

## Experimental phasing, including recent developments in SAD phasing

*Utilizing The Phenix Software Suite For Protein Crystallography*

*August 6-7, 2015*

University of Kansas

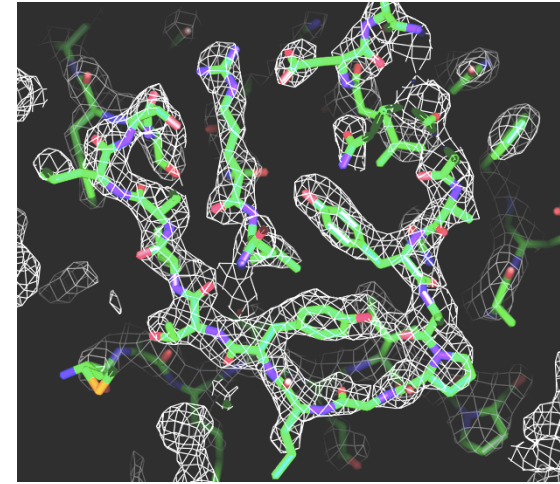
Tom Terwilliger

Los Alamos National Laboratory



# ***Steps in Single Wavelength Anomalous Diffraction (SAD) Structure Determination***

- **Plan the experiment**
- **Measure the data**
- **Scale the data**
- **Evaluate the accuracy of the anomalous differences**
- **Find the anomalous sub-structure**
- **Identify hand of sub-structure**
- **Calculate experimental phases and a map**
- **Improve the map with density modification**
- **Build and refine a model**



## Planning a SAD experiment

Maximizing the anomalous signal and the anomalous correlation

The **anomalous correlation** is a measure of the accuracy of each anomalous difference

The **anomalous signal** is a measure of how much total information is present in the anomalous differences

# Anomalous correlation: accuracy of anomalous differences

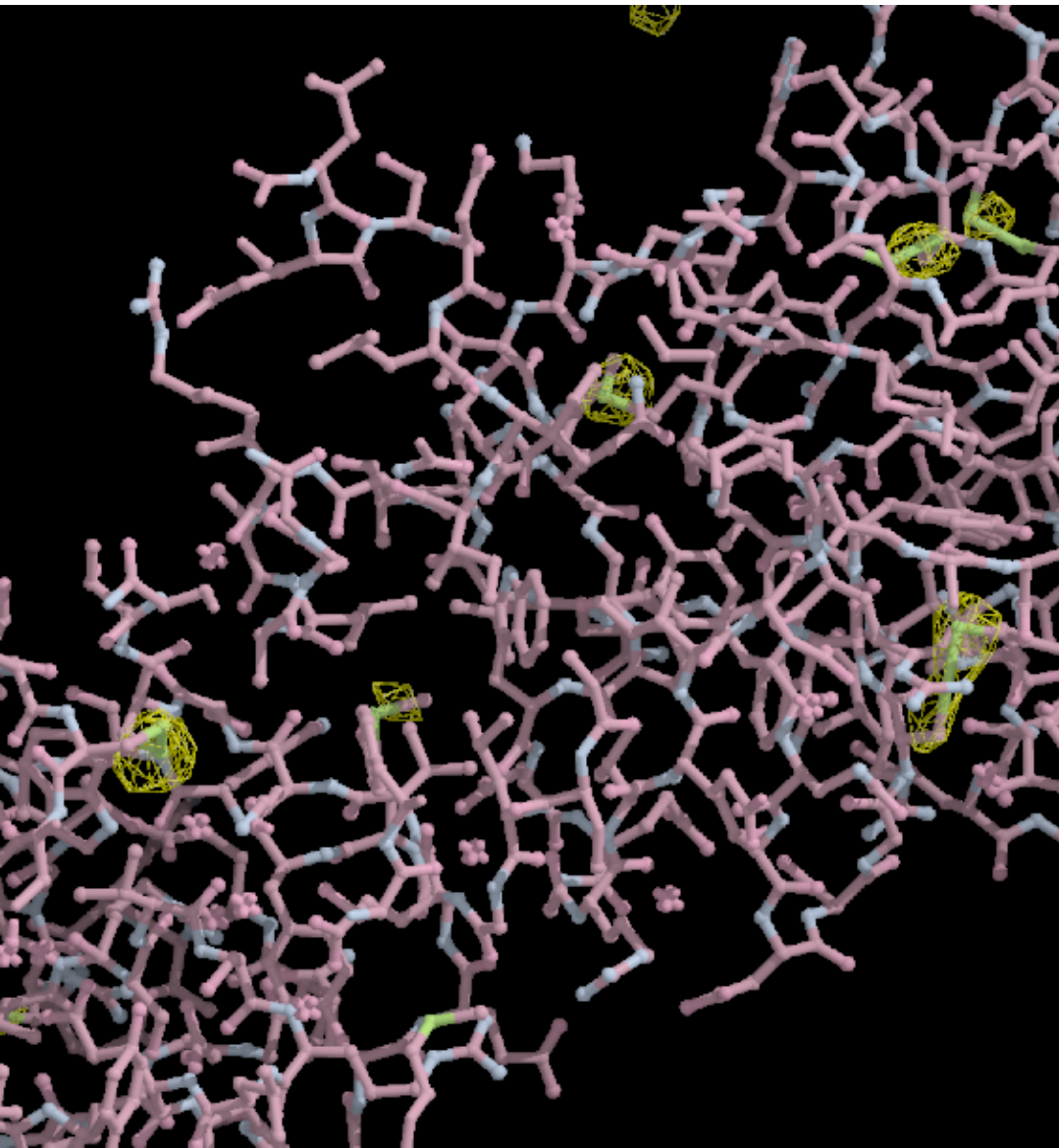
Correlation of observed and **sub-structure** anomalous differences

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

$CC_{ano}$  indicates how much of each anomalous difference is useful (on average)



# Anomalous signal: peak height at coordinates of anomalously-scattering atoms



$$S_{ano} = \frac{\langle \rho_{ano}(x_j) \rangle}{\langle \rho_{ano}^2 \rangle^{1/2}}$$

Typical values of  $S_{ano}$  for solved datasets: 10-20

*Anomalous difference Fourier with observed data and model phases*

## How big will my anomalous signal be?

Expected value of  
anomalous signal  $S_{ano}$

$$\langle S_{ano} \rangle = CC_{ano} \frac{N_{refl}^{1/2}}{f^{1/2} n_{sites}^{1/2}}$$

$f$  is 2<sup>nd</sup> moment of the  
anomalous scattering factor

$$f = \frac{\langle (f^h)^2 \rangle}{\langle f^h \rangle^2}$$

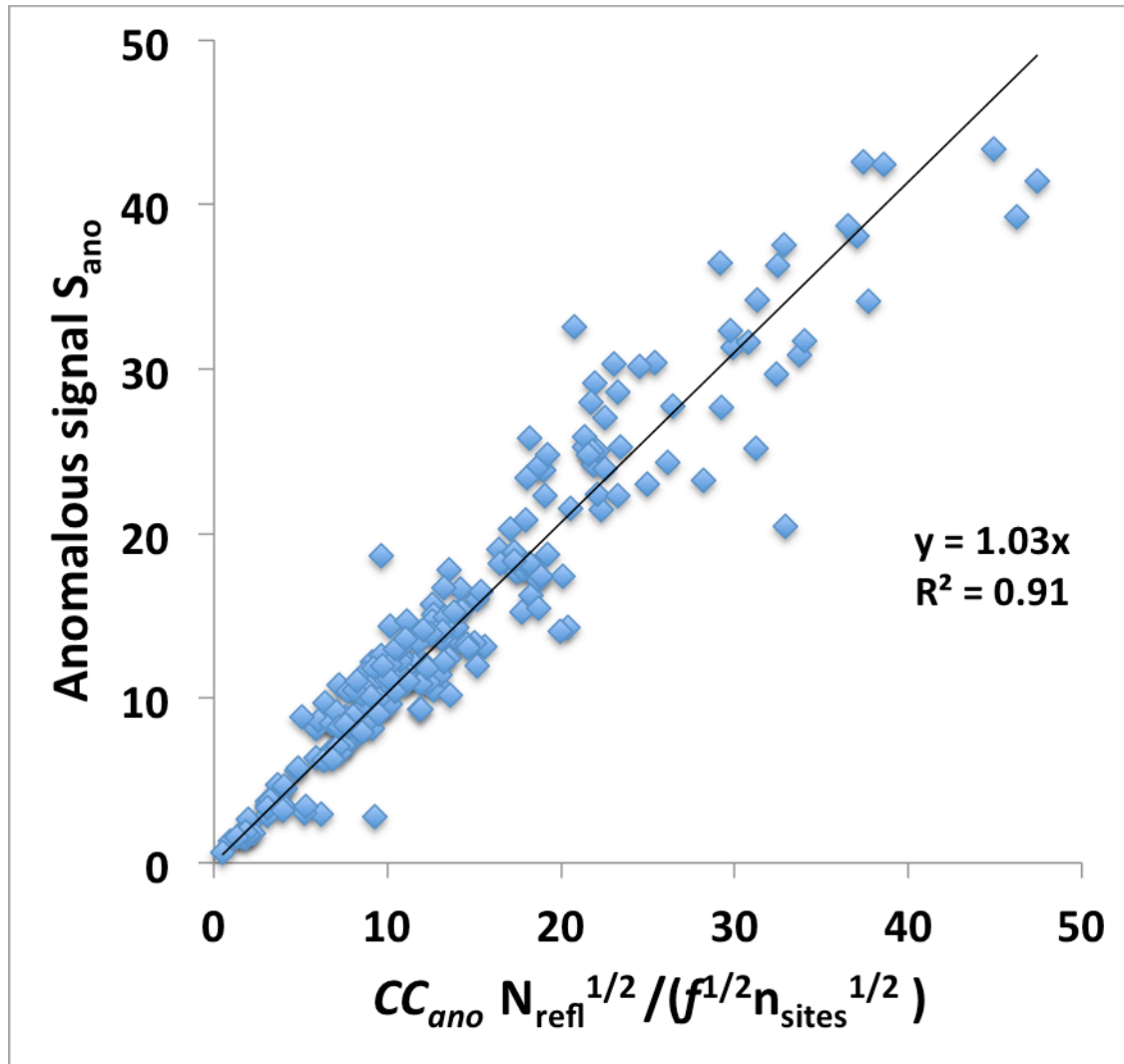
Anomalous scattering factor

$$f^h \equiv f'' e^{-B (\sin^2 \theta_h / \lambda^2)}$$

Perfect data (20,000 reflections, 8 sites):  $S_{ano} = (20000/8)^{1/2} = 50$

Good data (overall  $CC_{ano} = 0.36$   $f = 2.0$ ):  $S_{ano} = 12.6$

# Checking our simple model for anomalous signal



$$\langle S_{ano} \rangle = CC_{ano} \frac{N_{refl}^{1/2}}{f^{1/2} n_{sites}^{1/2}}$$

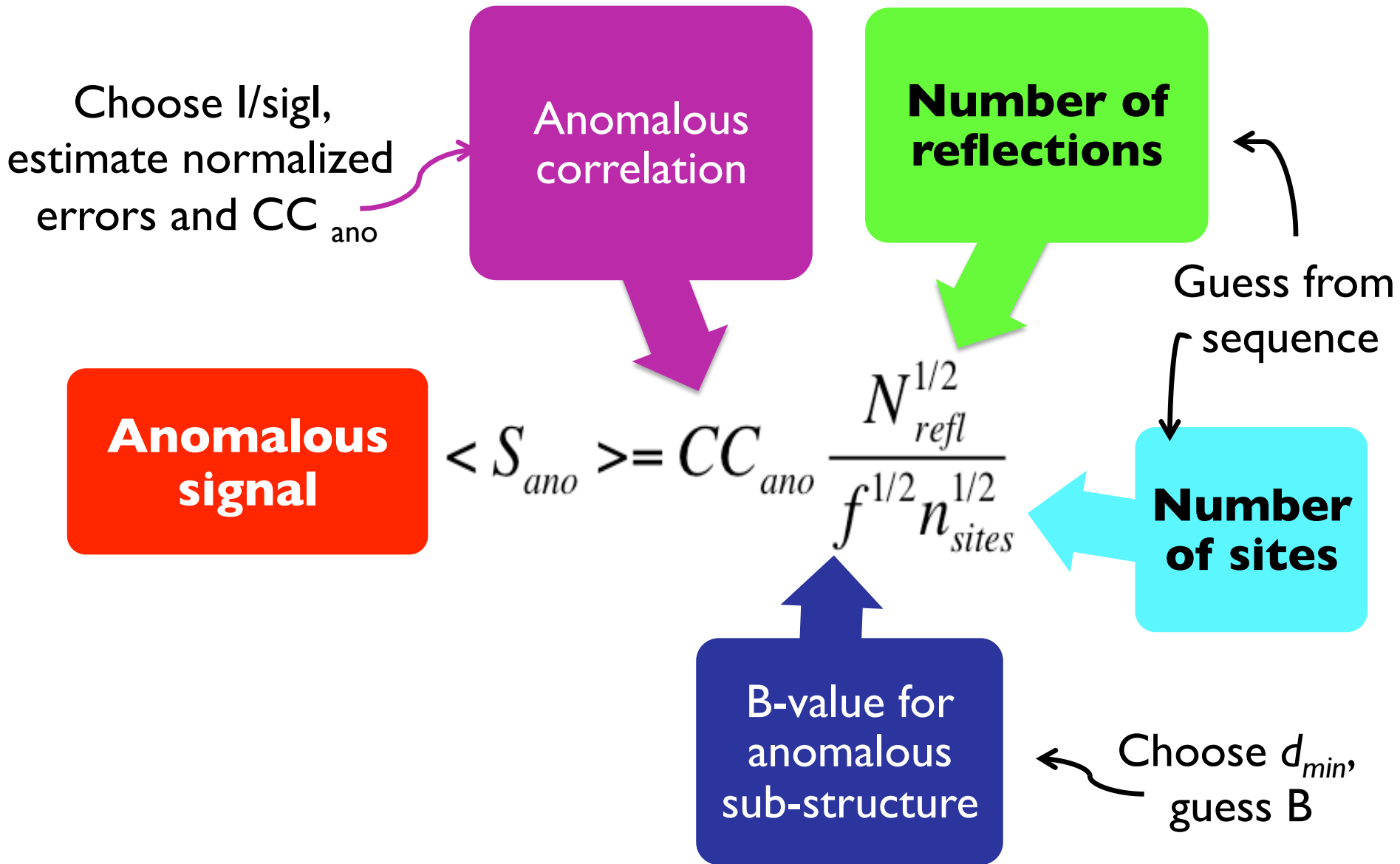
$CC_{ano}$ : Correlation of anomalous differences with model differences

