



# Finding NEMOs

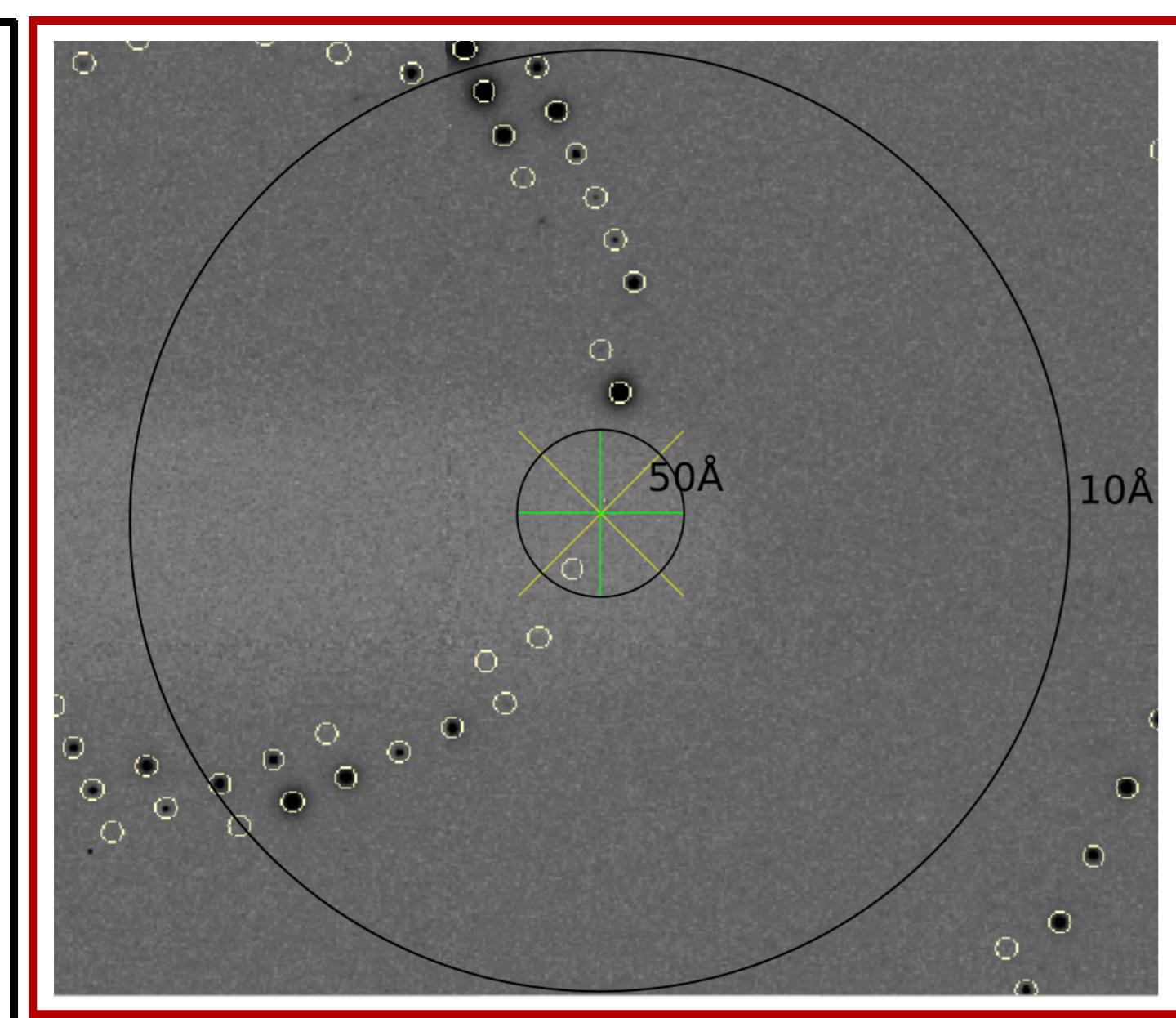
Yunyun Gao<sup>a</sup>, Helen Ginn<sup>a,b</sup>, Andrea Thorn<sup>a</sup>

a. Institut für Nanostruktur und Festkörperphysik, Universität Hamburg, Germany

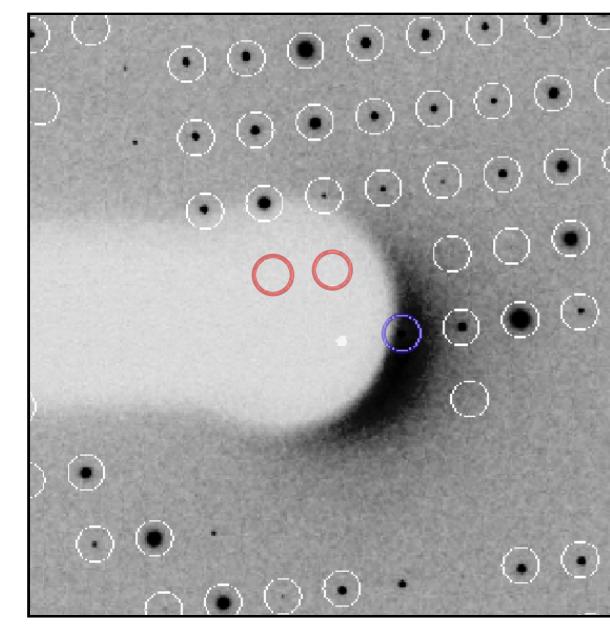
b. Center for Free-Electron Laser Science, CFEL, Deutsches Elektronen-Synchrotron DESY

