# DIALS: integration

James Parkhurst

DIALS-6 workshop May 2015

DIALS
Diffraction Integration for Advanced Light Sources

*DIALS: integration*

# DATA FLOW

# Internals: top-level

**Tasks in dials.integate:**

Calculate the bounding box parameters from strong reflections
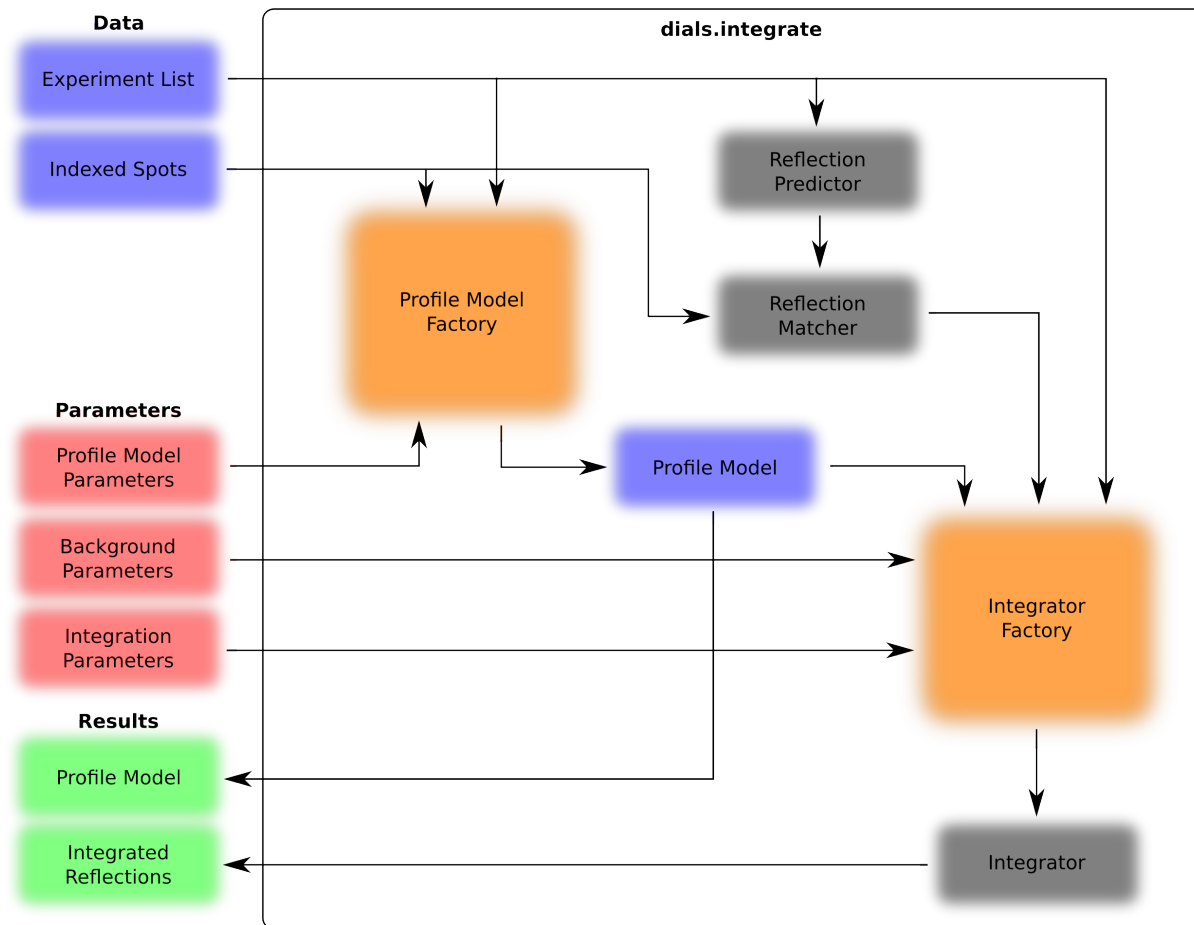
Predict the positions of reflections on the images

Build reference profiles across all images

Integrate the reflections and save output

# Internals: top-level

# Internals: types of integrator

**3D Integrator**
Each integration job reads a block of images and extracts reflections into 3D shoeboxes for processing.

**Flattened 3D Integrator**
Each integration job reads a block of images and extracts 3D shoeboxes which are "flattened" for processing.

**2D Integrator**
Each integration job reads a block of images and extracts partial reflections into 2D shoeboxes for processing.

**Single Frame 2D Integrator**
Each integration job reads a single image and extracts partials reflections into 2D shoeboxes for processing.
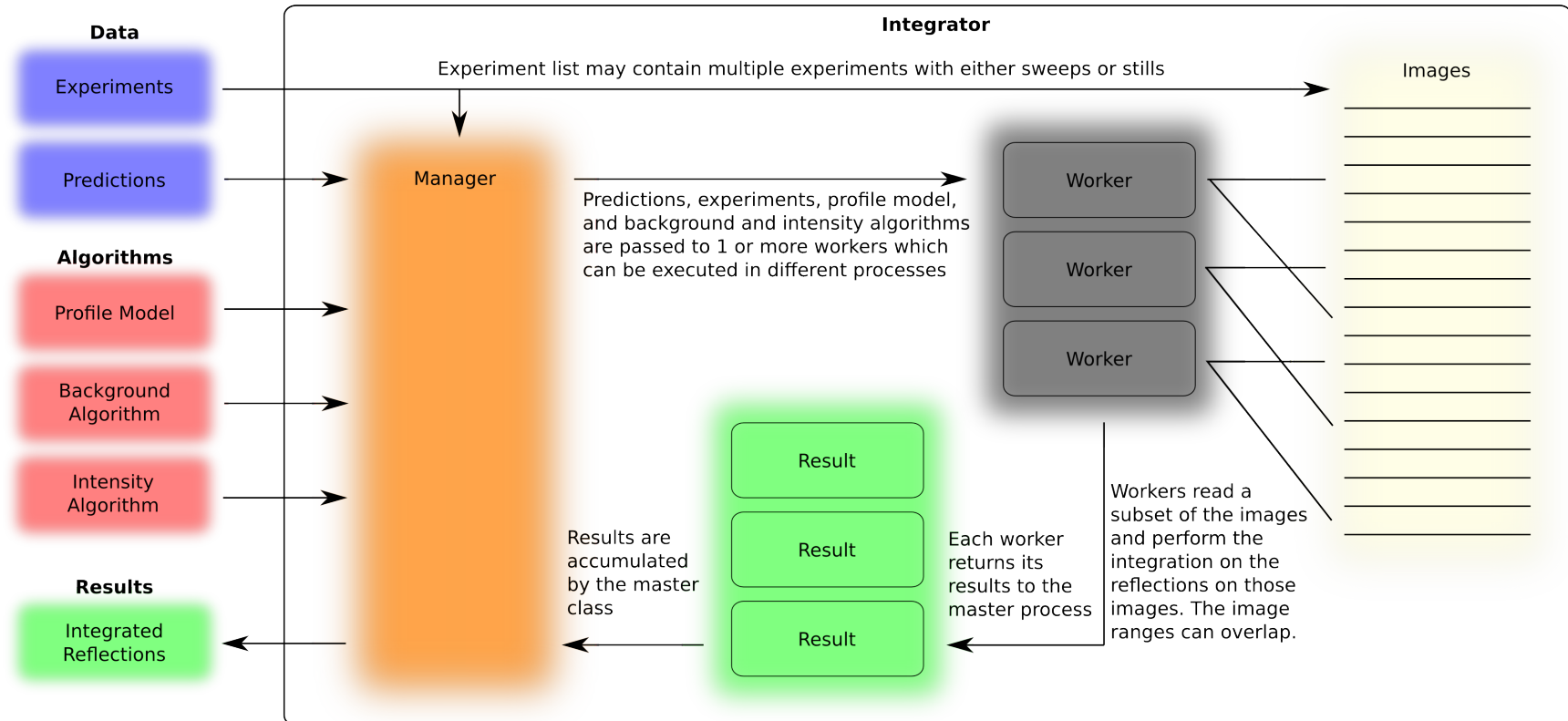
**Stills Integrator**
Same as the single frame 2d integrator but specialized to accept still experiments rather than rotation experiments.

