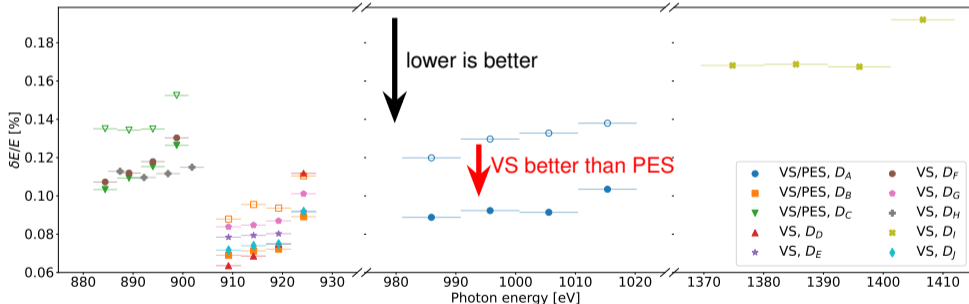# Enhancing non-invasive X-ray diagnostics: method

**Training**
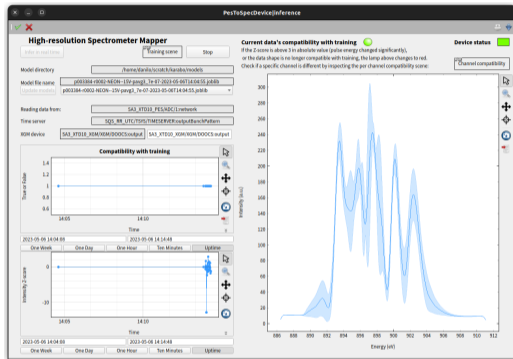


**Inference**

# Virtual Spectrometer's resolution

■ Systematic resolution studies under several conditions done.

■ Comparison with PES show better resolution.

■ Resolution calculated after training to inform scientists.



PES: open symbols; VS: full symbols.

European XFEL

# Deployment *online* and outlook

- Deployed in control system with simple interface to retrieve data and integrate ML projects.
- Combines advantages of low- and high-resolution devices:
  - Non-invasive.
  - Pulse-resolved.
  - Automated calibration.
  - Improved resolution.
- Adheres to self-defined guidelines:
  - Embedded **quality control**.
  - Resolution and uncertainty estimate ⇒ **interpretability**.
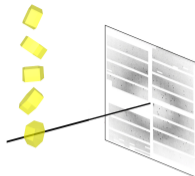  - SASE principle guides denoising ⇒ **context-aware**.



- Outlook:
  - Expand project for hard photons.
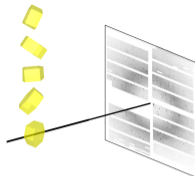  - Interpolate conditions to avoid pre-training stage.

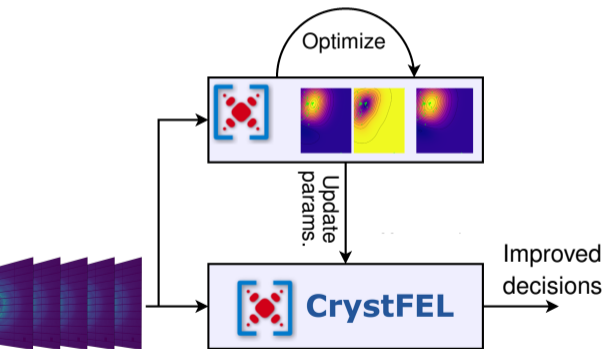# Automated data analysis

# Streamlining data analysis using ML

- Often data analysis pipelines have parameters.
- **Idea:** Simplify data analysis for non-experts.
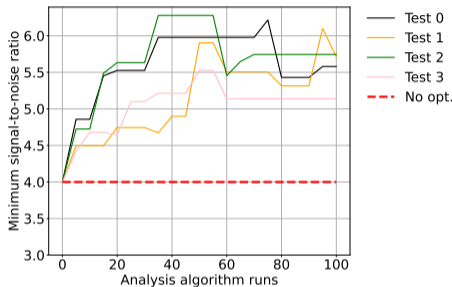
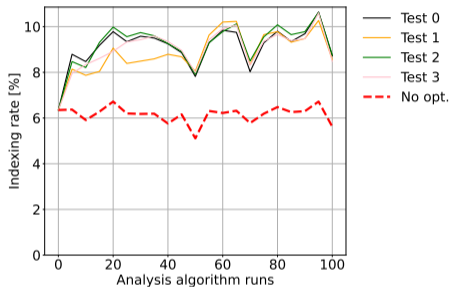# Streamlining data analysis using ML

- Often data analysis pipelines have parameters.
- **Idea:** Simplify data analysis for non-experts.