## A Story about Lost in Autoprocessing (and the Redemption)

Try me: How will the pipeline/you mask the beamstop



the Data Reduction Workflow  
Indexing->Prediction->Integration



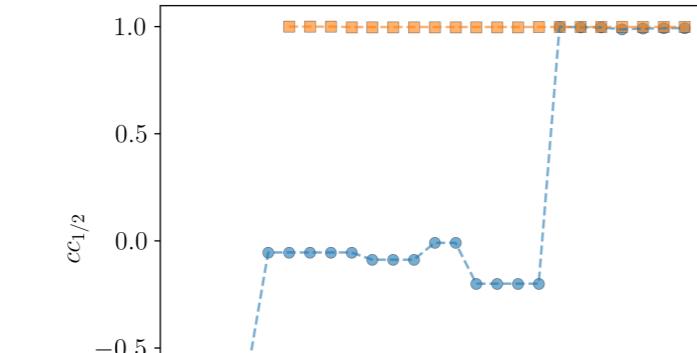
Two types of low-angle outliers related to beamstop. Predicted locations of reflections are circled. The blue circle and red circle highlight the Type 1 and Type 2 outliers, respectively. Type 1 outliers are typically strong and detectable using Wilson statistics. Detection and exclusion of the Type 2 outliers remains more difficult.

**NEMOs**  
Not-Excluded-unMasked-Outliers  
cluster(s) of weak observations  
at the low-resolution end  
in processed datasets

How NEMOs Escape



"Binning-based" Data Statistics

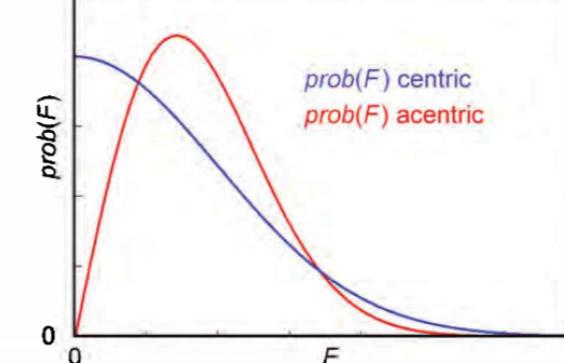


Perfect Beamstop Mask for the Entire Dataset Does That Even Exist?

