

# Ensemble reweighting using cryo-EM particles



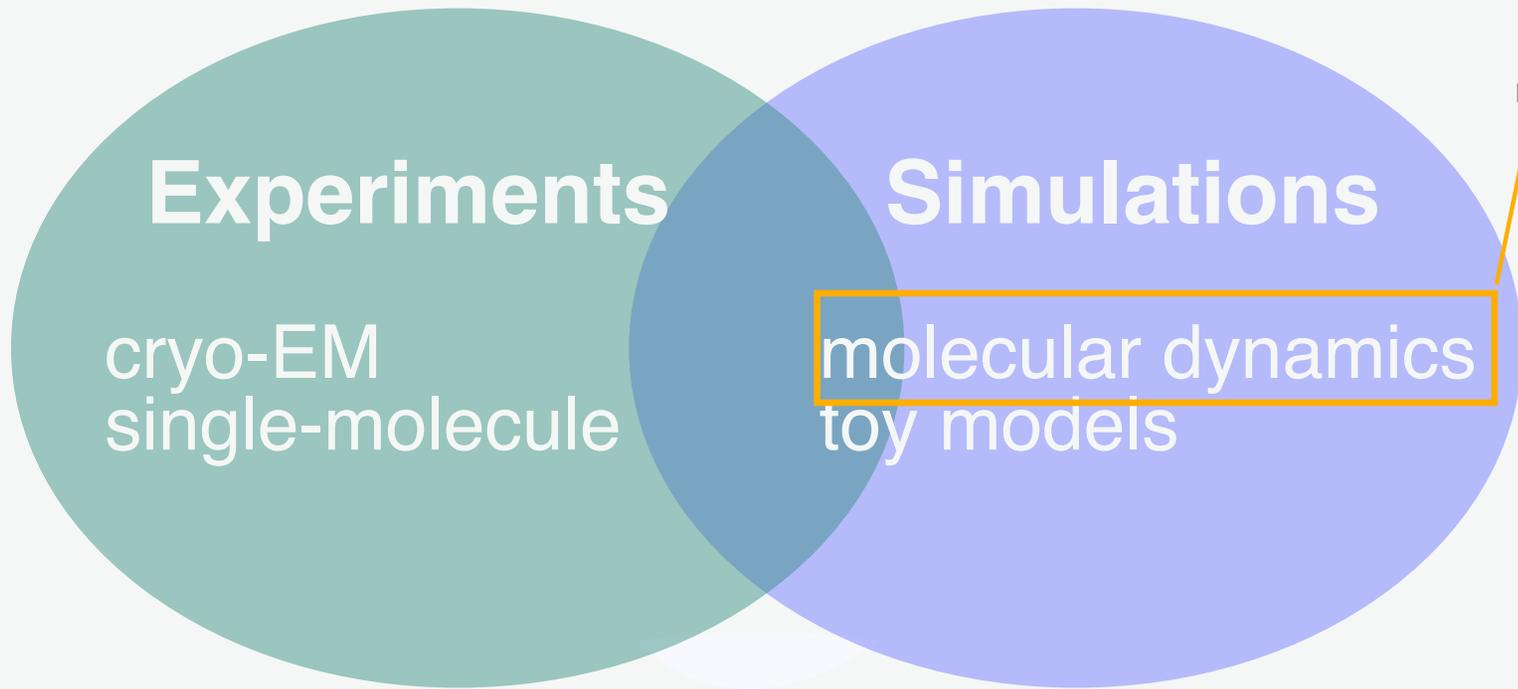
*Pilar Cossio*  
*Center for Computational Mathematics*  
*Flatiron Institute, NY*



## Outline:

1. Motivation
2. Cryo-EM as a single-molecule technique
3. Free-energies profiles along a path
4. Generalization: Ensemble Reweighting
5. Conclusions and perspectives: biasing MD...

# Structural and Molecular Biophysics\* at the Flatiron Institute, NY.



Check out:  
**Palacio-Rodriguez et al. JPCL, 2022**  
CV efficiency and rates from biased simulations

Theory, simulation, and application to extract **detailed molecular mechanisms: free energies and dynamics.**

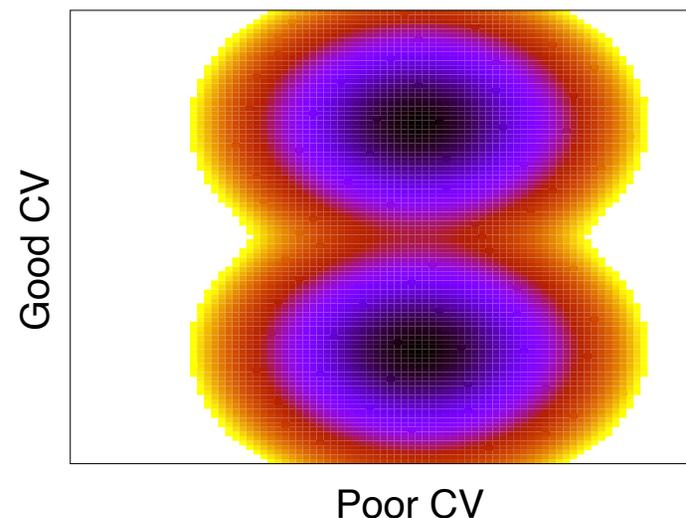
\* Co-leadership with Sonya Hanson

## Transition Rates and Efficiency of Collective Variables from Time-Dependent Biased Simulations

Karen Palacio-Rodriguez,<sup>†</sup> Hadrien Vroylandt,<sup>†</sup> Lukas S. Stelzl, Fabio Pietrucci, Gerhard Hummer, and Pilar Cossio\*

