

Structure solution from weak anomalous data

Diffraction Methods in Structural Biology

Gordon Research Conference

Bates College, Maine

July 30, 2014

Gábor Bunkóczi, Airlie McCoy, Randy Read (Cambridge University)
Nat Echols, Ralf Grosse-Kunstleve, Paul Adams, James Holton (Lawrence
Berkeley National Laboratory)
Tom Terwilliger (Los Alamos National Laboratory)



Structure solution from weak anomalous data

Problems with weak signal

Quantifying the anomalous signal

Solving the anomalous sub-structure with weak signal

Solving structures with weak signal

Estimating the anomalous signal from the data

Weak anomalous signal

Reasons:

*Few anomalous scatterers, sulfur SAD, weak diffraction,
wavelength far from peak*

Consequences:

Substructure identification is difficult

Phasing is poor

*Iterative density modification, model-building and refinement
works poorly*

Quantifying the anomalous signal I

CC_{ano} : How accurate are the anomalous differences?

Anomalous differences measured with errors ϵ_j

$$\Delta_{ano,j}^{obs} = \Delta_{ano,j} + \epsilon_j$$

Correlation of observed and true anomalous differences

$$CC_{ano} \equiv \frac{\langle \Delta_{ano,j} \Delta_{ano,j}^{obs} \rangle}{\langle \Delta_{ano}^2 \rangle^{1/2} \langle \Delta_{ano}^{2,obs} \rangle^{1/2}}$$

Fraction of observed anomalous differences that is noise

$$E_{ano}^2 = \frac{\langle \sigma_{ano}^2 \rangle}{\langle \Delta_{ano}^{2,obs} \rangle}$$

Expected value of CC_{ano}

$$CC_{ano} \sim [1 - E_{ano}^2]^{1/2}$$

Anomalous signal S_{ano} : How accurate are maps based on the anomalous differences?

Anomalous difference Fourier with model phases

$$\rho(x) = \frac{1}{V} \sum_h \Delta_{ano,h}^{obs} e^{i(\varphi_h^c - \frac{\pi}{2})} e^{-2\pi i(h \cdot x)}$$

Peak height at coordinates of anomalously-scattering atoms

$$S_{ano} \equiv \frac{\langle \rho(x_k) \rangle}{\langle \rho^2 \rangle^{1/2}}$$

Expected value of signal S_{ano}

$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

Example of anomalous signal S_{ano}
Holton Challenge data

bl831.als.lbl.gov/~jamesh/challenge/anom/

Simulated diffraction data from 3dk0 to 1.8 Å
(useful to 2.3 Å)

0% to 100% occupancy of Se in selenomethionine

“Impossible.mtz” has fraction Se of 0.21

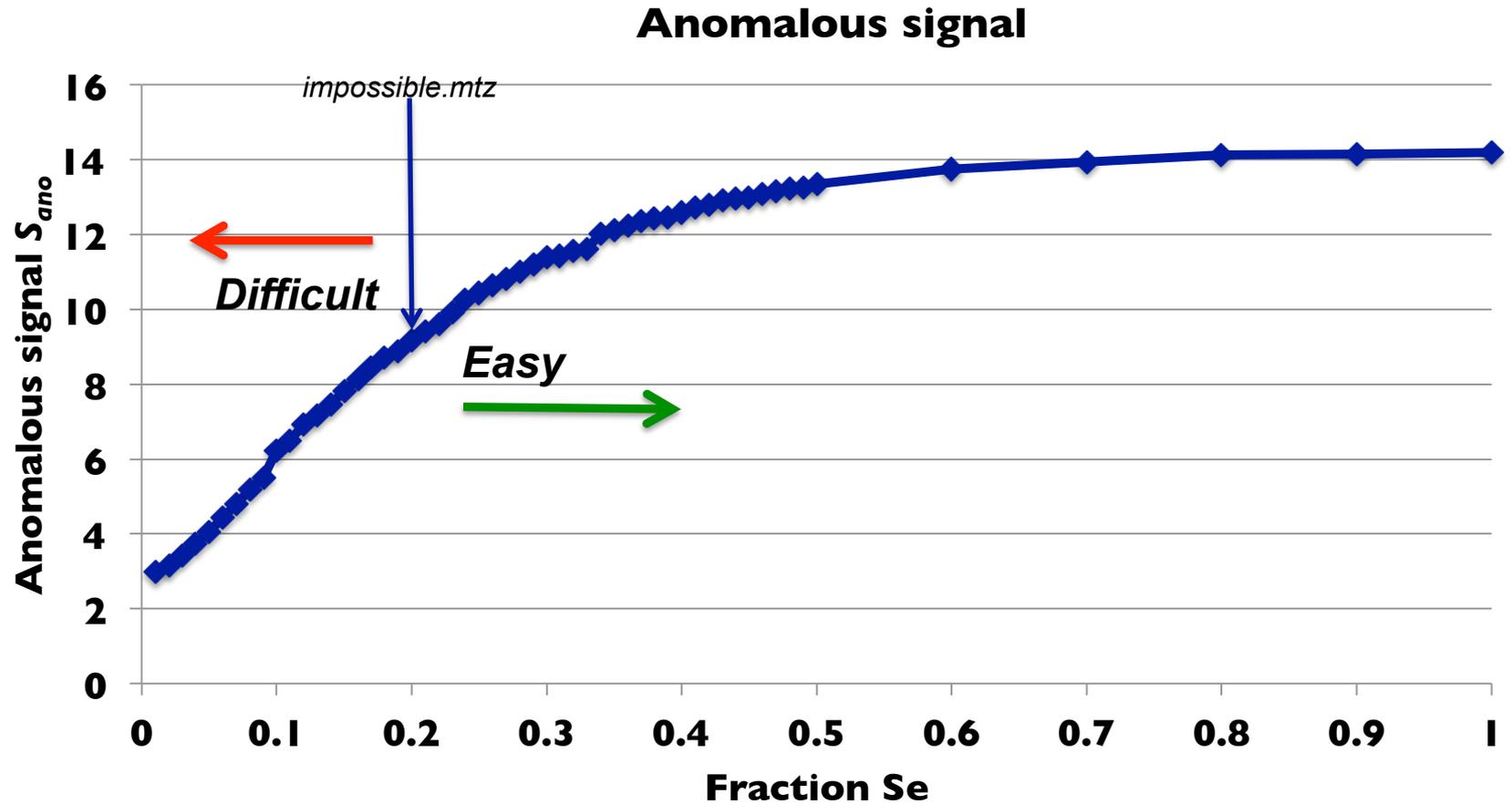
21% SeMet incorporation



22% SeMet incorporation



Example of anomalous signal S_{ano} Holton Challenge data



Finding the anomalous sub-structure with weak anomalous signal

Current approaches

***Dual-space methods (Shelxd, HySS, Crunch2)
Difference Fourier (Solve)***

Limitation of these approaches

***Anomalous differences are only approximately
proportional to the structure factors for anomalously-
scattering atoms***

Finding the anomalous sub-structure with weak anomalous signal

Most powerful source of information about substructure before phases are known is the SAD likelihood function:

The likelihood of measuring the observed anomalous data given a partial model

Using the SAD likelihood function to find the anomalous sub-structure

Start with guess about the anomalous sub-structure

From anomalous difference Patterson

Random

Any other source

Find additional sites that increase the likelihood

*LLG completion based on log-likelihood gradient maps**

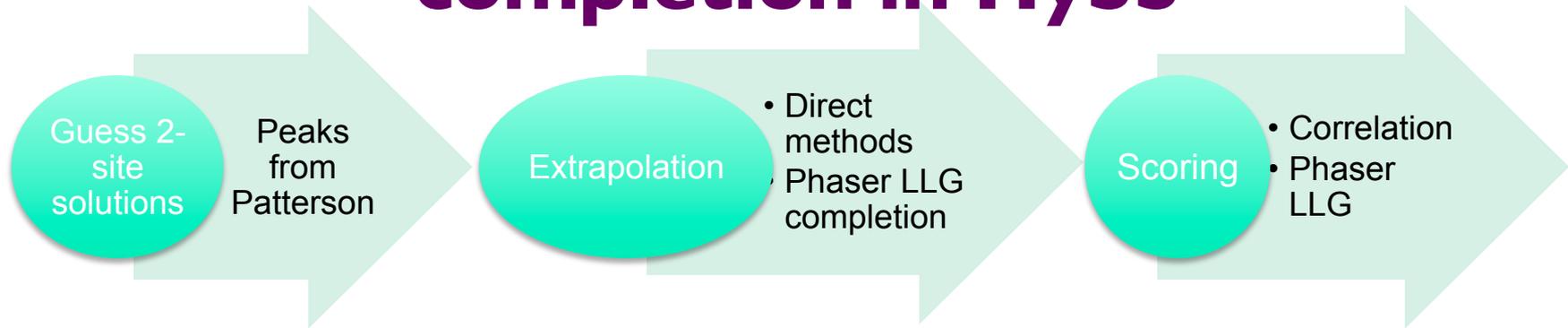
Iterative addition of sites

Related to using a difference Fourier—but much better

*La Fortelle, E. de & Bricogne, G. (1997). *Methods Enzymol.* 276, 472-494
McCoy, A. J. & Read, R. J. (2010). *Acta Cryst.* D66, 458-469.



Using LLG completion and dual-space completion in HySS



Guess 2-site solutions

Peaks from Patterson

Extrapolation

- Direct methods
- Phaser LLG completion

Scoring

- Correlation
- Phaser LLG

- Range of resolution
- Variable number of Patterson solutions

Adjustable
LLGC_SIGMA
(cut-off for peak height)

Use LLG score to compare solutions

Terminate early if same solution found several times

Run quick direct methods first

Using LLG completion in HySS

Test cases

164 SAD datasets from PDB (largely JCSG MAD data)

Using peak, remotes, inflection as available to include data
with low anomalous signal

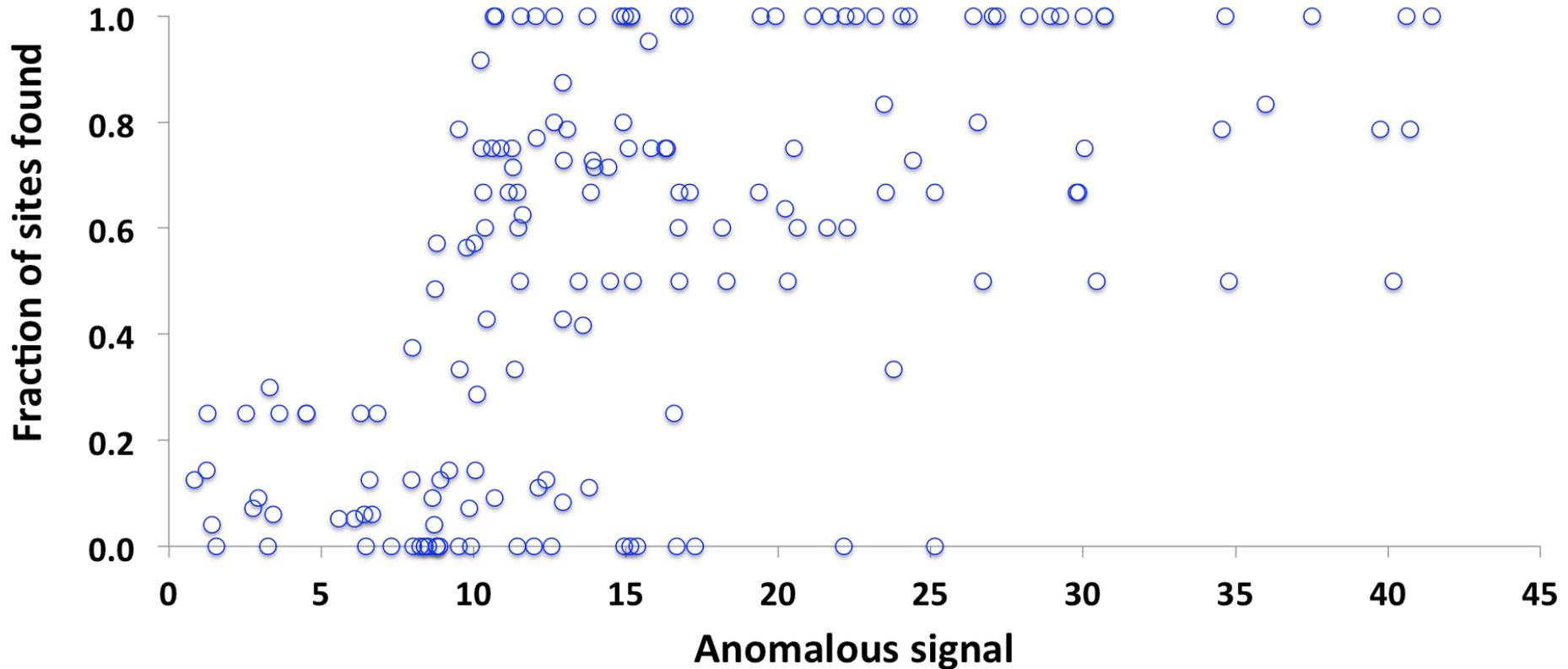
Setting up test data on 165 datasets

- *phenix.fetch_pdb 2o7t*
- *phenix.python \$PHENIX/phenix/phenix/autosol/sad_data_from_pdb.py 2o7t*
- Splits out each wavelength (peak, edge, remote etc) for MAD and run separately
- Run HySS with dual-space methods or LLG completion

