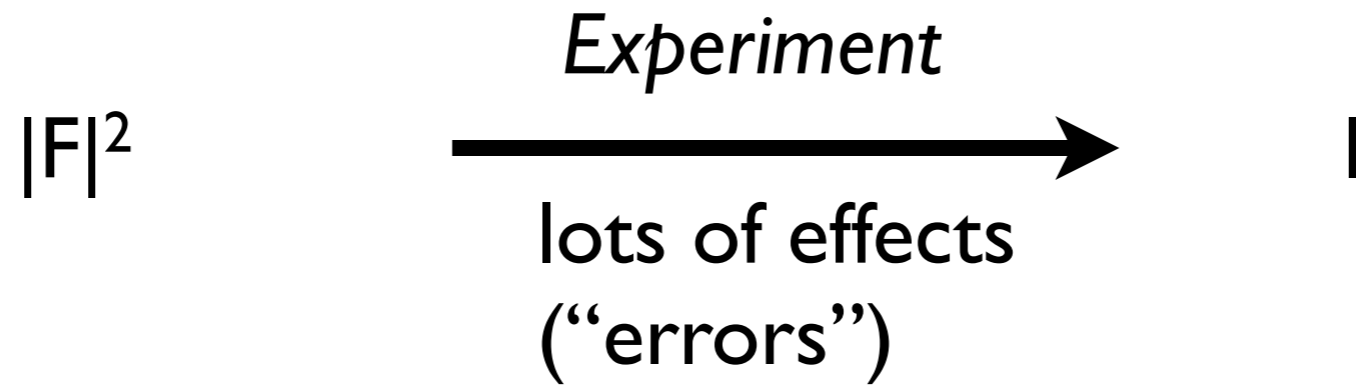


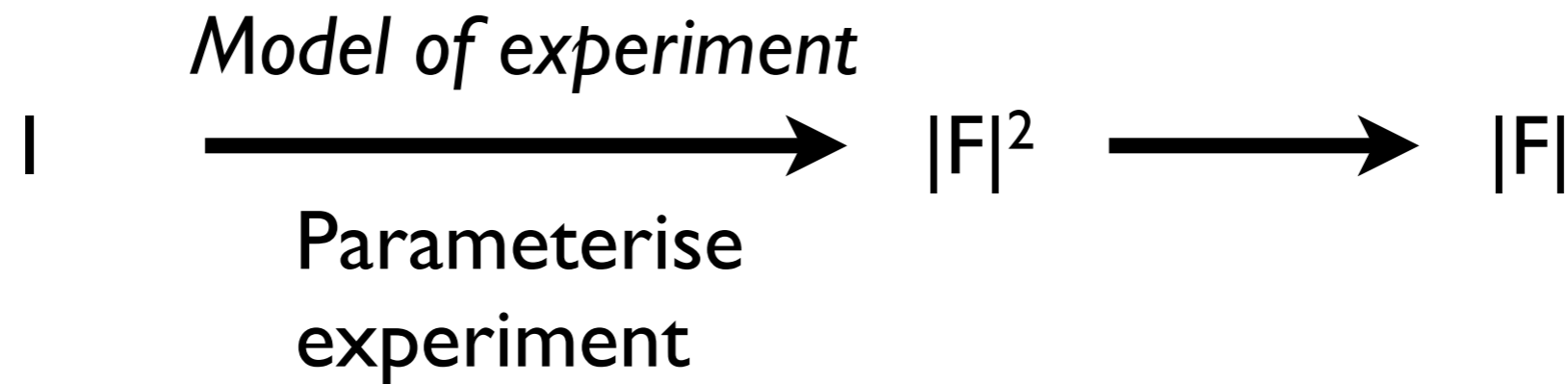
Data Reduction

Space Group Determination, Scaling and Intensity Statistics

Scaling and Merging



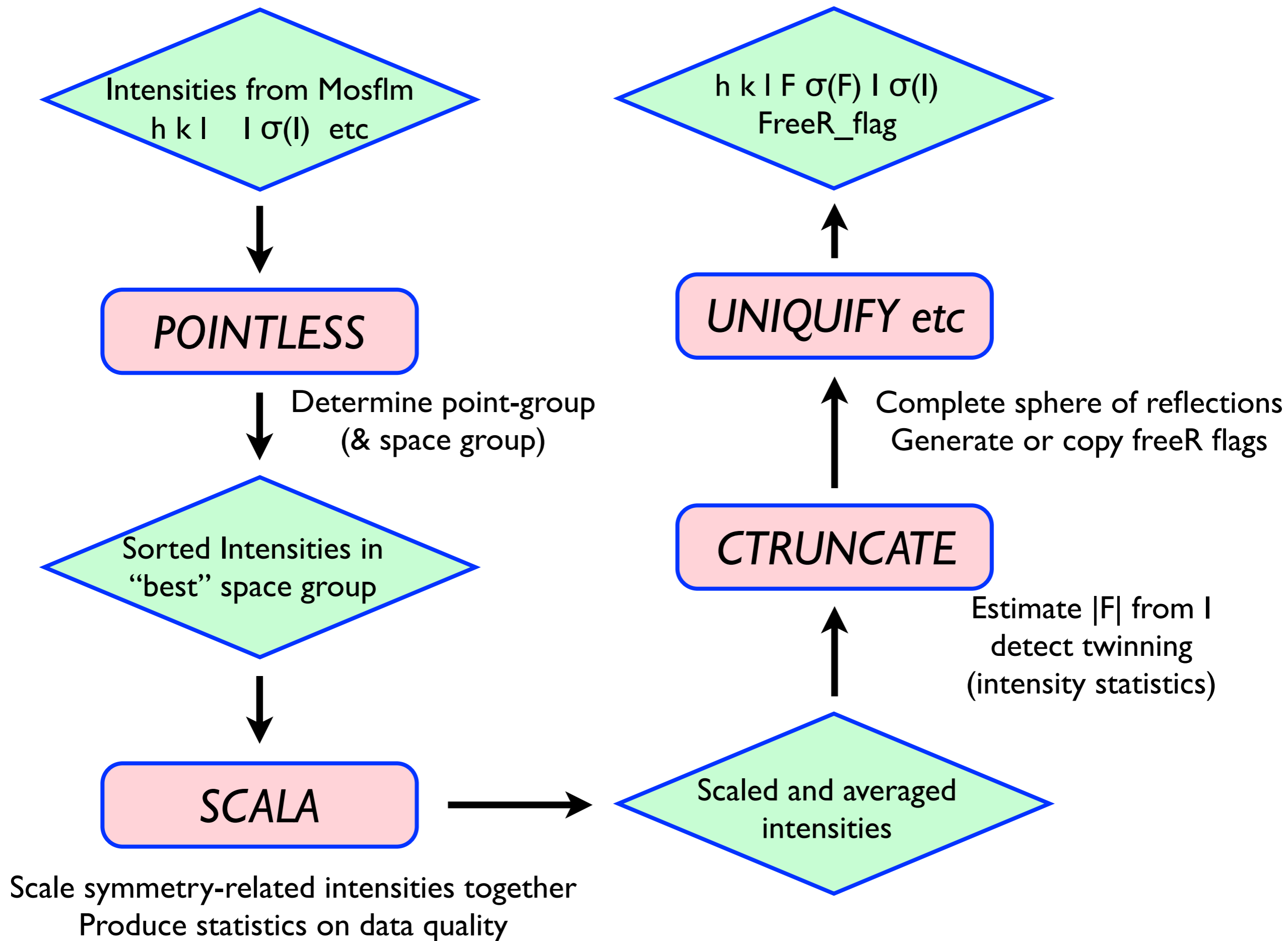
Our job is to invert the experiment: we want to *infer* $|F|$ from our measurements of intensity I

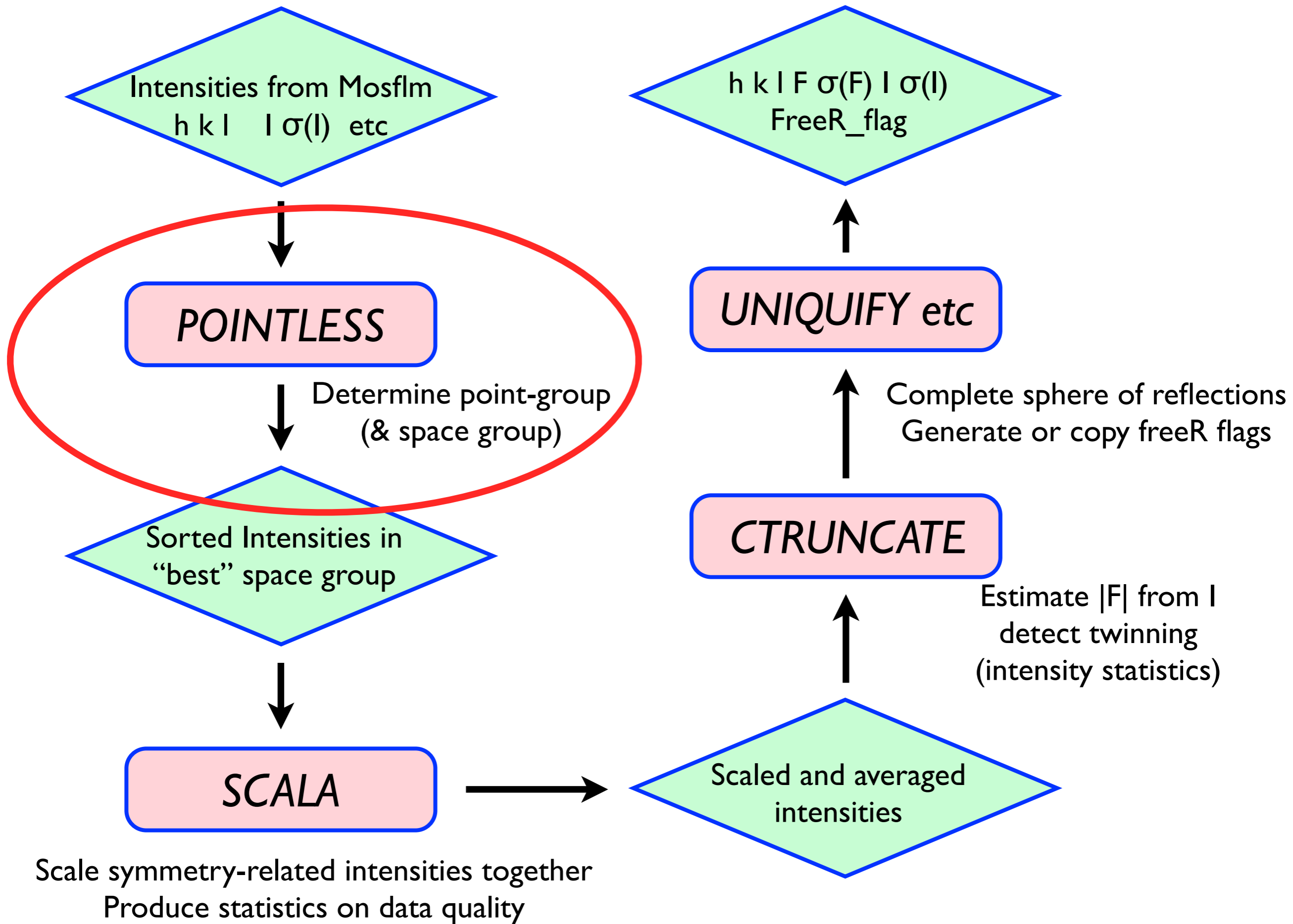


Data reduction can be done in an automated pipeline such as XIA2 (along with integration, *ie* go from images to a list of hkl F ready for structure determination)

XIA2 is used at Diamond (see Graeme Winter for more information)

This works pretty well, but in difficult cases you may need finer control over the process





Determination of Space group

The space group symmetry is only a **hypothesis** until the structure is solved, since it is hard to distinguish between true crystallographic and approximate (non-crystallographic) symmetry.

By examining the symmetry of the diffraction pattern we can get a good idea of the likely space group

It is also useful to find the likely symmetry as early as possible, since this affects the data collection strategy

Lattice symmetry imposes constraints on the cell dimensions (eg $\alpha=\beta=\gamma=90^\circ$ for an orthorhombic lattice), but the converse is not true: cell dimensions can have special relationships accidentally. Indexing in eg Mosflm only considers lattice *geometry* not symmetry (cubic, hexagonal/trigonal, tetragonal, orthorhombic, monoclinic, or triclinic, + lattice centring P, C, I, R, or F)

The Laue group (Patterson group) is the symmetry of the diffraction pattern, so can be determined from the observed intensities. It corresponds to the space group without any translations, and with an added centre of symmetry from Friedel's law.

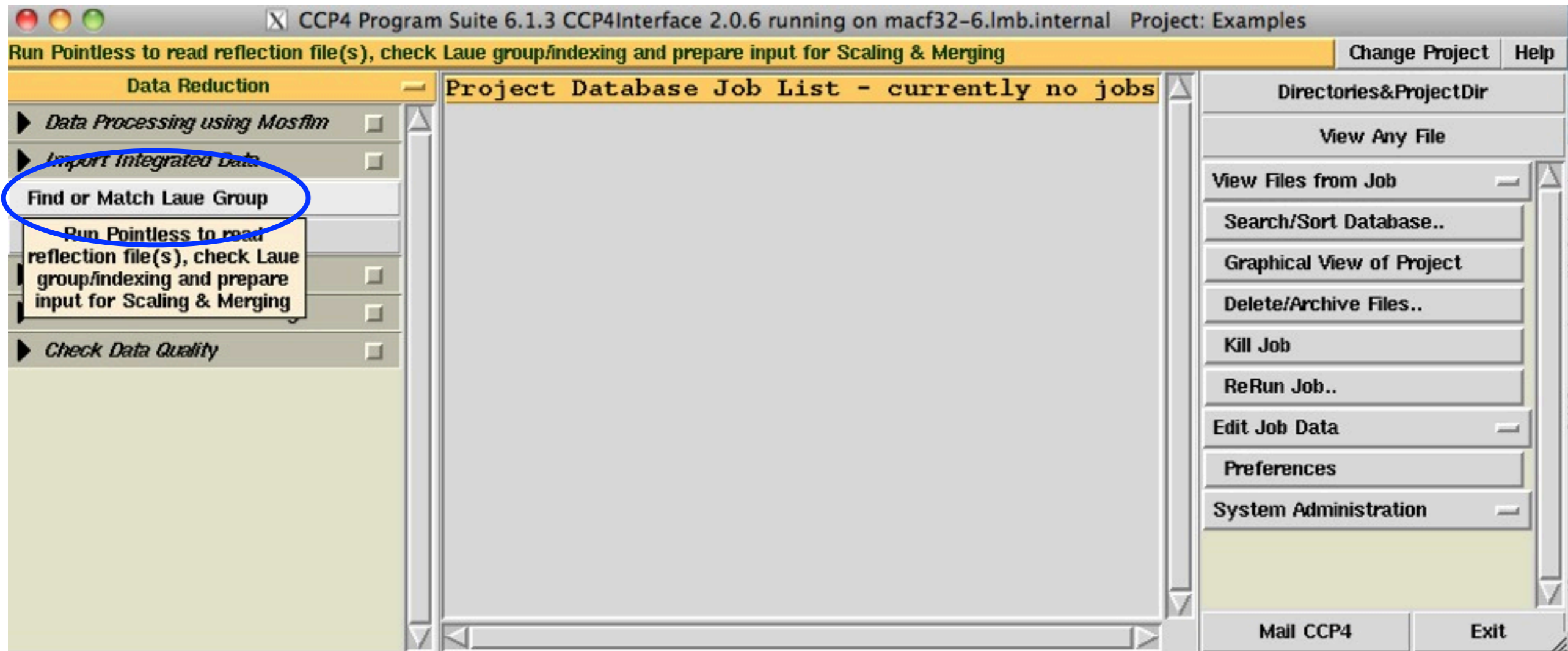
The space group is the point group + lattice centring + translations (eg screw dyad rather than pure dyad). Only visible in diffraction pattern as systematic absences along axes – these are not very reliable indicators as there are few axial reflections and there may be accidental absences.

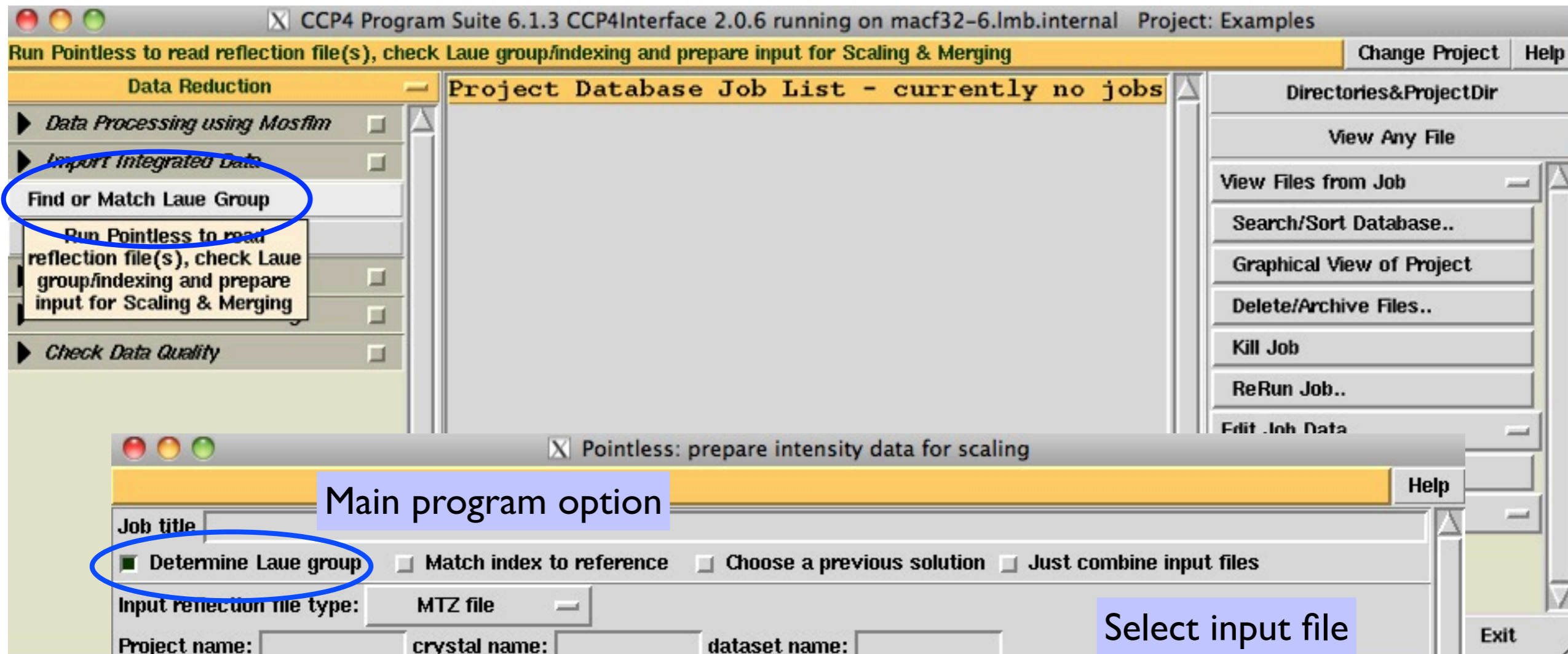
Protocol for space group determination (program *POINTLESS*)

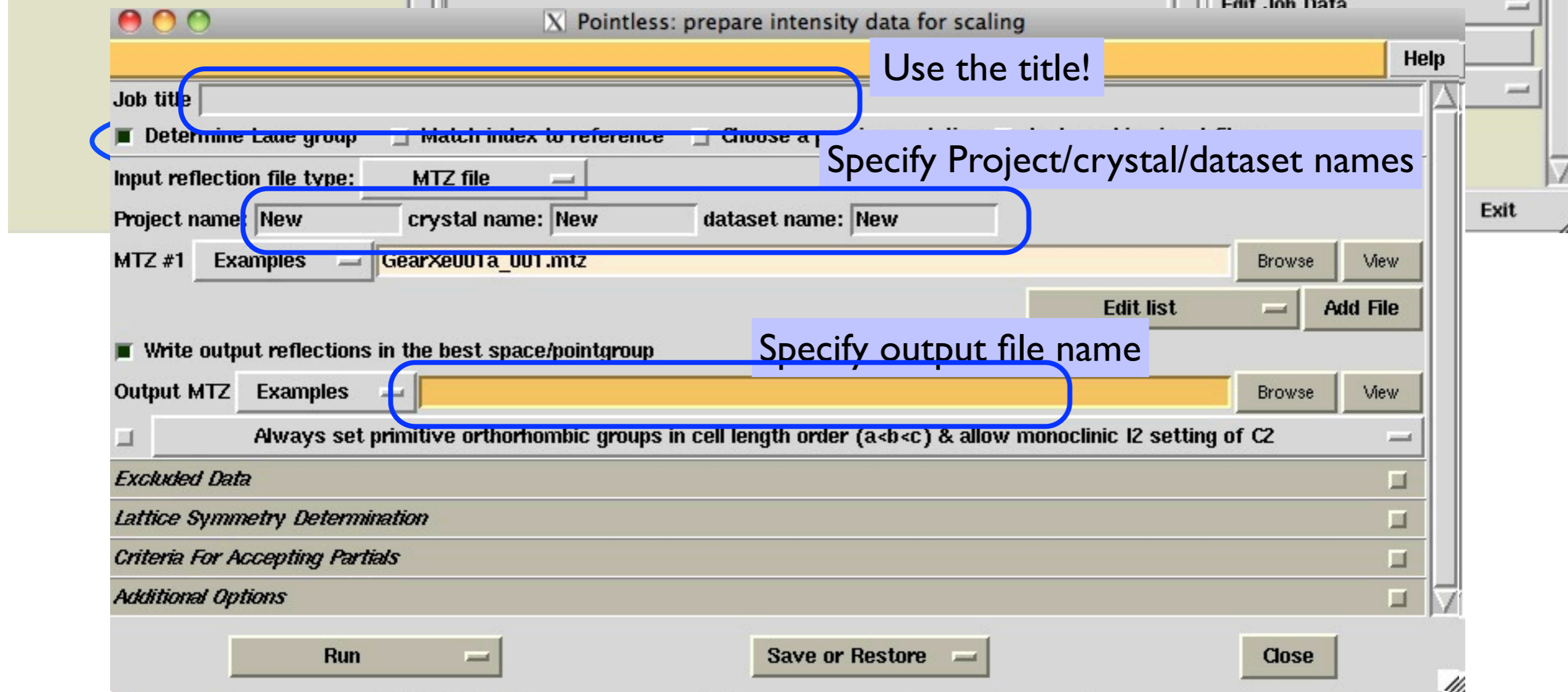
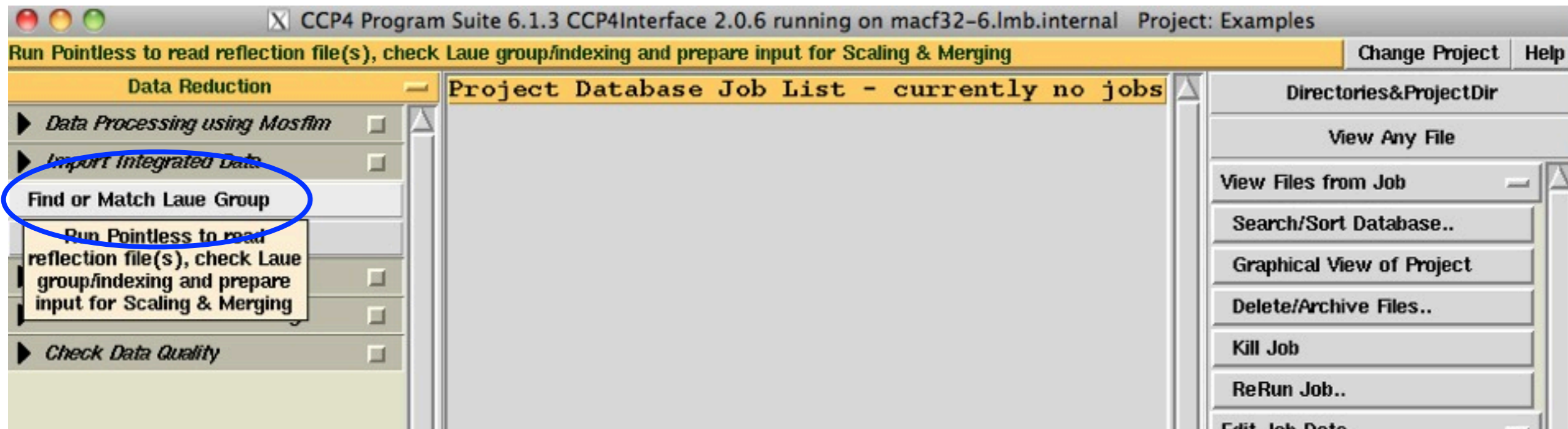
1. From the unit cell dimensions, find the highest compatible lattice symmetry (within a tolerance)
2. Score each symmetry element (rotation) belonging to lattice symmetry using all pairs of observations related by that element
3. Score combinations of symmetry elements for all possible sub-groups (Laue groups) of lattice symmetry group.
4. Score possible space groups from axial systematic absences