✓ Cite This: *J. Phys. Chem. Lett.* 2022, 13, 7490–7496

🌐 Read Online



↗ We introduce a measure of the **efficiency of CV** ( $\gamma$ ):

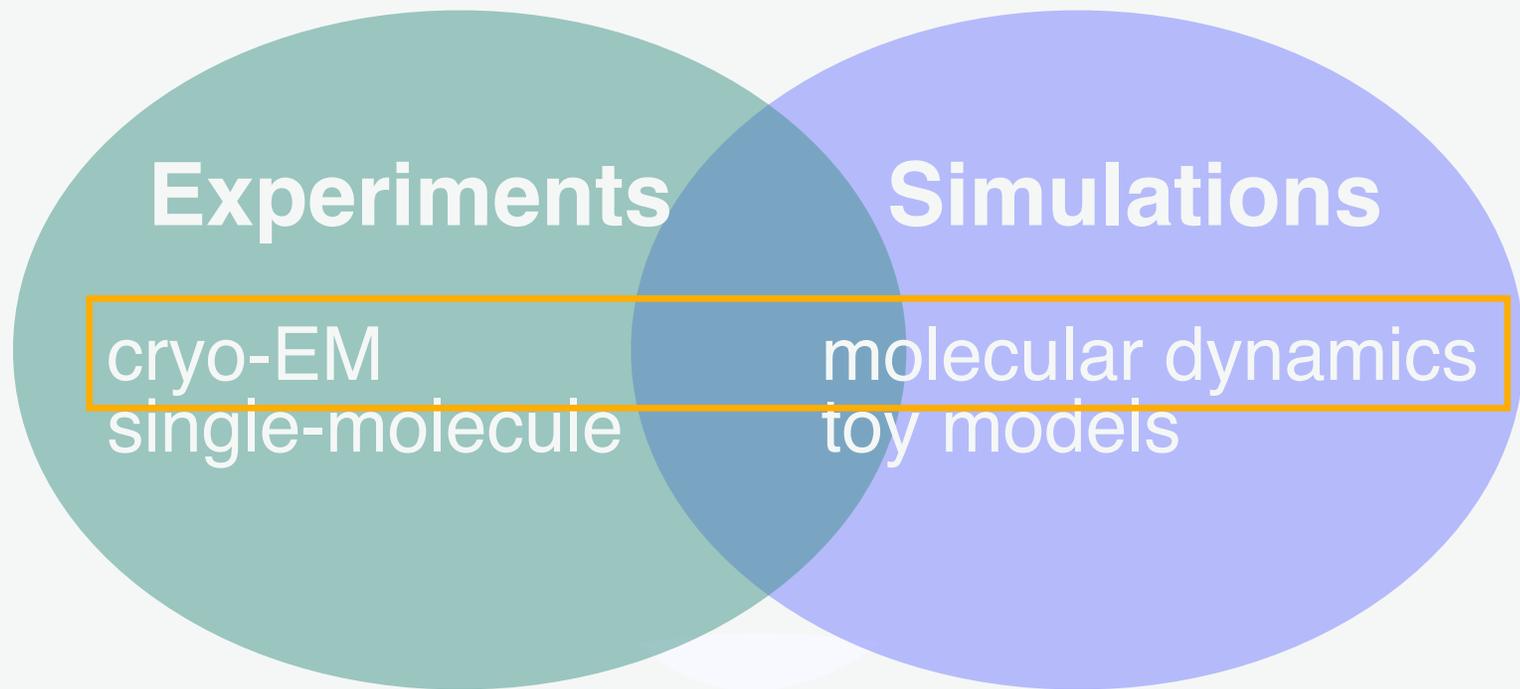
$$k(t) = k_{\text{pre}} e^{-\beta \Delta G^\ddagger + \beta \gamma V_{\text{MB}}(t)} = k_0 e^{\beta \gamma V_{\text{MB}}(t)}$$

Scales the bias

↖ ↗  
Fitting parameters:  
 $\gamma$  in  $[0,1]$  &  $\gamma=1 \rightarrow$  good CV

We can extract the unbiased rate (for any CV) and have a new measure of the efficiency of the CV.

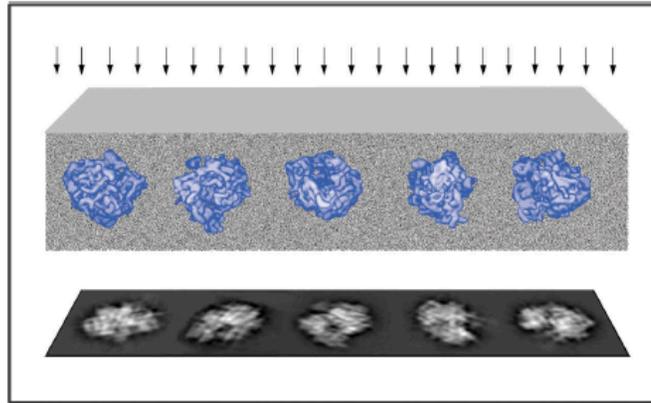
# Structural and Molecular Biophysics\* at the Flatiron Institute, NY.