Direct methods vs LLG completion
164 SAD datasets from PDB



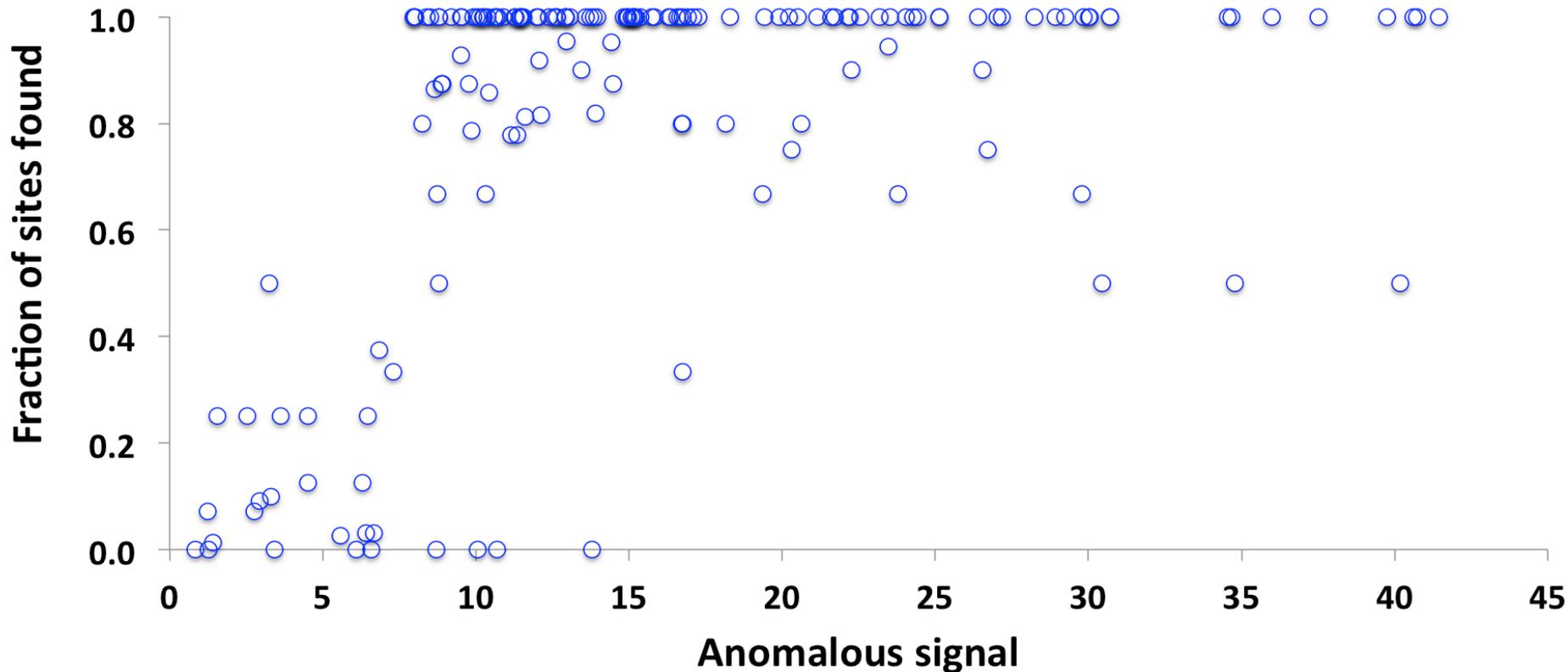
HySS Direct Methods



Direct methods vs LLG completion
164 SAD datasets from PDB

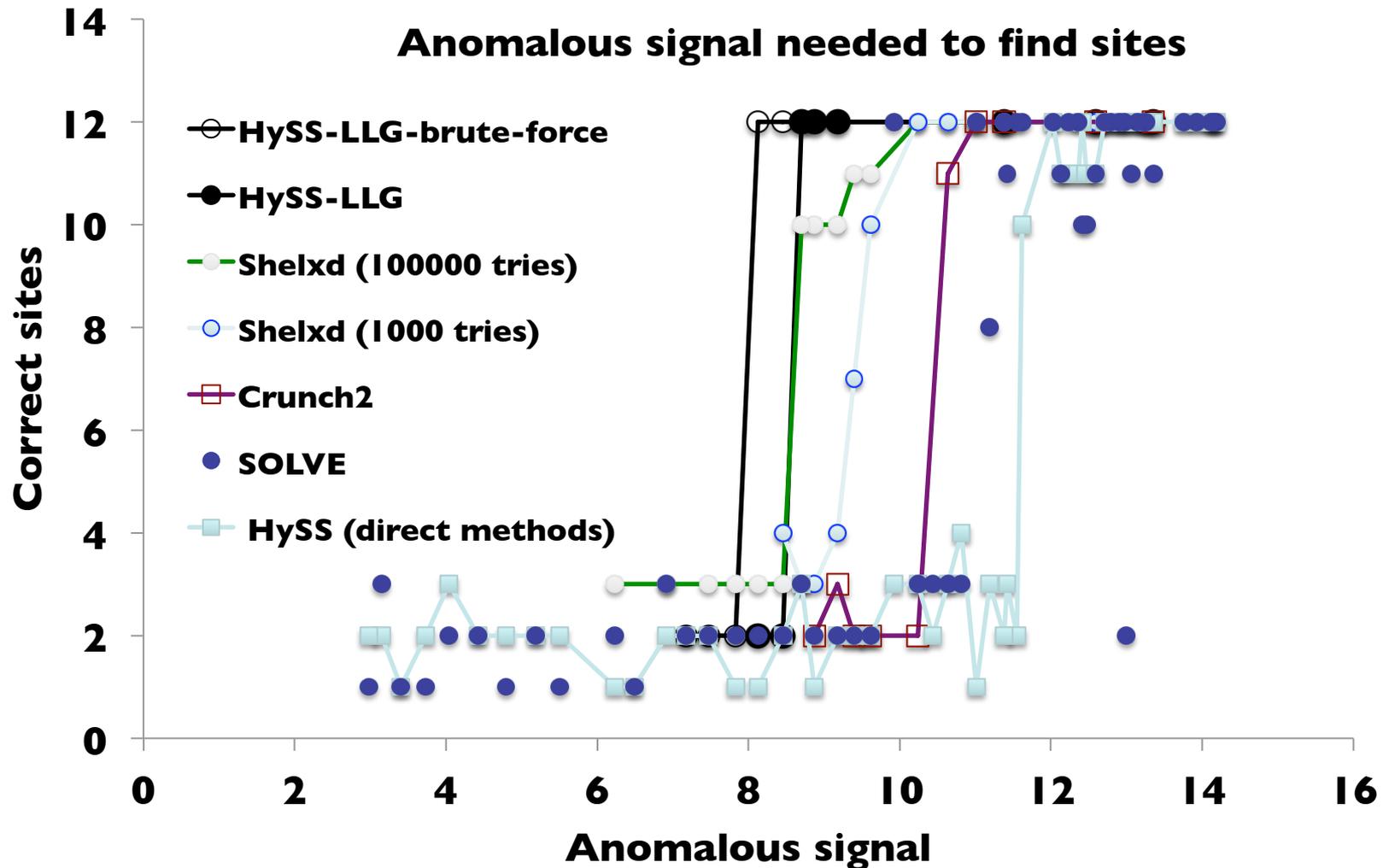


HySS LLG Completion



Holton Challenge data

Correct sites found vs anomalous signal S_{ano}



CysZ multi-crystal sulfur-SAD data

Qun Liu, Tassadite Dahmane, Zhen Zhang, Zahra Assur, Julia Brasch, Lawrence Shapiro, Filippo Mancina, Wayne Hendrickson (2012). Science 336,1033-1037

Data from 7 crystals collected at 1.74 Å

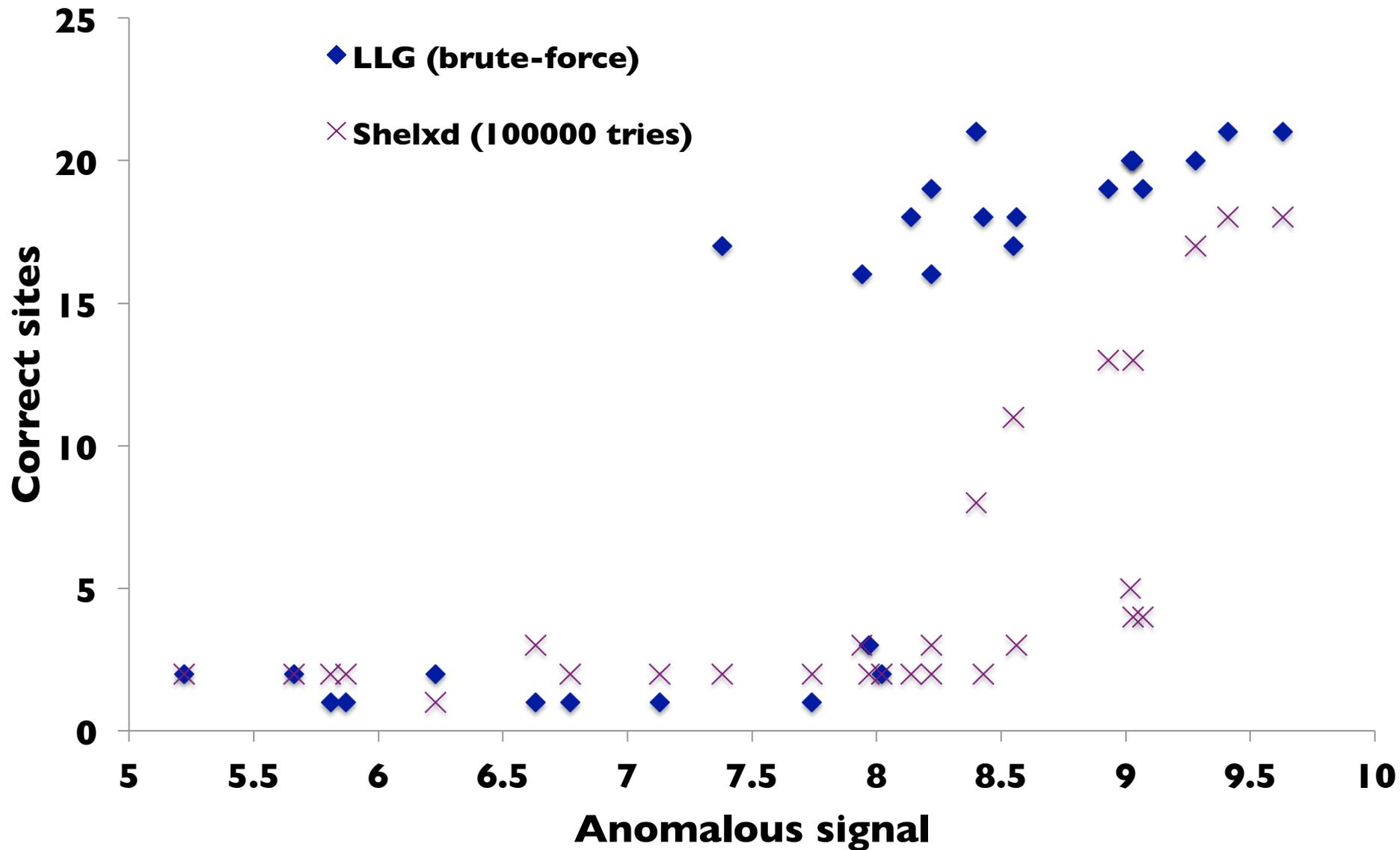
Only merged data could be solved

What is the minimum number of crystals that could have been used?

CysZ multi-crystal sulfur-SAD data

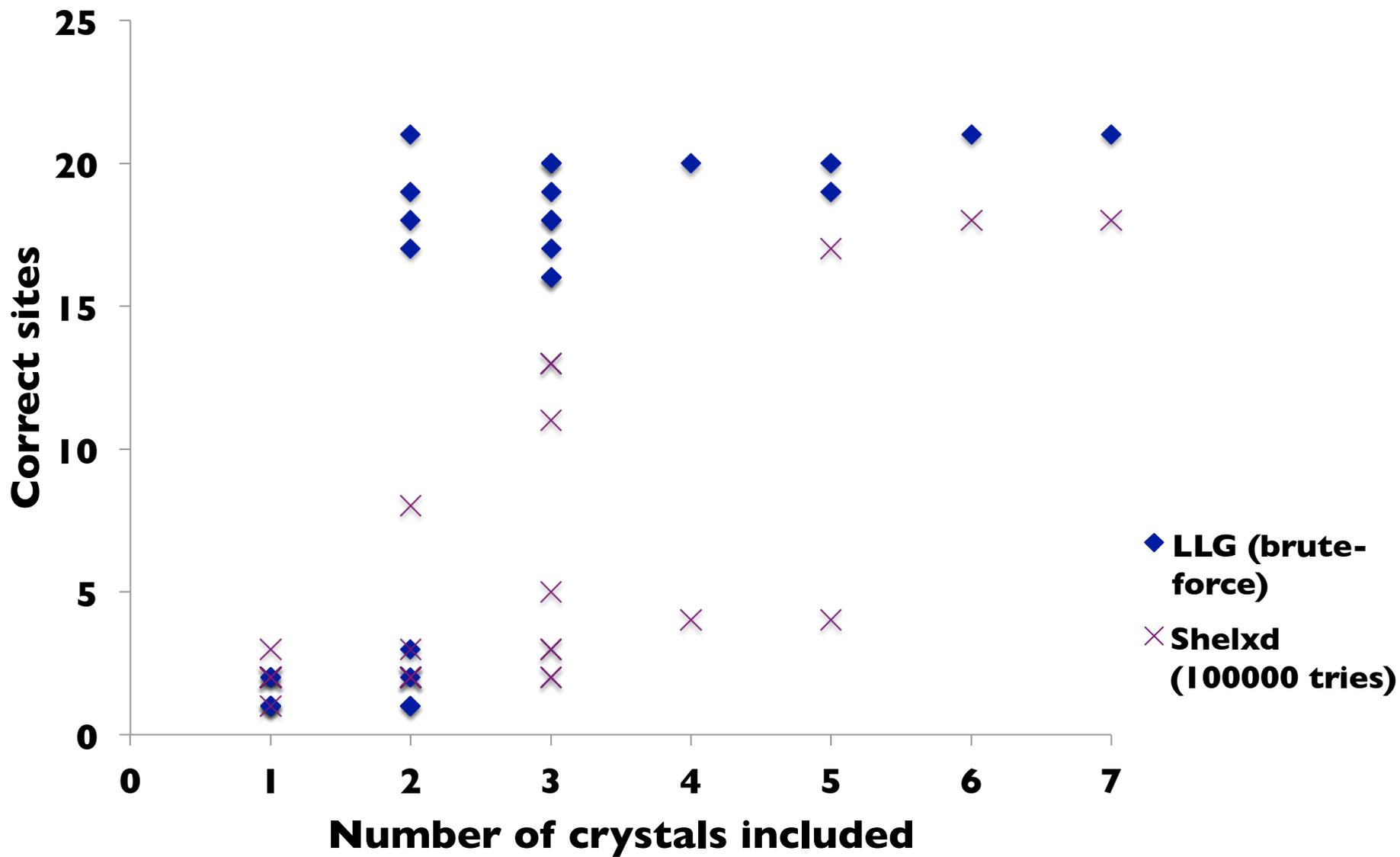
Datasets	Anomalous signal
5	5.22
1	5.66
4	5.81
2	5.87
6	6.23
7	6.63
3	6.77
56	7.13
561	7.94
67	8.22
273	9.02
2734	9.03
27345	9.07
27346	9.28
273456	9.41
2734561	9.63

