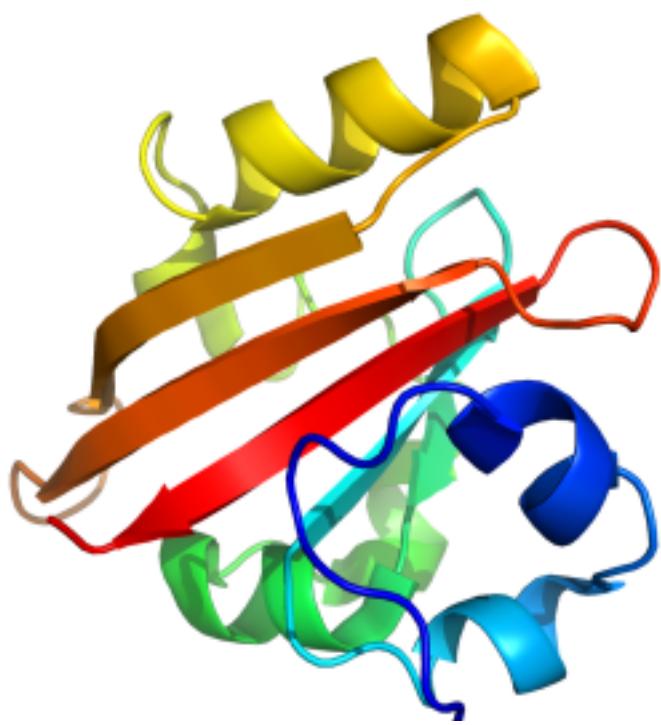
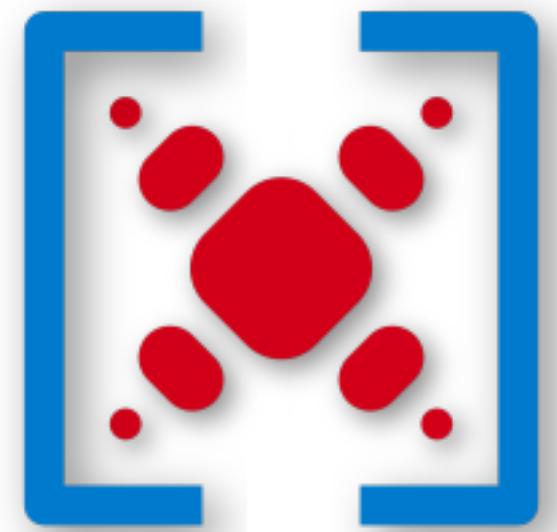