Integrator

DIALS
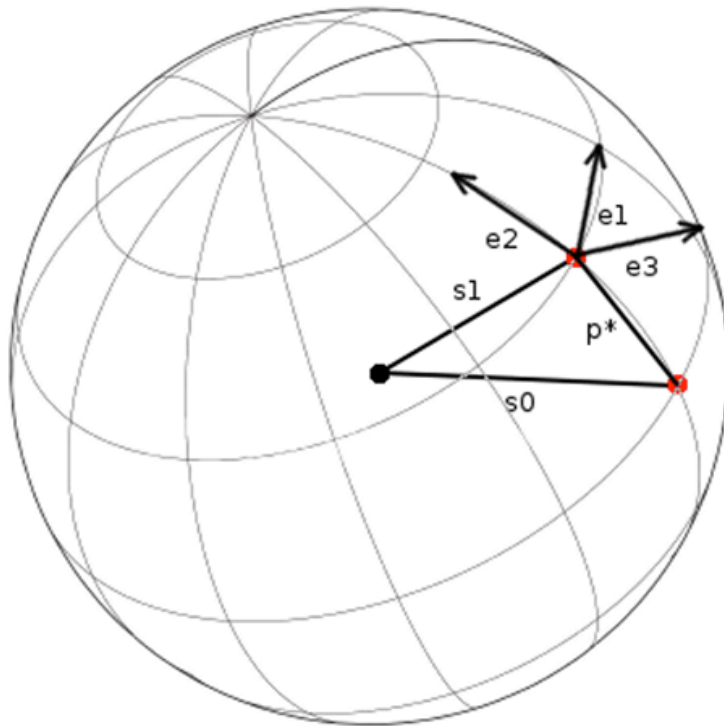Diffraction Integration for Advanced Light Sources

# Internals: integrator

**Data**

Experiments

Predictions

**Algorithms**

Profile Model

Background Algorithm

Intensity Algorithm

**Results**

Integrated Reflections

**Integrator**

Experiment list may contain multiple experiments with either sweeps or stills

Images

Manager

Predictions, experiments, profile model, and background and intensity algorithms are passed to 1 or more workers which can be executed in different processes

Worker

Worker

Worker

Result

Result

Result

Results are accumulated by the master class

Each worker returns its results to the master process

Workers read a subset of the images and perform the integration on the reflections on those images. The image ranges can overlap.

**DIALS**
Diffraction Integration for Advanced Light Sources

*DIALS: integration*

# REFLECTION SHOEBOXES

# Computing reflection shoeboxes



**Profile coordinate system**

Use the kabsch model of a normal distribution on the surface of the Ewald sphere

$$\exp(-\epsilon{\downarrow}1\,{\uparrow}2\,/2\sigma{\downarrow}D{\uparrow}2\,)\,\exp(-\epsilon{\downarrow}2\,{\uparrow}2\,/2c$$

$$e_1 = S_1 \times S_0/|S_1 \times S_0|$$
$$e_2 = S_1 \times e_1/|S_1 \times e_1|$$
$$e_3 = (S_1 + S_0)/|S_1 + S_0|$$

# Computing reflection shoeboxes



$\sigma_D$ is calculated from the spread of angles between the predicted diffracted beam vector and the vector for each strong pixel in the spot

$\sigma_M$ is calculated by maximum likelihood method assuming a normal distribution of phi residuals for each strong pixel in the spot

*DIALS: integration*

# BACKGROUND MODELLING

# Models

- Options to model the background under the peak as either
  - A constant across each image
  - A constant across all images
  - A plane across each image
  - A hyper-plane across all images
- Computed using simple linear least squares

# Outliers

- Large valued outliers can cause the background to be over-estimated
- This then causes the reflection intensity to be under-estimated
- Outliers in the background can come from:
  - Intensity from neighbouring spots
  - Hot pixels
  - Zingers
  - Unpredicted reflections
  - Ice rings
  - etc

# Outliers

# Simple outlier rejection



outlier.algorithm=nsigma

Reject pixels N sigma from the mean
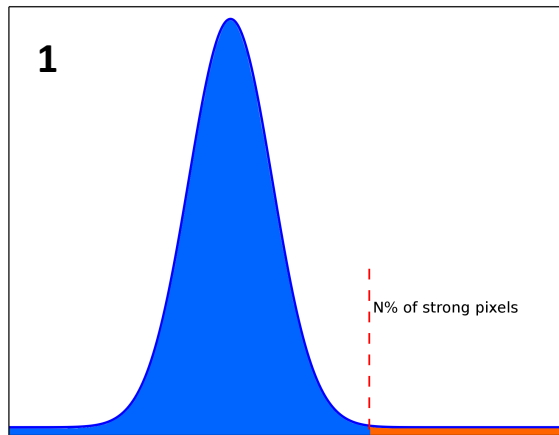


outlier.algorithm=truncated

Reject N% of the highest and lowest valued pixels
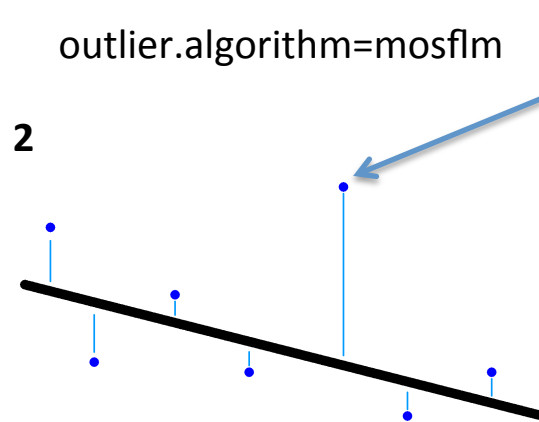


outlier.algorithm=tukey
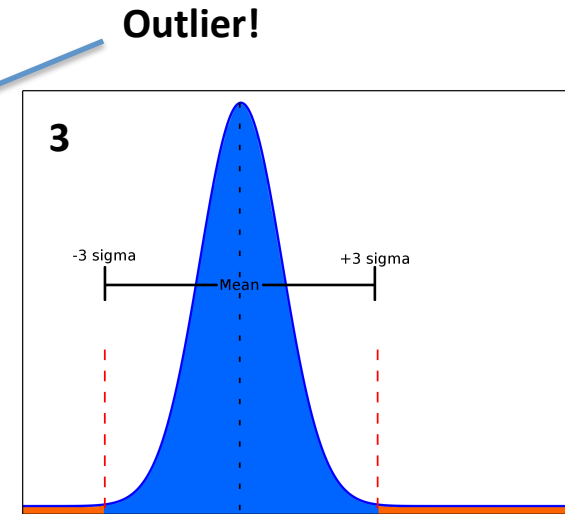
Reject pixels based on the interquartile range

# Mosflm-style outlier rejection

outlier.algorithm=mosflm

**Outlier!**

1

N% of strong pixels

2

3

-3 sigma    +3 sigma
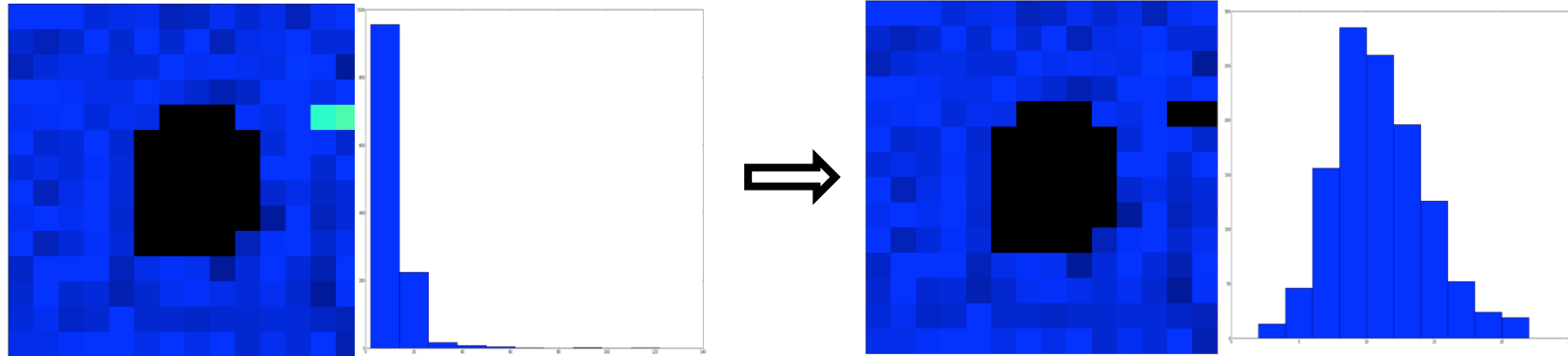
Mean

Remove N% of strongest pixels and compute the background plane

Compute the residuals of all background pixels to the plane

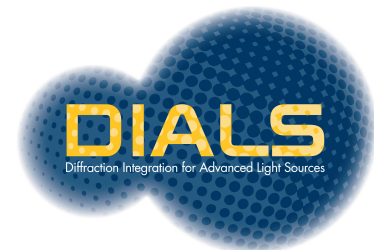Remove pixels whose residuals are greater than N sigma from the plane