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

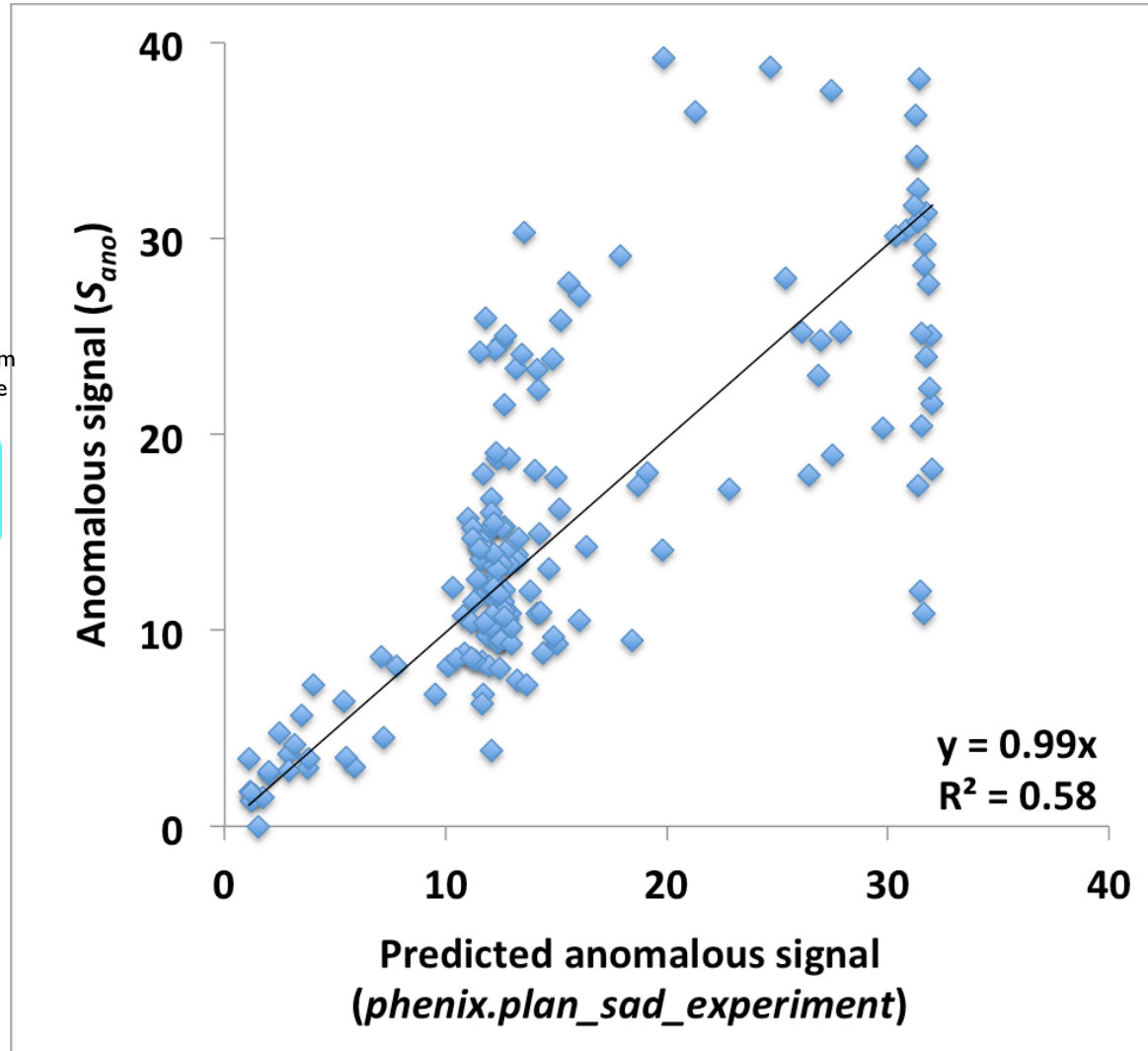
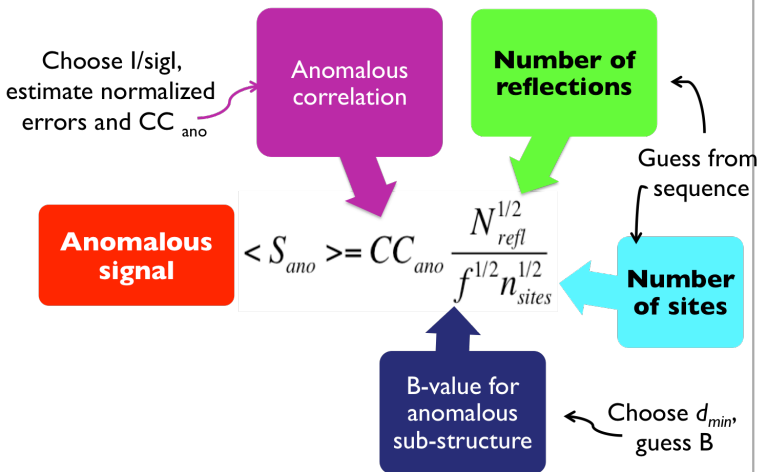
218 SAD datasets 1.2 – 4.5 Å

# *phenix.plan\_sad\_experiment*

Design an experiment that will give you enough anomalous signal



# Estimating the anomalous signal before collecting the data



# Optimizing scaling and merging of SAD data

*(phenix.scale\_and\_merge)*

# Why $F^+, F^-$ differ from one crystal to another

Crystal 1  
 $F^+, F^-$

Errors in measurement ( $\sigma_{\text{obs}}$ )

Crystals really are different  
( $\sigma_{\text{crystal}}$ )

Crystal 2  
 $F^+, F^-$

## Optimizing estimates of $F^+$ , $F^-$

Crystal 1  
 $F^+$ ,  $F^-$

Local scaling to reduce  
systematic errors

Use of  $\sigma_{\text{crystal}}$  in weighting

Crystal 2  
 $F^+$ ,  $F^-$



# Applying inter-dataset variances in weighting

Crystal I

$\Delta_{\text{ano}}$

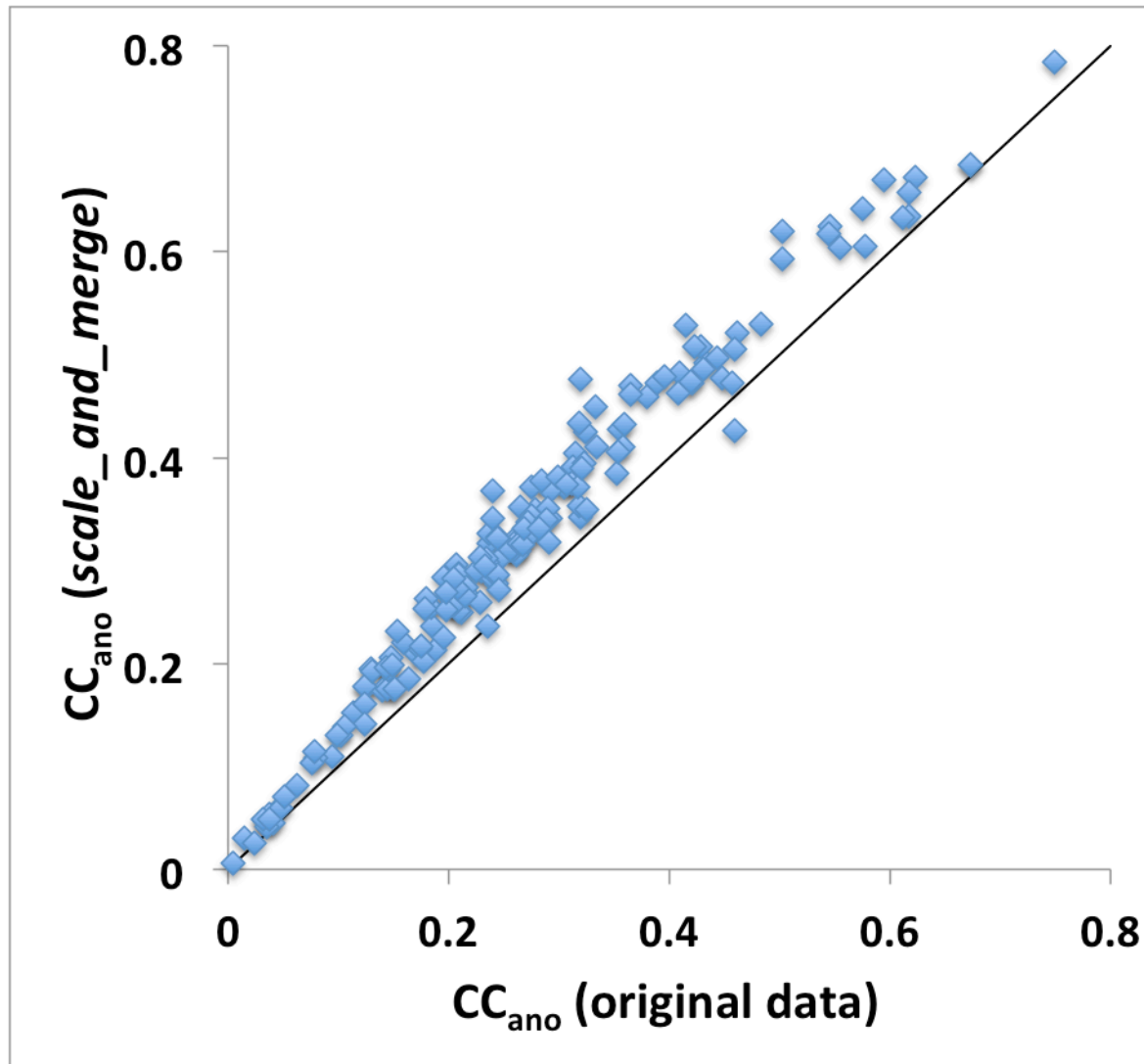
Weighting for data from an individual crystal:

$$\sigma^2_{\text{total}} \approx \sigma^2_{\text{obs}} + \sigma^2_{\text{crystal}}$$

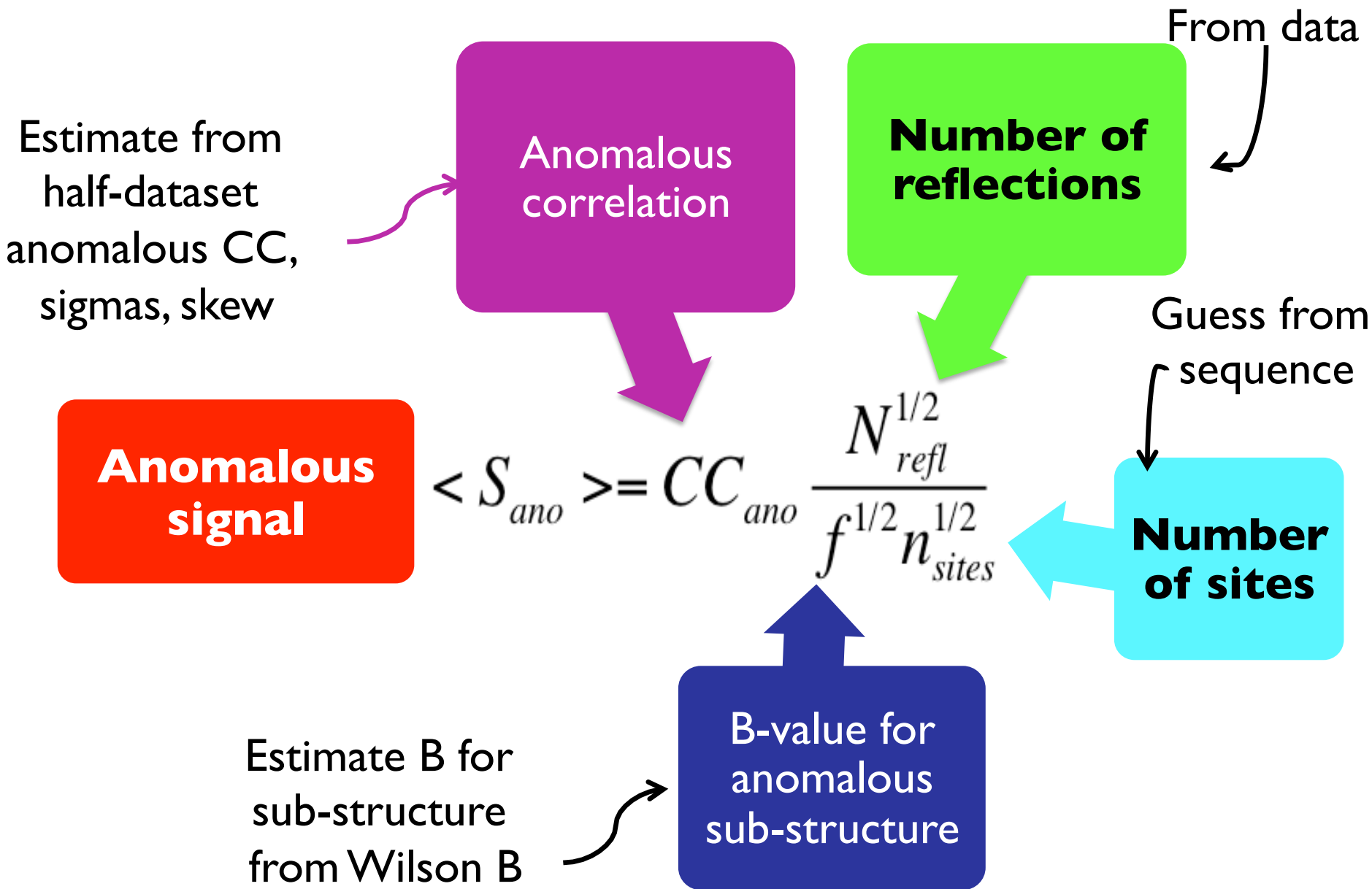
Average of all crystals

$\Delta^{\text{AVG}}$

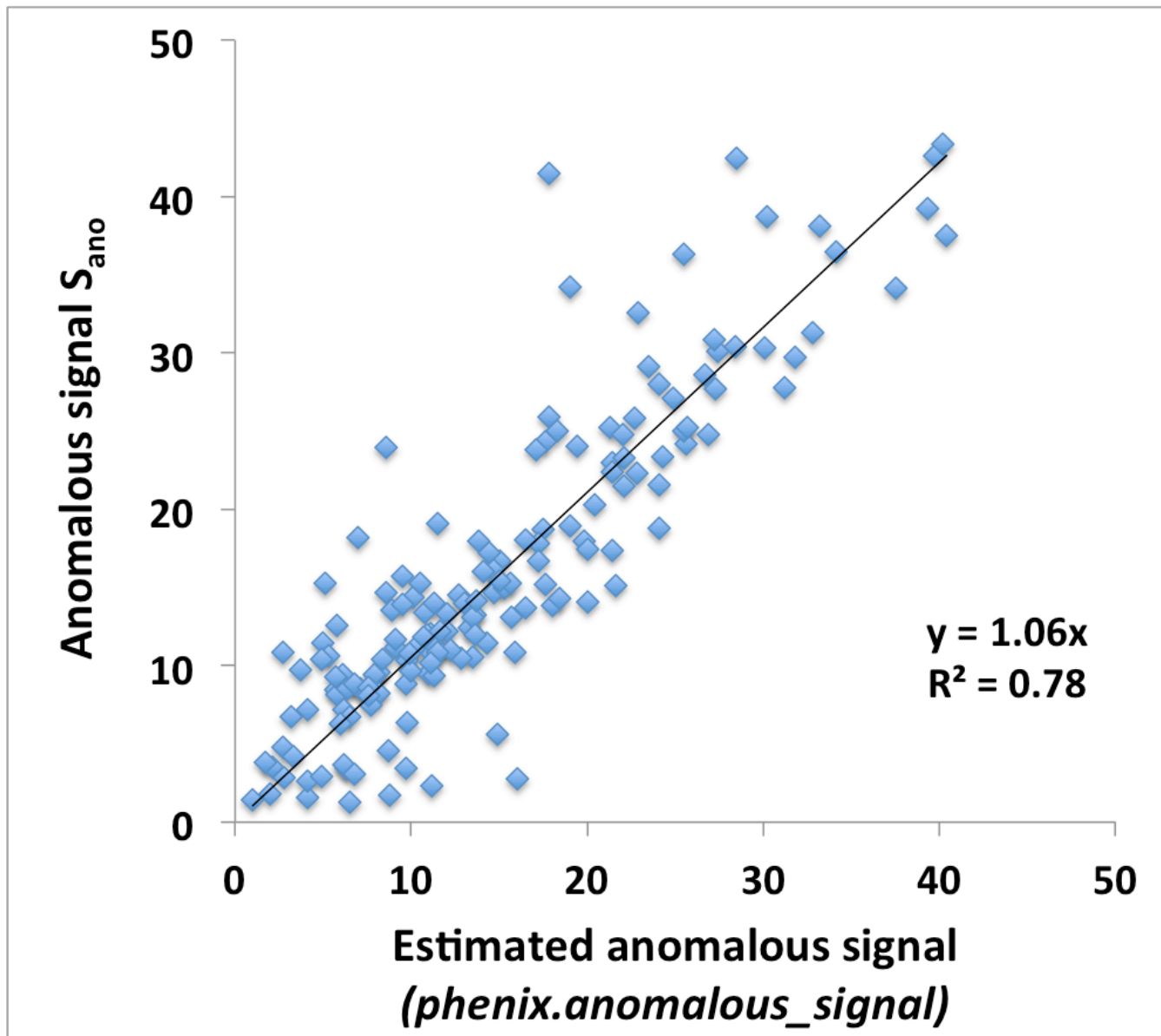
# Improvement in anomalous correlation using local scaling in *phenix.scale\_and\_merge*



# Estimating the anomalous signal after collecting the data



# Estimating the anomalous signal after collecting the data



## **Finding the anomalous sub-structure with the SAD likelihood function**

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

Most powerful source of information about the  
sub-structure before phases are known

# 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 an anomalous difference Fourier—but 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.