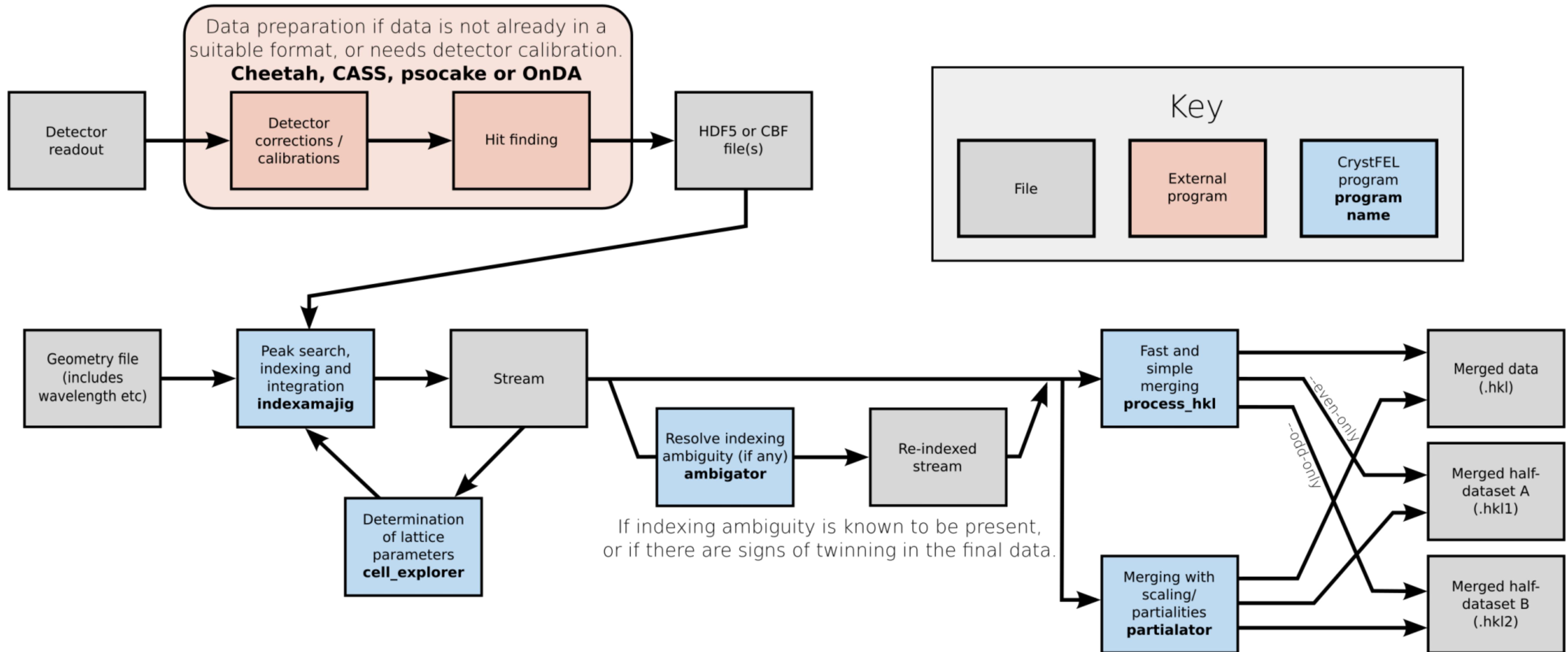
Data processing

Theory of data processing for serial crystallography

Thomas White
Gothenburg CrystFEL workshop
27. January 2020



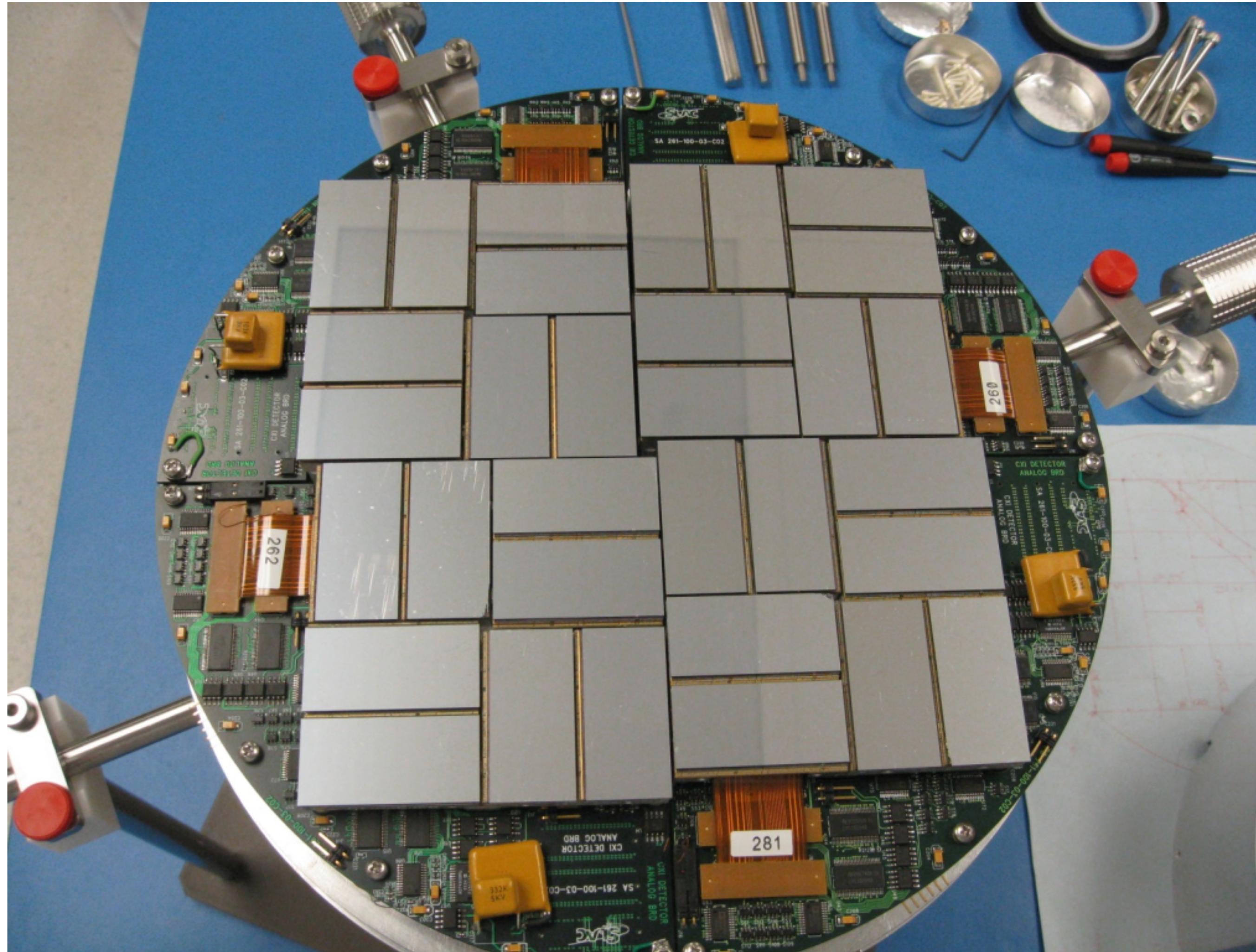
Data processing pipeline



The past and future

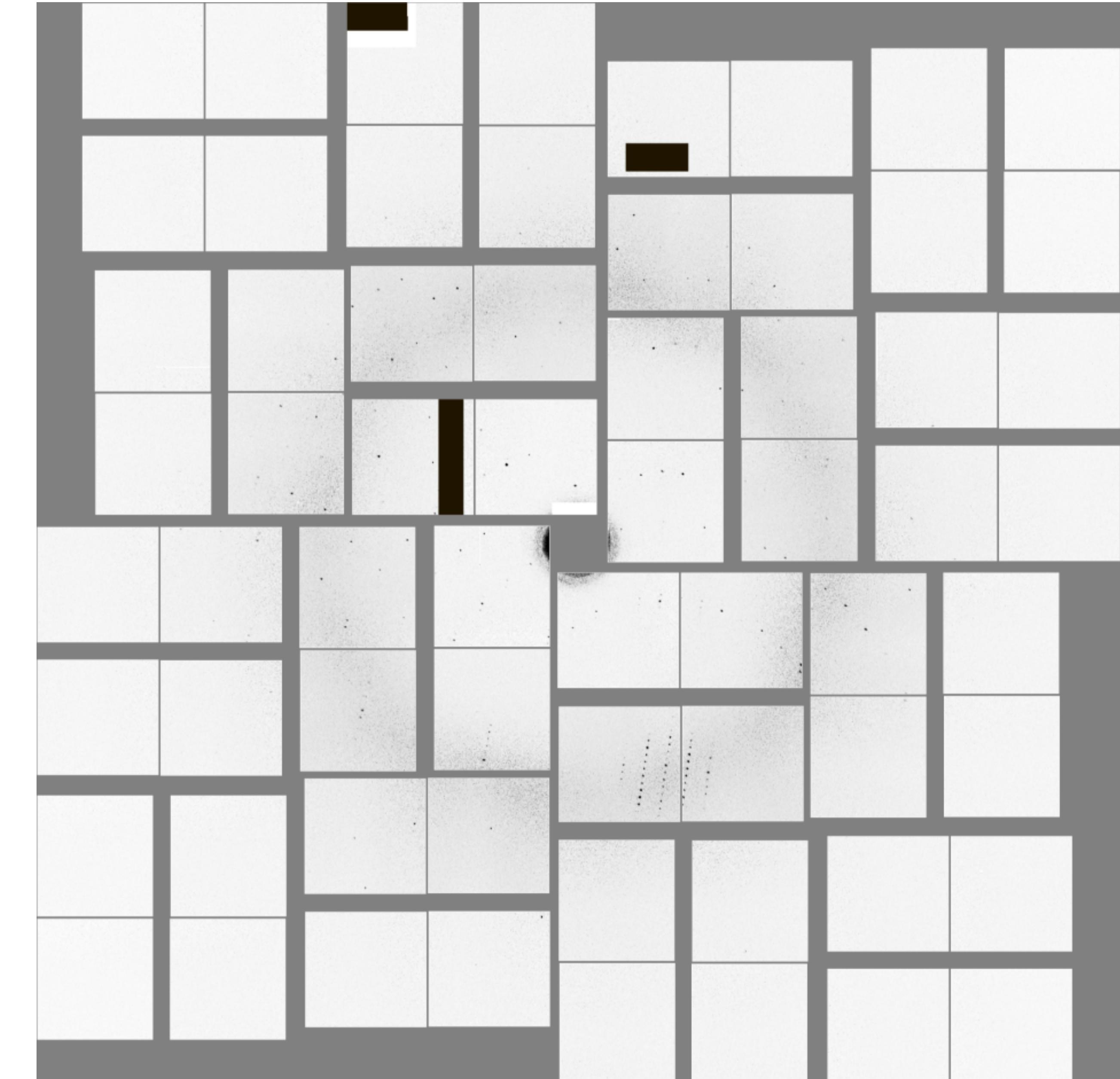
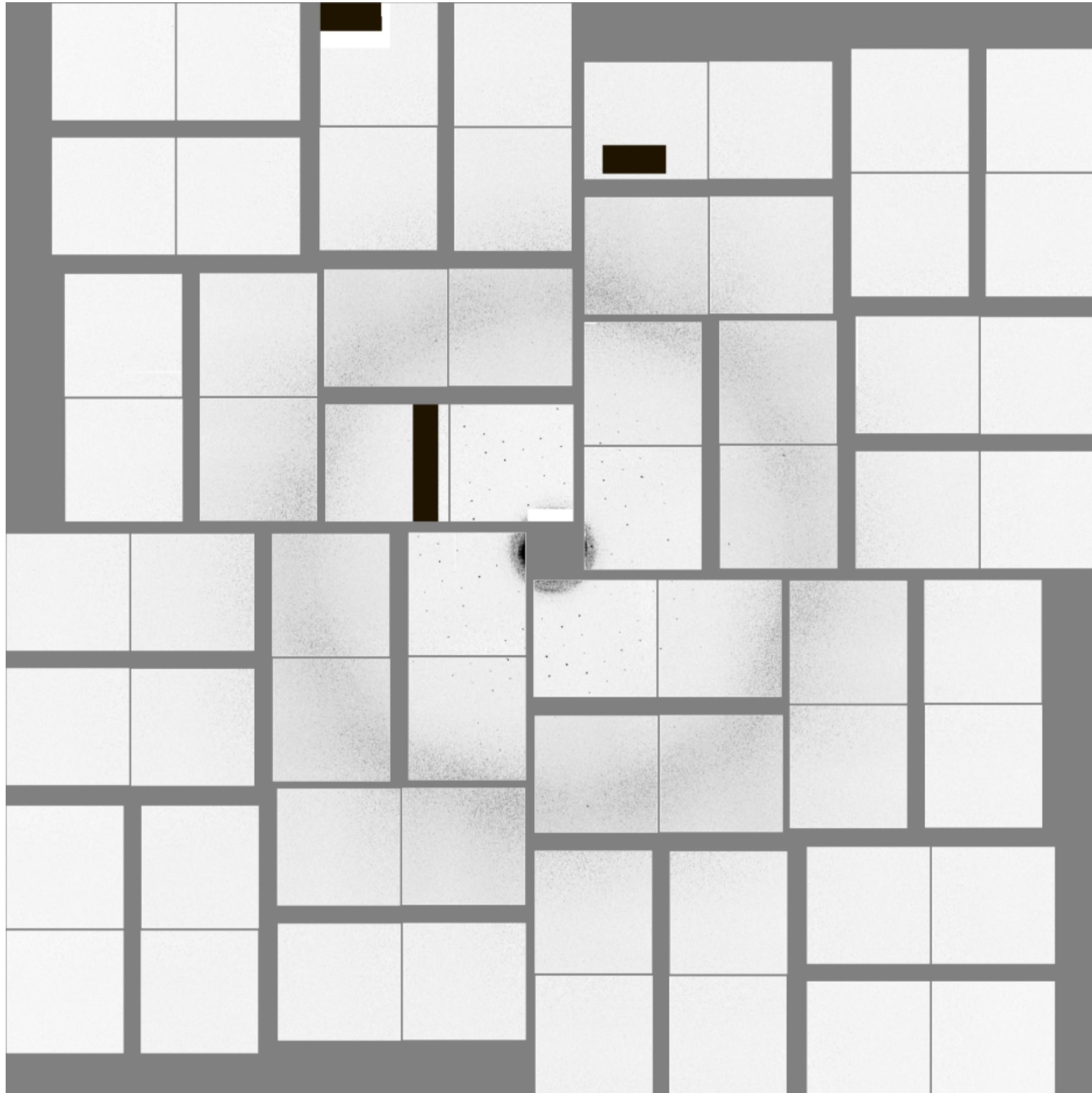
The oscillation method (Arndt & Wonacott, 1977) is now generally accepted as the most efficient means for photographic data collection from crystals with unit-cell dimensions greater than about 100 Å. For very large unit cells (average dimensions greater than about 250 Å), the usual method must be modified to realize its full efficiency – in particular, to cope with situations where only one good photograph, instead of a series of contiguous ones, can be taken within the lifetime of a crystal. Essentially, the resulting ‘one crystal–one photograph’ condition requires a different treatment of partially recorded reflections. This paper describes general aspects of the new procedure and gives a detailed account of its application to data collection from crystalline tomato bushy stunt virus (TBSV: space group $I23$, $a = 383.2$ Å) to 2.9 Å resolution.

Multi-panel detector geometry

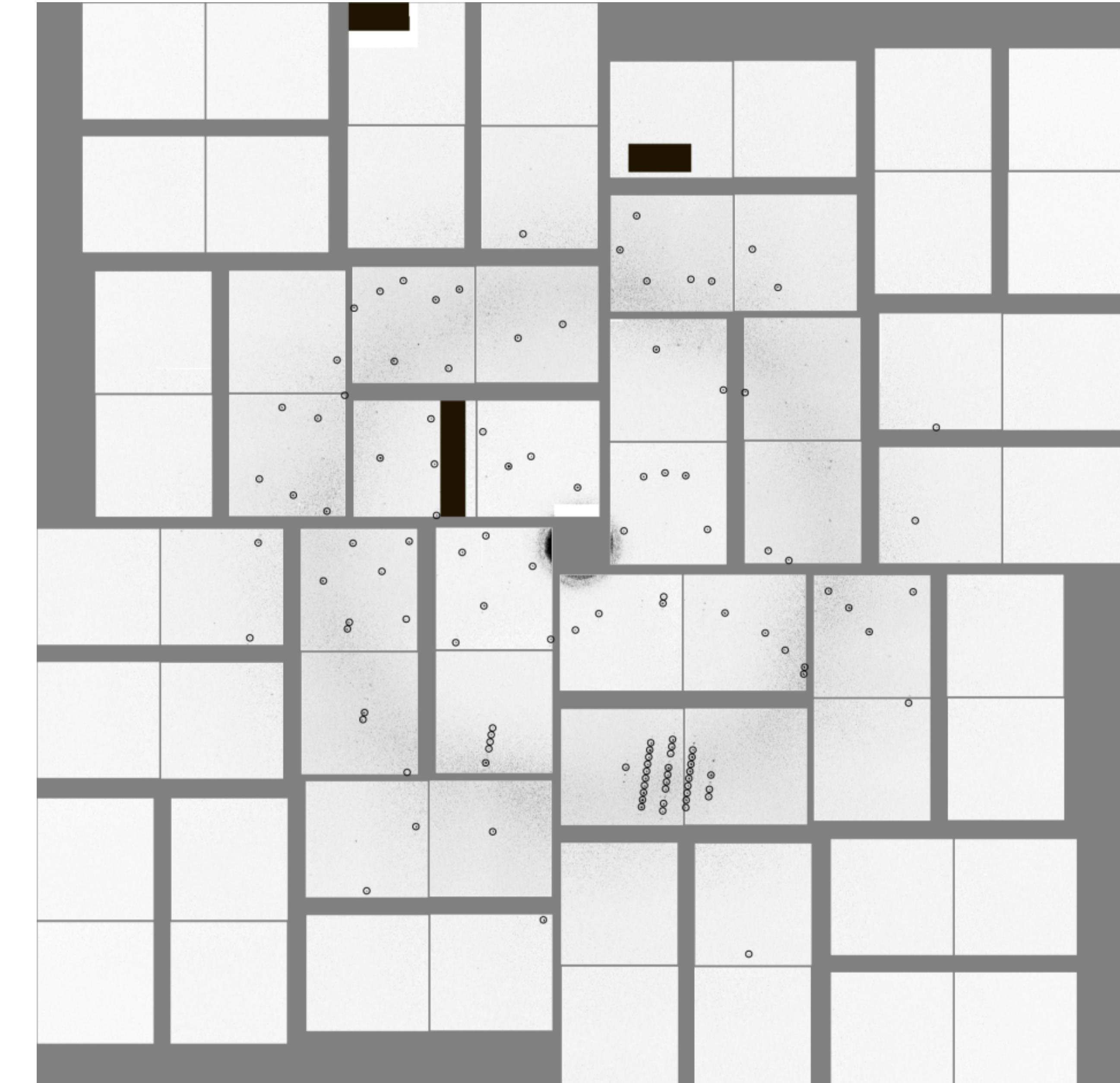
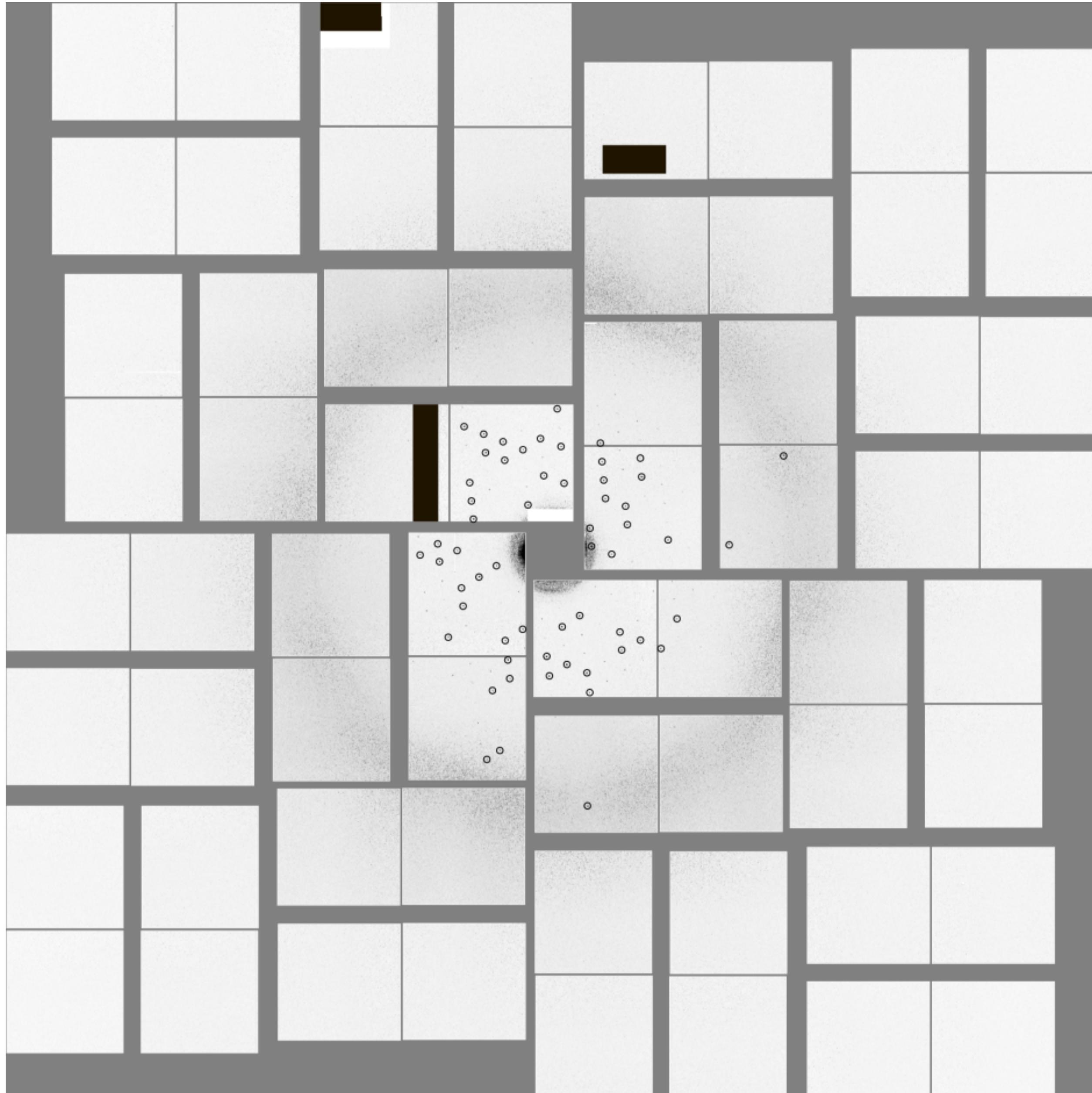


Cornell-SLAC Pixel Array Detector ("CSPAD")
Image: SLAC National Accelerator Laboratory