Theory, simulation, and application to extract **detailed molecular mechanisms: free energies and dynamics.**

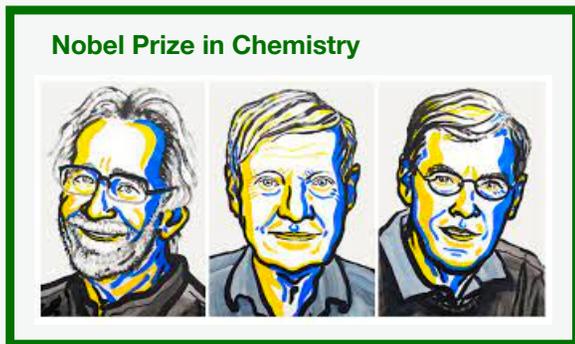
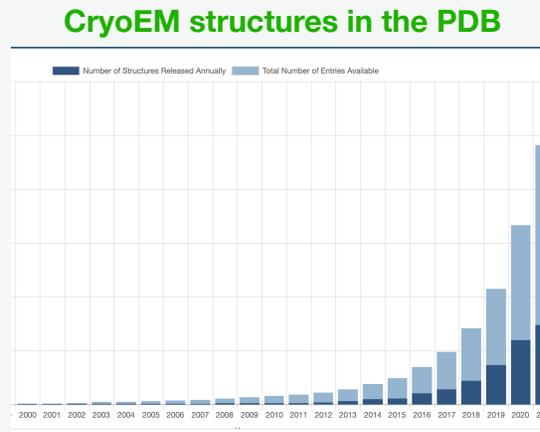
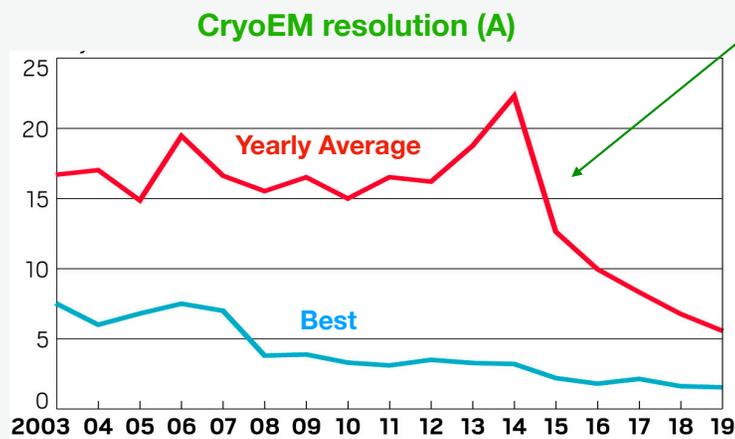
\* Co-leadership with Sonya Hanson