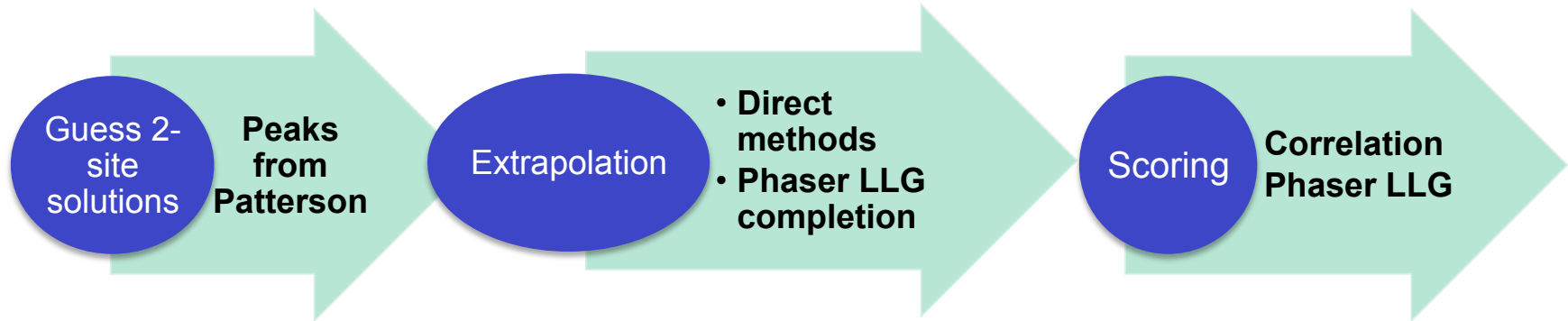
# LLG sub-structure searches in HySS

## Test cases

164 SAD datasets from PDB (largely JCSG MAD data)

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

# Finding anomalous substructure with LLG completion



- **Range of resolution**  
**Variable number of Patterson solutions**

**Adjustable  
LLGC\_SIGMA  
(cut-off for peak height)**

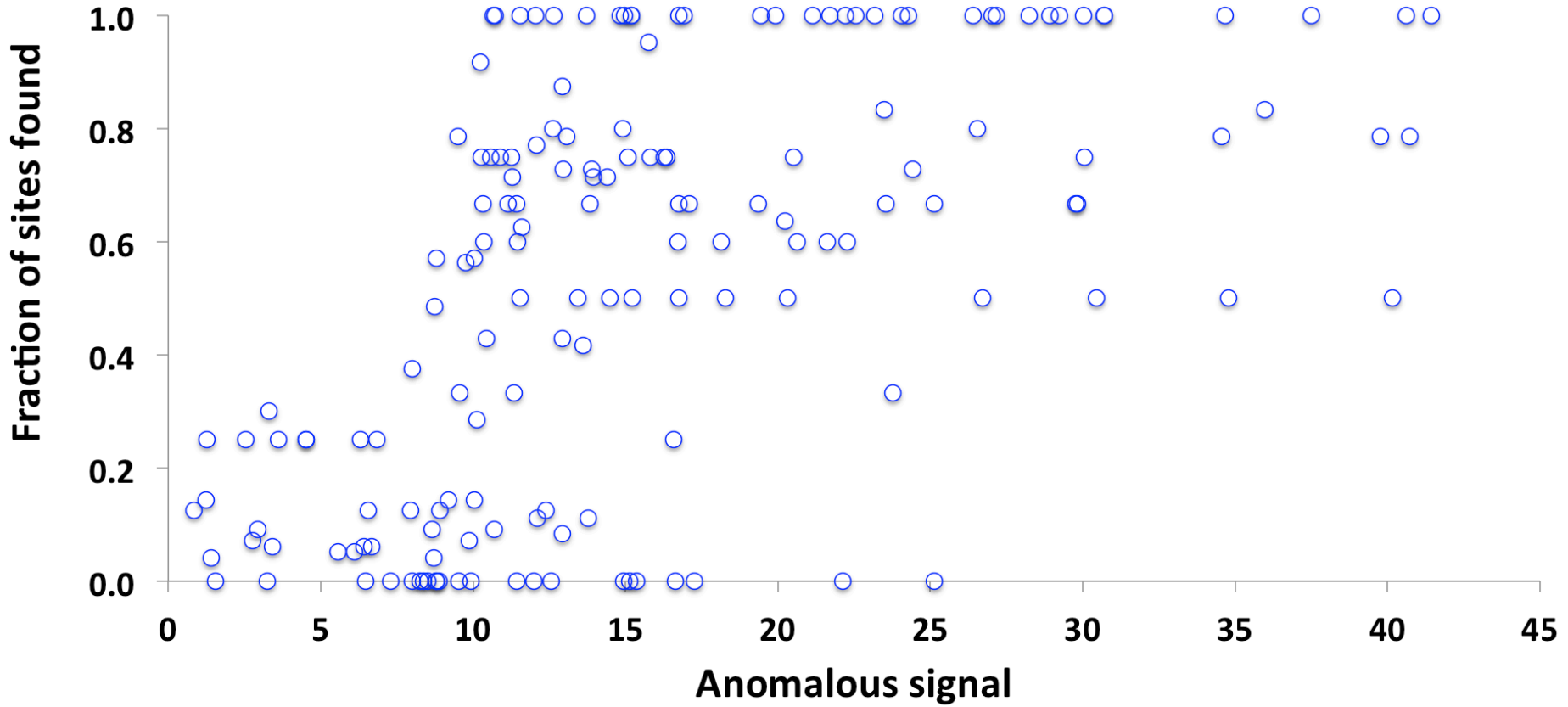
**Use LLG score to  
compare solutions**

**Terminate early if same  
solution found several  
times**

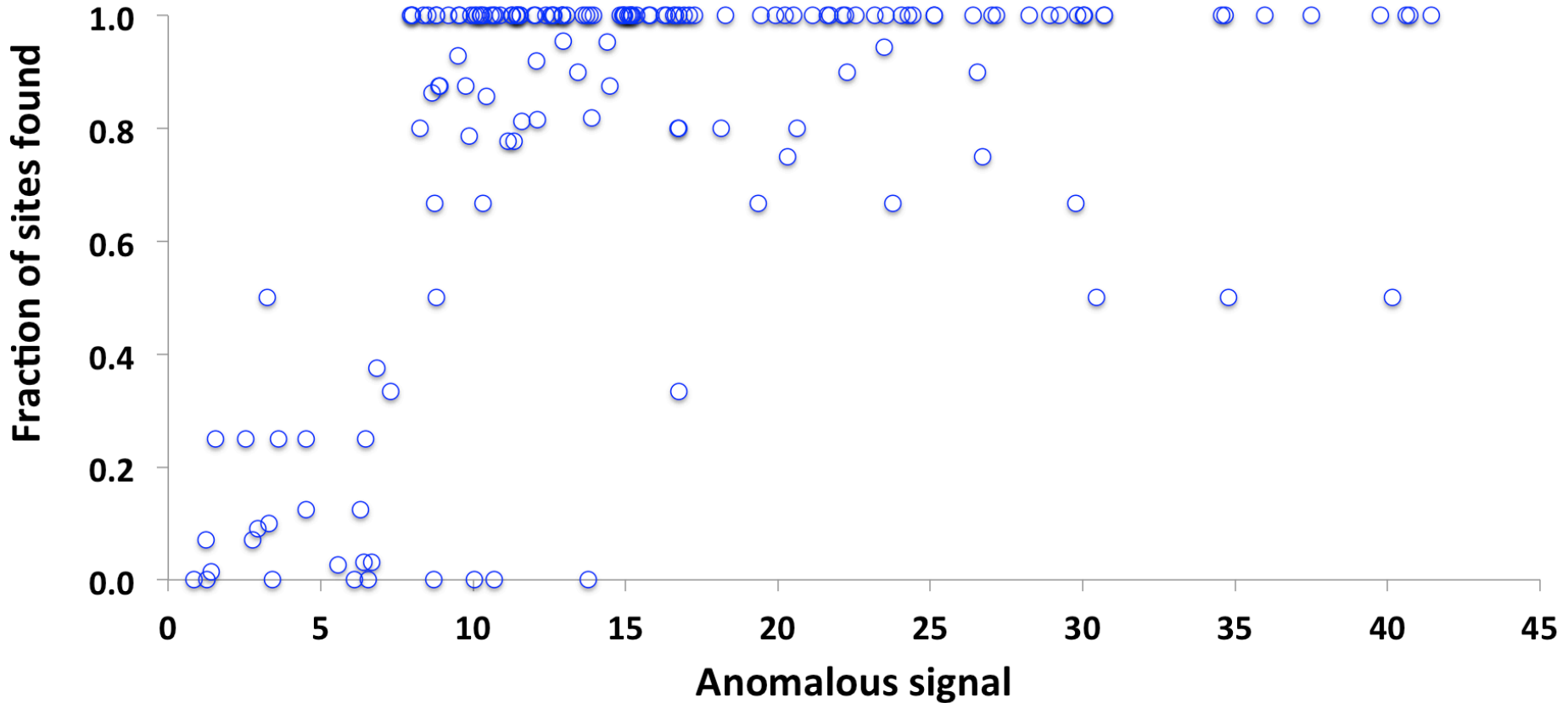
**Run quick direct  
methods first**



# Dual Space Sub-structure Completion

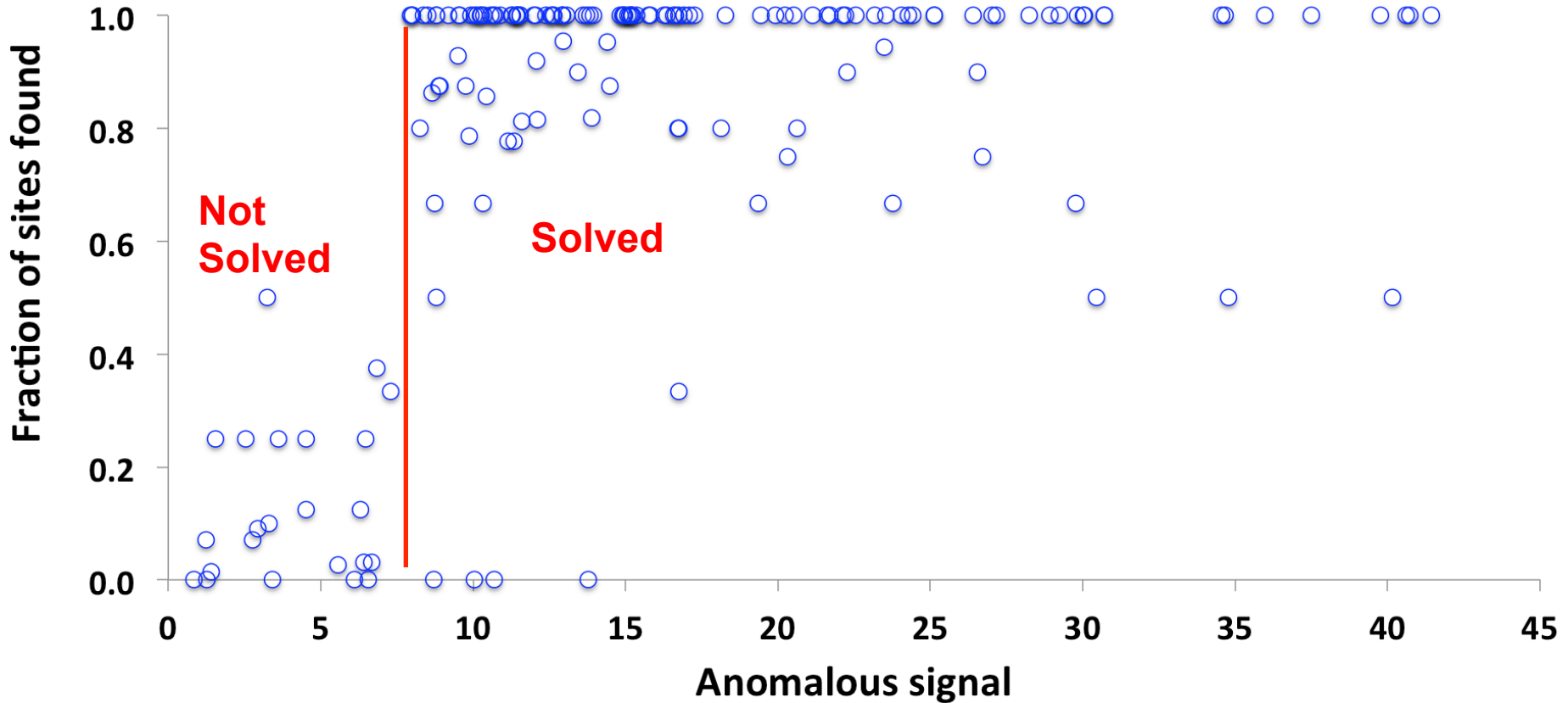


# LLG Sub-structure Search



***Bunkóczy et al., Nature Methods 12, 127–130 (2015).***

# Anomalous signal indicates if a dataset can be solved



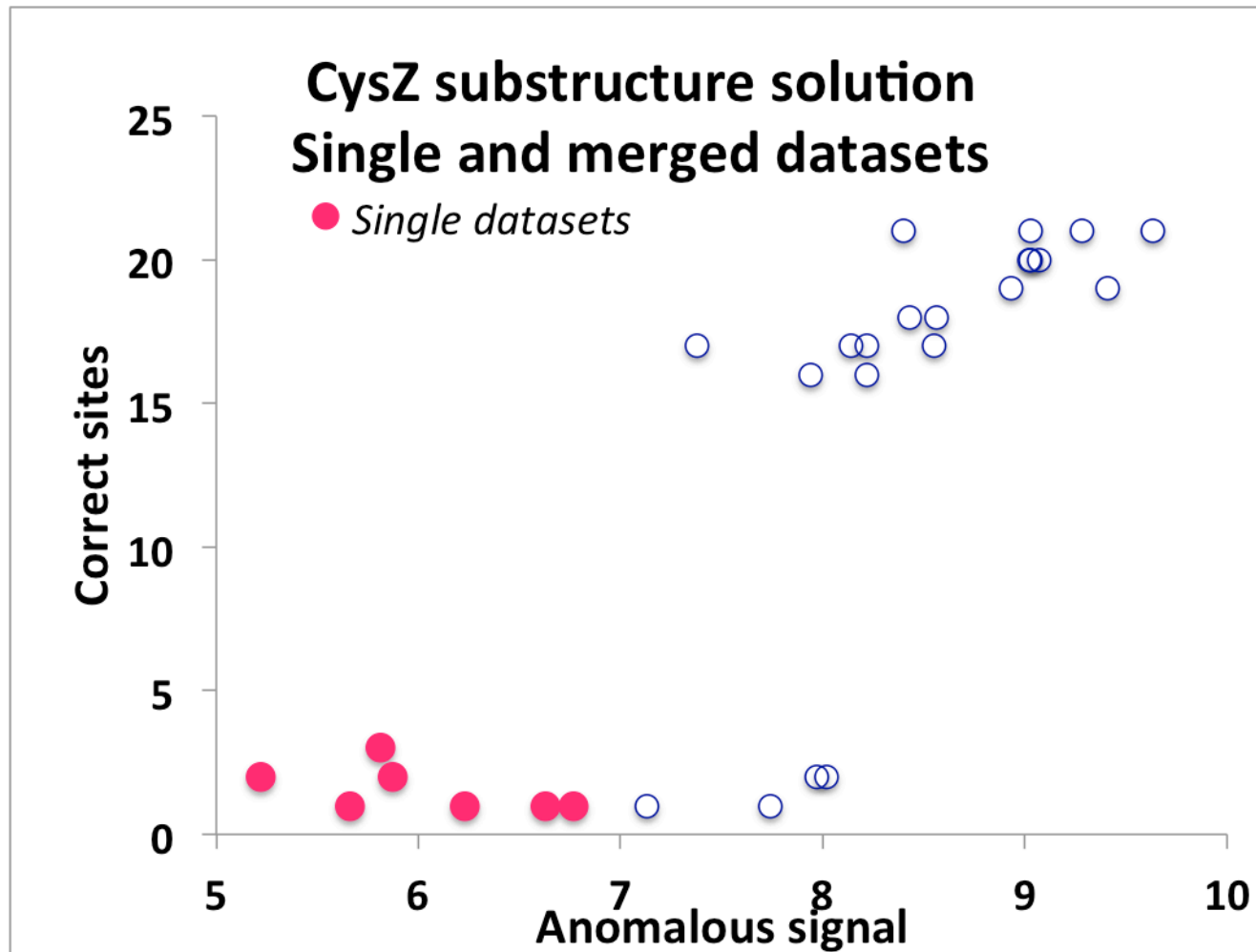
# **CysZ multi-crystal sulfur-SAD data**

