

# ***Ab initio* cryo-EM structure determination as a validation problem**

Pawel A. Penczek

**The University of Texas – Houston Medical School,  
Department of Biochemistry**



# ACKNOWLEDGMENTS

Francisco J. Asturias  
La Jolla, CA



Christian M.T. Spahn  
Charité, Berlin



NIH

# CONCLUSIONS

1. Validation should be an integral part of the structure determination process.
2. Any method should be permitted to fail under controlled circumstances as the failure can be as informative as success.
3. EM projection images are of very poor quality. Therefore, they should not be evaluated individually but as members of statistical assemblies.
4. Implementation in SPARX <http://sparx-em.org/sparxwiki/> with new additions of tools for the analysis of local variability (please see the poster).

# **Statistical cross-validation for detecting and preventing overfitting**

Problem of model selection

# EM DATA AND PARAMETER ERROR ESTIMATION

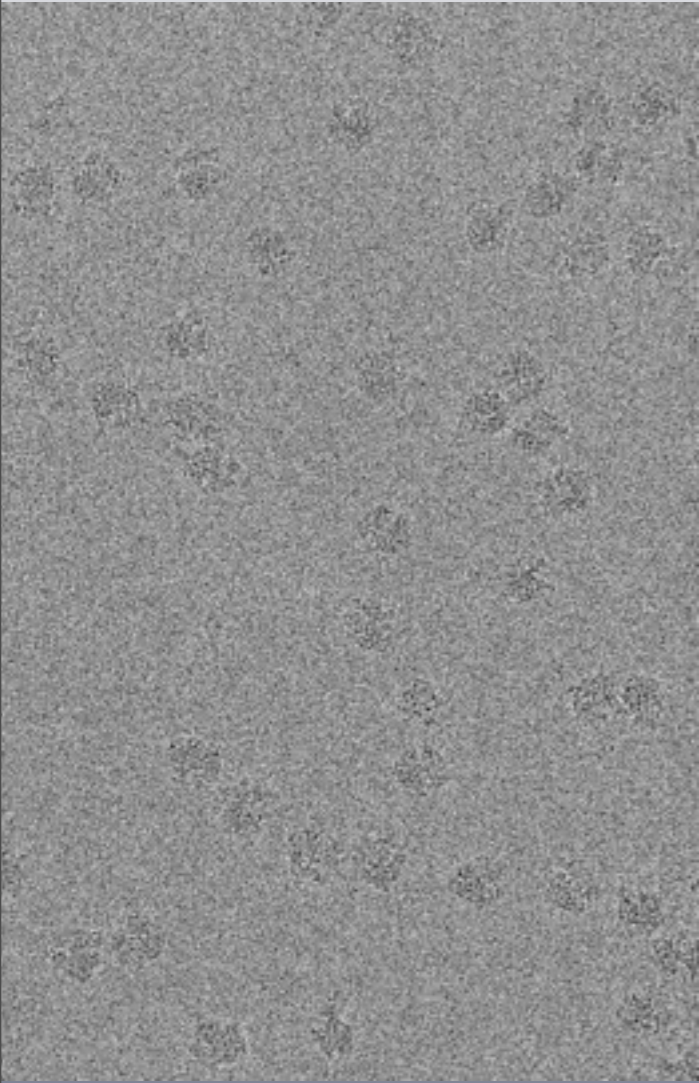
- A typical EM experiment generates a single dataset and it is not possible to derive an analytical expression to determine (alignment) parameter errors
- The challenge is then to estimate parameter errors in the absence of independent sample sets
- Statistical Resampling offers the best option for accurate estimation of parameter errors independent of assumptions about their statistical properties

# EM DATA AND PARAMETER ERROR ESTIMATION

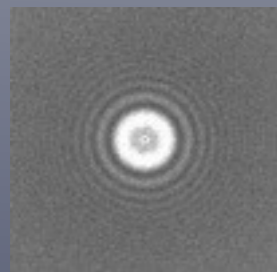
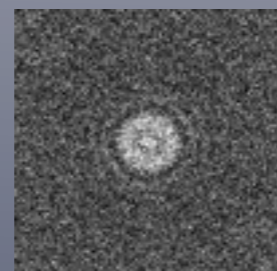
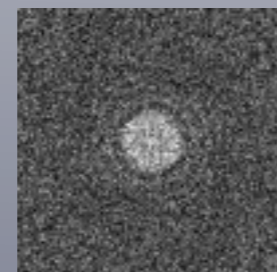
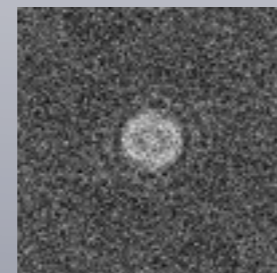
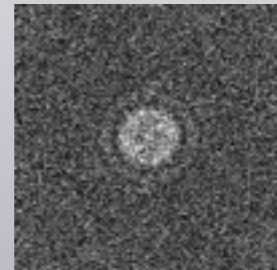
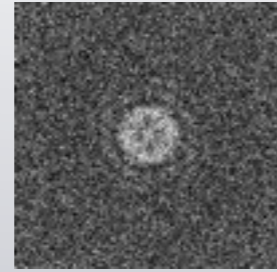
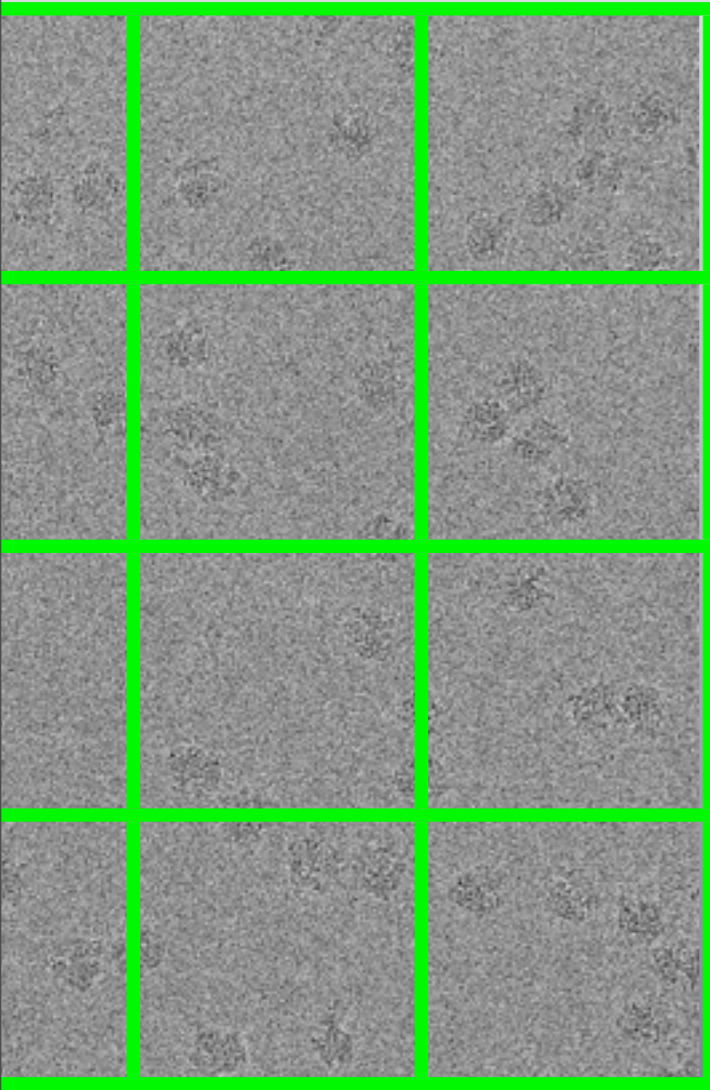
- A typical EM experiment generates a single dataset and it is not possible to derive an analytical expression to determine (alignment) parameter errors
- The challenge is then to estimate parameter errors in the absence of independent sample sets
- Statistical Resampling offers the best option for accurate estimation of parameter errors independent of assumptions about their statistical properties

If we treat the observed sample (EM dataset) as though it exactly represented the entire population, evaluating artificial variability generated through resampling allows us to accurately estimate variability of a sample statistic

# CTF parameter estimation and error assessment through bootstrap resampling (CTER)



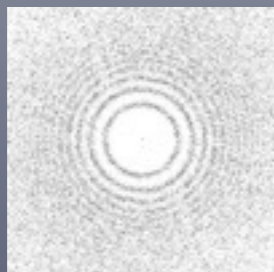
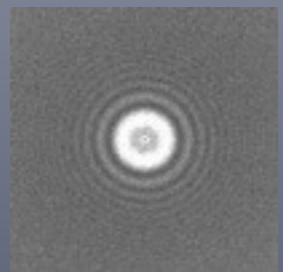
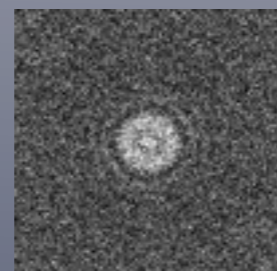
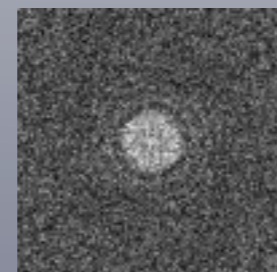
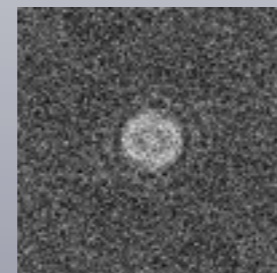
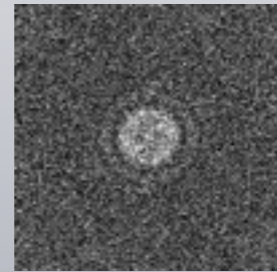
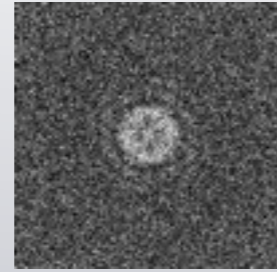
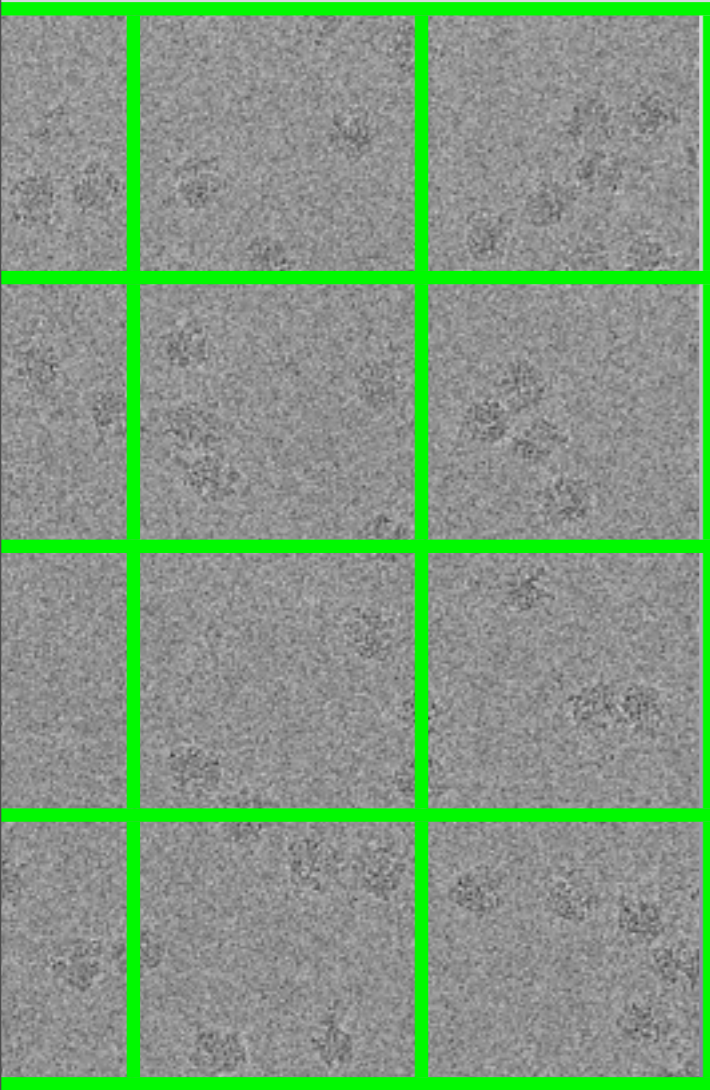
# CTF parameter estimation and error assessment through bootstrap resampling (CTER)