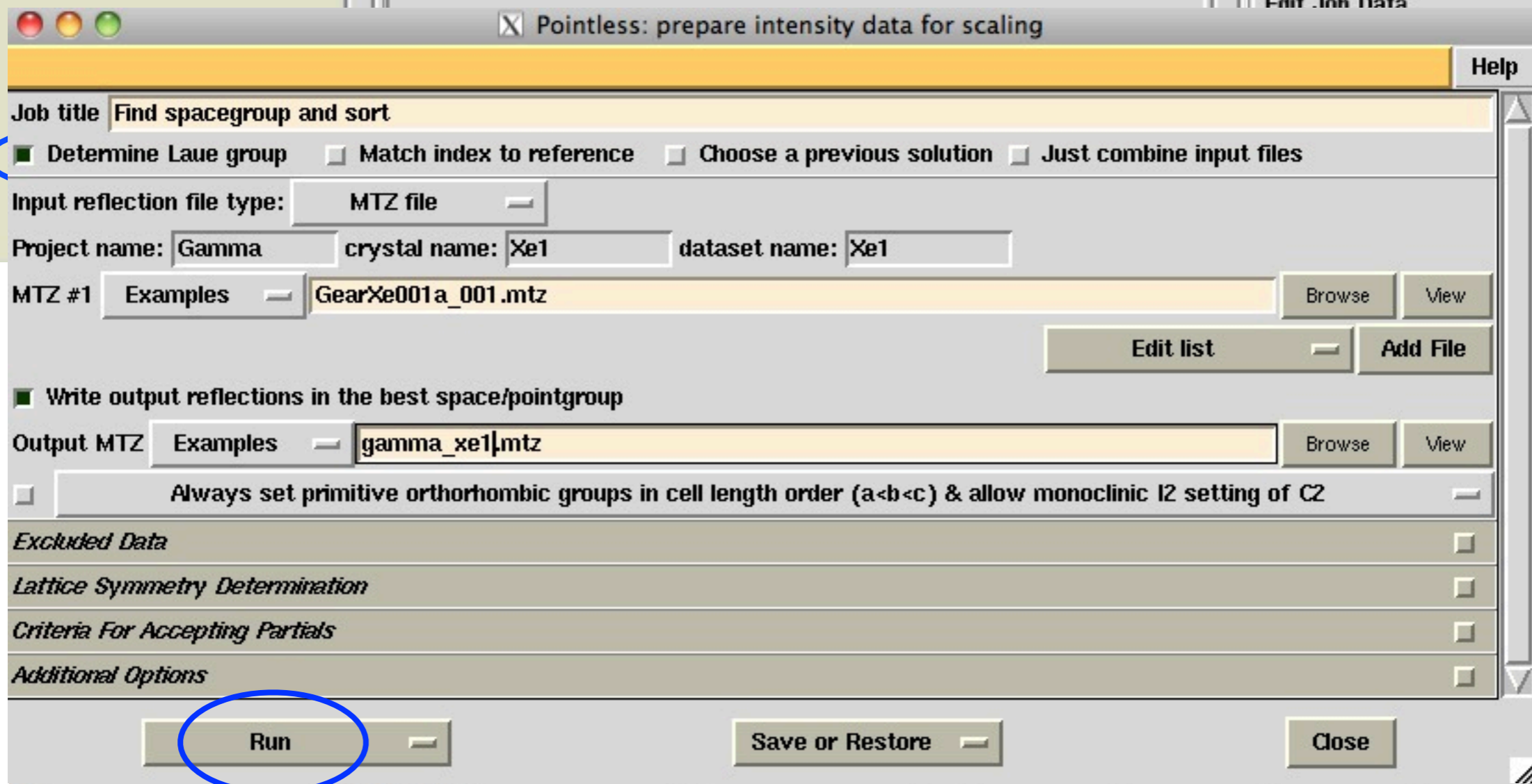
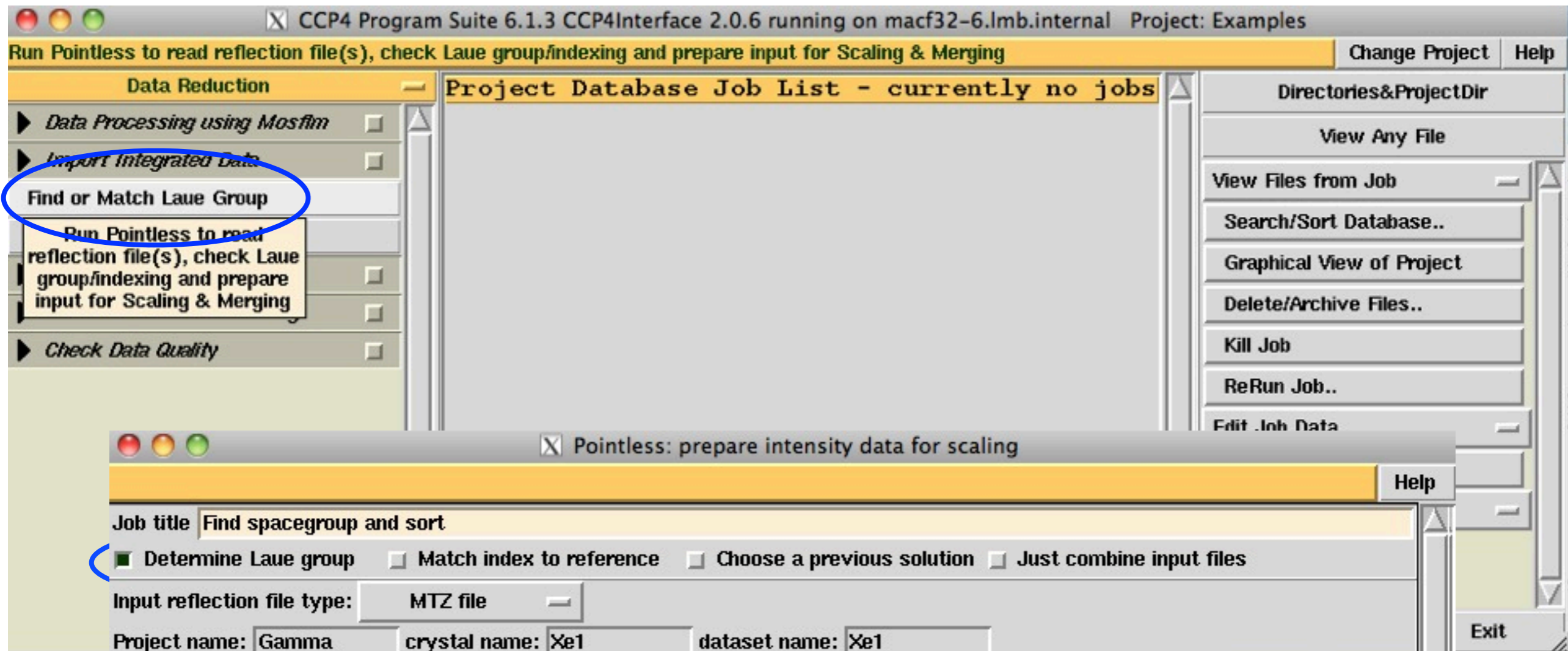
Scoring functions for rotational symmetry based on **correlation coefficient**, since this is relatively independent of the unknown scales.

R_{meas} values are also calculated

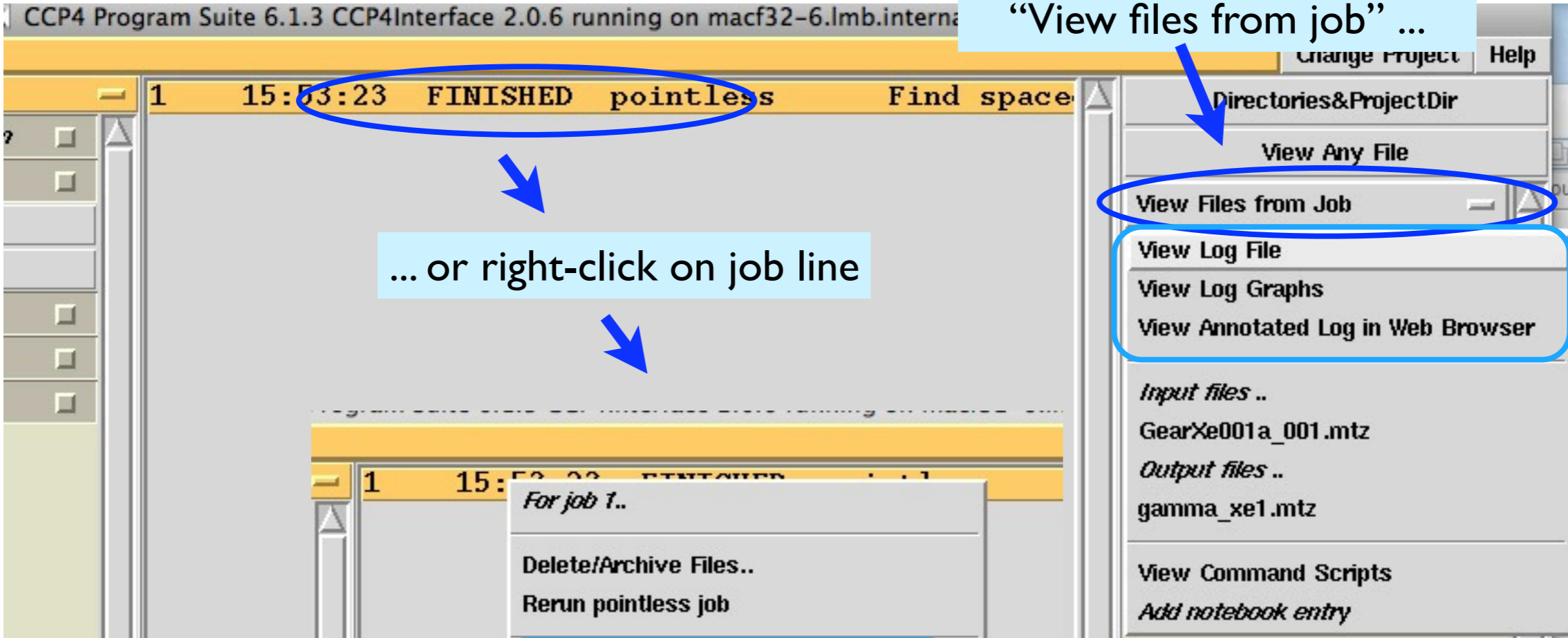




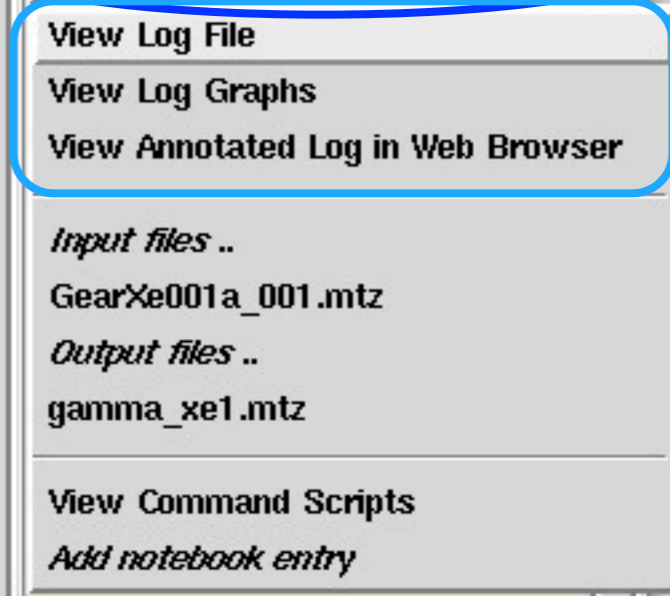
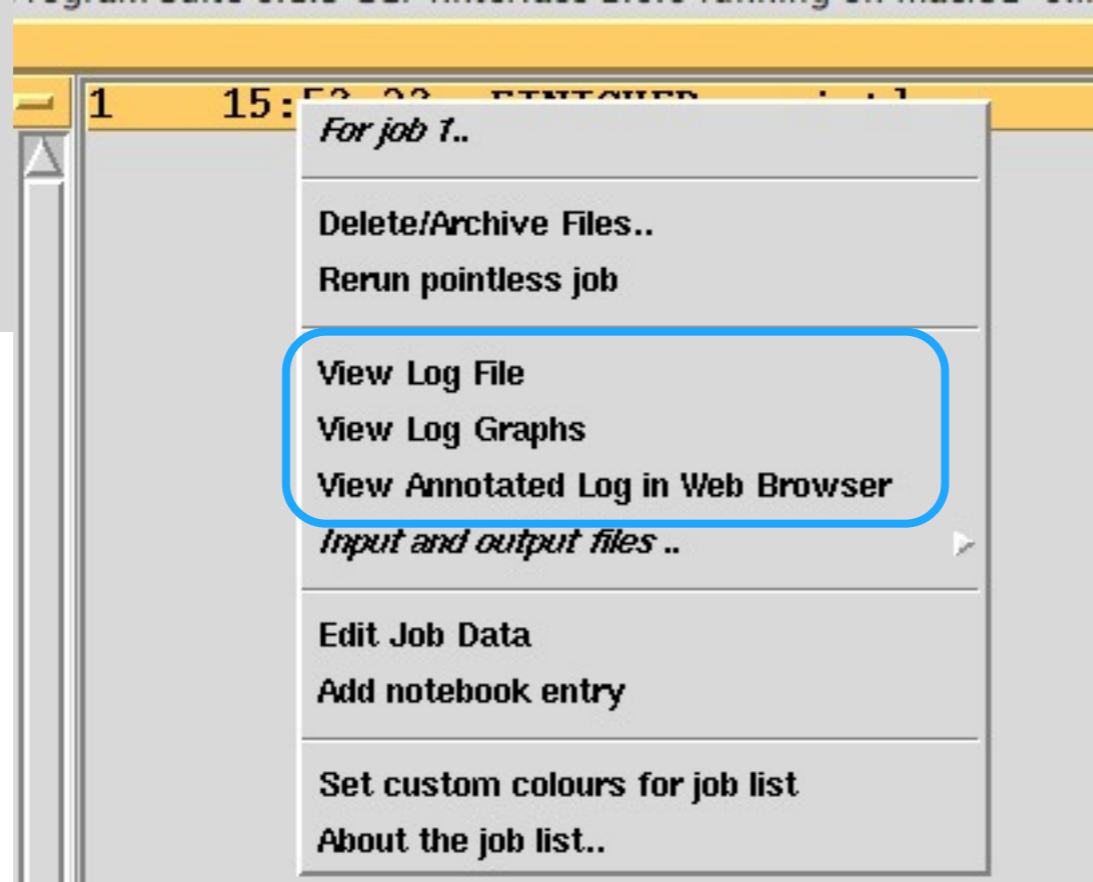




Examine output either from
"View files from job" ...



... or right-click on job line



Examine output either from
"View files from job" ...

1 15:53:23 FINISHED pointless Find space

... or right-click on job line

- Directories&ProjectDir
- View Any File
- View Files from Job**
- View Log File
- View Log Graphs
- View Annotated Log in Web Browser**
- Input files ..*
- GearXe001a_001.mtz
- Output files ..*
- gamma_xe1.mtz
- View Command Scripts
- Add notebook entry

- For job 1..
- Delete/Archive Files..
- Rerun pointless job
- View Log File
- View Log Graphs
- View Annotated Log in Web Browser**
- Input and output files ..*
- Edit Job Data
- Add notebook entry
- Set custom colours for job list
- About the job list..

View Annotated Log in Web Browser

1_pointless.log

Please consider citing the following papers:

- Pointless
 - P.R.Evans, 'Scaling and assessment of data quality' Acta Cryst. D62, 72-82 (2006).

Pointless Version 1.4.6 Run at 15:53:20 on 17/12/2009

Result:

Best Solution space group P 21 21 21

Reindex operator:	[h,k,l]
Laue group probability:	0.985
Systematic absence probability:	0.851
Total probability:	0.838
Space group confidence:	0.784
Laue group confidence:	0.982

Summary table:
probabilities and
confidence levels

Unit cell: 34.16 54.8 68 90 90 90

17.00 to 1.78 - Resolution range used for Laue group search

17.00 to 1.78 - Resolution range in file, used for systematic absence check

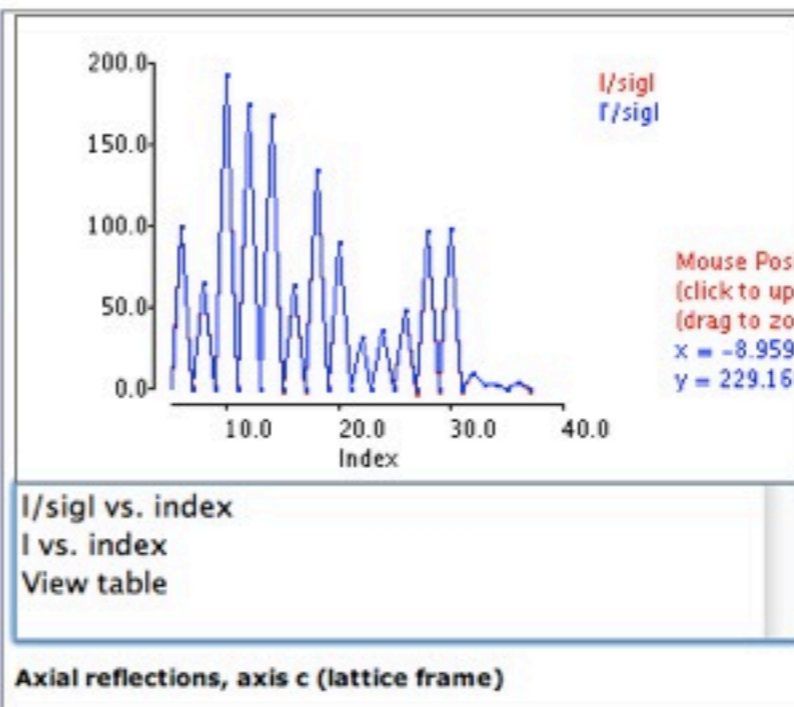
Number of batches in file: 100

The following tables were found in the logfile:

[Axial reflections, axis a \(lattice frame\)](#)

[Axial reflections, axis b \(lattice frame\)](#)

[Axial reflections, axis c \(lattice frame\)](#)



Graphs of axial
reflections for
systematic absences

[\[Show logfile summary\]](#) [\[Show full logfile\]](#)

[\[Documentation\]](#)

A straightforward orthorhombic case

Scoring the symmetry operators separately sometimes allows detection of pseudo-symmetry, eg if some rotation operators are much weaker than others

Analysing rotational symmetry in lattice group P m m m

Score	Probability	1 s	Correlation coefficient	ment	R-factor		
Netmt	Lklhd	Z-cc	CC	N	Rmeas	Symmetry & operator (in Lattice Cell)	
1	0.948	9.54	0.95	12122	0.097	identity	
2	0.942	9.44	0.94	18346	0.121	***	2-fold l { 0 0 1} {-h, -k, +l}
3	0.949	9.58	0.96	30259	0.097	***	2-fold h { 1 0 0} {+h, -k, -l}
4	0.912	9.15	0.92	17427	0.120	***	2-fold k { 0 1 0} {-h, +k, -l}

Separate scores for each symmetry operator in maximum possible lattice symmetry

	Laue Group	ReindexOperator	Lklhd	NetZc	Zc+	Zc-	CC	CC-	Rmeas	R-	Delta	
= 1	P m m m	***	0.985	9.35	9.35	0.00	0.94	0.00	0.11	0.00	0.0	[h, k, l]
2	P 1 2/m 1		0.006	0.38	9.56	9.18	0.96	0.92	0.10	0.12	0.0	[-k, -h, -l]
3	P 1 2/m 1		0.005	-0.01	9.38	9.39	0.94	0.94	0.11	0.11	0.0	[-h, -l, -k]
4	P 1 2/m 1		0.003	-0.13	9.31	9.44	0.93	0.94	0.11	0.11	0.0	[h, k, l]
5	P -1		0.000	0.22	9.54	9.32	0.95	0.93	0.10	0.11	0.0	[h, k, l]

Combined scores for all possible Laue (point) groups down to P1

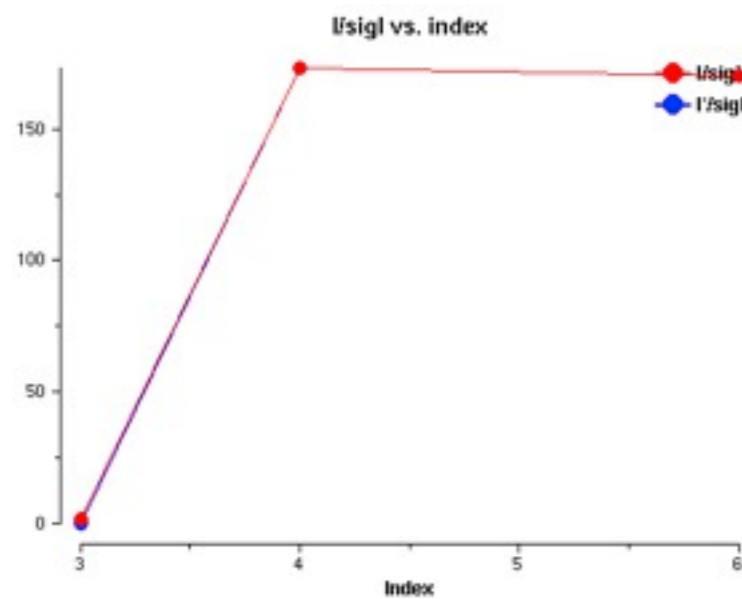
A clear indication that the Laue group is Pmmm (P222)

Possible axial systematic absences to determine space group

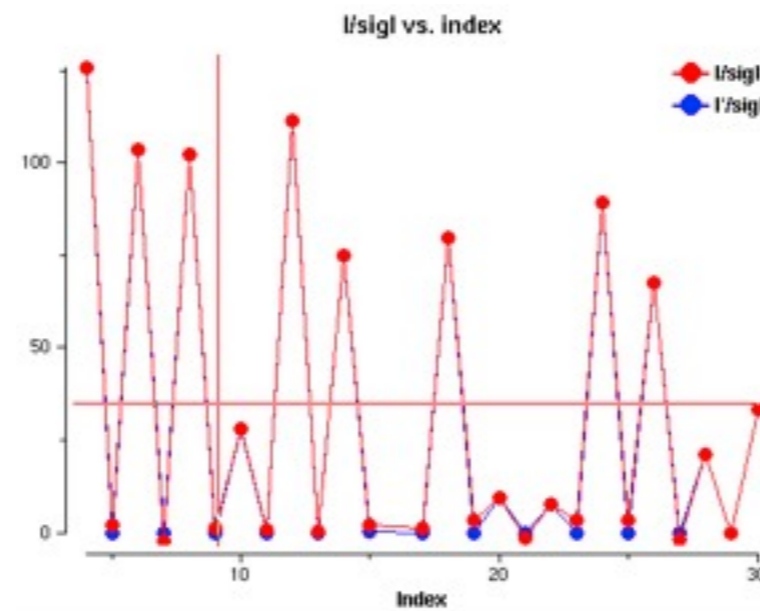
Zone	Number	PeakHeight	SD	Probability	ReflectionCondition
Zones for Laue group $P\ m\ m\ m$					
1 screw axis $2(1)$ [a]	3	1.000	0.296	** 0.889	$h00: h=2n$
2 screw axis $2(1)$ [b]	26	1.000	0.142	*** 0.971	$0k0: k=2n$
3 screw axis $2(1)$ [c]	46	0.997	0.097	*** 0.986	$00l: l=2n$

Fourier analysis of $I/\sigma(I)$

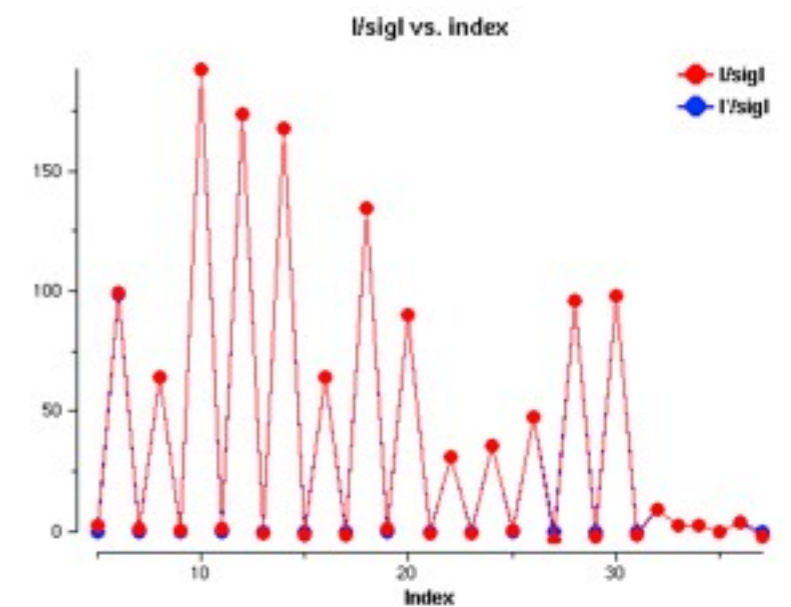
There are indications of 2_1 screw symmetry along all principle axes (though note there are only 3 observations on the a axis ($h00$ reflections))



Possible 2_1 axis along a



Clear 2_1 axis along b



Clear 2_1 axis along c

Possible spacegroups:

Indistinguishable space groups are grouped together on successive lines

'Reindex' is the operator to convert from the input hklin frame to the standard spacegroup frame.