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

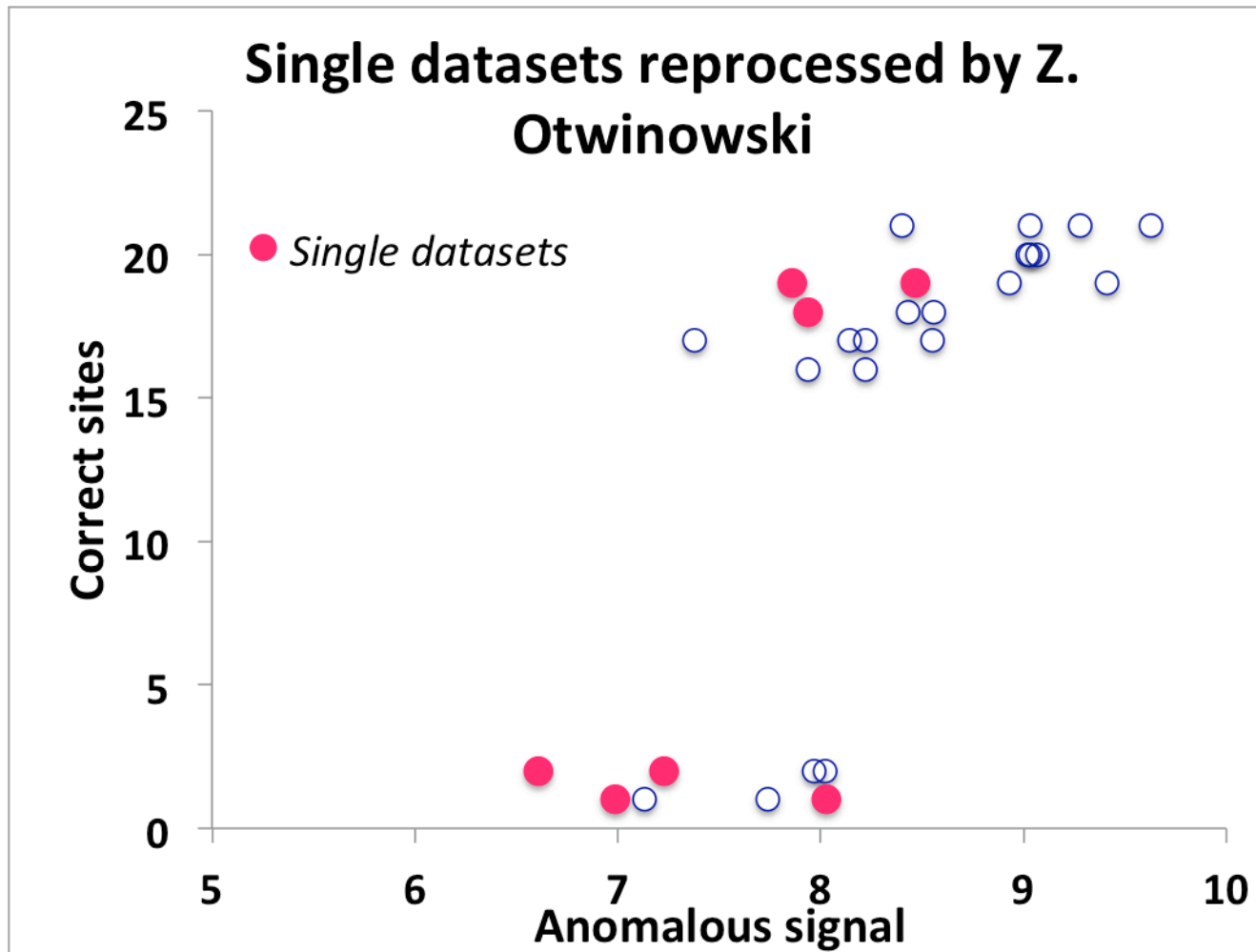
**Data from 7 crystals collected at wavelength of 1.74 Å to resolution of 2.3 Å**

**Can anomalous signal tell us which merged datasets will be solved?**

# CysZ multi-crystal sulfur-SAD data

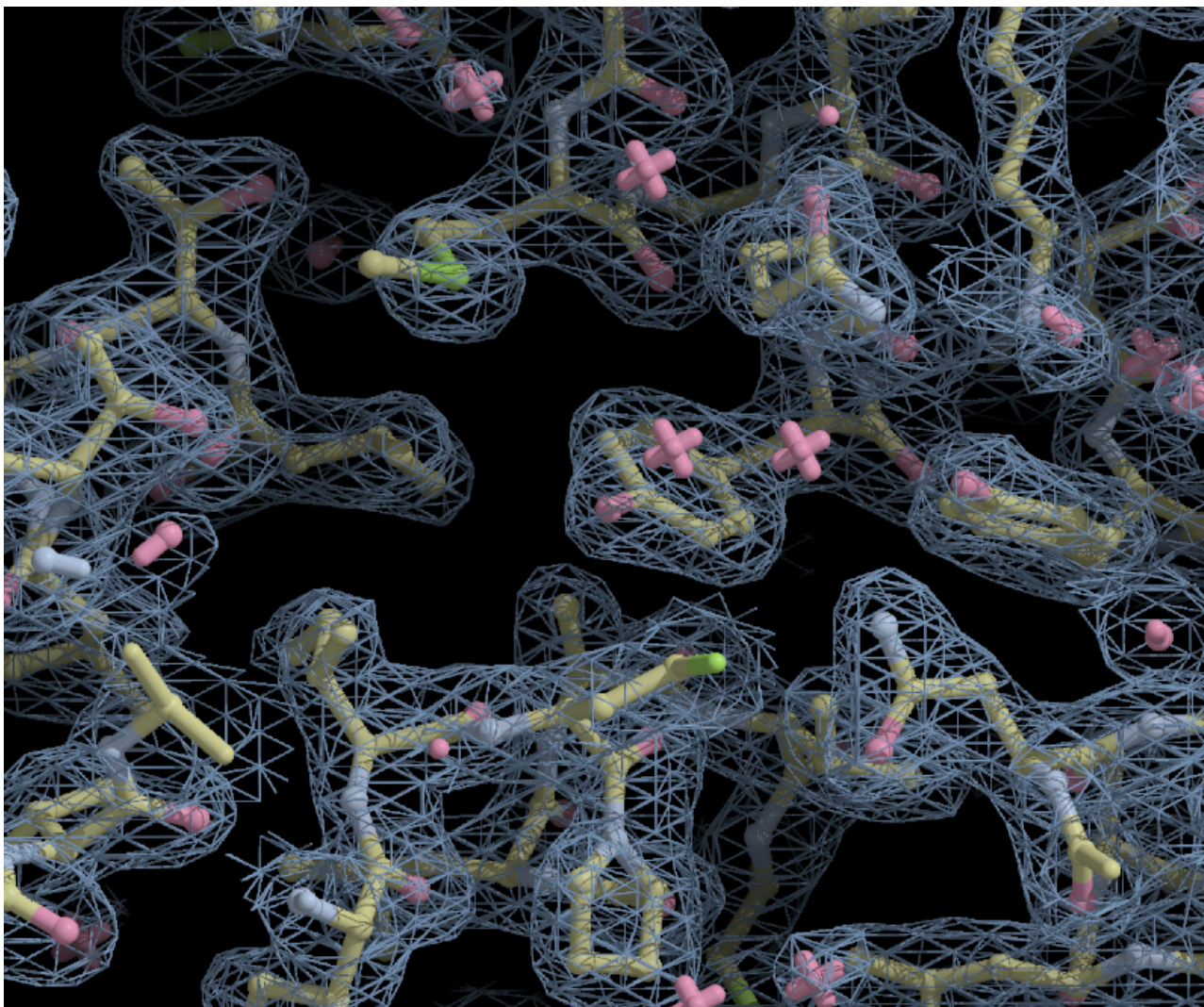


# CysZ multi-crystal sulfur-SAD data



# CysZ single-crystal sulfur-SAD data

**Crystal 6** *AutoSol R/Rfree=0.24/0.27*



## **Choosing the hand of the sub- structure**

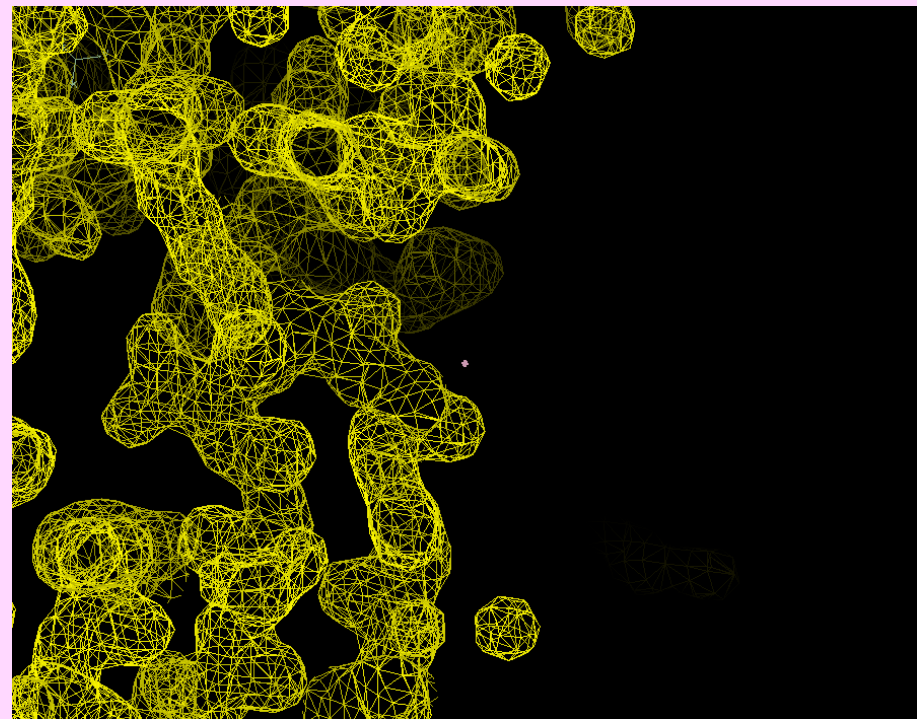
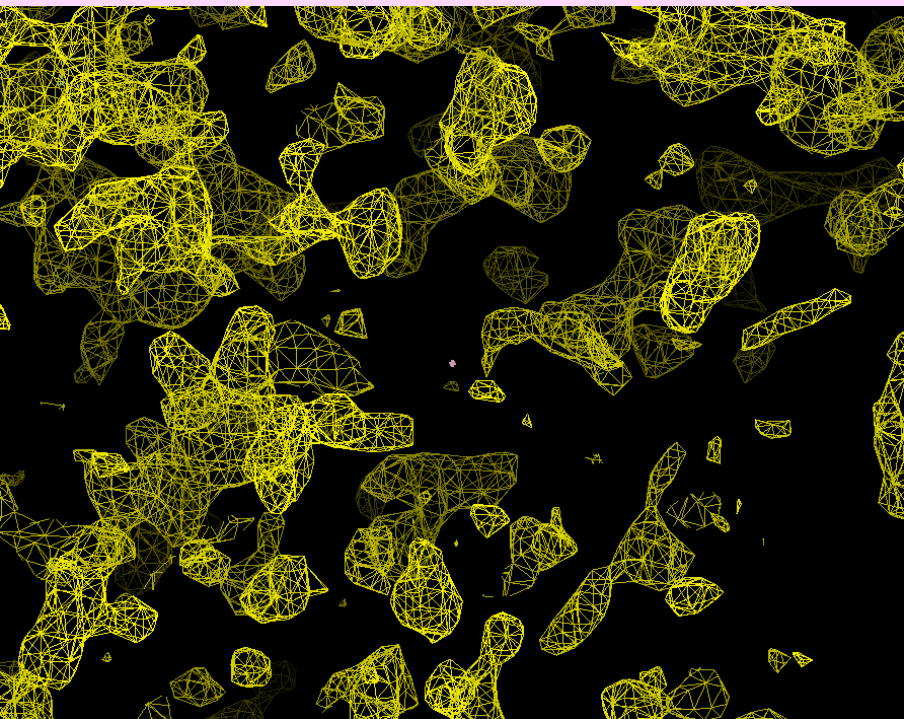
*Decision-making based on map quality*

**(*phenix.autosol*)**



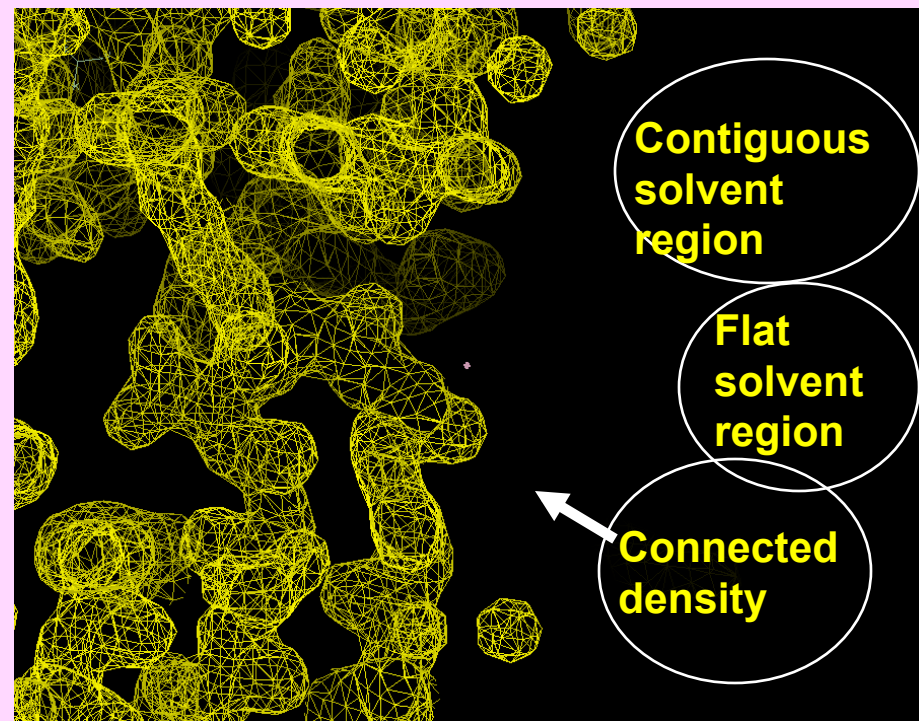
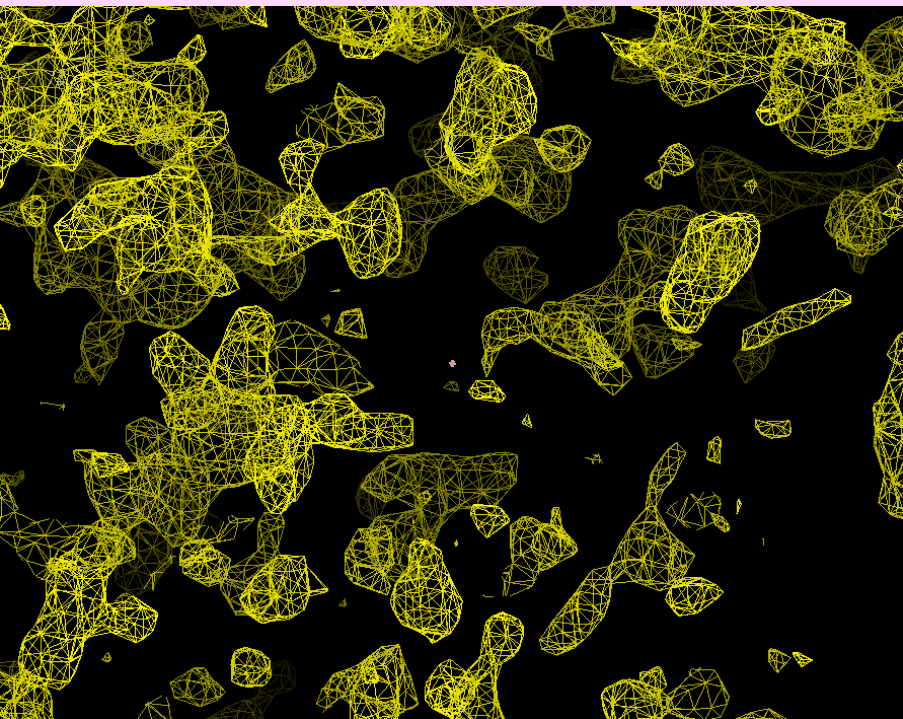
*Deciding what is good:*  
**Measures of the quality of an electron-density map:**