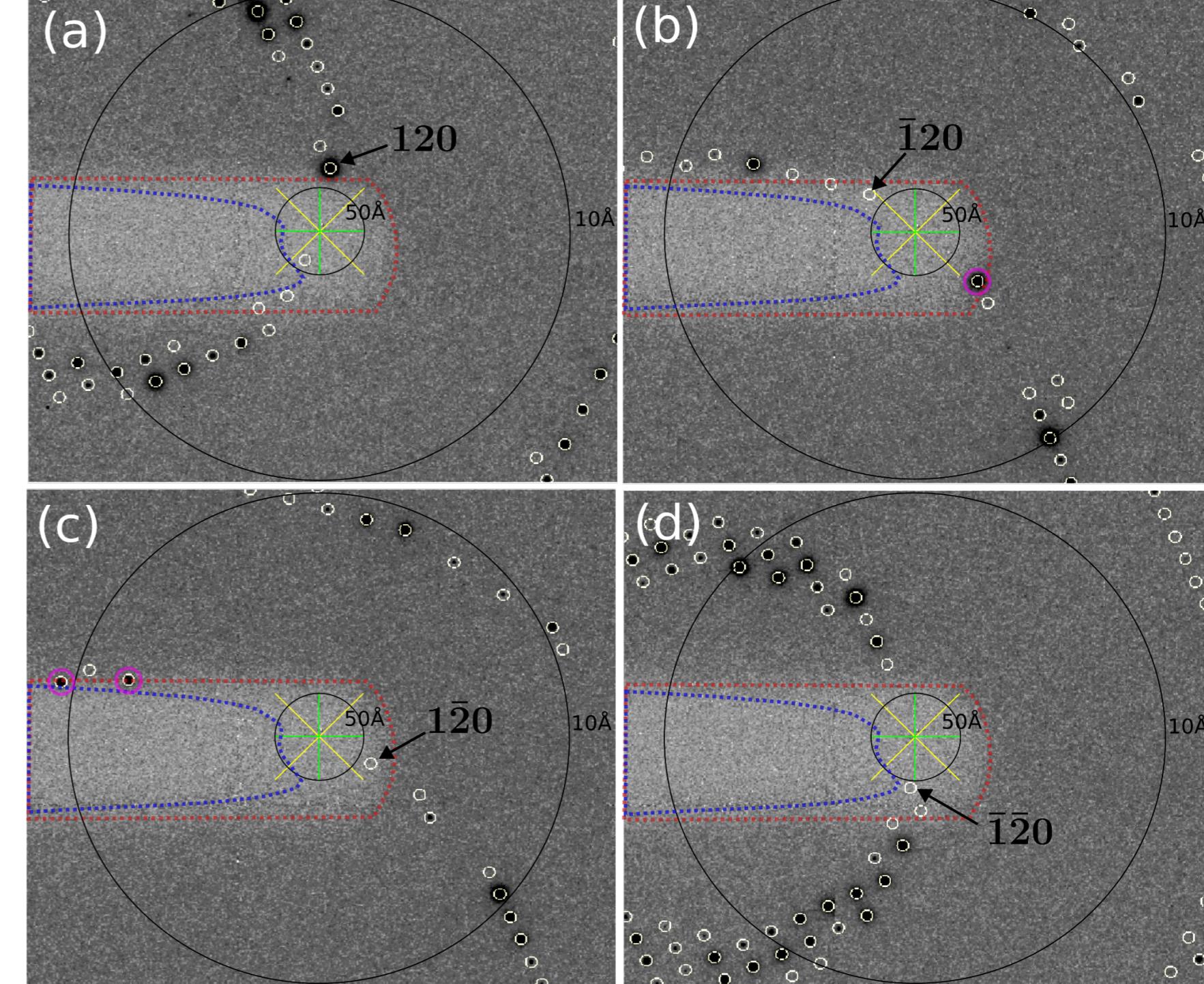
Crystallographical Statistics

$$p_a(E_0 < E_{0,\text{meas}}) = 1 - \exp(-E_{0,\text{meas}}^2)$$

$$p_c(E_0 < E_{0,\text{meas}}) = 1 - \text{erfc}\left(\frac{E_{0,\text{meas}}}{2^{1/2}}\right)$$