**DIALS**
Diffraction Integration for Advanced Light Sources
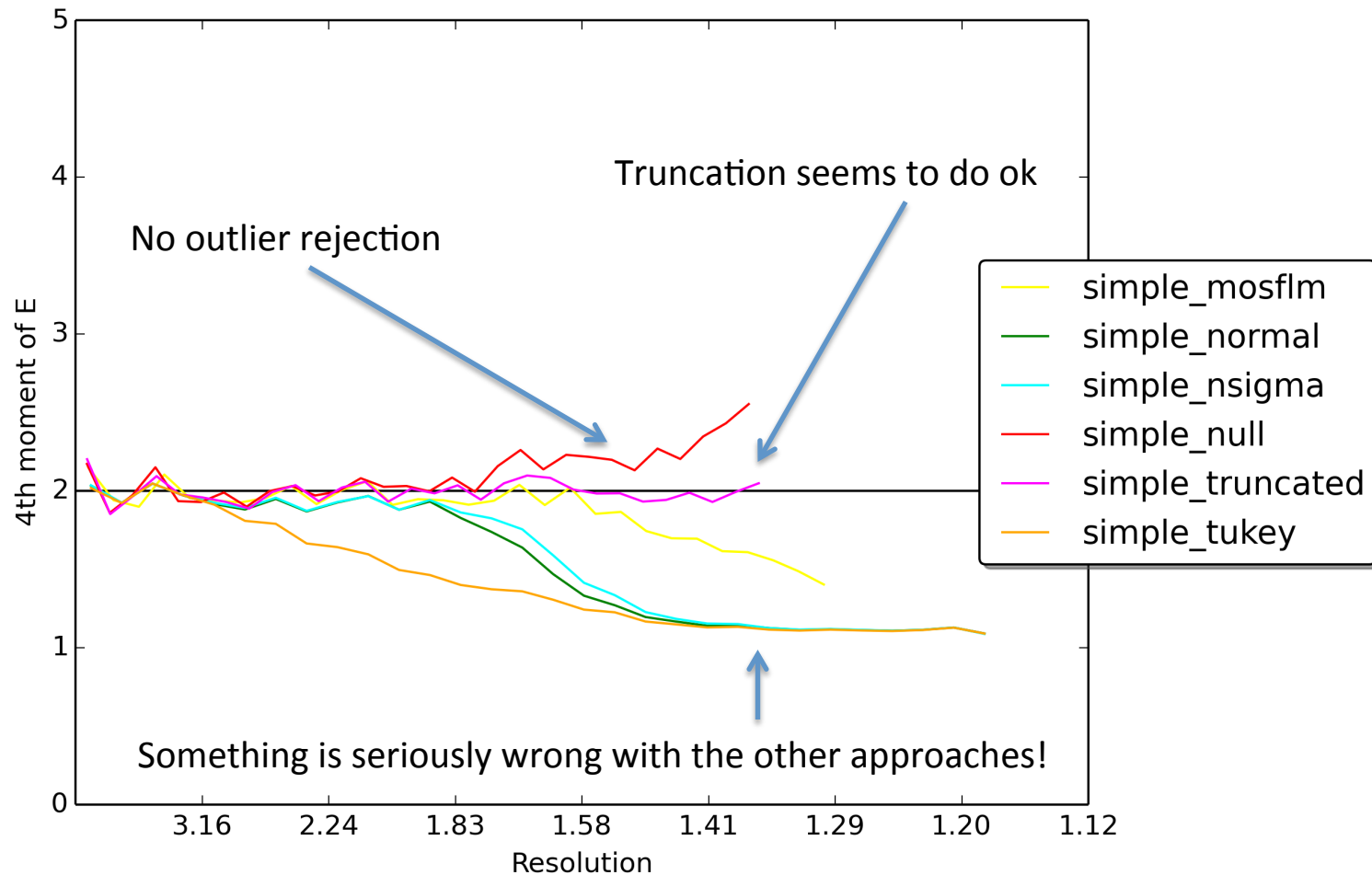
# XDS-style outlier rejection



Iteratively remove high valued pixels until the distribution of pixel counts resembles a normal distribution

# What effect does outlier rejection have

- Looked at two datasets
    - I04 Bag training. Good data with very few outliers.
    - PNAS data. Good data with some serious outliers.
    These outliers caused pointless to find the wrong point group when the data was processed without outlier rejection (pointless has now been fixed so this error no longer occurs).
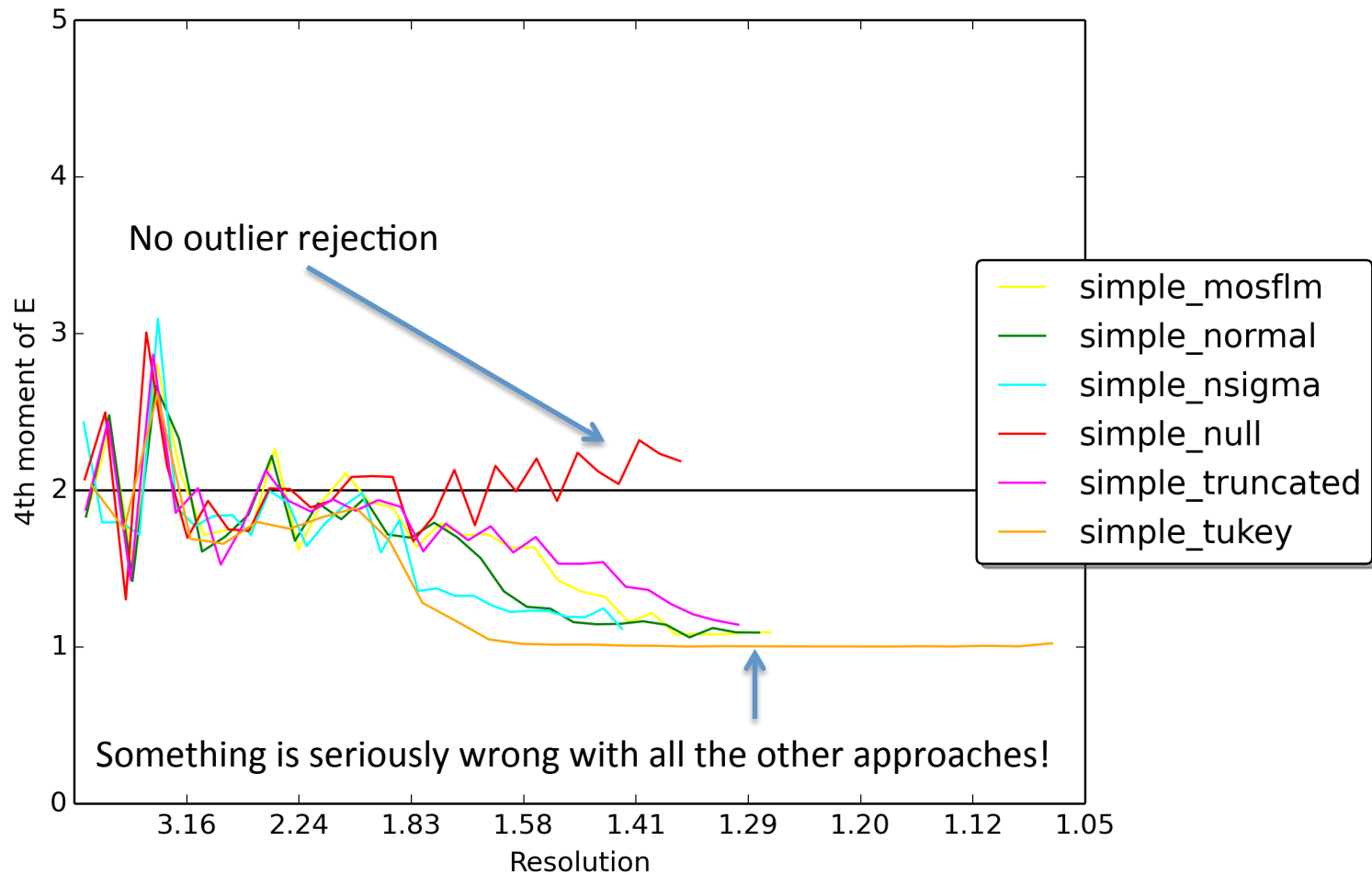
DIALS
Diffraction Integration for Advanced Light Sources