**Which solution is best?**  
**Are we on the right track?**

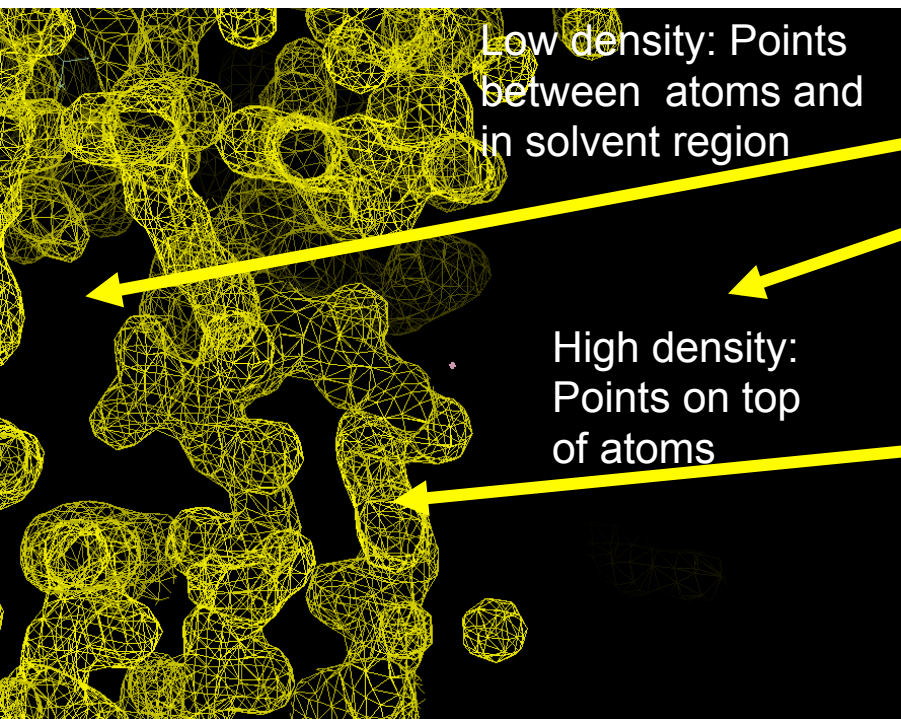


*Deciding what is good:*  
**Measures of the quality of an electron-density map:**

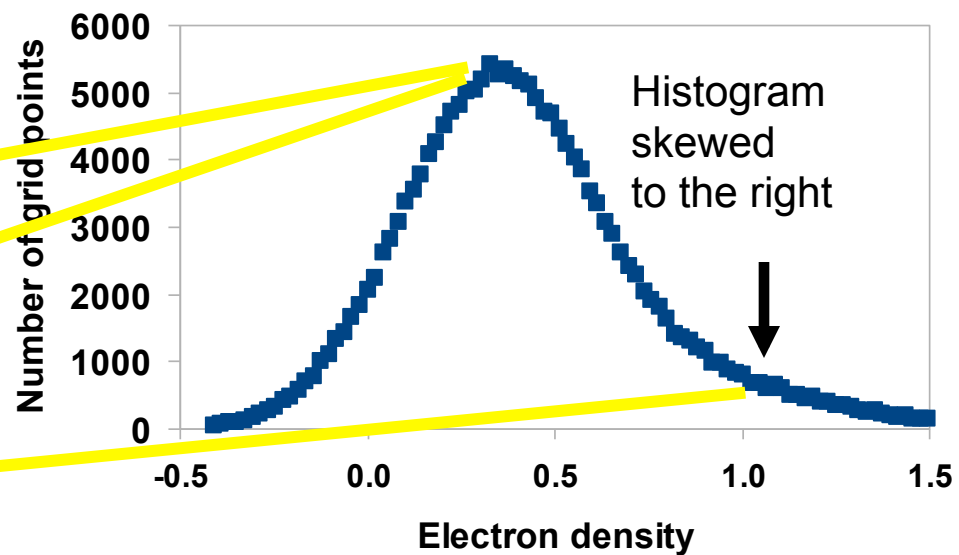
**Which solution is best?**  
**Are we on the right track?**



Histogram of electron density values has a positive “skew”

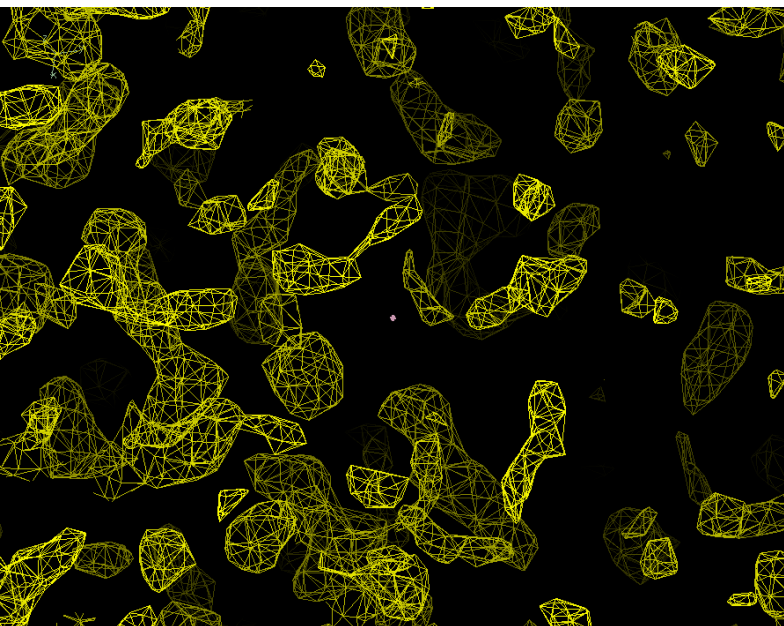


Typical histogram of electron density

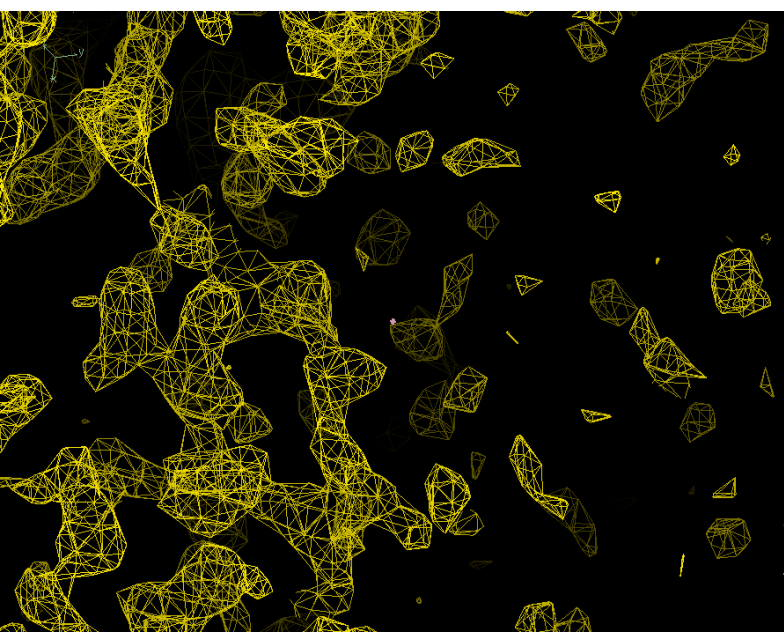




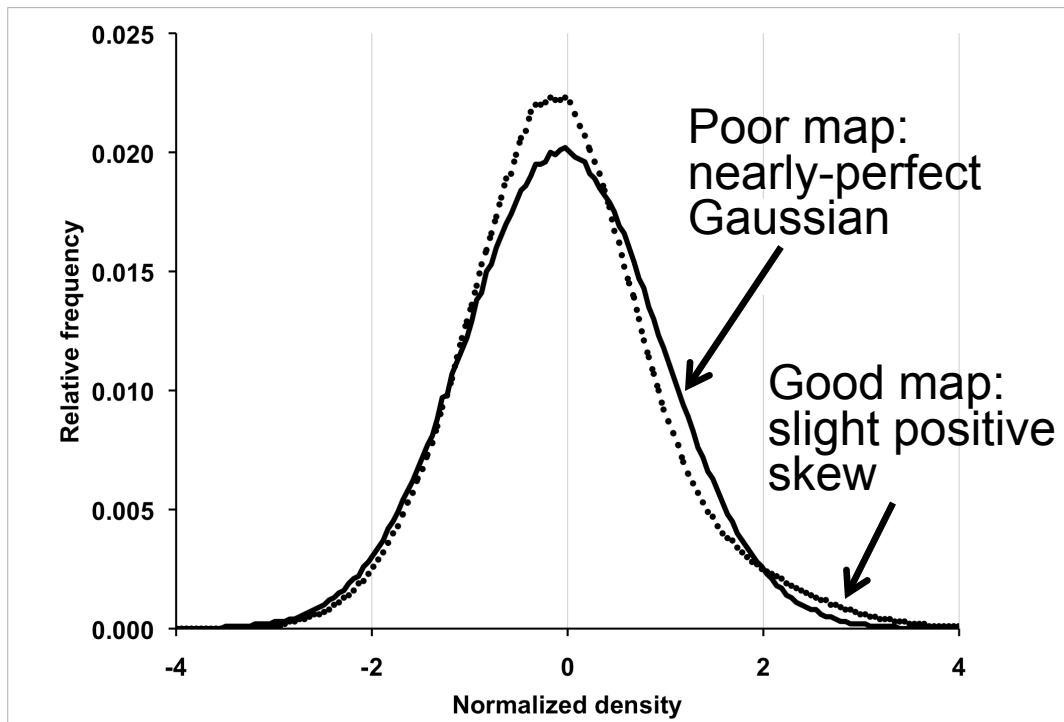
# ***Skew of electron density differentiates poor and good maps – even when the difference is barely visible***



Poor map  
(inverse hand)



Good map



# Evaluating electron density maps

<i>Basis</i>	<i>Good map</i>	<i>Random map</i>
Skew of density (Podjarny, 1977)	Highly skewed (very positive at positions of atoms, zero elsewhere)	Gaussian histogram
Connectivity of regions of high density (Baker, Krukowski, & Agard, 1993)	A few connected regions can trace entire molecule	Many very short connected regions
Correlation of local rms densities (Terwilliger, 1999)	Neighboring regions in map have similar rms densities	Map has uniform rms density
R-factor in 1 <sup>st</sup> cycle of density modification (Cowtan, 1996)	Low R-factor	High R-factor

*How well does the skew reflect map quality?*

Create real maps

Score the maps based on skew

Compare the scores with the actual quality of the maps

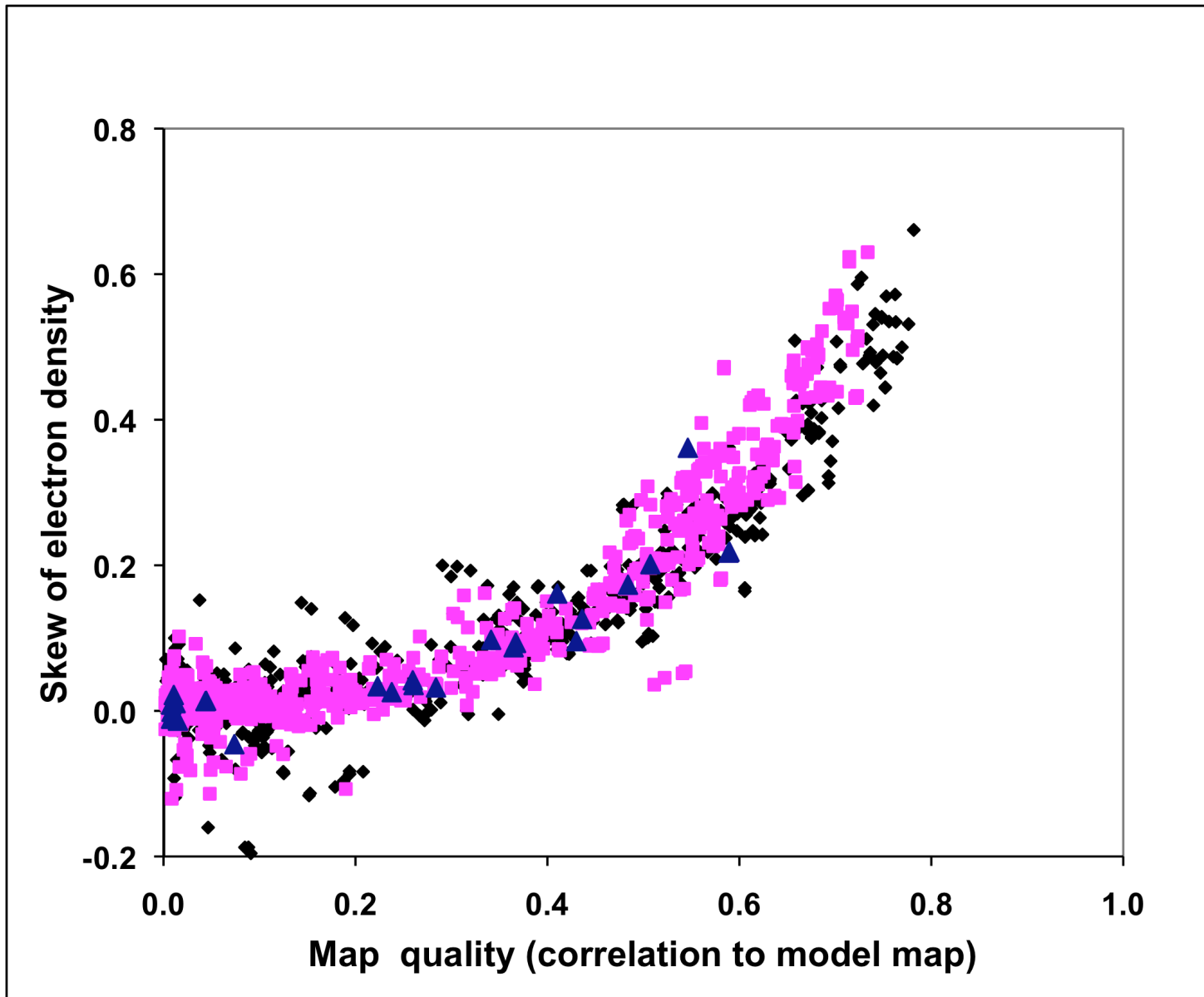
## Creating real maps

247 MAD, SAD, MIR datasets with final model available  
(*Phenix* library and JCSG publicly-available data)

Run *phenix.autosol* on each dataset.