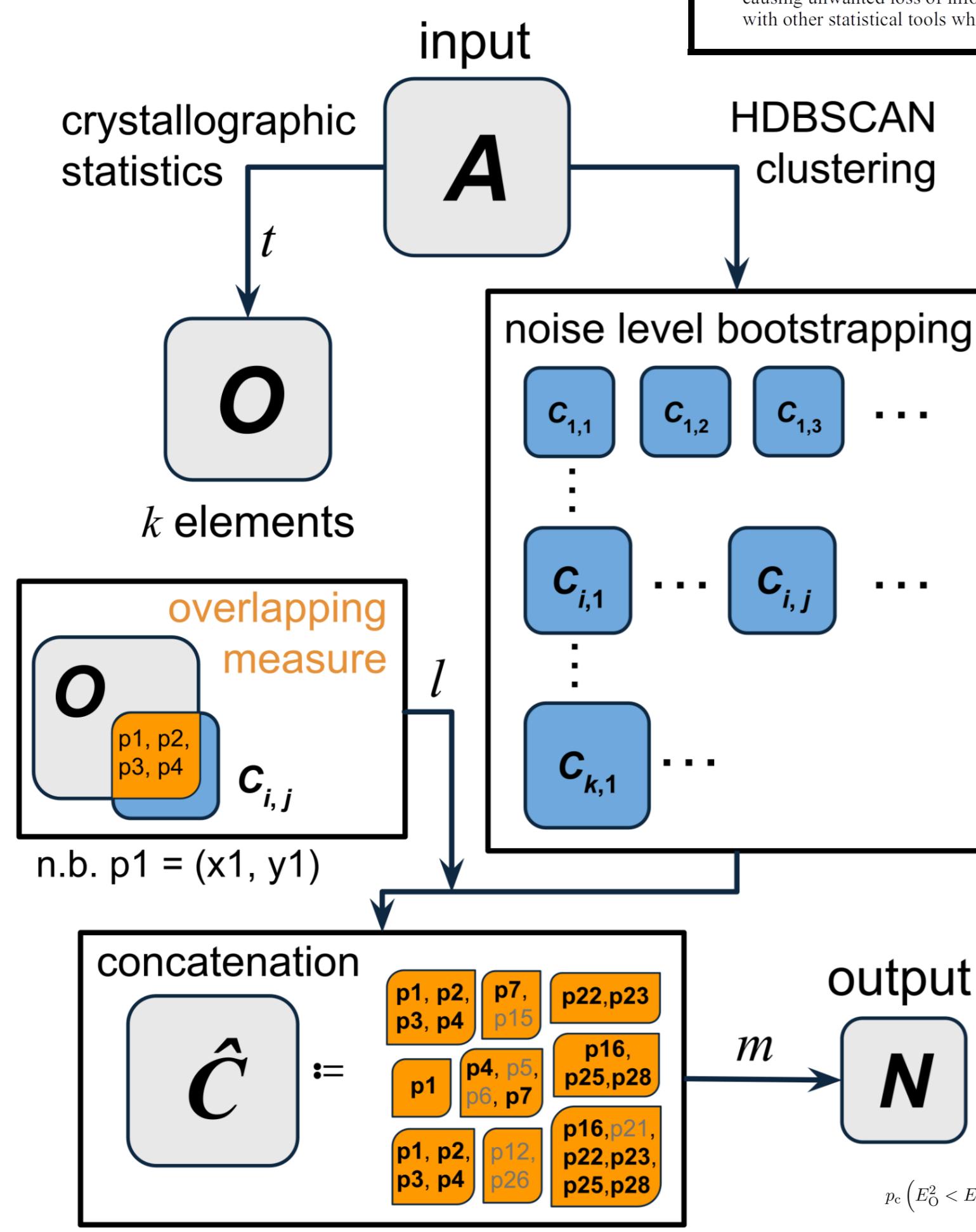


Outlier Rejection : Multiple Measurement of an Unique Reflection



Example of the experimental origin of NEMOs. The images are the corresponding raw detector frames of PDB entry 5bsf. The blue dashed polygon highlights the edge of excluded shaded regions recognised by the best possible XDS-DEFPX trial without masking unobstructed areas on the detectors. The red dashed polygon highlights the edge of how one intuitively would define a mask, manually created with the aid of *dials image viewer*. Using the blue mask, the similarity of observations in (b), in (c) and in (d) result