CysZ multi-crystal sulfur-SAD data



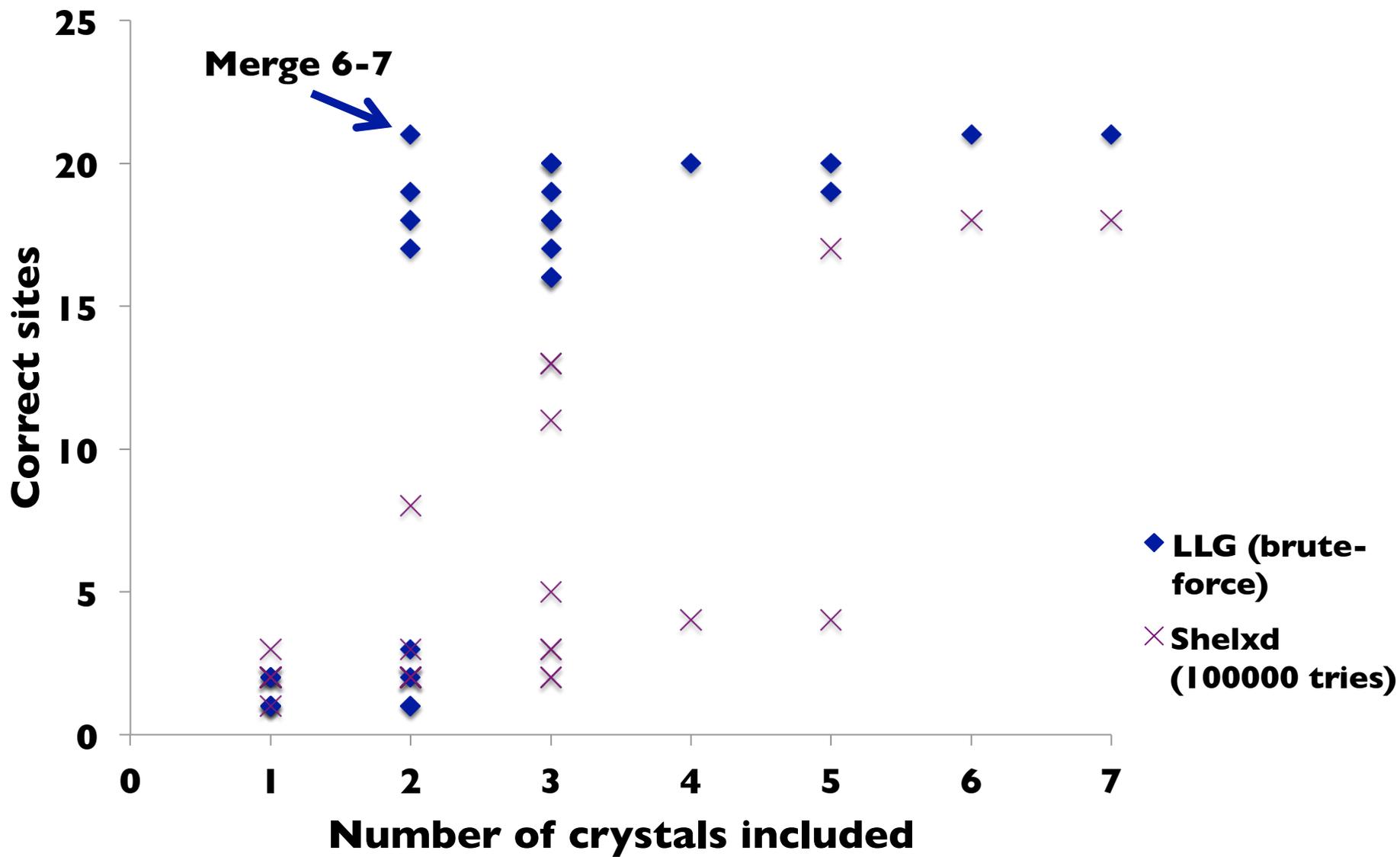


CysZ multi-crystal sulfur-SAD data





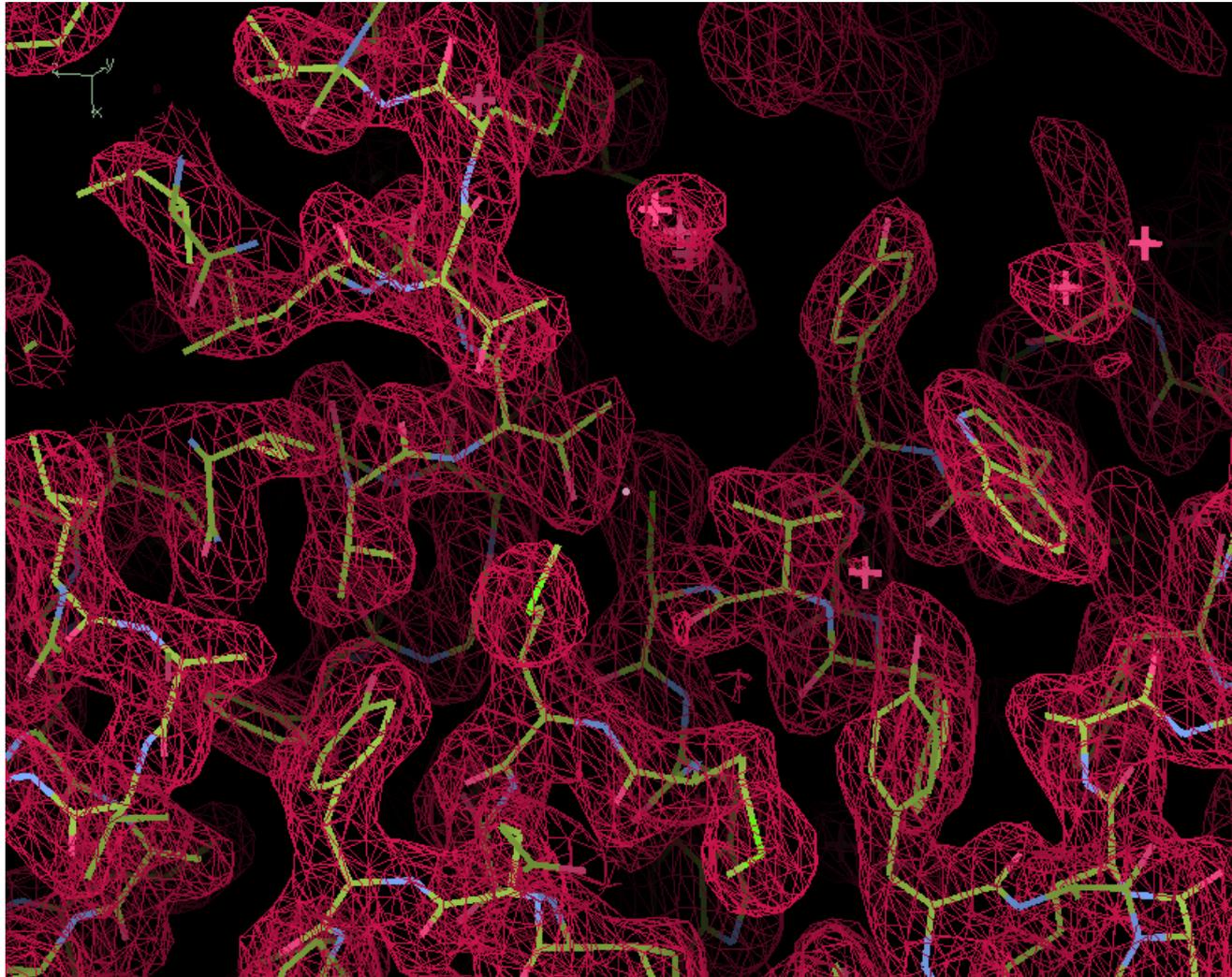
CysZ multi-crystal sulfur-SAD data



CysZ multi-crystal sulfur-SAD data

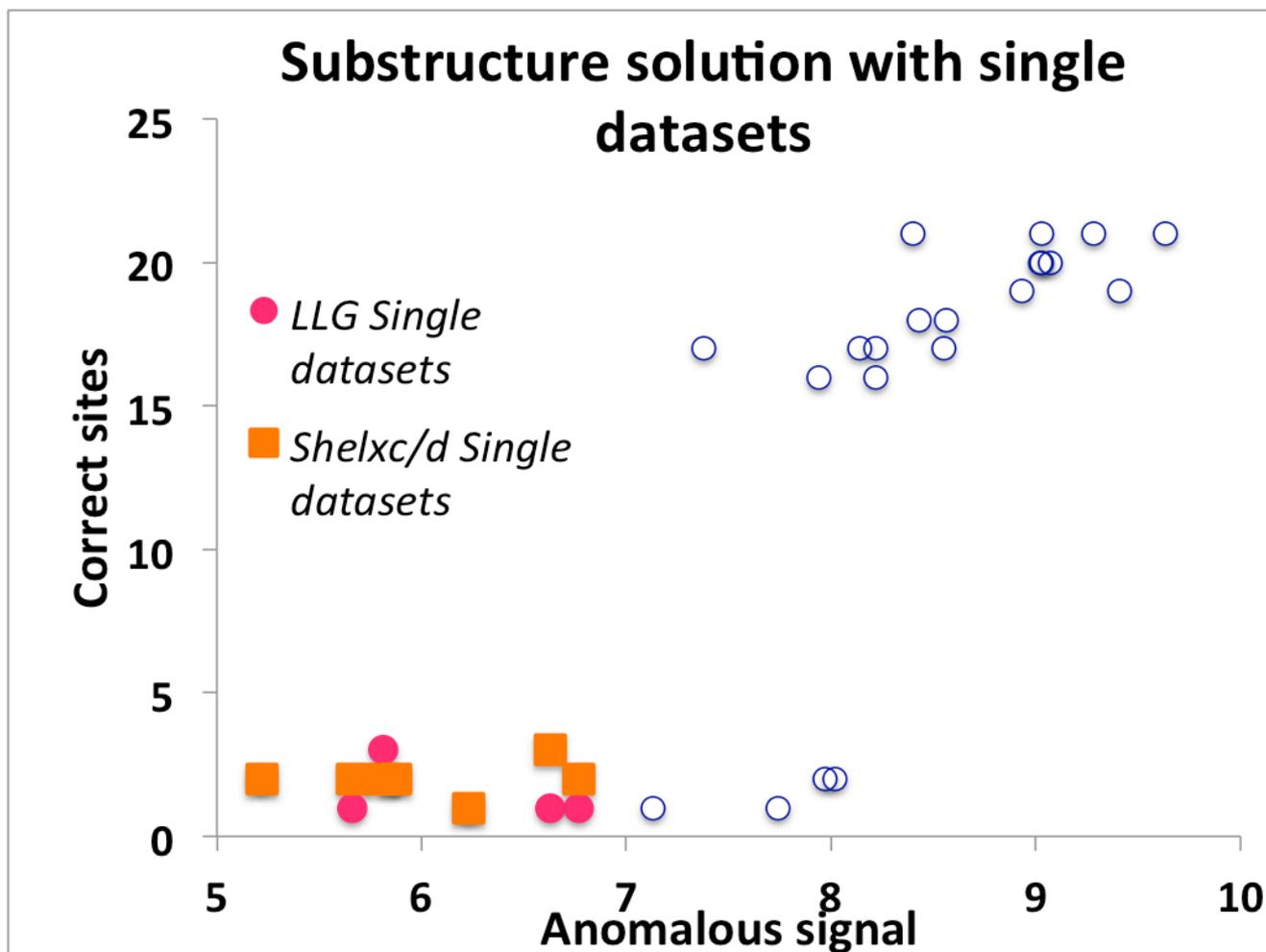
Merge of crystals 6, 7

AutoSol/Autobuild R/Rfree=0.22/0.26





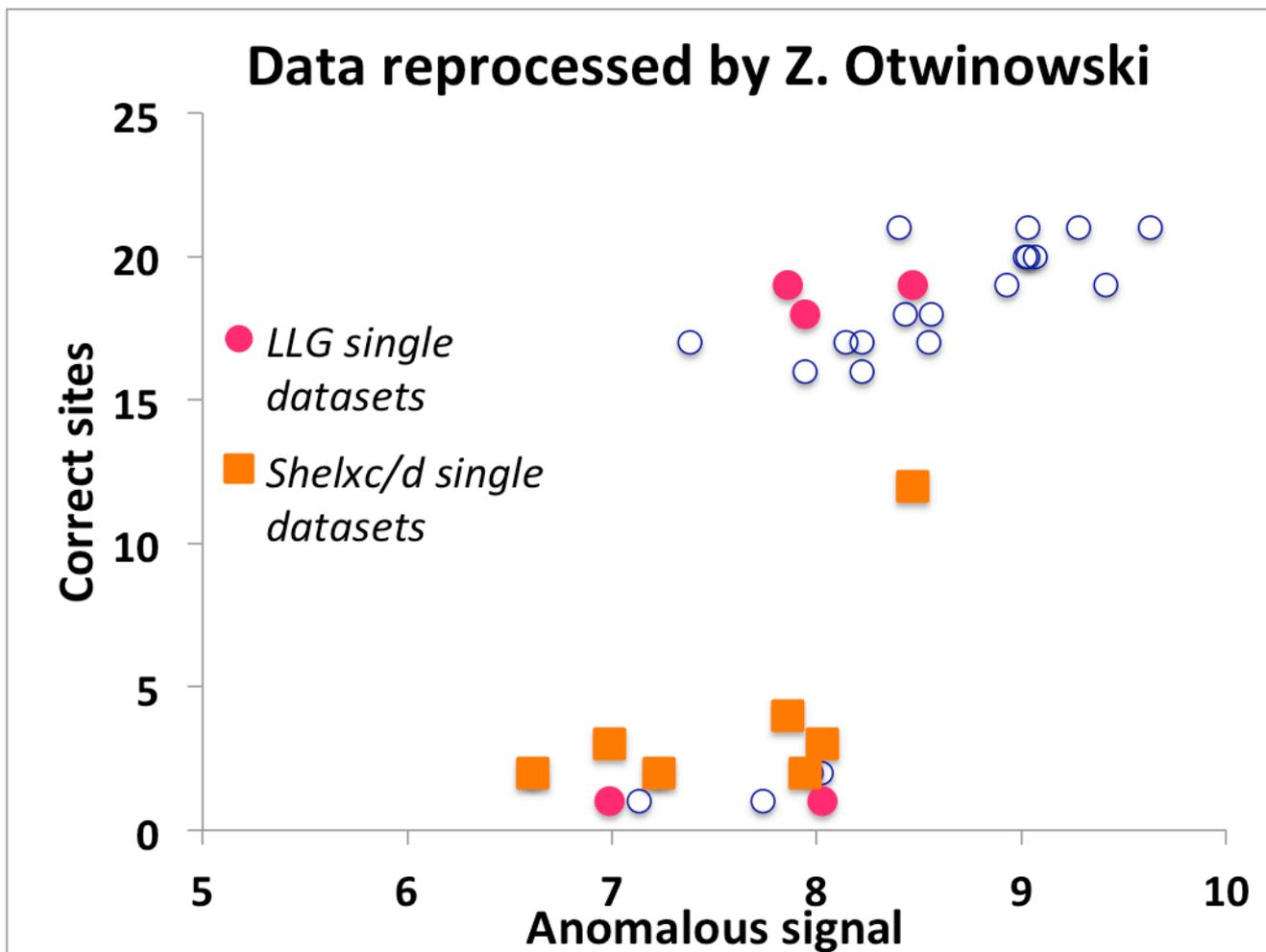
CysZ multi-crystal sulfur-SAD data





CysZ multi-crystal sulfur-SAD data

(The minimum number of datasets for this structure is 1)



Structure determination with weak anomalous signal

AutoSol:

