# Automated analysis workflows

James Wrigley for DA/MID/FXE/lots of lovely users

## `whoami`



■ Data analysis contact for MID.

■ Spends lots of time preparing analysis for experiments and helping users.

**European XFEL**

# Common pain points :(

(from an analysis perspective)

1.  Many experiments produce large amounts of data that require significant work to get something scientifically meaningful. Examples:

    - XPCS: 404 runs, 628TB

    - SFX: 286 runs, 161TB

    - Solution scattering: 435 runs, 441TB

2.  The raw data recorded by the facility is often difficult to understand (many data sources, confusing names, etc) and in a different representation than what's desired. Examples:

    - Nozzle temperature: `MID_EXP_UPP/CTRL/LSHORE.inputB.krdg`

    - Monochromator position: `FXE_XTD9_MONO-1/MOTOR/ACCM_PITCH.actualPosition`

**European XFEL**

## Key takeaway

With *high-level analysis* and *automation* it's possible to get a synchrotron analysis experience[1]!

---

[1]This is a lie, but it's close to being true.

**European XFEL**

# That synchrotron feeling

Concrete steps that we suggest doing:

**European XFEL**

## That synchrotron feeling

Concrete steps that we suggest doing:

- Most importantly, talk to us. Tell the instrument contact or local data contact about what you want to do or send an email to da@xfel.eu.

**European XFEL**

# That synchrotron feeling

Concrete steps that we suggest doing:

- Most importantly, talk to us. Tell the instrument contact or local data contact about what you want to do or send an email to da@xfel.eu.
- Some technical steps:

  - Use <u>extra-data's aliases</u> for sensible names to do e.g. `run.alias["mono-position"]`:

  ```yaml
  extra-data-aliases.yml

  nozzle-temperature: [MID_EXP_UPP/CTRL/LSHORE, inputB.krdg]
  mono-position:      [FXE_XTD9_MONO-1/MOTOR/ACCM_PITCH, actualPosition]
  ```
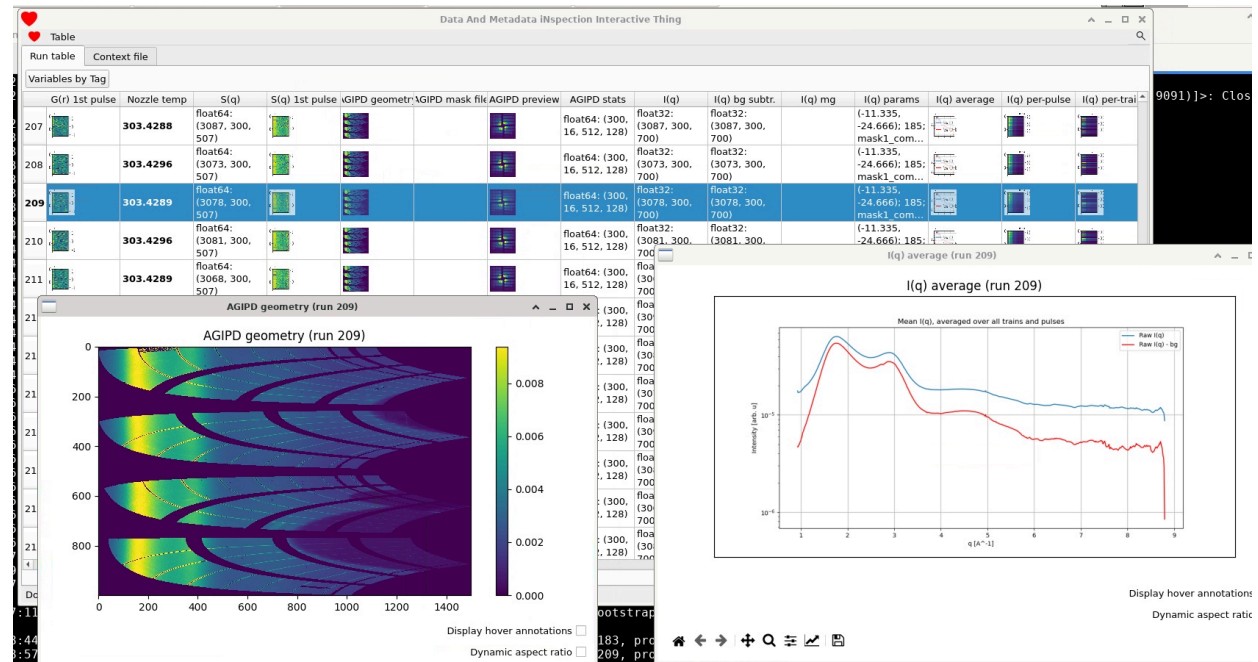
**European XFEL**