Cryo-EM images of biomolecules have two main unknowns: **pose and conformation**



Projection direction?  
~~Conformation?~~

# The CryoEM resolution revolution

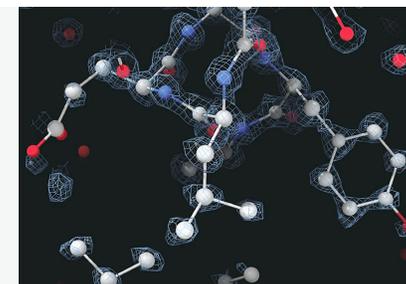


Dubochet, Frank & Henderson

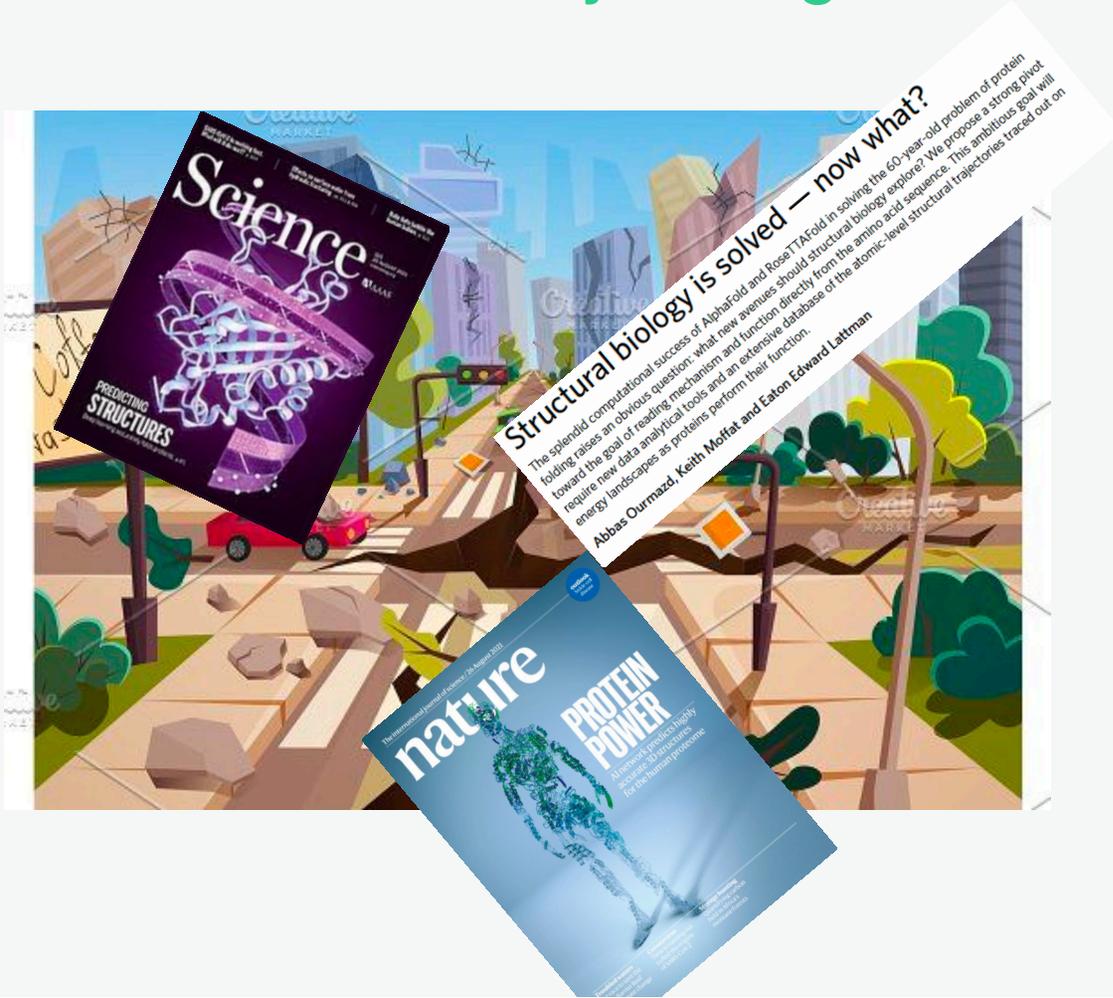
NEWS AND VIEWS | 21 October 2020

## Cryo-electron microscopy reaches atomic resolution

2020  
~1.2Å Resolution!



# However, CryoEM grounds have been shaken:

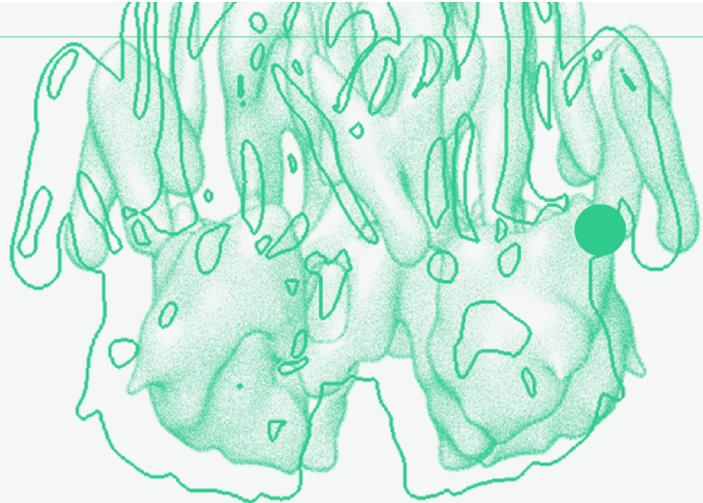
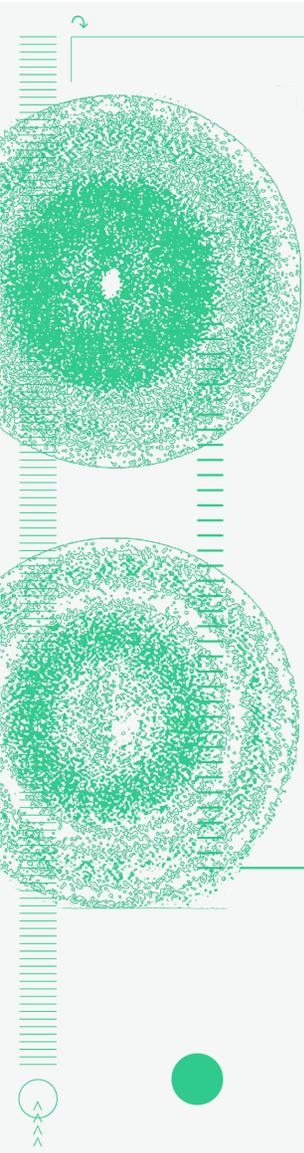


## Where is the field going?

➔ Faster, better and cheaper structure determination.

➔ Conformational variability, free energies and environment

➔ *In situ*



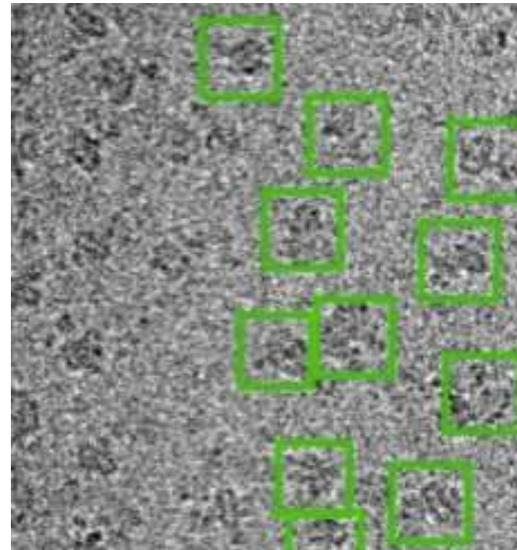
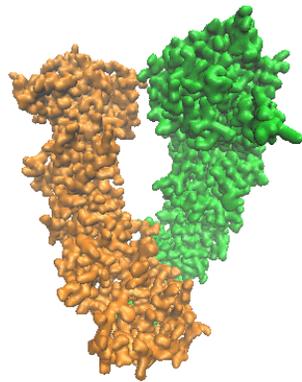
# CryoEM: conformational variability, free energies and environment



- ○ ○ Conformational heterogeneity from single-particle cryo-EM

**Freezing is done very rapidly\*.**

It's possible to trap individual conformations!!