- **Goal:** Tune parameters to maximize a *metric*.
- This example: maximize the fraction of indexed frames $f$.
- *Online*: fast feedback, higher success chances.
- *Offline*: improved scientific findings.

# Streamlining data analysis using ML: example



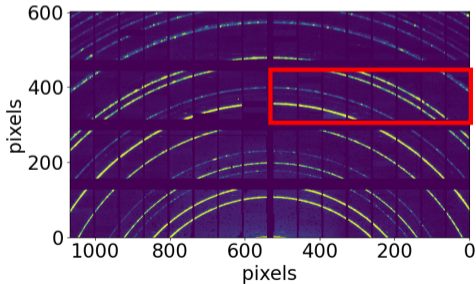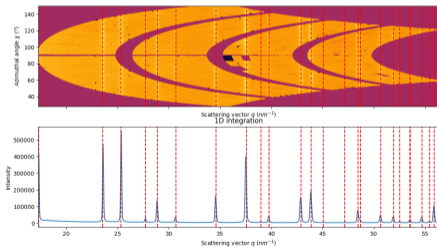- Hen Egg-White (HEW) Lysozyme.
- AGIPD detector at EuXFEL SPB/SFX.
- Web interface shows optimization progress ⇒ **interpretability**.
- **Quality metrics** also available.
- Optimizes clear science-based metrics ⇒ **context-aware**.

# Automated multi-modular geometry tuning

# Multi-modular geometry tuning

- ◼ Misalignment on module positions.
    - ◼ Manual alignment: requires lots of time.
    - ◼ Powder diffraction data are often the starting point for techniques requiring high-precision.
    - ◼ Powder diffraction-based methods require many parameters and often manual tuning.
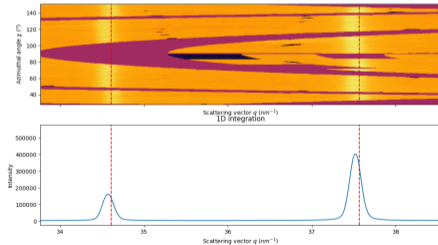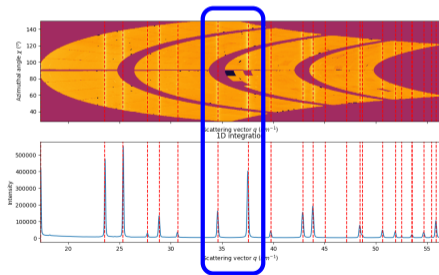- ◼ Let's start with powder diffraction: can we improve and *automate* it?

# Multi-modular geometry tuning

■ Misalignment on module positions.

   ■ Manual alignment: requires lots of time.
   ■ Powder diffraction data are often the starting point for techniques requiring high-precision.
   ■ Powder diffraction-based methods require many parameters and often manual tuning.

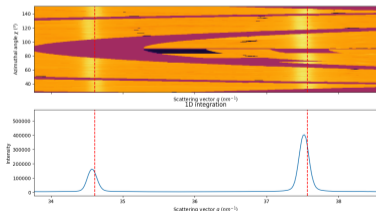■ Let's start with powder diffraction: can we improve and *automate* it?

# An information-theoretical approach

- ■ Optimizes the *mutual information* between radial distance and azimuthal angle → measures independence ⇒ **context-aware** and **interpretable**.
- ■ Pre-processing includes background subtraction and polar coordinate transformation.
- ■ Only first step in a long pipeline due to the limited experimental method resolution.
- ■ Validation tools available ⇒ quality metrics.