'TotProb' is a total probability estimate (unnormalised)

'SysAbsProb' is an estimate of the probability of the space group based on the observed systematic absences.

'Conditions' are the reflection conditions (absences)

Spacegroup	TotProb	SysAbsProb	Reindex	Conditions
$\langle P\ 21\ 21\ 21 \rangle$ (19)	0.838	0.851		$h00: h=2n, 0k0: k=2n, 00l: l=2n$ (zones 1,2,3)
$\langle P\ 2\ 21\ 21 \rangle$ (18)	0.104	0.106		$0k0: k=2n, 00l: l=2n$ (zones 2,3)
$\langle P\ 21\ 2\ 21 \rangle$ (18)	0.025	0.026		$h00: h=2n, 00l: l=2n$ (zones 1,3)
$\langle P\ 21\ 21\ 2 \rangle$ (18)	0.012	0.012		$h00: h=2n, 0k0: k=2n$ (zones 1,2)

Best Solution space group P 21 21 21

Reindex operator: [h,k,l]
Laue group probability: 0.985
Systematic absence probability: 0.851
Total probability: 0.838
Space group confidence: 0.784
Laue group confidence 0.982

Unit cell: 34.16 54.8 68 90 90 90

17.00 to 1.78 - Resolution range used for Laue group search

17.00 to 1.78 - Resolution range in file, used for systematic absence check

Number of batches in file: 100

Note high confidence in Laue group, but lower confidence in space group

Pseudo-cubic example

Cell: 79.15 81.33 81.15 90.00 90.00 90.00 $a \approx b \approx c$

Analysing rotational symmetry in lattice group $P m \bar{3} m$

Scores for each symmetry element

Nelmt	Lklhd	Z-cc	CC	N	Rmeas	Symmetry & operator (in Lattice Cell)		
1	0.955	9.70	0.97	13557	0.073	identity		
2	0.062	2.66	0.27	12829	0.488	2-fold	(1 0 1)	{+l,-k,+h}
3	0.065	2.85	0.29	10503	0.474	2-fold	(1 0 -1)	{-l,-k,-h}
4	0.056	0.06	0.01	16391	0.736	2-fold	(0 1 -1)	{-h,-l,-k}
5	0.057	0.05	0.00	17291	0.738	2-fold	(0 1 1)	{-h,+l,+k}
6	0.049	0.55	0.06	13758	0.692	2-fold	(1 -1 0)	{-k,-h,-l}
7	0.950	9.59	0.96	12584	0.100	*** 2-fold k	(0 1 0)	{-h,+k,-l}
8	0.049	0.57	0.06	11912	0.695	2-fold	(1 1 0)	{+k,+h,-l}
9	0.948	9.57	0.96	16928	0.136	*** 2-fold h	(1 0 0)	{+h,-k,-l}
10	0.944	9.50	0.95	12884	0.161	*** 2-fold l	(0 0 1)	{-h,-k,+l}
11	0.054	0.15	0.01	23843	0.812	3-fold	(1 1 1)	{+l,+h,+k} {+k,+l,+h}
12	0.055	0.11	0.01	24859	0.825	3-fold	(1 -1 -1)	{-l,-h,+k} {-k,+l,-h}
13	0.055	0.14	0.01	22467	0.788	3-fold	(1 -1 1)	{+l,-h,-k} {-k,-l,+h}
14	0.055	0.12	0.01	27122	0.817	3-fold	(1 1 -1)	{-l,+h,-k} {+k,-l,-h}
15	0.061	-0.10	-0.01	25905	0.726	4-fold h	(1 0 0)	{+h,-l,+k} {+h,+l,-k}
16	0.060	2.53	0.25	23689	0.449	4-fold k	(0 1 0)	{+l,+k,-h} {-l,+k,+h}
17	0.049	0.56	0.06	25549	0.653	4-fold l	(0 0 1)	{-k,+h,+l} {+k,-h,+l}

Only orthorhombic symmetry operators are present

Pseudo-cubic example

Cell: 79.15 81.33 81.15 90.00 90.00 90.00 $a \approx b \approx c$

	Laue	Group		Lklhd	NetZc	Zc+	Zc-	CC	CC-	Rmeas	R-	Delta	ReindexOperator
= 1	P	m m m	***	0.989	8.93	9.59	0.66	0.96	0.07	0.12	0.69	0.0	[-h,-l,-k]
2	P	1 2/m 1		0.003	7.85	9.65	1.80	0.97	0.18	0.09	0.60	0.0	[-h,-l,-k]
3	P	1 2/m 1		0.003	7.95	9.63	1.68	0.96	0.17	0.10	0.61	0.0	[l,h,k]
4	P	1 2/m 1		0.003	7.80	9.61	1.81	0.96	0.18	0.11	0.60	0.0	[h,k,l]
5	P	4/m m m		0.000	6.69	6.90	0.21	0.69	0.02	0.24	0.75	1.5	[-k,-h,-l]
6	P	4/m m m		0.000	4.55	5.41	0.85	0.54	0.09	0.34	0.68	0.1	[-l,-k,-h]
7		P 4/m		0.000	5.45	7.20	1.75	0.72	0.18	0.20	0.62	1.5	[-k,-h,-l]
8		P 4/m		0.000	4.72	6.53	1.81	0.65	0.18	0.25	0.60	0.1	[-l,-k,-h]
9		P -1		0.000	7.48	9.70	2.22	0.97	0.22	0.07	0.57	0.0	[-h,-l,-k]
10		P 4/m		0.000	4.03	5.96	1.92	0.60	0.19	0.29	0.59	1.4	[-h,-l,-k]
11	P	4/m m m		0.000	4.93	5.63	0.69	0.56	0.07	0.32	0.69	1.4	[-h,-l,-k]
12	C	m m m		0.000	4.97	6.67	1.70	0.67	0.17	0.24	0.62	1.5	[h-k,-h-k,-l]
13	C	1 2/m 1		0.000	4.80	6.99	2.19	0.70	0.22	0.21	0.57	1.5	[-h-k,-h+k,-l]
14	C	1 2/m 1		0.000	4.51	6.71	2.20	0.67	0.22	0.23	0.58	1.5	[h-k,-h-k,-l]
15	C	m m m		0.000	3.08	5.01	1.93	0.50	0.19	0.36	0.59	0.1	[-k-l,-k+l,-h]
16		P m -3		0.000	3.35	4.32	0.97	0.43	0.10	0.44	0.63	1.5	[h,k,l]
17	C	1 2/m 1		0.000	2.58	4.95	2.36	0.49	0.24	0.35	0.56	0.1	[k-l,-k-l,-h]
18	C	1 2/m 1		0.000	2.65	5.01	2.36	0.50	0.24	0.34	0.56	0.1	[-k-l,-k+l,-h]
19		H -3		0.000	2.17	4.56	2.39	0.46	0.24	0.40	0.55	1.5	[-k+l,-h-l,h-k-l]
20		H -3		0.000	2.09	4.48	2.39	0.45	0.24	0.40	0.55	1.5	[h-l,-h-k,-h+k-l]
21		H -3		0.000	2.15	4.54	2.39	0.45	0.24	0.39	0.55	1.5	[-h+k,-k-l,-h-k+l]
22		H -3		0.000	2.20	4.59	2.38	0.46	0.24	0.39	0.55	1.5	[k-l,h-k,-h-k-l]
23	C	1 2/m 1		0.000	3.10	5.42	2.32	0.54	0.23	0.31	0.56	1.4	[-h-l,h-l,-k]
24	C	1 2/m 1		0.000	3.36	5.67	2.31	0.57	0.23	0.30	0.56	1.4	[-h+l,-h-l,-k]
25	C	m m m		0.000	3.32	5.29	1.97	0.53	0.20	0.34	0.59	1.4	[-h-l,h-l,-k]
26		H -3 m		0.000	-0.01	2.66	2.67	0.27	0.27	0.52	0.54	1.5	[-h+k,-k-l,-h-k+l]
27		H -3 m		0.000	-0.03	2.65	2.68	0.26	0.27	0.52	0.54	1.5	[k-l,h-k,-h-k-l]
28		H -3 m		0.000	-0.13	2.58	2.71	0.26	0.27	0.53	0.53	1.5	[h-l,-h-k,-h+k-l]
29		H -3 m		0.000	-0.02	2.66	2.68	0.27	0.27	0.52	0.53	1.5	[-k+l,-h-l,h-k-l]
30	P	m -3 m		0.000	2.67	2.67	0.00	0.27	0.00	0.53	0.00	1.5	[h,k,l]

... symmetry is actually orthorhombic (P 2₁ 2₁ 2₁)

Combining multiple files (and multiple MAD datasets)

3 files
assigned to
same dataset

Job title: pk ip rm Se34

Determine Laue group Match index to reference Choose a previous solution Just combine input files

Input reflection file type: MTZ file

Project name: Brap crystal name: Se34 dataset name: pk

MTZ #1 Full path.. /Amb/home/pre/Projects/Brap/Se34/pk_1_001.mtz

MTZ #2 Full path.. /Amb/home/pre/Projects/Brap/Se34/pk_2_001.mtz

Assign to the same dataset as the previous file

MTZ #3 Full path.. /Amb/home/pre/Projects/Brap/Se34/pk_180_1_001.mtz

Assign to the same dataset as the previous file

MTZ #4 Full path.. /Amb/home/pre/Projects/Brap/Se34/ip_1_001.mtz

Assign to the same dataset as the previous file

Project name: Brap crystal name: Se34 dataset name: ip

MTZ #5 Full path.. /Amb/home/pre/Projects/Brap/Se34/rm_1_001.mtz

Assign to the same dataset as the previous file

Project name: Brap crystal name: Se34 dataset name: Rm

Write output reflections in the best space/pointgroup

Output MTZ Brap se34_pk_ip_rm.mtz

Test Laue group of 1st file before reading rest Assume all files have same indexing (faster)

Always set primitive orthorhombic groups in cell length order (a<b<c) & allow monoclinic I2 setting of C2

Excluded Data

Lattice Symmetry Determination

Criteria For Accepting Partial

Additional Options

Dataset 1, pk, 3 files

Dataset 2, ip, 1 file

Dataset 3, rm, 1 file

Combining multiple files (and multiple MAD datasets)

```
Alternative index test relative to first file
Alternative reindexing      CC      R(E^2)    Number Cell_deviation
[h,k,l]                    0.965    0.086    23592    0.00
[-k,h,l]                   0.789    0.205    22755    0.30
[l,k,-h]                   0.102    0.438    21060    0.76
[k,l,h]                    0.055    0.459    22714    0.66
[-h,l,k]                   0.048    0.461    23282    0.46
[l,h,k]                    0.043    0.457    21194    0.66
```

```
Alternative index test relative to files so far
Alternative reindexing      CC      R(E^2)    Number Cell_deviation
[h,k,l]                    0.933    0.124    40670    0.14
[-k,h,l]                   0.610    0.283    40494    0.43
[l,k,-h]                   0.061    0.463    40338    0.84
[-h,l,k]                   0.045    0.470    40635    0.43
[l,h,k]                    0.027    0.477    40352    0.68
[k,l,h]                    0.020    0.479    40461    0.77
```

```
Alternative index test relative to files so far
Alternative reindexing      CC      R(E^2)    Number Cell_deviation
[h,k,l]                    0.933    0.124    40670    0.14
[-k,h,l]                   0.610    0.283    40494    0.43
[l,k,-h]                   0.061    0.463    40338    0.84
[-h,l,k]                   0.045    0.470    40635    0.43
[l,h,k]                    0.027    0.477    40352    0.68
[k,l,h]                    0.020    0.479    40461    0.77
```

```
Alternative index test relative to files so far
Alternative reindexing      CC      R(E^2)    Number Cell_deviation
[h,k,l]                    0.960    0.095    22712    0.07
[-k,h,l]                   0.706    0.241    22712    0.36
[l,k,-h]                   0.084    0.455    22690    0.80
[k,l,h]                    0.050    0.468    22698    0.67
[-h,l,k]                   0.046    0.465    22701    0.44
[l,h,k]                    0.025    0.472    22693    0.72
```

Alternative indexing relative to first file(s):

	Reindex operator	CC	File name
2	[h,k,l]	0.965	pk_2_001.mtz
3	[h,k,l]	0.933	pk_180_1_001.mtz
4	[h,k,l]	0.960	ip_1_001.mtz
5	[h,k,l]	0.958	rm_1_001.mtz

Because of an indexing ambiguity (pseudo-cubic orthorhombic), we must check for consistent indexing between files

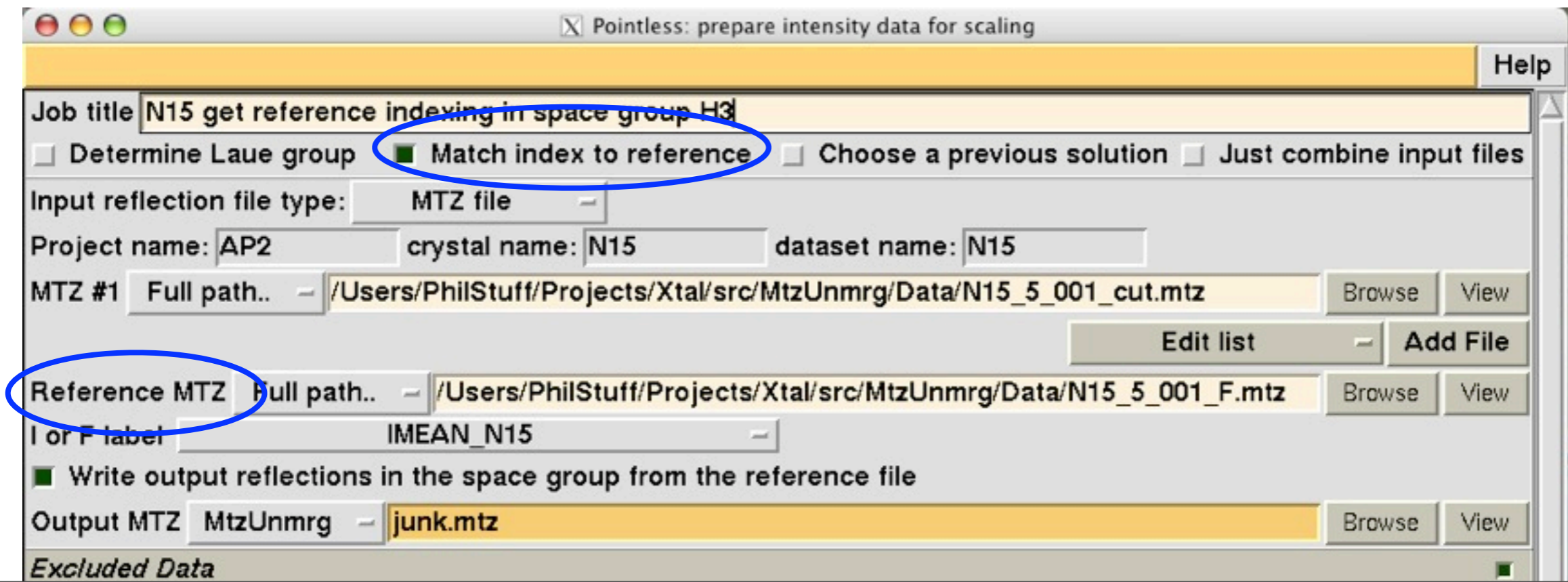
Alternative indexing

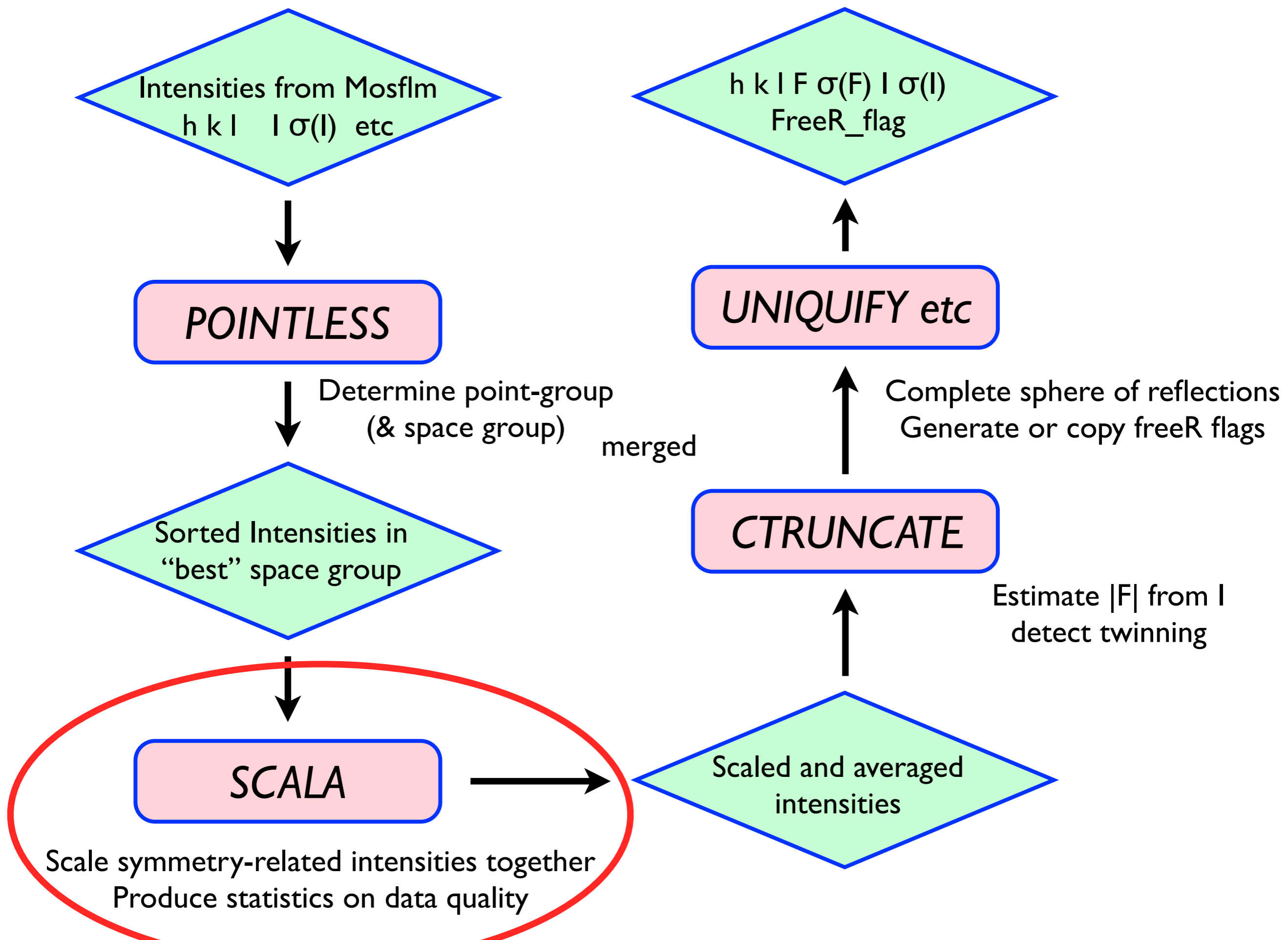
If the true point group is lower symmetry than the lattice group, alternative valid but non-equivalent indexing schemes are possible, related by symmetry operators present in lattice group but not in point group (*note that these are also the cases where merohedral twinning is possible*)

eg if in space group $P3$ (or $P3_1$) there are 4 different schemes
(h,k,l) or (-h,-k,l) or (k,h,-l) or (-k,-h,-l)

For the first crystal, you can choose any scheme

For subsequent crystals, the autoindexing will randomly choose one setting, and we need to make it consistent: *POINTLESS* will do this for you by comparing the unmerged test data to a reference dataset (merged or unmerged)





Scaling

Scaling tries to make symmetry-related and duplicate measurements of a reflection equal, by modelling the diffraction experiment, principally as a function of the incident and diffracted beam directions in the crystal. This makes the data **internally consistent**.

After scaling, the remaining differences between observations can be analysed to give an *indication* of data quality, though not necessarily of its absolute correctness.