**More information than just an average.**

\* Bock & Grubmuller, Nat Commun. 2022

# Dealing with cryoEM heterogeneity

## A typical case:

### EMPIAR-10278

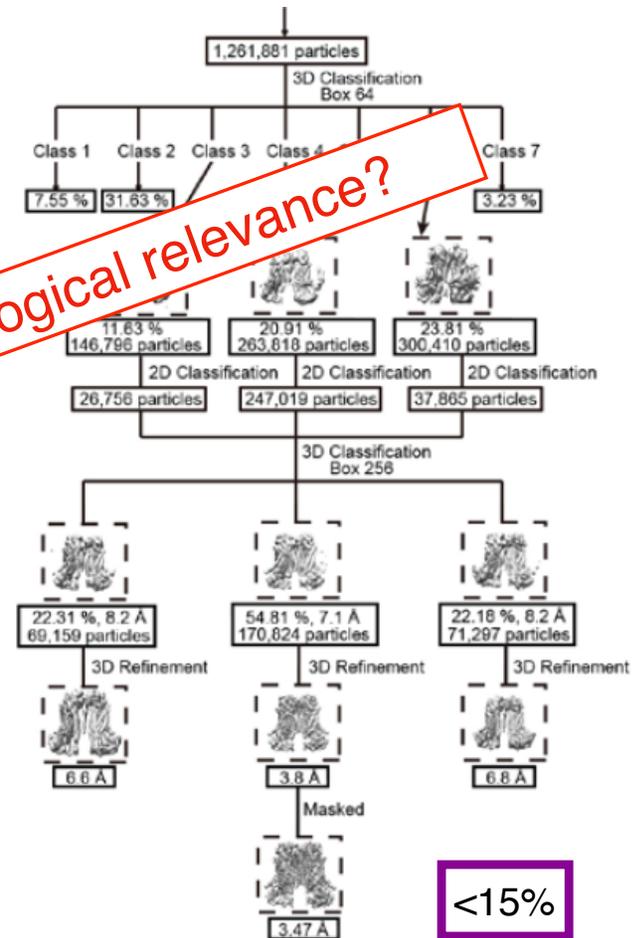
#### Cryo-EM structure of TMEM16F in digitonin with calcium bound

Publication:

Cryo-EM Studies of TMEM16F Calcium-Activated Ion Channel Suggest Features Important for Lipid Scrambling  
Feng S <sup>id</sup>, Han TW <sup>id</sup>, Ye W <sup>id</sup>, Jin P <sup>id</sup>, Cheng T <sup>id</sup>,  
Li J <sup>id</sup>, Jan LY <sup>id</sup>, Cheng Y <sup>id</sup>  
*Cell Rep* 28 567-579.e4 (2019)

Contains:

 picked particles



3D classification  
filter-out classes and  
reconstruct high-resolution  
density maps

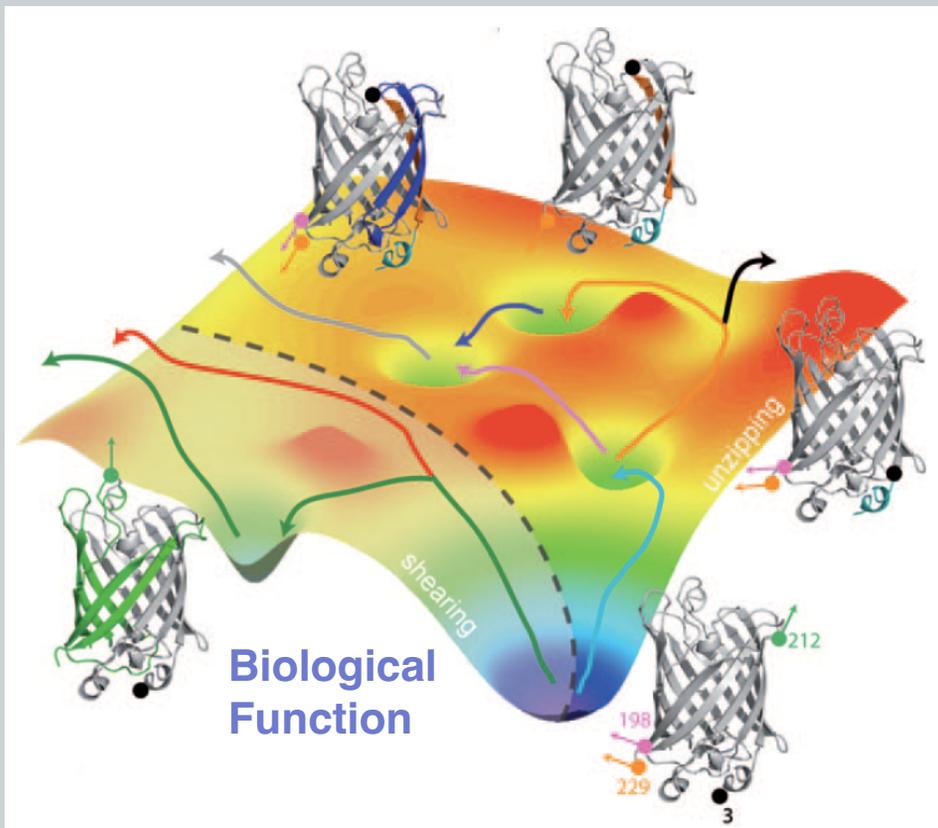
Why discard so much data? Biological relevance?