Average power spectrum and its variance

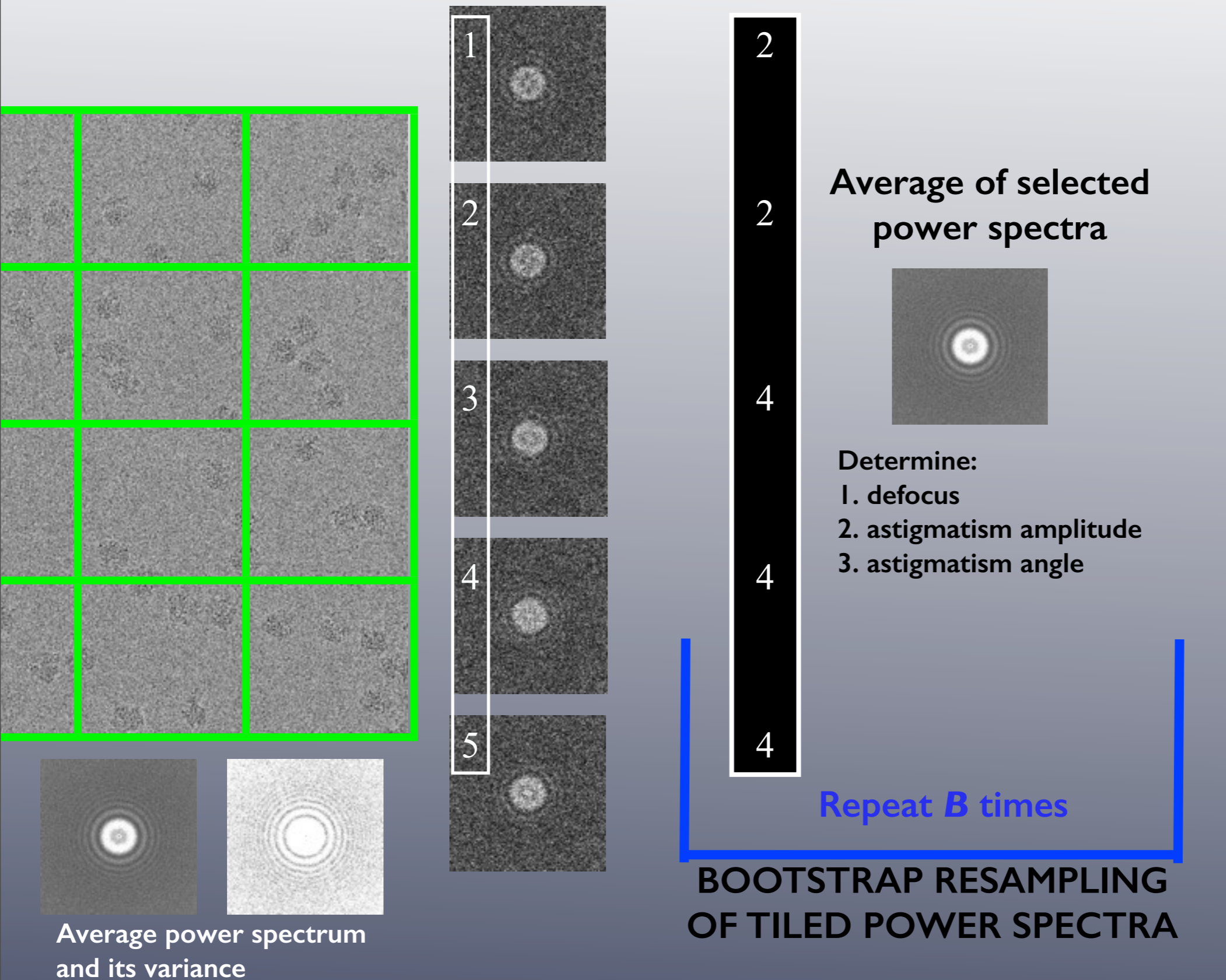


# CTF parameter estimation and error assessment through bootstrap resampling (CTER)

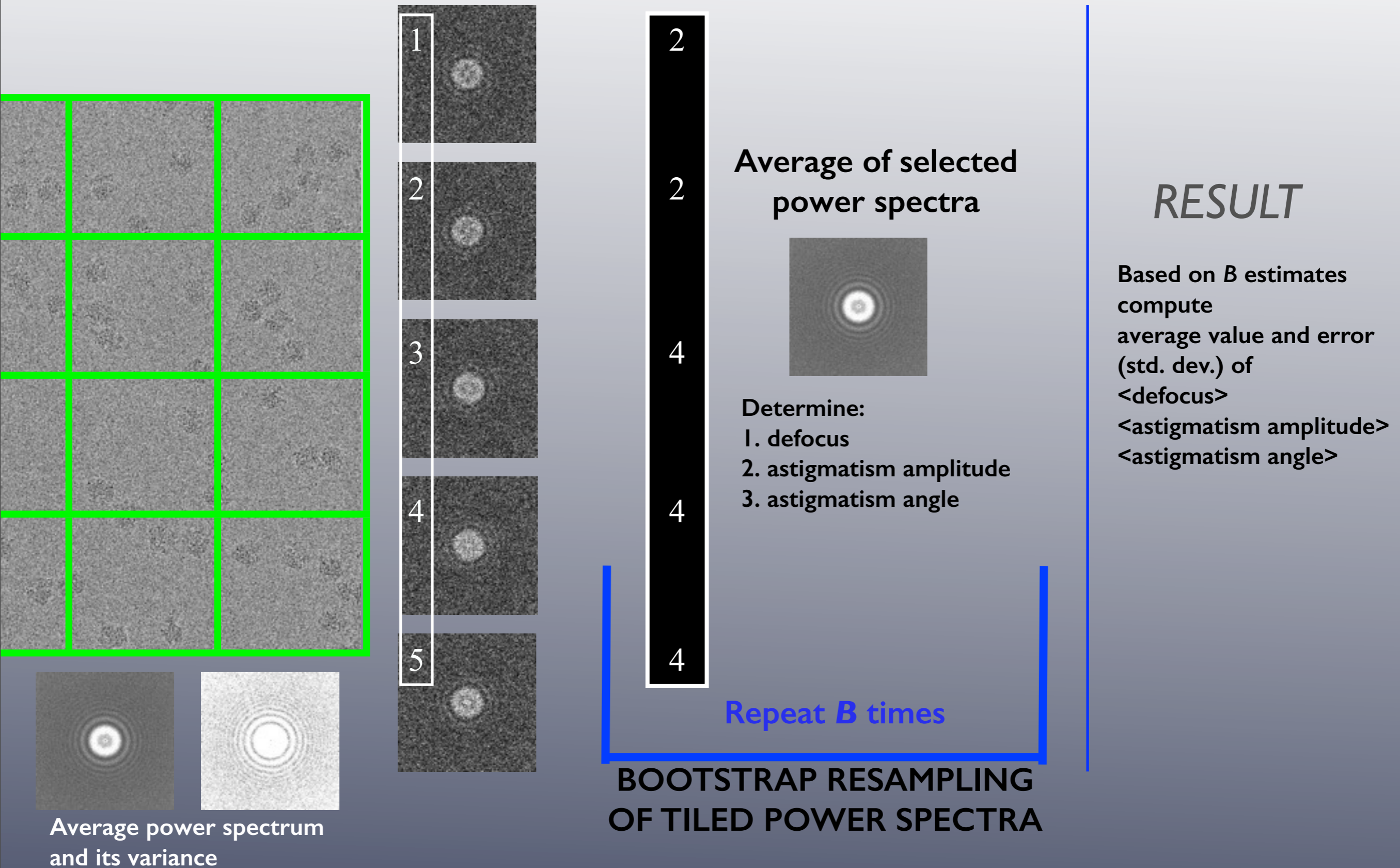


Average power spectrum and its variance

# CTF parameter estimation and error assessment through bootstrap resampling (CTER)



# CTF parameter estimation and error assessment through bootstrap resampling (CTER)



# ISAC: VALIDATION OF 2D MULTI-REFERENCE ALIGNMENT THROUGH STABILITY TESTING

1. If a set of images is homogeneous, the result from reference-free alignment is stable even for very low SNR data.
2. The converse is true, i.e., if a set of images is stable, it must be homogeneous.

2D alignment is **stable** if perturbation of initial alignment parameters does not produce dramatically different results.

# ISAC: VALIDATION OF 2D MULTI-REFERENCE ALIGNMENT THROUGH STABILITY TESTING

1. If a set of images is homogeneous, the result from reference-free alignment is stable even for very low SNR data.
2. The converse is true, i.e., if a set of images is stable, it must be homogeneous.

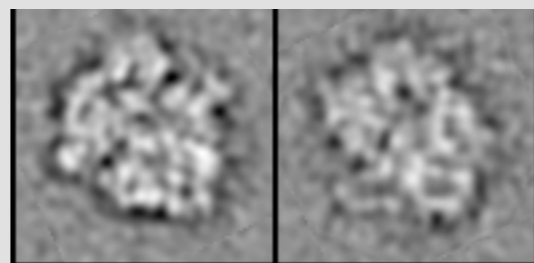
2D alignment is **stable** if perturbation of initial alignment parameters does not produce dramatically different results.

Assuming 1 and 2 are correct:

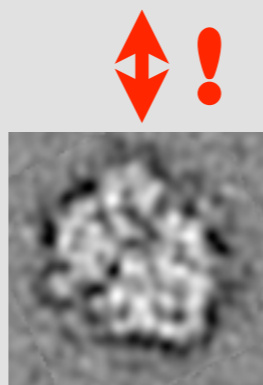
If we can find homogeneous subsets of images, we can solve the multi-reference alignment problem.

# STABLE VS. UNSTABLE CLASSES: A TEST CASE

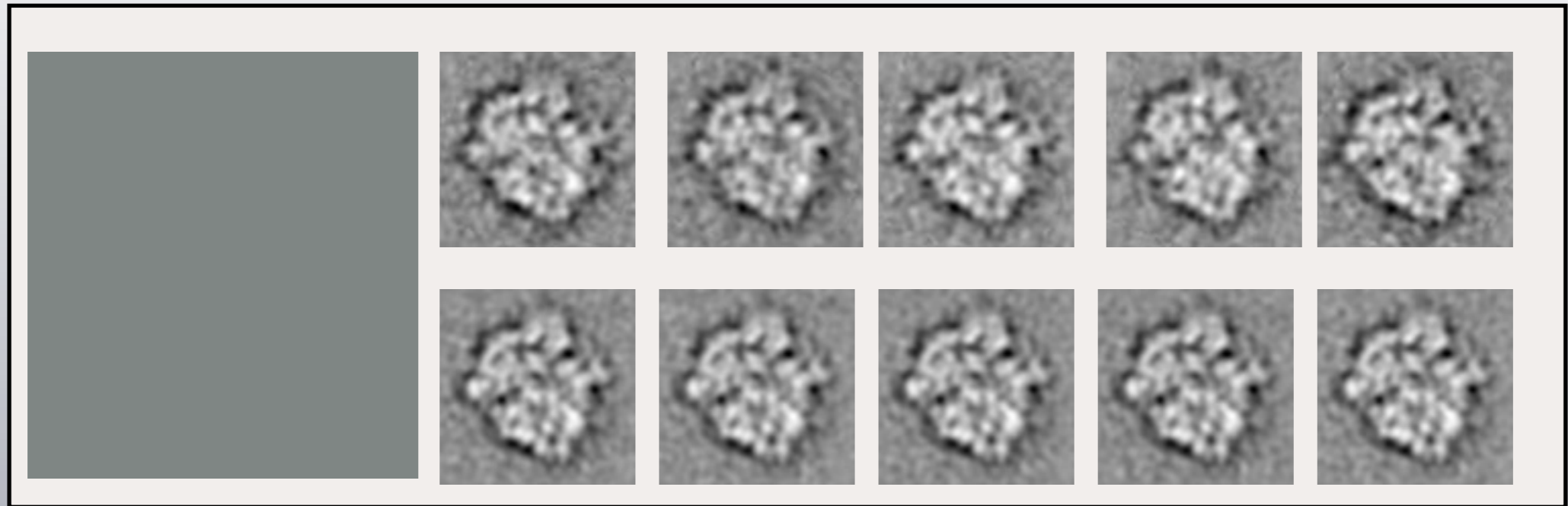
Two groups were mixed 50-50, their respective averages are:



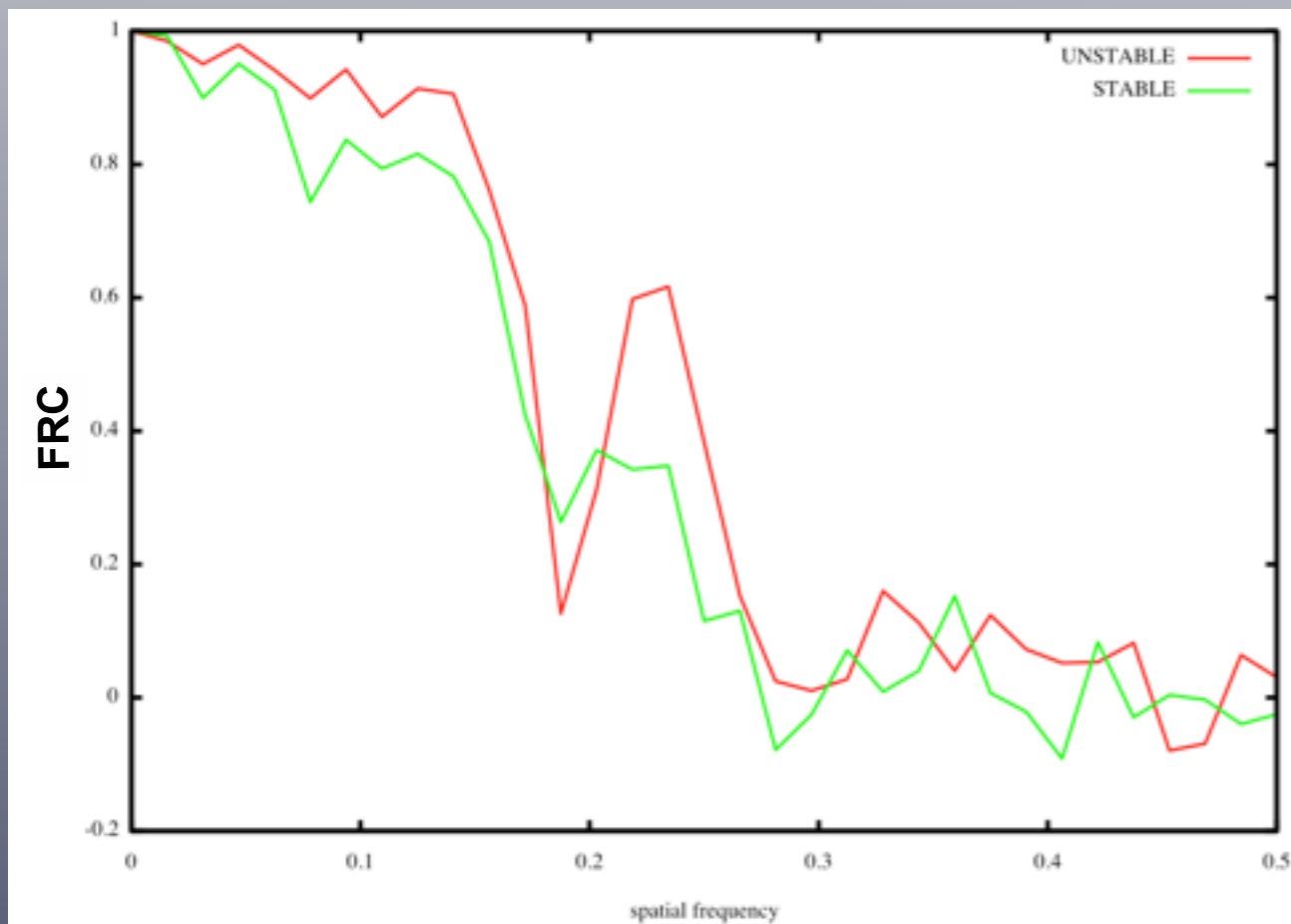
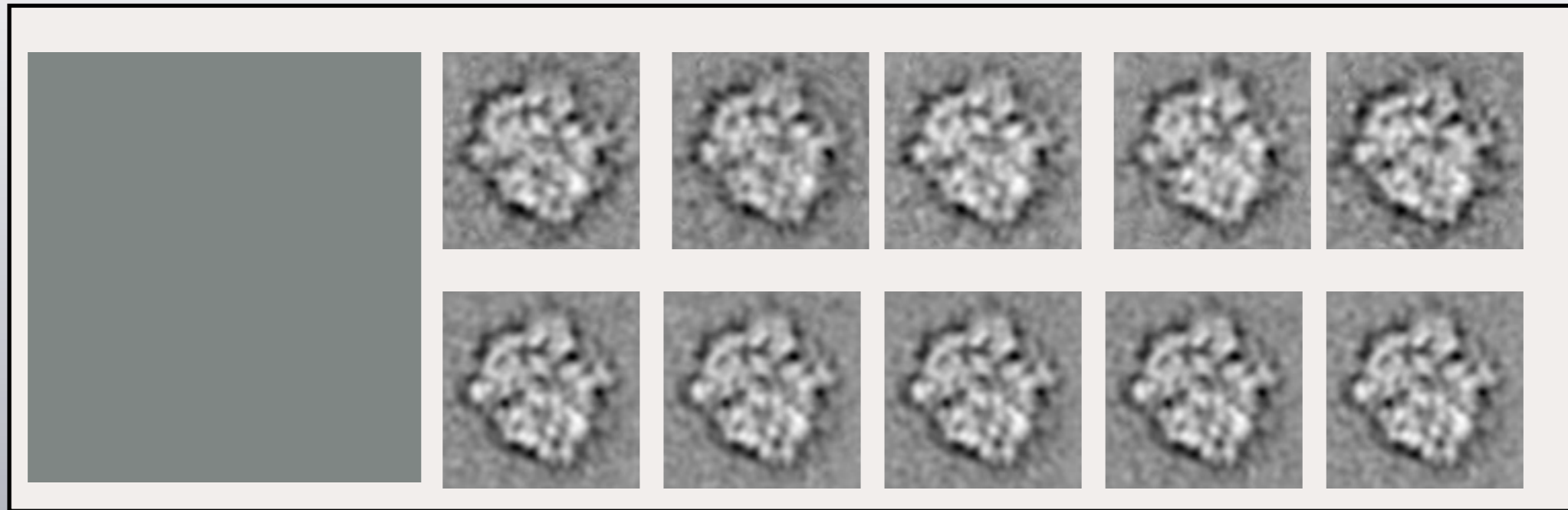
Sum of these two averages:



# STABLE VS. UNSTABLE CLASSES: TEST RESULTS

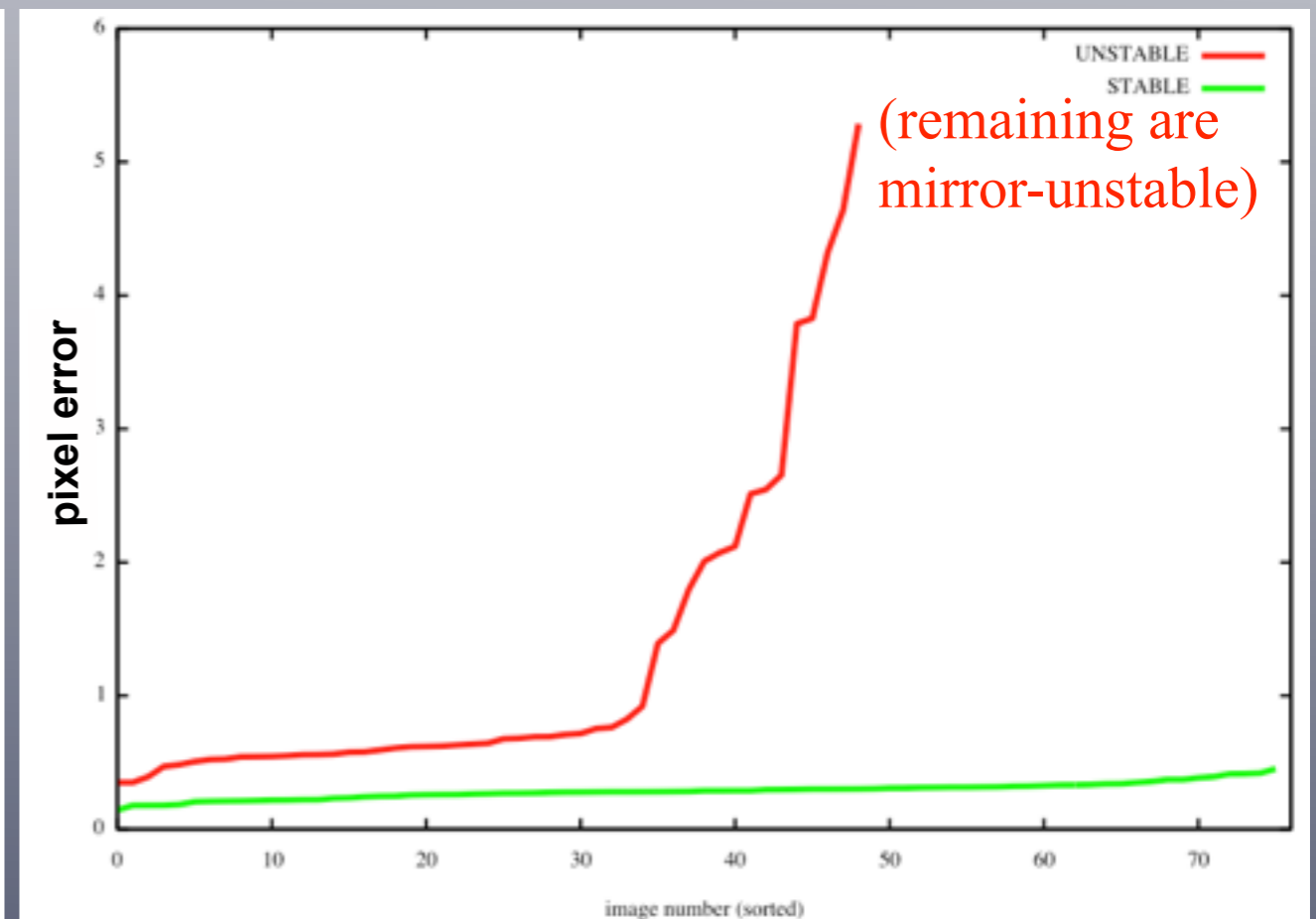
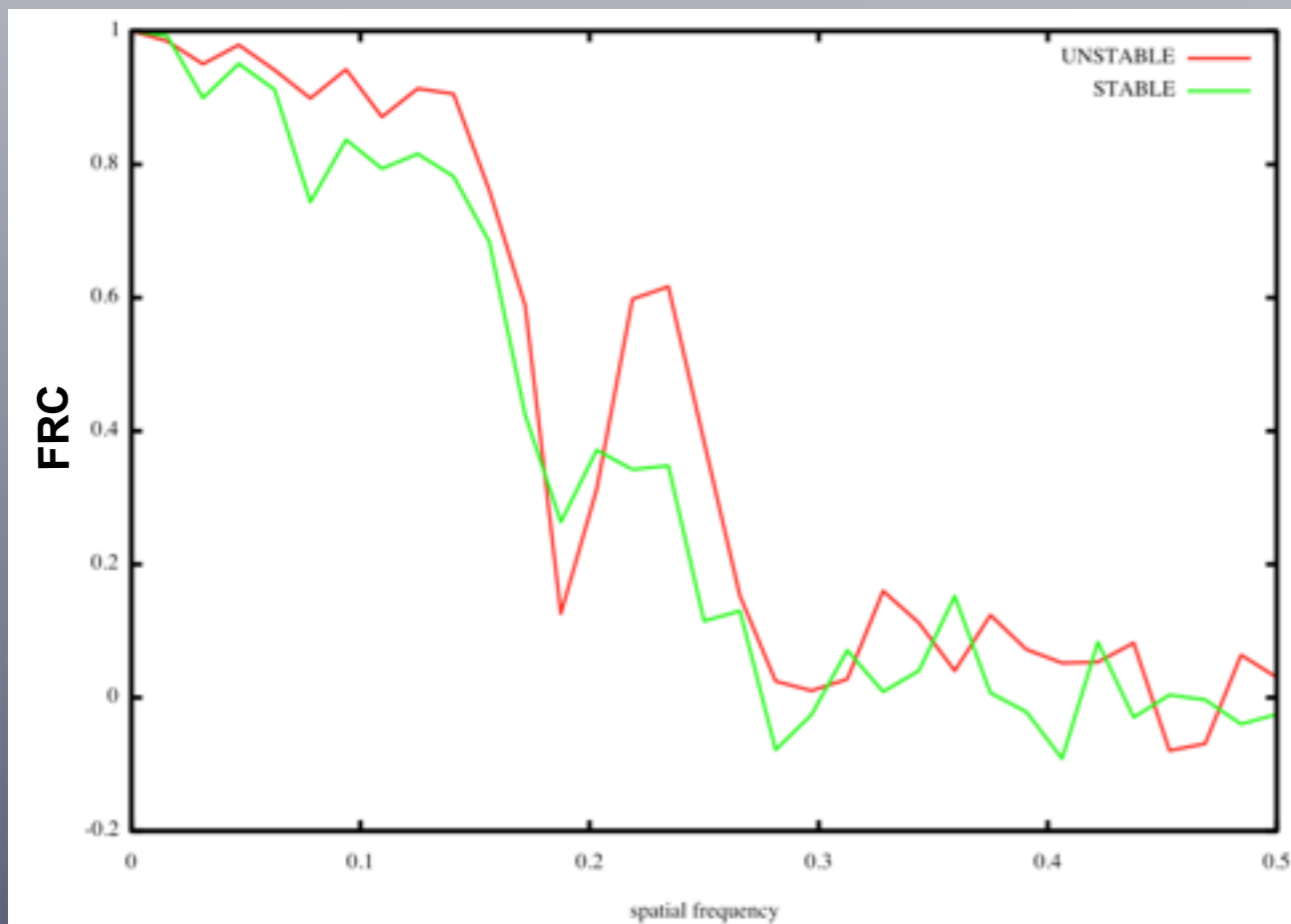
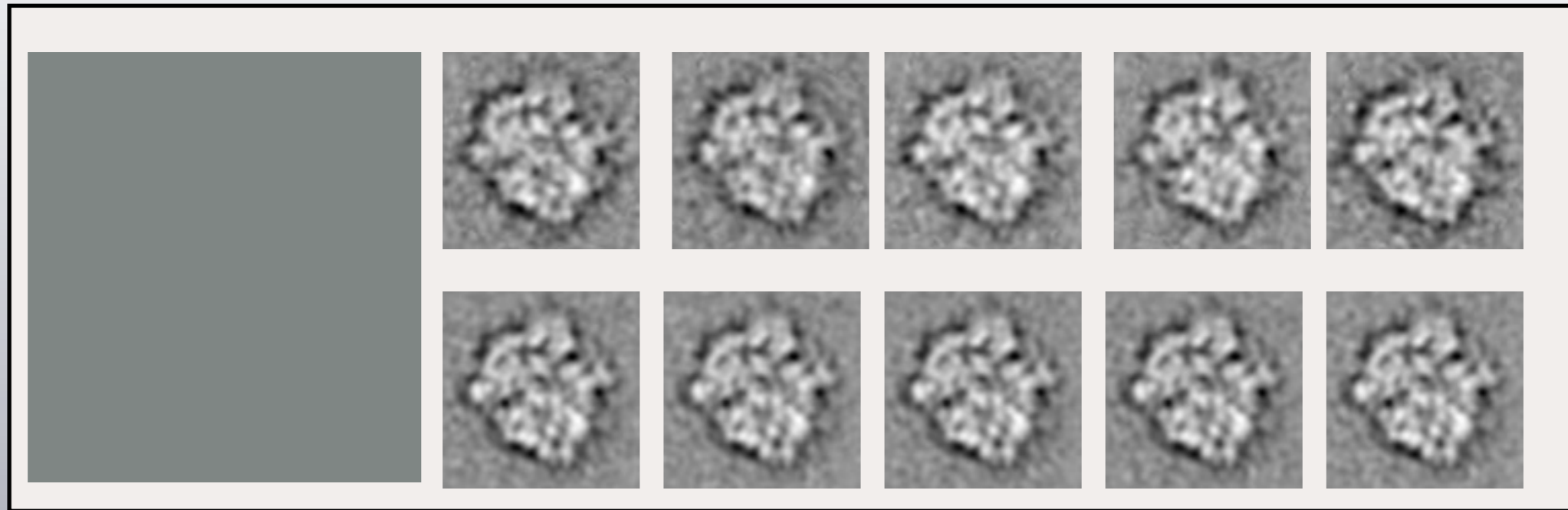


# STABLE VS. UNSTABLE CLASSES: TEST RESULTS





# STABLE VS. UNSTABLE CLASSES: TEST RESULTS



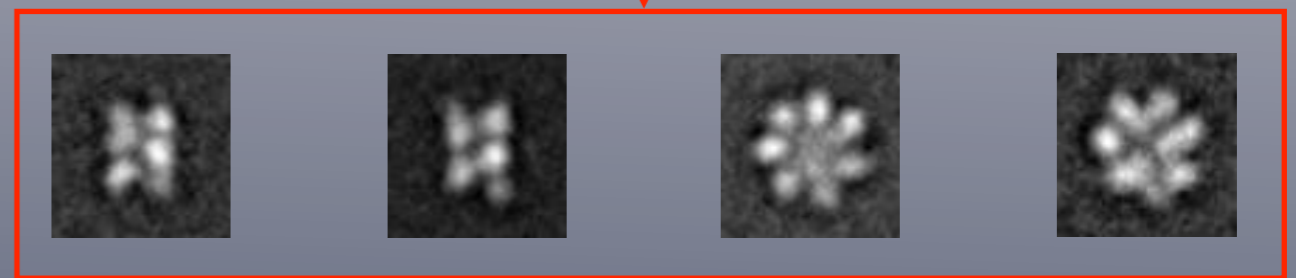
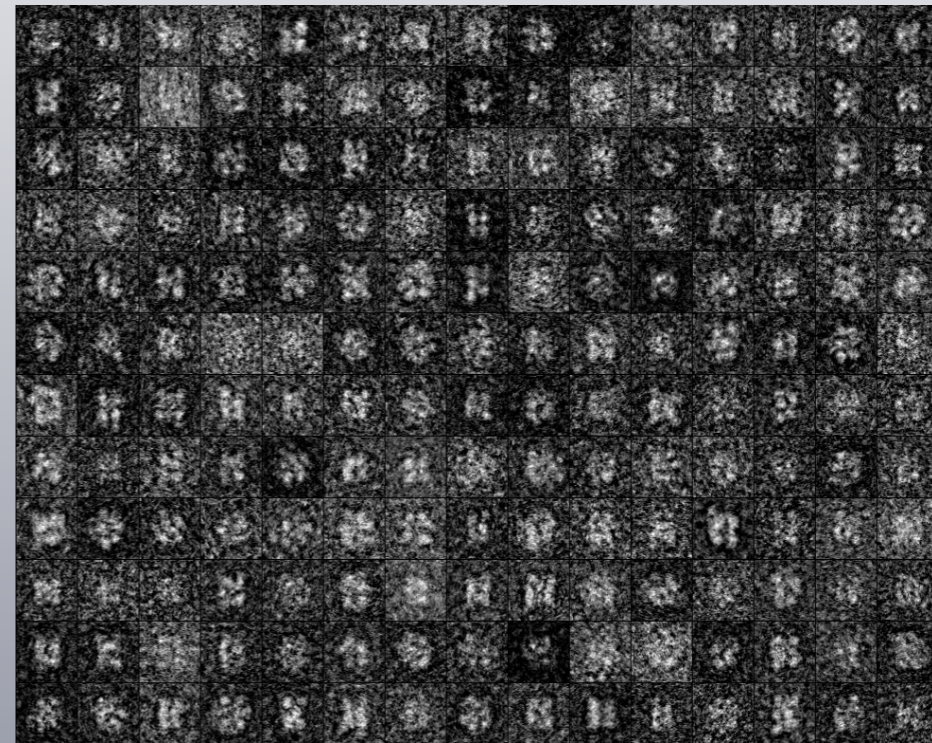
# 2D MULTI-REFERENCE ALIGNMENT (MRA)

MRA is equivalent to  $K$ -means clustering, with the distance between images defined as a maximum similarity over the permissible range of image rotations and translations.

$K$ -means results depend on the solution to another nontrivial problem: the alignment of a set of 2D images.

Because neither of these two problems can be easily solved, the difficulty is compounded.

$n$  images



$K$  averages (clusters)

# K-MEANS CLUSTERING

## KNOWN PROPERTIES:

- Very fast convergence guaranteed in a finite number of steps
- Converges only to a local minimum
- Unclear how to determine the appropriate number of classes ( $K$ )
- All images must be assigned to an average
- The solution (final averages) depends on the initial set of averages, and will change if clustering is repeated using different initial averages
- In EM, when alignment is added, classes tend to collapse

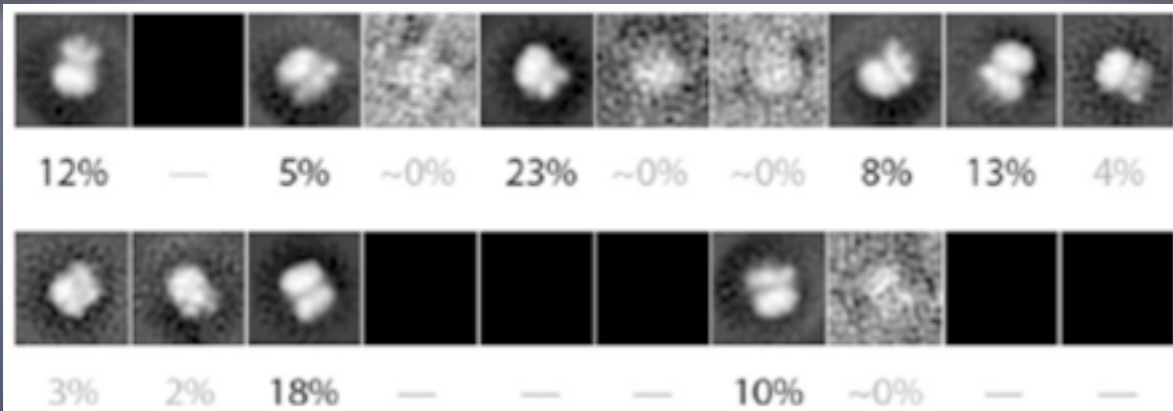
*K*-means group assignments  
minimum distance to a template within a row



# K-MEANS CLUSTERING

## KNOWN PROPERTIES:

- Very fast convergence guaranteed in a finite number of steps
- Converges only to a local minimum
- Unclear how to determine the appropriate number of classes ( $K$ )
- All images must be assigned to an average
- The solution (final averages) depends on the initial set of averages, and will change if clustering is repeated using different initial averages
- In EM, when alignment is added, classes tend to collapse