Measures of internal consistency:

R-factors & correlation coefficients:

$$R_{\text{merge}} (R_{\text{sym}}) = \sum | I_{hl} - \langle I_h \rangle | / \sum | \langle I_h \rangle |$$

traditional overall measures of quality, but increases with multiplicity although the data improves

$$R_{\text{meas}} = R_{\text{r.i.m.}} = \sum \sqrt{(n/n-1)} | I_{hl} - \langle I_h \rangle | / \sum | \langle I_h \rangle |$$

multiplicity-weighted, better (but larger)

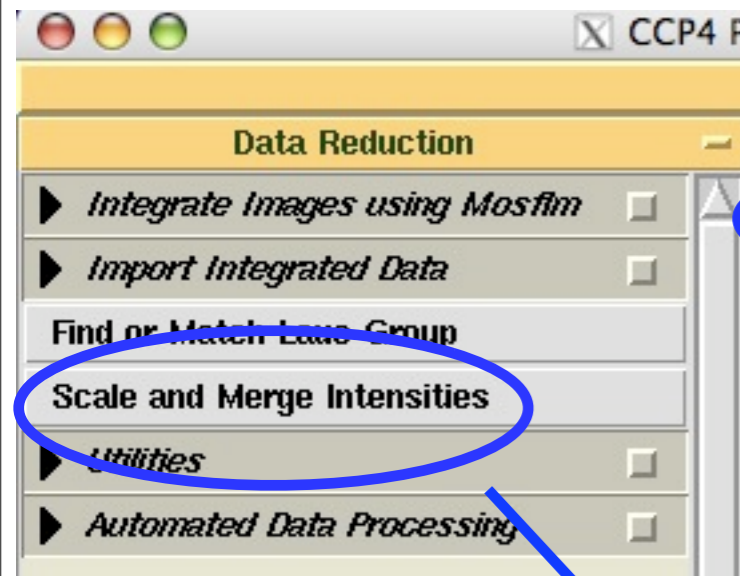
$$R_{\text{p.i.m.}} = \sum \sqrt{(1/n-1)} | I_{hl} - \langle I_h \rangle | / \sum | \langle I_h \rangle |$$

“Precision-indicating R-factor” gets better (smaller) with increasing multiplicity, ie it estimates the precision of the merged $\langle I \rangle$

CC pairwise correlation coefficients (see later)

Running SCALA from ccp4i interface

Click if you have anomalous scattering (changes the statistics and the outlier rejection)



A screenshot of the SCALA dialog box. The 'Job title' is 'Gamma ear Xe'. The 'Separate anomalous pairs for merging statistics' checkbox is checked and circled in blue. The 'MTZ in' field is 'Examples - gamma_xe1.mtz' and is circled in blue with a label 'Input file'. The 'MTZ out' field is 'Examples - gamma_xe1_scala1.mtz'. The 'Scaling Protocol' section is expanded, showing 'Scale on rotation axis with secondary beam correction' with 'isotropic' Bfactor scaling. The 'rotation interval' is 5 and the 'Secondary beam correction maximum number of spherical harmonics' is 6. The 'Independent Bfactors defined by rotation interval' is 20. The 'Apply tails correction' checkbox is unchecked. The 'Run' button is highlighted.

Usually use the default scaling options

Running SCALA from ccp4i interface

CCP4 P

Data Reduction

- Integrate Images using Mosfilm
- Import Integrated Data
- Find or Match Low Group
- Scale and Merge Intensities**
- Utilities
- Automated Data Processing

Scala - Scale Experimental Intensities

Job title

- Customise Scala process (default is to refine & apply scaling)
- Separate anomalous pairs for merging statistics
- Run to output Wilson plot and SFs after scaling and output a single MTZ file
- Ensure unique data & add FreeR column for fraction of data.
- Generate Patterson map and do peaksearch to check for pseudo-translations

MTZ in

- Override automatic definition of 'runs' to mark discontinuities in data
- Exclude data resolution less than Angstrom or greater than Angstrom

MTZ out

Convert to SFs & Wilson Plot

Use as identifier to append to column labels

Scala - Scale Experimental Intensities

Job title

- Customise Scala process (default is to refine & apply scaling)
- Separate anomalous pairs for merging statistics
- Run to output Wilson plot and SFs after scaling and output a single MTZ file
- Ensure unique data & add FreeR column for fraction of data.
- Copy FreeR from another MTZ

Extend reflector with Bfactor scaling

Use this option for your first dataset ...

... or this one for subsequent ones

Observations Used & Handling of Partial

What to look at?

View Annotated Log in Web Browser

Scala Version 3.3.15 Run at 17:31:58 on 21/12/2009 Finished with: **** Normal termination ****

Result:

Summary data for Project: Gamma Crystal: Xe1 Dataset: Xe1

	Overall	InnerShell	OuterShell
Low resolution limit	17.00	17.00	1.88
High resolution limit	1.78	5.63	1.78
Rmerge	0.034	0.025	0.196
Rmerge in top intensity bin	0.021	-	-
Rmeas (within I+/I-)	0.046	0.034	0.261
Rmeas (all I+ & I-)	0.059	0.056	0.264
Rpim (within I+/I-)	0.030	0.023	0.171
Rpim (all I+ & I-)	0.029	0.030	0.133
Fractional partial bias	-0.003	-0.002	-0.010
Total number of observations	44572	1443	4824
Total number unique	12130	435	1403
Mean((I)/sd(I))	18.0	30.0	5.6
Completeness	95.1	93.9	77.4
Multiplicity	3.7	3.3	3.4
Anomalous completeness	88.5	92.2	65.1
Anomalous multiplicity	2.1	2.1	2.0
DelAnom correlation between half-sets	0.539	0.762	-0.024
Mid-Slope of Anom Normal Probability	1.399	-	-

Outlier rejection and statistics assume that there is anomalous scattering, ie I+ differs from I-

Average unit cell: 34.16 54.80 68.00 90.00 90.00 90.00

Space group: P 21 21 21

Average mosaicity: 0.98

Minimum and maximum SD correction factors: Fulls 1.09 2.60 Partials 1.73 16.73

Dataset: Gamma/Xe1/Xe1
written as averaged data to output file /Users/pre/Projects/Xtal/Temp/Examples_2_1_mtz.tmp

Maximum resolution: 1.78A

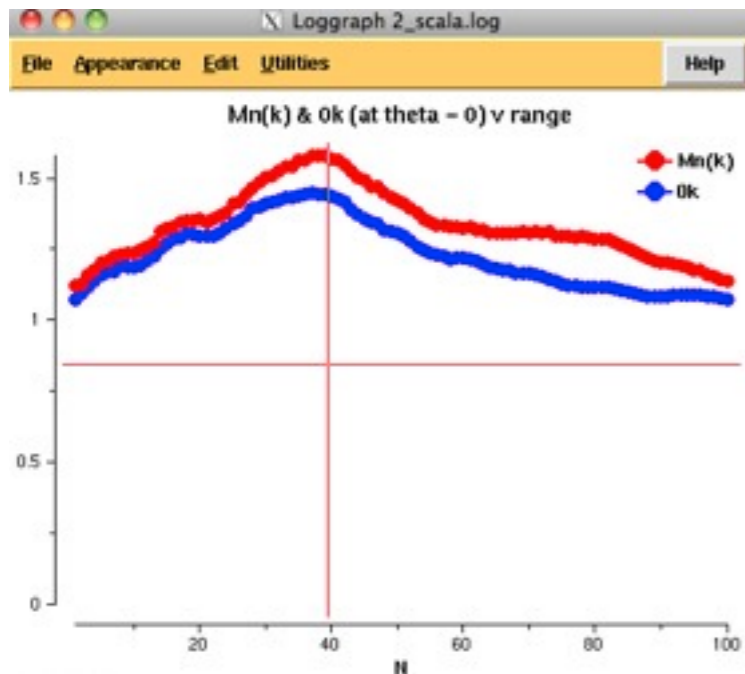
Summary

“Table I”

Graphs: Analyses against “batch” (image number or “time”)

- check for level of radiation damage
 - if you cut back from the end, there is a trade-off between damage and completeness
- check for bad images or regions

A good case



39.33,0.845

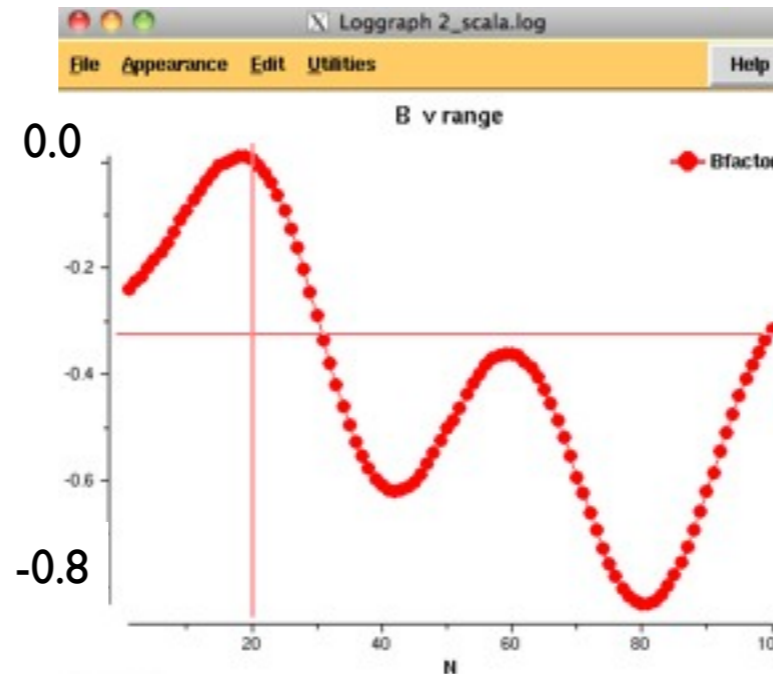
Tables in File

- >>> Scales v rotation range, %e1
- Analysis against all Batches for all runs, %e1
- Analysis against resolution, %e1
- Analysis against intensity, %e1
- Completeness, multiplicity, Rmeas v. resolution, %e1

Graphs in Selected Table

- Mn(k) & 0k (at theta = 0) v range
- B v range
- Number rejected v range

No great difference between average scale Mn(k) & scale at $\theta=0$



19.99,-0.32

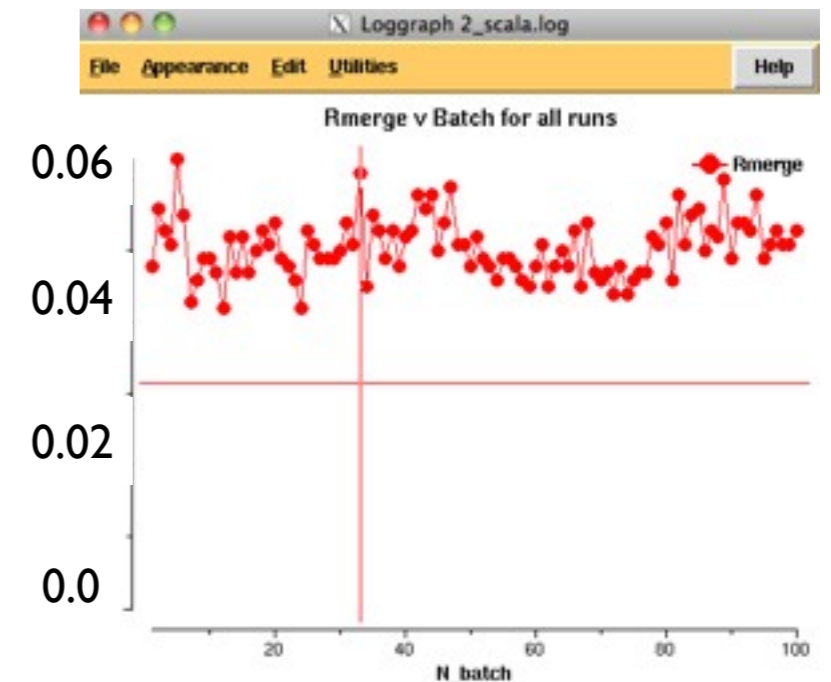
Tables in File

- >>> Scales v rotation range, %e1
- Analysis against all Batches for all runs, %e1
- Analysis against resolution, %e1
- Analysis against intensity, %e1
- Completeness, multiplicity, Rmeas v. resolution, %e1

Graphs in Selected Table

- Mn(k) & 0k (at theta = 0) v range
- B v range
- Number rejected v range

Small variation in relative B-factor



32.96,0.031

Tables in File

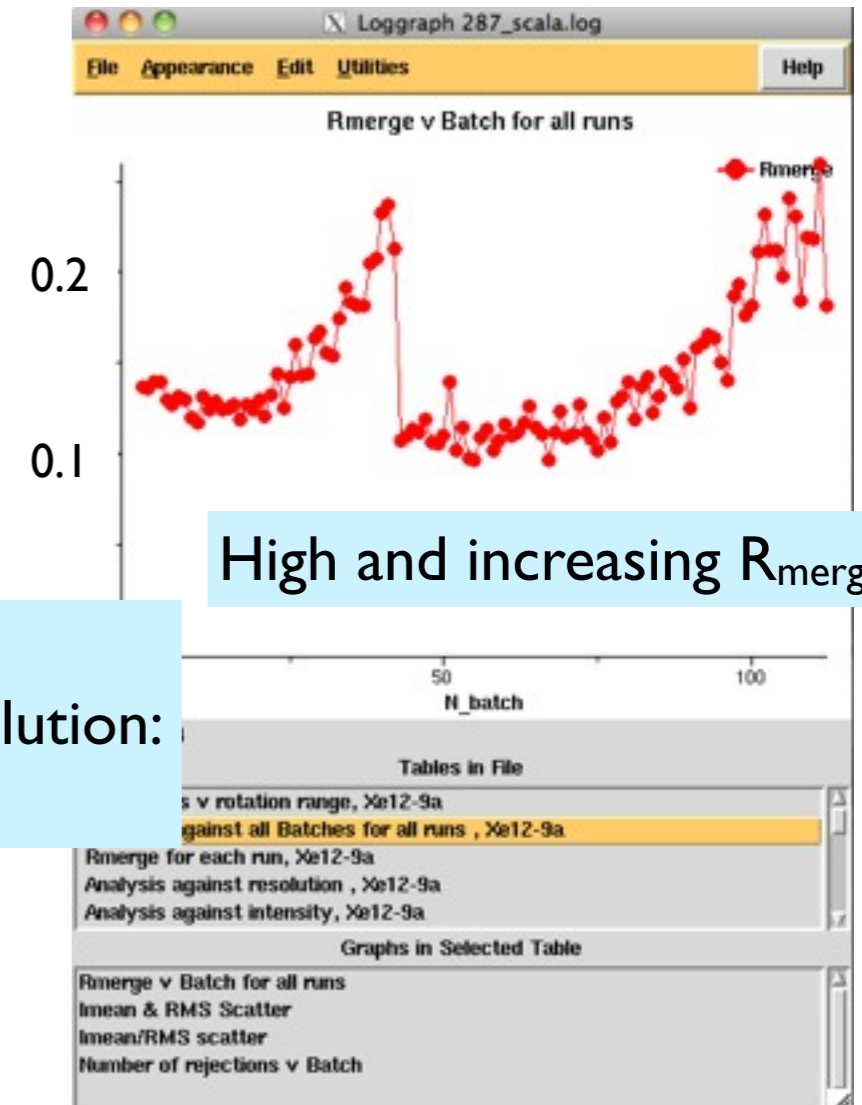
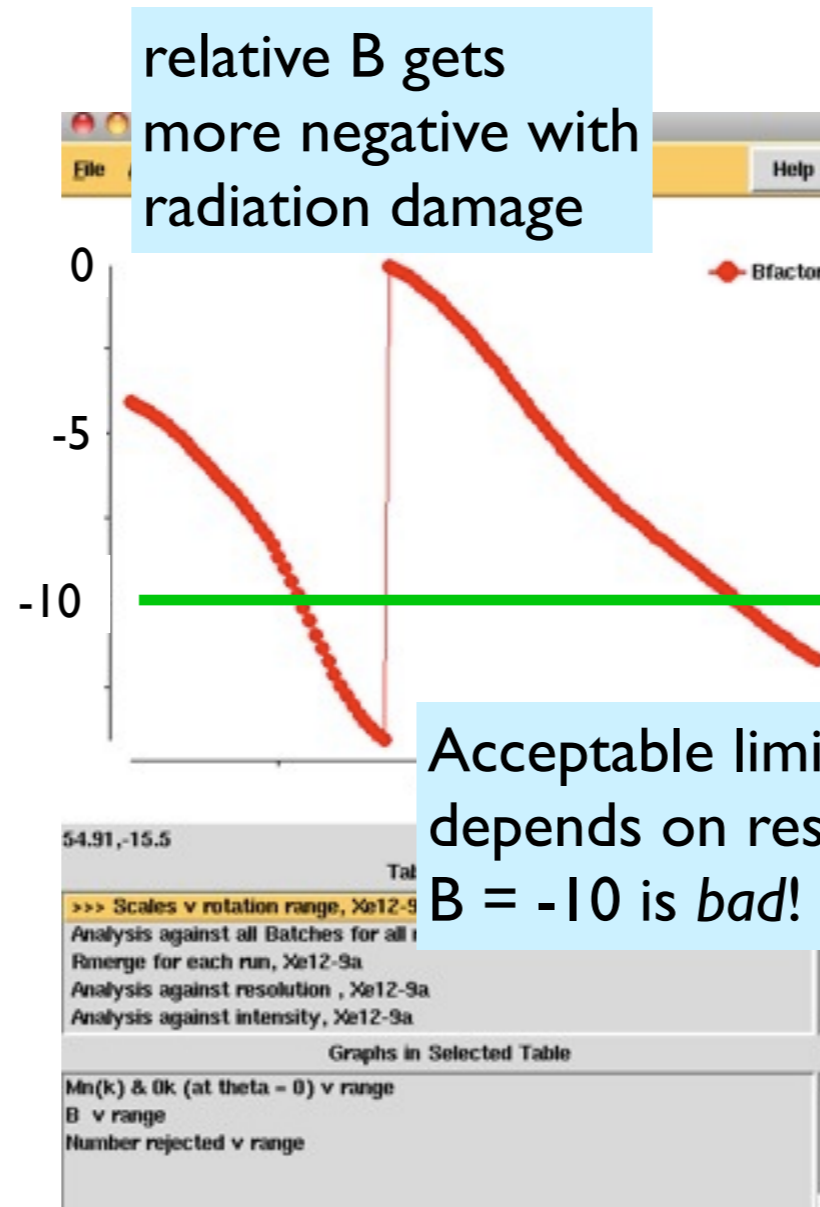
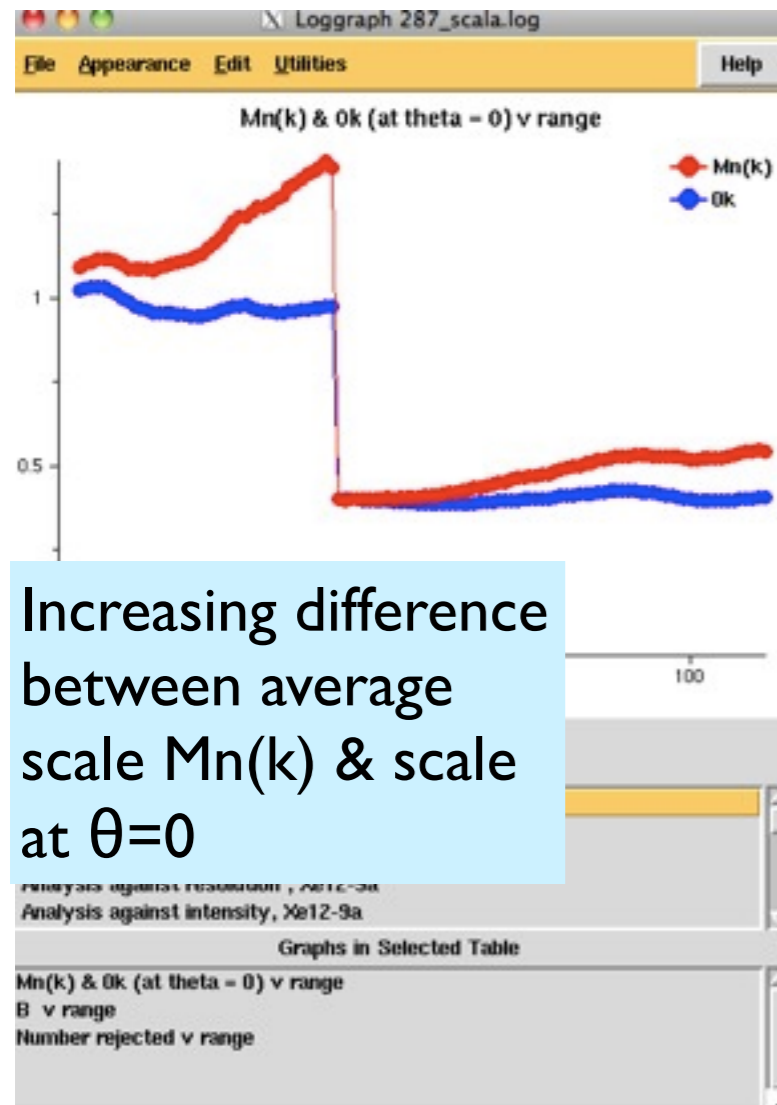
- >>> Scales v rotation range, %e1
- Analysis against all Batches for all runs, %e1
- Analysis against resolution, %e1
- Analysis against intensity, %e1
- Completeness, multiplicity, Rmeas v. resolution, %e1

Graphs in Selected Table

- Rmerge v Batch for all runs
- I mean & RMS Scatter
- I mean/RMS scatter
- Number of rejections v Batch

Uniform and low R_{merge}

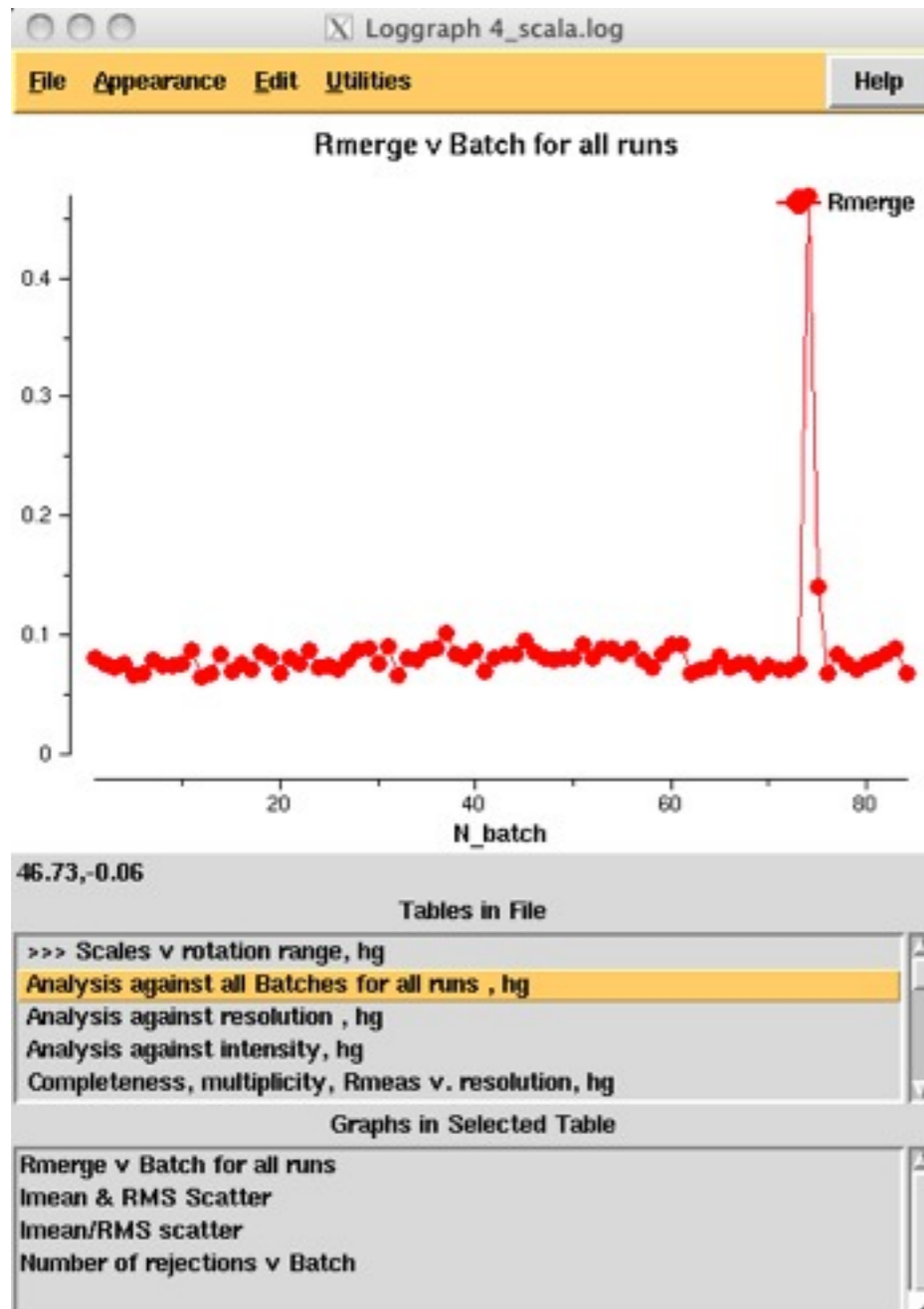
A bad case: two crystals, both dying, both incomplete