Calculate maps for each solution considered  
(opposing hands, additional sites, including various derivatives  
for MIR)

# Skew of electron density – positive skew of density values



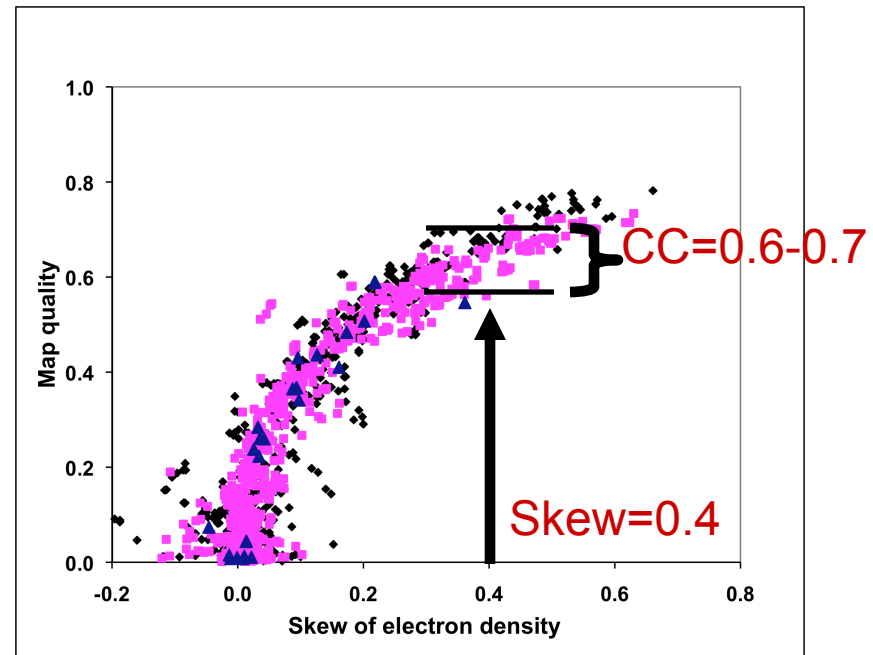
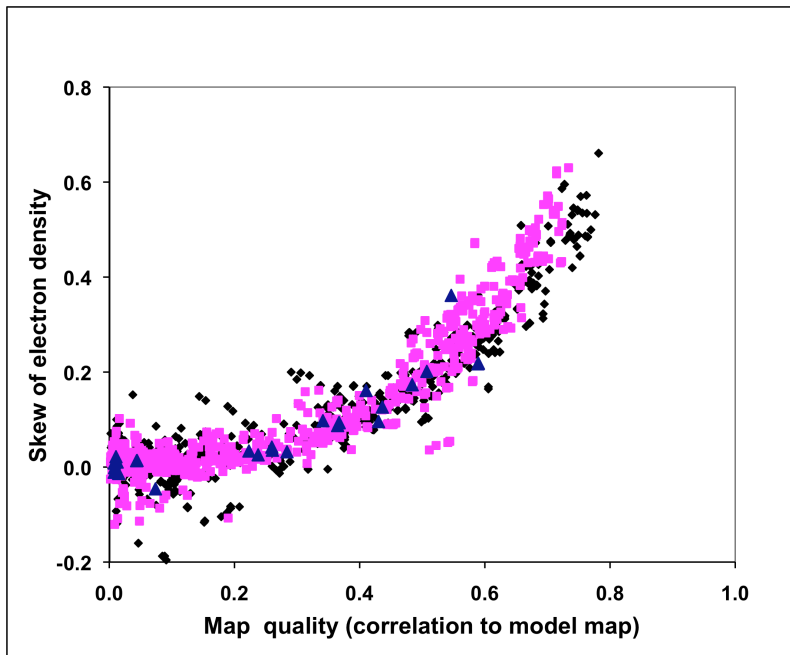


# Using scoring criteria to estimate the quality of a map

Skew depends on map quality



Estimate map quality from skew



## Estimated map quality in practice

Evaluating solutions to a 2-wavelength MAD experiment  
(JCSG Tm3681, 1VPM, SeMet 1.6 Å data)

Data for HYSS	Sites	Estimated CC $\pm 2SD$	Actual CC
Peak	12	$0.73 \pm 0.04$	0.72 ←
Peak (inverse hand)	12	$0.11 \pm 0.43$	0.04
$F_A$	12	$0.73 \pm 0.03$	0.72
$F_A$ (inverse)	12	$0.11 \pm 0.42$	0.04
Sites from diff Fourier	9	$0.70 \pm 0.17$	0.69

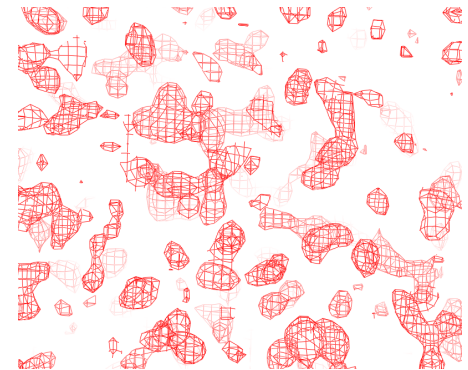
## Density modification

*Improving phase quality by including expectations  
about the map*

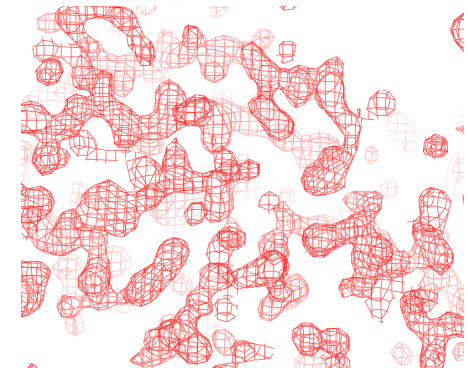
**(*phenix.autosol; phenix.autobuild*)**

# Statistical density modification

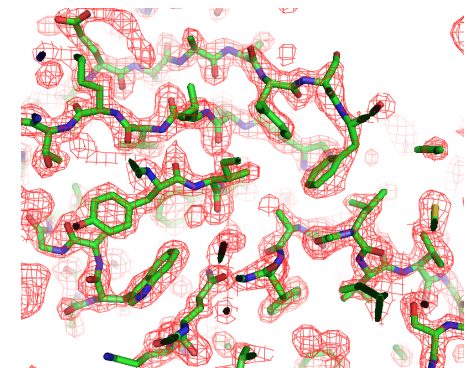
- Principle: phase probability information from probability of the map and from experiment:
- $P(\phi) = P_{\text{map probability}}(\phi) P_{\text{experiment}}(\phi)$
- “Phases that lead to a believable map are more probable than those that do not”
- **A believable map is a map that has...**
  - a relatively flat solvent region
  - NCS (if appropriate)
  - A distribution of densities like those of model proteins
- **Method:**
  - calculate how map probability varies with electron density  $\rho$
  - deduce how map probability varies with phase  $\phi$
  - combine with experimental phase information



Experimental map

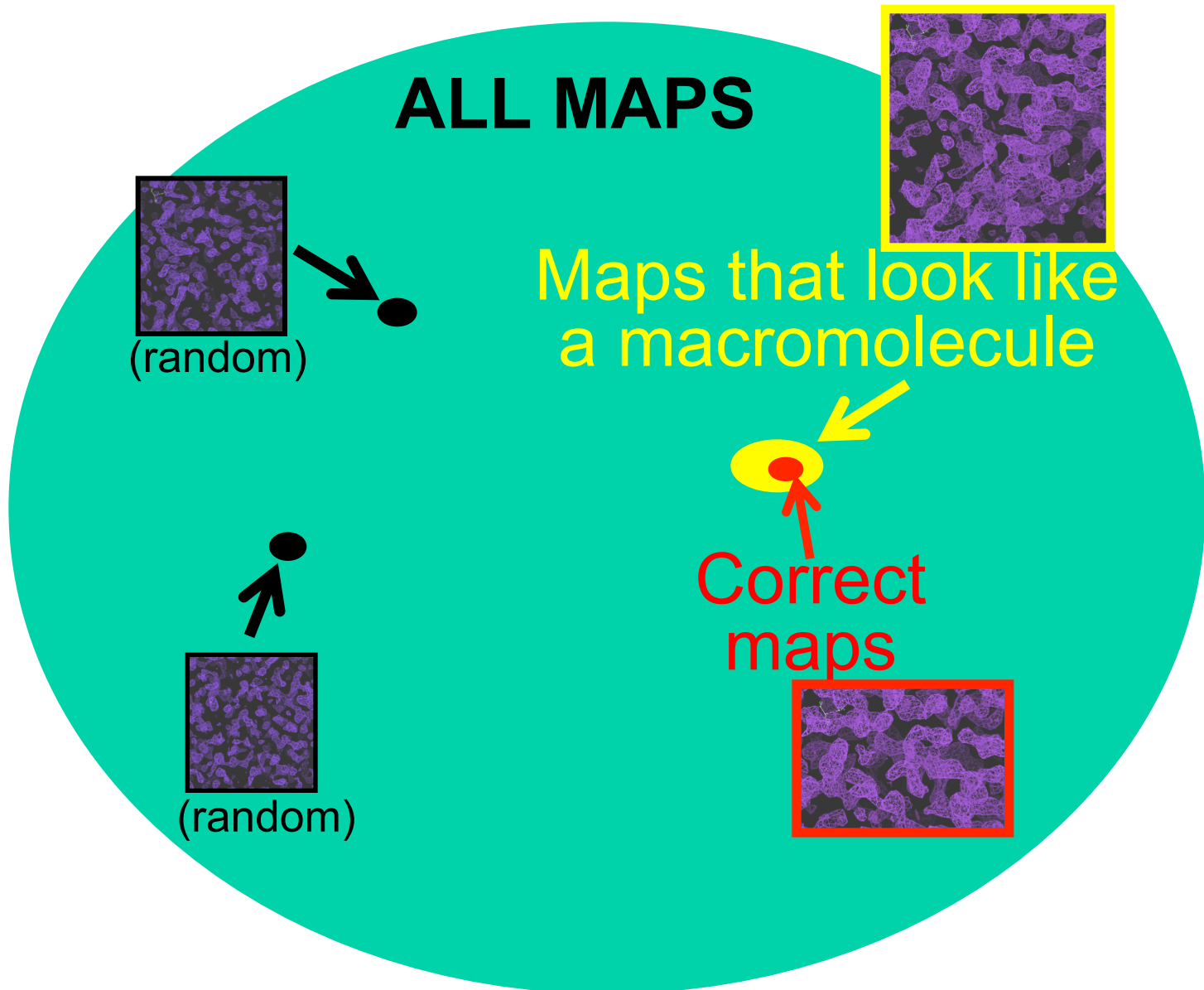


Density-modified



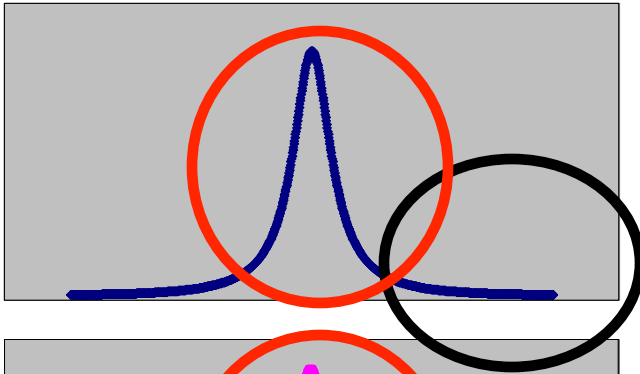
Interpreted

Maps that look like proteins are MUCH more likely to be correct than ones that do not

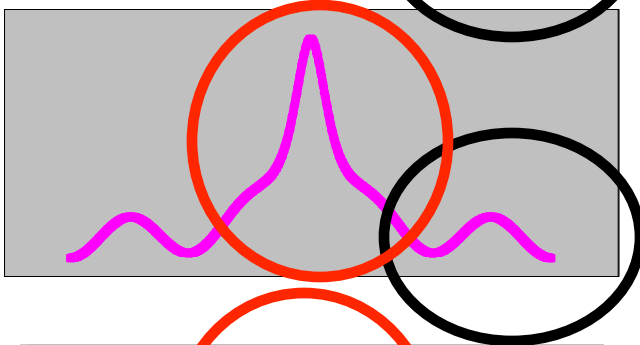


## Map probability phasing: Getting a new probability distribution for each phase given estimates of all others

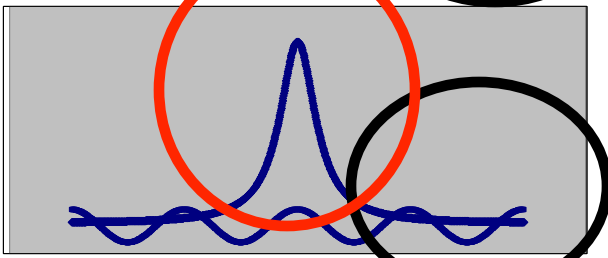
1. Identify expected features of map (flat far from center)
2. Calculate map with current estimates of all structure factors except one ( $k$ )
3. Test all possible phases  $\phi$  for structure factor  $k$  (for each phase, calculate new map including  $k$ )
4. Probability of phase  $\phi$  estimated from agreement of map with expectations
5. **Phase probability of reflection  $k$  from map is independent of starting phase probability because reflection  $k$  is omitted from the map**



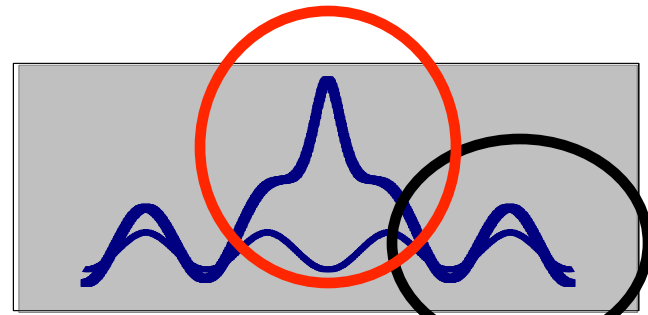
A function that is (relatively) flat far from the origin



Function calculated from estimates of all structure factors but one ( $k$ )



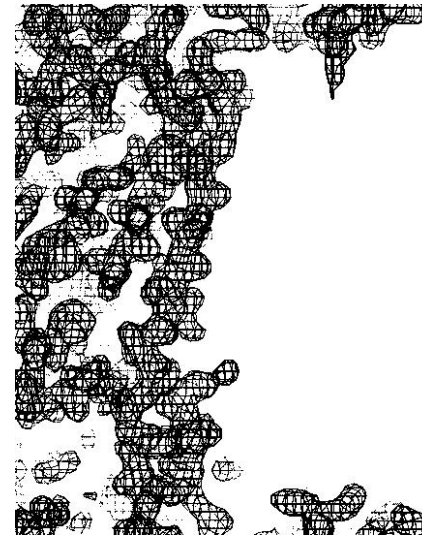
Test each possible phase of structure factor  $k$ .  $P(\phi)$  is high for phase that leads to flat region



## ***A map-probability function – allowing different weighting of information from different parts of the map***

Log-probability of the map is sum over all points in map of local log-probability

$$LL^{MAP}(\{\mathbf{F}_h\}) \approx \frac{N_{REF}}{V} \int_V LL(\rho(\mathbf{x}, \{\mathbf{F}_h\})) d^3\mathbf{x}$$



A map with a flat (blank) solvent region is a likely map

Local log-probability is believability of the value of electron density ( $\rho(\mathbf{x})$ ) found at this point

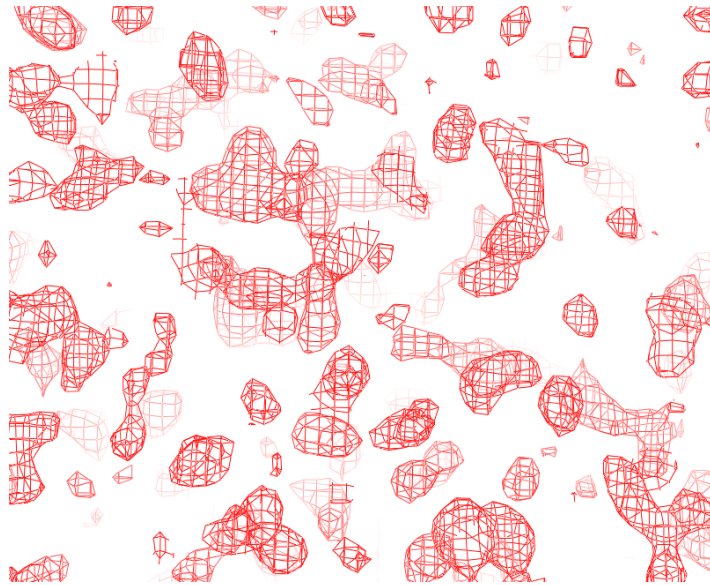
$$LL(\rho(\mathbf{x}, \{\mathbf{F}_h\})) = \ln[p(\rho(\mathbf{x})|PROT)p_{PROT}(\mathbf{x}) + p(\rho(\mathbf{x})|SOLV)p_{SOLV}(\mathbf{x})]$$

If the point is in the PROTEIN region, most values of electron density ( $\rho(\mathbf{x})$ ) are believable