*K*-means group assignments  
minimum distance to a template within a row



# $EQK$ (EQUAL GROUP SIZE) - MEANS CLUSTERING

Assign  $n$  images to  $K$  classes  
such that each class contains










$$\frac{n}{K} \text{ images}$$

# $EQK^{(\text{EQUAL GROUP SIZE})}$ -MEANS CLUSTERING

Assign  $n$  images to  $K$  classes  
such that each class contains

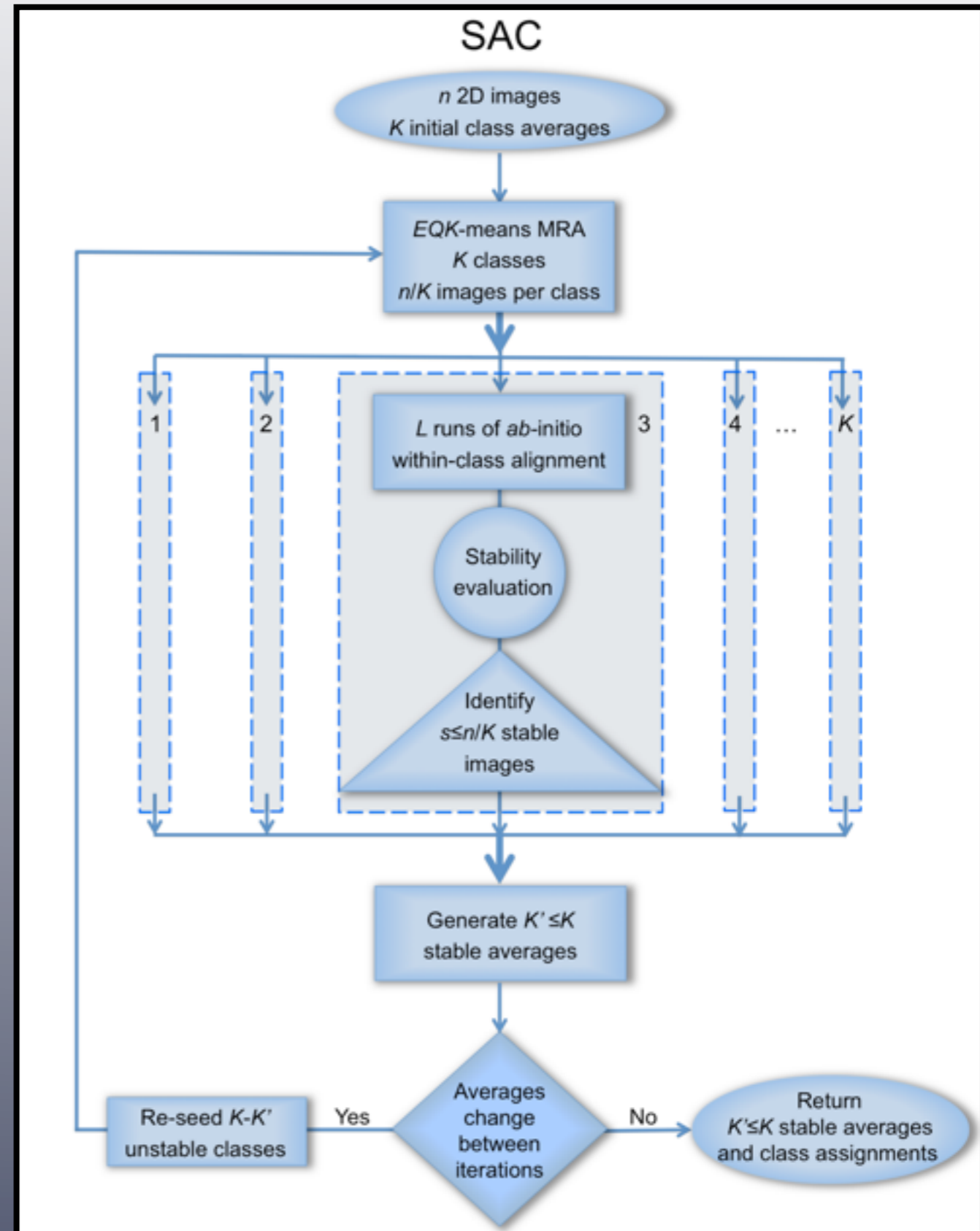
$$\frac{n}{K} \text{ images}$$

$EQK$ -means group assignments  
minimum distance to all templates, maximum number per group=3

|                                                                                       |  |  | ... |  |
|---------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|-----|-------------------------------------------------------------------------------------|
|    | $d_{11}^2$                                                                          | $d_{12}^2$                                                                          | ... | $d_{1K}^2$                                                                          |
|    | $d_{21}^2$                                                                          | $d_{22}^2$                                                                          | ... | $d_{2K}^2$                                                                          |
|  | $d_{31}^2$                                                                          | $d_{32}^2$                                                                          | ... | $d_{3K}^2$                                                                          |
|  | $d_{41}^2$                                                                          | $d_{42}^2$                                                                          | ... | $d_{4K}^2$                                                                          |
|  | $d_{51}^2$                                                                          | $d_{52}^2$                                                                          | ... | $d_{5K}^2$                                                                          |
| ⋮                                                                                     | ⋮                                                                                   | ⋮                                                                                   | ⋮   | ⋮                                                                                   |
|  | $d_{n1}^2$                                                                          | $d_{n2}^2$                                                                          | ... | $d_{nK}^2$                                                                          |

# A PROTOCOL FOR TESTING ALIGNMENT STABILITY

1. Run reference-free alignment  $L$ -times, using randomized initial orientation parameters
2. Bring all  $L$  sets of solutions into register by simultaneous minimization of the variance of orientation parameters (similar but not equivalent to alignment of resulting averages)
3. Compute pixel error for each image using orientation parameters for  $L$  positions it adopted
4. The set is called stable if the average of pixel errors for all images in  $L$  alignments is less than a predefined threshold (usually one pixel).



# CANDIDATE CLASS AVERAGES

The image displays a grid of 300 small, low-resolution images arranged in 15 columns and 20 rows. The top-right corner of the grid is blacked out. One cell in the 15th column and 14th row is highlighted with a red border. The images appear to be small, blurry, and low-contrast, possibly representing candidate class averages for a classification task. The grid is organized into 15 columns and 20 rows, with the top-right corner being blacked out.



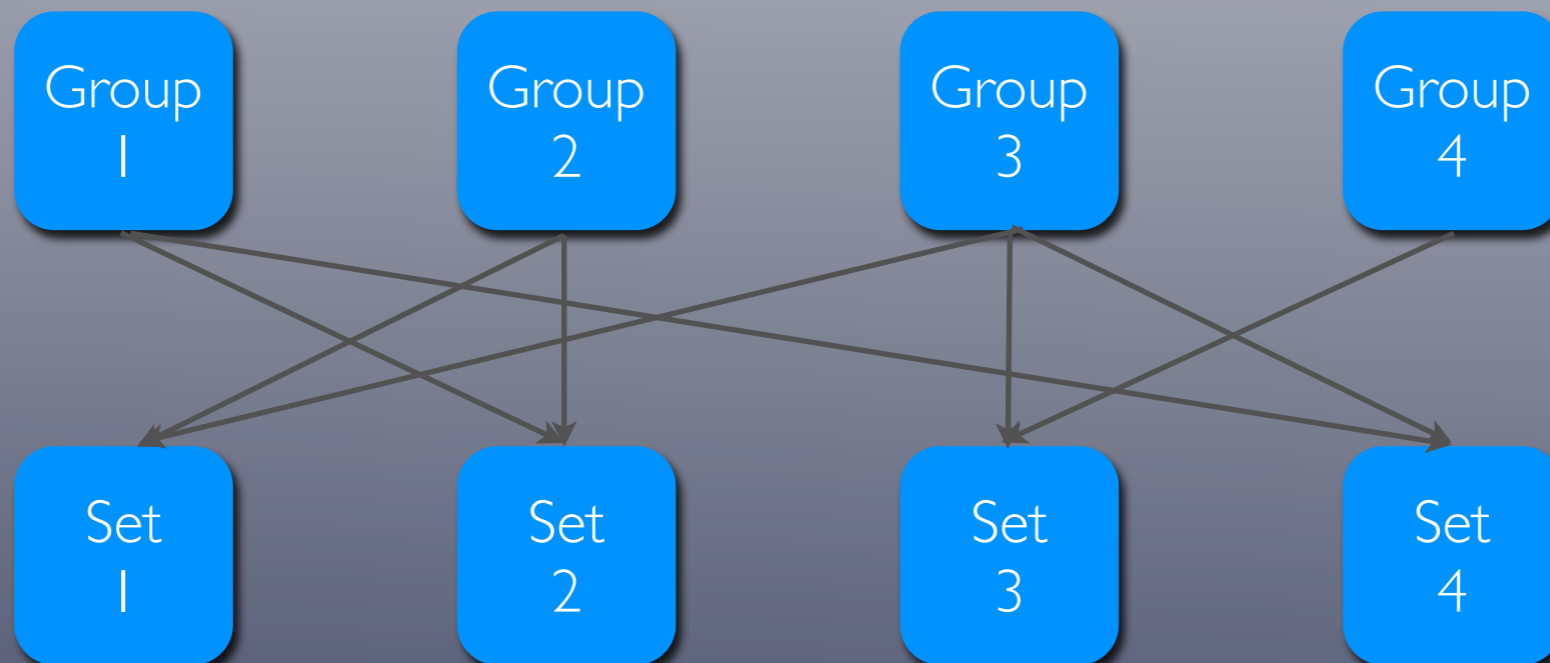
# CANDIDATE CLASS AVERAGES

- All images are accounted for (assigned to class averages)
- No validation
- The candidate class averages are used as initial templates for proper ISAC

# REPRODUCIBILITY

- Since *EQK*-means, even if combined with an alignment stability test, does not guarantee an optimum solution (global minimum) and stable groups can be fake, we require the solution to be reproducible over a number of quasi-independent runs.
- We have  $m=4$  *EQK*-means runs analyzing the data in parallel. Once all runs produce their respective averages, we compare assignments of images to class averages and select as reproducible subsets shared among quasi-independent runs.

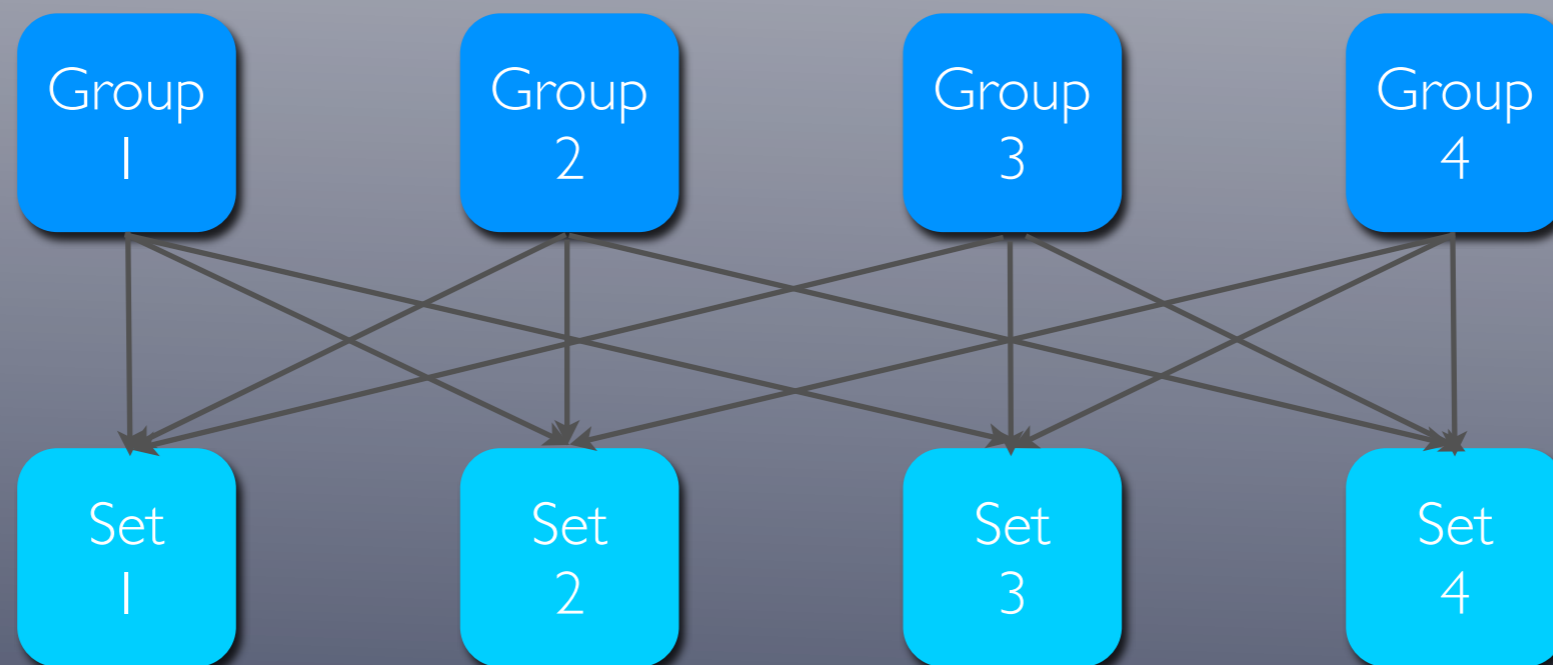
$m=2$



# REPRODUCIBILITY

- Since *EQK*-means, even if combined with an alignment stability test, does not guarantee an optimum solution (global minimum) and stable groups can be fake, we require the solution to be reproducible over a number of quasi-independent runs.
- We have  $m=4$  *EQK*-means runs analyzing the data in parallel. Once all runs produce their respective averages, we compare assignments of images to class averages and select as reproducible subsets shared among quasi-independent runs.

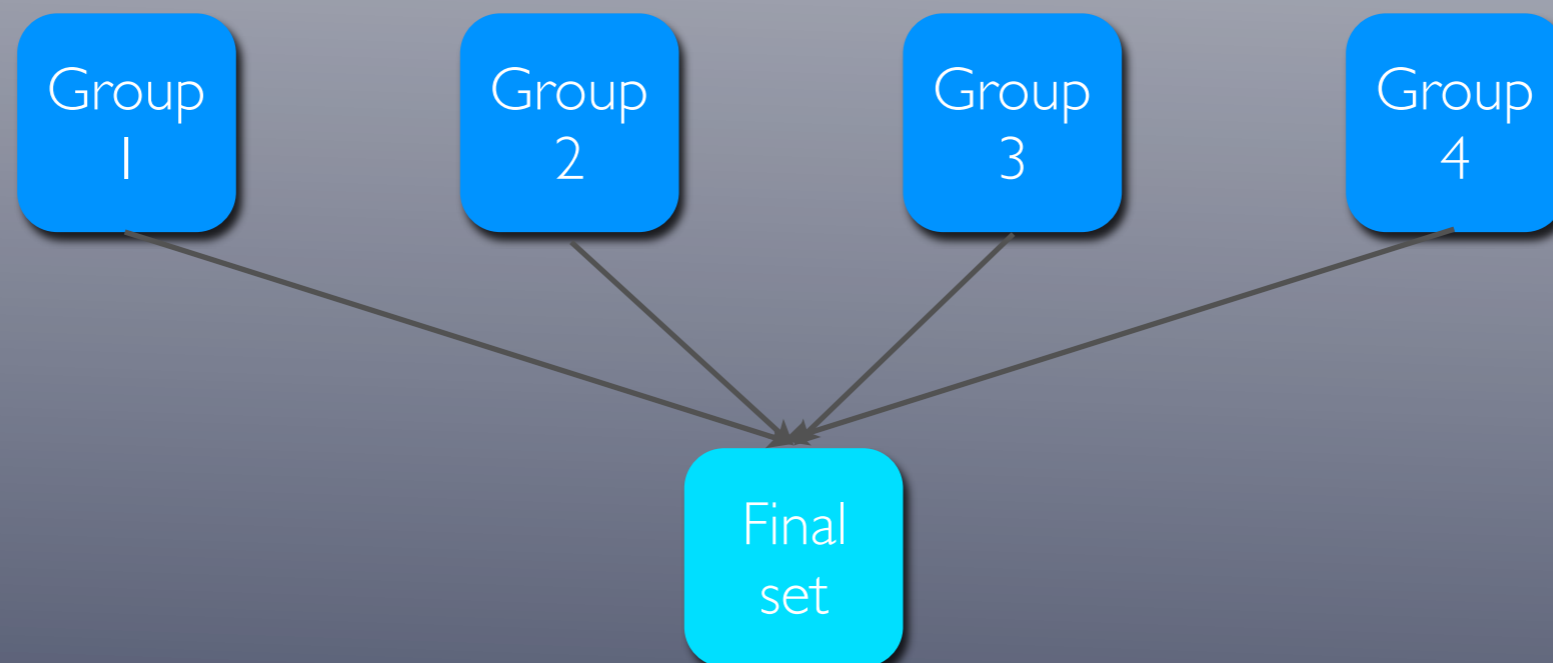
$m=3$



# REPRODUCIBILITY

- Since *EQK*-means, even if combined with an alignment stability test, does not guarantee an optimum solution (global minimum) and stable groups can be fake, we require the solution to be reproducible over a number of quasi-independent runs.
- We have  $m=4$  *EQK*-means runs analyzing the data in parallel. Once all runs produce their respective averages, we compare assignments of images to class averages and select as reproducible subsets shared among quasi-independent runs.