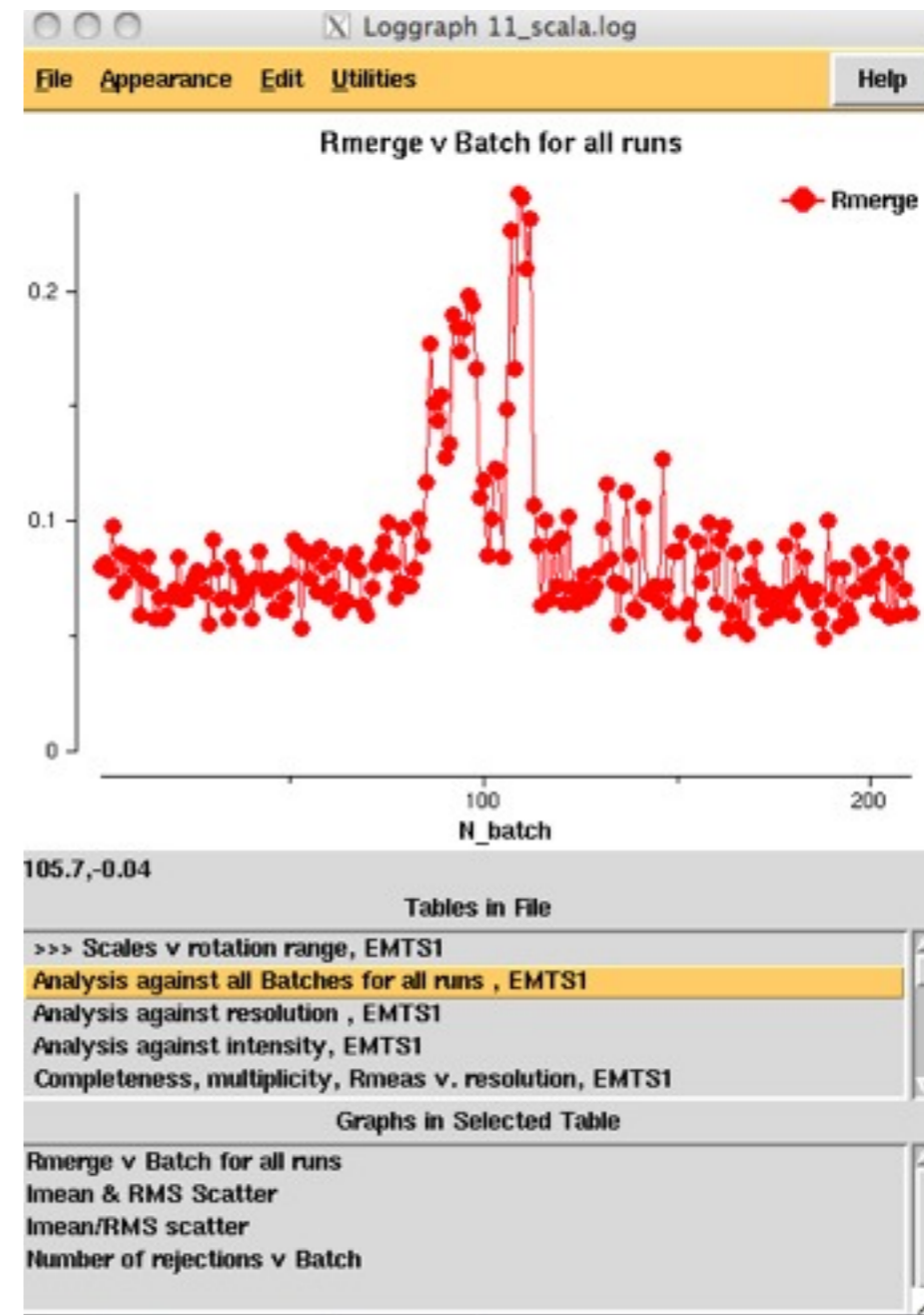
The relative B-factor gives a resolution-dependent scale factor as a function of “time” (dose): average radiation damage decay is greater at high resolution

$$k(\text{time}) = \exp[-2B(\text{time}) \sin^2\theta/\lambda^2]$$

Graph of R_{merge} vs batch may also detect individual bad images, or bad regions, that should be investigated or rejected

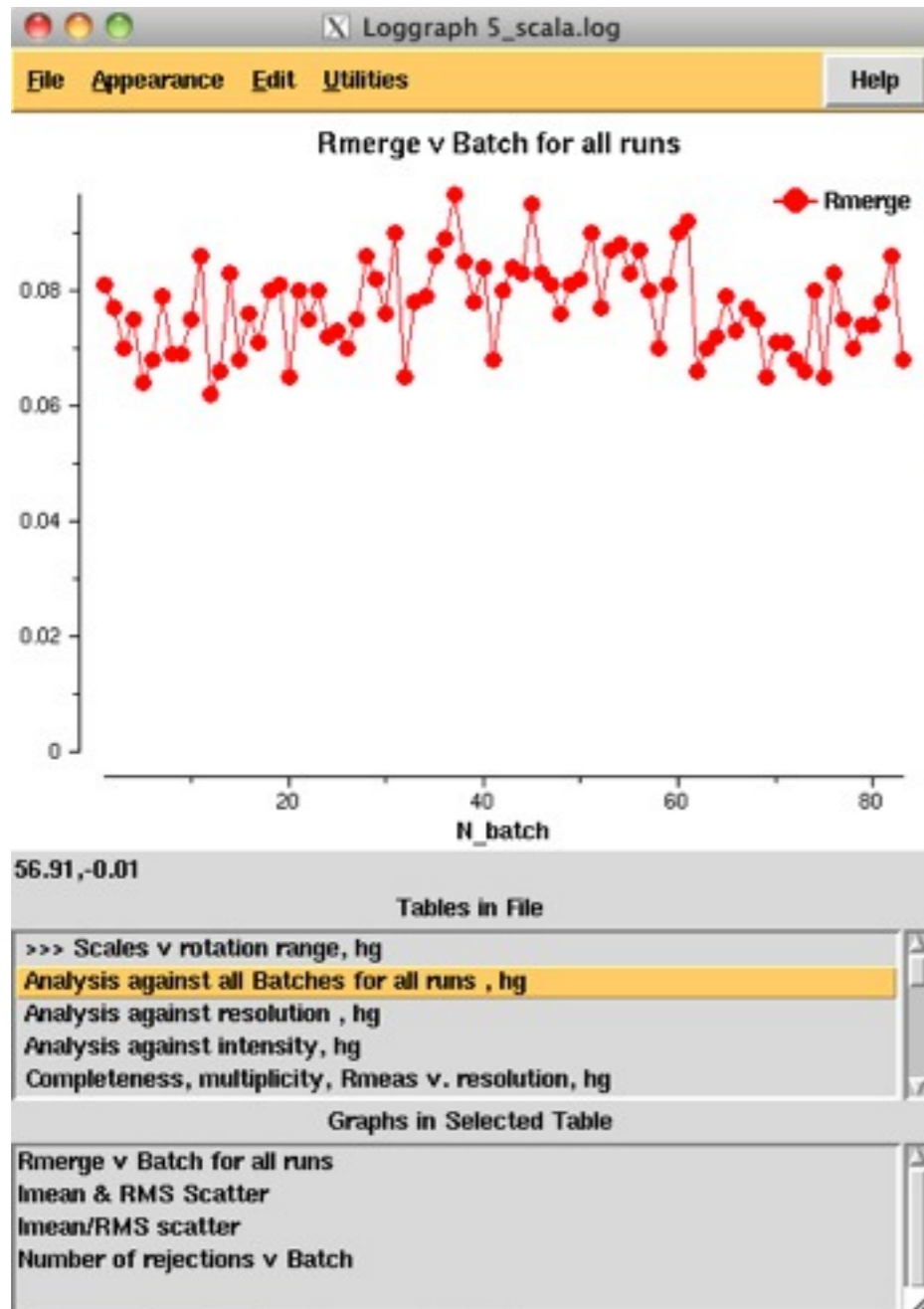


One bad (weak) image

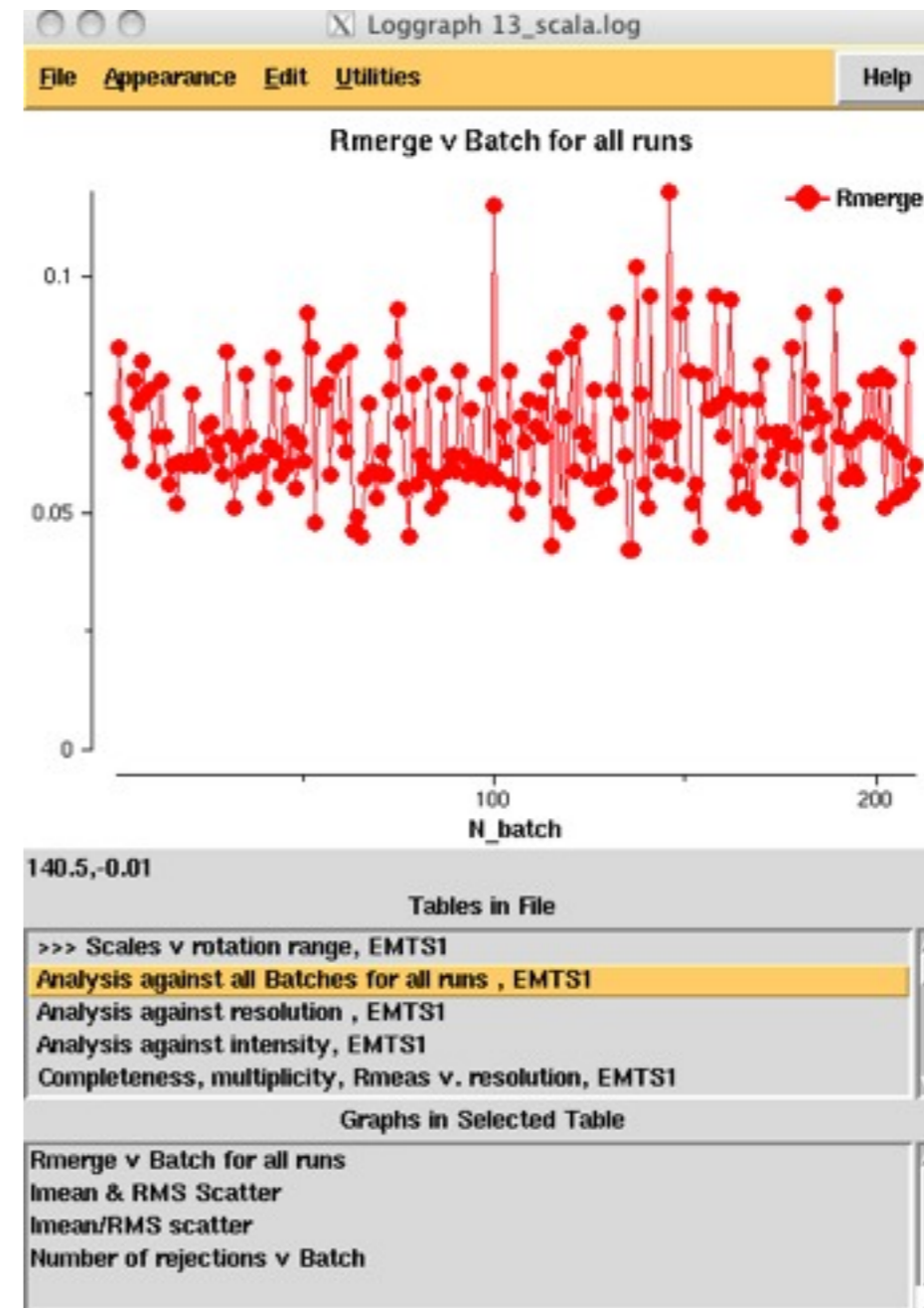


Bad region where integration had gone wrong

Graph of R_{merge} vs batch may also detect individual bad images, or bad regions, that should be investigated or rejected

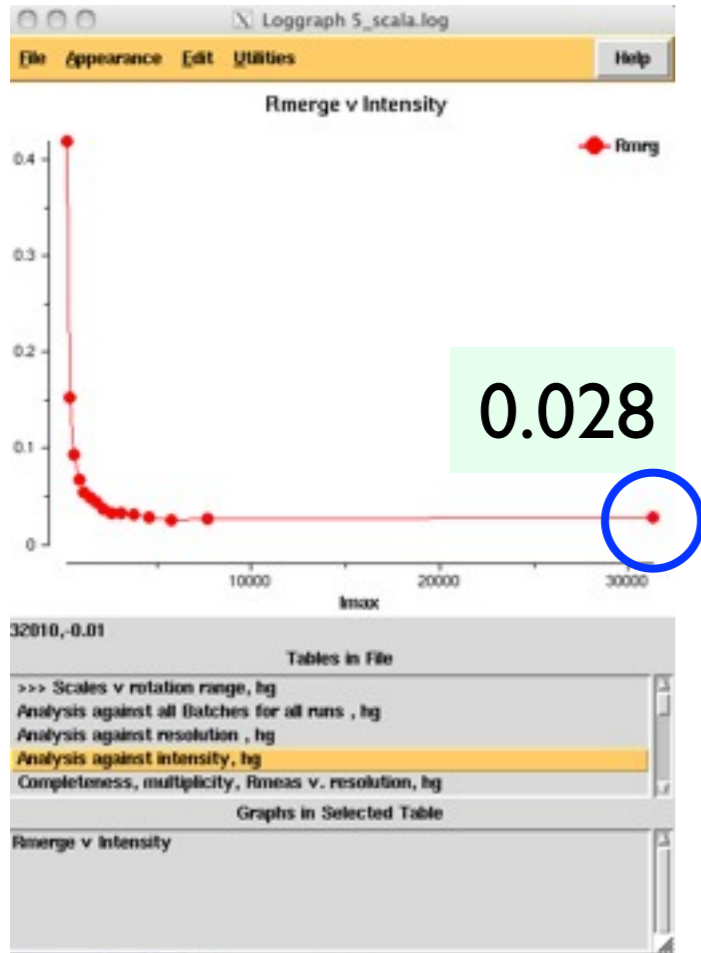


Omitting bad image

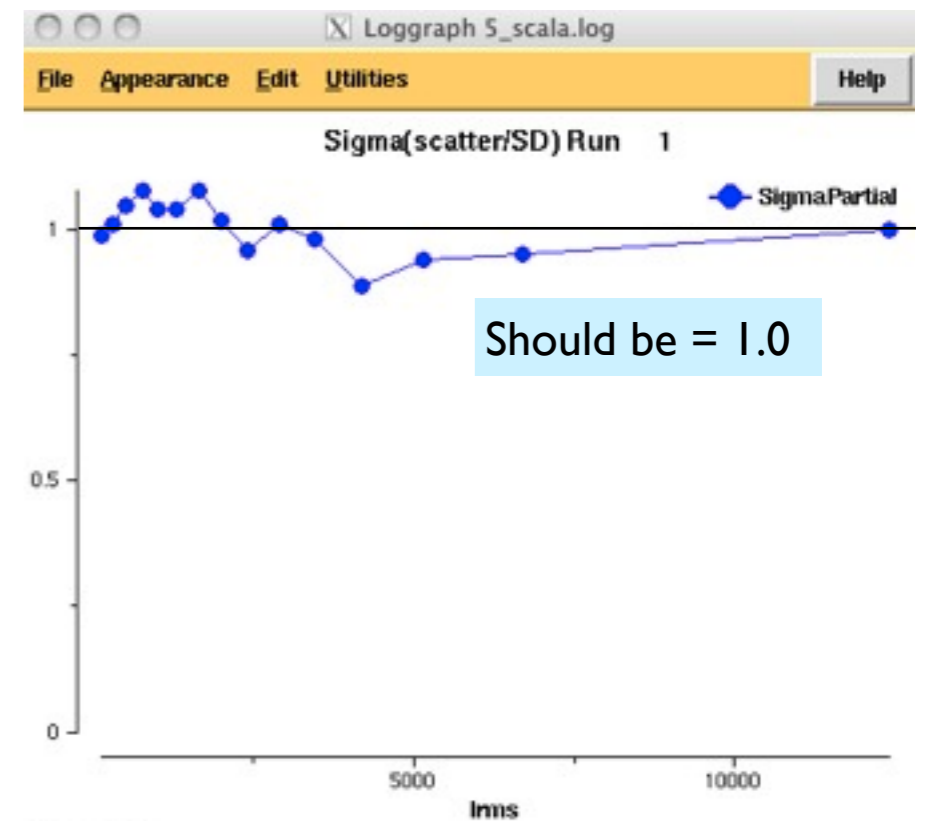


Reprocessed

Analyses against intensity



R_{merge} vs. I not generally useful (since R is a fractional measure, it will always be large for small I), but the value in the top intensity bin should be small



Improved estimate of $\sigma(I)$

The error estimate $\sigma(I)$ from the integration program is too small particularly for large intensities. A “corrected” value may be estimated by increasing it for large intensities such that the mean scatter of scaled observations on average equals $\sigma'(I)$, in all intensity ranges

$$\text{Corrected } \sigma'(I)^2 = \text{SDFac}^2 [\sigma^2 + \text{SdB} \langle I_h \rangle + (\text{SdAdd} \langle I_h \rangle)^2]$$

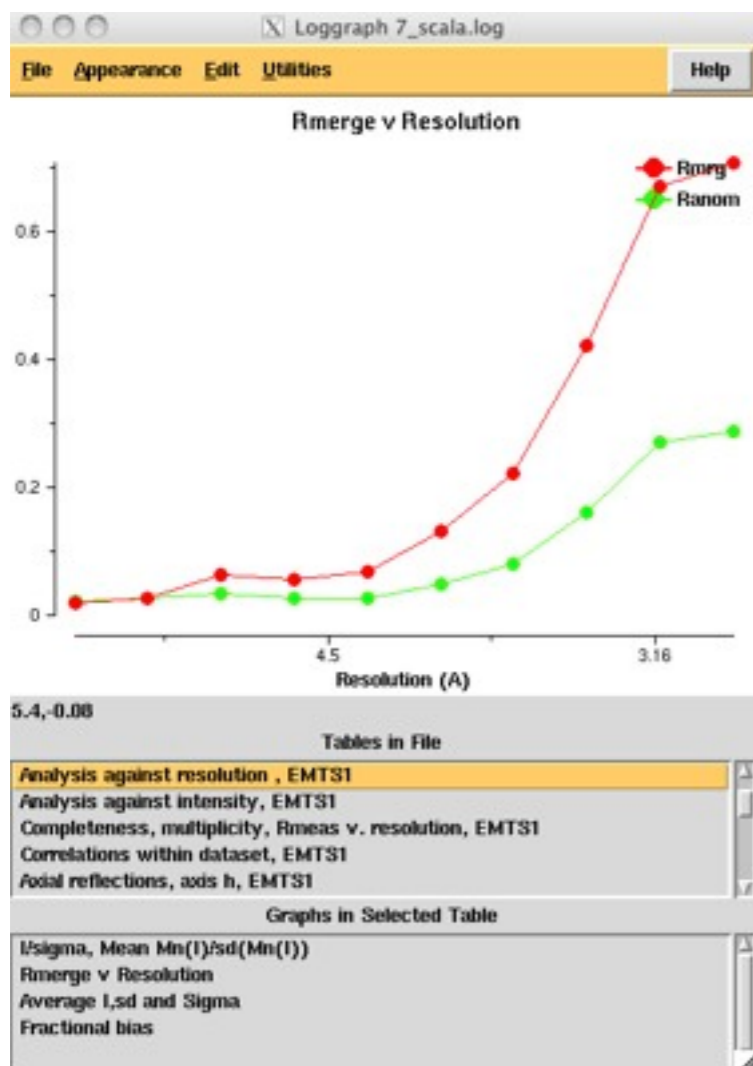
SDFac, **SdB** and **SdAdd** are adjustable parameters

Analyses against resolution

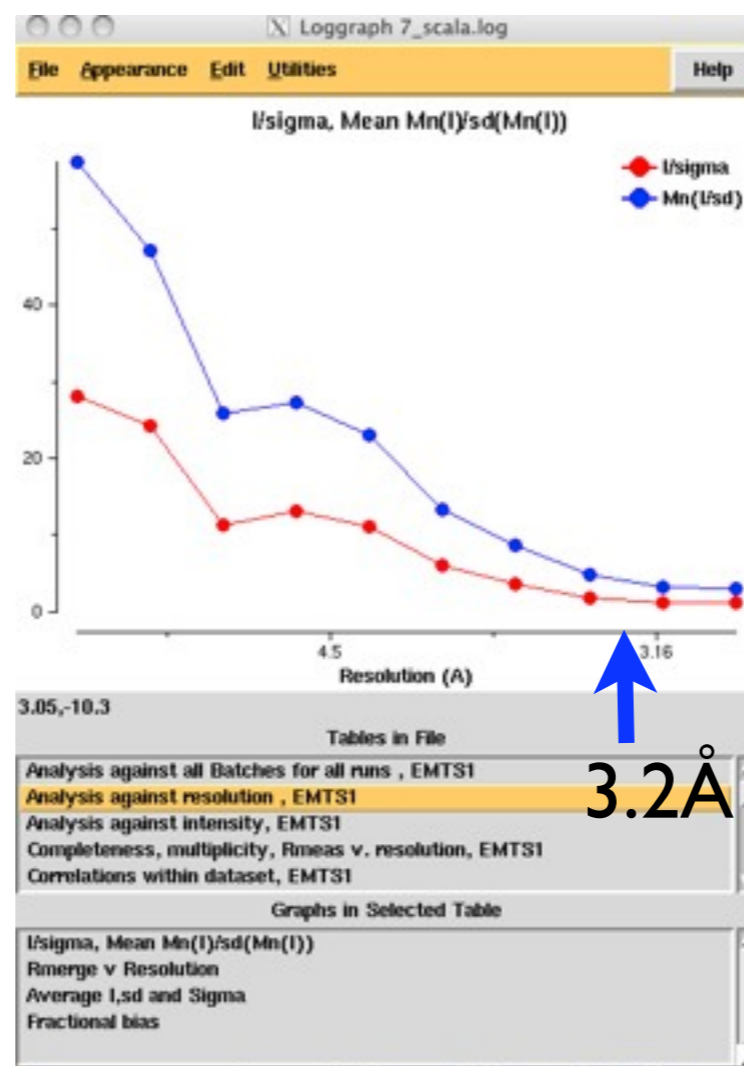
What is the real resolution? not an easy question to answer

May depend on what you want the data for: more stringent for experimental phasing than for refinement