Diffraction patterns

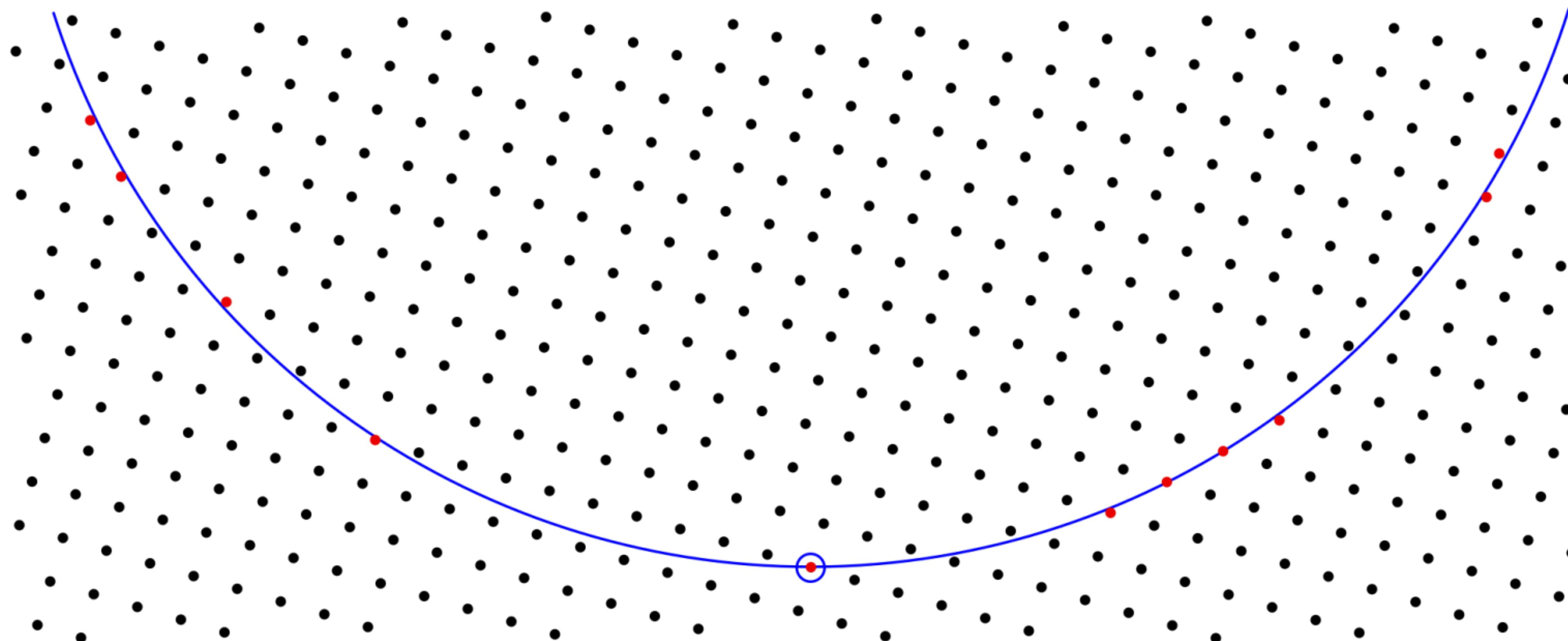


Diffraction patterns

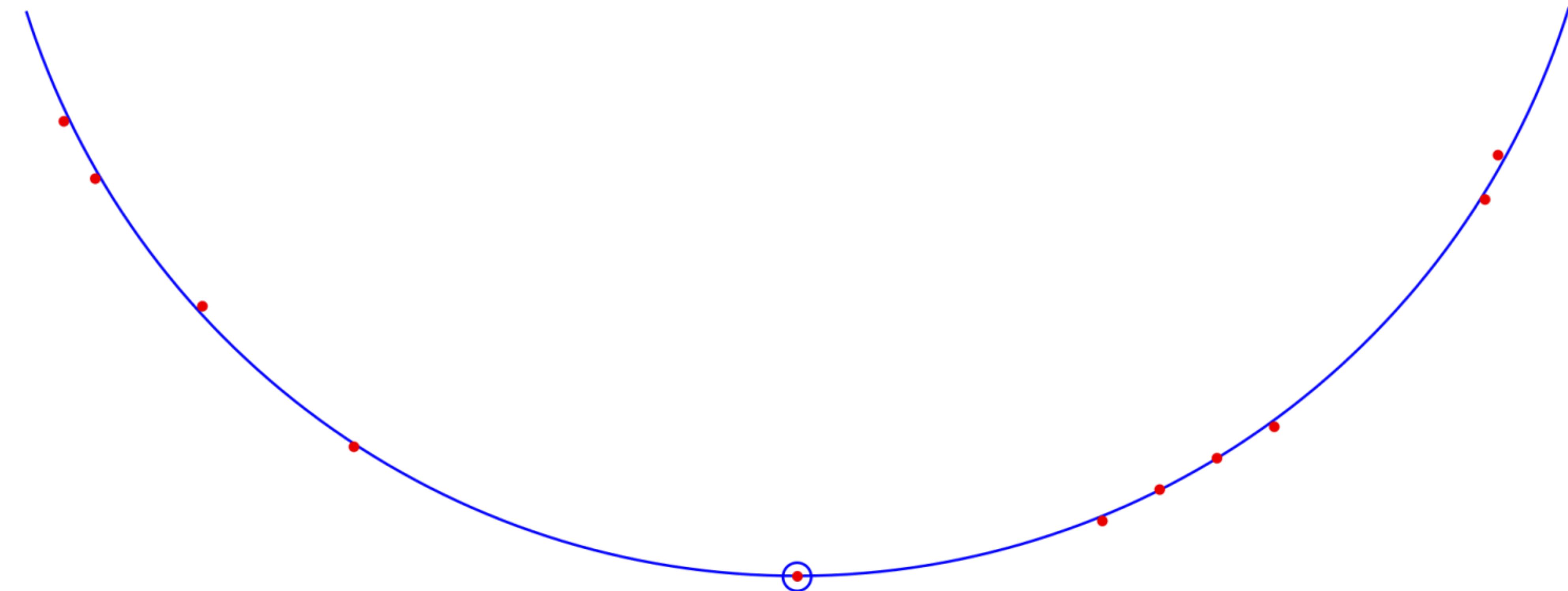


Indexing the diffraction pattern

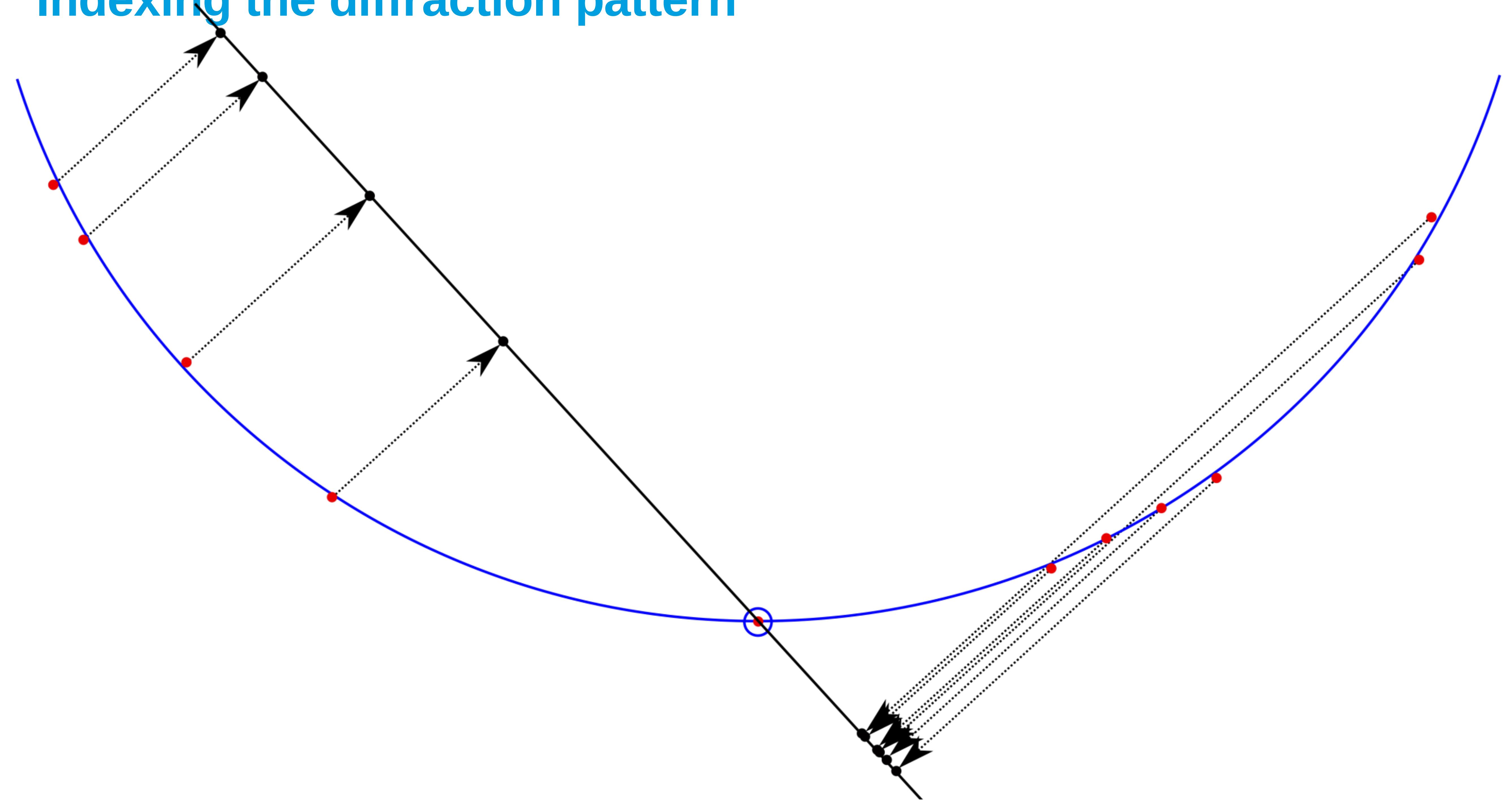
$$2d \sin \theta = \lambda$$



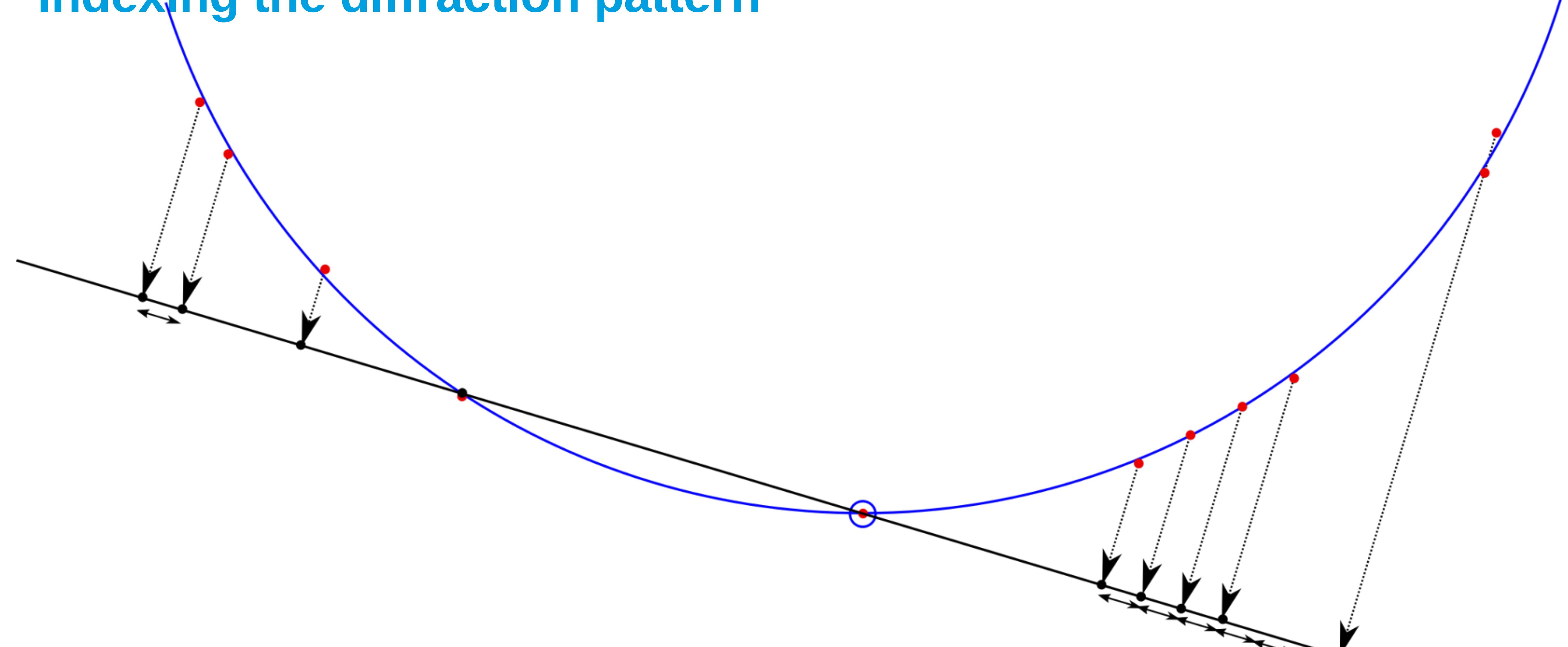
Indexing the diffraction pattern



Indexing the diffraction pattern

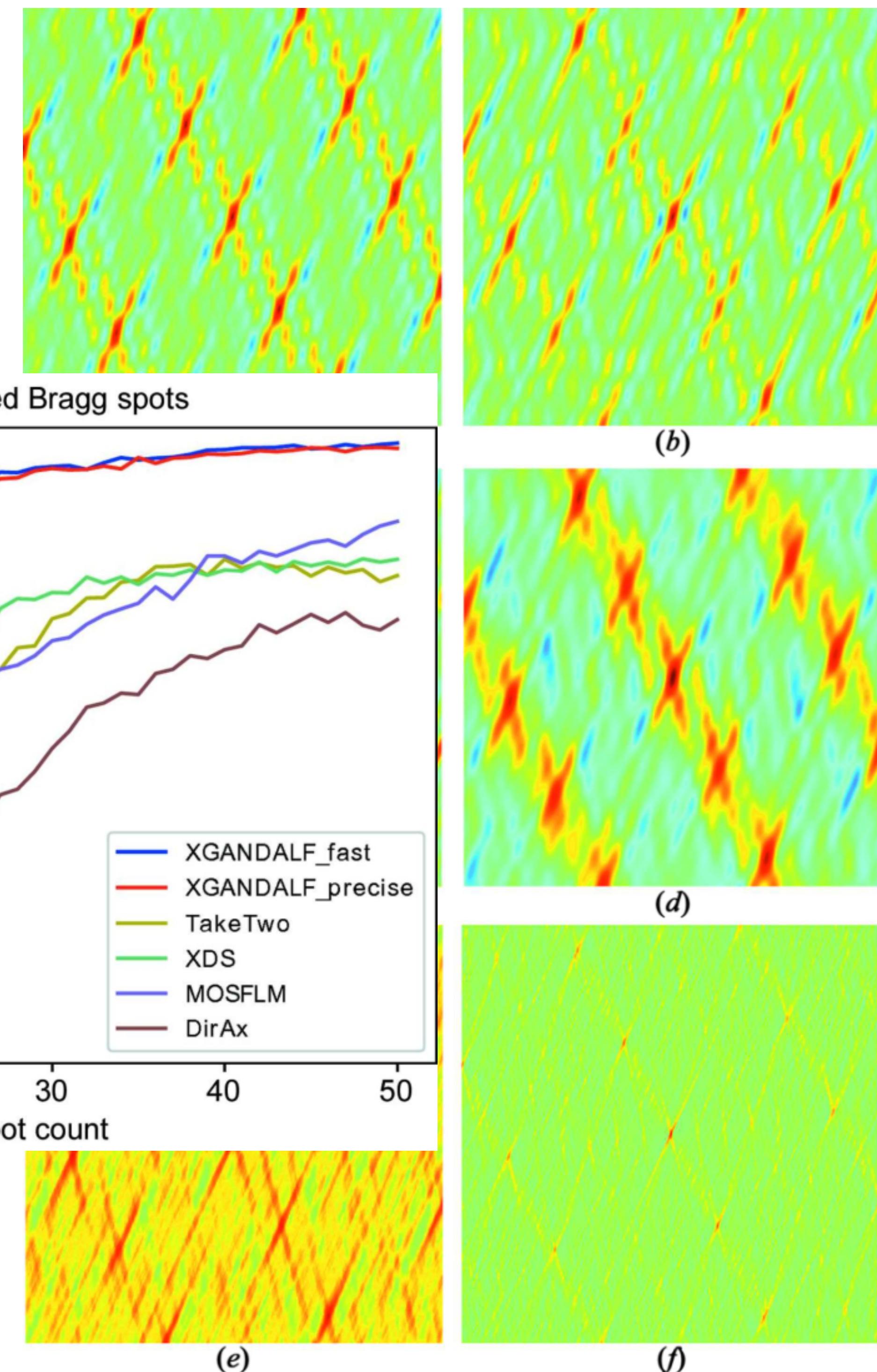