<15%



# Extracting free-energies from cryo-EM (?)

*Related to the probability distribution*



Bertz et al. *Angew. Chem.* 2008, 47, 8192–8195

The free energy (FE) is related to the probability distribution of the configurations.

A configuration  $x \in \mathbb{R}^{3N}$  and the Boltzmann factor gives the probability

$$\rho(x) = \frac{1}{Z} e^{-\beta H(x)}$$

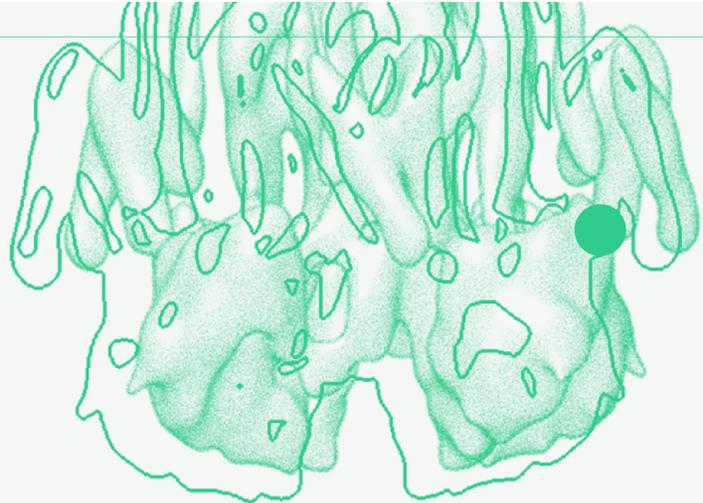
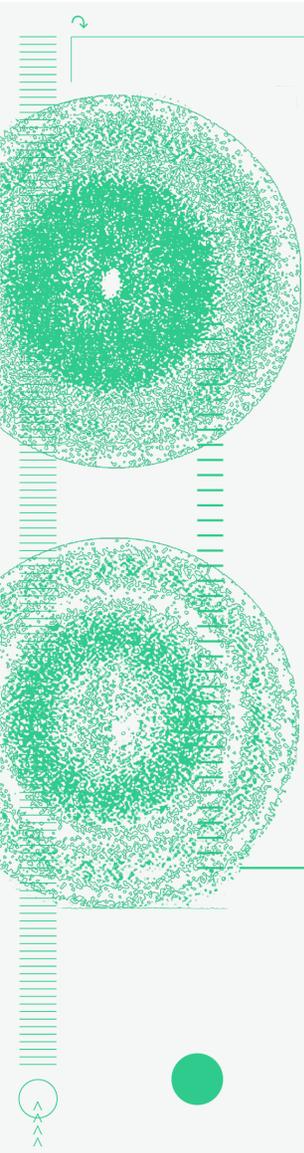
Hamiltonian

Too complex! Project along some collective variables ( $s$ ):

$$\rho(s) = \int \delta(S(x) - s) \rho(x) dx = \frac{1}{Z_1} e^{-\beta G(s)}$$

Free energy



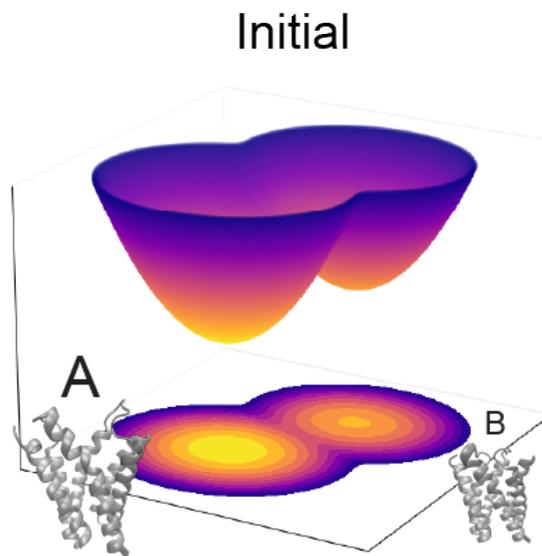


# General framework



# Free energies from CryoEM

Inference of free energies given cryoEM particles.



Conformational space  $x \in \mathbb{R}^{3\Lambda}$

Ensemble of models (e.g.  
AlphaFold, MD trajectory)

# Theory

## Bayes' theorem

$$p(H|Y) \propto p(H)p(Y|H)$$

we would like to know

**Too complex!**

a set of *iid* cryo-EM particles  
(not an averaged observable)

$$p(H|Y) \propto p(H) \prod_i \int \underbrace{p(y_i|x_i)}_{\text{Likelihood of image } y_i \text{ generated by conformation } x_i^*} \underbrace{p(x_i|H)}_{\text{Probability of } x_i} dx_i$$