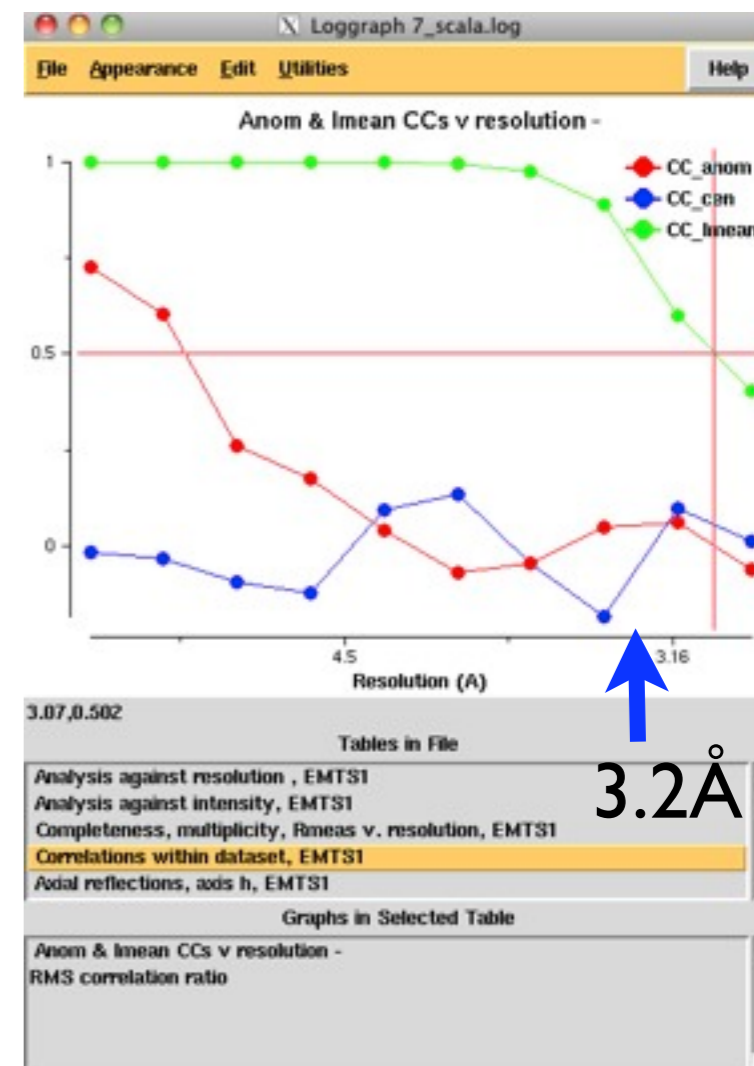
Anisotropic data needs a less stringent overall cut-off to keep best data



R_{merge} is not particularly useful: it gets higher at high resolution



$\langle I/\sigma(I) \rangle$ after merging (blue line) should be $> \sim 1-2$

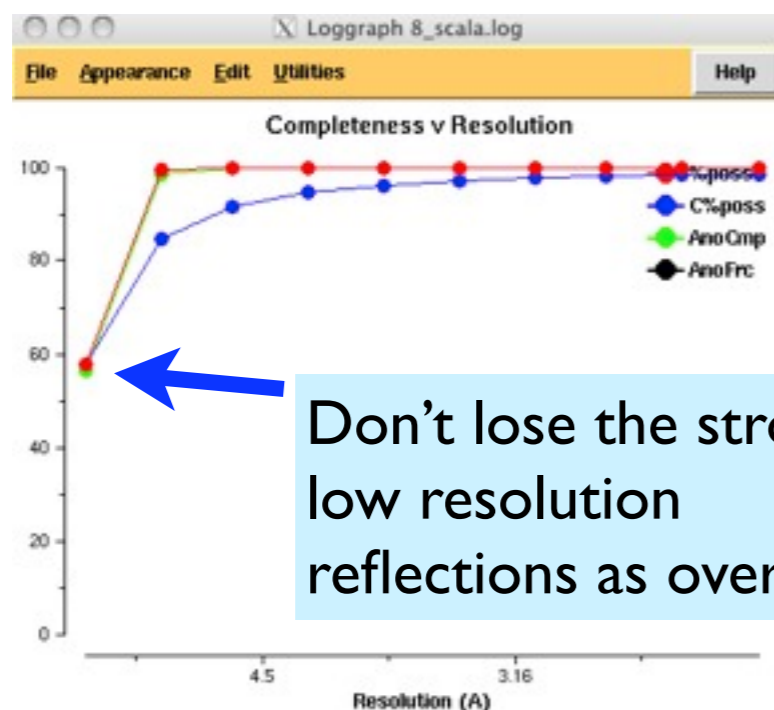


CC between random half-datasets should be $> \sim 0.5$

Completeness

Data completeness is important, preferably in all resolution shells, though you can probably get away with some incompleteness at the outer edge.

See James Holton's movies for an illustration of the importance of completeness <http://ucxray.berkeley.edu/~jamesh/movies/>



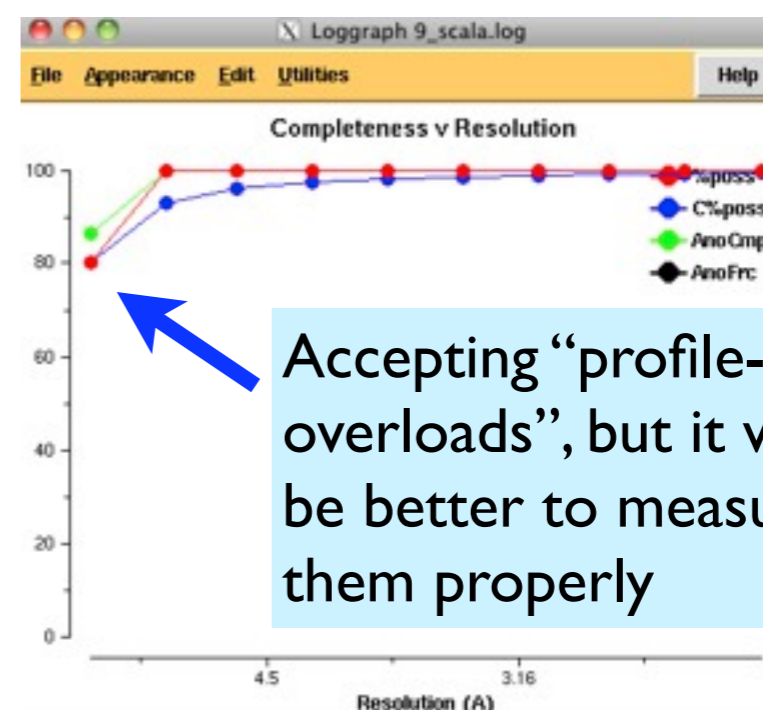
Don't lose the strong low resolution reflections as overloads

3.16,-17.5

Tables in File
>>> Scales v rotation range, N16
Analysis against all Batches for all runs , N16
Analysis against resolution , N16
Analysis against intensity, N16
Completeness, multiplicity, Rmeas v. resolution, N16

Graphs in Selected Table

Completeness v Resolution
Multiplicity v Resolution
Rpim (precision R) v Resolution
Rmeas, Rsym & PCV v Resolution



Accepting "profile-fitted overloads", but it would be better to measure them properly

3.27,-16.8

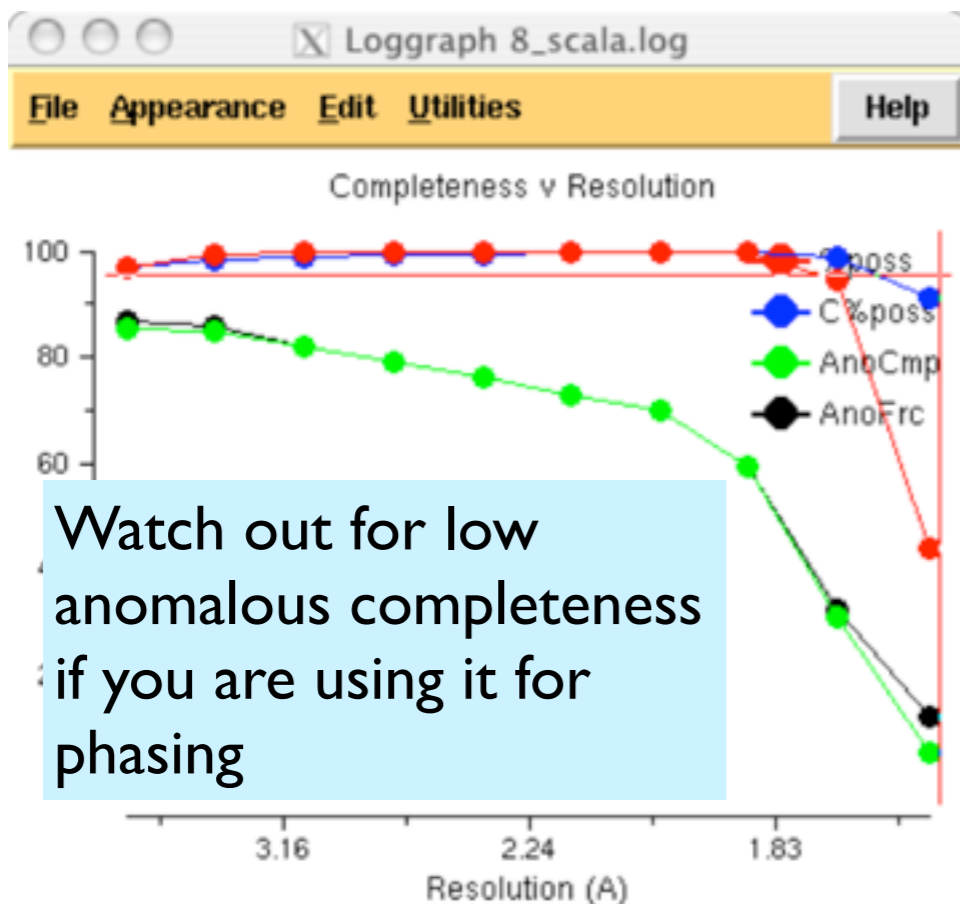
Tables in File
>>> Scales v rotation range, New
Analysis against all Batches for all runs , New
Analysis against resolution , New
Analysis against intensity, New
Completeness, multiplicity, Rmeas v. resolution, New

Graphs in Selected Table

Completeness v Resolution
Multiplicity v Resolution
Rpim (precision R) v Resolution
Rmeas, Rsym & PCV v Resolution

Completeness

Data completeness is important, preferably in all resolution shells, though you can probably get away with some incompleteness at the outer edge.



Watch out for low anomalous completeness if you are using it for phasing

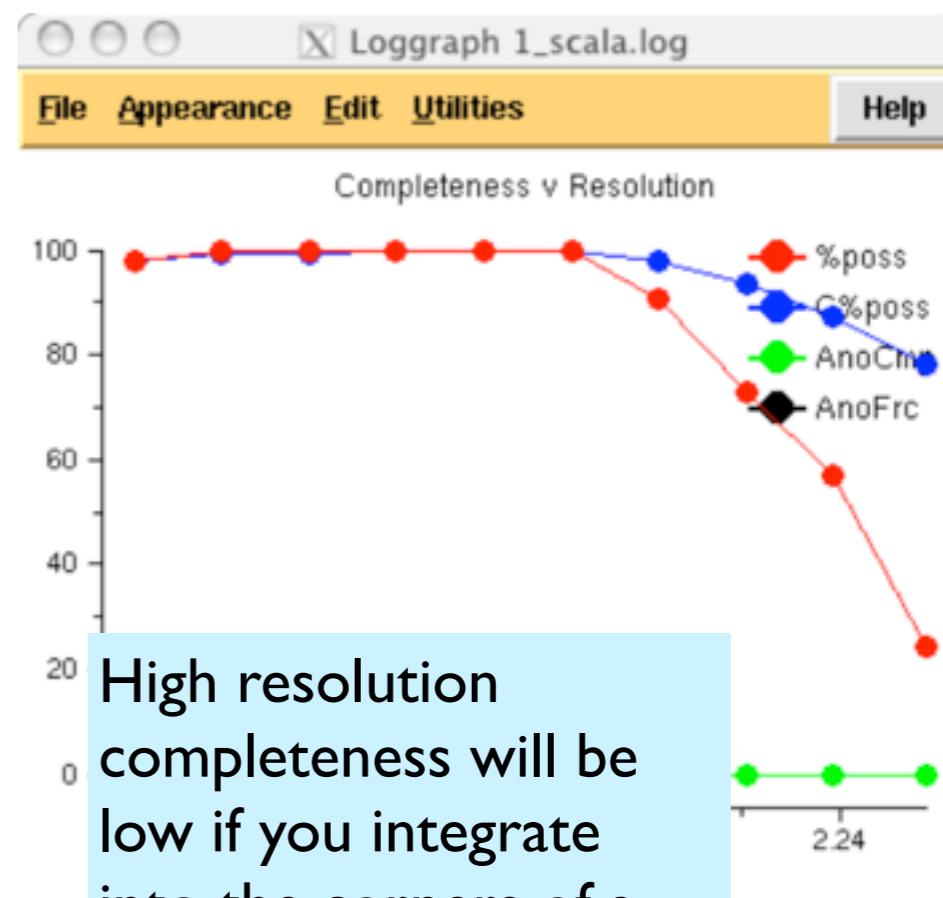
1.65,95.67

Tables in File

- >>> Scales v rotation range, N
- Analysis against Batch, N
- Analysis against resolution, N
- Analysis against intensity, N
- Completeness, multiplicity, Rmeas v. resolution, N**

Graphs in Selected Table

- Completeness v Resolution
- Multiplicity v Resolution
- Rpim (precision R) v Resolution
- Rmeas, Rsym & PCV v Resolution



High resolution completeness will be low if you integrate into the corners of a square detector

2.34

Analysis against intensity, p2x

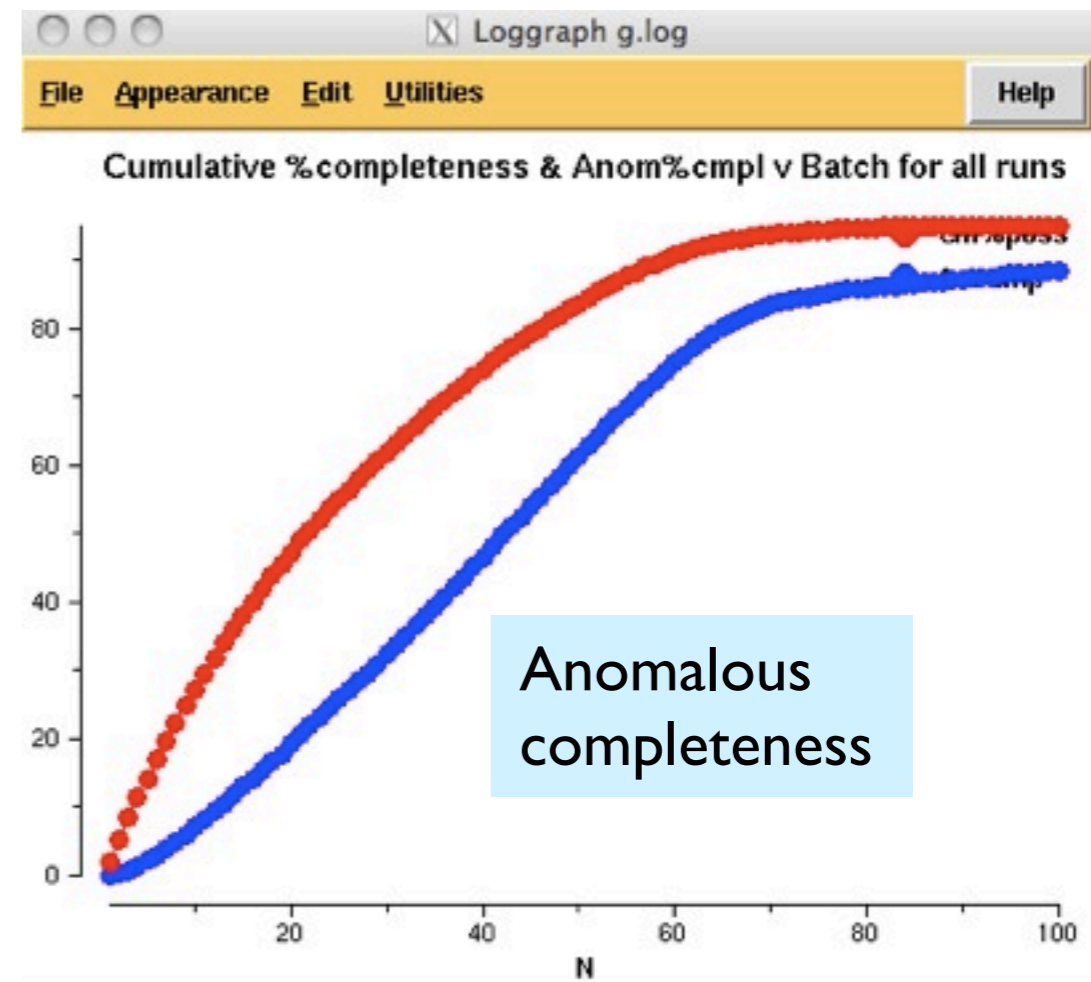
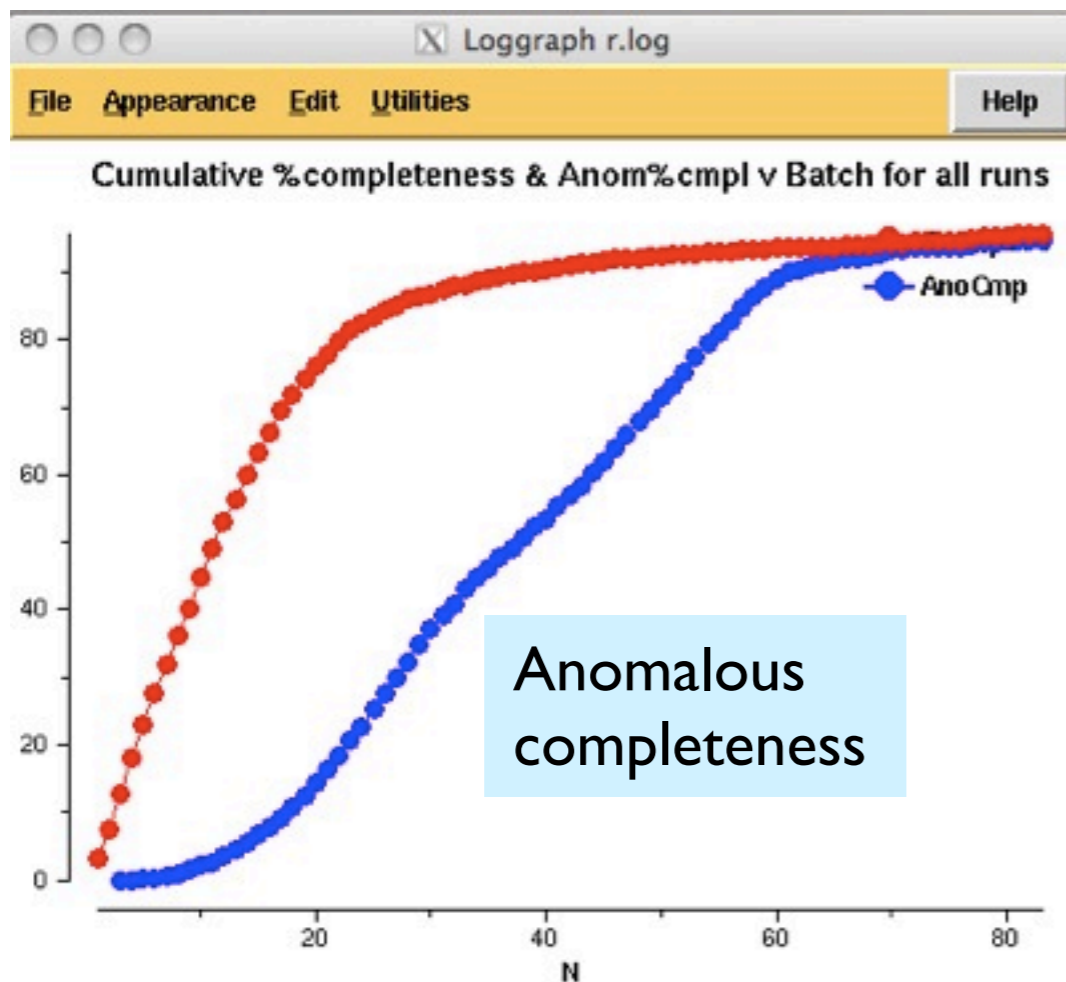
- Completeness, multiplicity, Rmeas v. resolution, p2x**
- Correlations within dataset, p2x
- Axial reflections, axis k, p2x
- Axial reflections, axis l, p2x

Graphs in Selected Table

- Completeness v Resolution
- Multiplicity v Resolution
- Rpim (precision R) v Resolution
- Rmeas, Rsym & PCV v Resolution

Completeness

Data completeness is important, preferably in all resolution shells, though you can probably get away with some incompleteness at the outer edge.



84.93,109.0

Tables in File

- >>> Scales v rotation range, New
- Analysis against all Batches for all runs, New
- Analysis against resolution, New
- Analysis against resolution, with & without anomalous (Ov), New
- Analysis against intensity, New

Graphs in Selected Table

- Rmerge v Batch for all r
- Cumulative %complemen
- lmean & RMS Scatter
- lmean/RMS scatter
- Number of rejects

0,0

Tables in File

- >>> Scales v rotation range, Xe1
- Analysis against all Batches for all runs, Xe1
- Analysis against resolution, Xe1
- Analysis against resolution, with & without anomalous (Ov), Xe1
- Analysis against intensity, Xe1

Graphs in Selected Table

- lmean/RMS scatter
- Number of rejects

Cumulative completeness against batch

Graph not yet available!

Outliers

Detection of outliers is easiest if the multiplicity is high

Removal of spots behind the backstop shadow does not work well at present: usually it rejects all the good ones, so tell Mosflm where the backstop shadow is.

Reasons for outliers

- outside reliable area of detector (eg behind shadow)

specify backstop shadow, calibrate detector

- ice spots

do not get ice on your crystal!