# That synchrotron feeling

Concrete steps that we suggest doing:

- Most importantly, talk to us. Tell the instrument contact or local data contact about what you want to do or send an email to da@xfel.eu.
- Some technical steps:

    - Use extra-data's aliases for sensible names to do e.g. `run.alias["mono-position"]`:

        extra-data-aliases.yml

        ```
        nozzle-temperature: [MID_EXP_UPP/CTRL/LSHORE, inputB.krdg]
        mono-position:      [FXE_XTD9_MONO-1/MOTOR/ACCM_PITCH, actualPosition]
        ```

    - Move code from notebooks into an automated pipeline with DAMNIT.

**European XFEL**

# DAMNIT

- Automatically creates a run table from custom, user-defined functions.
- Results saved into a database for display and accessible through a <u>Python API</u>.

```
from extra.damnit import Damnit

db = Damnit(1234)
i_q = db[200]["i_q"].read()
```



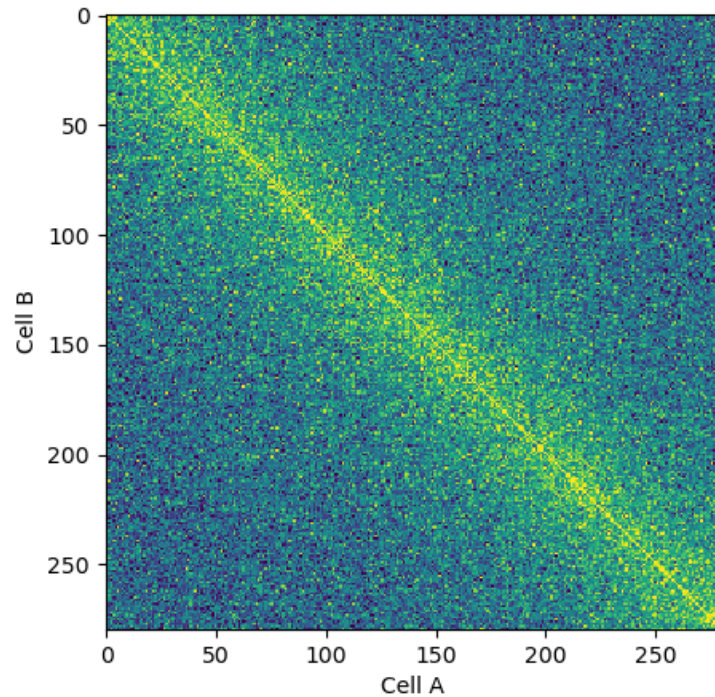Also see poster 58, "DAMNIT, a tool for Automated Experiment Overview" by T. Michelat.

**European XFEL**

# Example: XPCS



- X-ray Photon Correlation Spectroscopy: technique for investigating sample dynamics based on correlating laser speckle across time.
- Key quantities are Two-Time Correlation Functions (TTCFs) and $g_2$ functions.



XPCS experiment diagram by M. Dargasz.

**European XFEL**

# Example: XPCS



Example TTCF



$g_2$'s for multiple TTCFs at different $q$ values

**European XFEL**

# Example: XPCS

Mature offline pipeline, developed over a few years by multiple people from DA/MID/University of Siegen.

Performance: processing a 5 minute run takes ~20 minutes using multiple processes.

Fully integrated with DAMNIT



```python
from extra.damnit import Damnit

db = Damnit(1234)
xpcs = db[300]["xpcs_mean_dataset"].read()
```



Also see poster 11, "A High-Throughput Data Pipeline for MHz-XPCS: Offline Analysis" by A. Leonau.

Also see the "A High-Throughput Data Pipeline for MHz XPCS: Analyzing Protein Dynamics via Solution Scattering" talk by A. Leonau in the Data Science session on the 22nd.
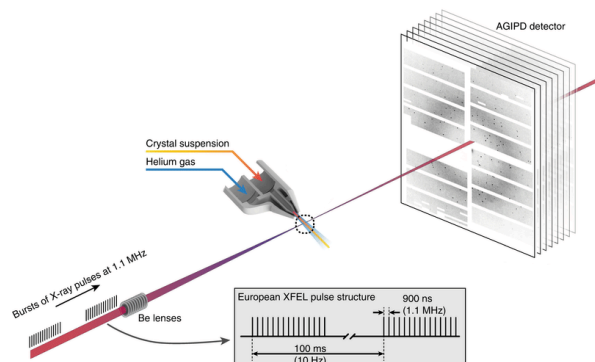
**European XFEL**