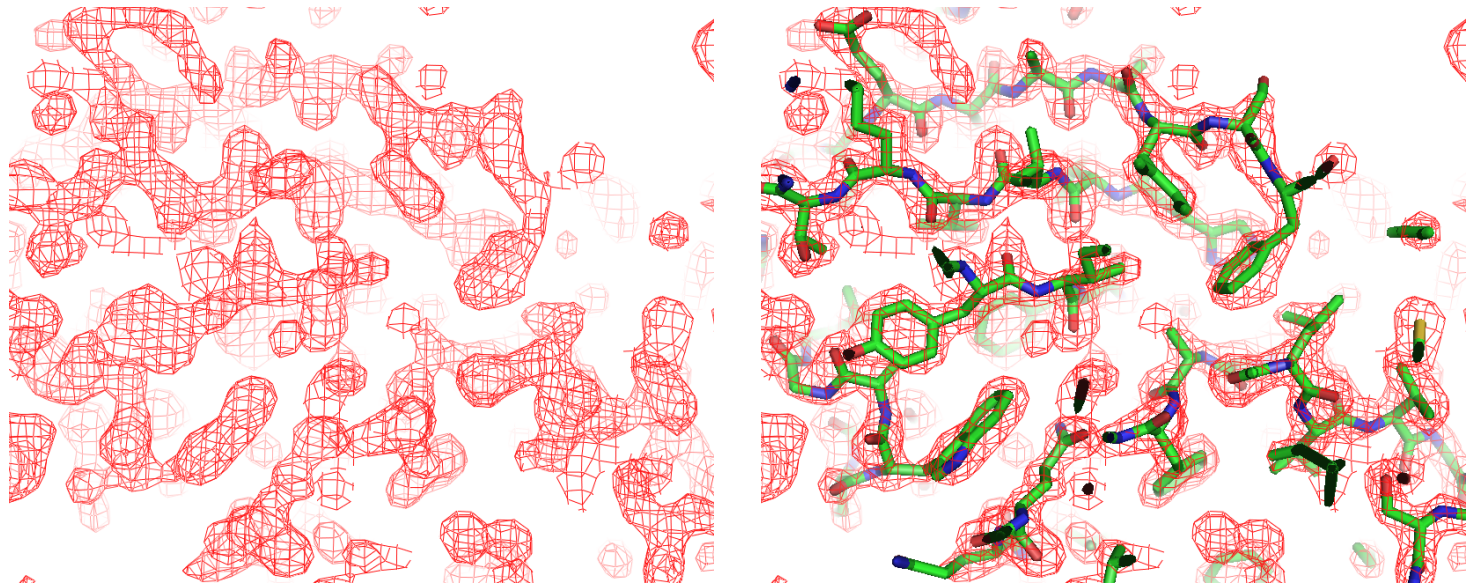
If the point is in the SOLVENT region, only values of electron density near zero are believable

# Statistical density modification (nsf-N SAD map , 2Å, no NCS, 50% solvent)

Phaser SAD map  
(CC=0.43)

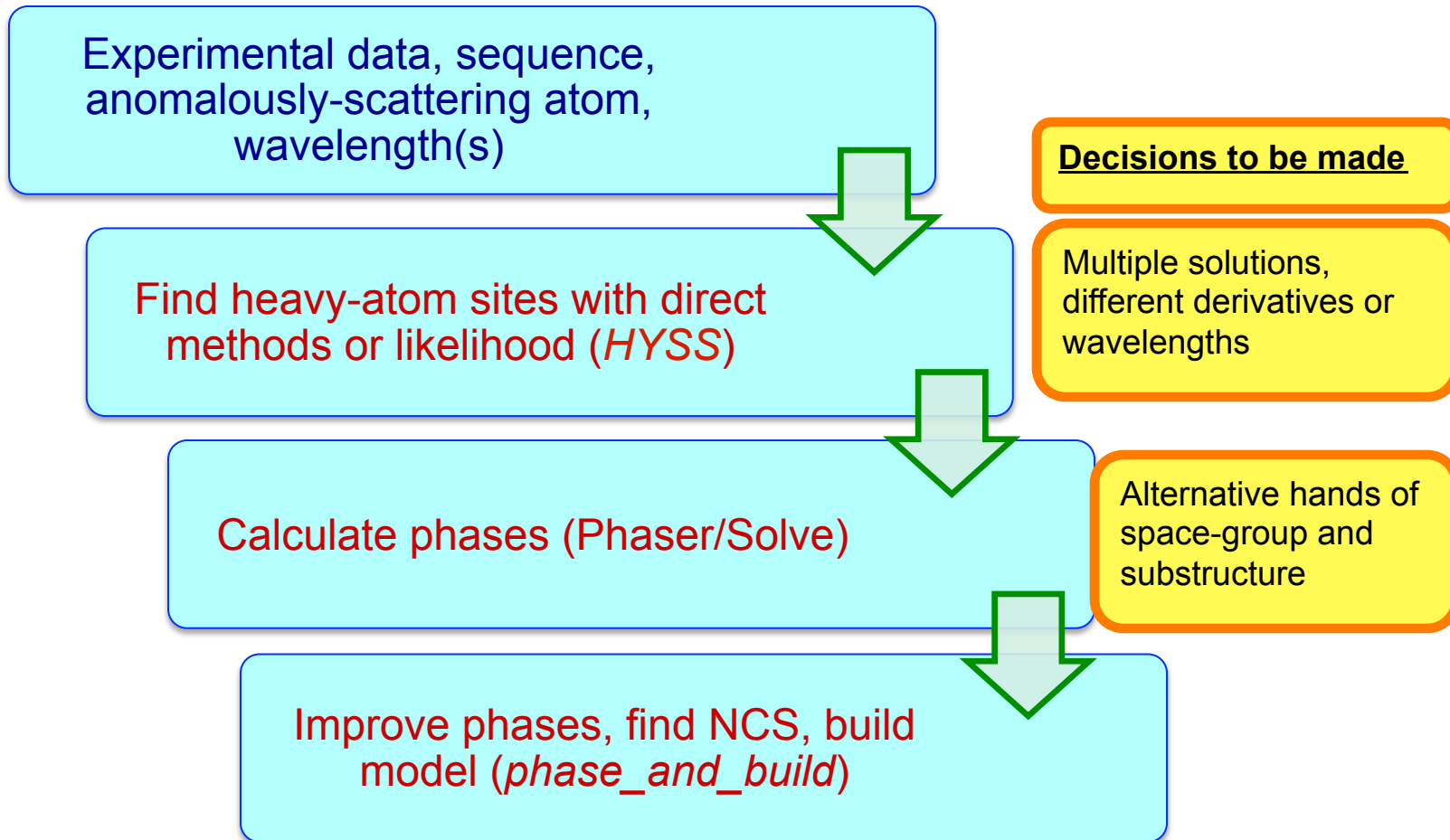


Phaser +RESOLVE  
(CC=0.79)

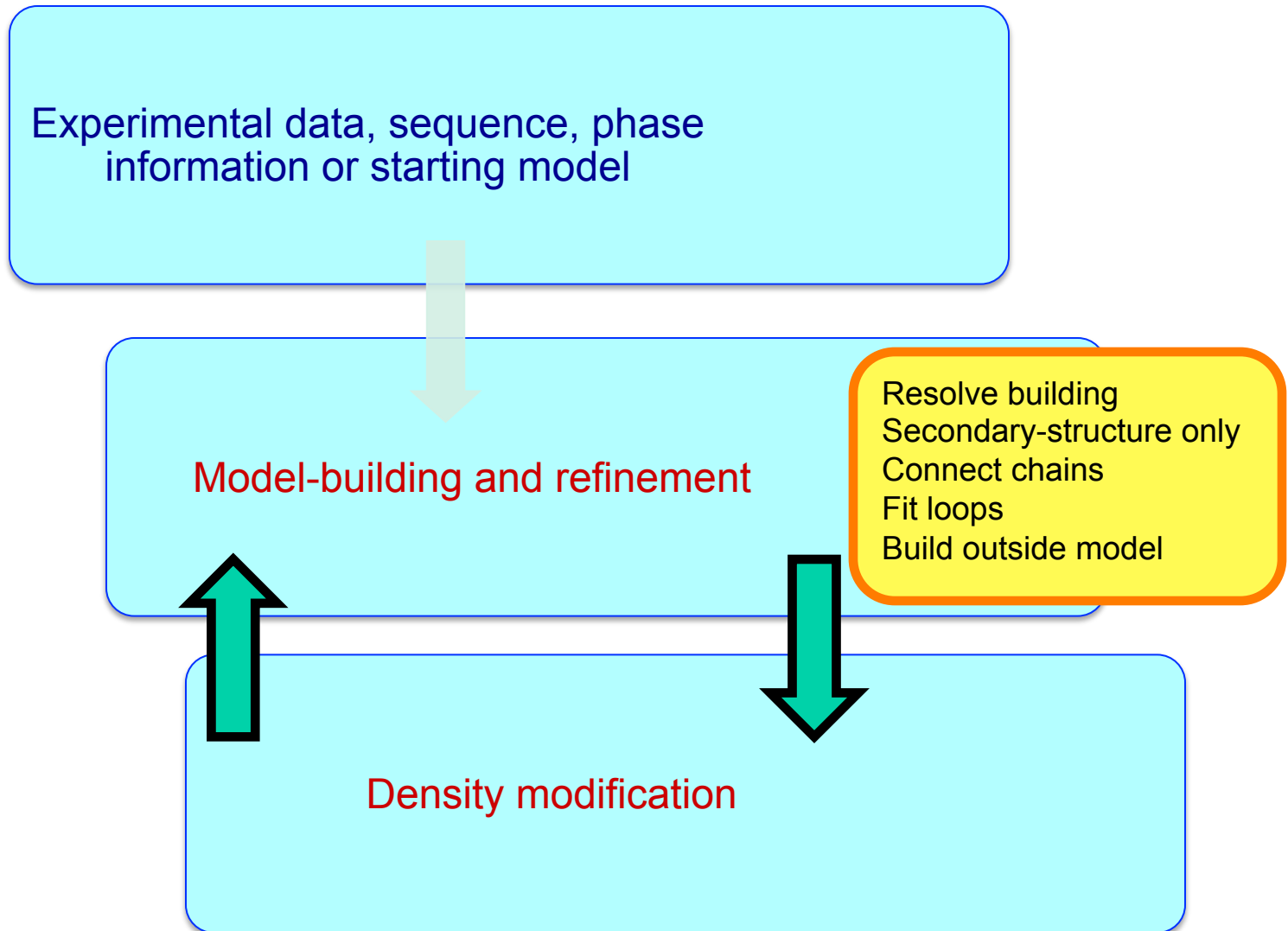




# Structure solution with *phenix.autosol*

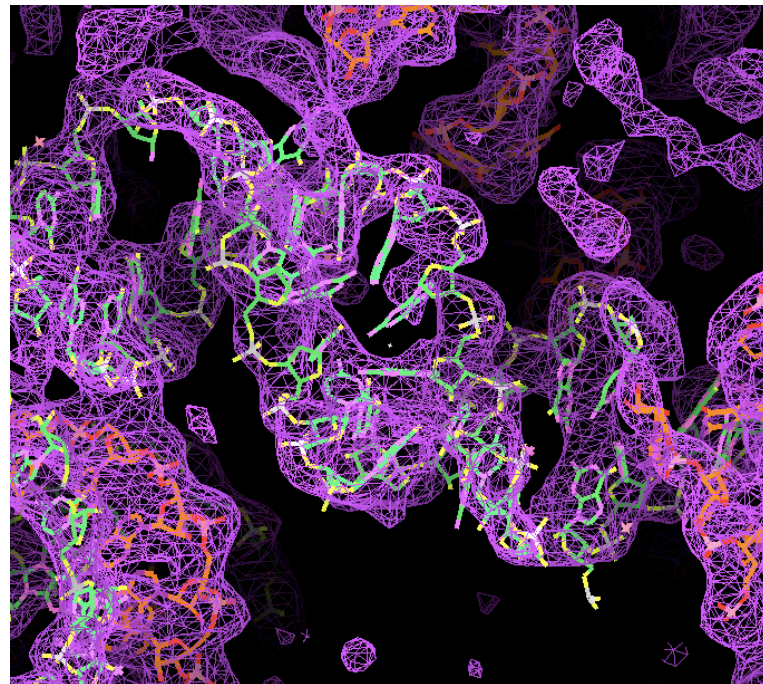
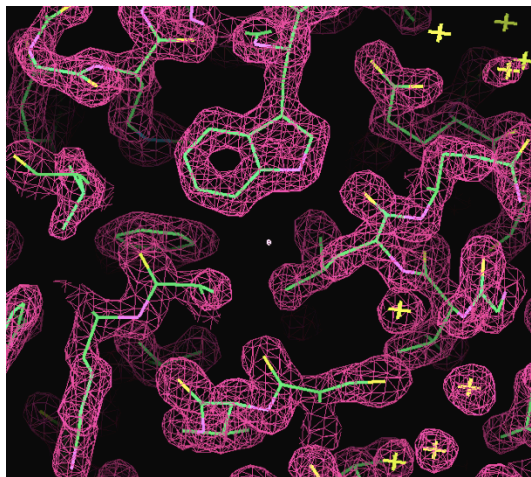


# Iterative density modification, model-building and refinement with *phenix.autobuild*

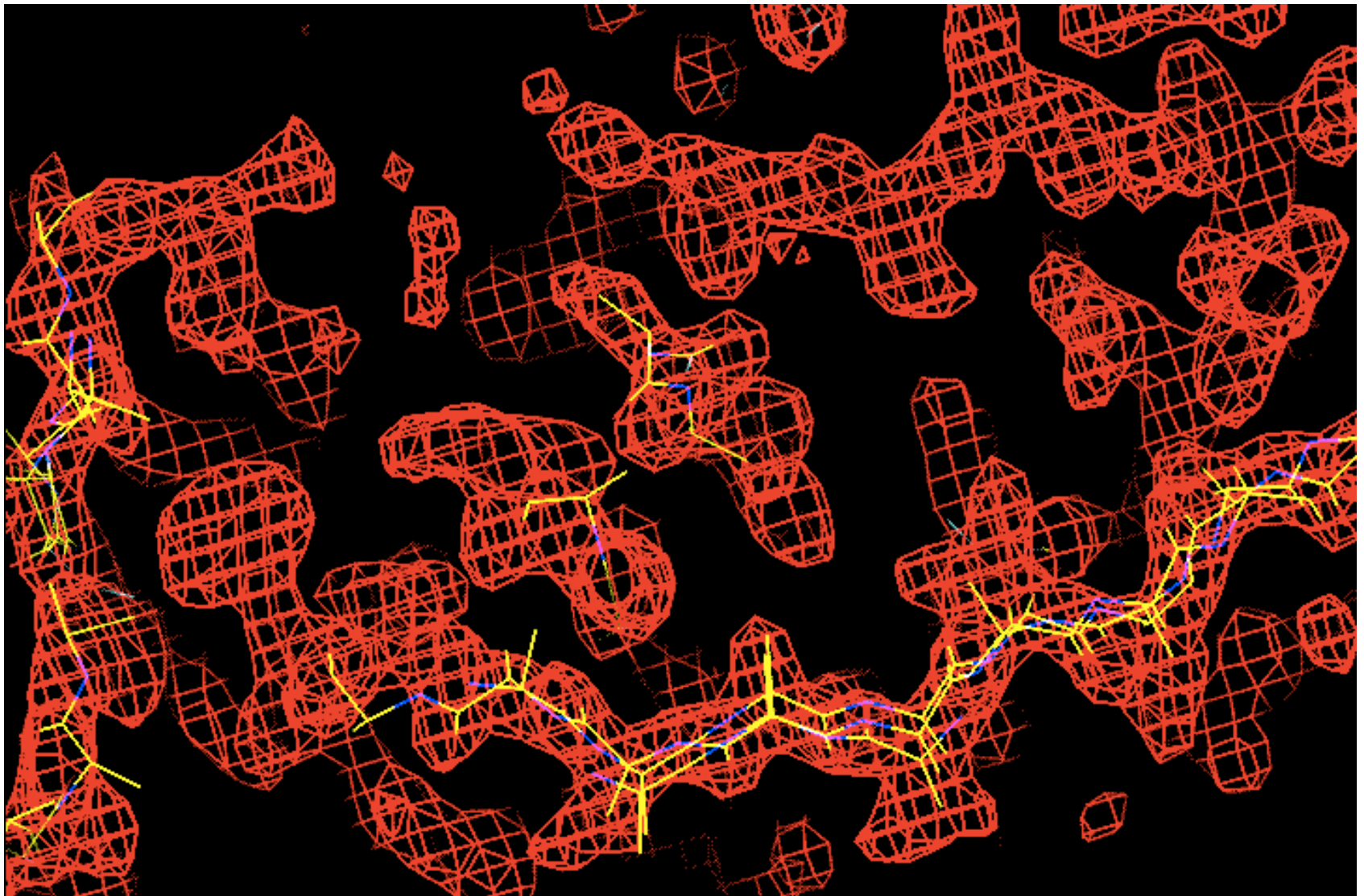


# Model-building at moderate or high resolution

- FFT-based identification of regular secondary structure
- Extension with short fragments from high-resolution structures
- Probabilistic sequence alignment