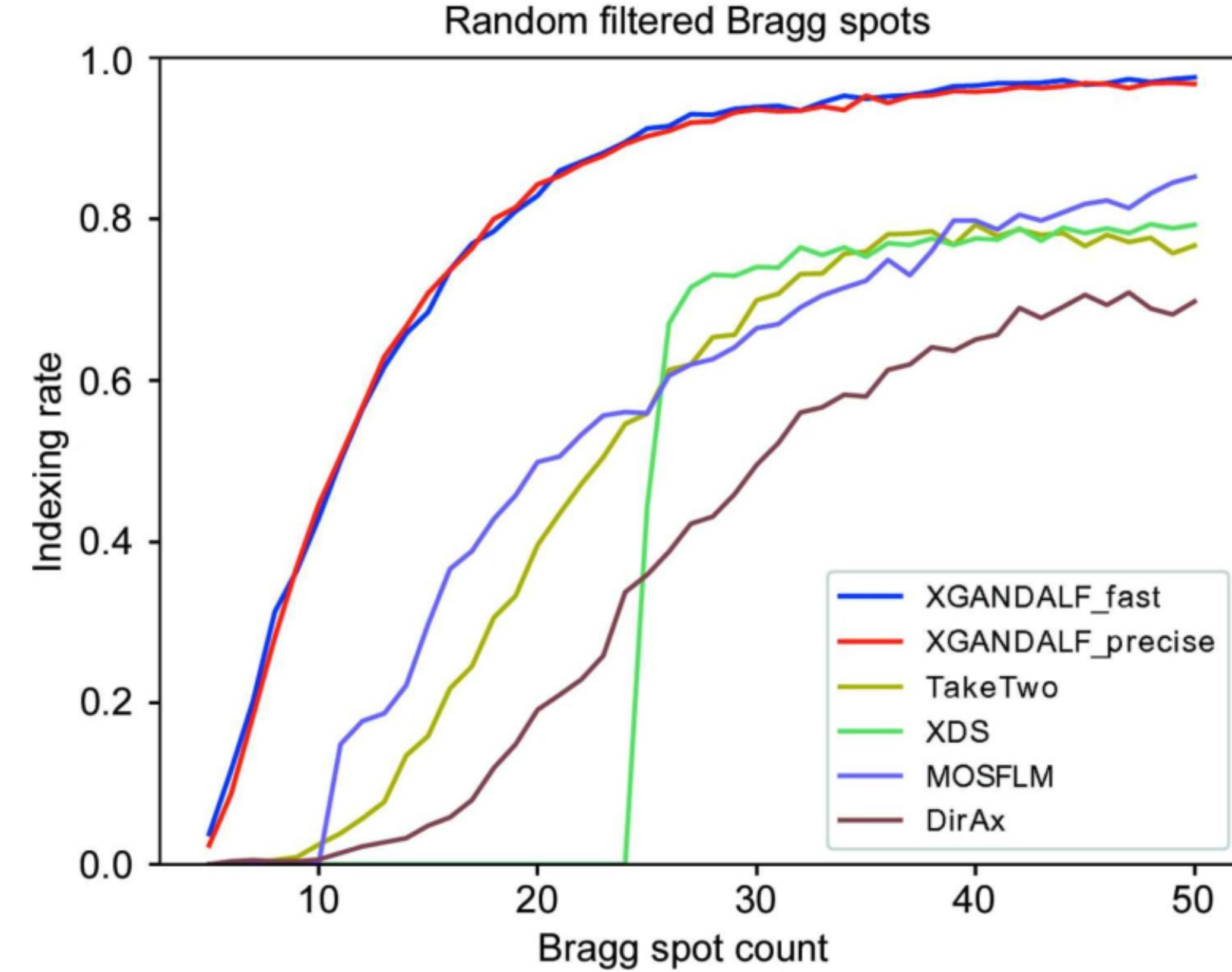
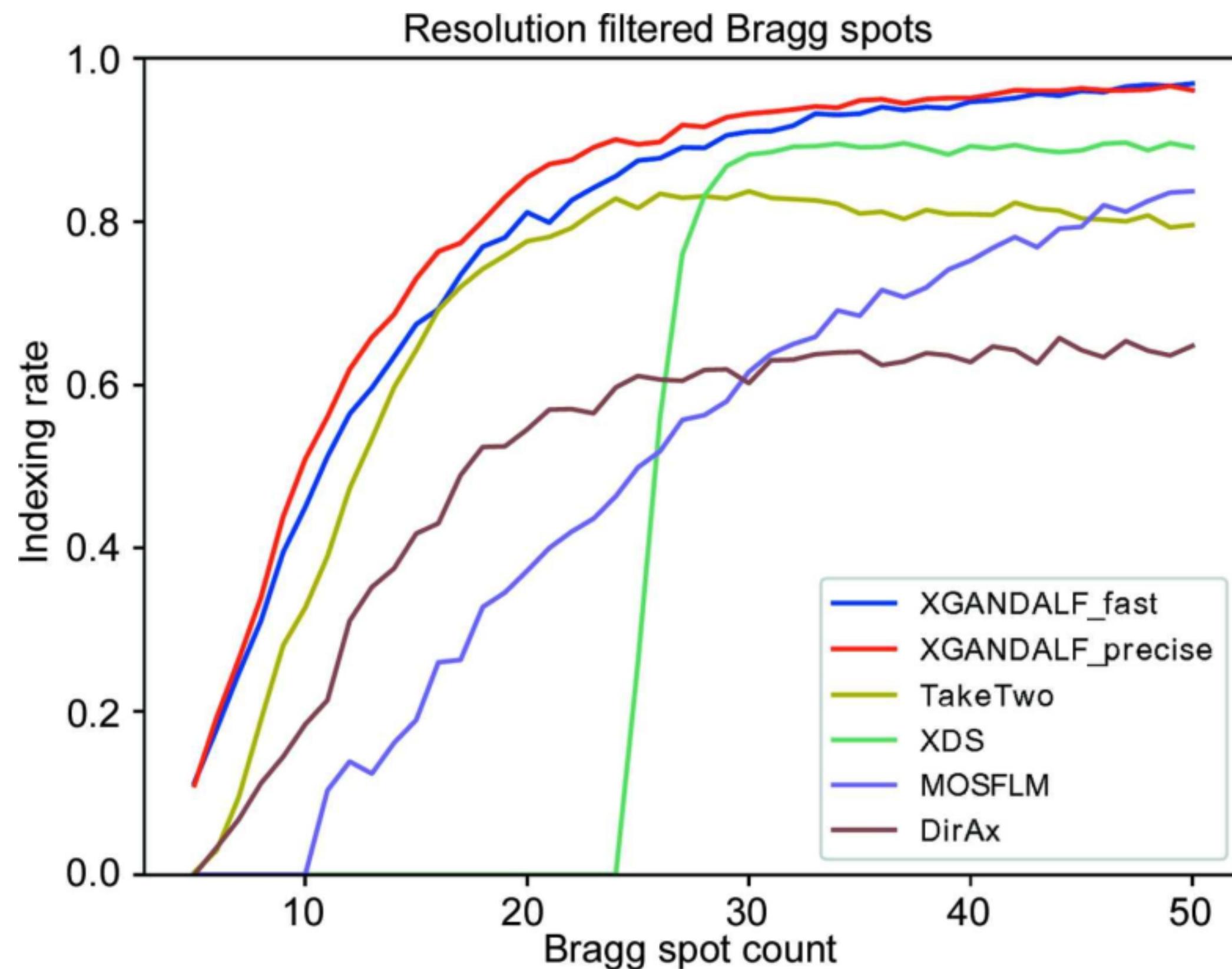


Indexing the diffraction pattern



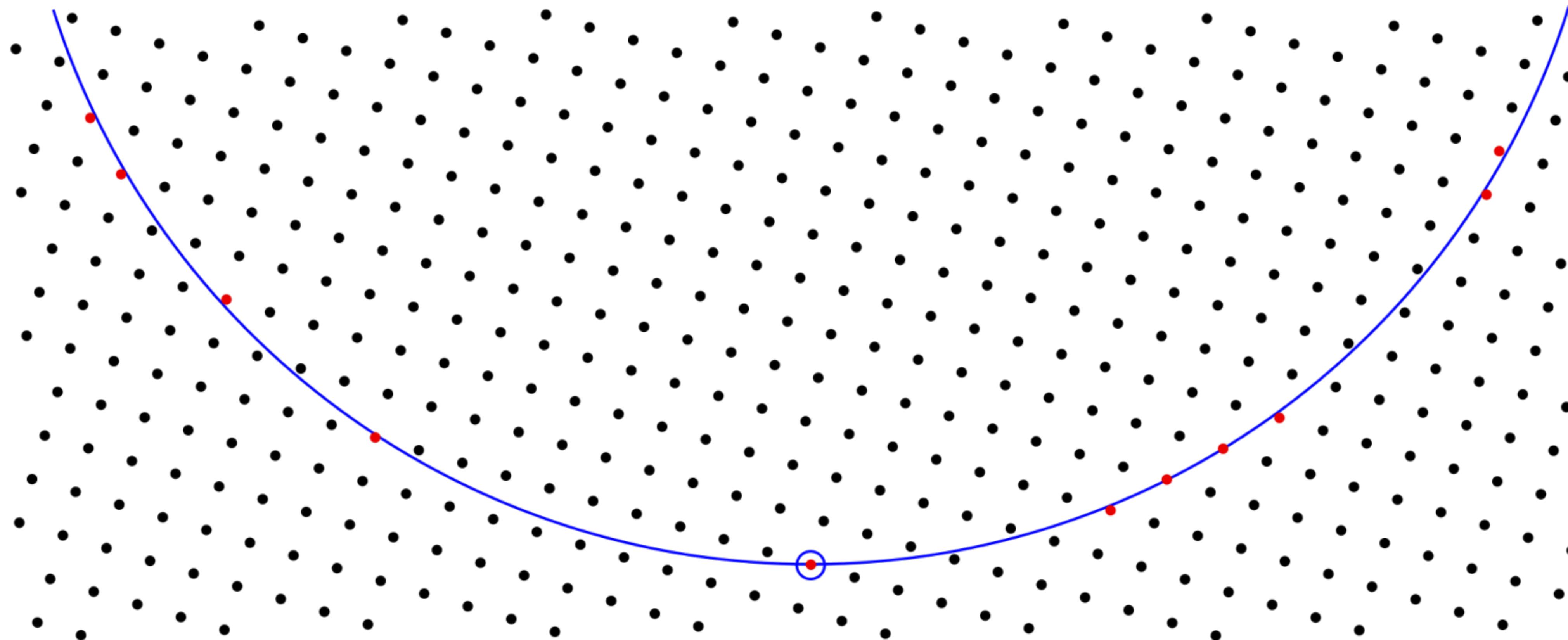
$$\begin{pmatrix} a_x^* & a_y^* & a_z^* \\ b_x^* & b_y^* & b_z^* \\ c_x^* & c_y^* & c_z^* \end{pmatrix}^T = \begin{pmatrix} a_x & a_y & a_z \\ b_x & b_y & b_z \\ c_x & c_y & c_z \end{pmatrix}^{-1}$$