# Handling outliers badly



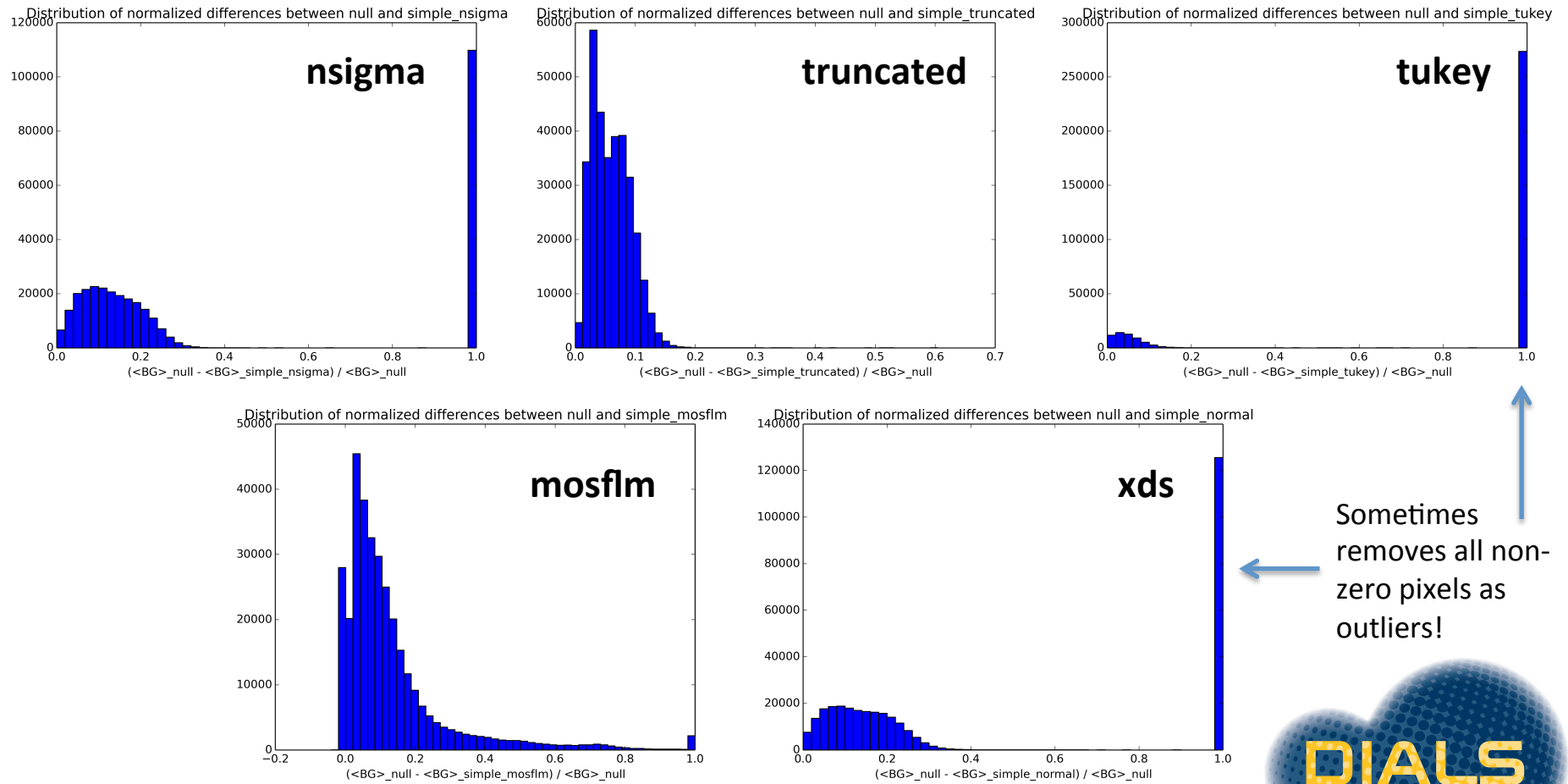I04 bag training data processed with xia2

# Handling outliers badly



PNAS data processed with xia2

# Handling outliers badly

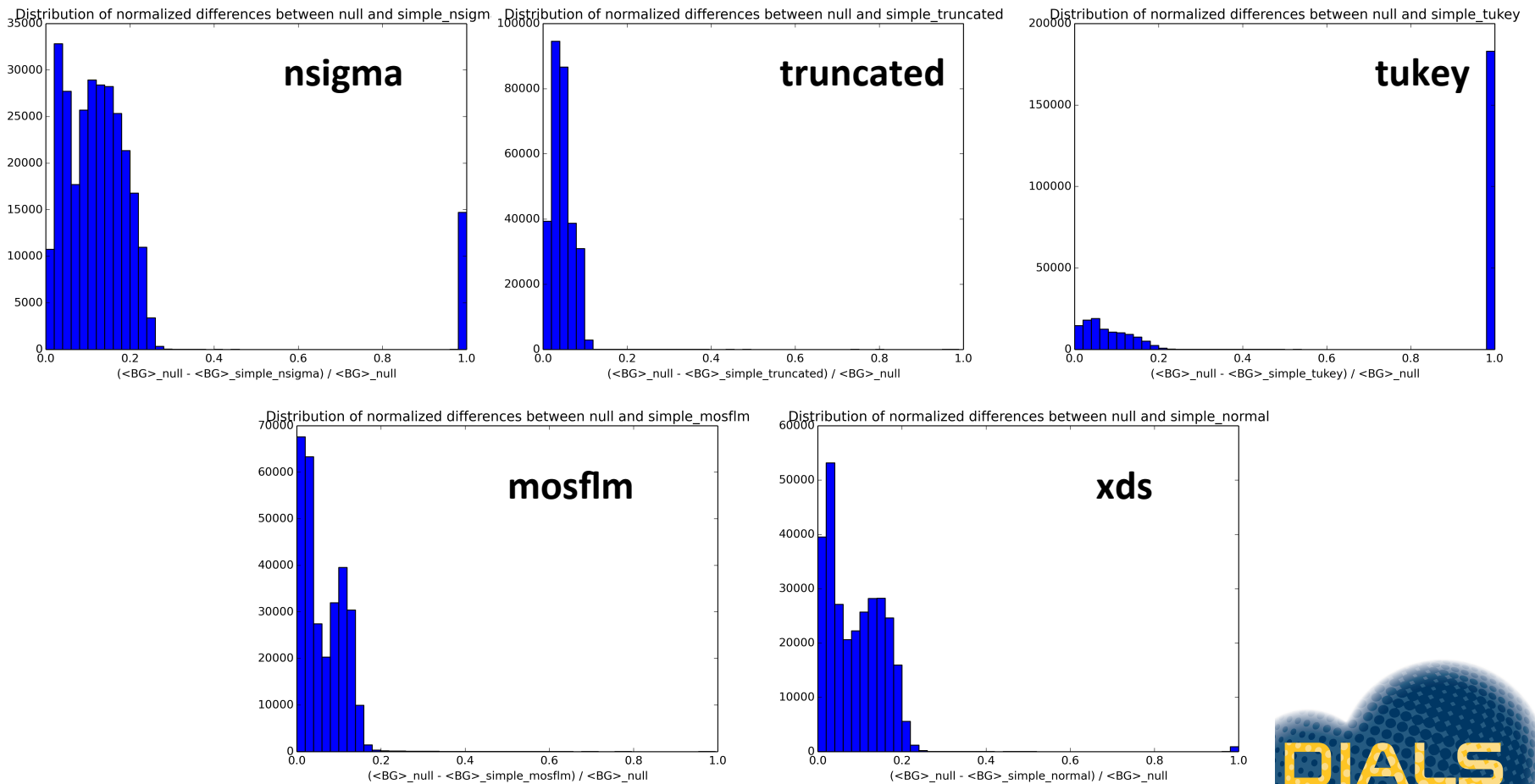Histograms of background differences vs no outlier rejection for I04 bag training data



Sometimes removes all non-zero pixels as outliers!