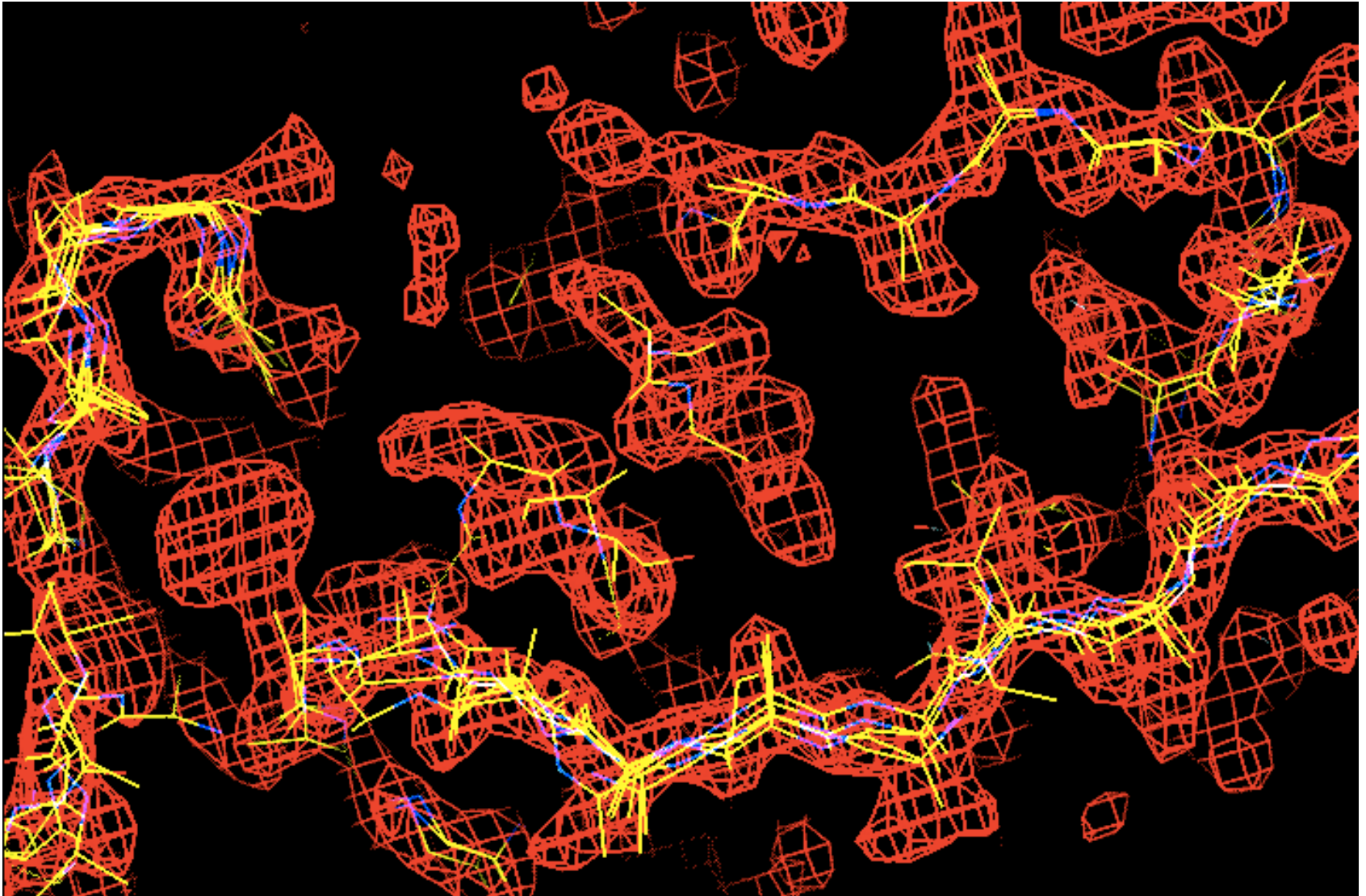


## *Initial model-building – strand fragments*

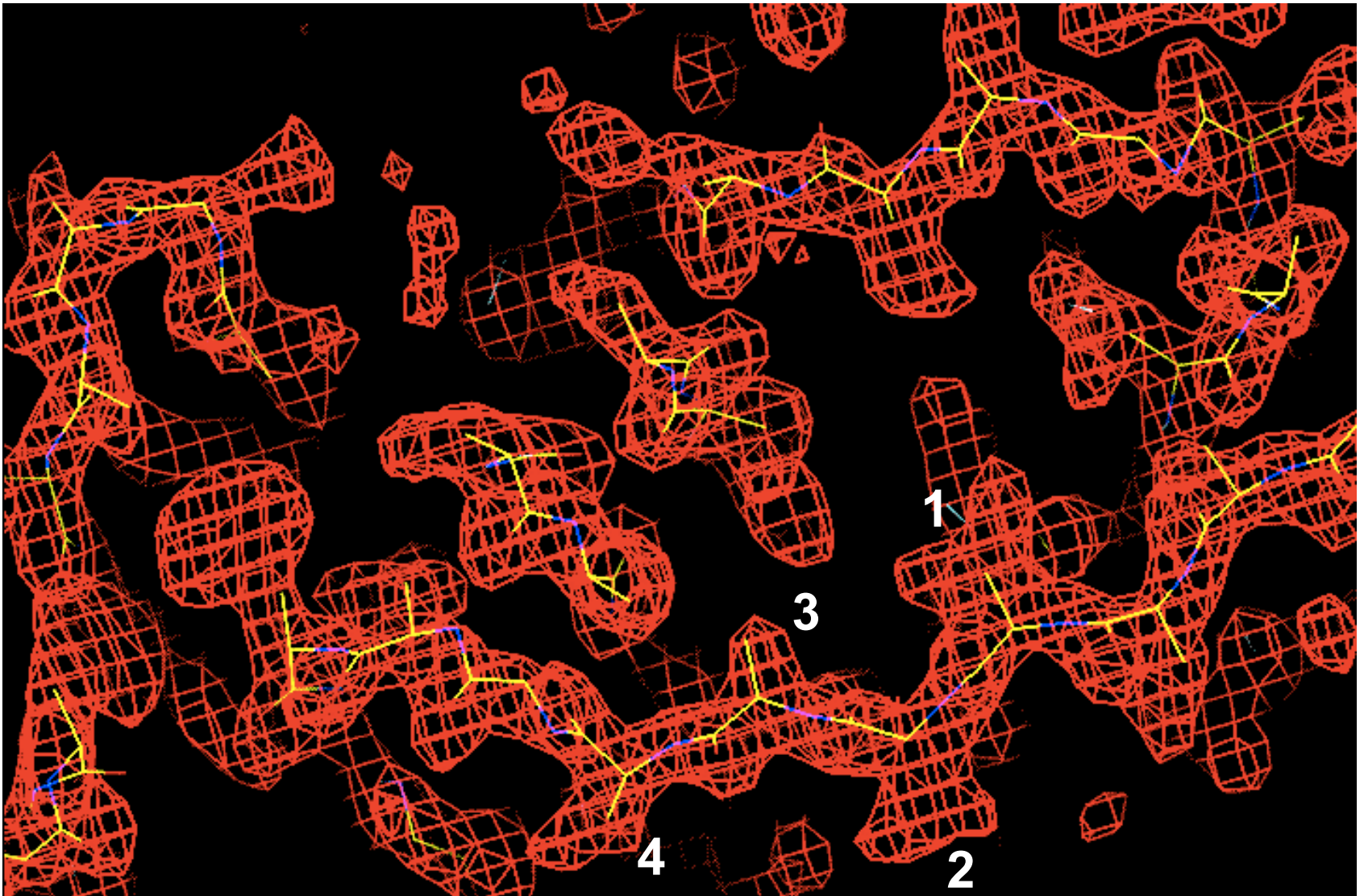




*Chain extension  
(result: many overlapping fragments)*



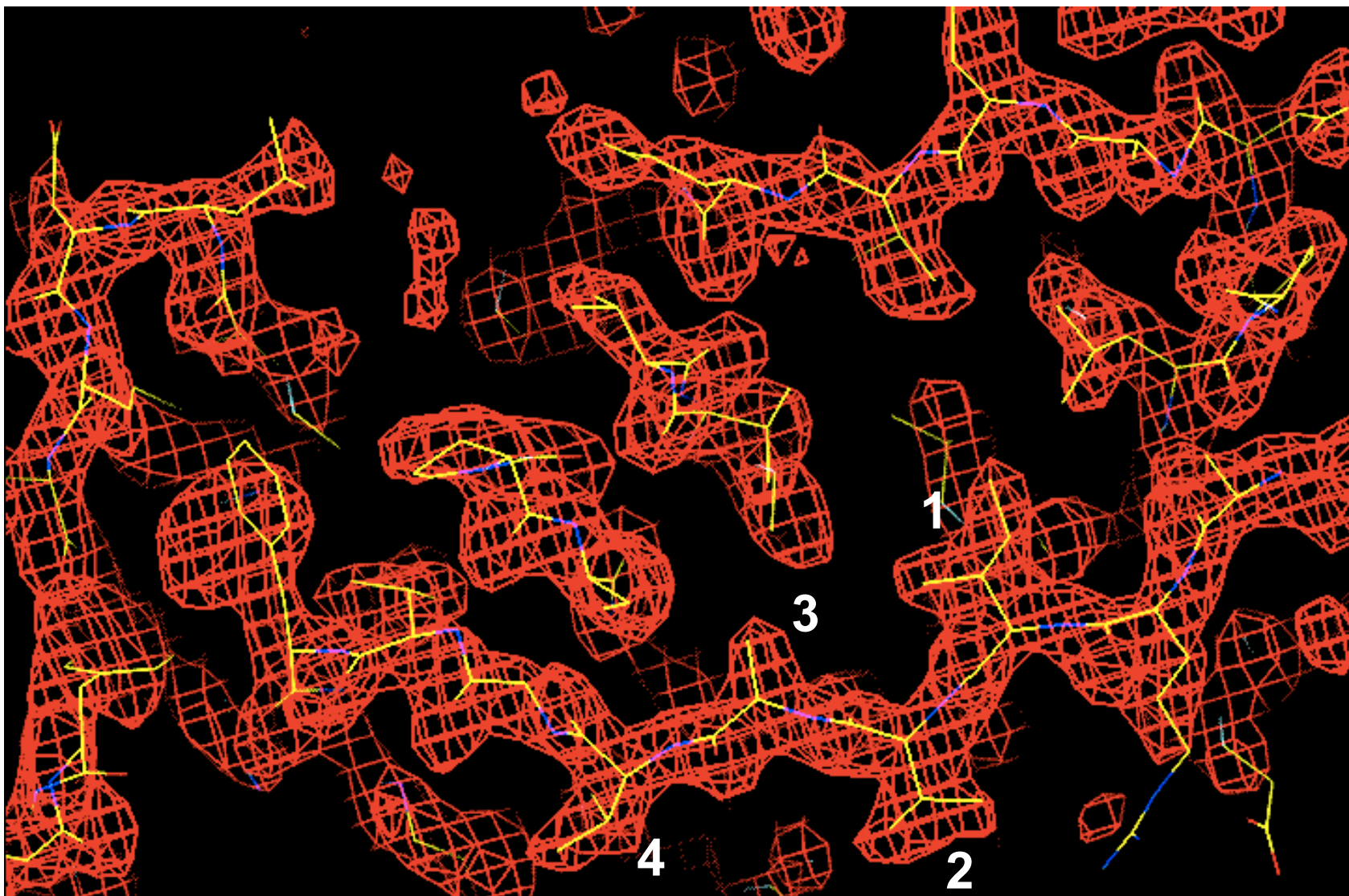
*Main-chain as a series of fragments  
(choosing the best fragment at each location)*







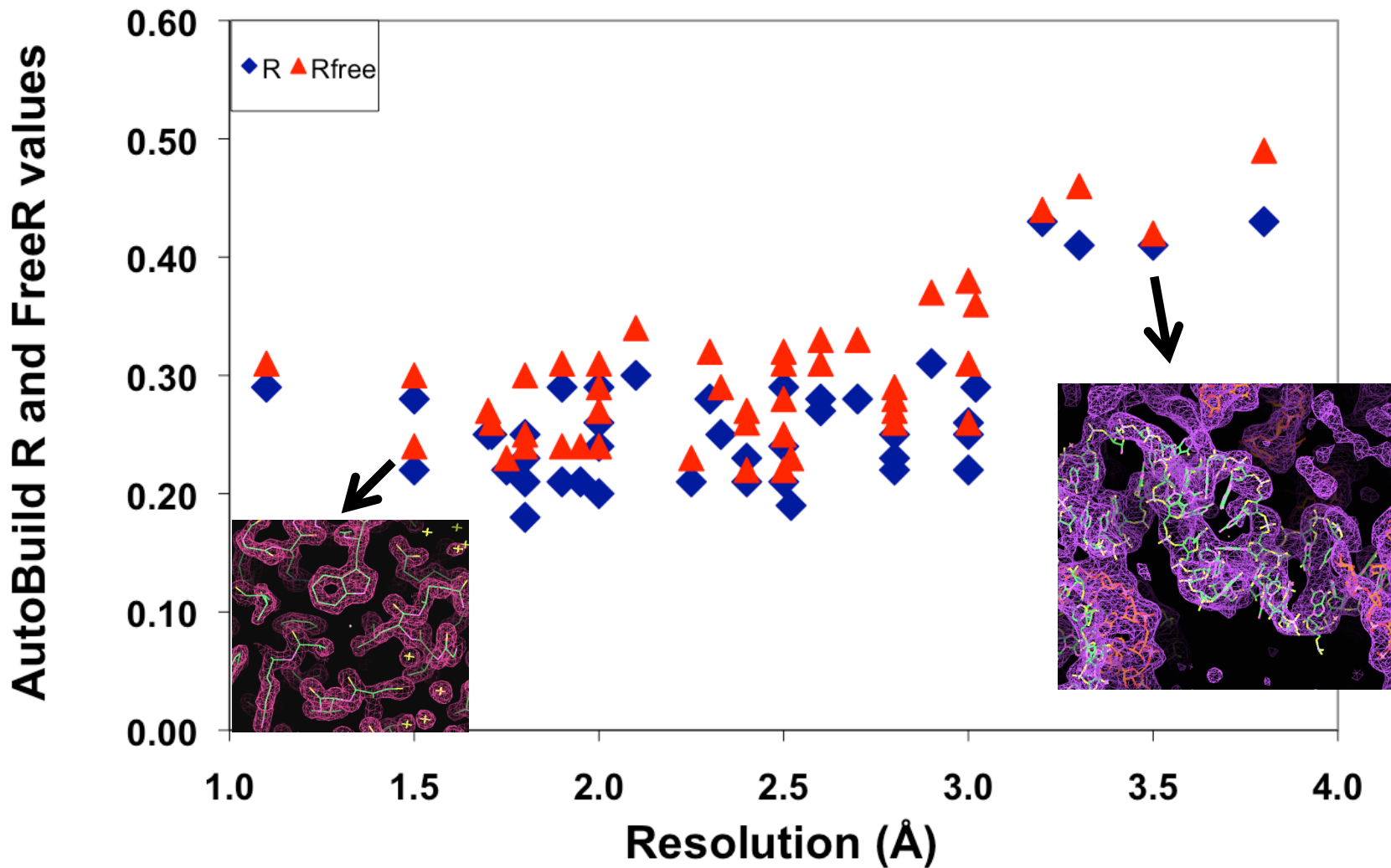
# *Addition of side-chains to fixed main-chain positions*





# AutoBuild – tests with structure library

Fully automated iterative model-building, final R/Rfree



## The Phenix Team

### Lawrence Berkeley Laboratory

Paul Adams, Pavel Afonine, Nigel Moriarty, Nicholas Sauter, Oleg Sobolev, Billy Poon



### Los Alamos National Laboratory

Tom Terwilliger, Li-Wei Hung



Randy Read, Airlie McCoy, Gabor Bunkóczi, Robert Oeffner

### Cambridge University



*An NIH/NIGMS funded  
Program Project*



### Duke University

Jane & David Richardson, Christopher Williams, Bradley Hintze