# Example: XPCS (online)

- Quantities are computed in the same way as offline, modulo certain filtering methods.
- Provides real-time feedback at 10Hz.
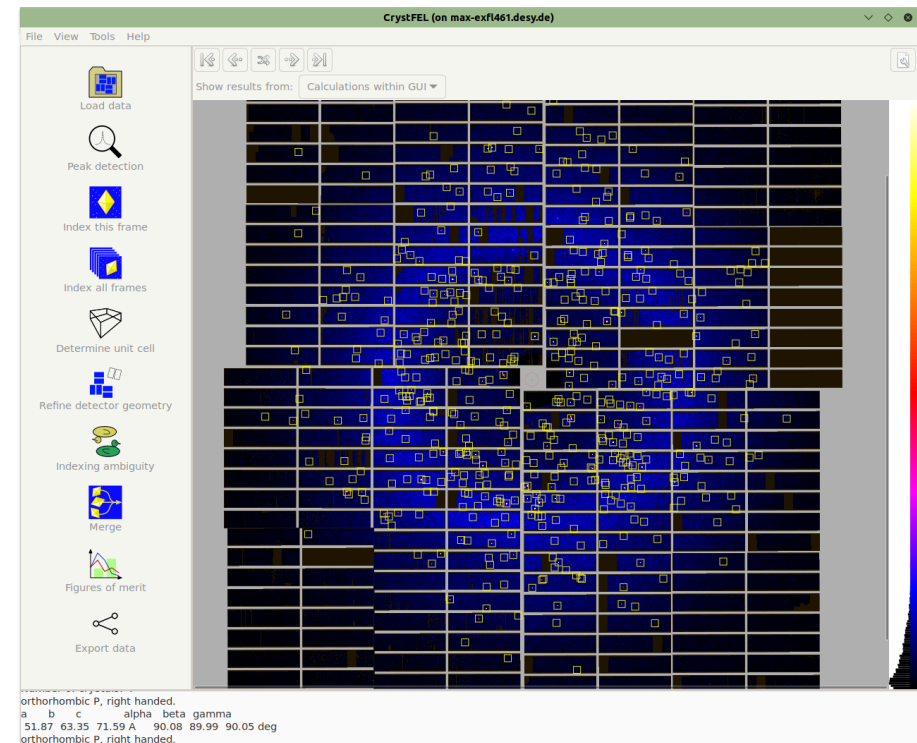- Also integrated with DAMNIT to make online results available as soon as the run finishes.



Also see poster 10, "A high throughput Data Pipeline for MHz XPCS: Online Analysis" by M. Jakobsen.

**European XFEL**

# Example: SFX

- Serial Femtosecond Crystallography: technique for finding the atomic structure of a sample.
- Users come with many small crystals (e.g. from a protein), which are fed into the beam to diffract onto the detector.
- Analysis requires multiple steps (e.g. hit finding, indexing) to end up with an electron density map.



SFX experiment diagram adapted from Wiedorn et al. (2018), Nat. Comm. 9.

**European XFEL**

# Example: SFX

■ <u>Extra-Xwiz</u>: a wrapper around <u>CrystFEL</u>.

■ Goal is to simplify parallelization, data handling, and configuration. End result is the structure factors.

■ Integrated with DAMNIT to analyze runs automatically.