**Before inter-module tuning**

Tuning only detector-sample distance



**After inter-module tuning**

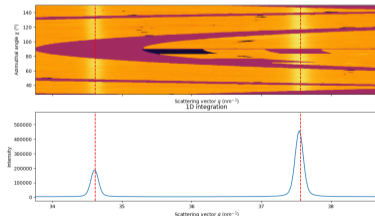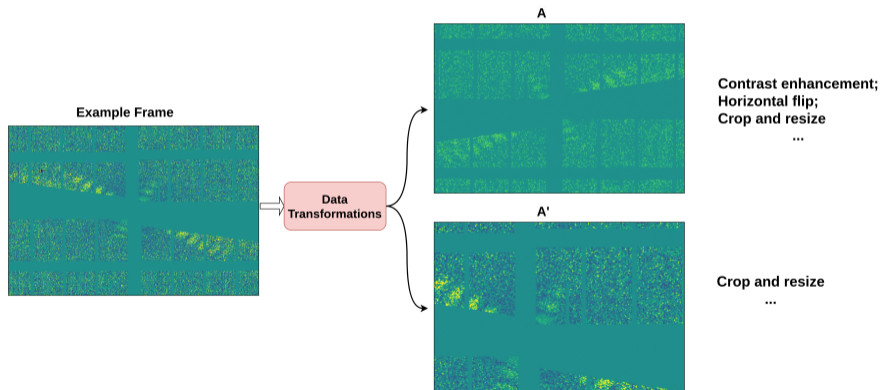Additionally, inter-module separation tuning



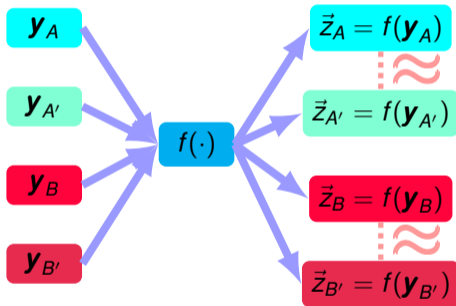Minimize
$MI(r, \phi)$

# Image clustering

# How do we Google data?

- How can we make data findable as soon as we collect it?
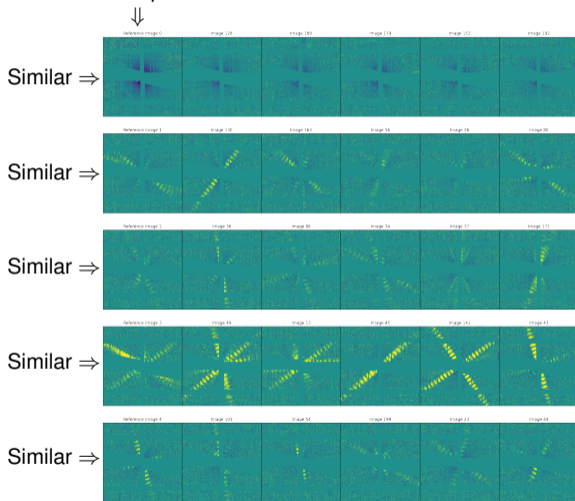- **Concept**: *Change* the data *view* and enforce their similarity.



**Example Frame**

**Data Transformations**

**A**

Contrast enhancement;
Horizontal flip;
Crop and resize
...

**A'**

Crop and resize
...

# Creating a similarity metric

- *Equivalent views* → variations to ignore, based on the science ⇒ **context-aware**.

$$\vec{z}_A = f(\mathbf{y}_A)$$

$$\approx$$

$$\vec{z}_{A'} = f(\mathbf{y}_{A'})$$

$$\vec{z}_B = f(\mathbf{y}_B)$$

$$\approx$$

$$\vec{z}_{B'} = f(\mathbf{y}_{B'})$$

$\mathbf{y}_A$  $\mathbf{y}_{A'}$  $\mathbf{y}_B$  $\mathbf{y}_{B'}$  →  $f(\cdot)$