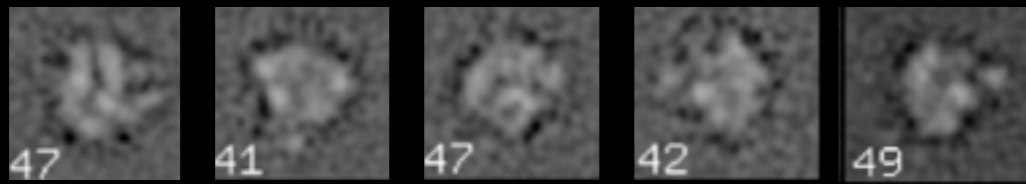
$m=4$



# ISAC: ITERATIVE STABLE ALIGNMENT AND CLUSTERING

- We use 4 CPU groups to analyze the data set simultaneously
- Irreproducible averages are eliminated

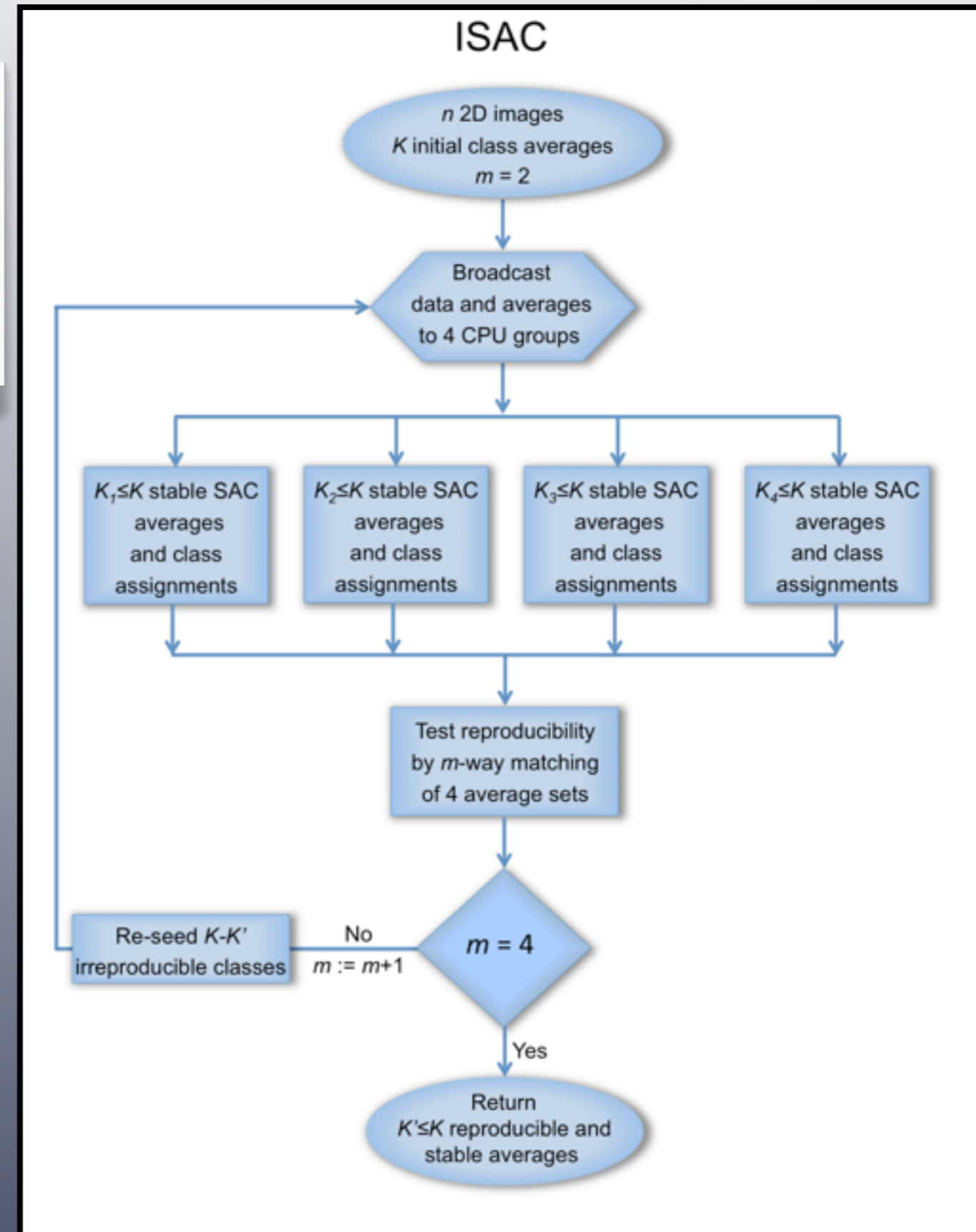
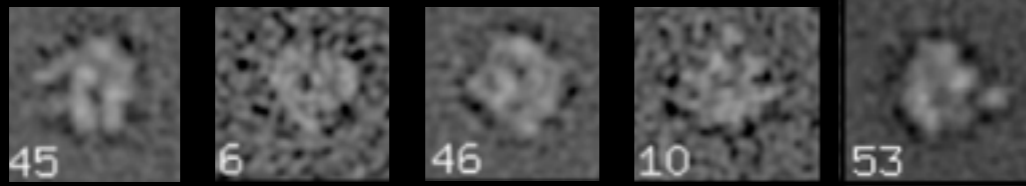
$m=2$



$m=3$



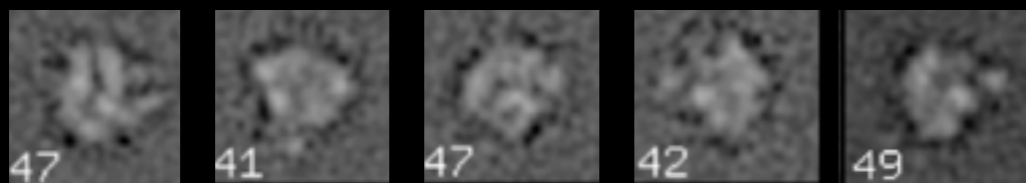
$m=4$



# ISAC: ITERATIVE STABLE ALIGNMENT AND CLUSTERING

- We use 4 CPU groups to analyze the data set simultaneously
- Irreproducible averages are eliminated

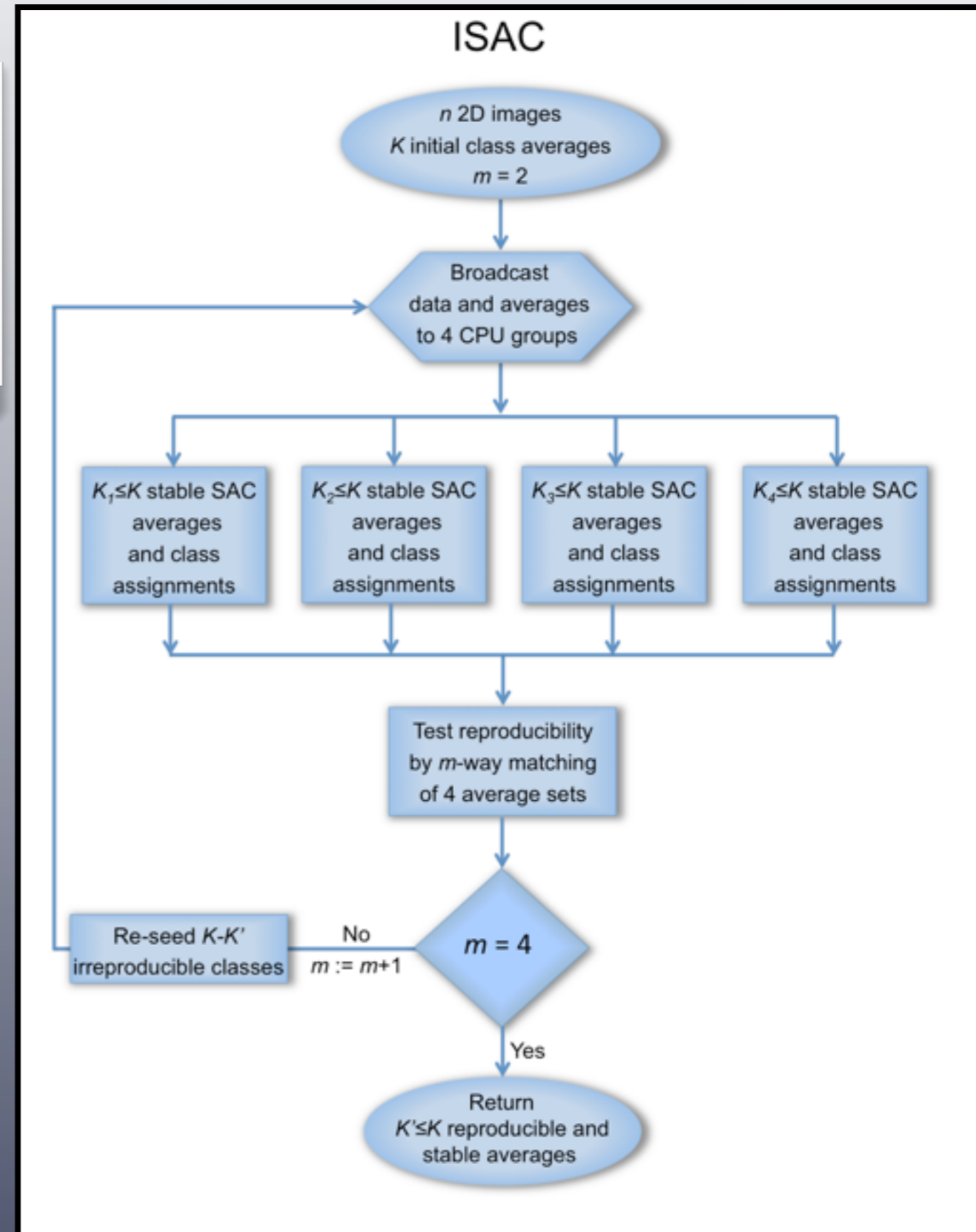
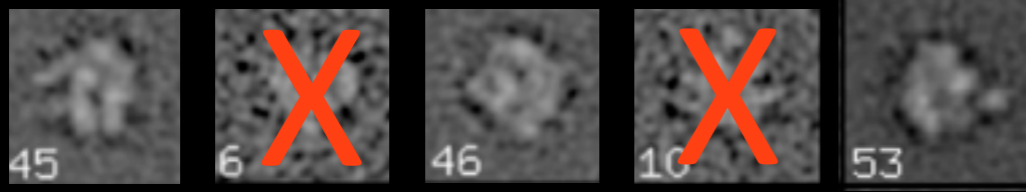
$m=2$