Likelihood of image  $y_i$   
generated by  
conformation  $x_i^*$

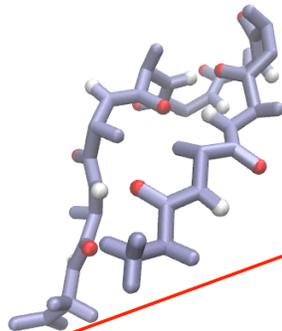
Probability of  $x_i$

data



CryoEM  
images  
 $Y=\{y_i\}$

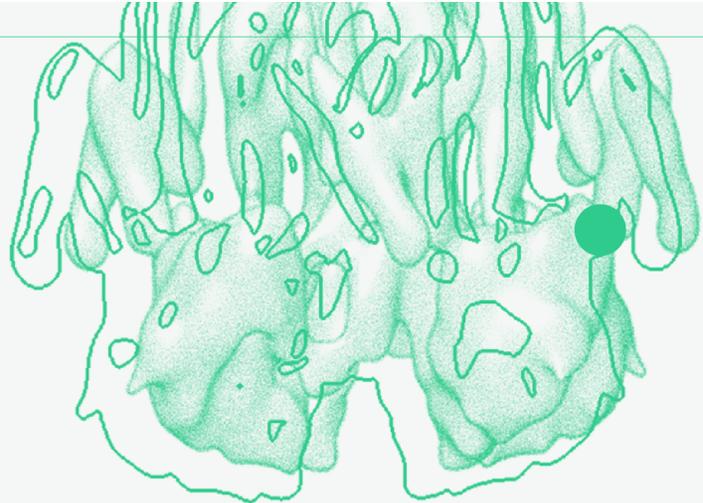
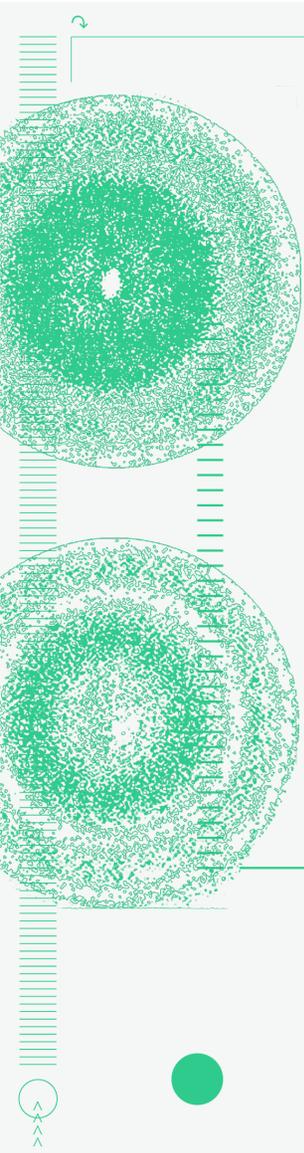
Conformation  
 $x \in \mathbb{R}^{3N}$



distribution

$$p(x|H) = \frac{1}{Z} e^{-\beta H(x)}$$

Hamiltonian



# Free-energy profile along a path



# Cryo-BIFE

## Cryo-EM Bayesian Inference of Free Energy

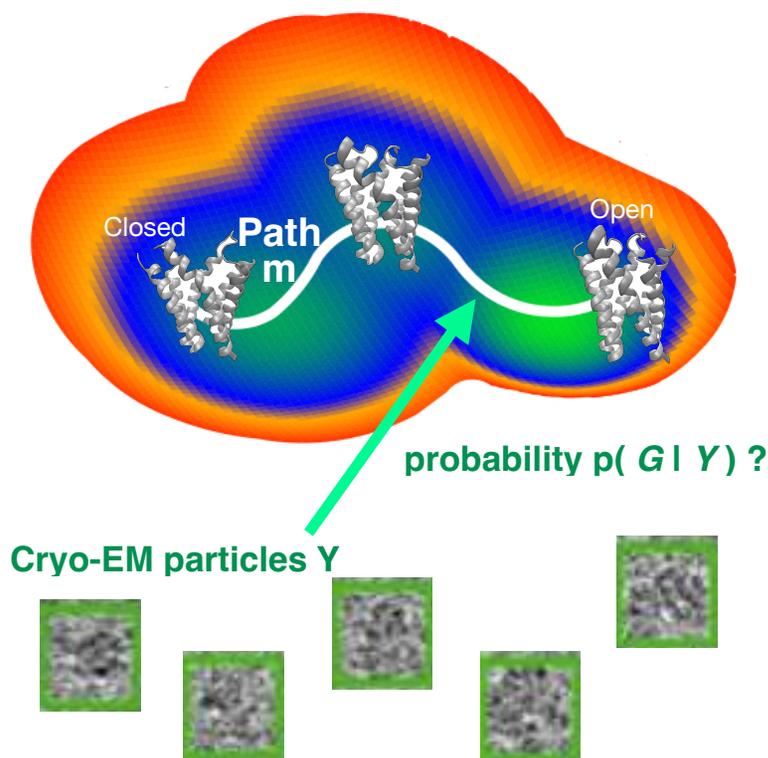
Article | [Open Access](#) | [Published: 01 July 2021](#)

### A Bayesian approach to extracting free-energy profiles from cryo-electron microscopy experiments

[Julian Giraldo-Barreto](#), [Sebastian Ortiz](#), [Erik H. Thiede](#), [Karen Palacio-Rodriguez](#), [Bob Carpenter](#), [Alex H. Barnett](#) & [Pilar Cossio](#) ✉

We assume a low-temperature approximation and project along a path collective variable\*

$$p(x|G) = \sum_m \delta(x - x_m) \frac{e^{-\beta G(s_m)}}{Z_1}$$



\*Inspired by MD methods like the string method (Weinan, Ren, Vandenn-Eijnden - Physical Review B, 2002)