Substructure solution, phasing, density modification, preliminary model-building

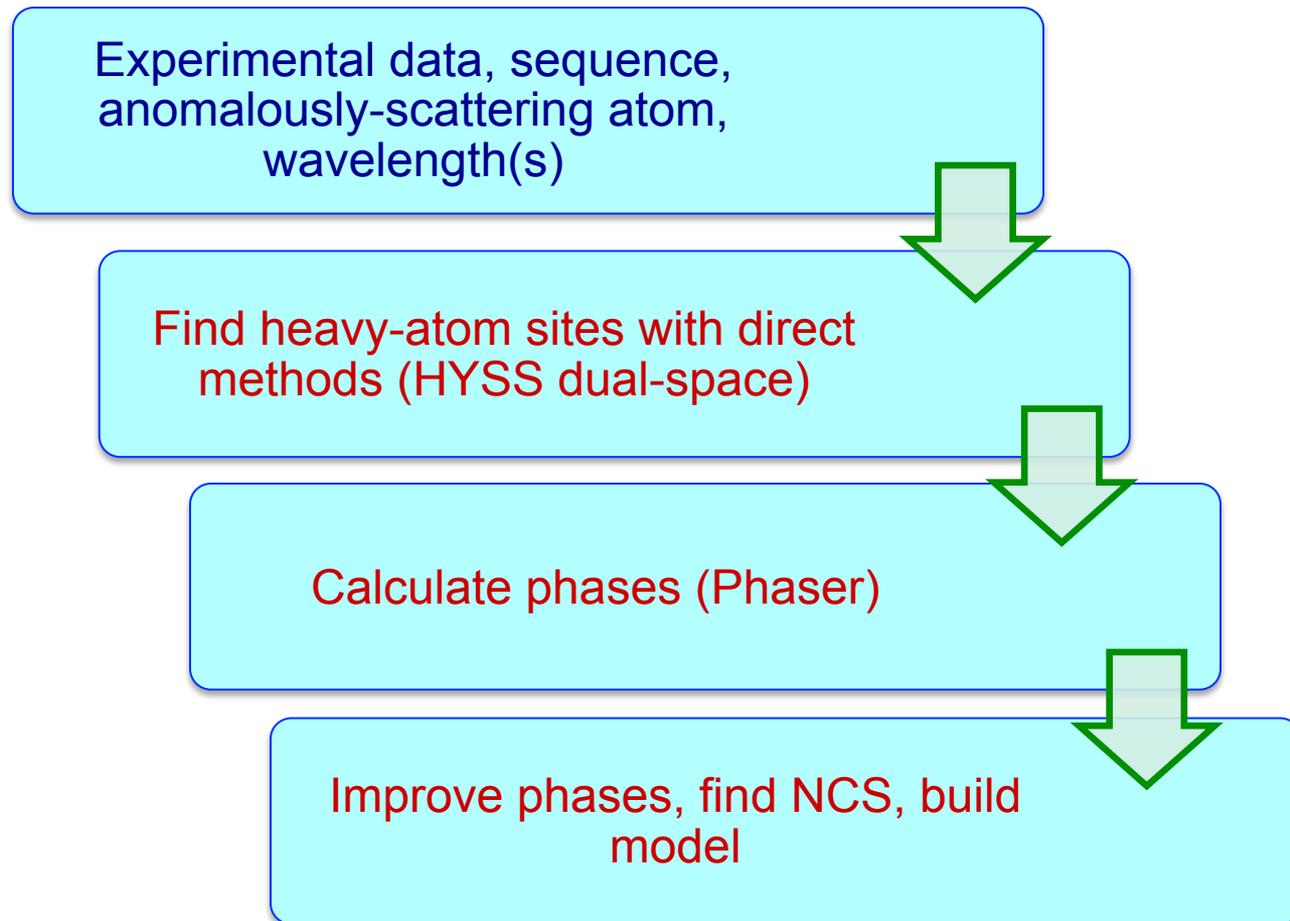
AutoBuild

Iterative model-building, refinement, density modification

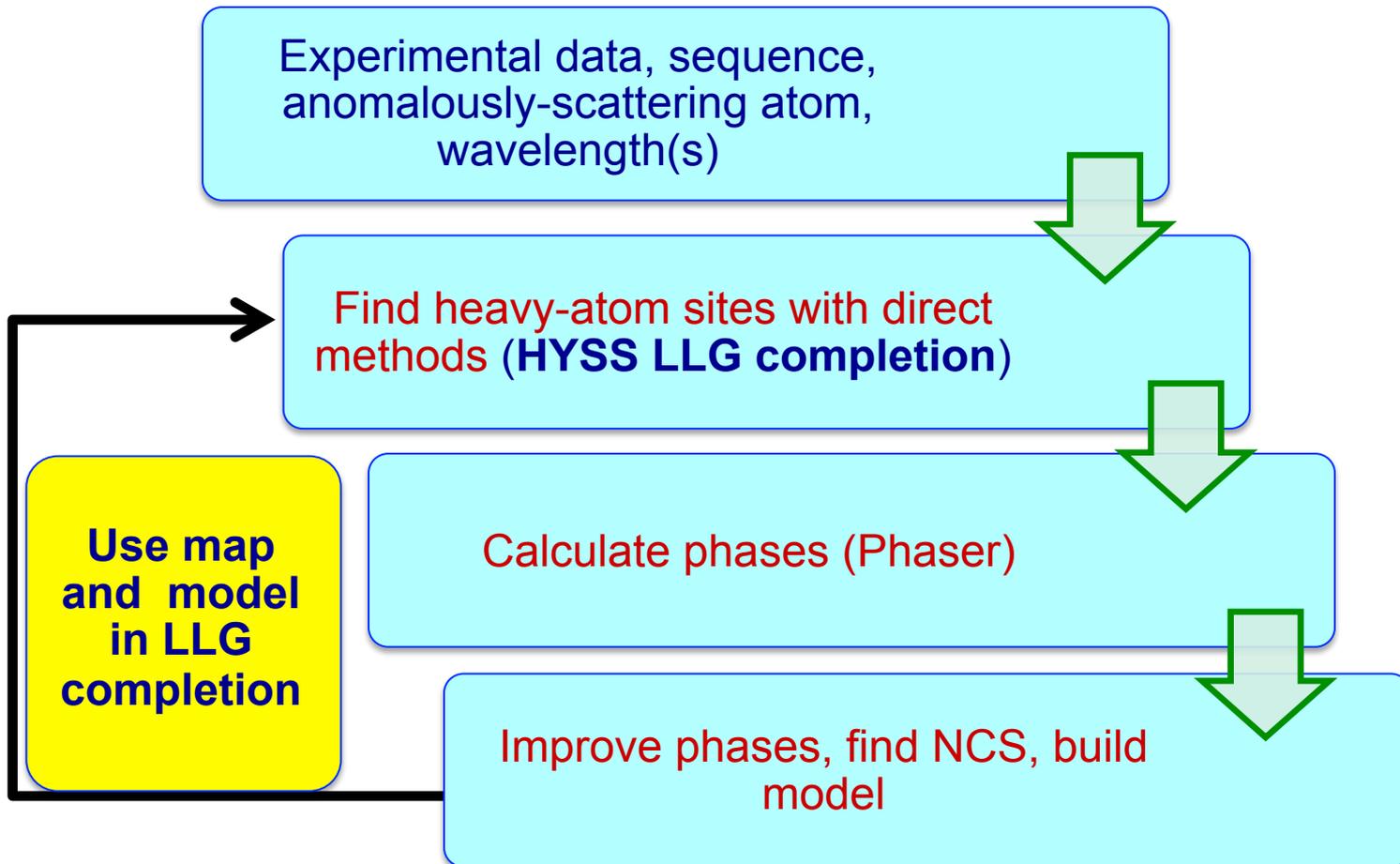
Parallel AutoBuild

Parallel runs of AutoBuild with map averaging and picking best models

Structure solution with phenix.autosol

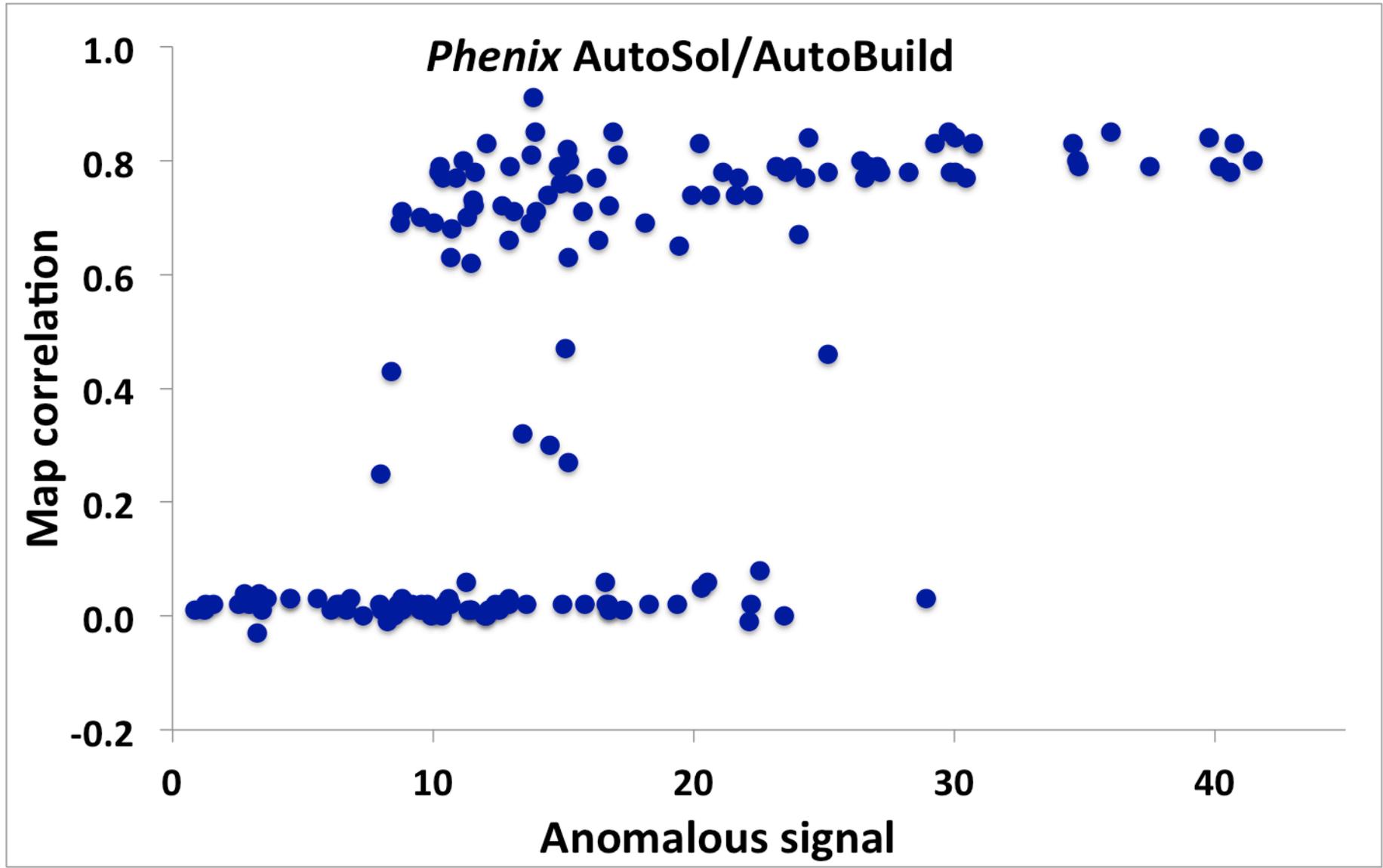


Structure solution with phenix.autosol: enhancements for weak SAD data

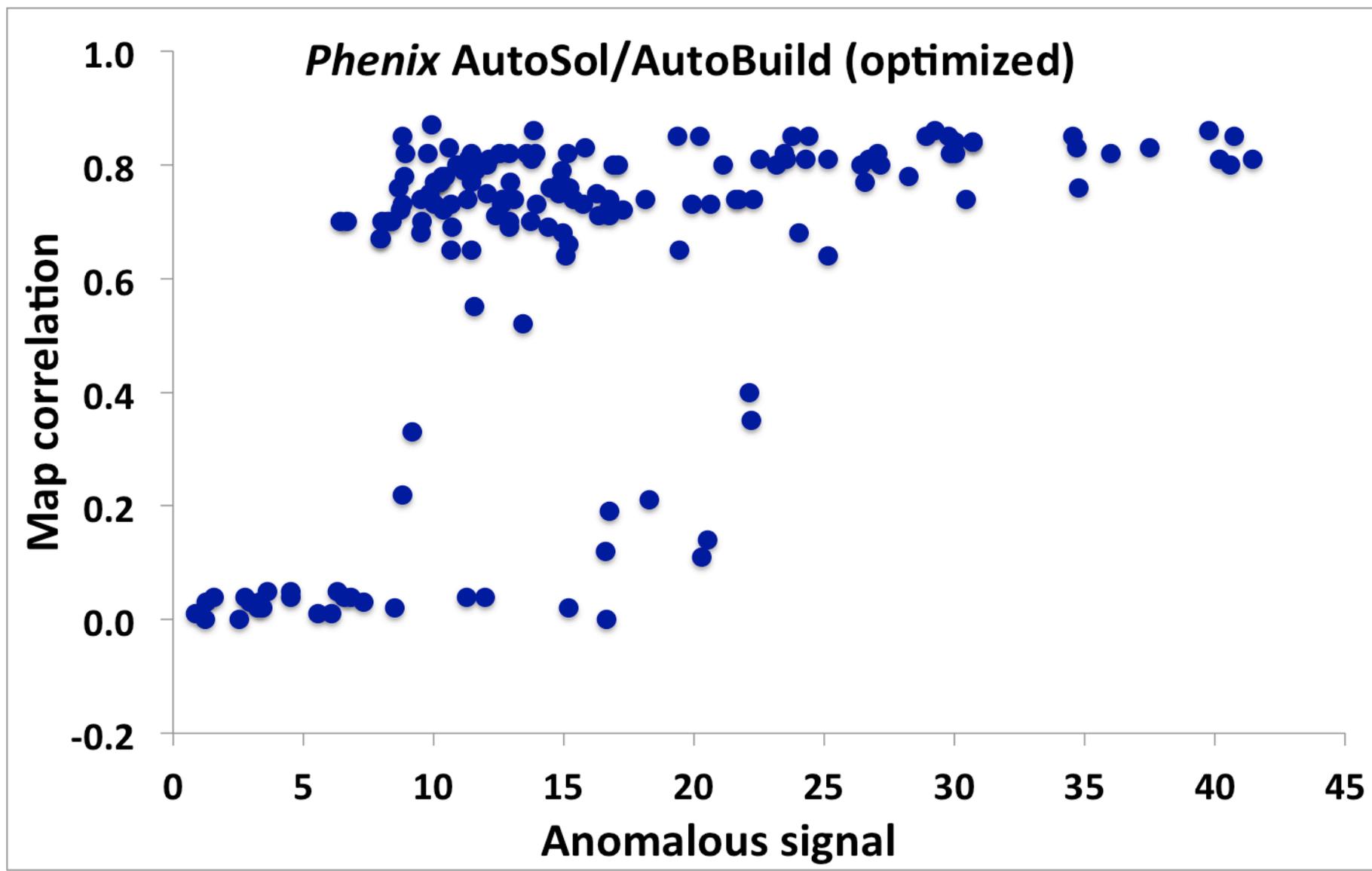


AutoSol structure solution
164 SAD datasets from PDB

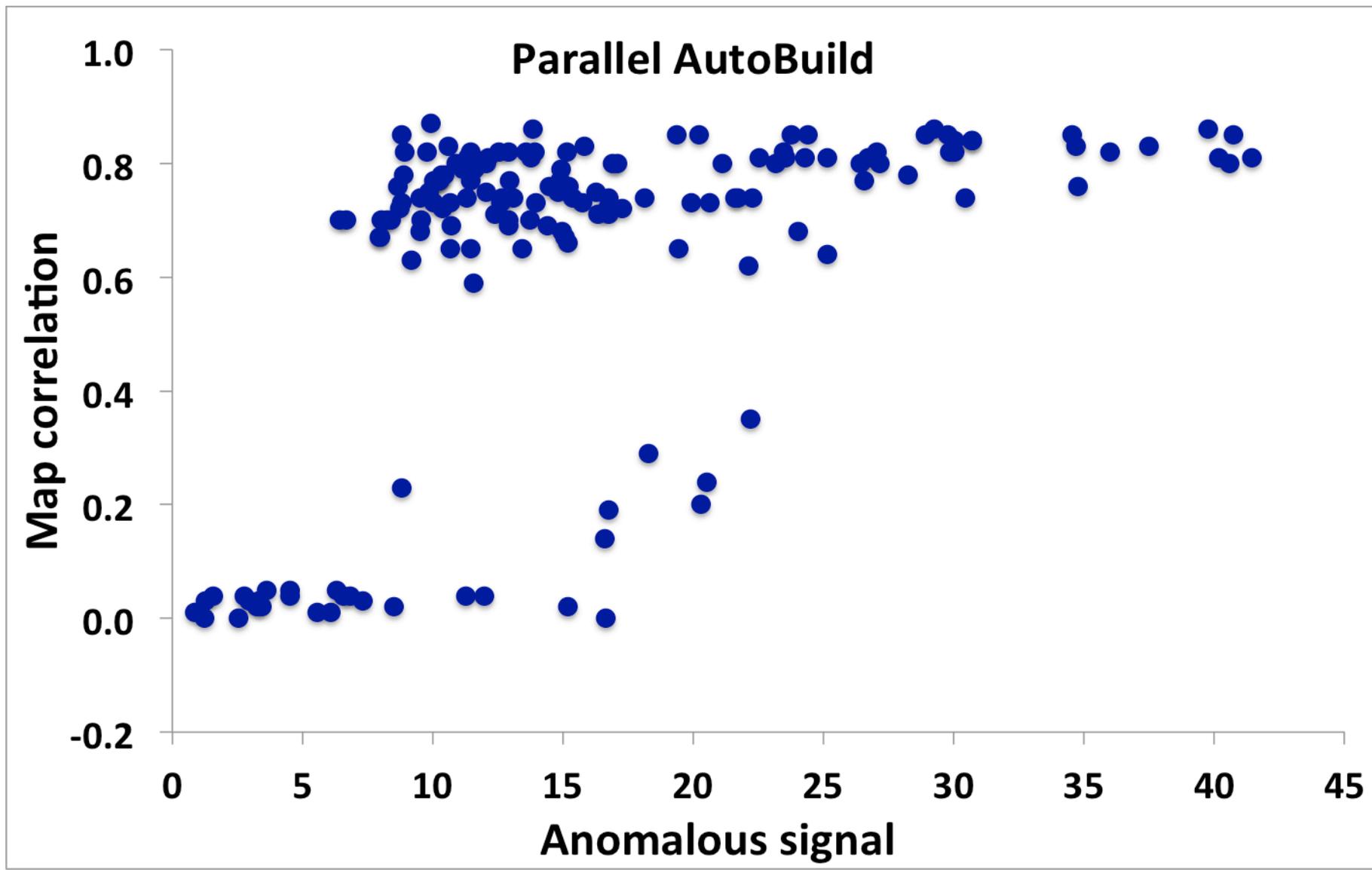
(including inflection/remote datasets not previously used as SAD data)



AutoSol structure solution
164 SAD datasets from PDB



AutoBuild model-building
164 SAD datasets from PDB



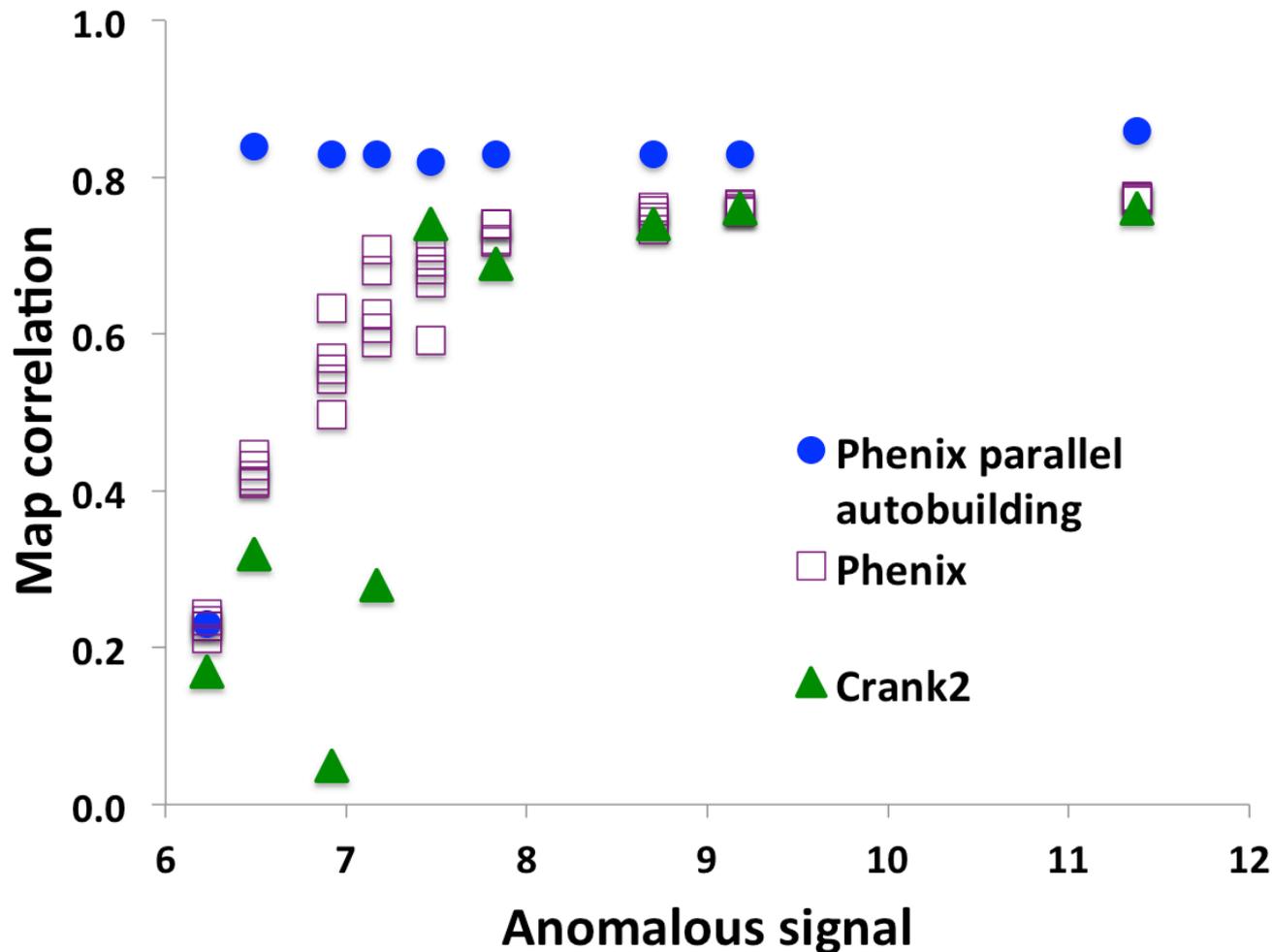


Holton Challenge data

Starting point: known sites.

*Calculate phases, carry out iterative density modification,
model-building and refinement.*

Final map correlation vs anomalous signal-to-noise.