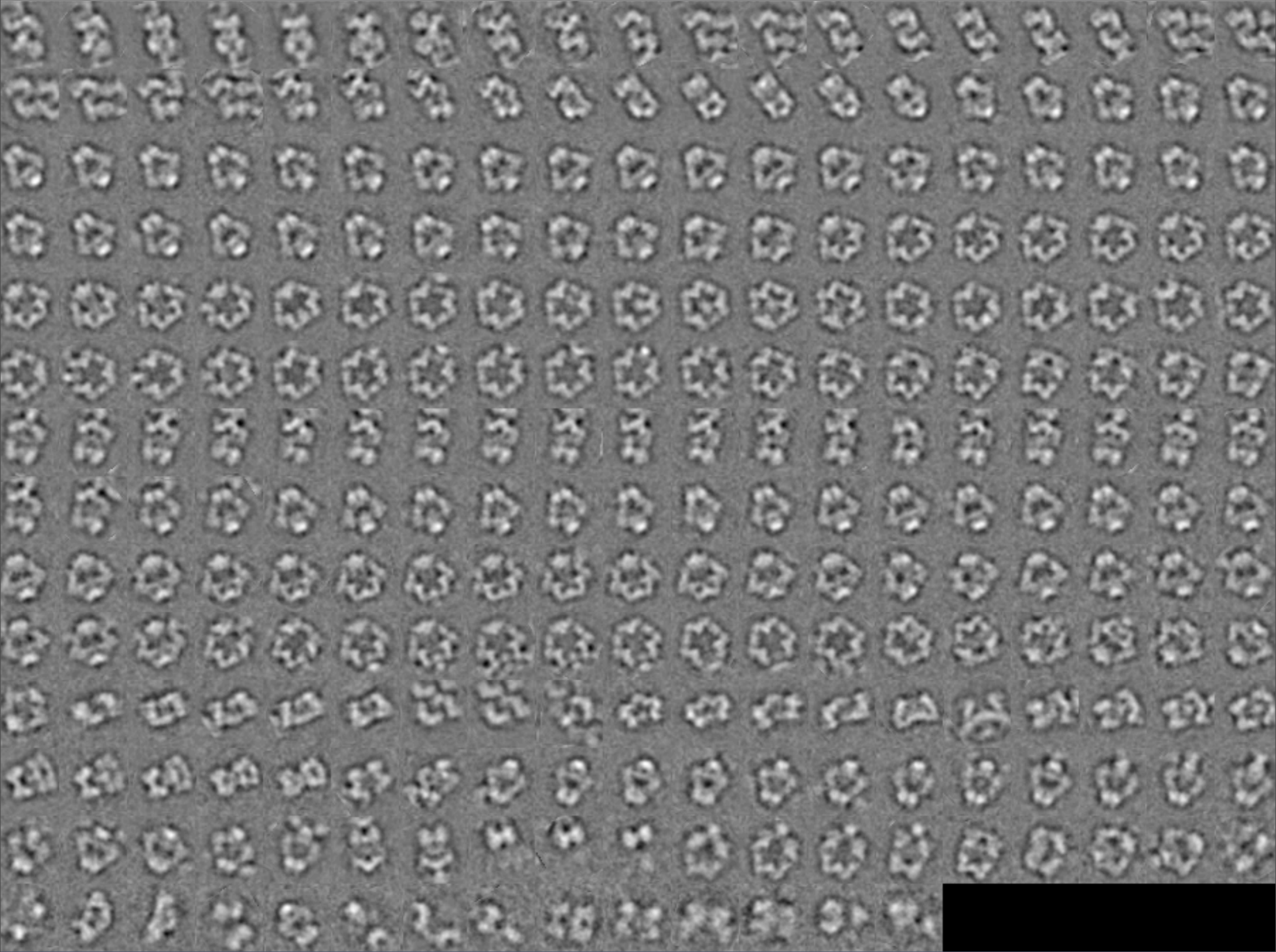


$m=3$



$m=4$





# ISAC

*Validated and reproducible class averages*

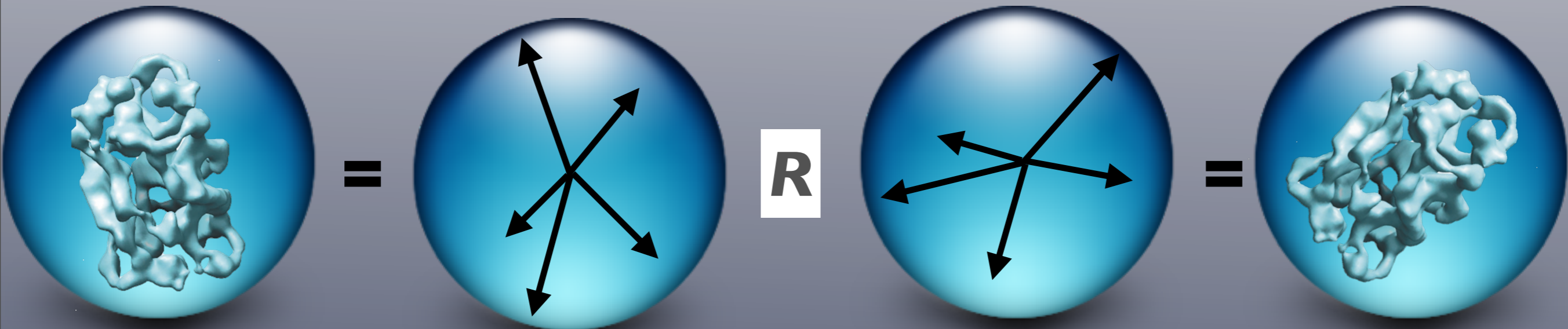


# Constructive validation: from *ab initio* EM map determination to map refinement



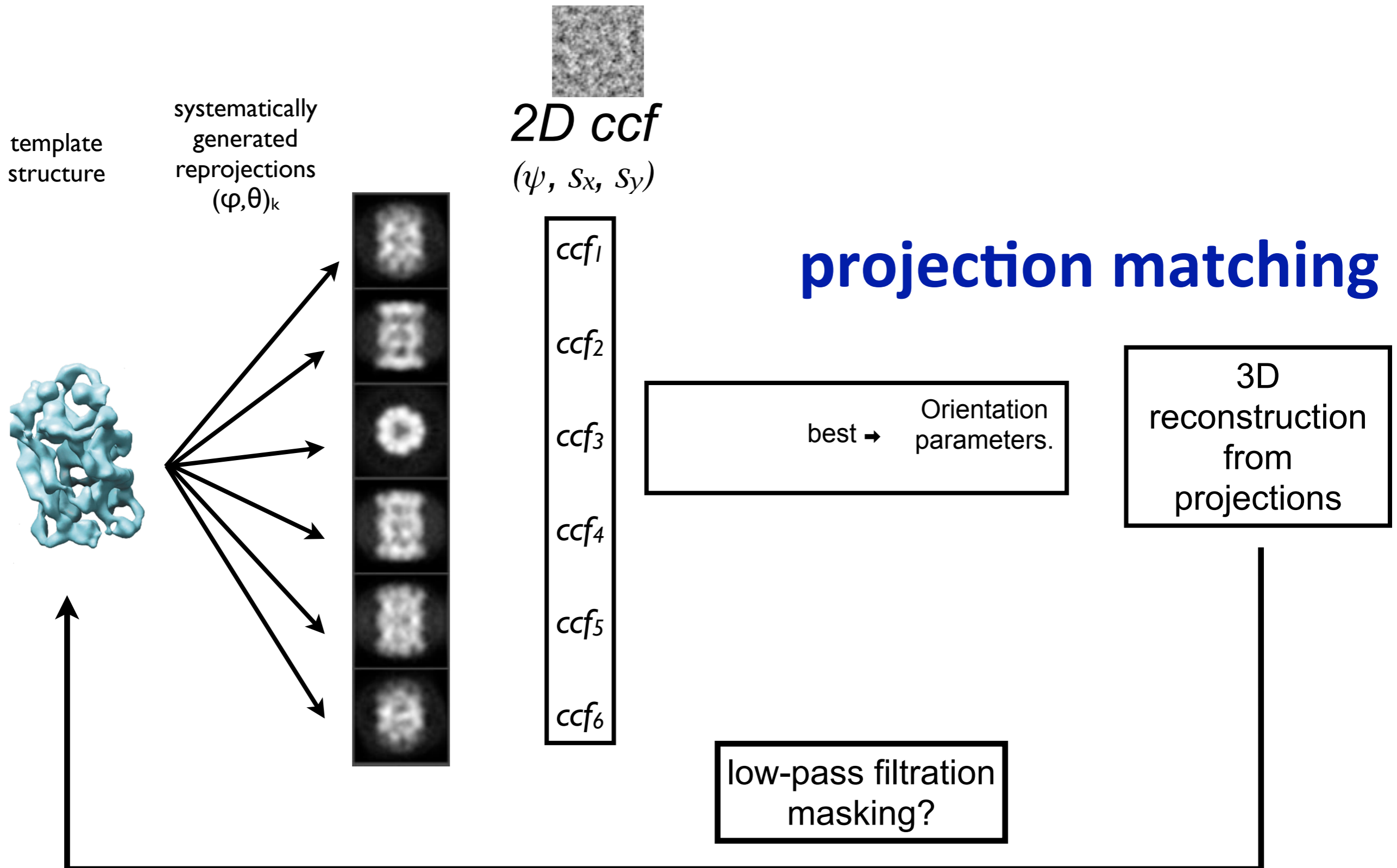
$$(\phi, \theta), \psi, S_x, S_y$$

$$\tau, \psi, S_x, S_y$$

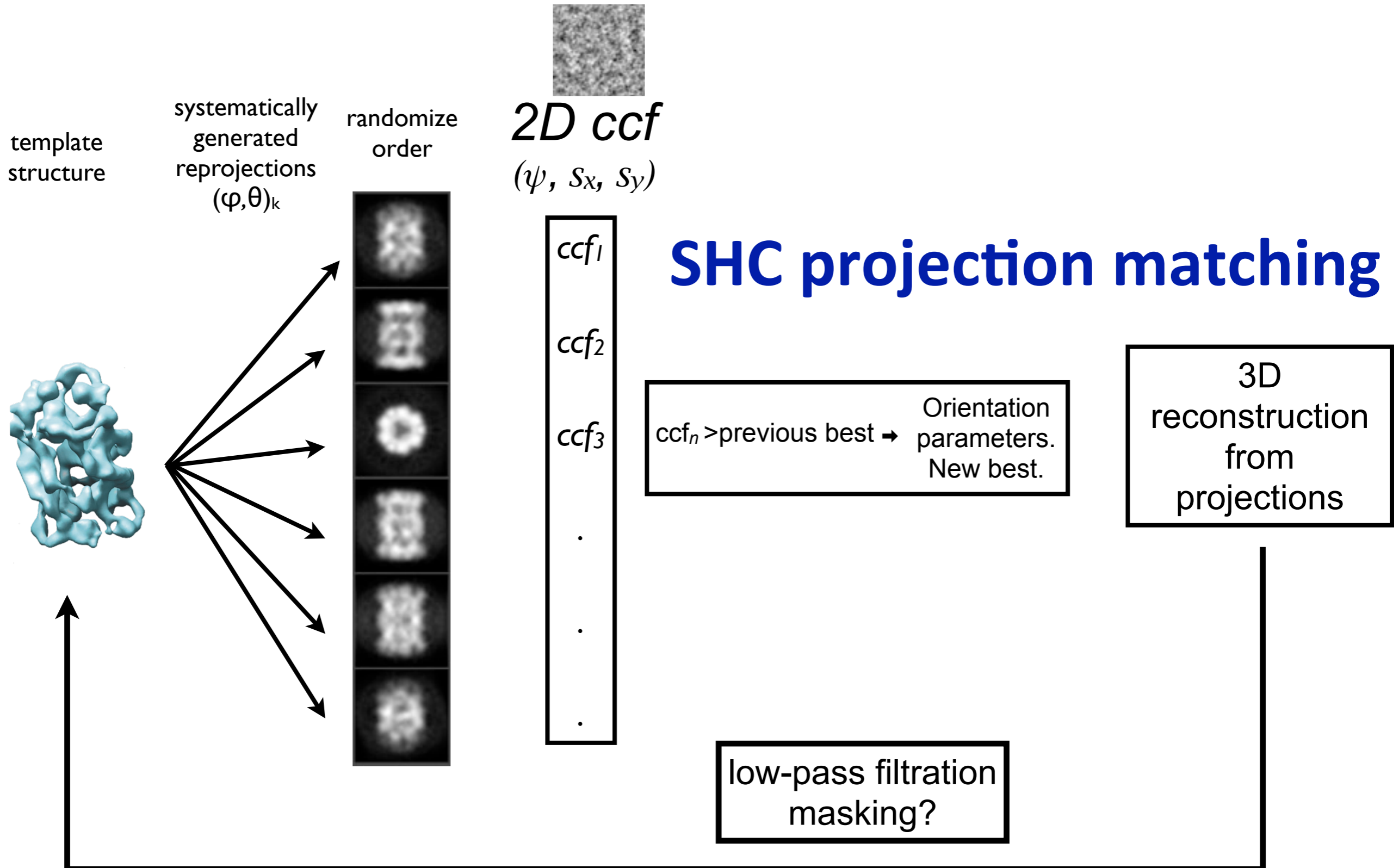


$$\max_{\mathbf{R}} \sum_{n=1}^N \tau_n^{1'} \mathbf{R} \tau_n^2$$

# STEP 1: GENERATING A MAP

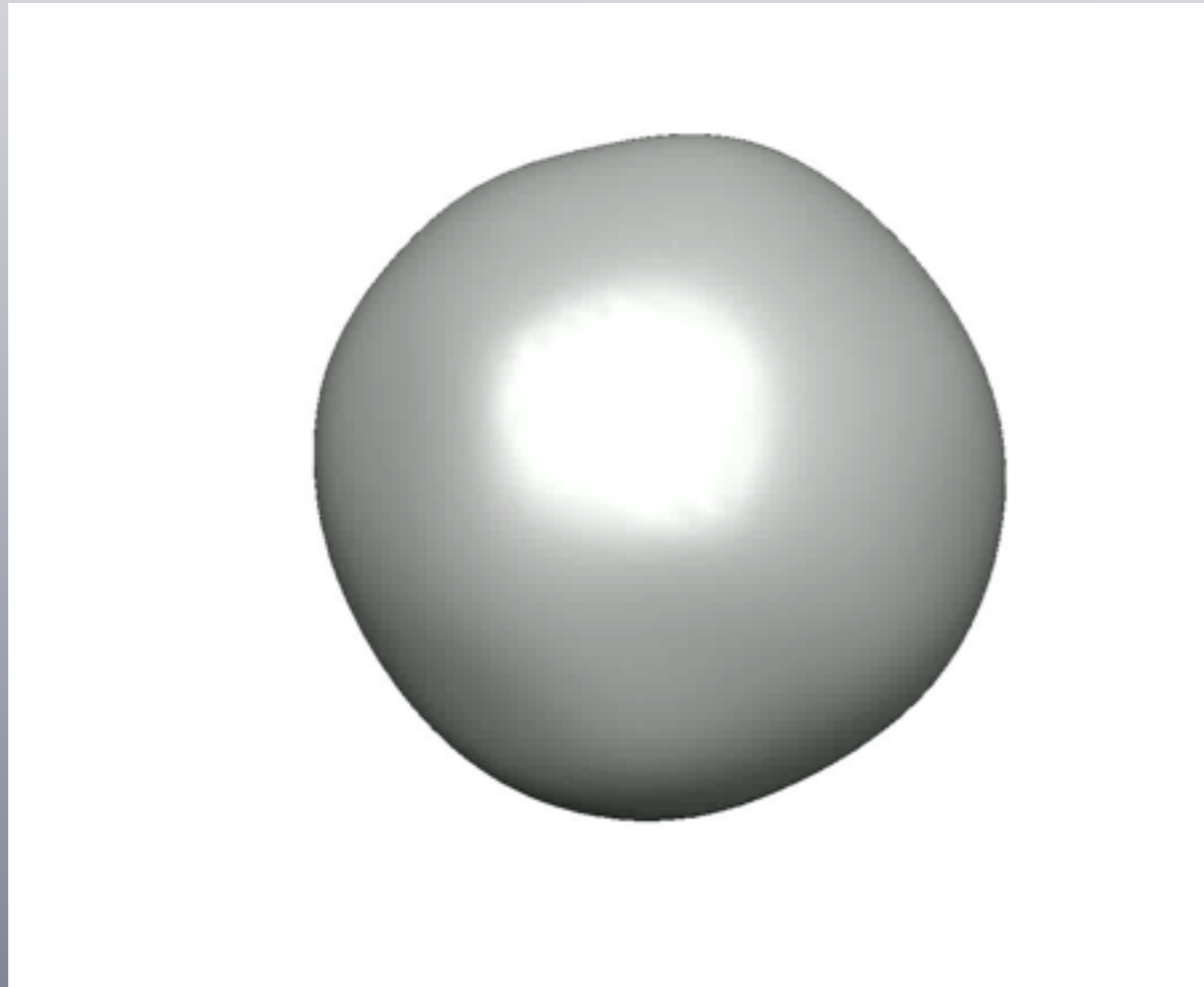


# STEP 1: GENERATING A MAP

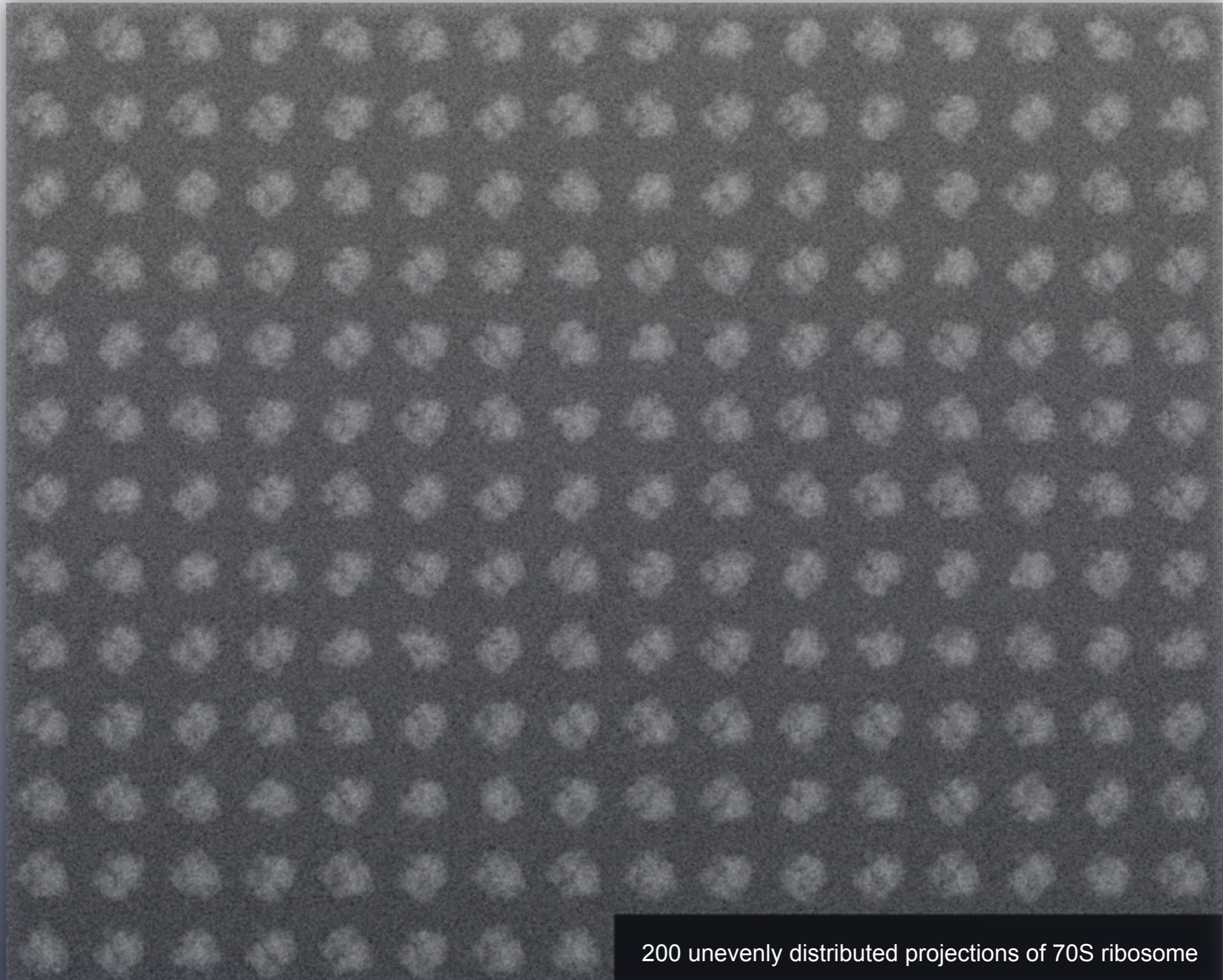


# SHC - CONVERGENCE

# SHC - CONVERGENCE

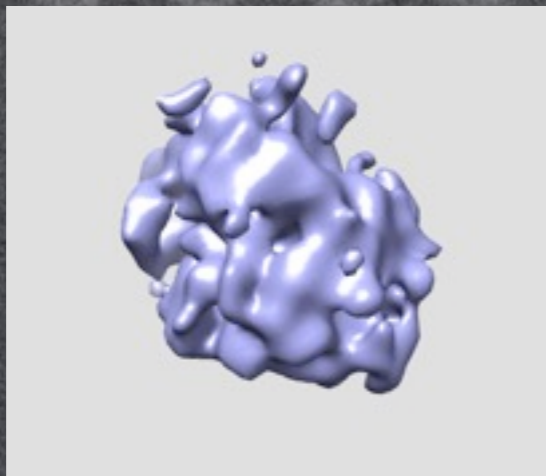
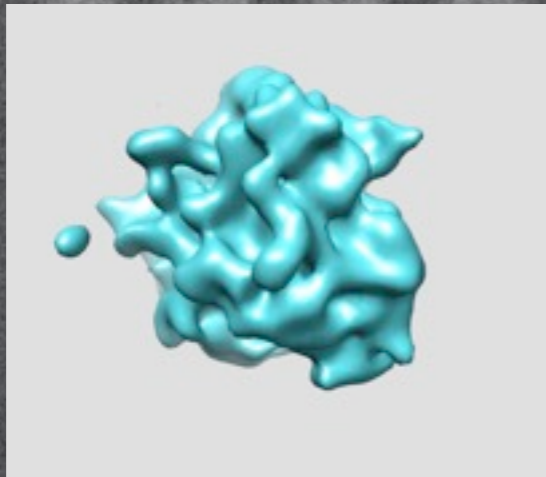
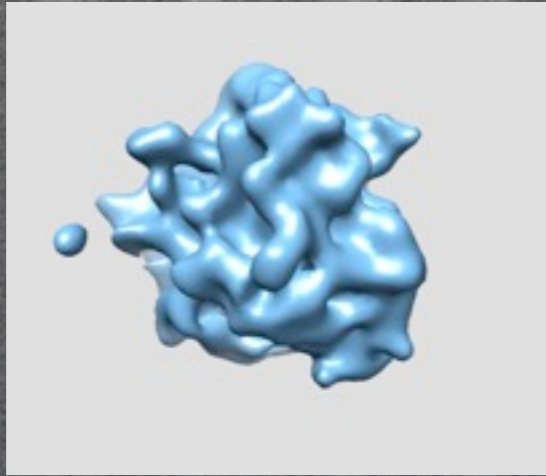


# OVERCOMING SHC CONVERGENCE LIMITATIONS BY MONITORING PARAMETER REPRODUCIBILITY



200 unevenly distributed projections of 70S ribosome

# OVERCOMING SHC CONVERGENCE LIMITATIONS BY MONITORING PARAMETER REPRODUCIBILITY



## **GOOD:**

- No bias towards the initial structure, in normal use always randomized start
- Often converges to a plausible solution
- Very good for structure refinement