In *Automated phase transition discovery*: Sun, Y., Brockhauser, S., Hegedüs, P., Plückthun, C., Gelisio, L., Ferreira de Lima, D. E., Sci. Rep., (2023).
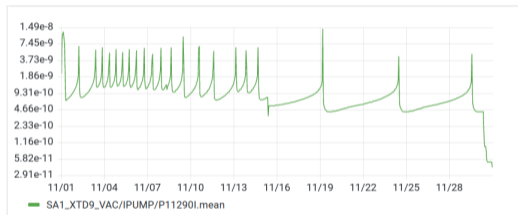
Example
⇓

Similar ⇒

Similar ⇒

Similar ⇒

Similar ⇒

Similar ⇒



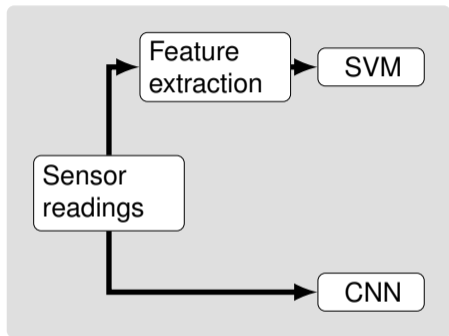**European XFEL**

# Predictive Maintenance

# Predictive Maintenance: ion vacuum pump use-case

- Faults may lead to loss of beam time.
- Important to detect them early.
- Difficulty: complex system makes it hard for humans to monitor everything.
- Example: Ion pump faults have lead to significant downtime.
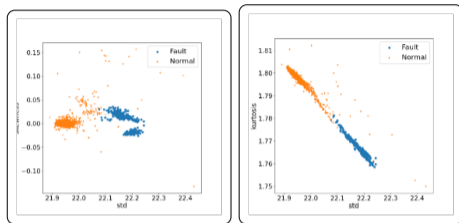- Detection mechanism: frequent surges in pressure level.



(Amna Majid)

# How can we detect it?



| Method | Accuracy [%] | Precision | Recall |
|--------|-------------|-----------|--------|
| SVM | 99.98 | 1.00 | 0.96 |
| CNN | 99.95 | 0.99 | 0.99 |

(Amna Majid)

European XFEL



- Two methods researched with similar performance.

- SVM makes a linear cut in the feature space of peak characteristics → easy **interpretation** and based on **context**.

- CNN uses all information.

- Prefer interpretable method!

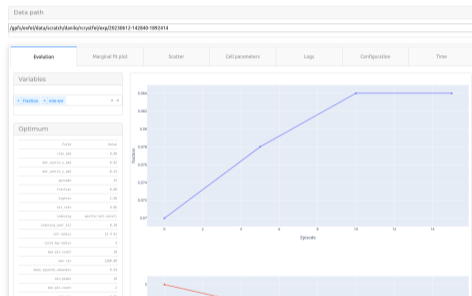- Web interface for monitoring ⇒ **quality control**.

# Summary

- Several approaches to enhance automation at the EuXFEL.
- **Interpretability**, **context-awareness** and **quality control** are seen as assets to guide towards adequate solutions.
- Control system allows for integration and deployable methods.
  - Interface design is simple, but highlights those characteristics to guide users.
  - Aim for a holistic approach to integrate those features in all applications.
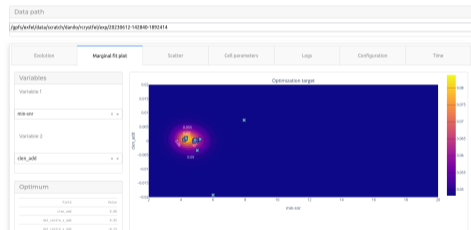
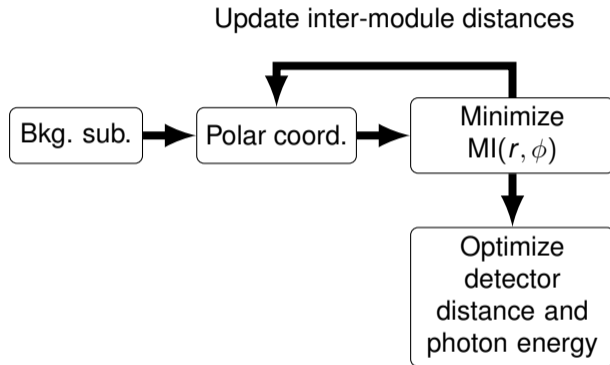Thank you!