Estimating the anomalous signal from the data

Gold standards for the anomalous information:

Correlation of true and observed differences:

$$CC_{ano} \equiv \frac{\langle \Delta_{ano,j} \Delta_{ano,j}^{obs} \rangle}{\langle \Delta_{ano}^2 \rangle^{1/2} \langle \Delta_{ano}^{2,obs} \rangle^{1/2}}$$

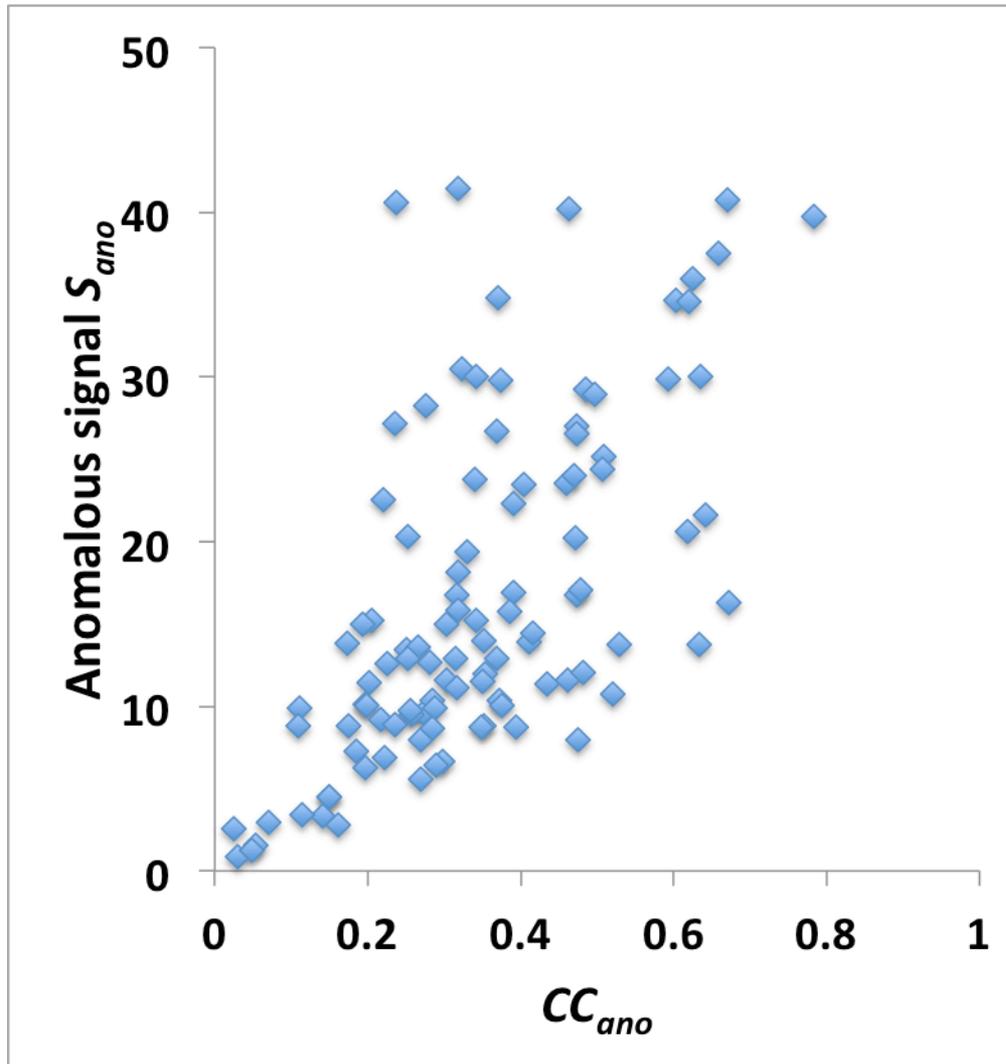
Peak height in model-phased Difference Fourier:

$$S_{ano} \equiv \frac{\langle \rho(x_k) \rangle}{\langle \rho^2 \rangle^{1/2}}$$

Relationship between CC_{ano} and S_{ano}

$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

Checking the relationship between CC_{ano} and S_{ano}

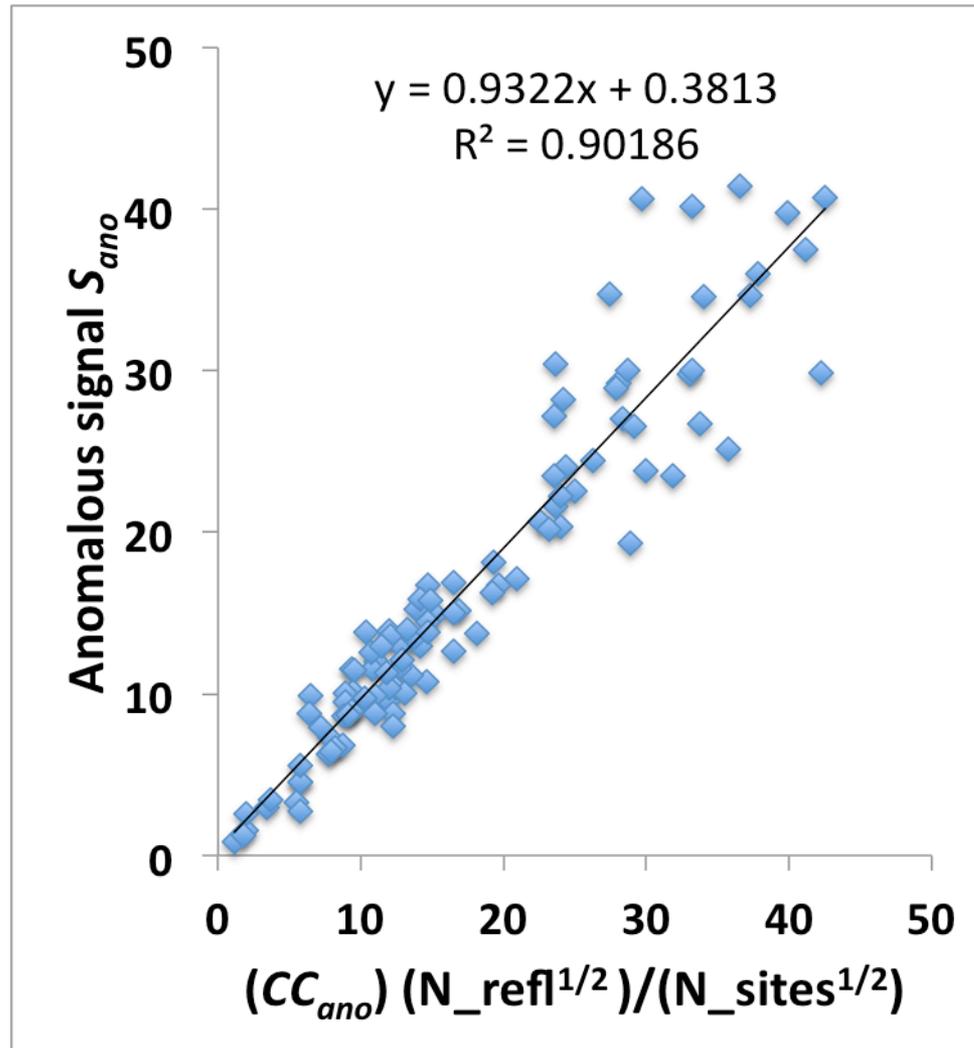


$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

CC_{ano} : Correlation of anomalous differences with model differences

S_{ano} : Peak height in model-phased difference Fourier

Checking the relationship between CC_{ano} and S_{ano}



$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

CC_{ano} : Correlation of anomalous differences with model differences

S_{ano} : Peak height in model-phased difference Fourier

Estimating the anomalous correlation CC_{ano} from the data

CC_{ano} estimates based on simple theory:

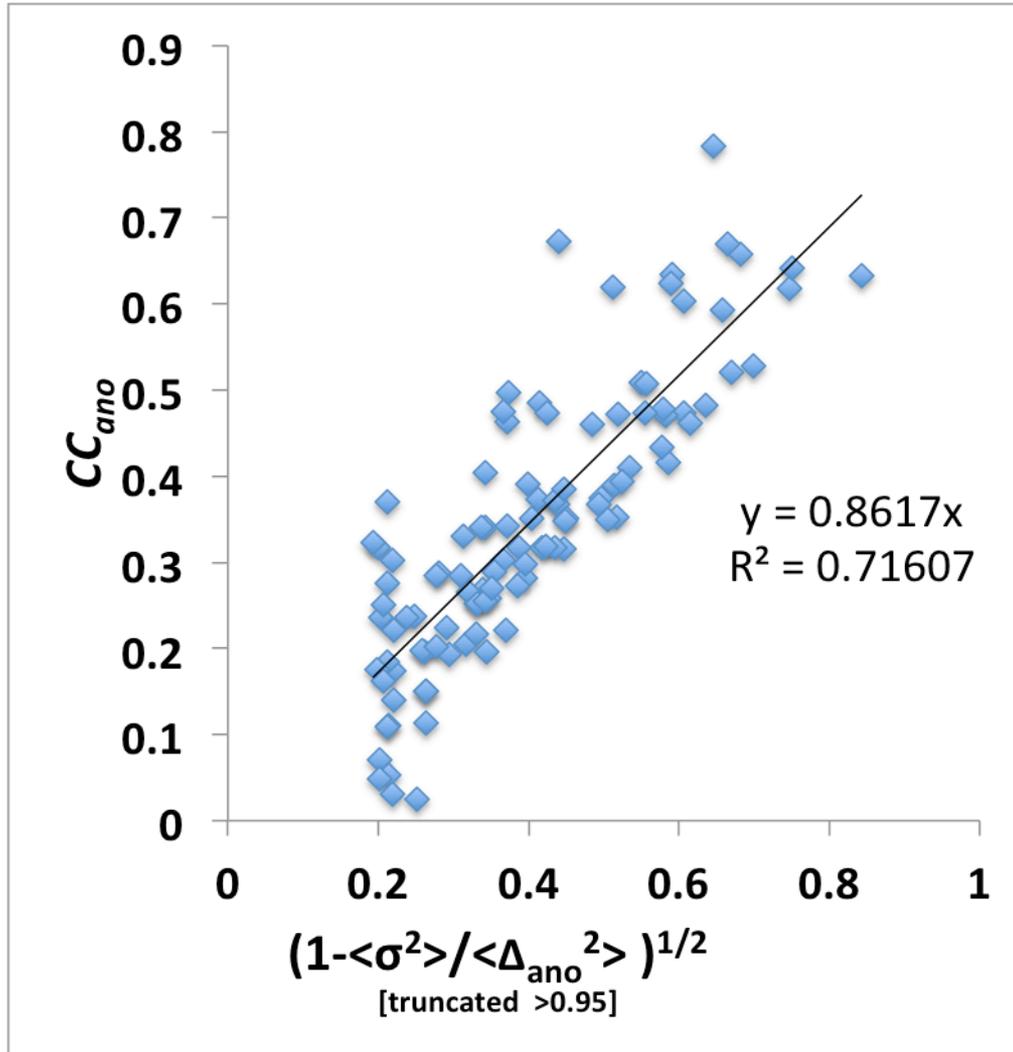
Estimated from experimental uncertainties and anomalous differences

Estimated from half-dataset correlation of experimental anomalous differences

$$E_{ano}^2 = \frac{\langle \sigma_{ano}^2 \rangle}{\langle \Delta_{ano}^{2,obs} \rangle}$$
$$CC_{ano} \sim [1 - E_{ano}^2]^{1/2}$$

$$CC_{ano}^* = \left[\frac{2CC_{ano}^{half_dataset}}{1 + CC_{ano}^{half_dataset}} \right]^{1/2}$$
$$CC_{ano} \sim CC_{ano}^*$$

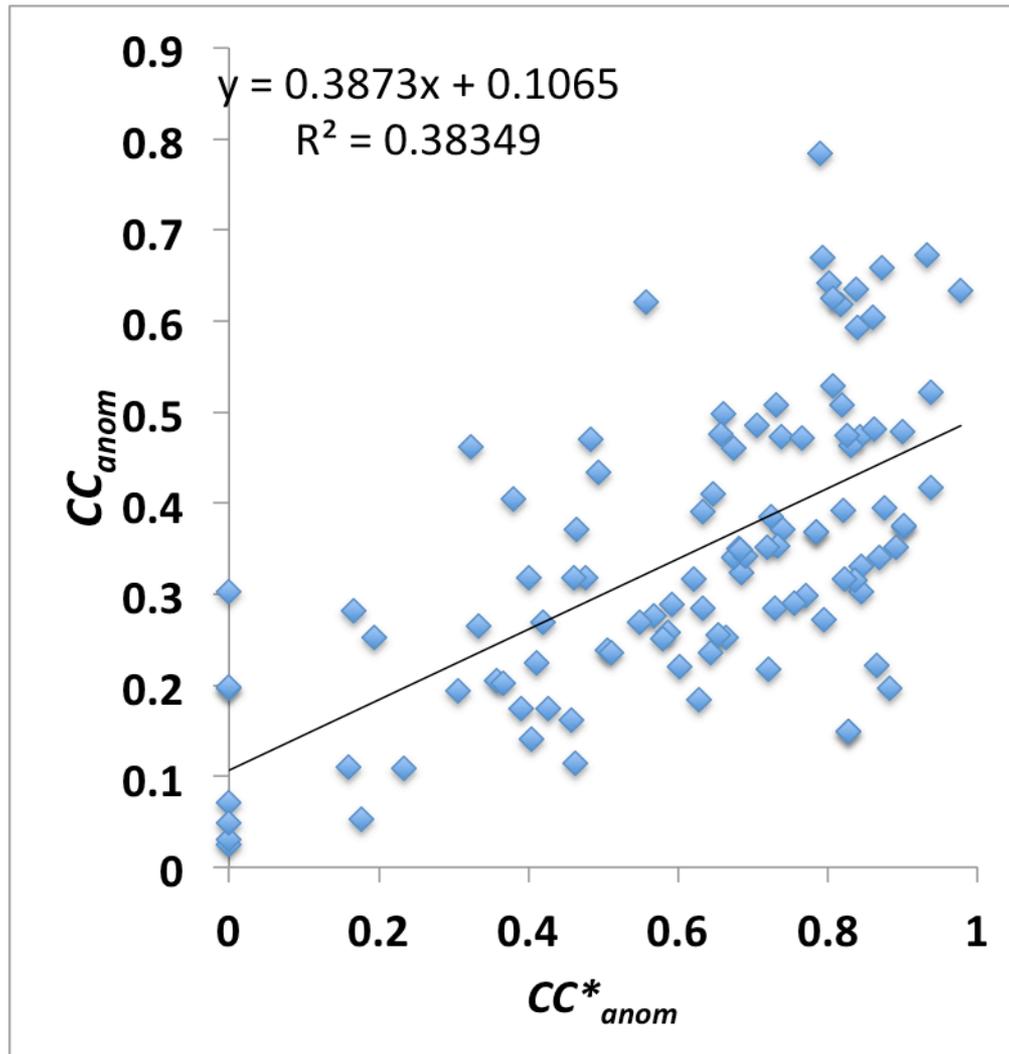
Estimating CC_{ano} from experimental uncertainties and anomalous differences



$$E_{ano}^2 = \frac{\langle \sigma_{ano}^2 \rangle}{\langle \Delta_{ano}^{2,obs} \rangle}$$

$$CC_{ano} \sim [1 - E_{ano}^2]^{1/2}$$

Estimating CC_{ano} from the half-dataset anomalous correlation.



$$CC_{ano}^* = \left[\frac{2CC_{ano}^{half_dataset}}{1 + CC_{ano}^{half_dataset}} \right]^{1/2}$$

$$CC_{ano} \sim CC_{ano}^*$$

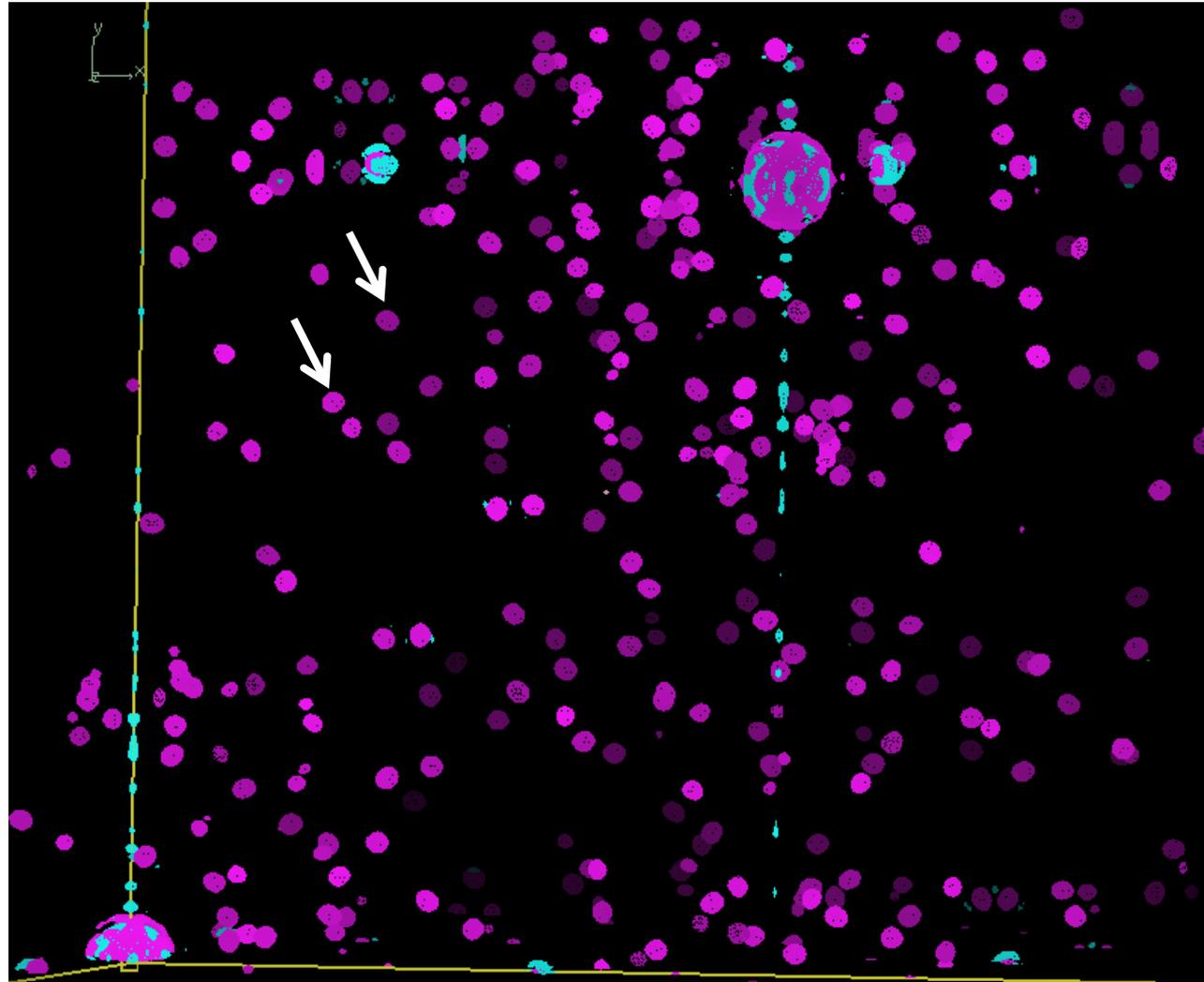


Skew of anomalous difference Patterson

Anomalous difference Patterson for 2a3n (14 Se sites, 1.3 Å)