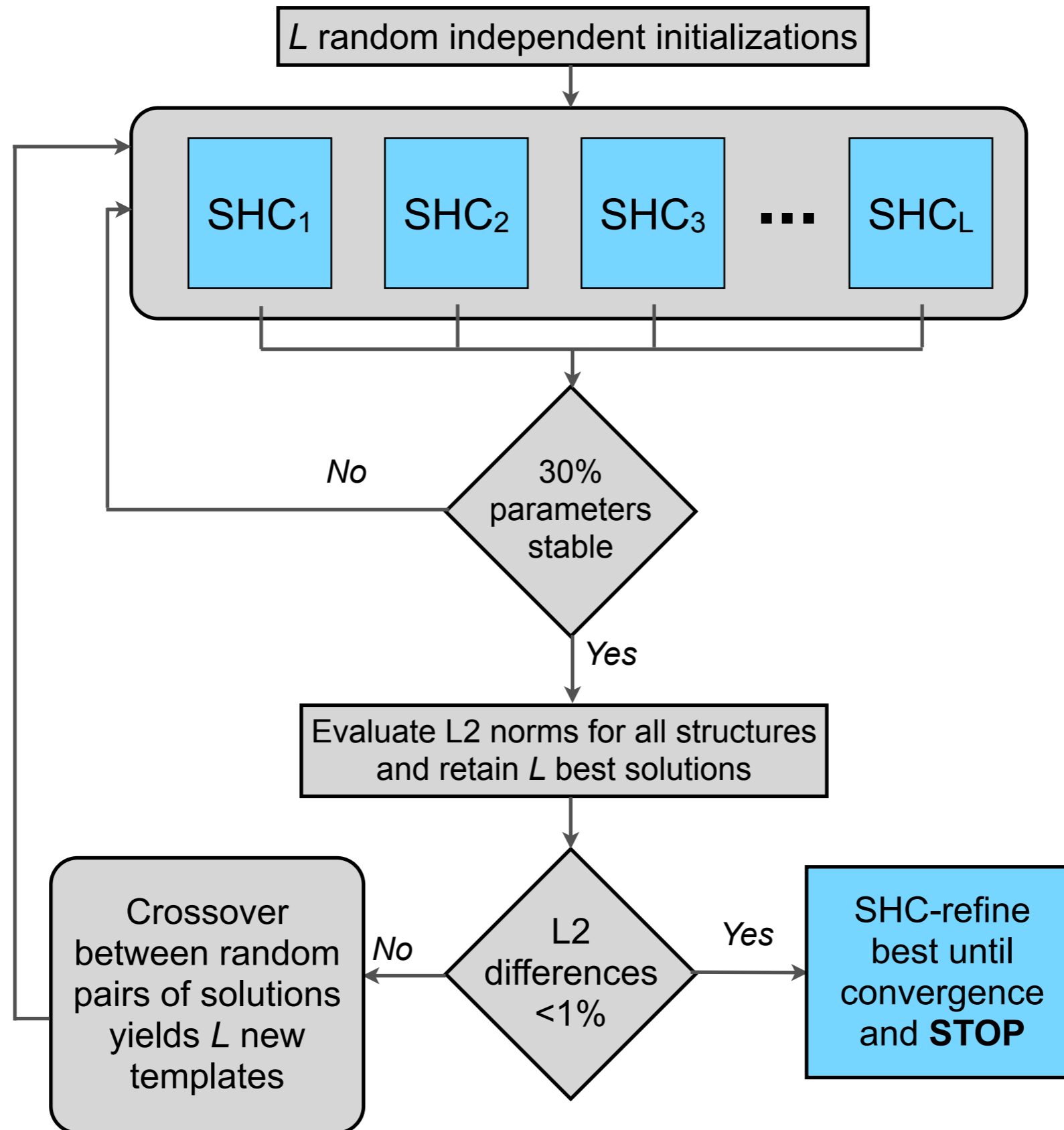
## **NOT SO GOOD:**

- Convergence properties poorly characterized/understood, unclear how often it converges and what does it depend on
- Sometimes gets stuck in a completely wrong solution
- Plausible solutions somewhat different

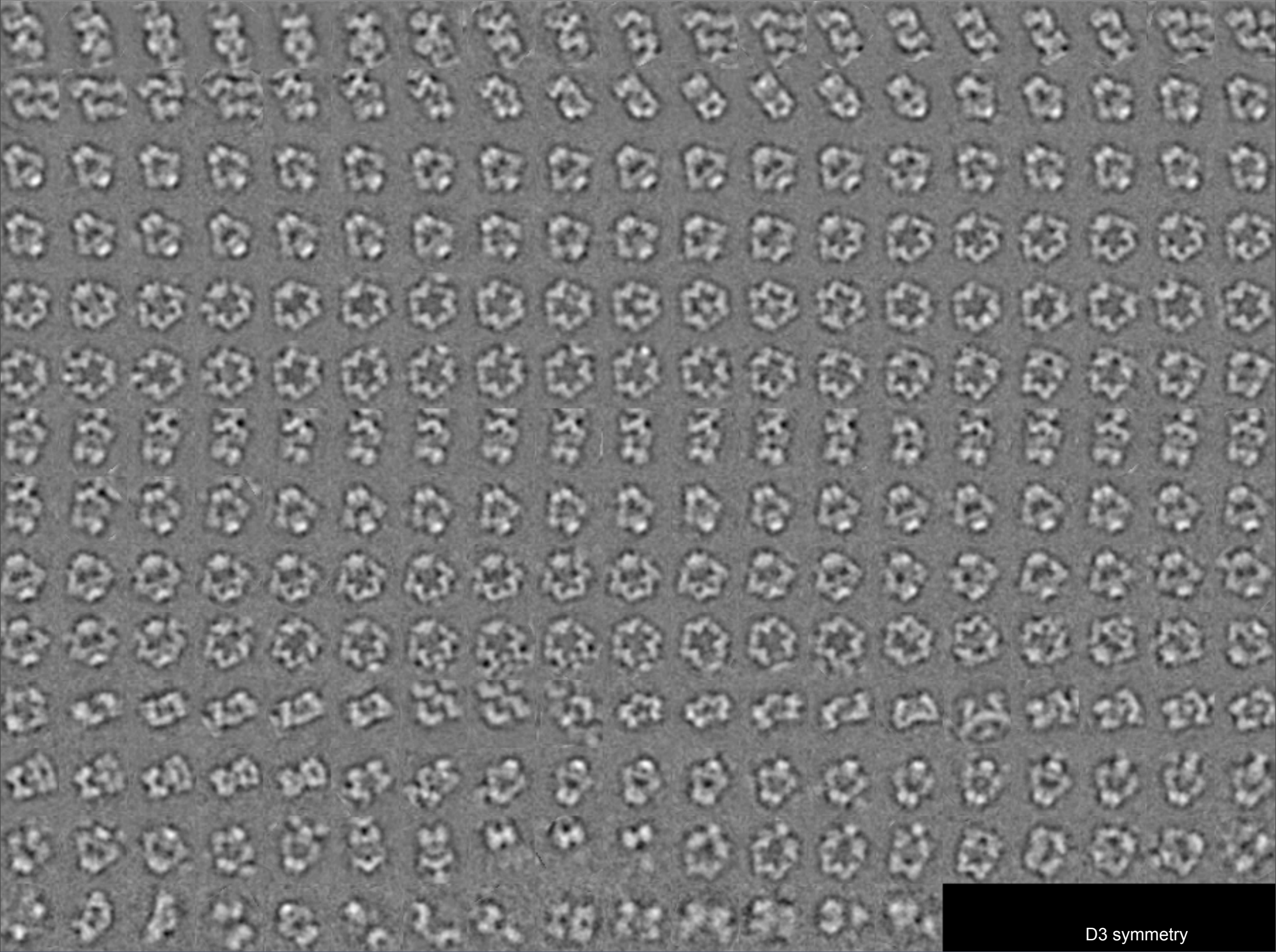
200 unevenly distributed projections of 70S ribosome

# STEP 2: VIPER

(Validation of Individual Parameter Reproducibility)

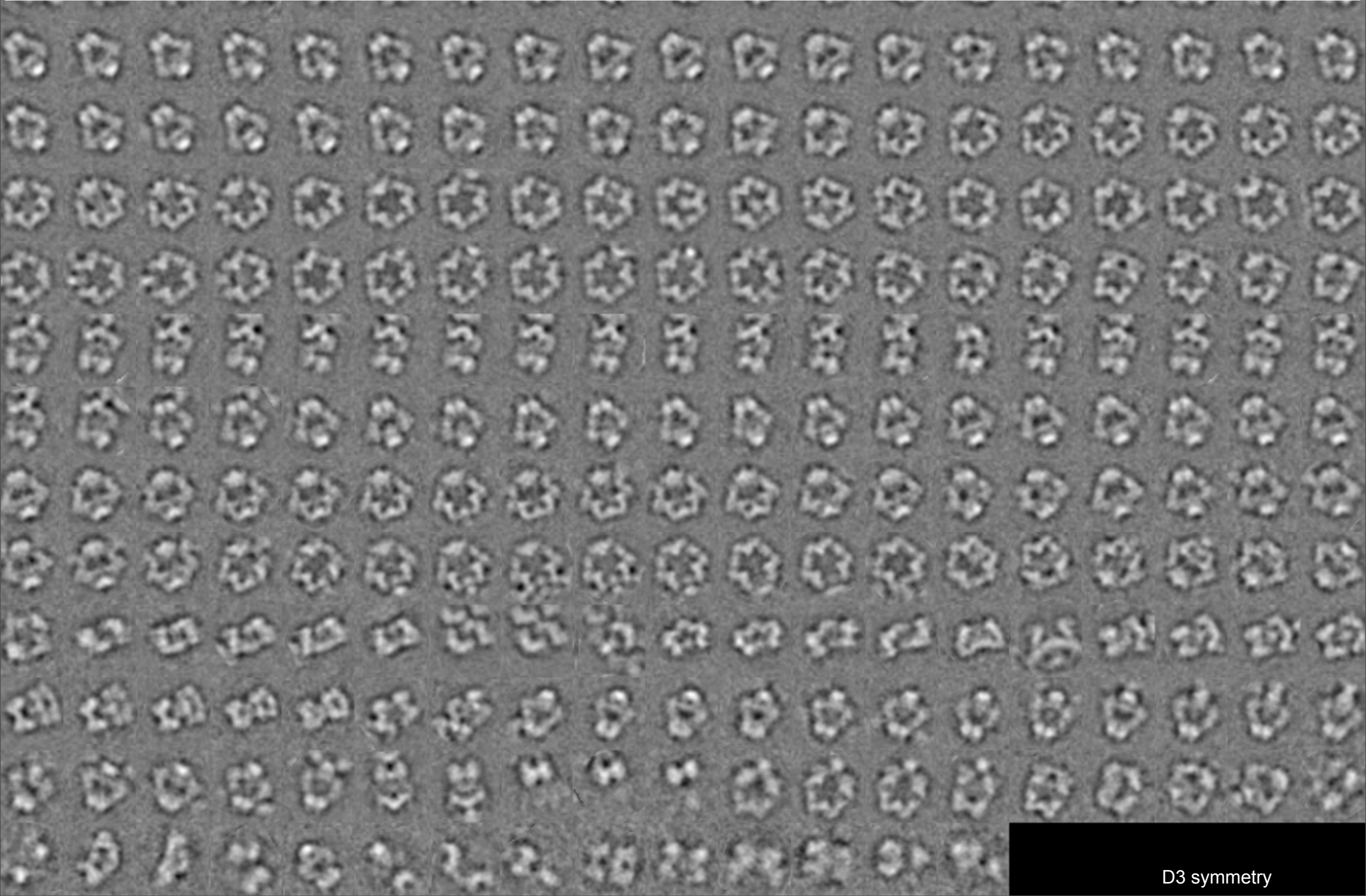






D3 symmetry

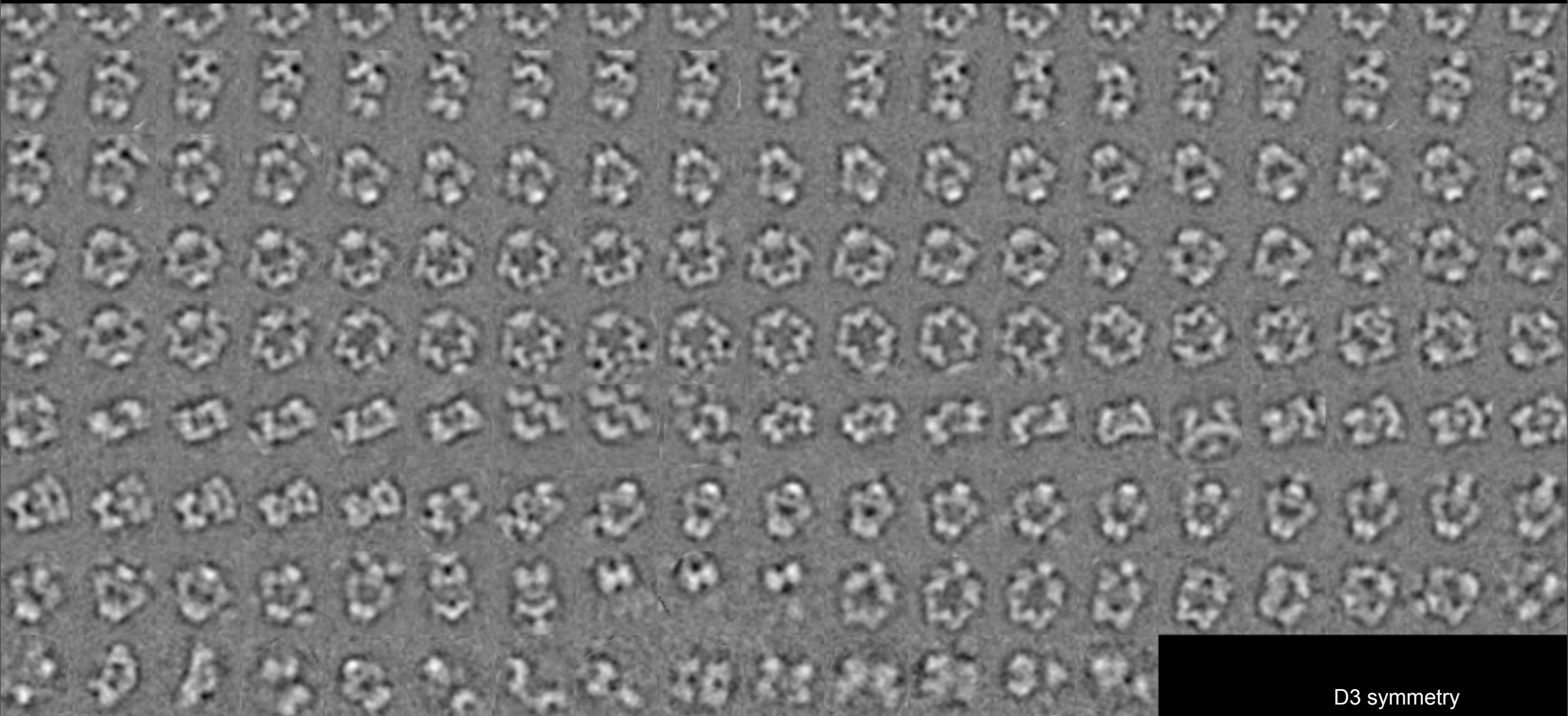
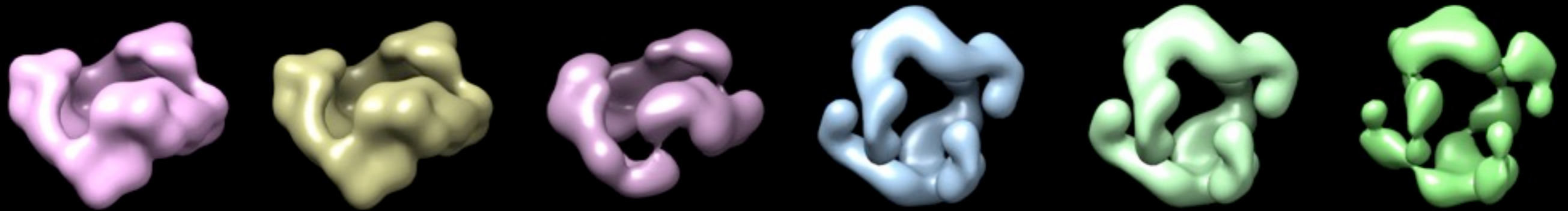
*Ab initio* structure determination with ISAC/VIPER:  
only the correct averages, only the correct structure



D3 symmetry

*Ab initio* structure determination with ISAC/VIPER:  
only the correct averages, only the correct structure