**Background is systematically under-estimated**

# Handling outliers badly

Histograms of background differences vs no outlier rejection for PNAS data
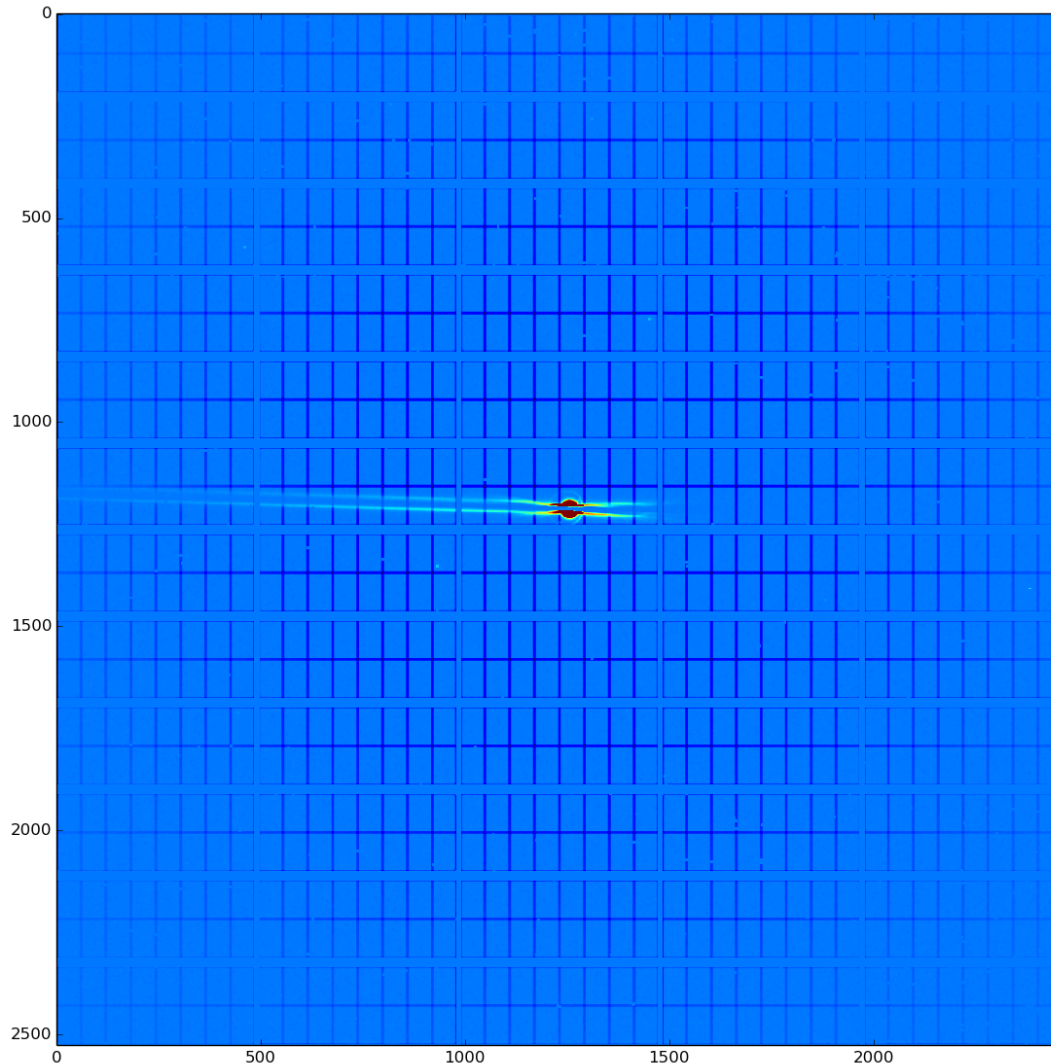


**Background is systematically under-estimated**

# Handling outliers better

- Most of our methods assume a normal distribution of counts – not a good approximation for data with low background

- Instead of rejecting outliers could we use a robust estimation method?

- Most methods also focus on normally distributed data – leads to under-estimation

- Could we use a robust generalized linear model approach assuming a Poisson distribution?

DIALS
Diffraction Integration for Advanced Light Sources
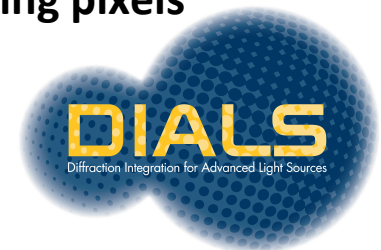
# Is the background Poisson distributed?



- Analysed 9000 blank images
- Local index of dispersion computed at each pixel
- Average index of dispersion computed for each pixel

**Background data is Poisson distributed**

**Virtual pixels show under-dispersion due to correlations with neighbouring pixels**

**~7.2% of pixels are affected**

# Robust GLM algorithm

Eva Cantoni and Elvezio Ronchetti (2001), "*Robust Inference for Generalized Linear Models*", Journal of the American Statistical Association, Vol. 96, No. 455

Solve $\sum i=1\uparrow n\blacksquare[\psi\downarrow c\ (r\downarrow i\ )w(x\downarrow i\ )\mu\downarrow i\ /\sqrt{V(\mu\downarrow i)}\ -a(\beta)]=0$

$r\downarrow i=\ y\downarrow i\ -\mu\downarrow i\ /\sqrt{V(\mu\downarrow i)}$  Pearson residuals

$V(\mu)=\mu$  Variance function

$w(x)=1$  Weights for explanatory variables

$\psi\downarrow c\ (r)=\{\blacksquare r,\qquad |r|\leq cc\ sign(r),\ |r|>c$  Weights for dependant variables

$c=1.345$  Tuning constant

$a(\beta)=\ 1/n\ \sum i=1\uparrow n\blacksquare E[\psi\downarrow c\ (r\downarrow i\ )]\ \ w(x\downarrow i\ )\mu\downarrow i\ /\sqrt{V(\mu\downarrow i)}$  Consistency correction
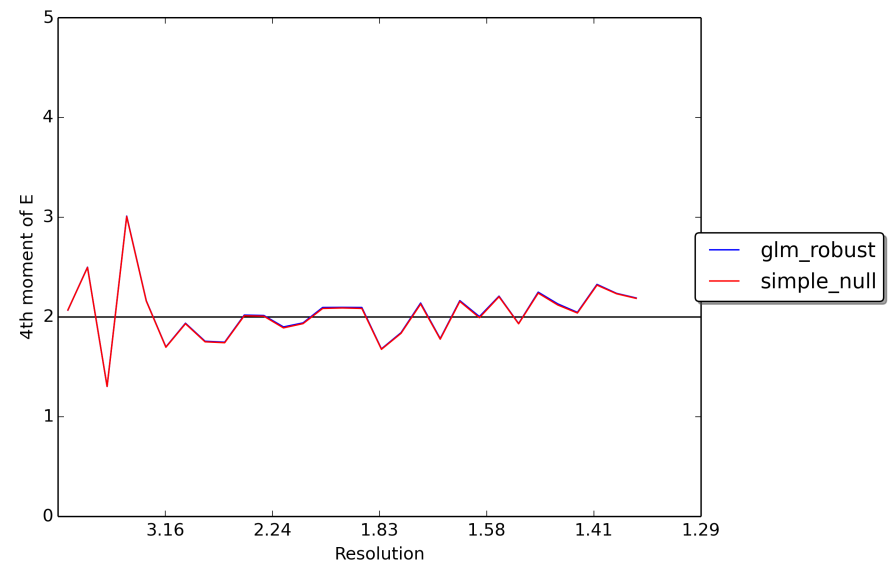
# Handling outliers better (?)