in the exclusion of observation in (a). The resulting merged unique reflection then would have a value close to 0, but an "algorithm acceptable" uncertainty. Using the red mask, the unique reflection can be properly recorded, as the other symmetry equivalent observations are completely masked. However, in (b) and (c) such a mask results in the masking of unbiased observations (highlighted with magenta circles).

## Find NEMOs



But Keep Dori

2. HDBSCAN on multiple noise level

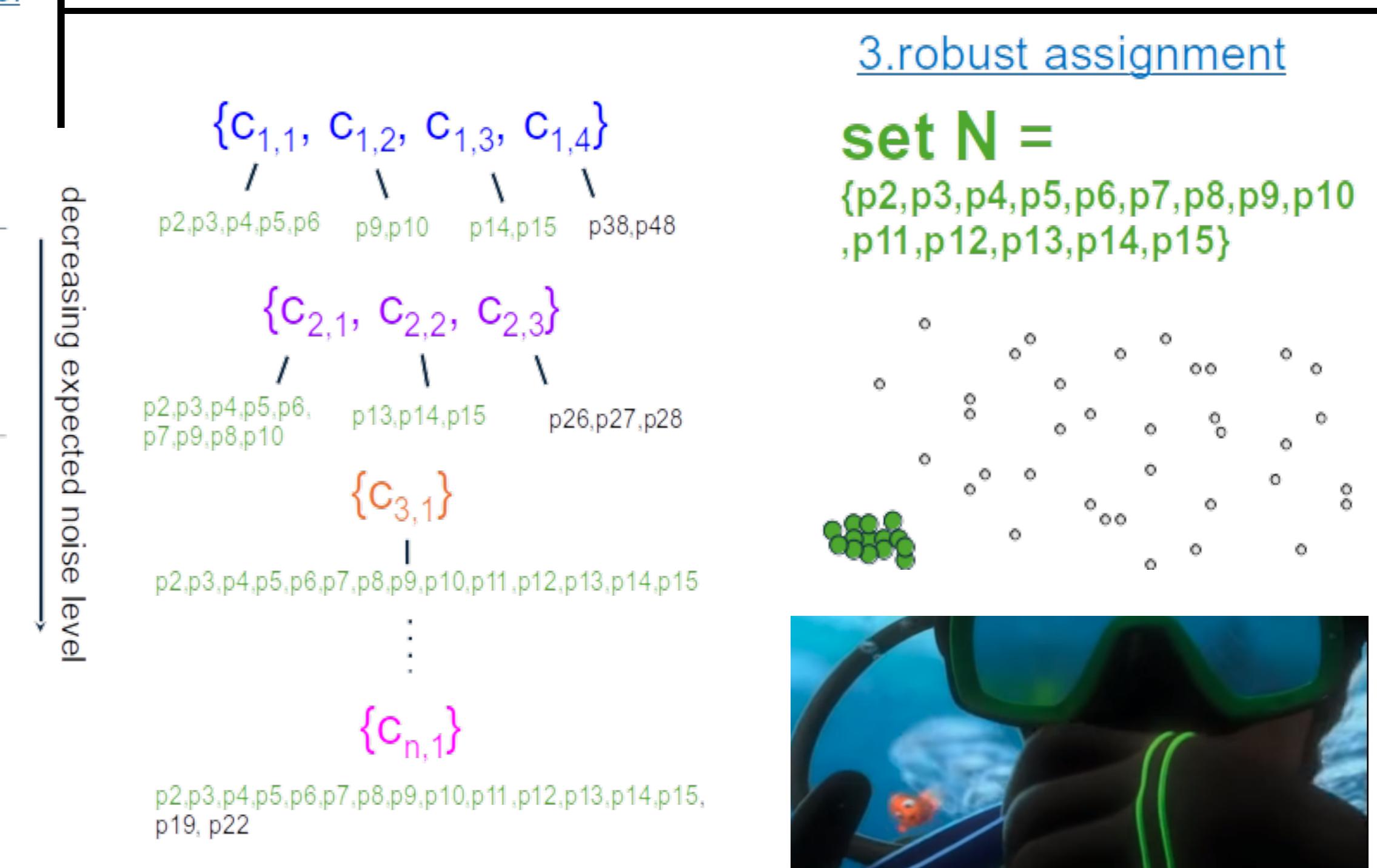
1. Constructing a pool for potential beamstop shadow outlier  
 $P(E < E_0) < t, d > 10 \text{ \AA}$