Contours at $\pm 4\sigma$.
Positive pink, negative blue

Model anomalous differences





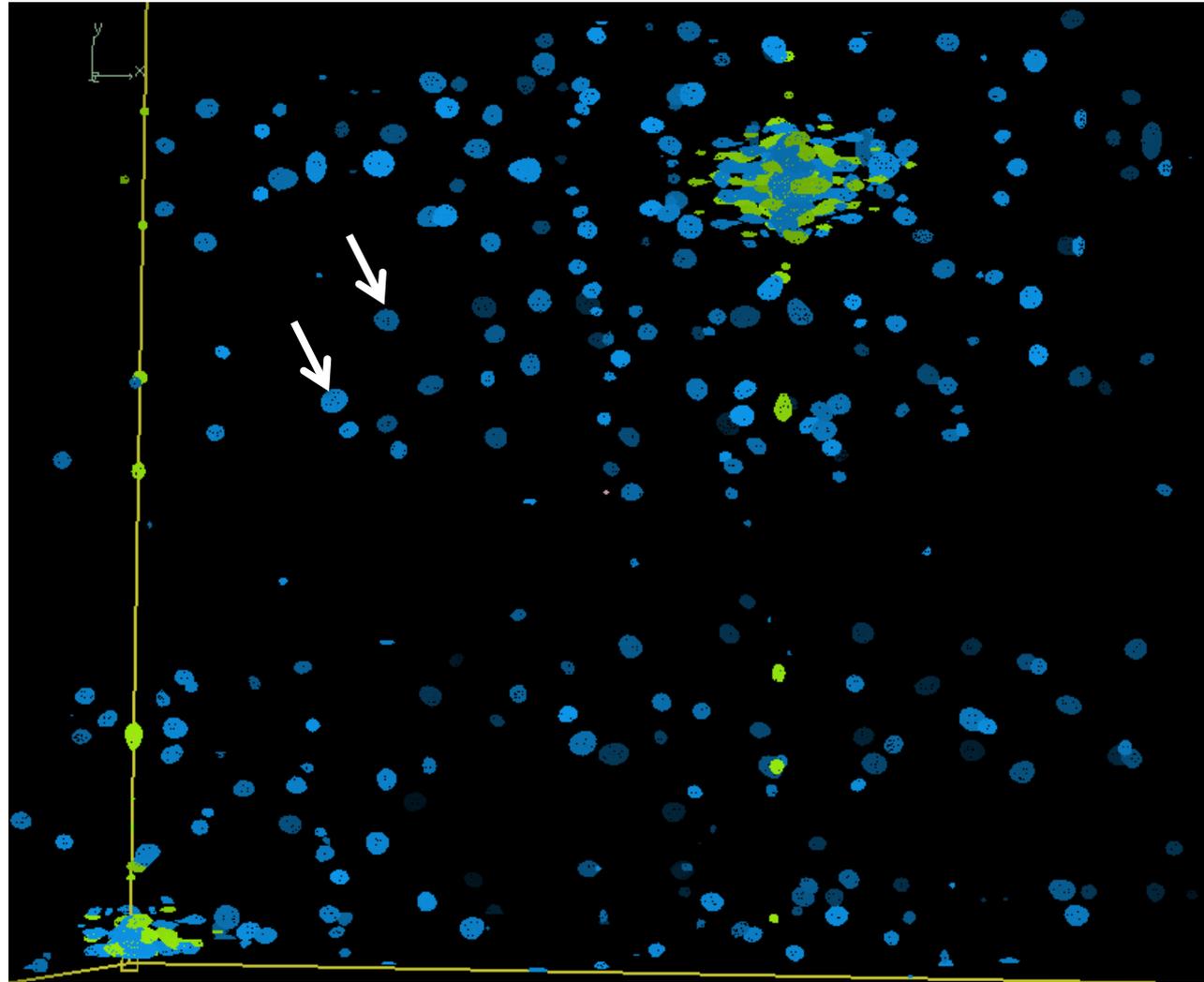
Skew of anomalous difference Patterson

Anomalous difference Patterson for $2a3n$ (14 Se sites, 1.3 Å)

Contours at 4σ .

Positive blue, negative green

Measured anomalous differences

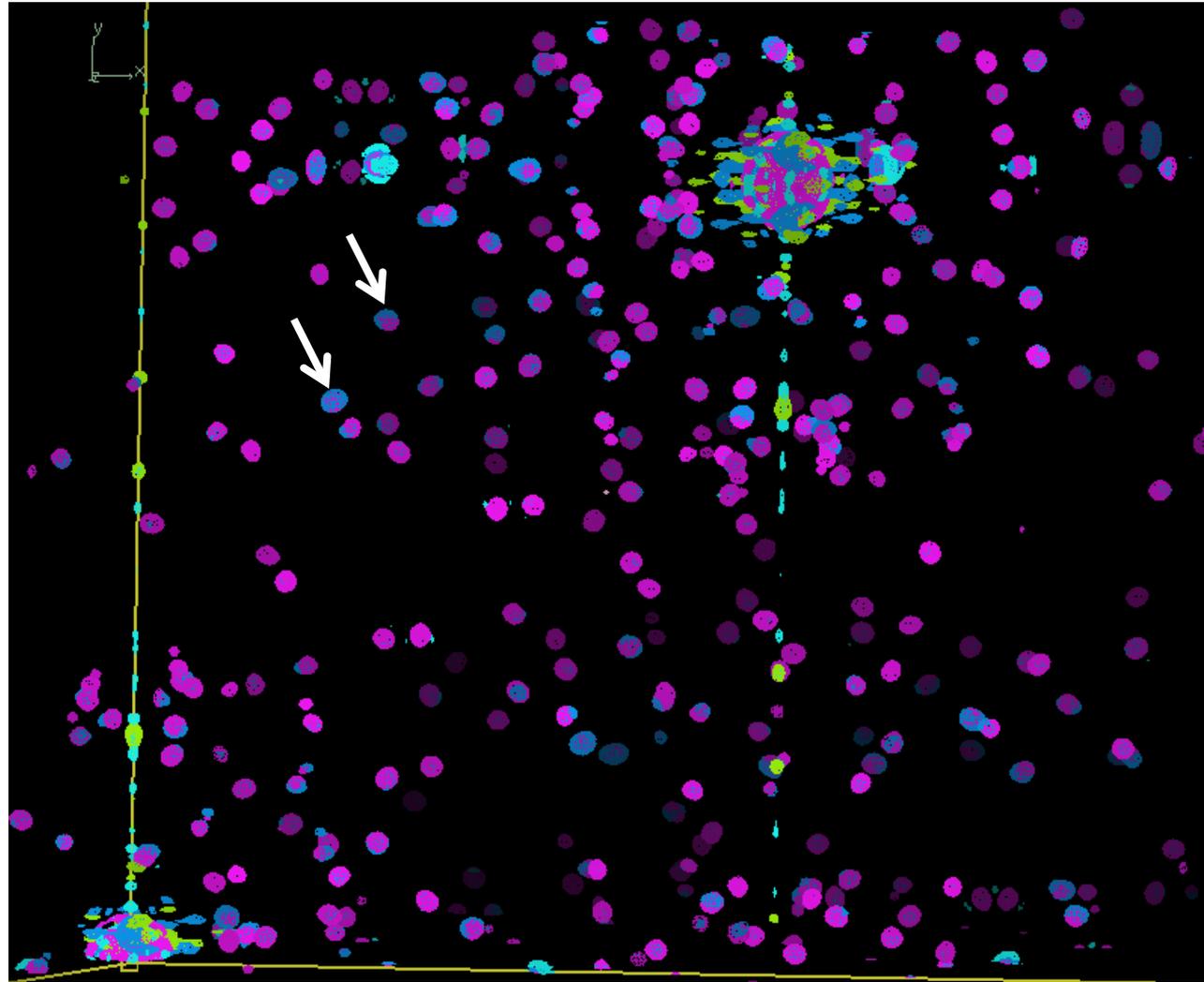




Skew of anomalous difference Patterson

Anomalous difference Patterson for 2a3n (14 Se sites, 1.3 Å)
Contours at 4σ .

Model (pink) and experimental (blue) anomalous differences





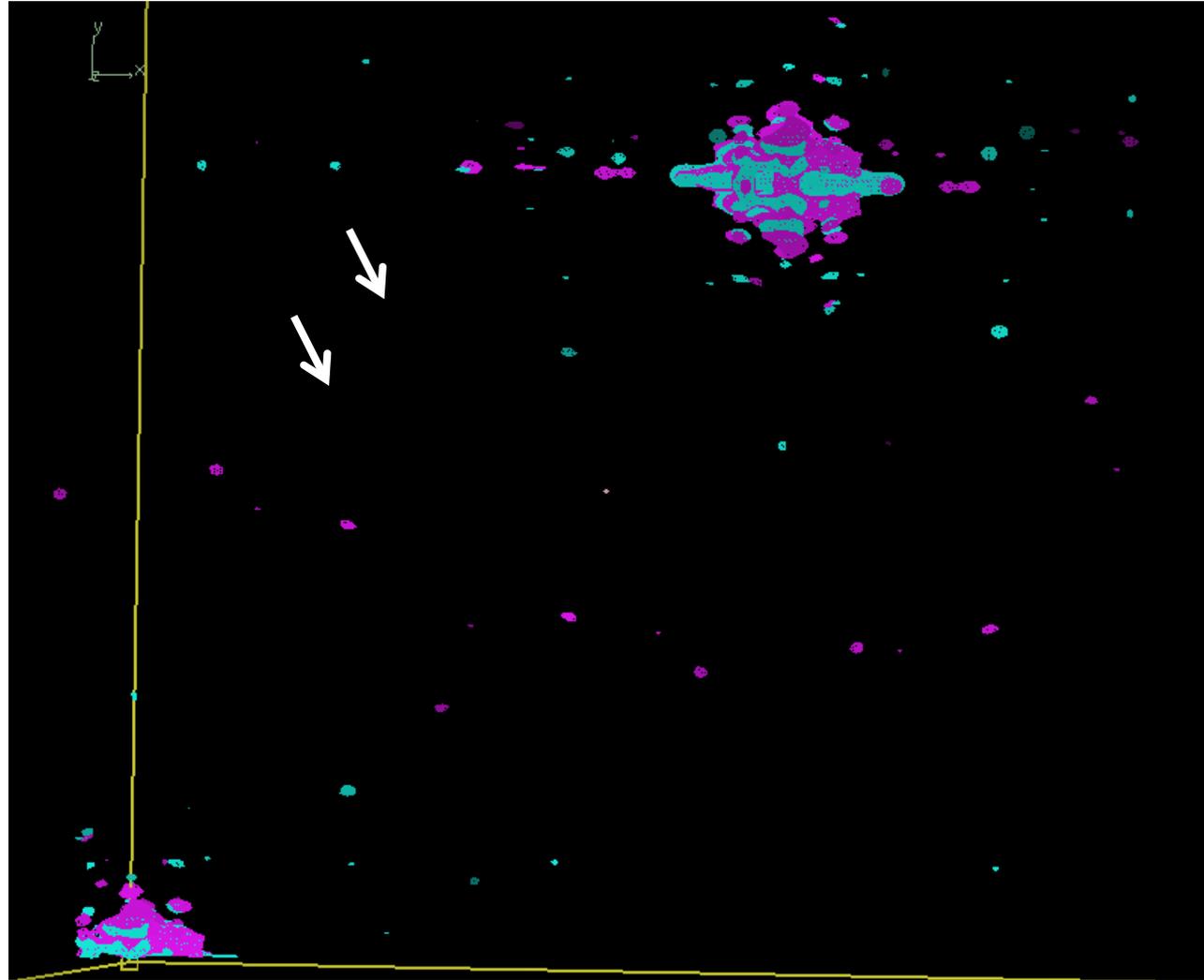
Skew of anomalous difference Patterson

Anomalous difference Patterson for 2a3n (14 Se sites, 1.3 Å)

Contours at 4σ .

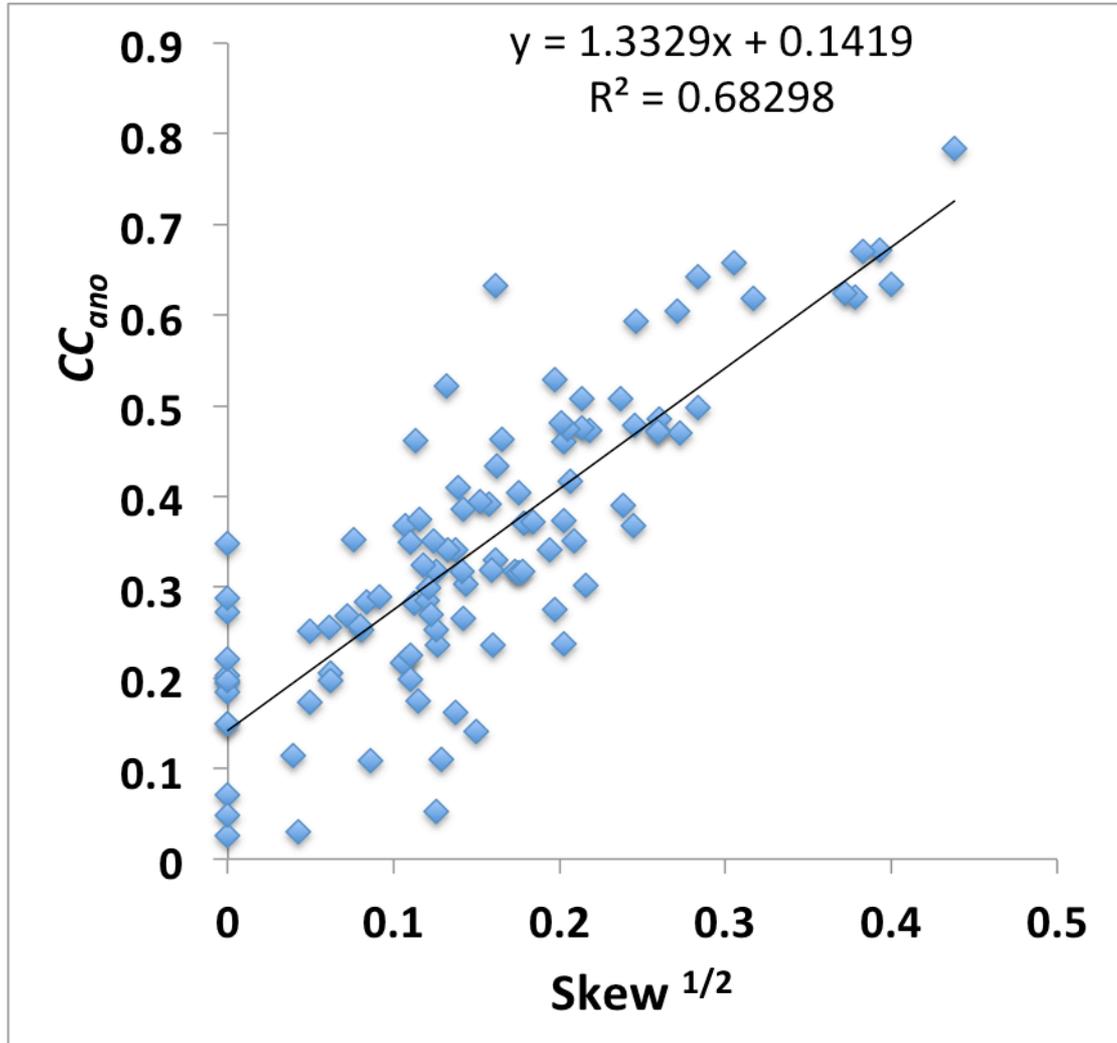
Positive blue, negative pink.