XGANDALF indexer for snapshots

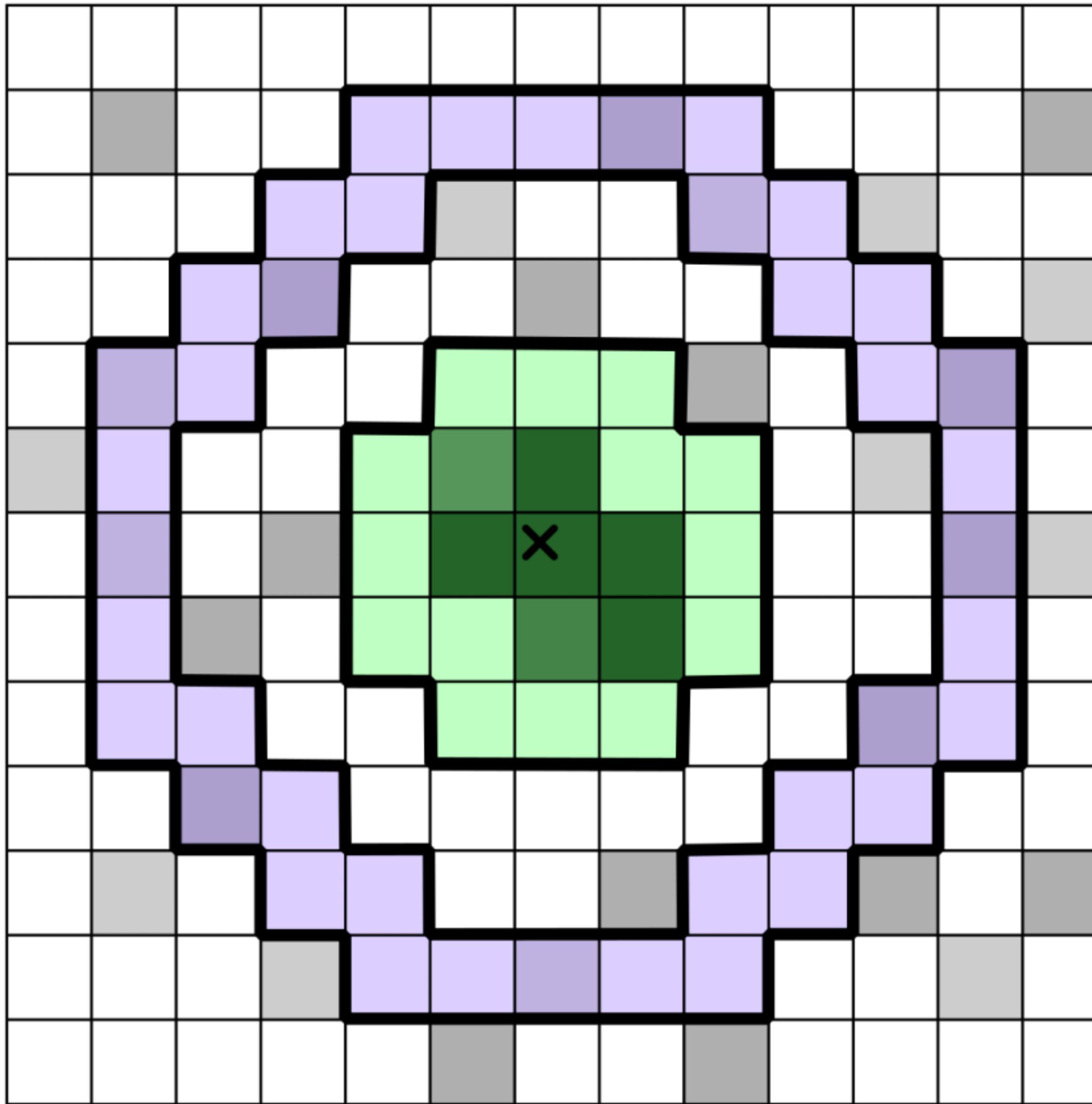


Reflection "prediction"

Important difference from before: all lattice points get projected, not just the ones with peaks

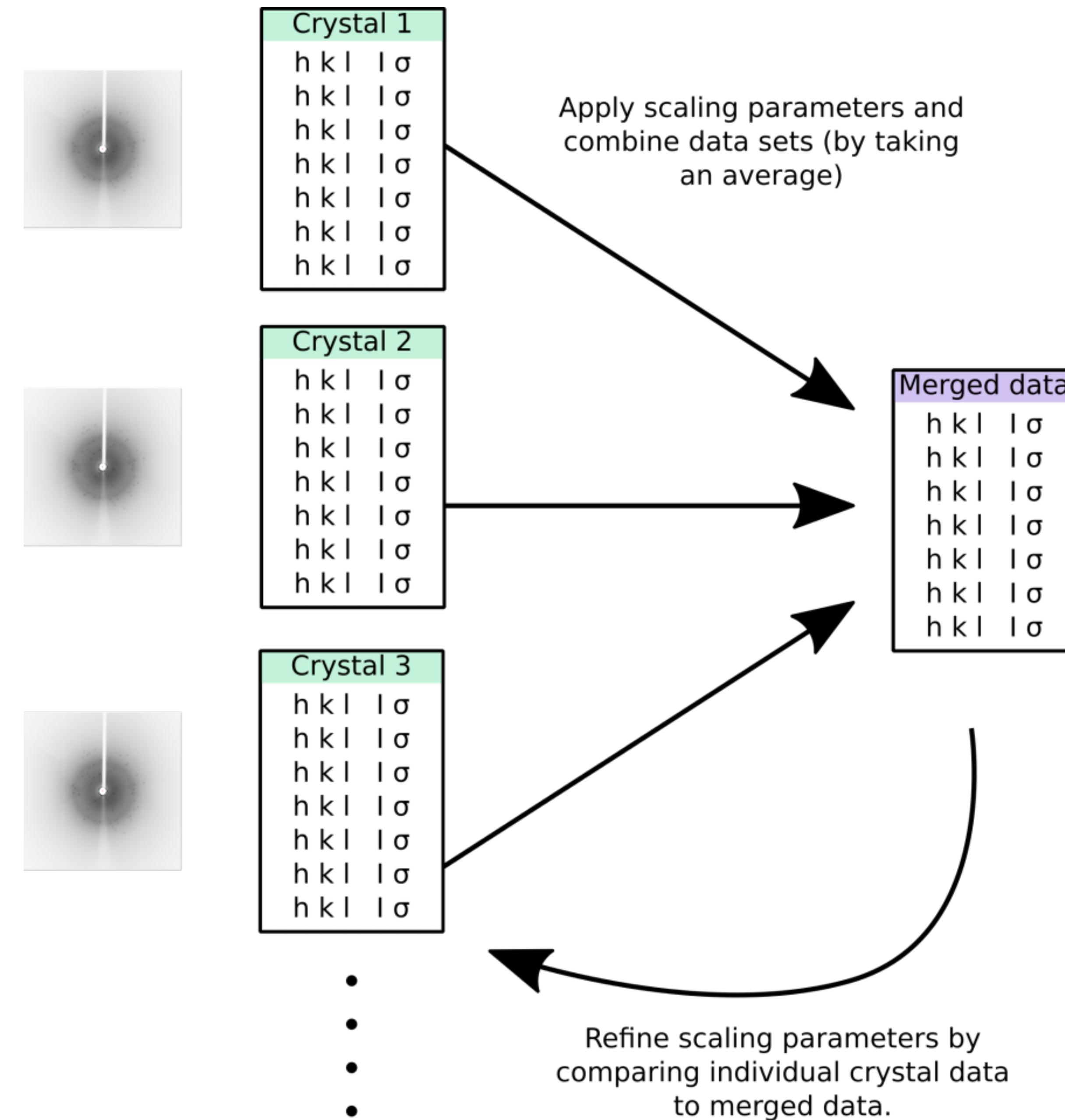


Reflection integration



- ✖ Integration fiducial
- █ Peak region
- █ Background region

Merging the measurements



Indexing ambiguities

Clustering algorithm:

Brehm and Diederichs

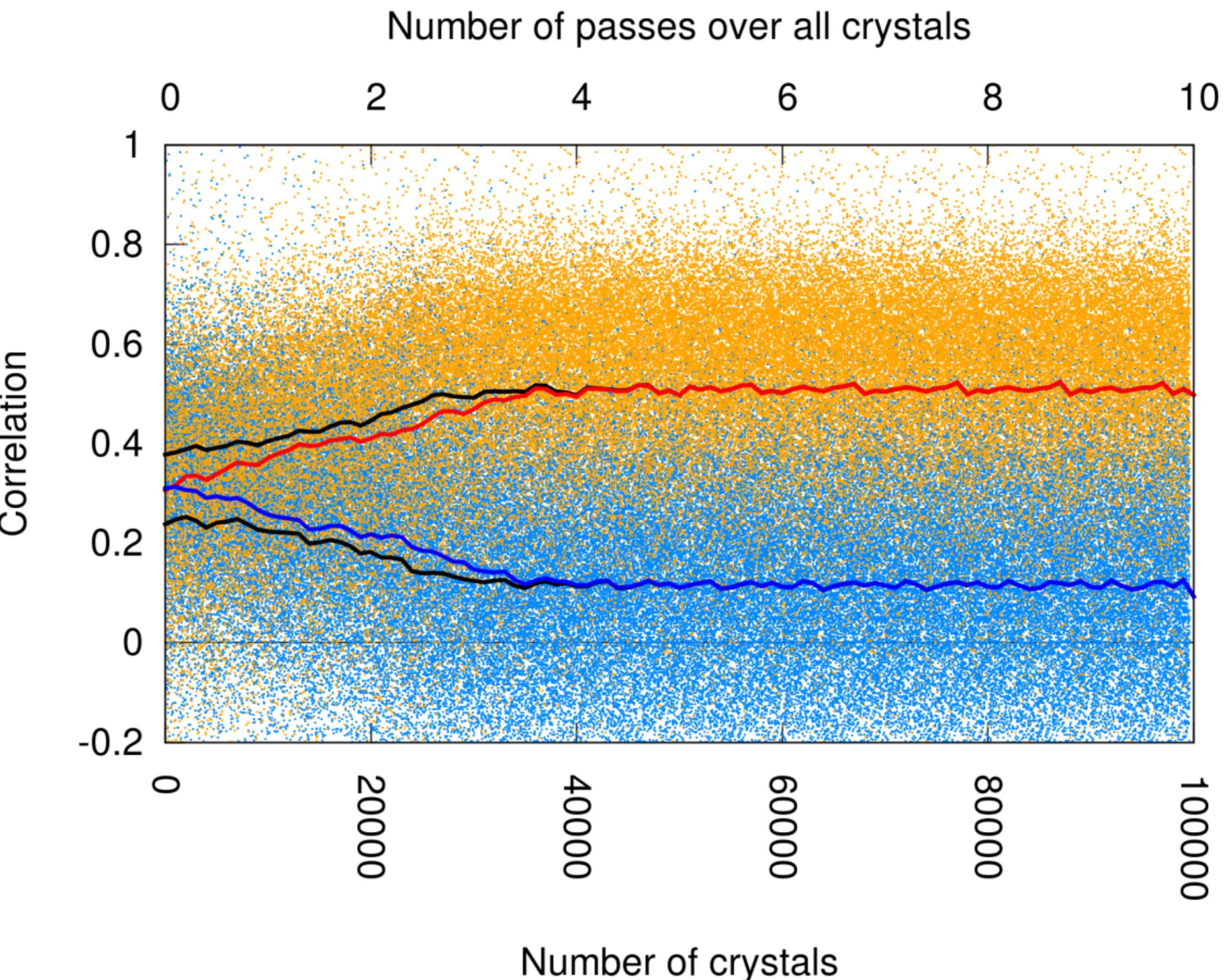
Acta Cryst. D70 (2014) p101-109

Simplified (1D) version:

White, ... Brehm, ... Diederichs, ...