# Additional material

# Interpretability in a Web interface

# Multi-modular geometry tuning: concept



Update inter-module distances

Bkg. sub. → Polar coord. → Minimize MI($r, \phi$) → Optimize detector distance and photon energy
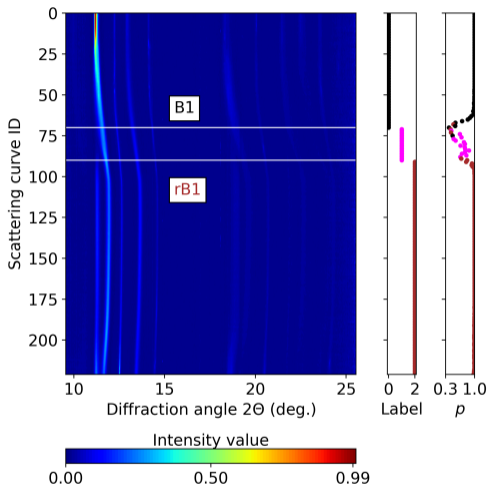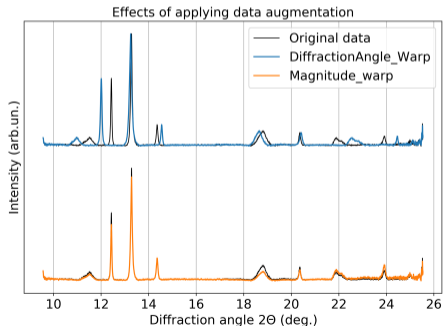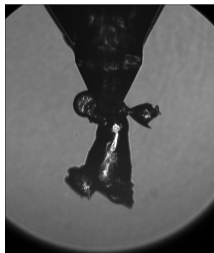
# Googling phase transitions

- ■ Ignore all changes in the spectra, *but* phase transitions.

- ■ Learn to map irrelevant variations into the same $\vec{z}$.



Effects of applying data augmentation



Sun, Y., Brockhauser, S., Hegedüs, P. , Plückthun, C., Gelisio, L., Ferreira de Lima, D.E., Sci. Rep., (2023).

European XFEL

# Protecting detectors from damage



- Ice can form on the tip of nozzles, and scatter X-rays that can destroy detector pixels.
- **Idea:** Use computer vision techniques to detect:
    - jet instabilities, reducing beamtimes efficiency;
    - ice formation.
- Information can be used to alert operators or even intervene.
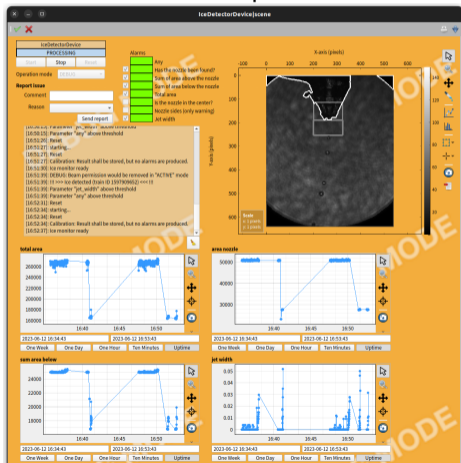
# Protecting detectors from damage



Experiment's side camera.

Two consecutive frames are separated by 0.1 s.