- multiple lattices

find single crystal

- zingers

- bad prediction (spot not there)

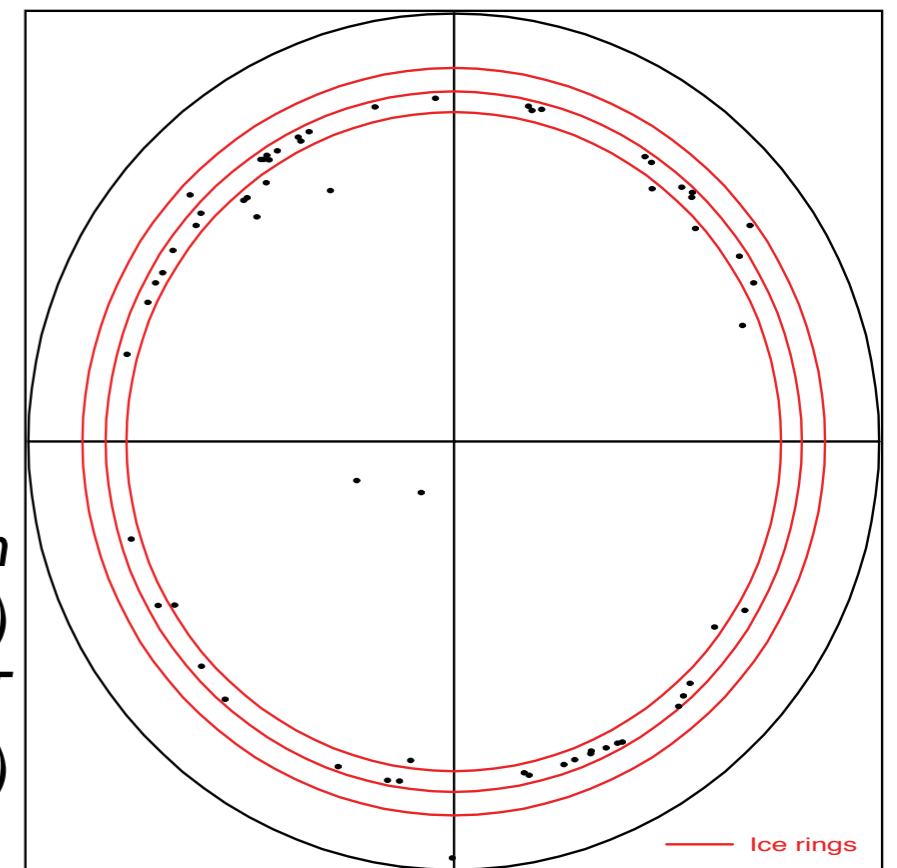
improve prediction

- spot overlap

lower mosaicity, smaller slice, move detector back

deconvolute overlaps

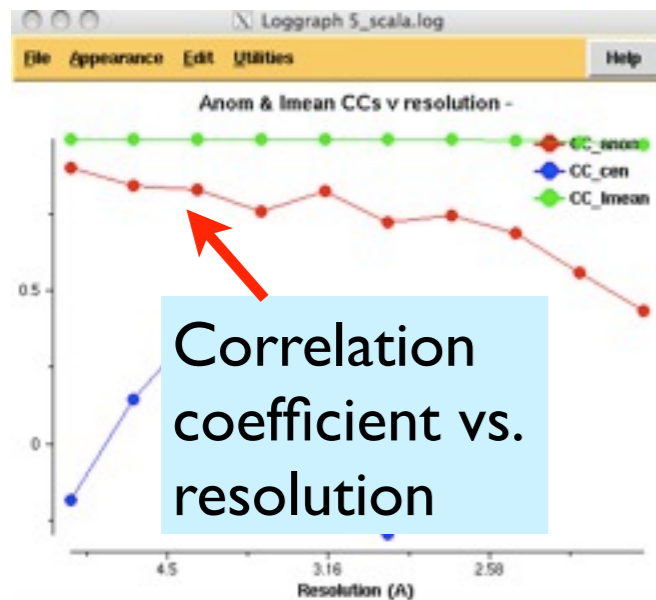
Rejects lie on
ice rings (red)
(ROGUEPLOT
in Scala)



Position of rejects on detector

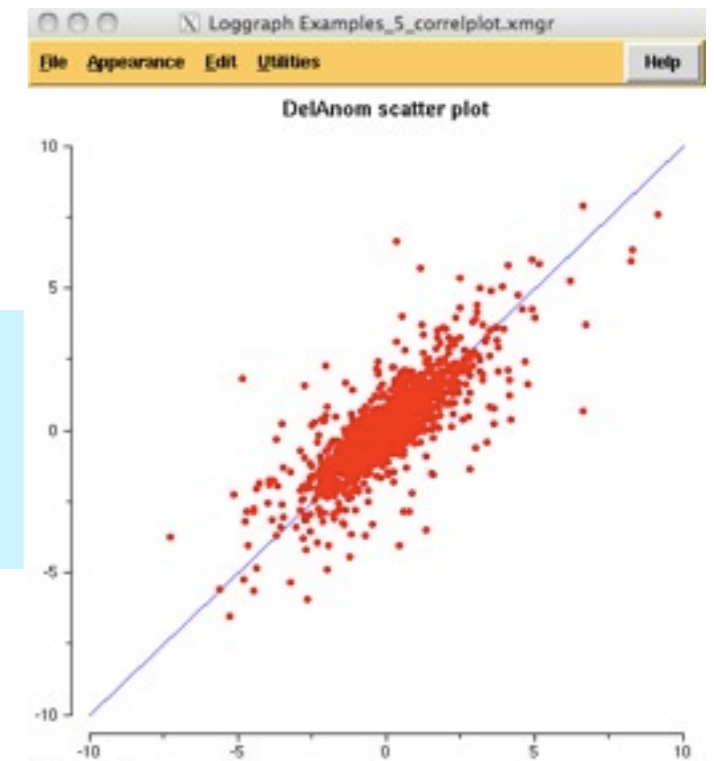
Detecting anomalous signals

The data contains both I+ (hkl) and I- (-h-k-l) observations and we can detect whether there is a significant difference between them.

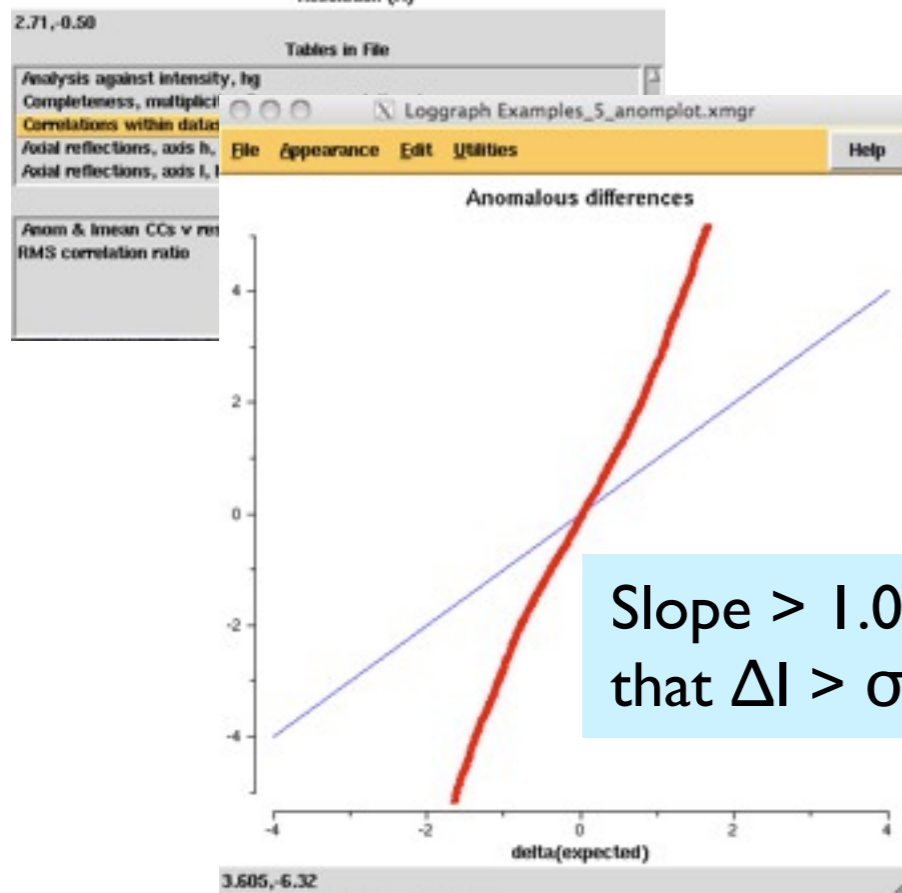


Split one dataset randomly into two halves, calculate correlation between the two halves

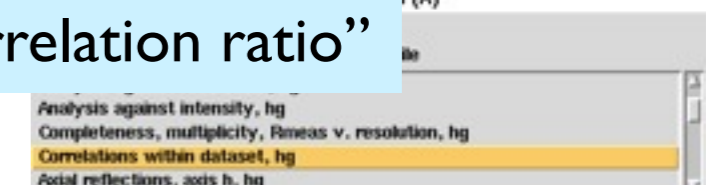
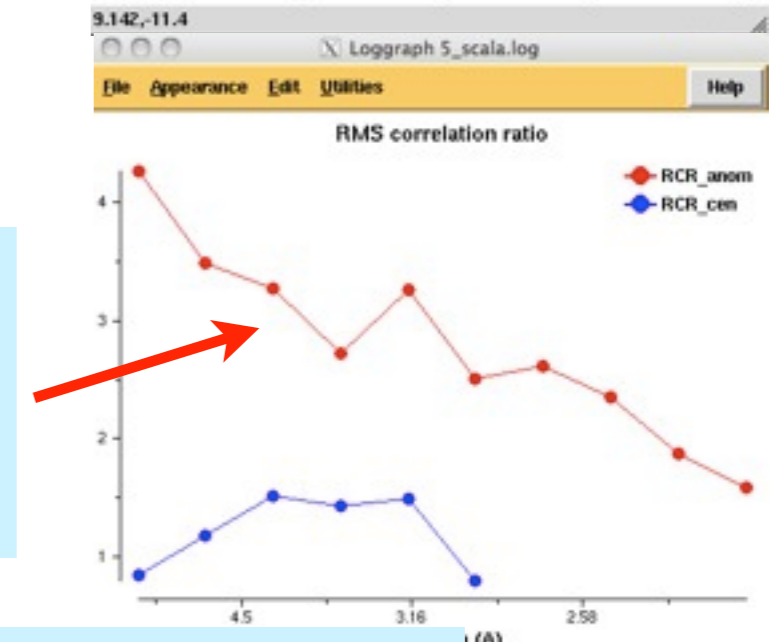
Plot ΔI_1 against ΔI_2 should be elongated along diagonal



Strong anomalous signal

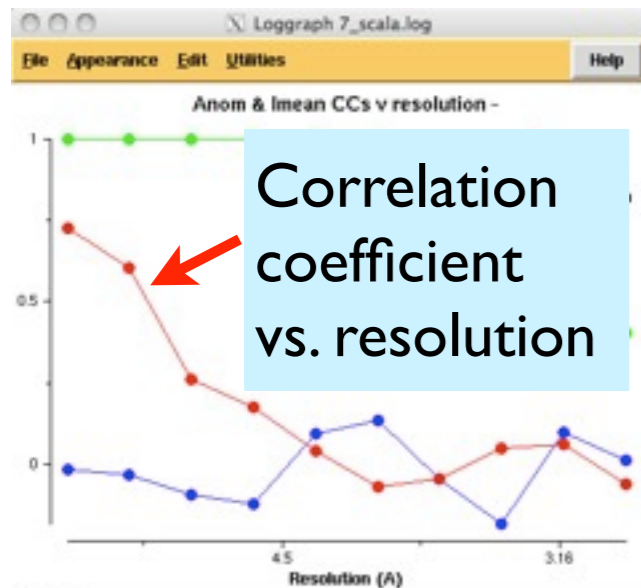


Ratio of width of distribution along diagonal to width across diagonal



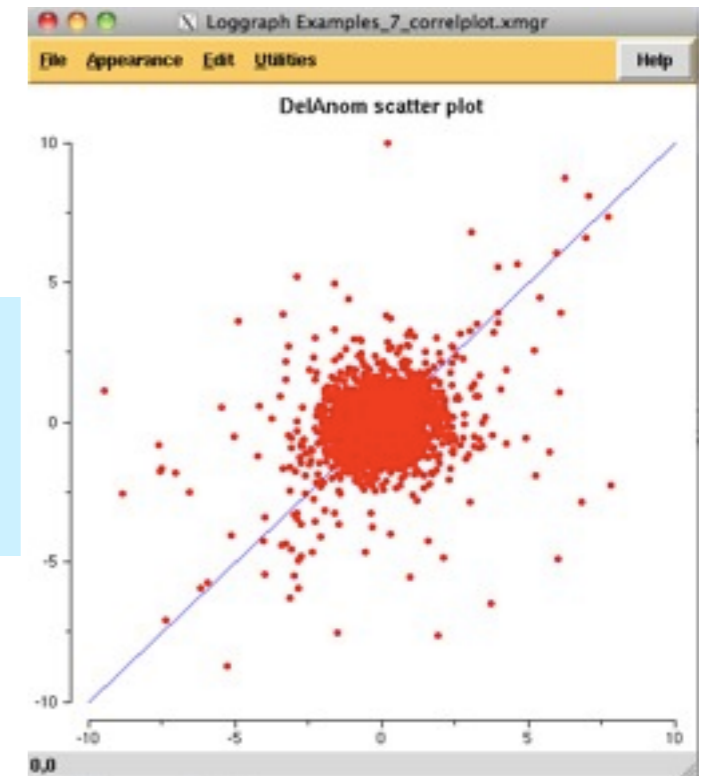
Detecting anomalous signals

The data contains both I+ (hkl) and I- (-h-k-l) observations and we can detect whether there is a significant difference between them.

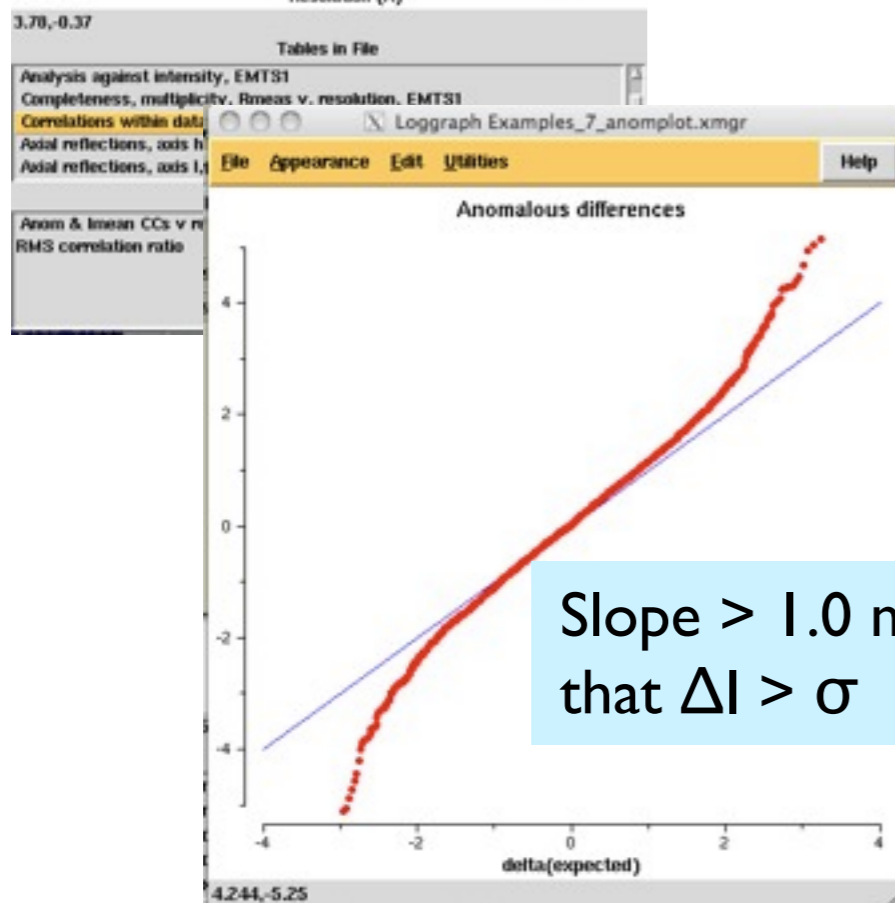


Split one dataset randomly into two halves, calculate correlation between the two halves

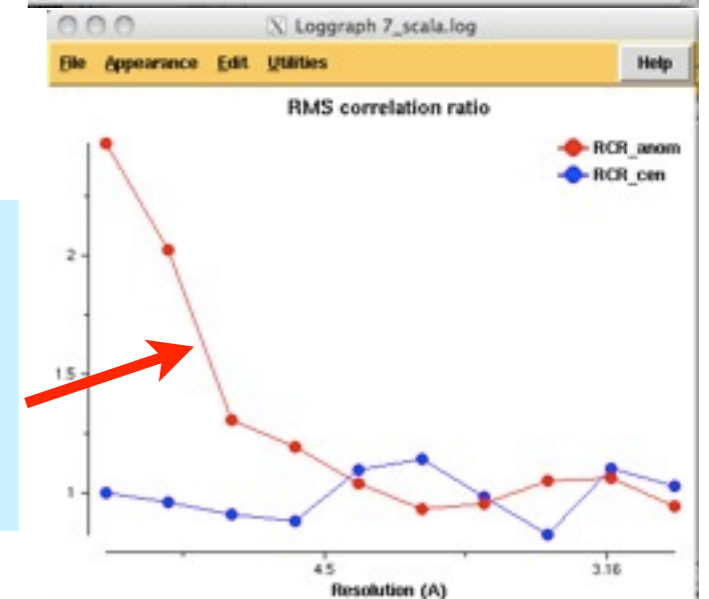
Plot ΔI_1 against ΔI_2 should be elongated along diagonal



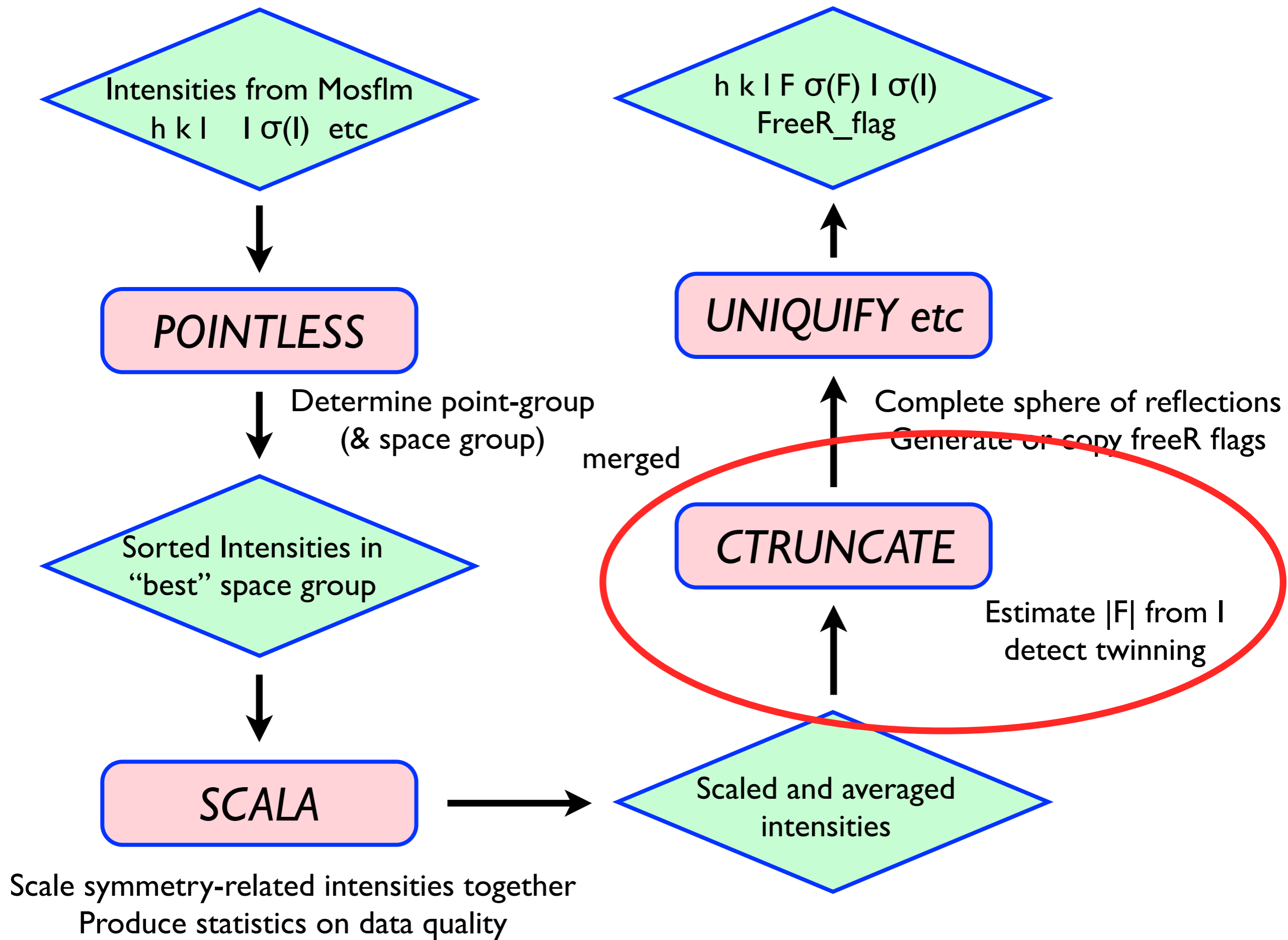
Weak but useful anomalous signal



Ratio of width of distribution along diagonal to width across diagonal



“RMS correlation ratio”



Estimation of amplitude $|F|$ from intensity I

If we knew the true intensity J then we could just take the square root

$$|F| = \sqrt{J}$$

But measured intensities I have an error $\sigma(I)$ so a small intensity may be measured as negative.