xwiz_lysozyme.toml

```
[data]
proposal = 1234
runs = [28, 33, 34]
list_prefix = "r0028_t12"

[geom]
file_path = "geom_with_masks.geom"

[proc_coarse]
resolution = 0.4
peak_method = "peakfinder8"
peak_threshold = 300
peak_snr = 7.0
peak_min_px = 1
peak_max_px = 20
peaks_hdf5_path = "entry_1/result_1"
index_method = "xgandalf"
```
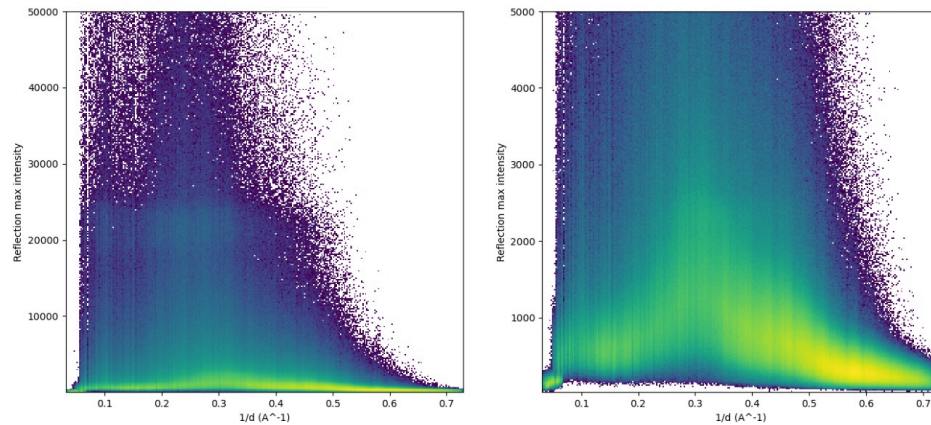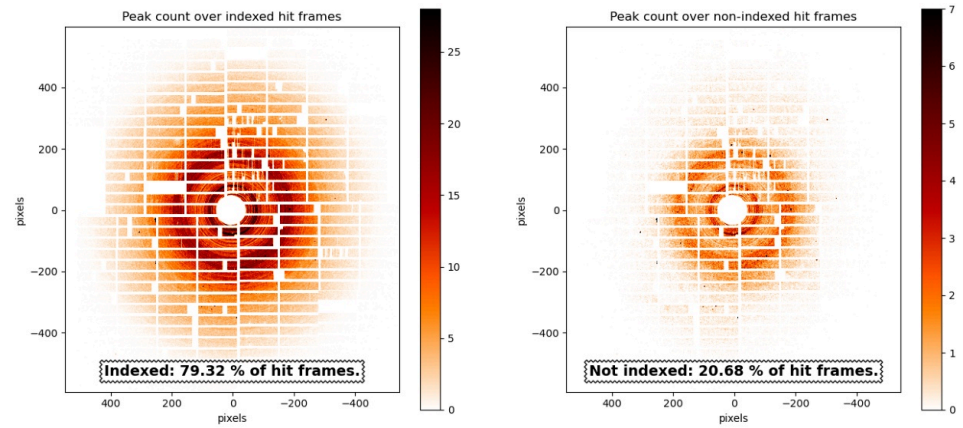
Also see poster 40, "Automatic data processing and results overview during SFX experiments" by O. Turkot.
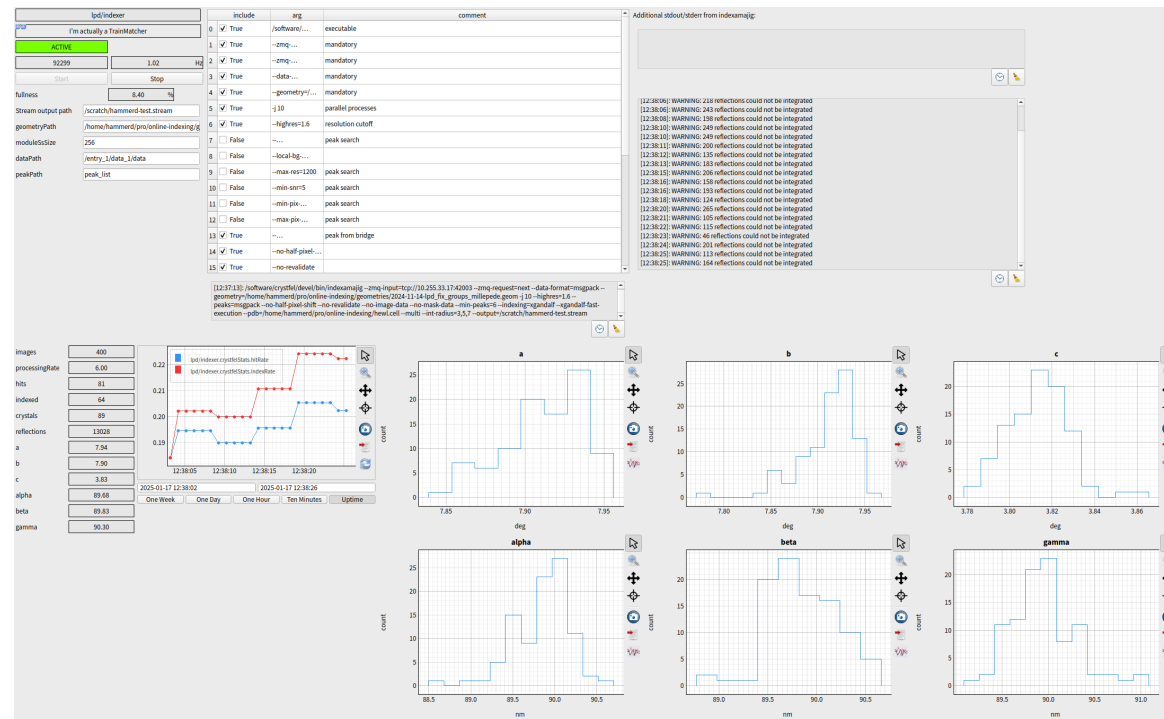
**European XFEL**