J. Appl. Cryst. 49 (2016) p680

(see section 4)



Figures of merit

From merged data:

Average I/sigma(I)

Number of measurements per merged reflection

Split pattern into two halves, merge each half separately:

R_split

CC^{1/2}

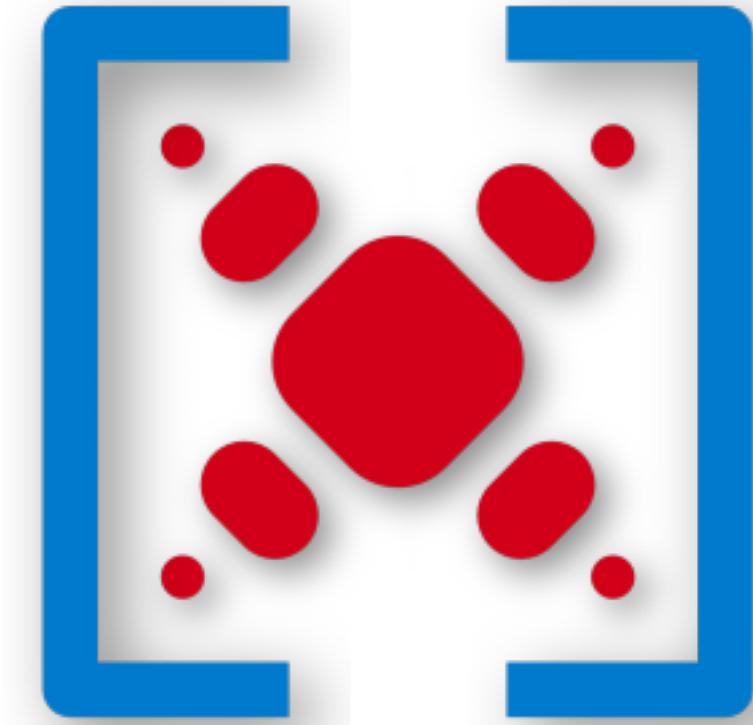
CC*

..... all measure pretty much the same thing

Go back and look at individual patterns again:

$\Delta CC^{1/2}$ ("delta CC half") - see Assmann et al. J. Appl. Cryst. 49 (2016) p1021

What is CrystFEL?



SFX data processing pipeline

- **indexamajig** - Indexing and integration
- **process_hkl** - Monte Carlo merging, simple scaling
- **partialator** - Scaling, partialities, post-refinement
- **ambiguator** - Resolve indexing ambiguities

Figures of merit

- **check_hkl** - Completeness, multiplicity, SNR
- **compare_hkl** - Compare reflection data

Simulation tools

- **partial_sim** - Calculate partial reflections
- **pattern_sim** - Fourier simulation

Visualisation tools

- **hdfsee** - Display diffraction patterns
- **render_hkl** - Visualise reflection data
- **cell_explorer** - Exampline unit cell distributions

Utilities / other

- **get_hkl** - Twin, example, add noise
- **list_events** - Create event lists
- **geoptimiser** - Refine detector geometry
- **whirligig** - Find rotation series
- **cell_tool** - Unit cell manipulations

Lots of supporting scripts (mostly in Python)

Documentation (manual pages, --help, example files)

Lots more on website (tutorial, demonstration data, citation lists, best practice guidelines, ...)

CrystFEL 0.9.0

- New `cell_tool`
- Interface to PinkIndexer
- Reflection prediction for electrons and wide bandwidth X-rays
- Arbitrary spectrum API
- ZMQ/MsgPack interface
- New cell comparison algorithm
- Faster deltaCChalf calculation
- Arbitrary polarisation/degree of polarisation (needed for LCLS-II)
- Improvements to multi-lattice indexing
- Handle CBF files natively

Advanced topics

Indexing ambiguities

... how to know there's an ambiguity, and how to fix it

Partiality modelling and post-refinement

Special considerations for TR-SX, MAD and MIR

Recognise when there's an ambiguity

Warning signs:

Point group 3, 4, 6, 312, 321 or 23

(432, 622, 32, 422, 222, 2 and 1 are OK)

.... but you'd have to know the structure already

Twin warnings from Phenix, CCP4 etc, or check_hkl --ltest :

```
$ check_hkl test.hkl -p 4ET8.pdb --ltest --ignore-negs
```

Using symmetry from reflection file: 4/mmm

Discarded 0 reflections (out of 3033) with I/sigma(I) < -inf

Discarded 672 reflections because they had negative intensities.

2631 pairs

$\langle |L| \rangle = 0.512$ (ideal untwinned 0.500, twinned 0.375)

$\langle L^2 \rangle = 0.344$ (ideal untwinned 0.333, twinned 0.200)

\$

Twin warnings mean...

1. There's an indexing ambiguity which you need to resolve
 2. The crystals are really (physically) twinned
 3. There's some weird data artifact (e.g. bad detector calibration)
- any combination of the above.

Symmetry classification for serial crystallography

Triclinic lattice

$\bar{1}$	1	$P\bar{1}$	$P1$
-----------	---	------------	------

Monoclinic lattice

m	$P2, P2_1, C2$	Pm, Pc, Cm, Cc $P2/m, P2_1/m, C2/m, P2/c, P2_1/c, C2/c$
2	$2/m$	

Orthorhombic lattice

$mm2$		$Pmm2, Pmc2_1, Pcc2, Pma2, Pca2_1, Pnc2, Pmn2_1, Pba2, Pna2_1, Pnn2, Cmm2, Cmc2_1, Ccc2, Amm2, Aem2, Ama2, Aea2, Fmm2, Fdd2, Imm2, Iba2, Im2$
222	mmm	$P222, P222_1, P2_12_12, P2_12_12_1, C222_1, C222, F222, I222, I2_12_12_1$ $Pmmm, Pnnn, Pccm, Pban, Pmma, Pnna, Pmna, Pcca, Pbam, Pccn, Pbcm, Pnnm, Pmmn, Pbcn, Pbca, Pnma, Cmcm, Cmce, Cmmm, Cccm, Cmme, Ccce, Emmm, Fddd, Imm, Ibam, Ibca, Imma$

Tetragonal lattice

4	$\bar{4}$			4mm	$P4, P4_1, P4_2, P4_3, I4, I4_1$	$\bar{P}4, I\bar{4}$			$P4mm, P4bm, P4_2cm, P4_{2}nm, P4cc, P4nc, P4_{2}mc, P4_{2}bc, I4mm, I4cm, I4_1md, I4_1cd$
	$\bar{4}2m$	$\bar{4}m2$	$4/m$			$P\bar{4}2m, P\bar{4}2c, P\bar{4}2_1m, P\bar{4}2_1c, I\bar{4}2m, I\bar{4}2d$	$P\bar{4}m2, P\bar{4}c2, P\bar{4}b2, P\bar{4}n2, I\bar{4}m2, I\bar{4}c2$	$P4/m, P4_2/m, P4/n, P4_{2}/n, I4/m, I4_1/a$	
422	$4/mmm$				$P422, P42_12, P4_122, P4_12_12, P4_22, P4_2_12, P4_322, P4_32_12, I422, I4_122$	$P4/mmm, P4/mcc, P4/nbm, P4/nnc, P4/mbm, P4/mnc, P4/nmm, P4/ncc, P4_2/mmc, P4_2/mcm, P4_2/nbc, P4_2/nmm, P4_2/mbc, P4_2/mnm, P4_2/nmc, P4_2/ncm, I4/mmm, I4/mcm, I4_1/AMD, I4_1/ACD$			

Rhombohedral lattice

3	$\bar{3}$	3m	$R3 (H3)$	$R\bar{3} (H\bar{3})$
32		$\bar{3}m$	$R32 (H32)$	$R\bar{3}m (H\bar{3}m)$

Hexagonal lattice

3	$\bar{3}$	$3m1 \quad \bar{6} \quad 31m$	$6mm \quad 6/m$	$P3, P3_1, P3_2$	$\bar{P}3 \quad P3m1, P3c1 \quad P\bar{6} \quad P31m, P31c$	$P6mm, P6cc, P6_3cm, P6_3mc$
6	$312 \quad 321$	$\bar{3}m1 \quad \bar{6}m2 \quad \bar{6}2m \quad \bar{3}1m$	$622 \quad 6/mmm$	$P6, P6_1, P6_5, P6_2, P6_4, P6_3, P312, P3_12, P3_21, P3_212$	$P\bar{3}m1, P\bar{3}c1 \quad P\bar{6}m2, P\bar{6}c2 \quad P\bar{6}2m, P\bar{6}2c \quad P\bar{3}1m, P\bar{3}1c$	$P6/mmm, P6/mcc, P6_3/mcm, P6_3/mmc$

Cubic lattice

23	$\bar{4}3m$	$m\bar{3}$	$P23, F23, I23, P2_13, I2_13$	$P\bar{4}3m, F\bar{4}3m, I\bar{4}3m, P\bar{4}3n, F\bar{4}3c, I\bar{4}3d$	$Pm\bar{3}, Pn\bar{3}, Fm\bar{3}, Fd\bar{3}, Im\bar{3}, Pa\bar{3}, Ia\bar{3}$
432		$m\bar{3}m$	$P432, P4_232, F432, F4_132, I432, P4_32, P4_132, I4_132$	$Pm\bar{3}m, Pn\bar{3}n, Pm\bar{3}n, Pn\bar{3}m, Fm\bar{3}m, Fm\bar{3}c, Fd\bar{3}c, Im\bar{3}m, Ia\bar{3}d$	

Spotting "accidental" ambiguities

```
$ cell_tool 4ET8.pdb --find-ambi -y 422
```

Input unit cell: 4ET8.pdb

-----> The input unit cell:

tetragonal P, unique axis c, right handed.

a	b	c	alpha	beta	gamma
---	---	---	-------	------	-------

79.00	79.00	38.00	A	90.00	90.00	90.00	deg
-------	-------	-------	---	-------	-------	-------	-----

Looking for ambiguities up to 3x each lattice length.

This will take about 30 seconds. Please wait...

[....]

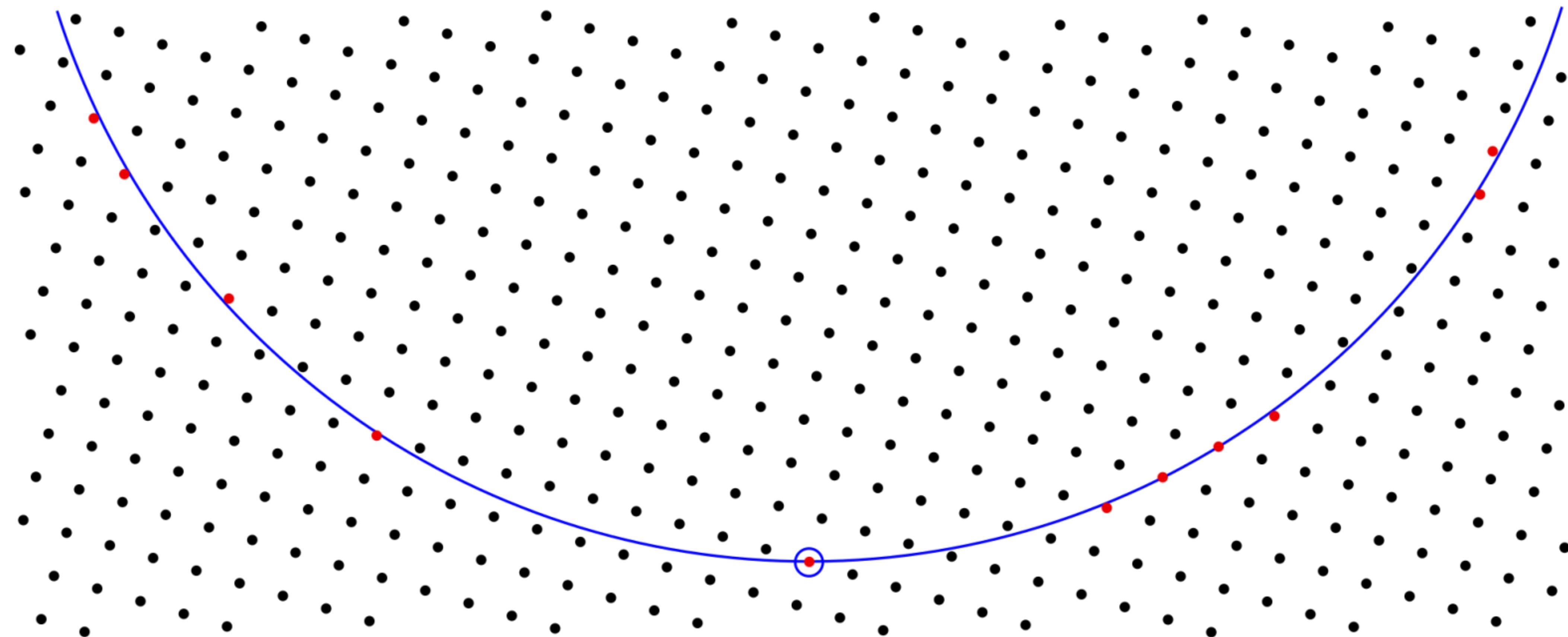
Observed symmetry operations:

Observed : hkl -h,-k,l -h,k,-l -k,-h,-l k,-h,l -k,h,l k,h,-l h,-k,-l

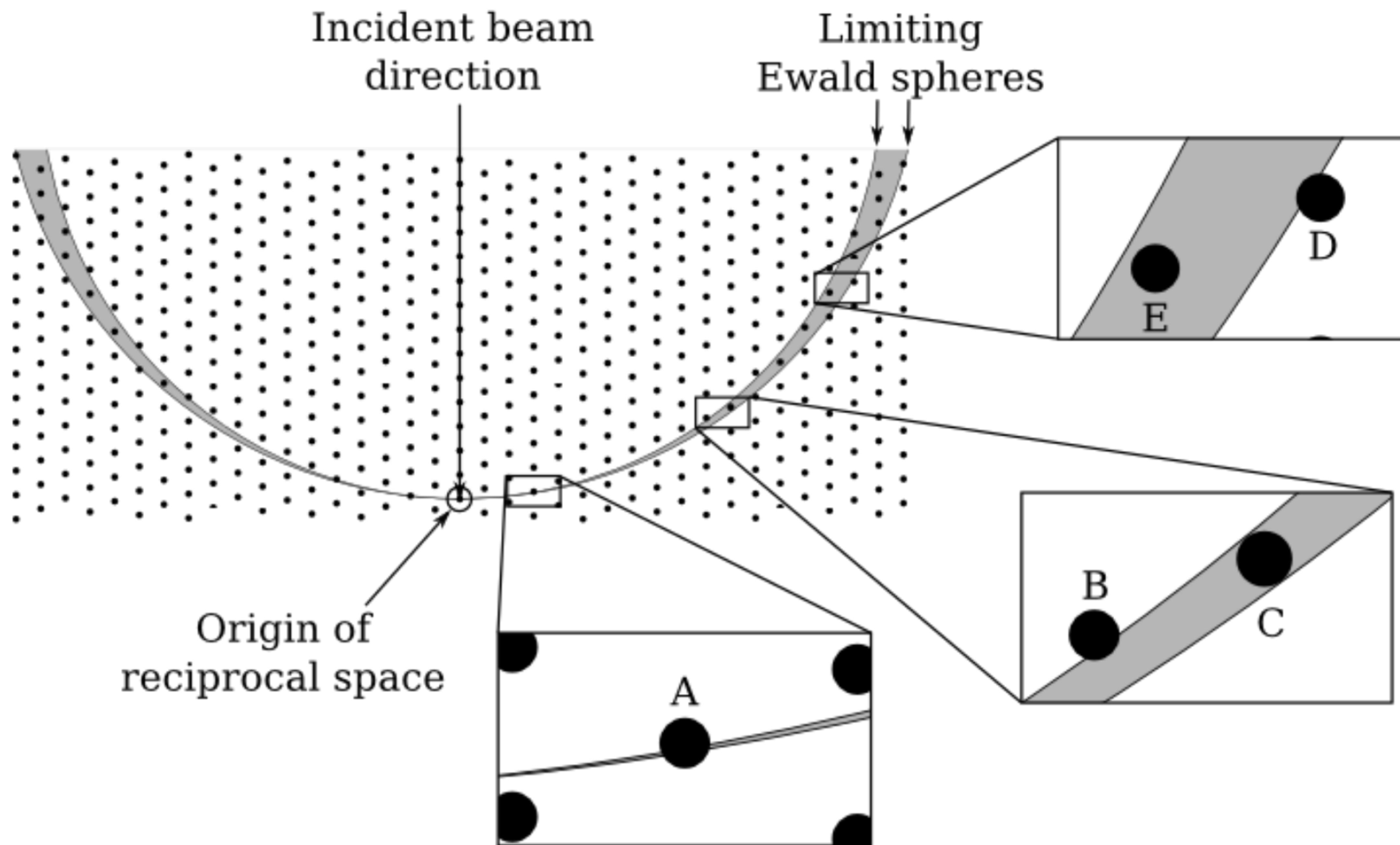
Ambiguity operations:

Observed -> 422 :

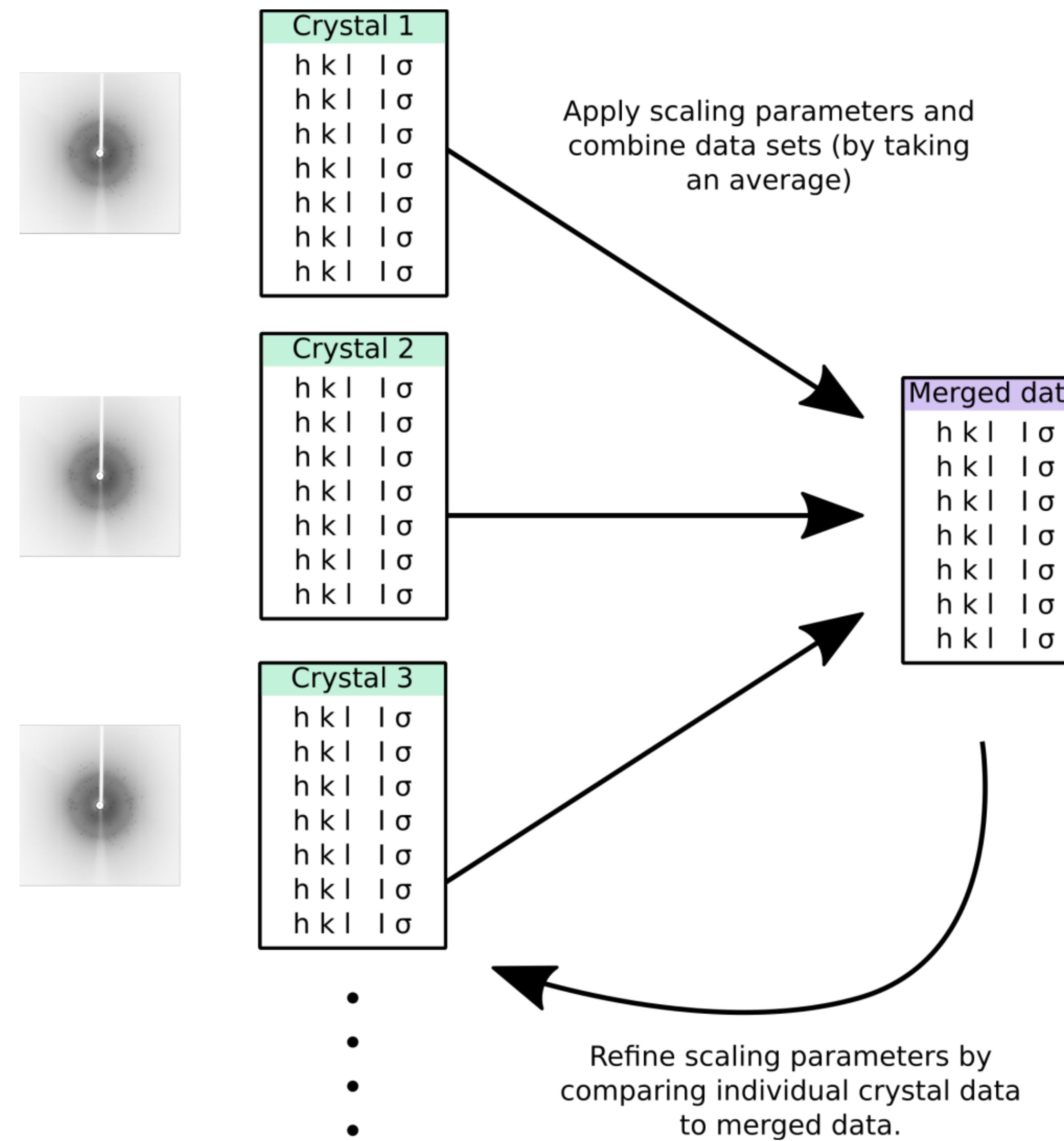
Reflection partiality



Reflection partiality



Post-refinement



Partiality models available in CrystFEL (v0.9.0)

partialator --model=unity

Set all partialities to 1 (no modelling)

partialator --model=xsphere

Numerical integration across arbitrary spectrum

See Ginn et al, Acta Cryst. D71, p1400 (2015)

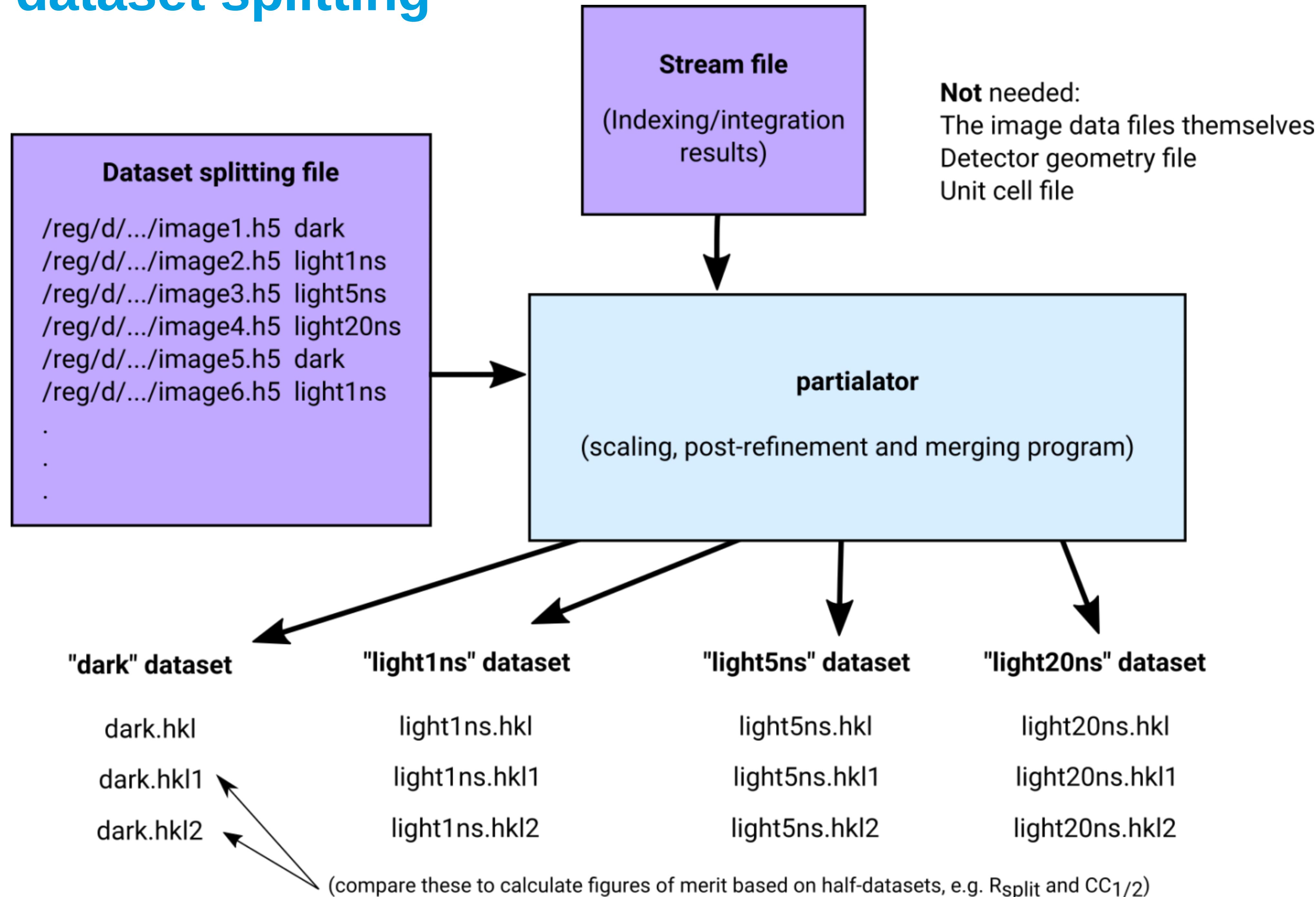
partialator --model=offset

Thin Ewald sphere, like Kabsch, Acta Cryst. D70, p2204 (2014)

partialator --model=ggpm

Analytically-defined model based on Gaussian overlap integrals
(Needs X-ray spectrum to be an analytical function)