The “best” estimate of $|F|$ larger than \sqrt{I} for small intensities ($< \sim 3 \sigma(I)$) to allow for the fact that we know that $|F|$ must be positive

[c]truncate estimates $|F|$ from I and $\sigma(I)$ using the average intensity in the same resolution range: this give the prior probability $p(J)$

$$E(F ; I, \sigma(I)) = \int_0^{\infty} F p(I ; J, \sigma(I)) p(J) dJ$$

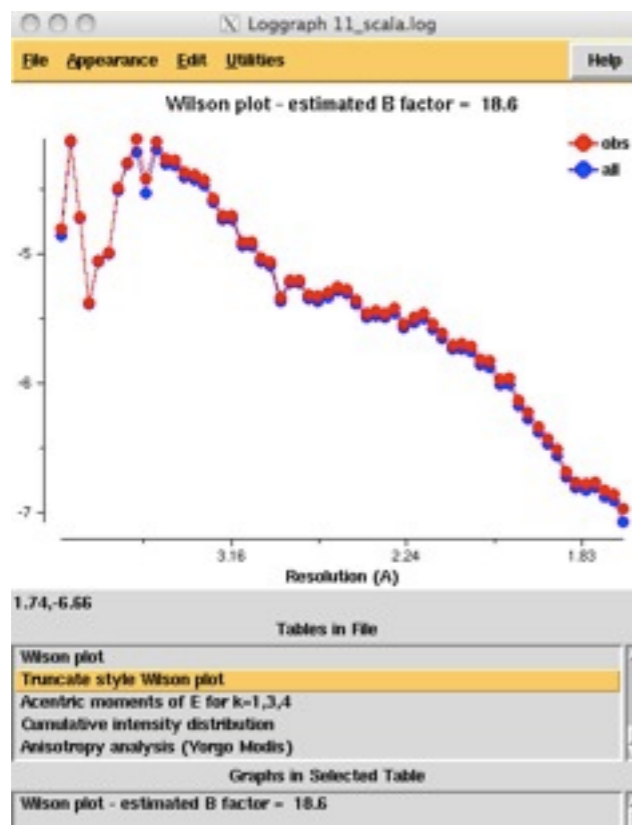
Intensity statistics

We need to look at the distribution of intensities to detect twinning

Assuming atoms are randomly placed in the unit cell, then

$$\langle I \rangle(s) = \langle F F^* \rangle(s) = \sum_j g(j, s)^2$$

where $g(j, s)$ is the scattering from atom j at $s = \sin\theta/\lambda$



Average intensity falls off with resolution, mainly because of atomic motions (B-factors)

For the purposes of looking for crystal pathologies, we are not interested in the variation with resolution, so we can use “normalised” intensities which are independent of resolution

$$\langle I \rangle(s) = C \exp(-2 B s^2)$$

Wilson plot: $\log(\langle I \rangle(s))$ vs s^2

This would be a straight line if all the atoms had the same B-factor

Normalised intensities: relative to average intensity at that resolution

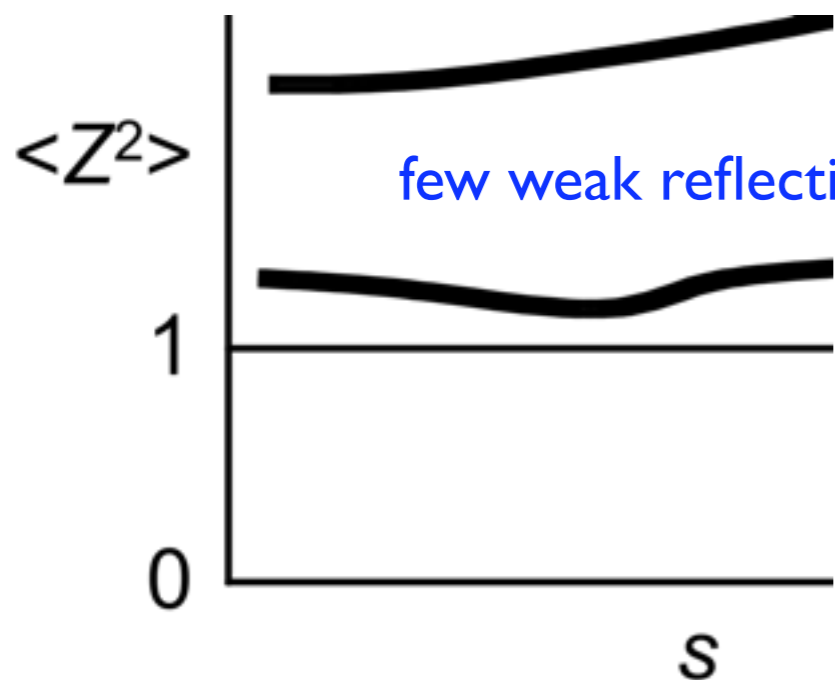
$$Z(h) = I(h)/\langle I(s) \rangle \approx |E|^2$$

$$\langle Z(s) \rangle = 1.0 \text{ by definition}$$

$$\langle Z^2(s) \rangle > 1.0 \text{ depending on the distribution}$$

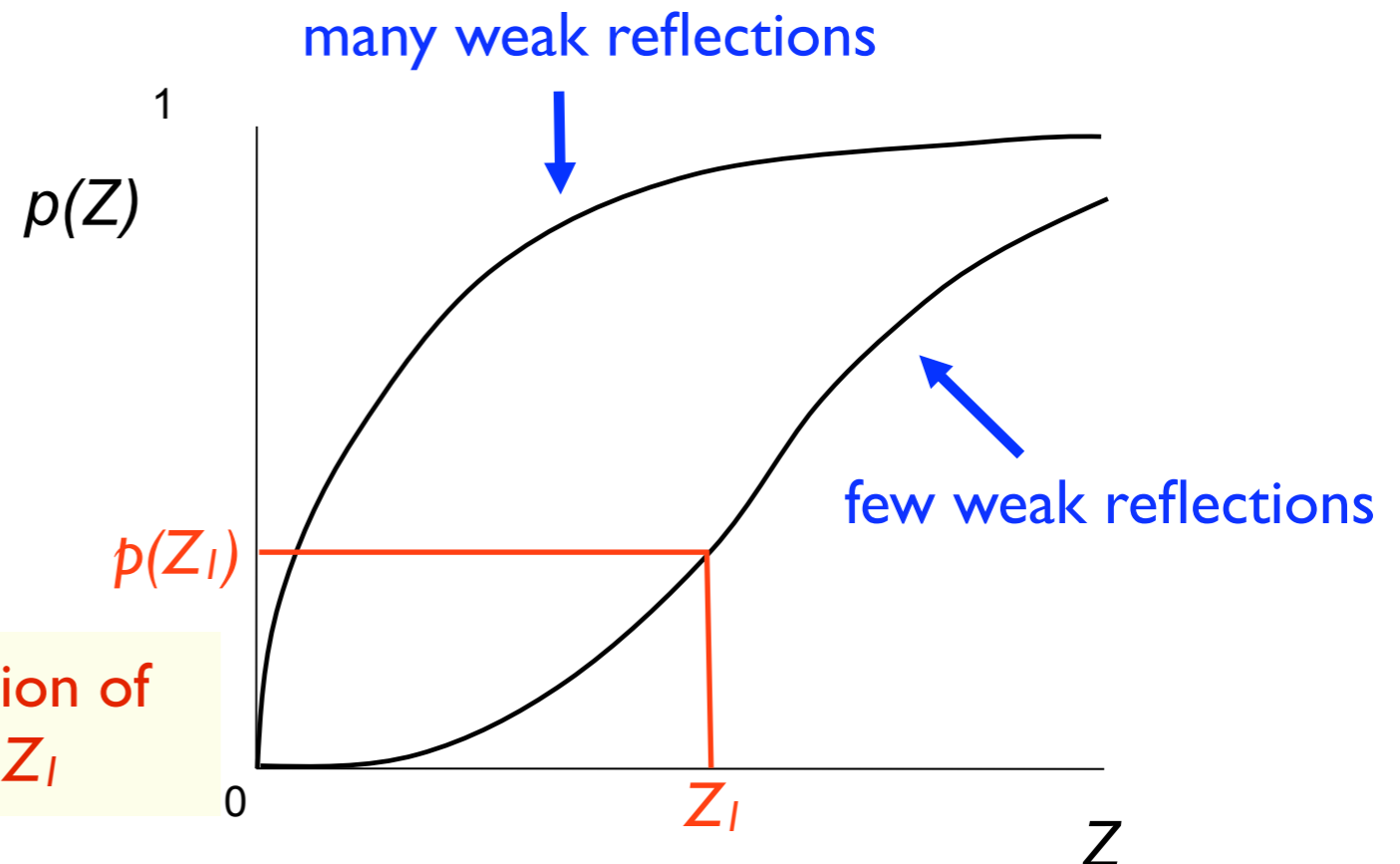
$\langle Z^2(s) \rangle$ is larger if the distribution of intensities is wider: it is the 2nd moment ie the *variance* (this is the 4th moment of E)

many weak reflections



few weak reflections

Cumulative distribution of Z: $p(Z)$ vs. Z



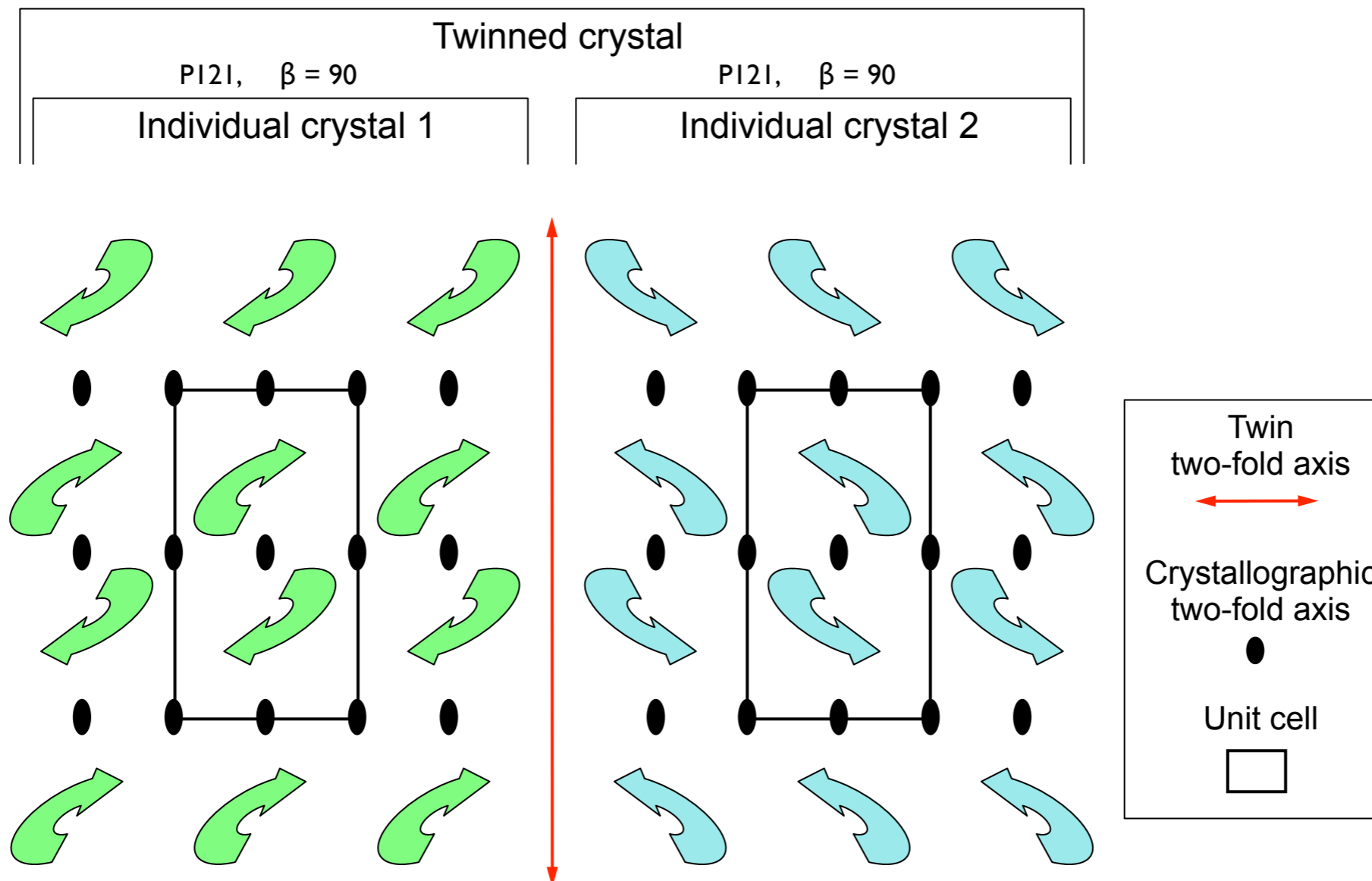
$p(Z_1)$ is the proportion of reflections with $Z < Z_1$

Twinning by (pseudo)merohedry

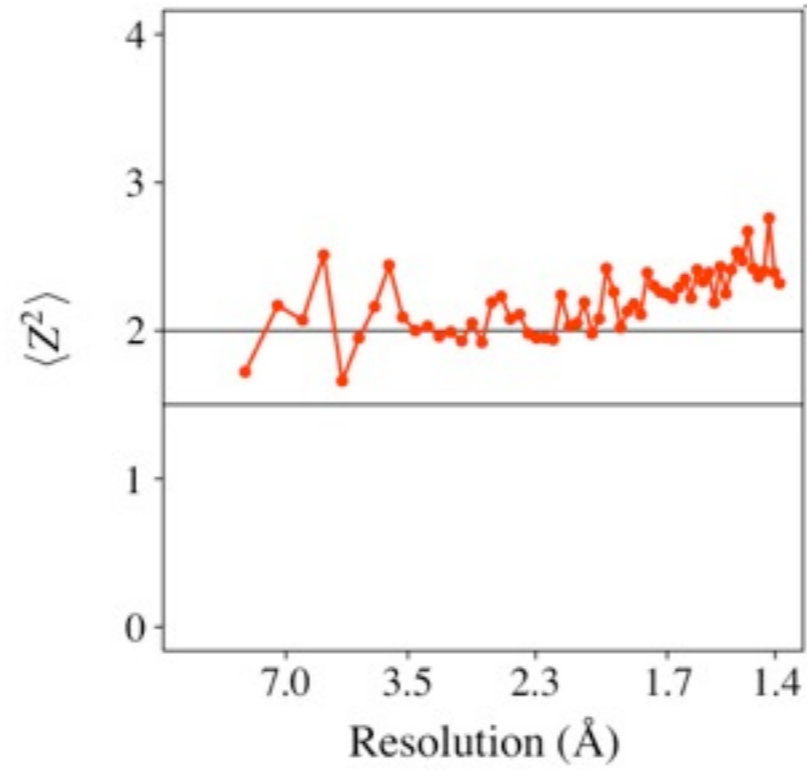
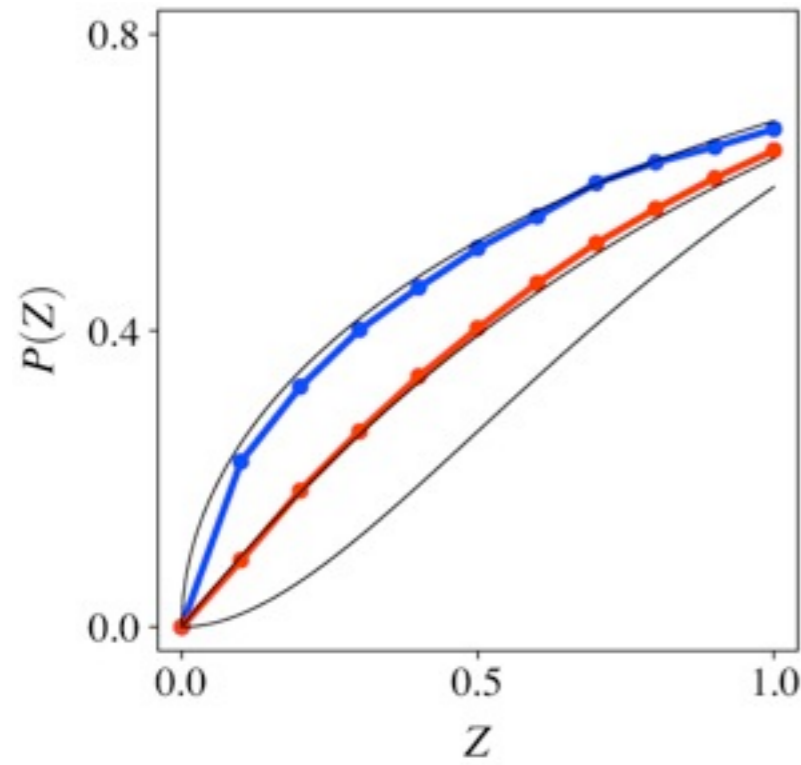
Two crystals whose lattices overlap (nearly) exactly: this can happen when the true symmetry is lower than the lattice symmetry

Measured intensities are the *sum* of two different reflections related by the twin operator, so a weak intensity is likely to be inflated by a stronger one

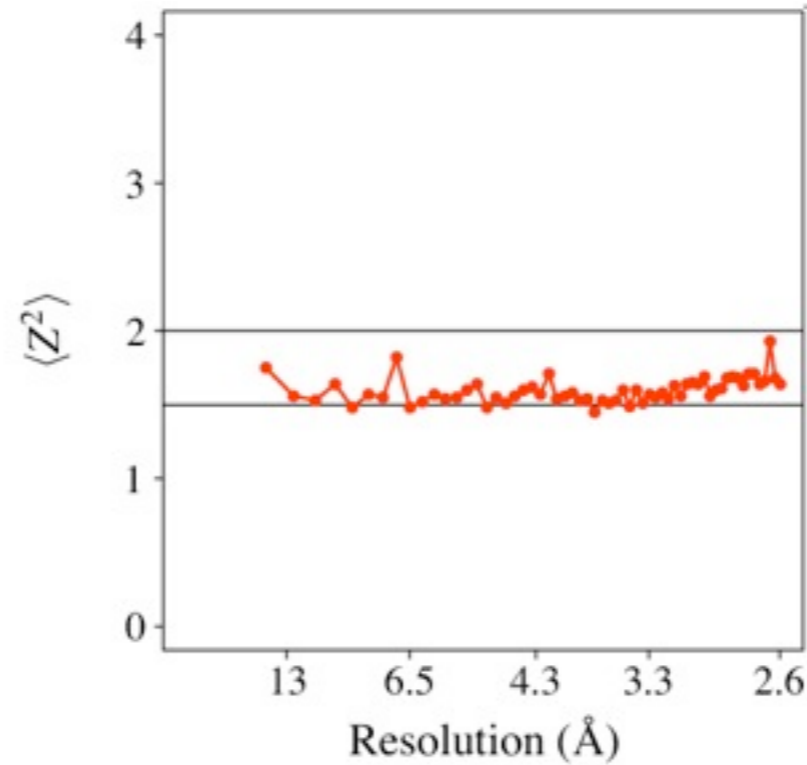
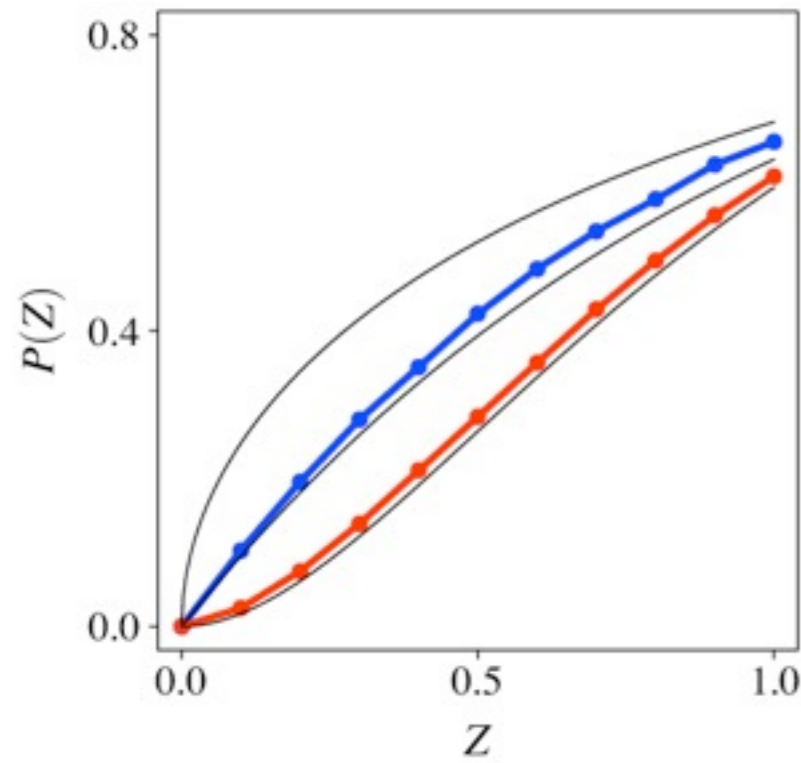
too few weak intensities



Examples



PDB entry 1ilj
single crystal



C-terminal domain of gp2
protein from phage SPP1
(unpublished)
perfect twin

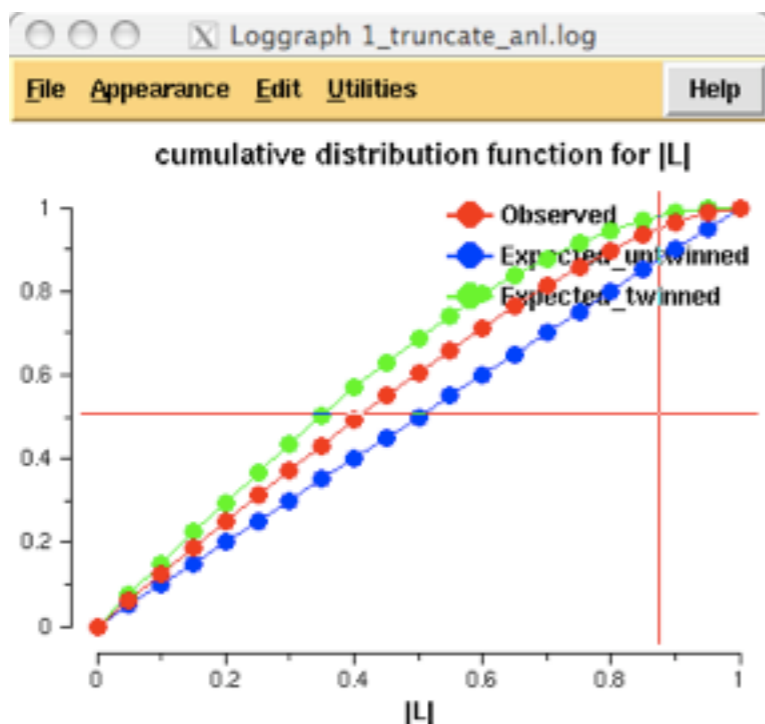
Cumulative intensity distribution

2nd moment of Z or $\langle E^4 \rangle$

Andrey Lebedev

Truncate: L- and H-tests

Cumulative distribution of L
(L-test)



0.873,0.509

Tables in File

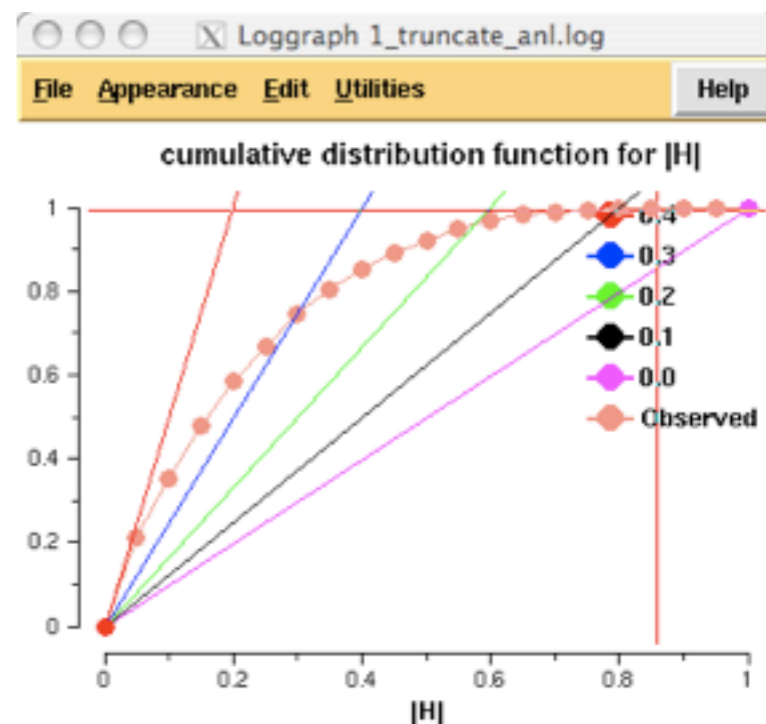
- Acentric moments of E using Truncate method
- H test for twinning (operator k, h, -l)
- L test for twinning**
- Wilson plot
- Truncate style Wilson plot

Graphs in Selected Table

L-test is probably the most reliable test for twinning

$L = (I_1 - I_2)/(I_1 + I_2)$
I1 & I2 close in reciprocal space

Cumulative distribution of H
(H-test)
(Partial twinning test)



0.857,0.995

Observed 0.85 , 0.999

Tables in File

- Acentric moments of E using Truncate method
- H test for twinning (operator k, h, -l)**
- L test for twinning
- Wilson plot
- Truncate style Wilson plot

Graphs in Selected Table

cumulative distribution function for |H|

$H = (I_1 - I_2)/(I_1 + I_2)$
I1 & I2 related by twin symmetry

Other features of the intensity distribution which may obscure or mimic twinning

Translational non-crystallographic symmetry:

whole classes of reflections may be weak

eg h odd with a NCS translation of $\sim 1/2, 0, 0$

$\langle I \rangle$ over all reflections is misleading, so Z values are inappropriate

The reflection classes should be separated (not yet done)

Anisotropy: $\langle I \rangle$ is misleading so Z values are wrong

ctruncate applies an anisotropic scaling before analysis

Overlapping spots: a strong reflection can inflate the value of a weak neighbour, leading to too few weak reflections

this mimics the effect of twinning

Summary

Questions & Decisions

- What is the point group (Laue group)?
- What is the space group?
- Is there radiation damage: should data be cut away from the end (possibly at the expense of resolution)?
- What is the best resolution cut-off?
- Is there anomalous signal (if you expect one)?
- Are the data twinned?
- Is this dataset better or worse than ones you have already?

Acknowledgements

Andrew Leslie	many discussions
Harry Powell	many discussions
Ralf Grosse-Kunstleve	cctbx
Kevin Cowtan	clipper, simplex minimiser, C++ advice
Martyn Winn & CCP4 gang	ccp4 libraries
Peter Briggs	ccp4i
Airlie McCoy	C++ advice, code etc
Graeme Winter	testing & bug finding
Clemens Vonrhein	testing & bug finding
Eleanor Dodson	many discussions
Andrey Lebedev	intensity statistics & twinning
Norman Stein	ctruncate
George Sheldrick	discussions on symmetry detection