# Software deployed
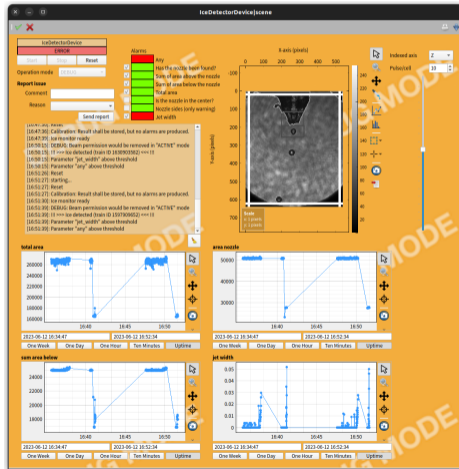
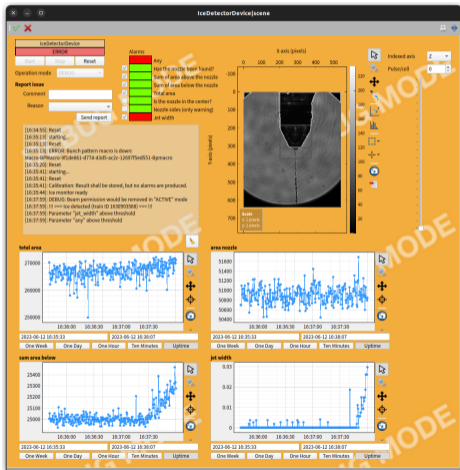Standard operation

Ice detected



...age. Human monitoring slow.

# Drive control system responsibly



- Ice formation may lead to detector damage. Human monitoring slow.
- *Interpretability*: Interface informs on alarm source and low-quality data.
- Operators are still in control: we only guide them on request.
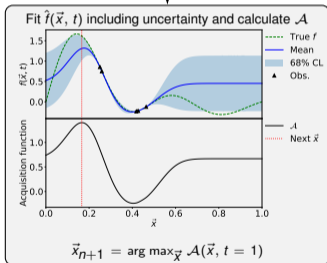
# Baysian Optimization: How does it work?

Dynamic Bayesian optimization.

- A Gaussian process is used to fit the objective function $f$

$$C(\vec{x}_1, t_1, \vec{x}_2, t_2 | S, \vec{L}, T, \sigma) = S^2 \, e^{-\frac{(\vec{x}_1 - \vec{x}_2)^2}{2\,\vec{L}^2}} \, e^{-\frac{(t_1 - t_2)^2}{2\,T^2}} + \sigma^2$$

- The acquisition function is defined as

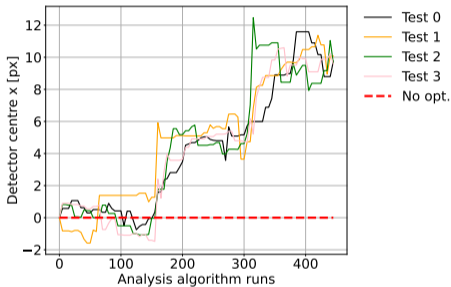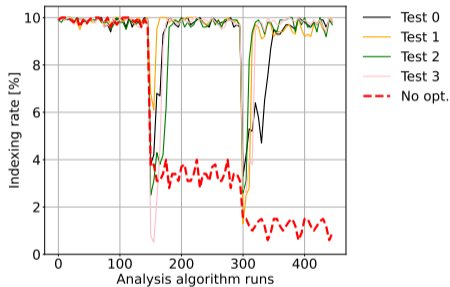$$\mathcal{A}(\vec{x}) = \bar{f}(\vec{x}, t = \text{current}) + \sqrt{\beta} \, \delta f(\vec{x}, t = \text{current})$$

Initially run analysis ($\times n$) and store:

| Parameter | Objective |
|-----------|-----------|
| $\vec{x}_1$ | $f(\vec{x}_1, \mathcal{D}(t=1))$ |
| . | . |
| . | . |
| $\vec{x}_n$ | $f(\vec{x}_n, \mathcal{D}(t=1))$ |



Fit $\hat{f}(\vec{x}, t)$ including uncertainty and calculate $\mathcal{A}$

- - - True $f$
- Mean
- 68% CL
- ▲ Obs.
- $\mathcal{A}$
- Next $\vec{x}$

$\vec{x}_{n+1} = \arg\max_{\vec{x}} \mathcal{A}(\vec{x}, t = 1)$

Run analysis at $\vec{x}_{n+1}$ and store:

| Parameter | Objective |
|-----------|-----------|
| $\vec{x}_1$ | $f(\vec{x}_1, \mathcal{D}(t=1))$ |
| . | . |
| . | . |
| $\vec{x}_{n+1}$ | $f(\vec{x}_{n+1}, \mathcal{D}(t=2))$ |

# Bayesian Optimization: Simulating a detector shift



- Hen Egg-White (HEW) Lysozyme.
- Simulated AGIPD data, X-ray beam pointing shifting twice.

# Enforce self-consistency mapping

**Idea**: different views of the data containing the same information must be understood as the same.