Custom dataset splitting



Three-way ambiguity in human melatonin receptor

Tetragonal P

$a = 122.3 \text{ \AA}$

$b = 122.3 \text{ \AA}$

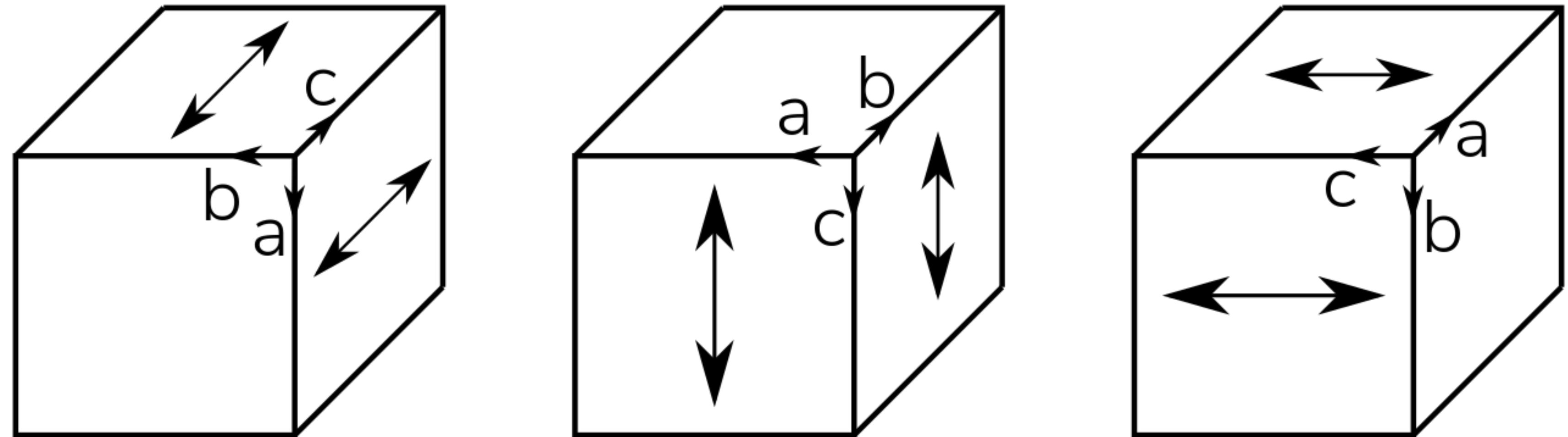
$c = 122.8 \text{ \AA}$

$\alpha = 90^\circ$

$\beta = 90^\circ$

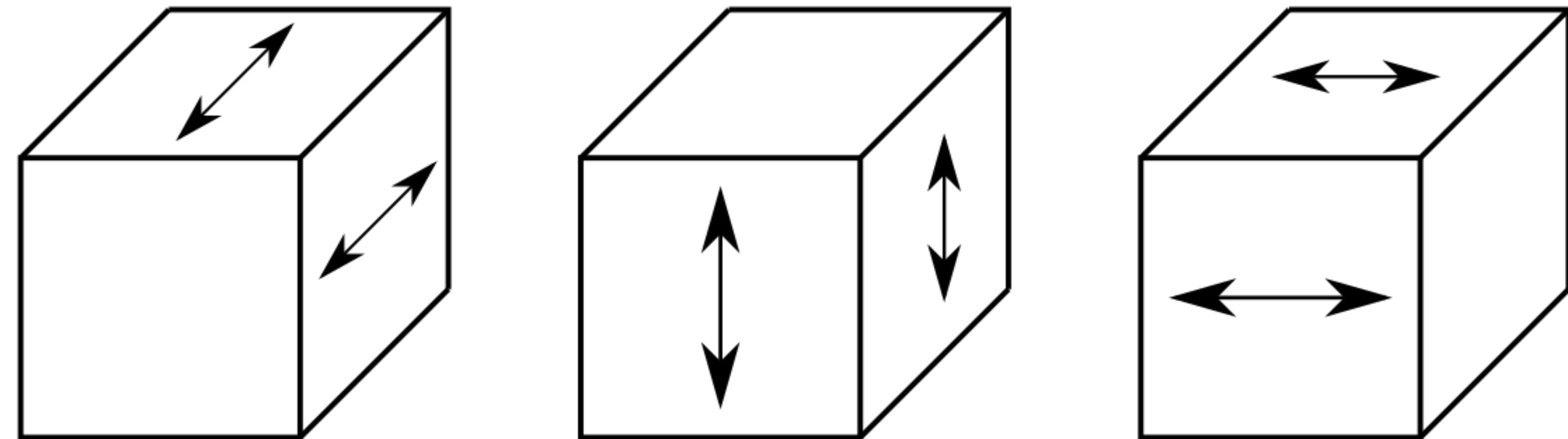
$\gamma = 90^\circ$

Indexing ambiguity: $hkl \rightarrow klh \rightarrow lhk \rightarrow hkl \rightarrow \dots$
One indexing ambiguity, three possibilities



B. Stauch, L. C. Johannson, J. D. McCorry, N. Patel, G. Won Han, X.-P. Huang, C. Gati, A. Batyuk, S. T. Slocum, A. Ishchenk, W. Brehm, T. A. White, N. Michaelian, C. Madsen, L. Zhu, T. D. Grant, J. M. Grandner, A. Shiriaeva, R. H. J. Olsen, A. R. Tribo, S. Yous, R. C. Stevens, U. Weierstall, V. Katritch, B. L. Roth, W. Liu, V. Cherezov.
Nature 569 (2019) p284

Resolving a three-way cyclic ambiguity



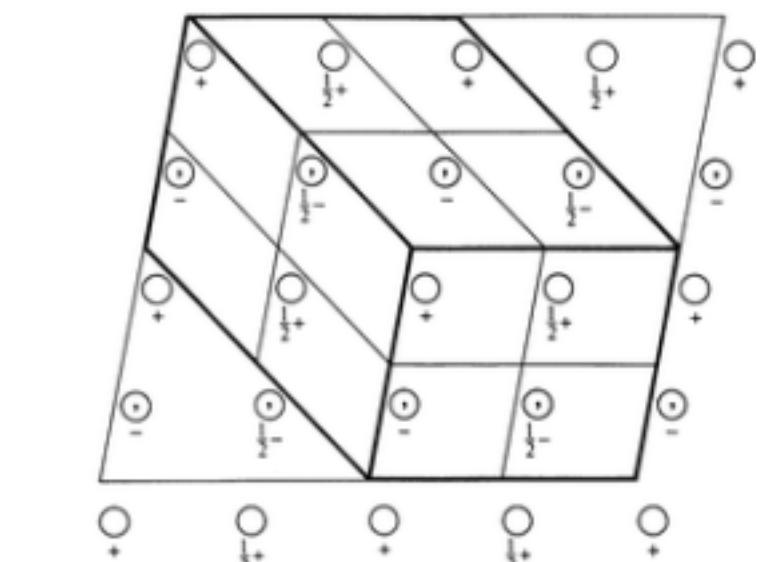
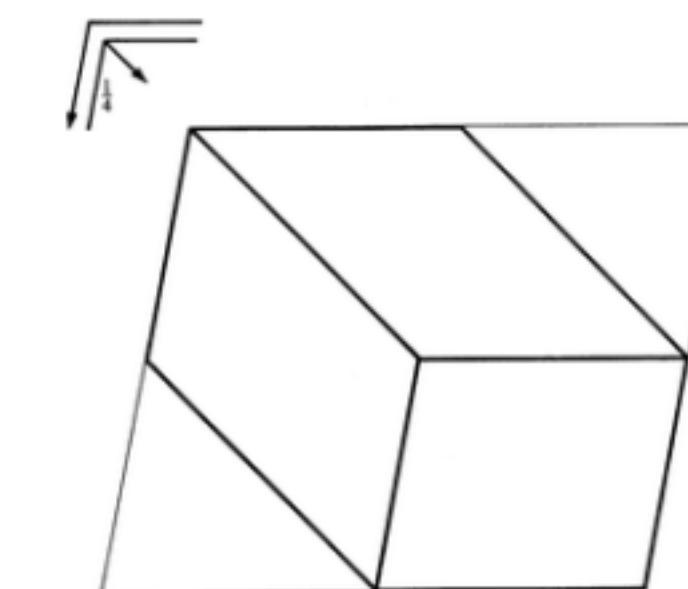
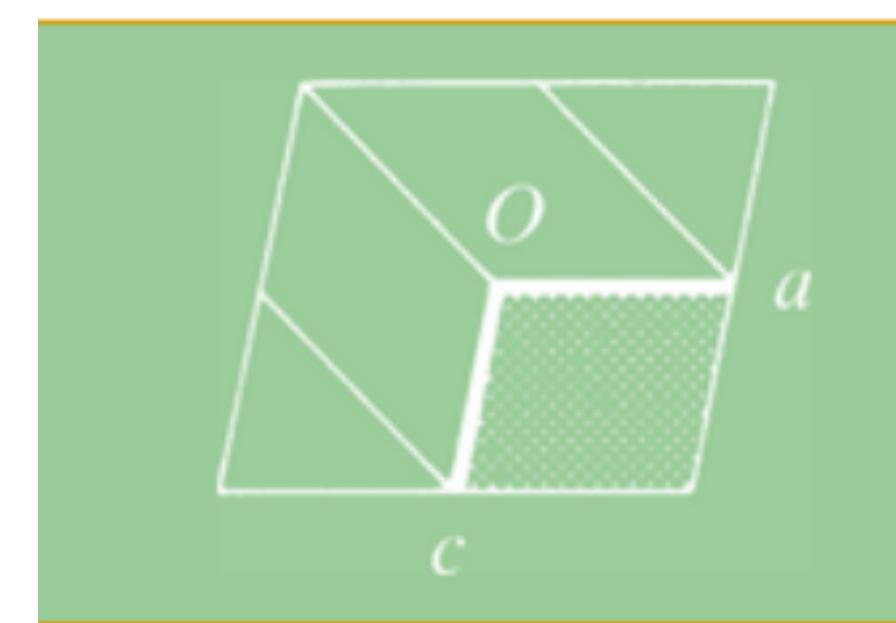
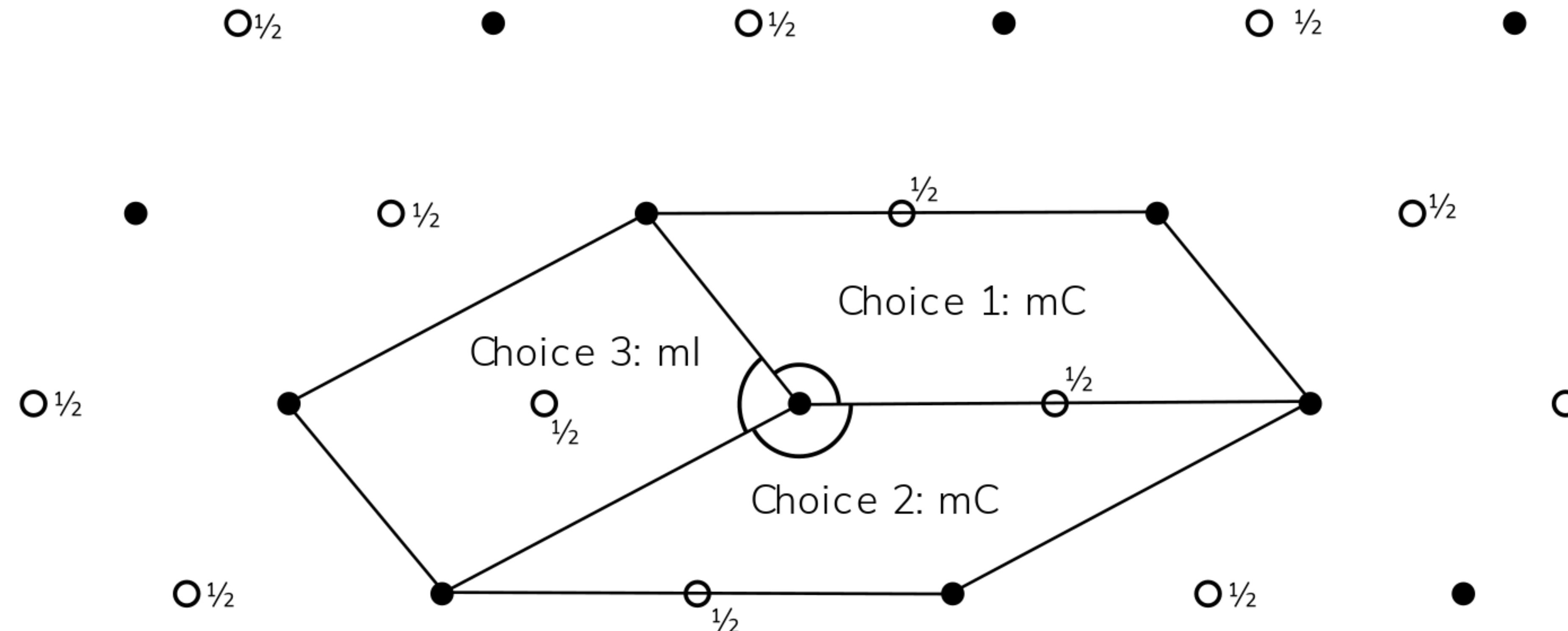
First run of ambig:

```
graph TD; A1[1/3] --> B1[2/3]; A1 --> C1[1/6]; A1 --> D1[1/6]; B1 --> E1[5/6]; B1 --> F1[1/12]; B1 --> G1[1/12];
```

Second run of ambig:

```
graph TD; B1[2/3] --> E1[5/6]; B1 --> F1[1/12]; B1 --> G1[1/12]; C1[1/6] --> F1[1/12]; C1 --> G1[1/12]; D1[1/6] --> F1[1/12]; D1 --> G1[1/12];
```

Cell choices for monoclinic C-centered structures



CrystFEL unit cell tool

Find possible indexing ambiguities:

```
$ cell_tool -p mystructure.cell --find-ambi
```

Calculate radii of all powder rings (useful for calibration):

```
$ cell_tool -p mystructure.cell --rings
```

Calculate all unit cell choices for monoclinic C cell:

```
$ cell_tool -p mystructure.cell --cell-choices
```

Transform a unit cell:

```
$ cell_tool -p mystructure.cell --transform=b,a-2b,c
```

Calculate primitive unit cell from centered cell:

```
$ cell_tool -p mystructure.cell --uncenter
```

Find transformation relating two cells:

```
$ cell_tool -p mystructure.cell --compare other.cell
```

Acknowledgements

Richard Kirian (Arizona State U.)

Kenneth Beyerlein

Andrew Aquila

Andrew Martin

Lorenzo Galli

Chun Hong Yoon

Karol Nass

Nadia Zatsepin (Arizona State U.)

Anton Barty

Thomas Grant (Hauptman-Woodward Inst., Buffalo)

Cornelius Gati

Aleksandra Tolstikova

Wolfgang Brehm (U. Konstanz, now CFEL)

Valerio Mariani

Parker de Waal (Van Andel Inst.)

Takanori Nakane (U. Tokyo)

Keitaro Yamashita (SPring-8)

Oleksandr Yefanov

Helen Ginn (U. Oxford)

Steve Aplin

Nicolas Riebesel (Technical U. Hamburg-Harburg)

Yaroslav Gevorkov

... and many more for feedback, ideas and guidance!

Funding: Helmholtz Association (Program-Oriented Funding); PIER Helmholtz Graduate School; European Union's 2020 Research and Innovation Programme (Marie Skłodowska-Curie grant agreement 637295 "X-Probe"); X-ray Free Electron Laser Priority Strategy Programme (MEXT, Japan); SACLAC HPC and Mini-K; BMBF German-Russian Cooperation "SyncFELMed" (grant 05K14CHA); BioStruct-X (Seventh Framework Programme of the European Commission)