# Example: SFX

Example debugging plots:



Peakogram



Peak counts for indexed/non-indexed frames

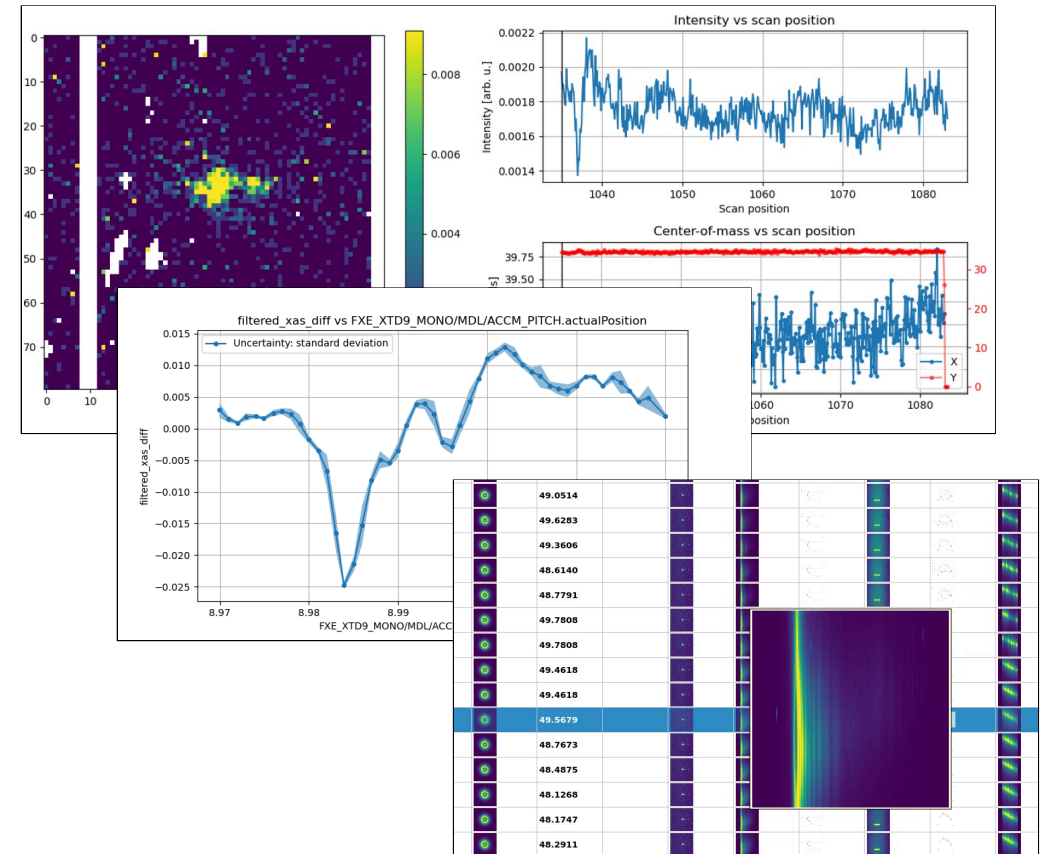**European XFEL**

# Example: SFX (online)

What we can currently do online: hit finding and indexing.

# Others

These techniques are in various stages of maturity and coverage:

- Bragg XPCS
- Single particle imaging
- Single crystal diffraction
- X-ray absorption spectroscopy
- X-ray emission spectroscopy
- X-ray scattering



**European XFEL**

# Conclusion

This only works by collaborating with amazeballs users and instrument scientists <3

**European XFEL**

# Conclusion

- This only works by collaborating with amazeballs users and instrument scientists <3
- Analysis at our facility is probably never going to be quite as easy as at a synchrotron, but we can get close.

**European XFEL**

# Conclusion

- This only works by collaborating with amazeballs users and instrument scientists <3
- Analysis at our facility is probably never going to be quite as easy as at a synchrotron, but we can get close.

*Questions?*

(or contact us at da@xfel.eu)

**European XFEL**