In both cases robust GLM method gives sensible results for the 4$^{th}$ moment of E plots

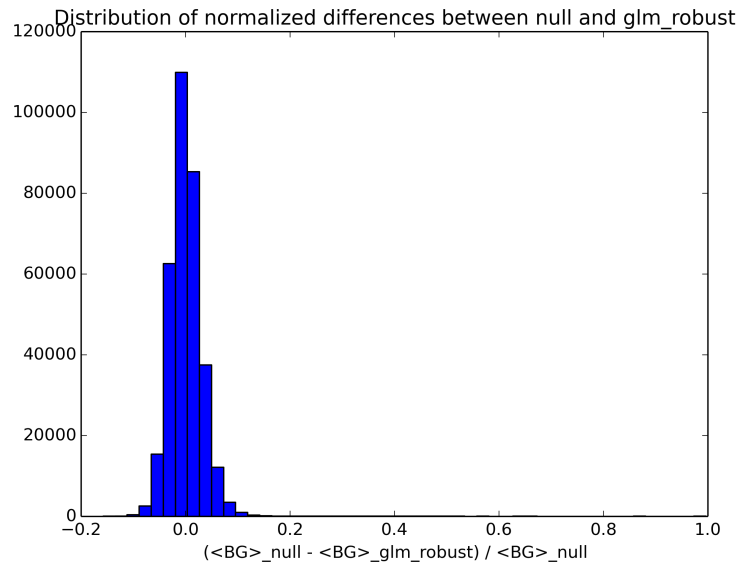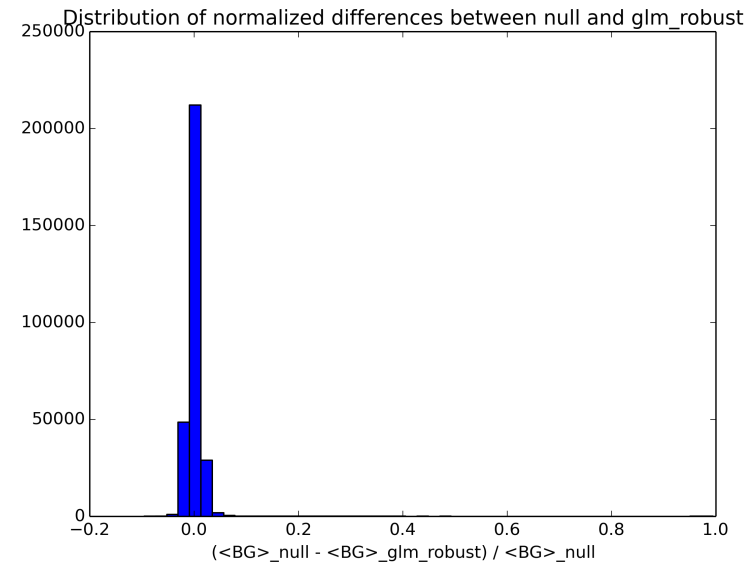

**I04 bag training data**

**PNAS data**

# Handling outliers better (?)

Histograms of background differences vs no outlier rejection for I04 bag training data
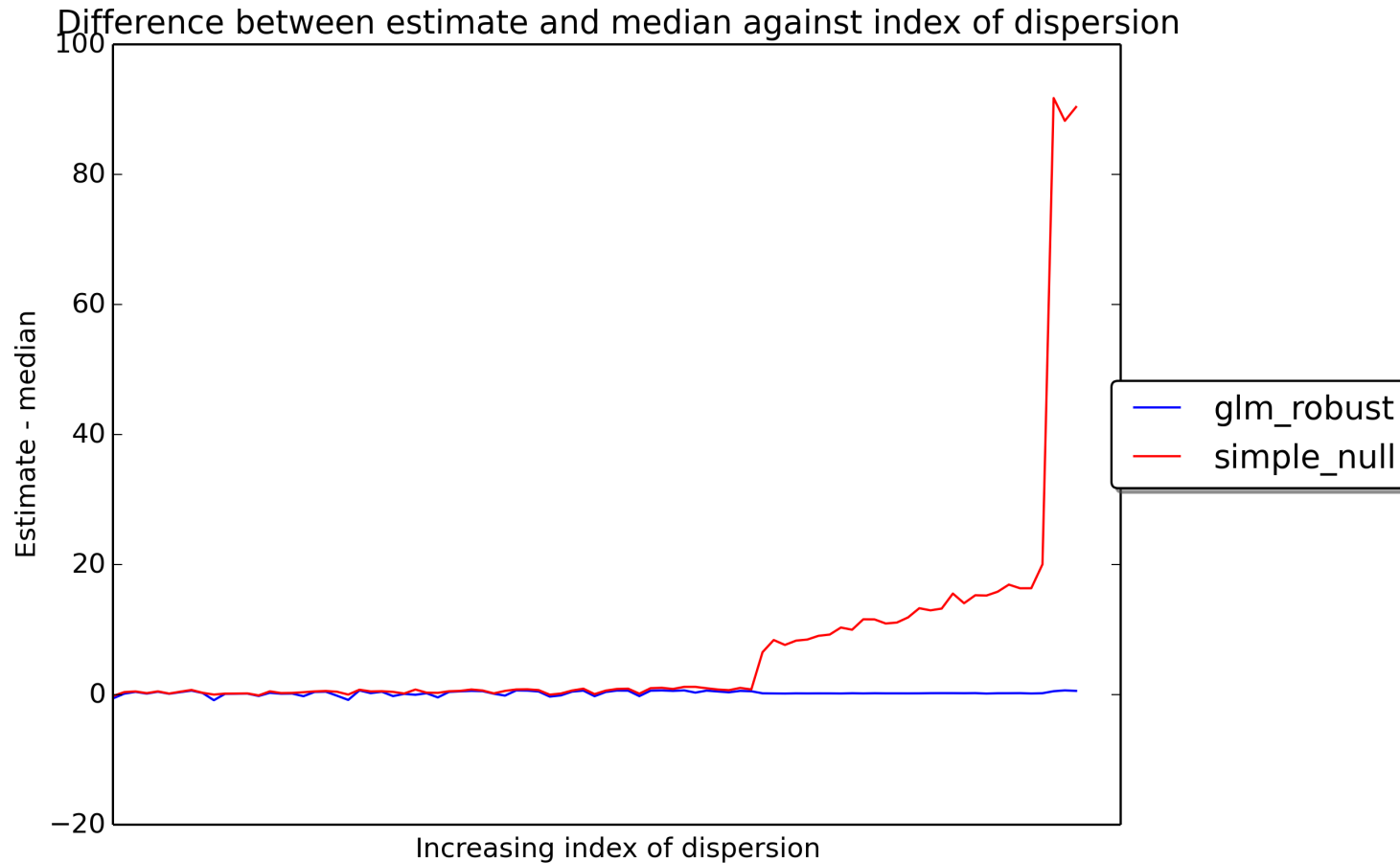
Histograms of background differences vs no outlier rejection for PNAS data



**No systematic difference in the background!**

# Are we actually doing anything?



Difference between estimate and median against index of dispersion

**Looks like we're handling outliers ok**

# Robust algorithm

- Algorithm requires a good seed value or it won't converge – use the median

- Can't represent a straight line so currently using constant background – plan to look at more sophisticated background if needed

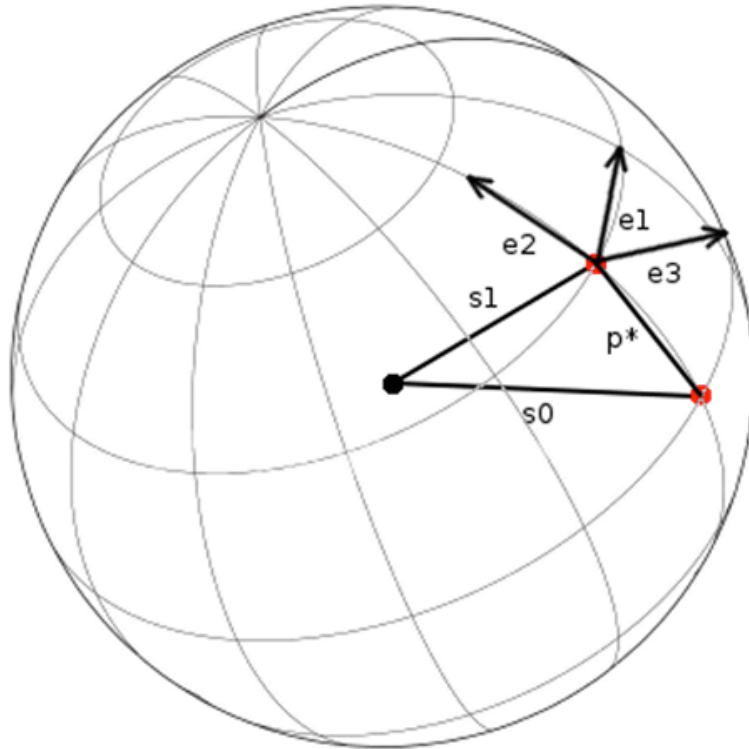- Current results look promising

*DIALS: integration*

# SIGNAL INTEGRATION

# Integration

- Integration algorithm options:
  - Summation
  - 3D profile fitting (as in XDS)
  - 2D profile fitting (future)