Randomized anomalous differences





Estimating CC_{ano} from skew of the anomalous difference Patterson



$$CC_{ano} \sim skew_{Patterson}^{1/2}$$

Estimating the anomalous correlation



$$CC_{ano}$$

Estimated fraction of observed anomalous differences that is noise

$$E_{ano}^2 = \frac{\langle \sigma_{ano}^2 \rangle}{\langle \Delta_{ano}^{2,obs} \rangle}$$

$$CC_{ano} \sim [1 - E_{ano}^2]^{1/2}$$

Half-dataset CC of anomalous differences

$$CC_{ano}^* = \left[\frac{2CC_{ano}^{half_dataset}}{1 + CC_{ano}^{half_dataset}} \right]^{1/2}$$

$$CC_{ano} \sim CC_{ano}^*$$

Skew of anomalous difference Patterson

$$CC_{ano} \sim skew_{Patterson}^{1/2}$$

Estimating the anomalous signal S_{ano}

Estimation of S_{ano} requires the value of CC_{ano} and the number of sites

$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

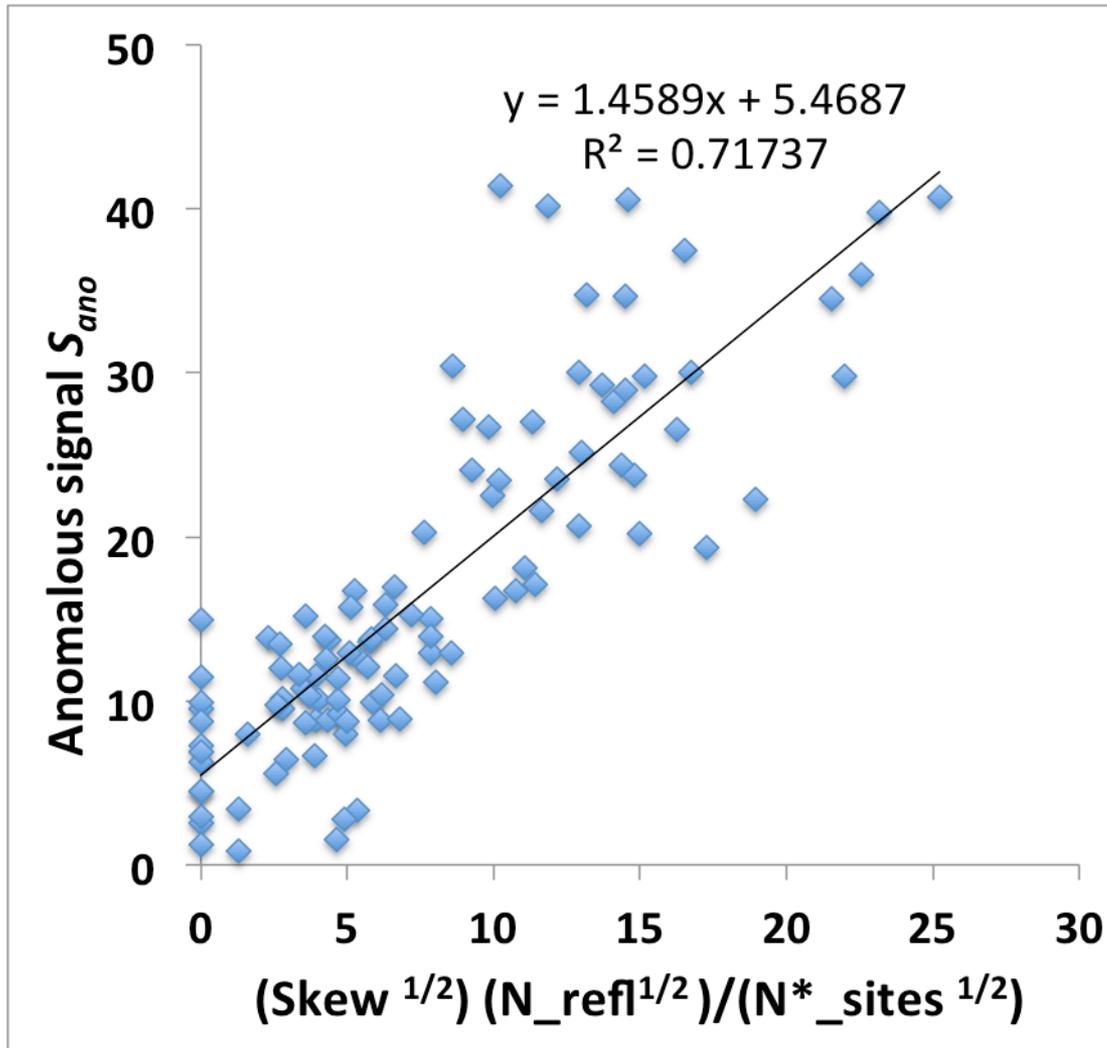
Use *phenix.autosol* estimate of number of sites

Based on sequence, asymmetric unit volume

Guess of number of NCS copies

Guess of number of sites for atoms other than S, Se
(typically 1-2 per 100 residues)

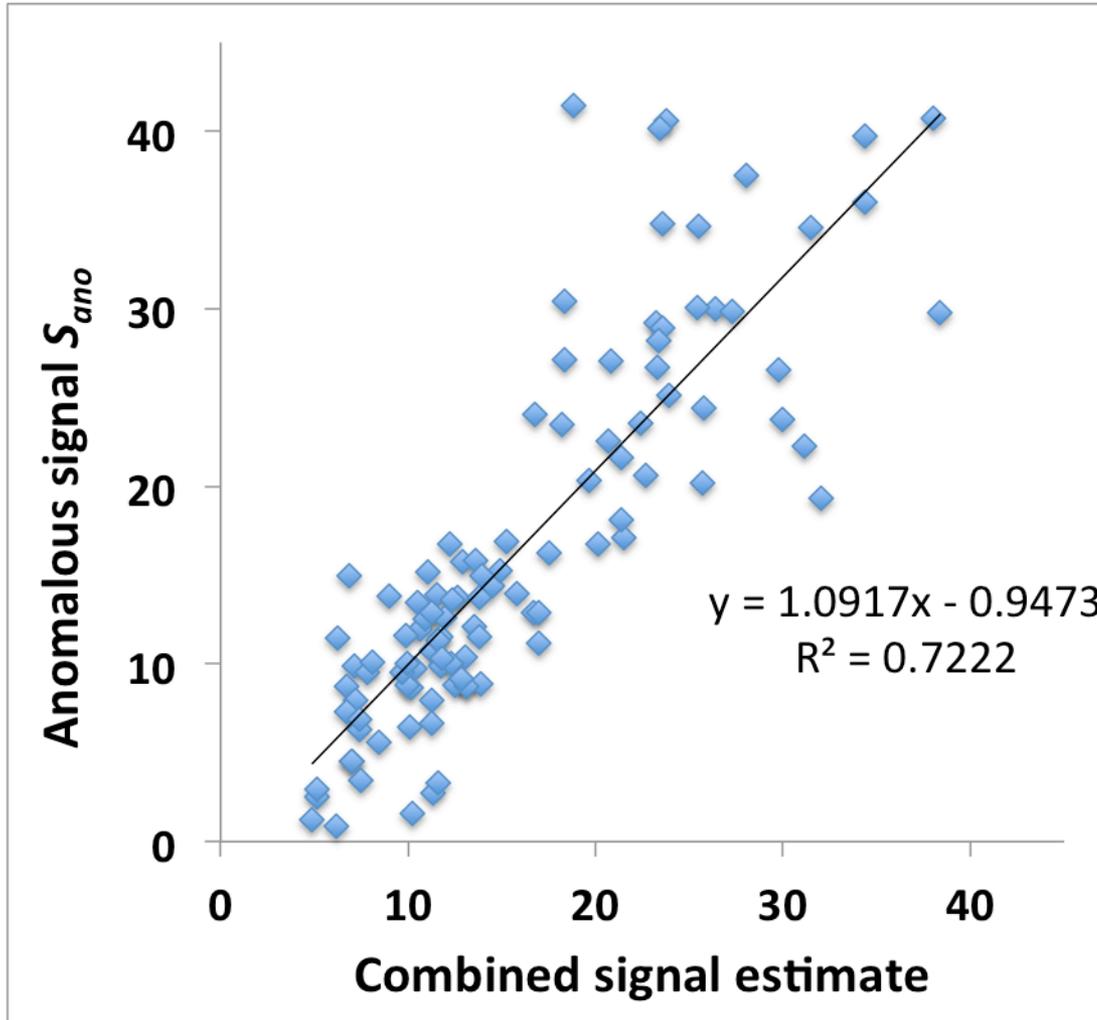
Estimating S_{ano} from skew of the anomalous difference Patterson



$$CC_{ano} \sim skew_{Patterson}^{1/2}$$

$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

Estimating S_{ano} from all 3 measures of anomalous correlation



$$CC_{ano} \sim [1 - E_{ano}^2]^{1/2}$$

$$CC_{ano} \sim \left[\frac{2CC_{ano}^{half_dataset}}{1 + CC_{ano}^{half_dataset}} \right]^{1/2}$$

$$CC_{ano} \sim skew_{Patterson}^{1/2}$$

$$S_{ano} \sim CC_{ano} \frac{N_{refl}^{1/2}}{N_{sites}^{1/2} \left(\frac{5}{4}\right)^{1/2}}$$

Using the anomalous signal S_{ano} and correlation
 CC_{ano}

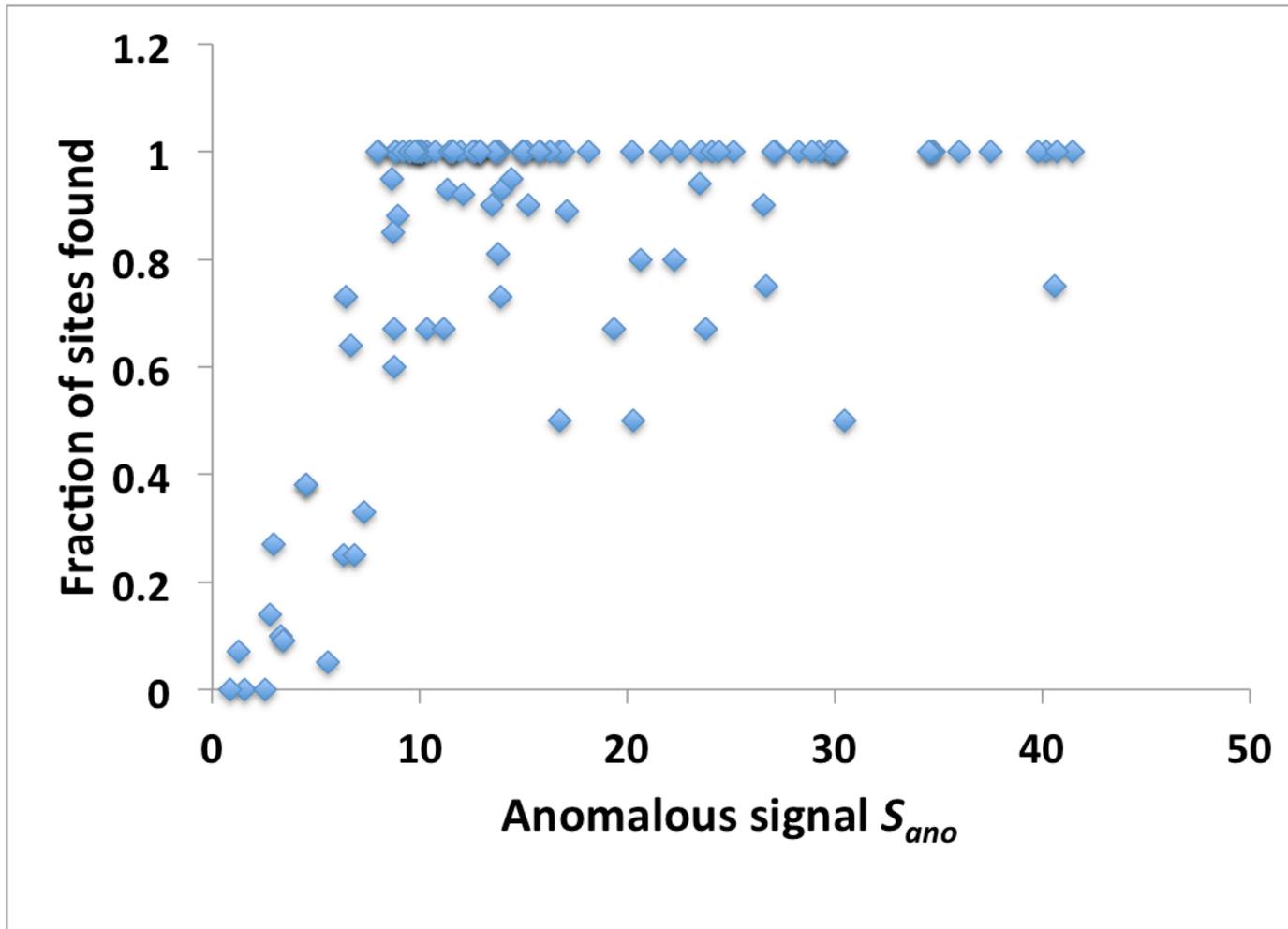
What do we expect:

Finding sites may be most closely related to map quality (S_{ano})

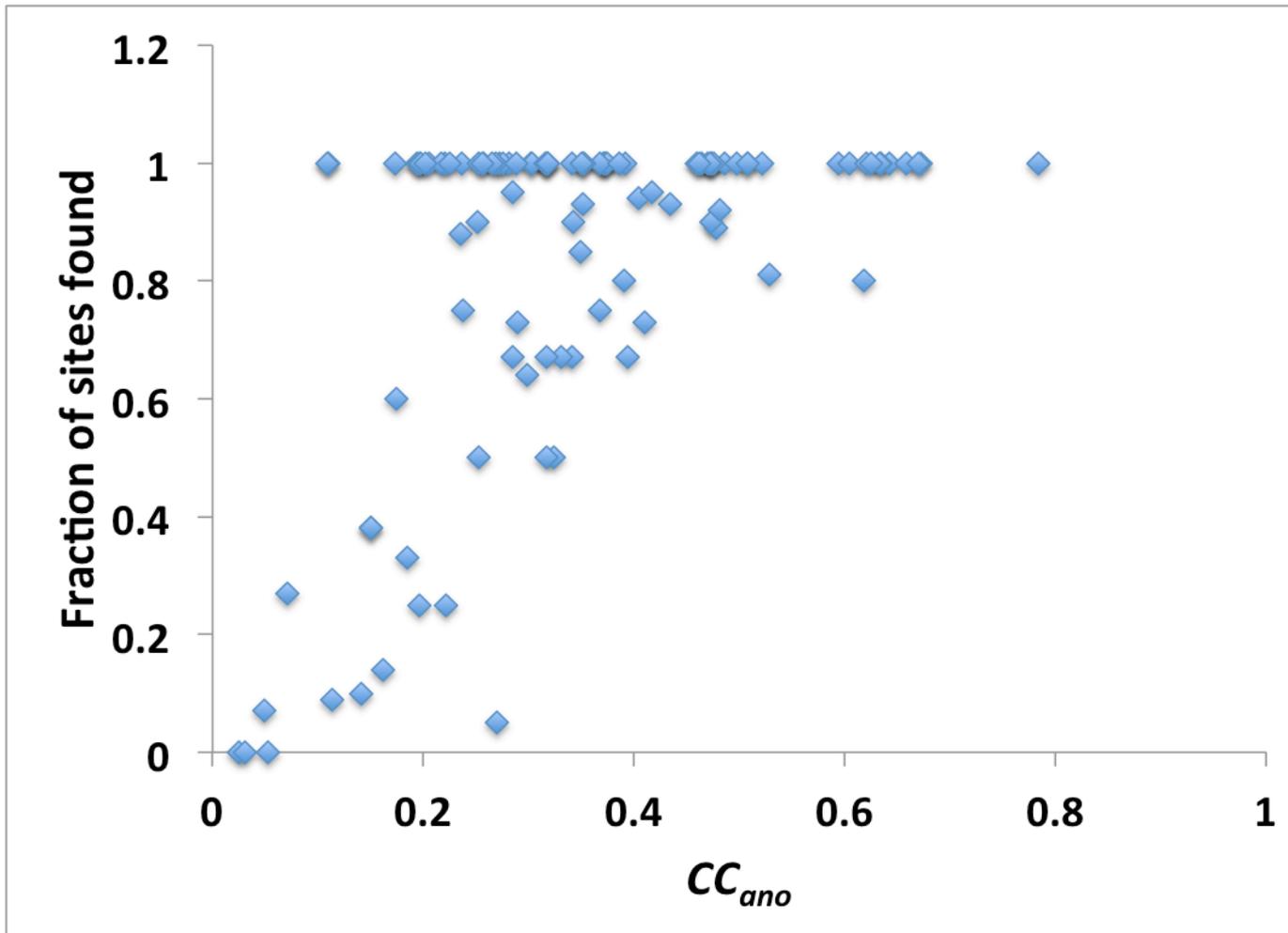
Experimental phase quality may be most closely related to the accuracy of the anomalous differences (CC_{ano})