$t > 10^{-2}$

set O =  
(p2,p3,p4,p5,p6,p7,p8,p9,p10,  
p11,p13,p14,p15,  
p21,p22,...,  
p26,c28,...,  
p38,...,  
p49,...)

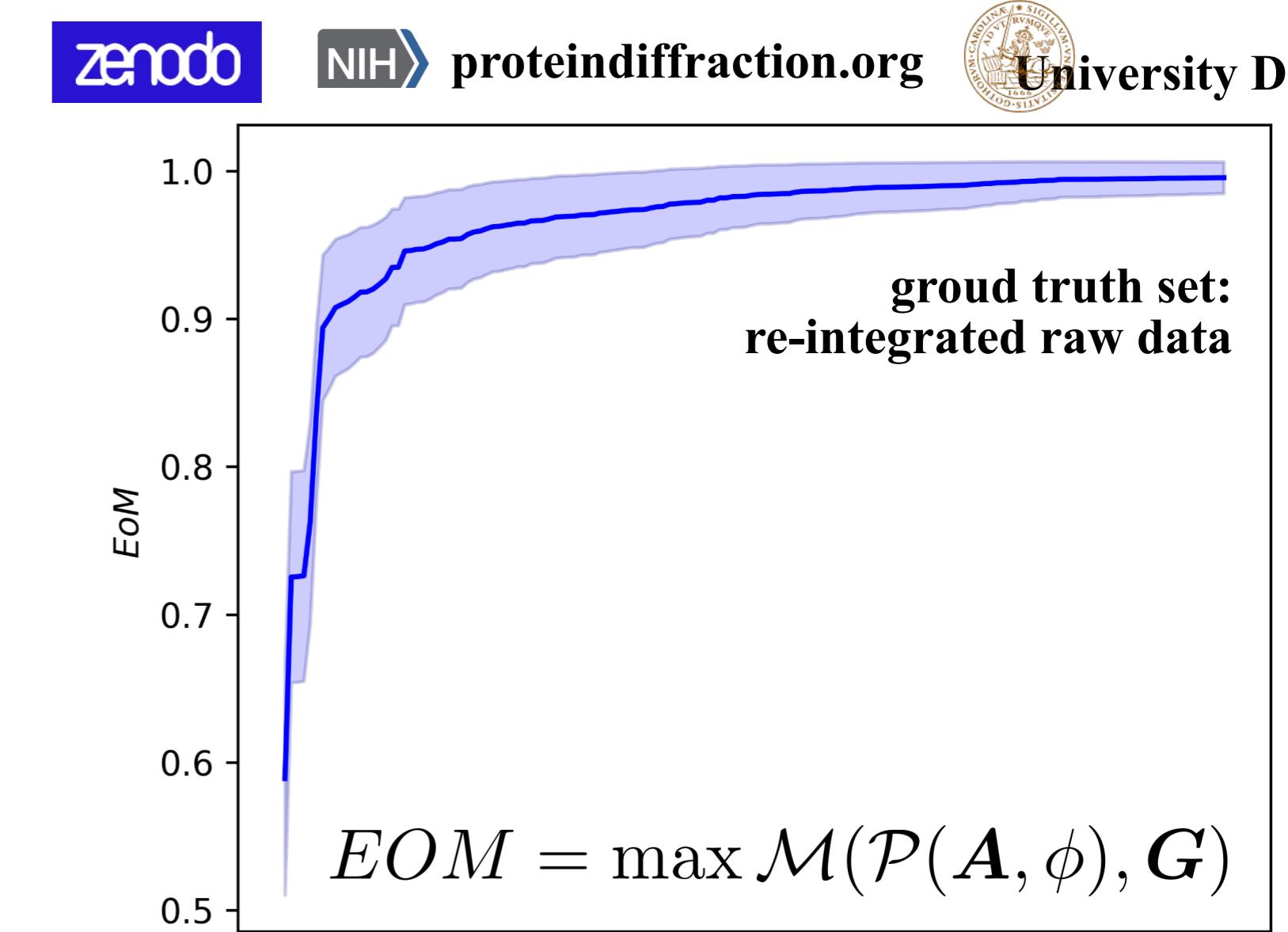
$$p_a(E_0 < E_{0,\text{meas}}) = \frac{1}{2} \left[ \text{erfc}\left(\frac{E_{0,\text{meas}}}{2^{1/2} \sigma_{E_0}}\right) - \exp\left(-\frac{\pi(E_{0,\text{meas}} - E_{0,\text{min}})^2}{2\sigma_{E_0}^2}\right) \text{erfc}\left(\frac{\pi(E_{0,\text{meas}} - E_{0,\text{min}})}{2^{1/2} \sigma_{E_0}}\right) \right]$$