# 3D profile fitting coordinate system



**Profile coordinate system**

Use Kabsch coordinate system
- Corrects for geometrical distortions
- Makes spots appear to have taken shortest path through Ewald sphere
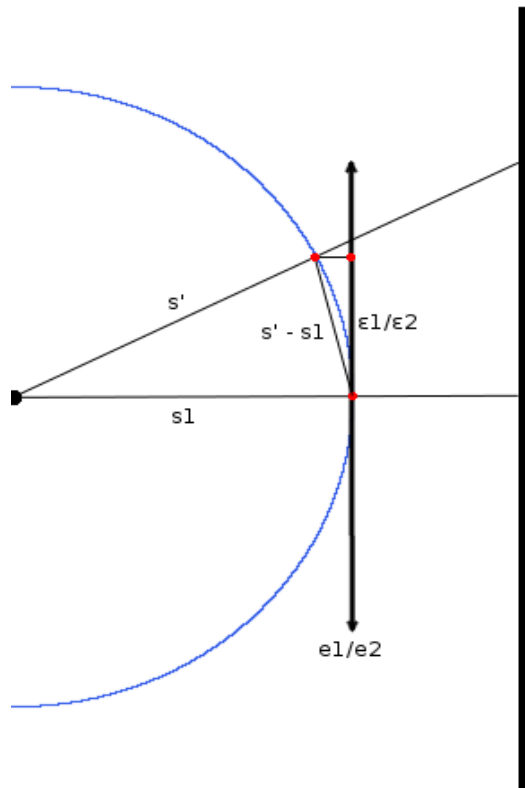- Model assumes a Gaussian profile in Kabsch coordinate system

$$e_1 = S_1 \times S_0 / |S_1 \times S_0|$$
$$e_2 = S_1 \times e_1 / |S_1 \times e_1|$$
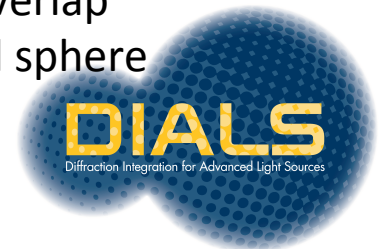$$e_3 = (S_1 + S_0) / |S_1 + S_0|$$

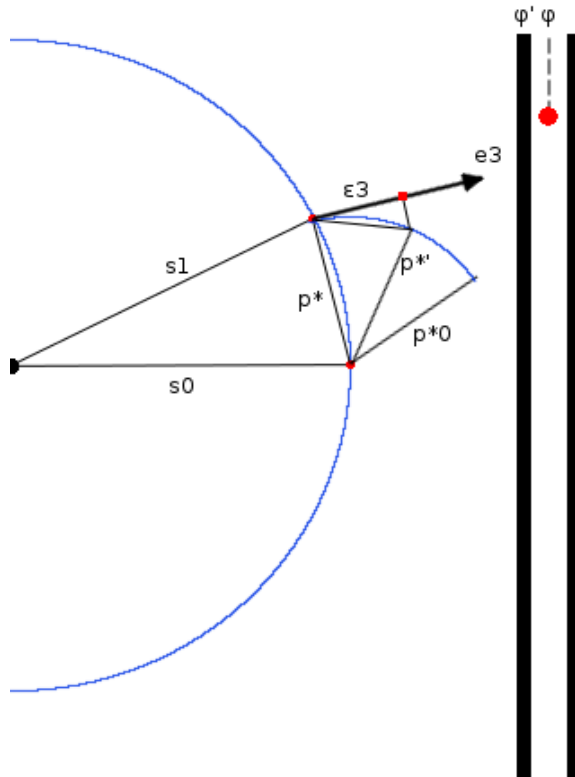# 3D profile fitting pixel gridding



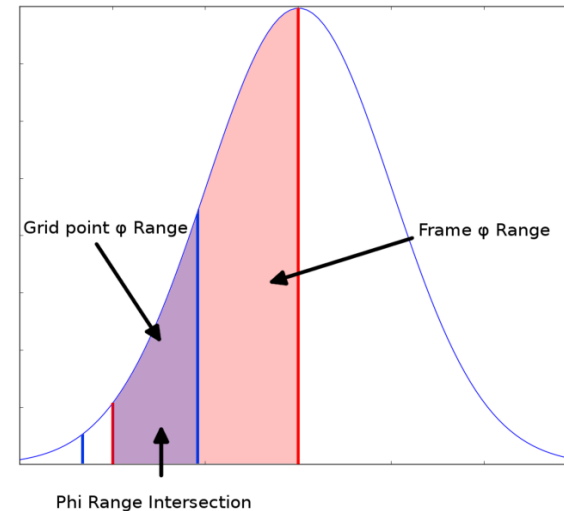Pixels are mapped to the Ewald sphere



Counts are redistributed to Ewald sphere grid by computing fractional overlap of each pixel and Ewald sphere grid point

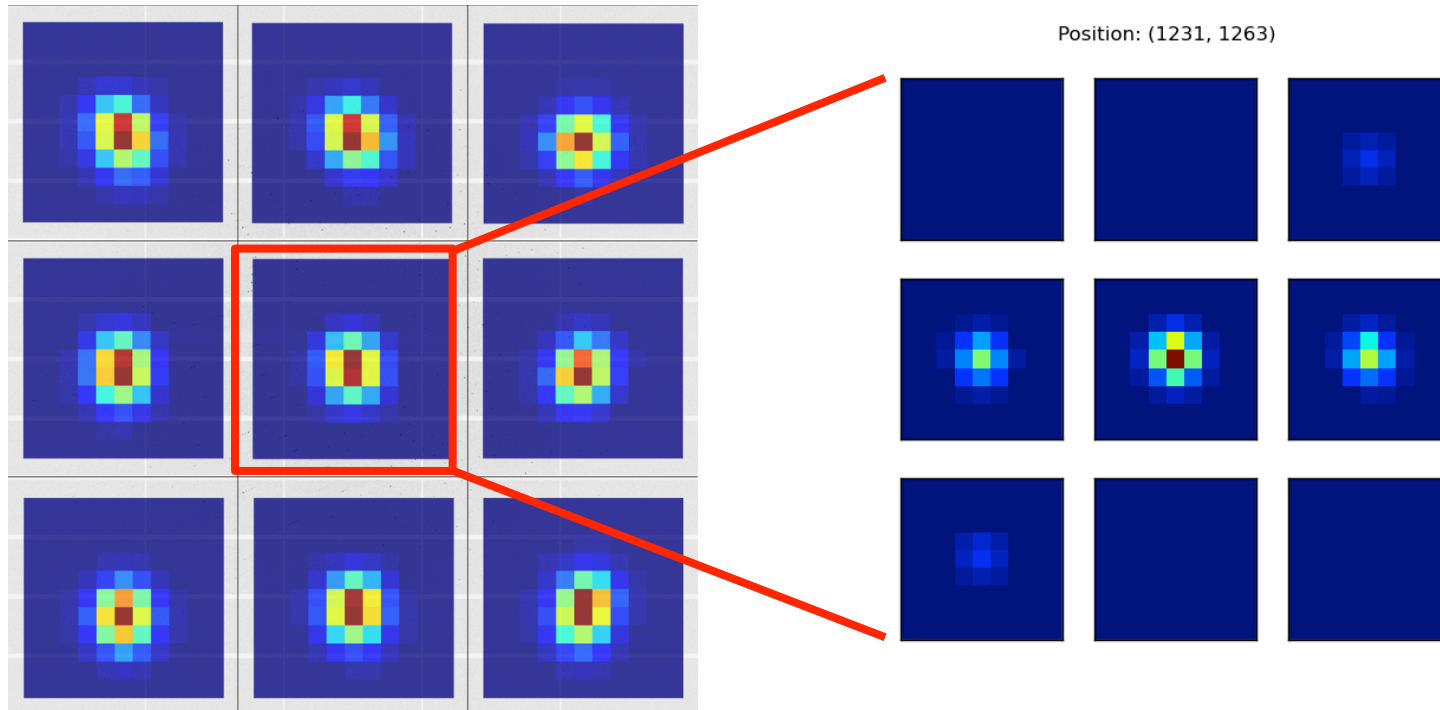# 3D profile fitting phi gridding



Frames are transformed to make reflection appear as if it took the shortest path through the Ewald sphere