Can I find the substructure:
Using the anomalous signal S_{ano} to guess

Best possible case: using known signal S_{ano}

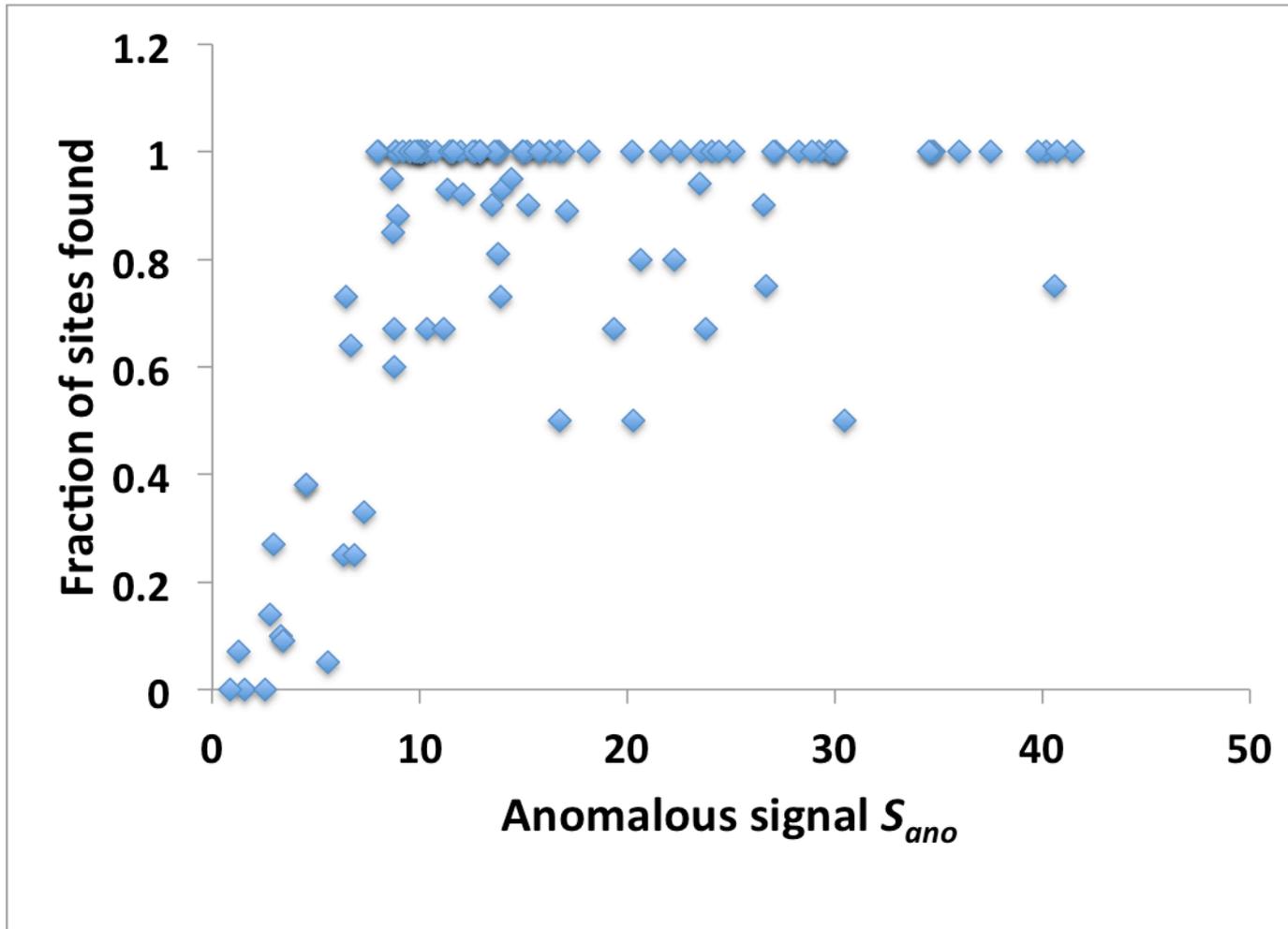


Can I find the substructure:
 Using the anomalous correlation CC_{ano}
Best possible case: using true CC_{ano}



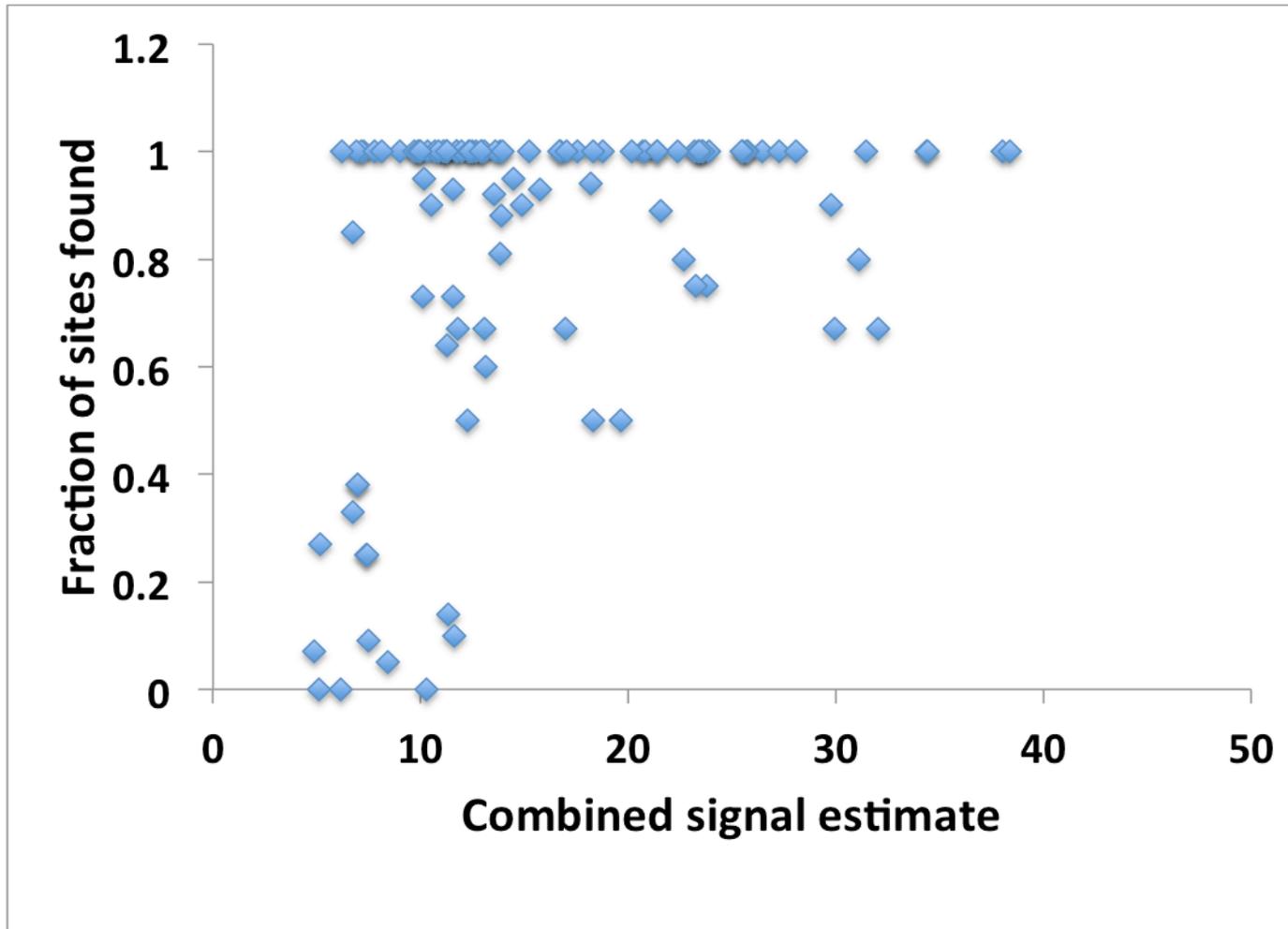
Can I find the substructure:
Using the anomalous signal S_{ano} to guess

Best possible case: using known signal S_{ano}

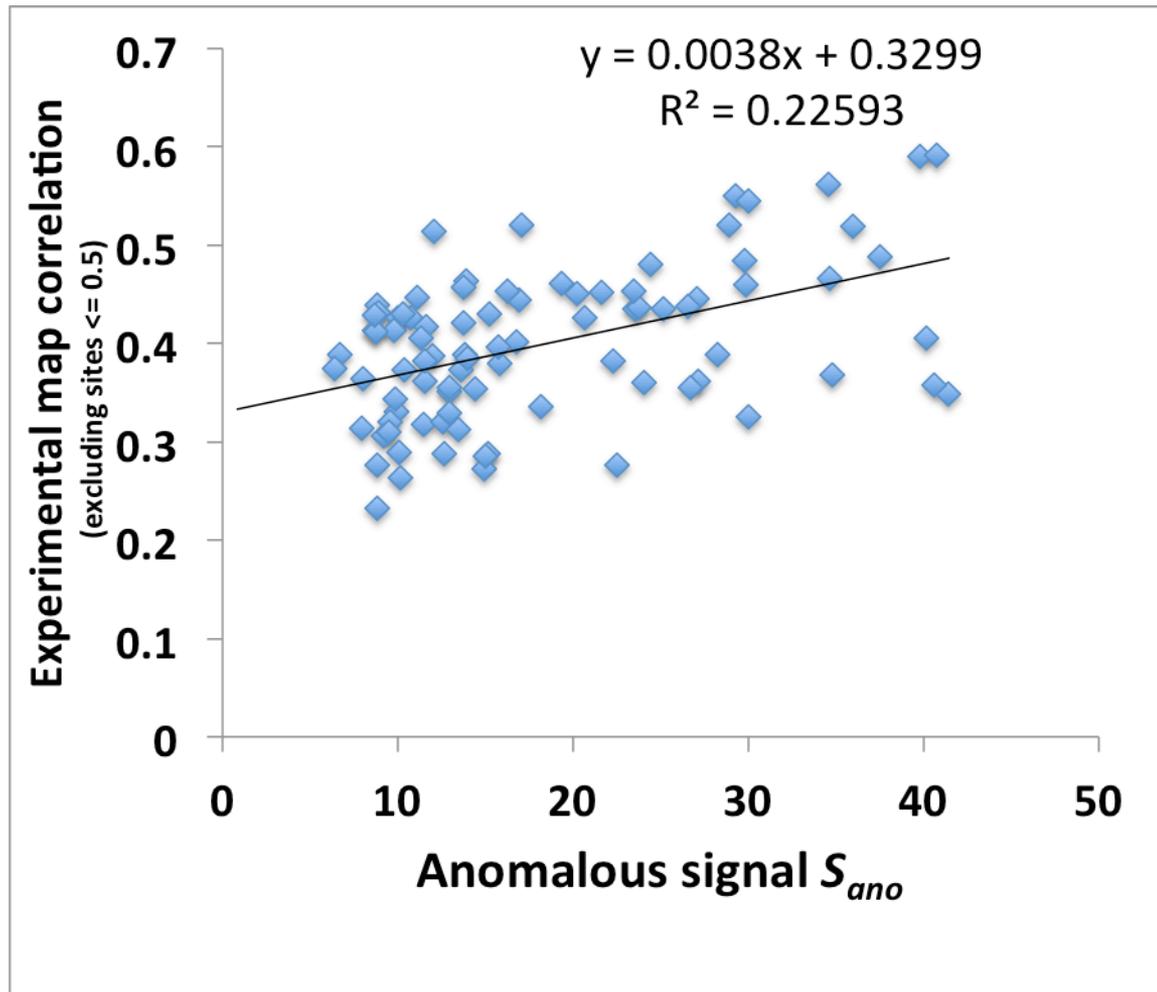


Can I find the substructure:
Using the anomalous signal S_{ano} to guess

Real-world case: S_{ano} estimated from the data

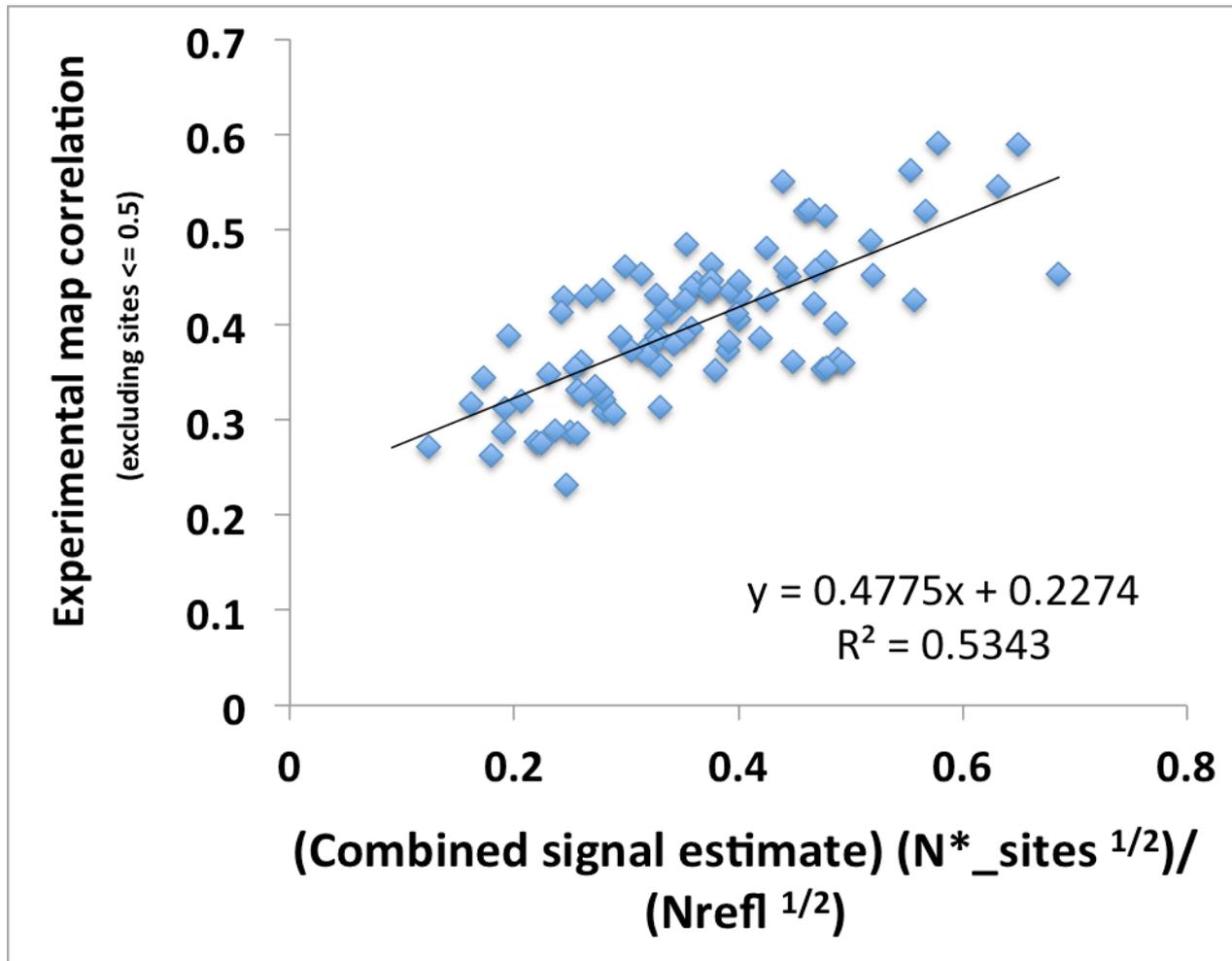


How good will the phasing be:
Could we use the anomalous signal S_{ano} ?



How good will the phasing be:
 Using the anomalous correlation CC_{ano} to
 guess

Real-world case: CC_{ano} estimated from the data



Structure solution from weak anomalous data: Perspectives

Anomalous signal and anomalous correlation are useful measures of quality and can be estimated from the data

Likelihood-based methods for finding the anomalous substructure are powerful even with weak signal

Structures can be solved with weak signal

The PHENIX Project



Lawrence Berkeley Laboratory

Paul Adams, Pavel Afonine, Nat Echols, Ralf Grosse-Kunstleve, Nigel Moriarty, Nicholas Sauter, Peter Zwart



Los Alamos National Laboratory

Tom Terwilliger, Li-Wei Hung



Randy Read, Airlie McCoy, Gabor Bunkoczi, Rob Oeffner, Richard Mifsud

Cambridge University



Duke University

Jane & David Richardson, Jeff Headd, Vincent Chen, Chris Williams, Bryan Arendall, Swati Jain, Bradley Hintze, Lizbeth Videau



An NIH/NIGMS funded Program Project