$$p_c(E_0 < E_{0,\text{meas}}) = \frac{E_{0,\text{meas}}}{E_{0,\text{min}}} \int_{E_{0,\text{min}}}^{\infty} \frac{1}{(2\pi)^{1/2} \sigma_{E_0}} \exp\left[-\frac{(E_0 - E_{0,\text{meas}})^2}{2\sigma_{E_0}^2}\right] \frac{1}{(2\pi)^{1/2} \sigma_{E_0}} \exp\left(-\frac{E_0^2}{2\sigma_{E_0}^2}\right) dE_0$$



Weaving the Fishing Net: Crystallographical Statistics + Hierarchical Density-Based Clustering + Semi-Supervised Hyperparameter Tuning

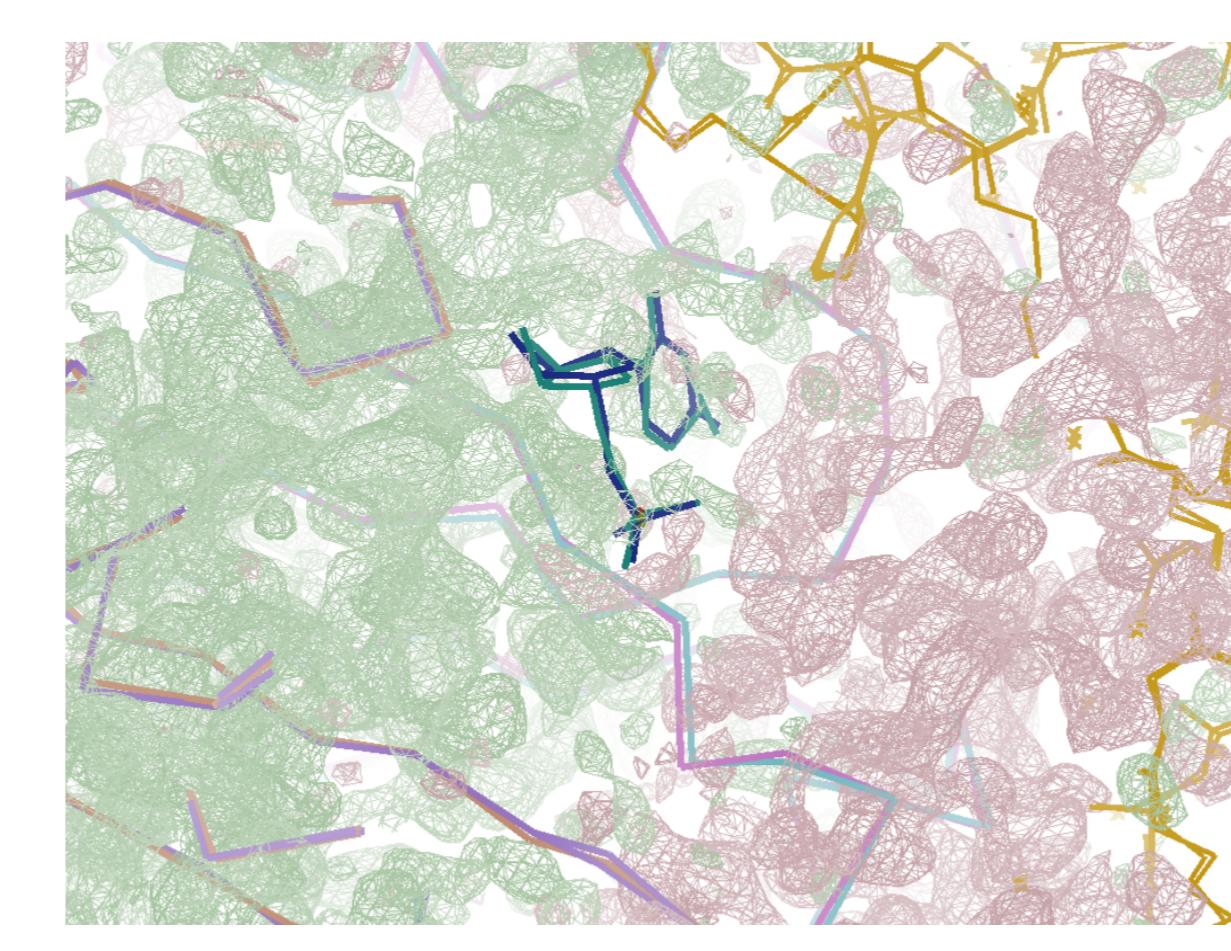
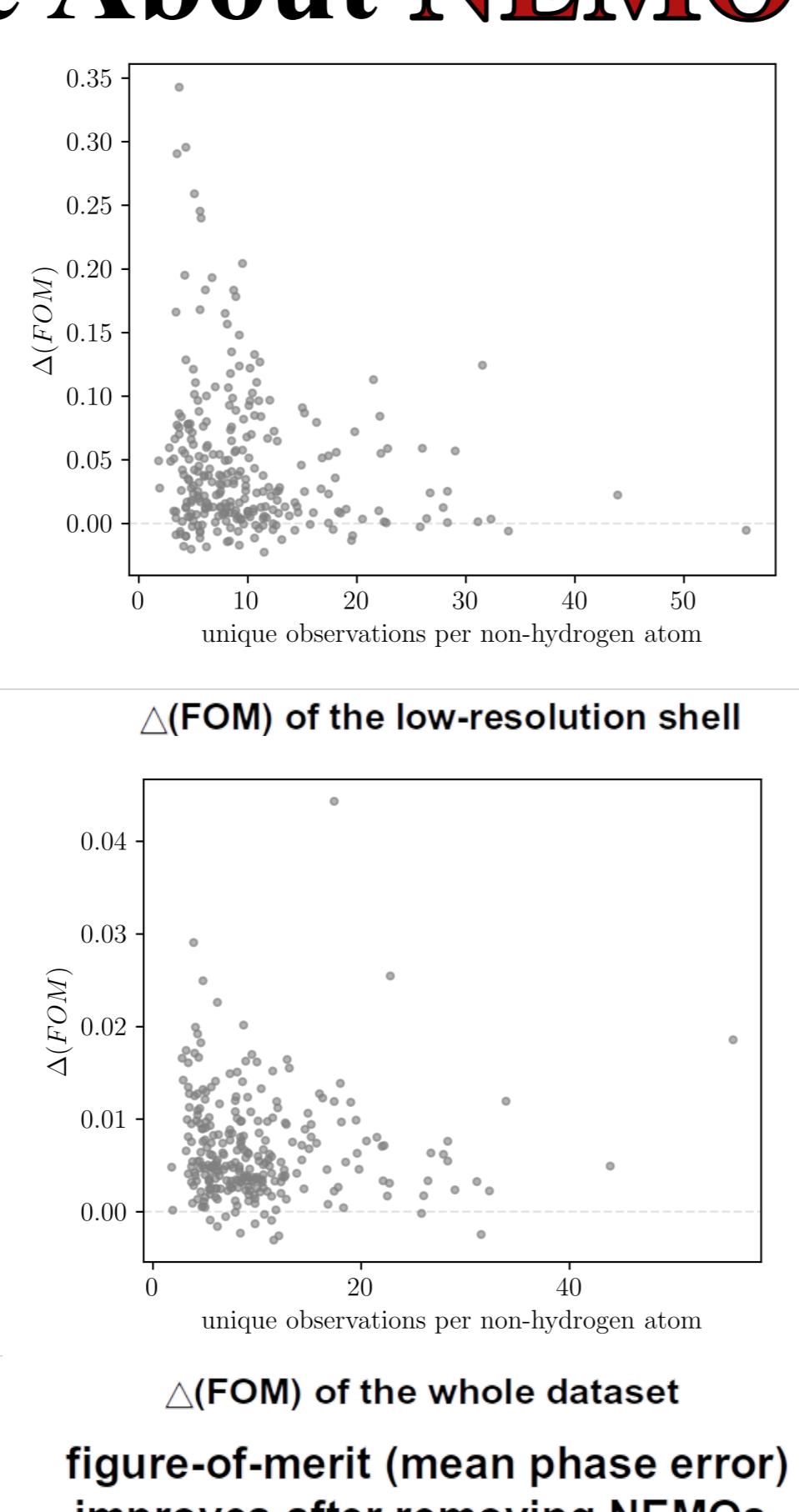
MORE RAW IMAGES/UNTRUNCATED PROCESSED DATA PLS!