Counts on each image are distributed by finding the angular overlap between each grid point and each image and integrating over the intersection

# Building reference profiles

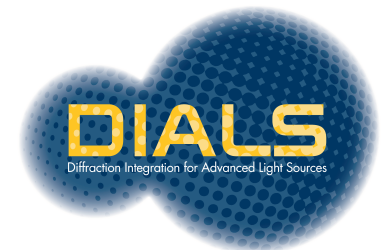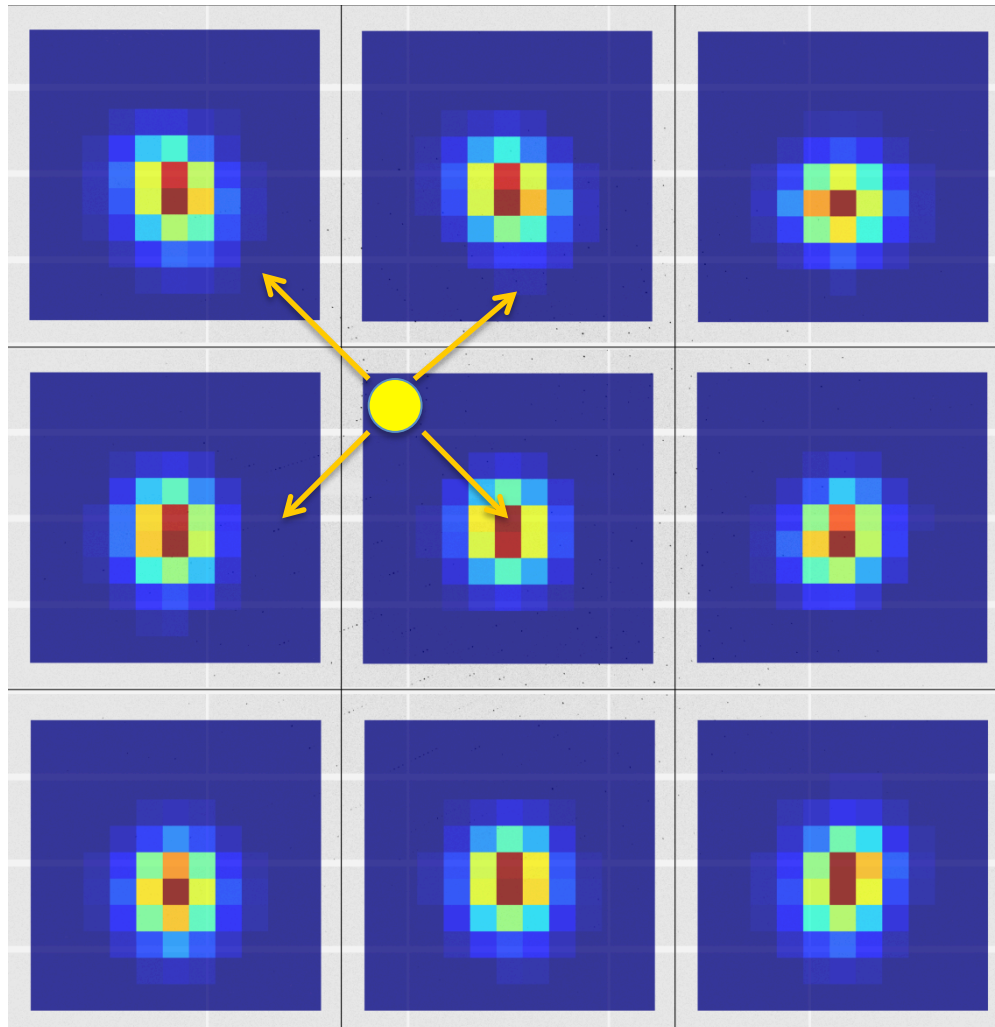

Position: (1231, 1263)

- Reference profiles are formed on a grid covering a given angular range
- Grid options include:
  - Rectangular grid (as in Mosflm)
  - Circular grid (as in XDS)
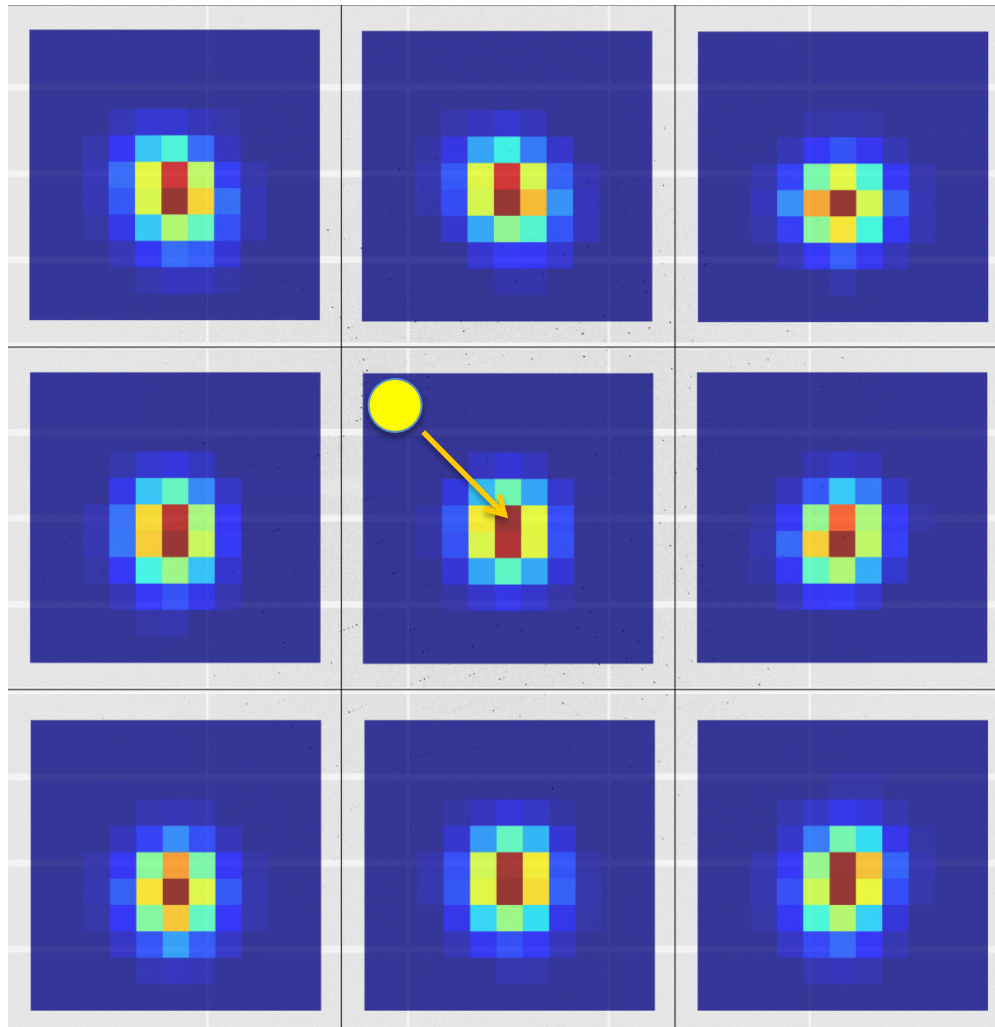  - Single reflection (currently for multi-panel detectors)

DIALS
Diffraction Integration for Advanced Light Sources

# Building reference profiles

Each strong spot contributes to building the profile at adjacent grid points

# Fitting reference profiles

Each reflection is fitted against its closest reference profile

# Thanks!