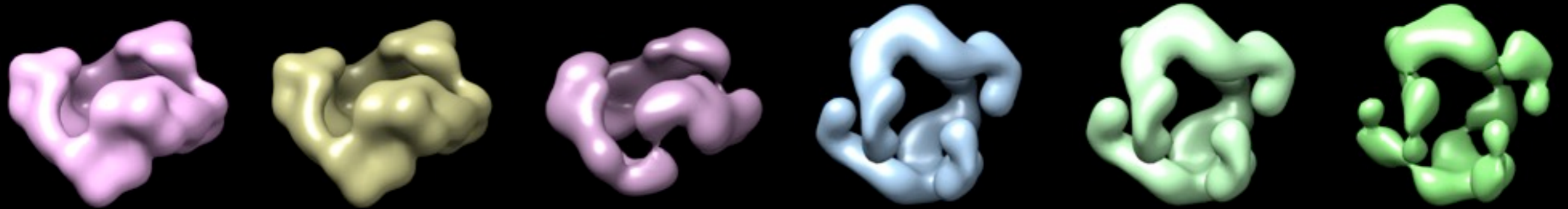
GA - generation I



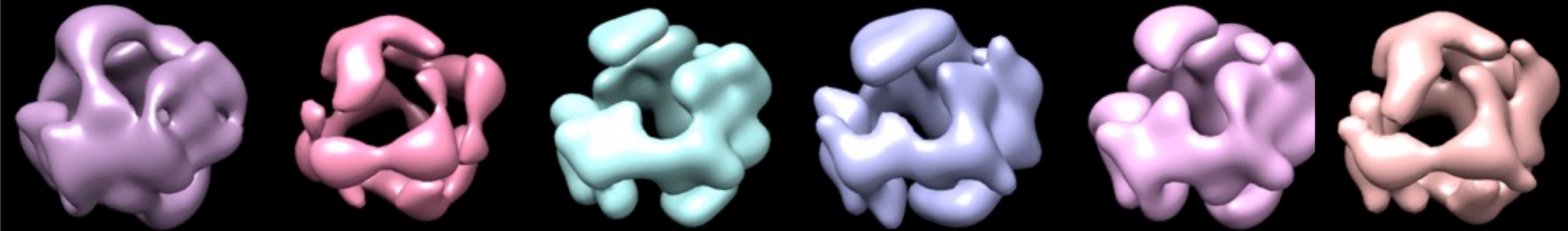
D3 symmetry

*Ab initio* structure determination with ISAC/VIPER:  
only the correct averages, only the correct structure

GA - generation I



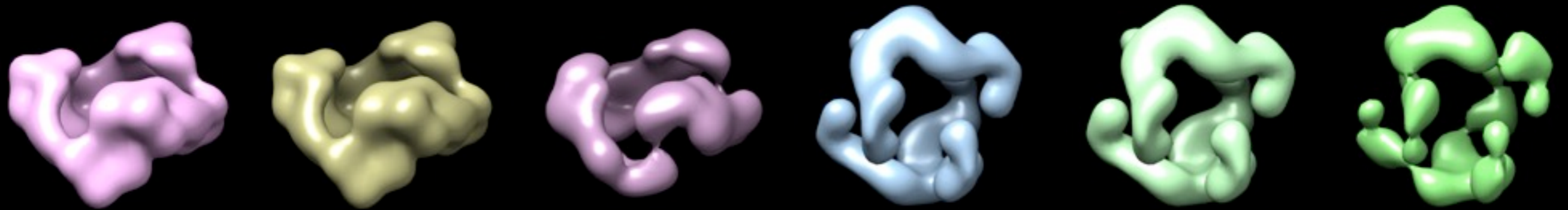
offsprings I



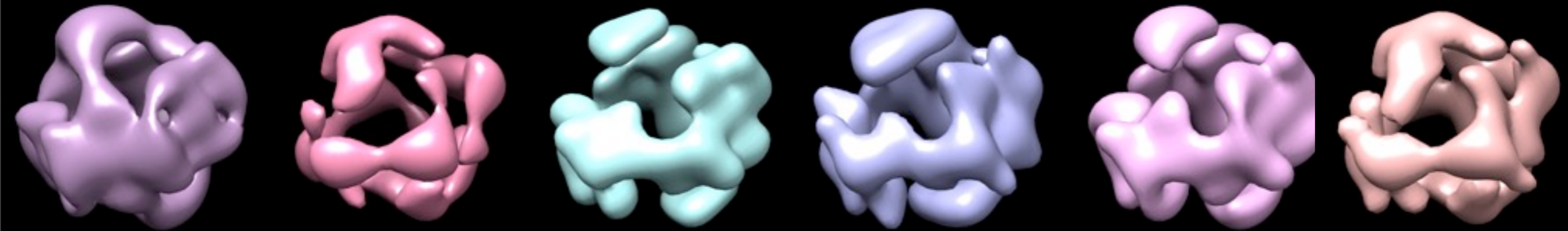
D3 symmetry

*Ab initio* structure determination with ISAC/VIPER:  
only the correct averages, only the correct structure

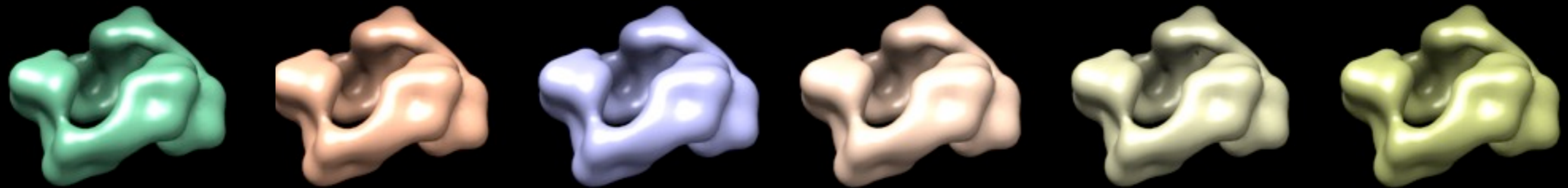
GA - generation I



offsprings I



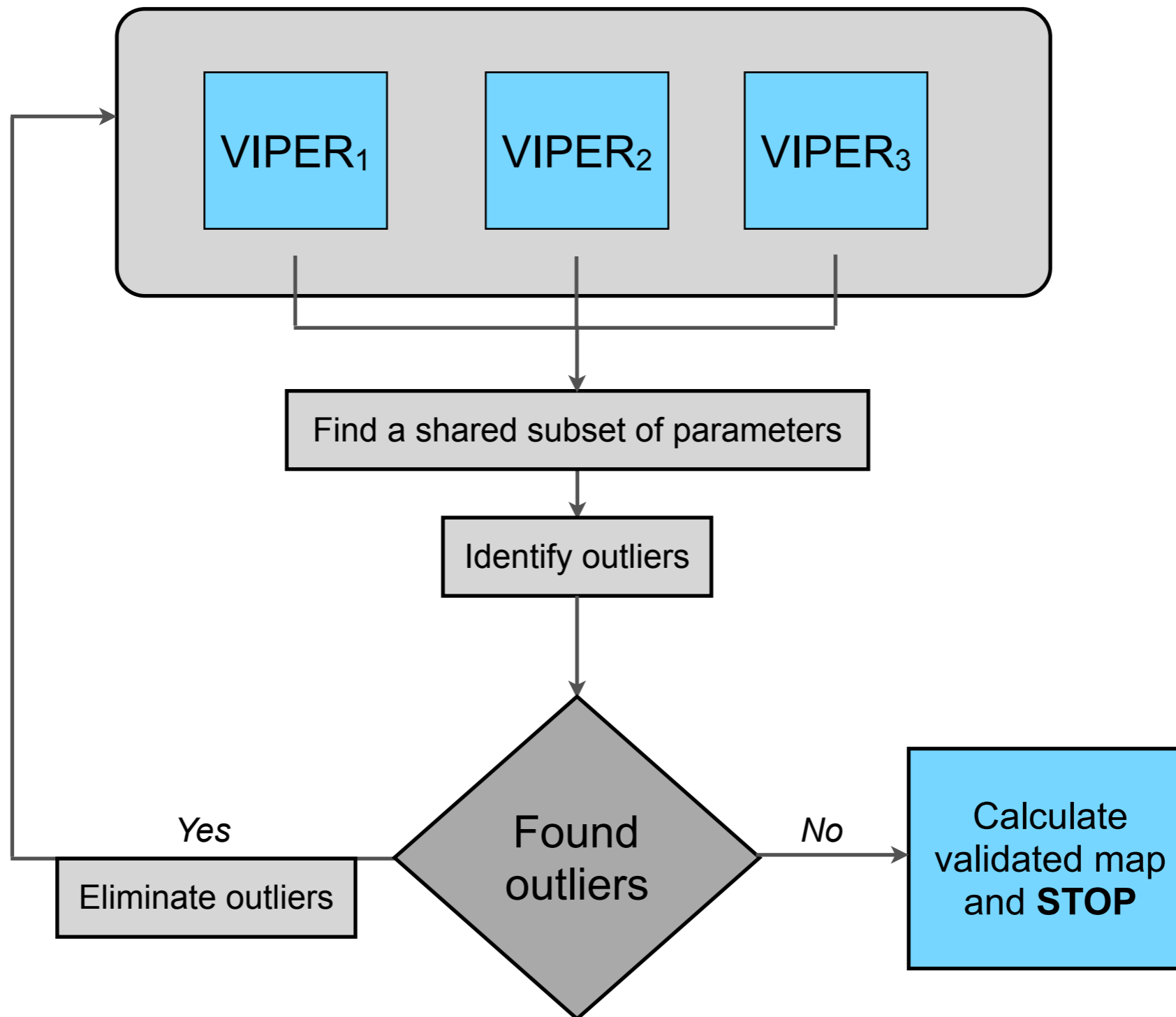
GA - generation II



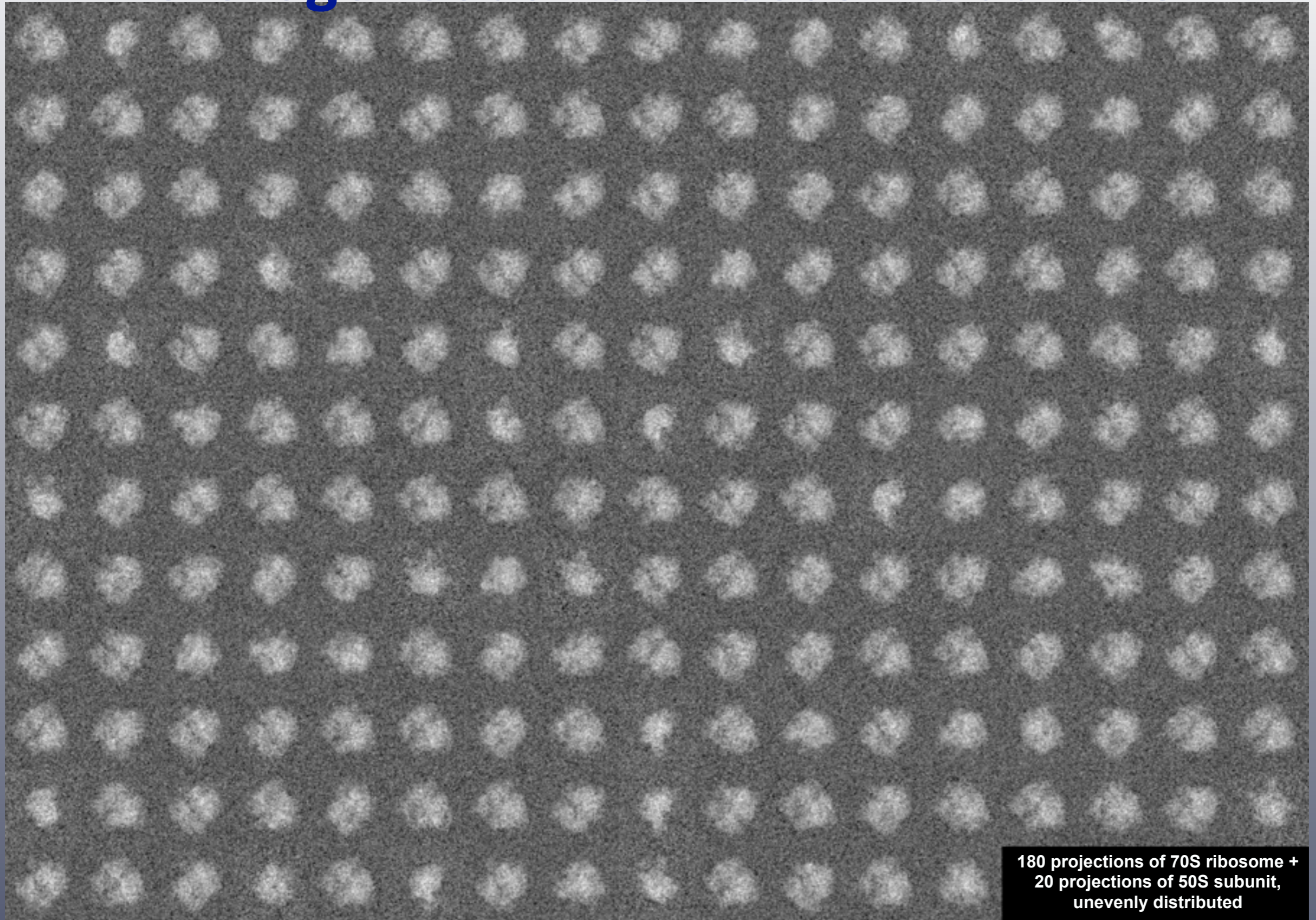
D3 symmetry

# STEP 3: VALIDATION WITH REPRODUCIBLE VIPER

## R-VIPER YIELDS A VALIDATED *ab initio* MAP

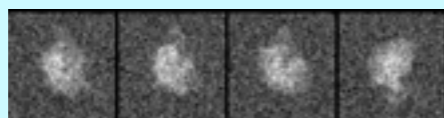
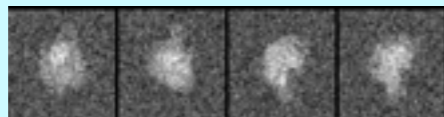
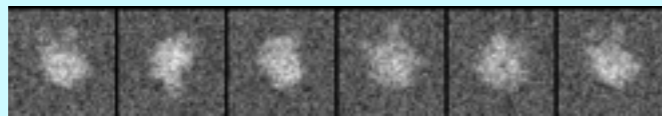
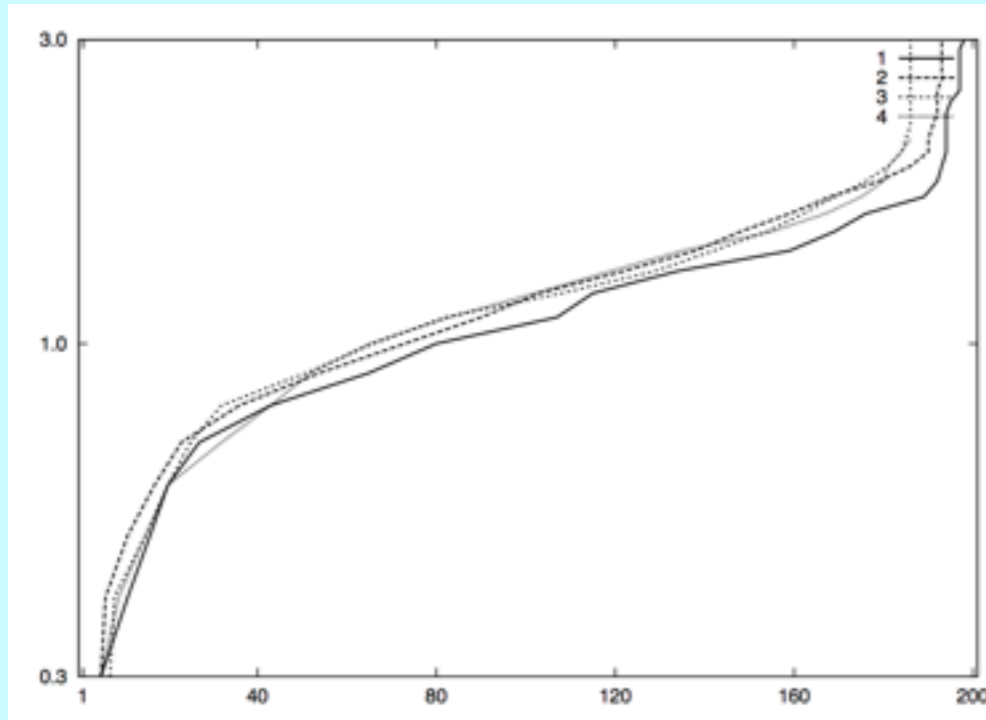


# Application of VIPER to a simulated heterogeneous 70S ribosome data set



180 projections of 70S ribosome +  
20 projections of 50S subunit,  
unevenly distributed

# Application of VIPER to a simulated heterogeneous 70S ribosome data set



180 projections of 70S ribosome +  
20 projections of 50S subunit,  
unevenly distributed