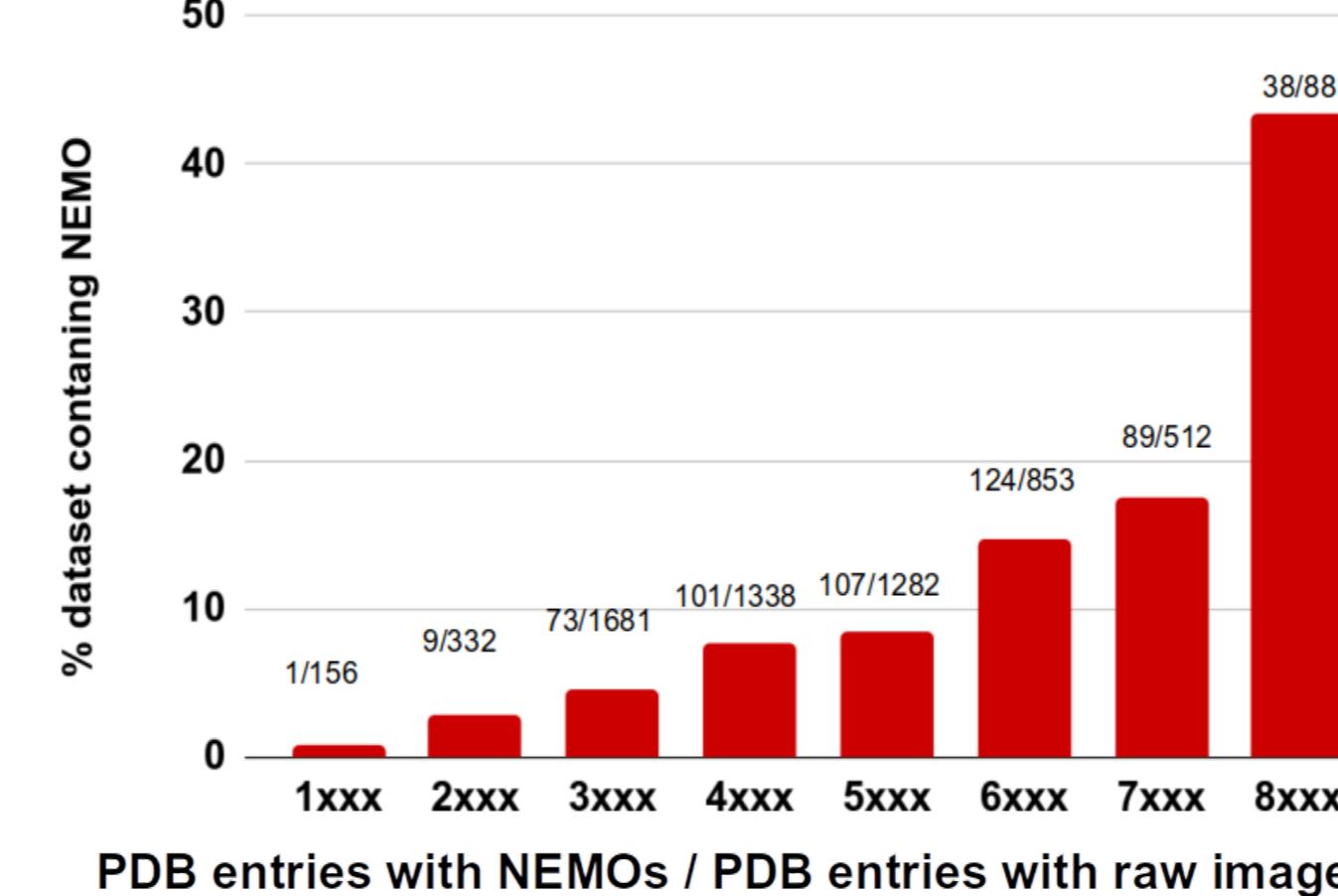
## Why You Should Care About NEMOs

	range
NEMOs %	0.75 - 0.0015
resolution	3.38 - 0.97
R-free	0.341 - 0.1191
# unique observations per non-hydrogen atoms	1.8 - 55.7
solvent %	32.25 - 81.83
synchrotron	NSLS-II,ESRF,BESSY, PETRA III,APS,Diamond, australian synchrotron, ALS,SSRL,LNLS
data reduction	XDS,HKL-3000, HKL-2000,DIALS, MOSFLM,autoPROC,DENZO

- 270 PDB entries
- data contain NEMOs
- confirmed beamstop outliers through re-integration
- no twinning or tNCS
- re-refine & rebuild by PDB-REDO 8.03
- model: deposited coordinates
  - + 0.25 shake with phenix.pdbtools
- data: original vs. NEMOs removed



NEMOs-containing – NEMOs-excluded (0.03 e $\text{\AA}^{-3}$ )



## NEMOs Found, What Now

Guide for Low-resolution Cutoff

Minimize Information Loss

Refine Processed Data post-mortem

Objective Metrics for Automatic Beamstop Mask Generation

(Central) Beamstop Mask Free When Reaching 1.0 Accuracy?

Thank Arwen Pearson, Dale Tronrud, Kay Diederichs and Randy Read for Helping us Find NEMOs

Ref.

Read, R. J. (1999). Acta Crystallogr. D Biol. Crystallogr. 55, 1759–1764.

Read, R. J. & McCoy, A. J. (2016). Acta Crystallogr. D Struct. Biol. 72, 375–387.

Campello, R. J. G. B., et al. (2013). Advances in Knowledge Discovery and Data Mining, Vol. 7819. Springer Berlin Heidelberg.

Mishra, S., Monath, N., Boratko, M., Kobren, A. & McCallum, A. (2022). AAAI. 36, 7788–7796.