# Cryo-BIFE

For multiple images  $Y$ , the posterior reduces to

$$p(G|Y) \propto p(G) \prod_i \left[ \sum_m p(y_i|x_m) \frac{e^{-\beta G(s_m)}}{Z_1} \right]$$

*Prior*
*Single image*

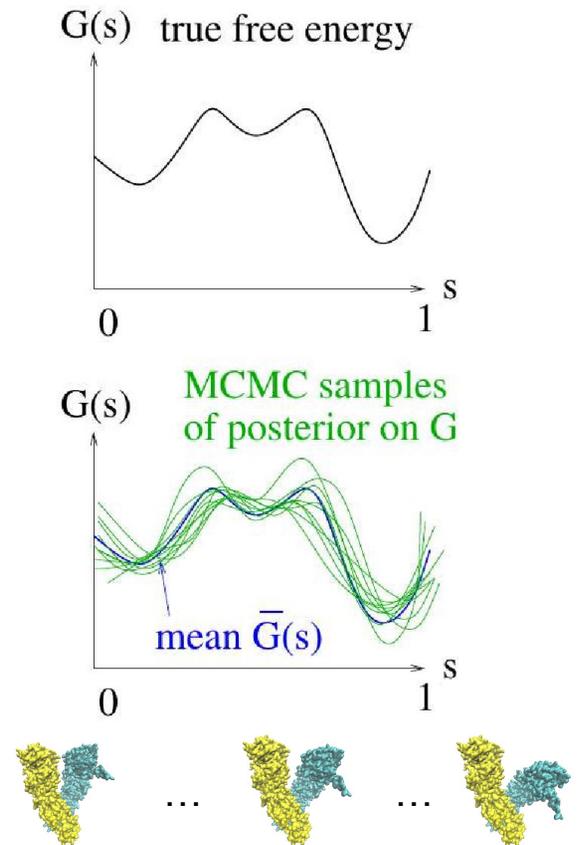
Probability of a  $G$  given the particles

## Main point:

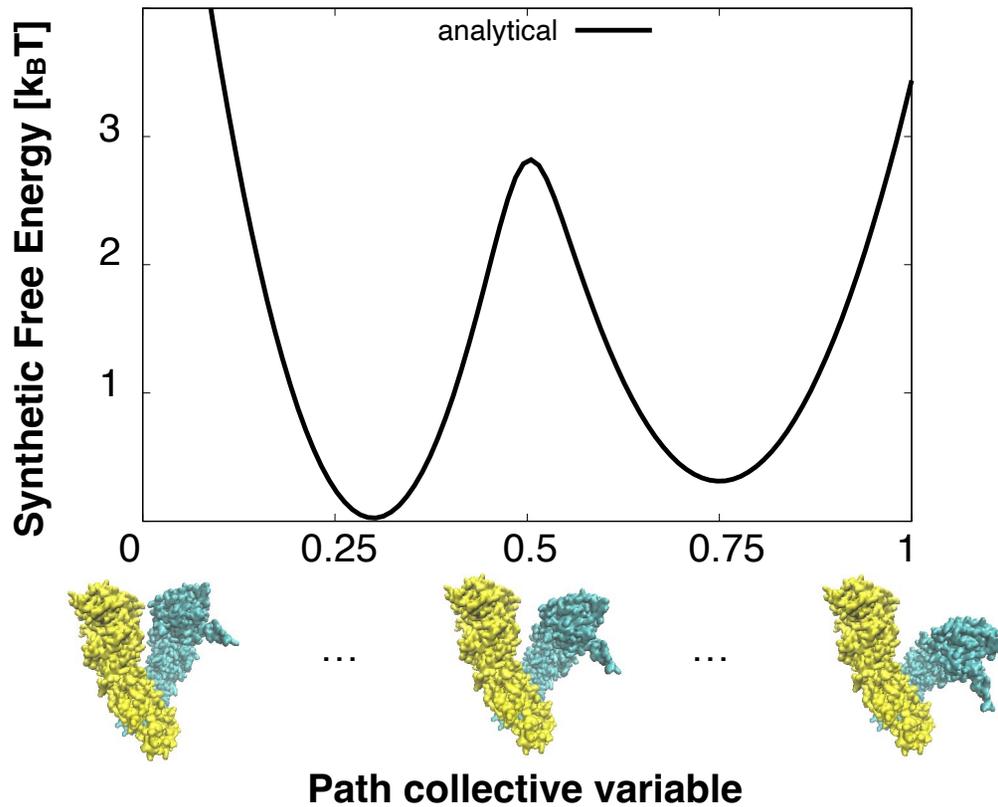
Sampling  $G$  from  $p(G|Y)$ , allows us to calculate its mean and **uncertainty** at each node given the cryo-EM particles.

## Key ideas:

- Particle reconstructions are **not** crucial.
- We assume that we have a sufficiently good transition path of 3D structures (by MD, AlphaFold etc).
- We use MCMC sampling (e.g. in STAN) to extract the expected  $G$  and its uncertainty.



# 1D Hsp90 synthetic data



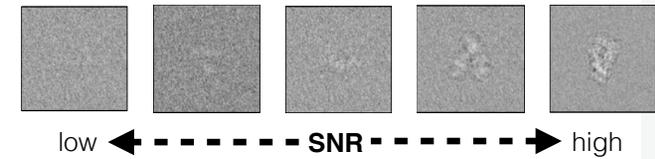
## Path calculation:

- 20 nodes (along the single degree of freedom)
- 2 orientation-rounds of BioEM

## Synthetic images:

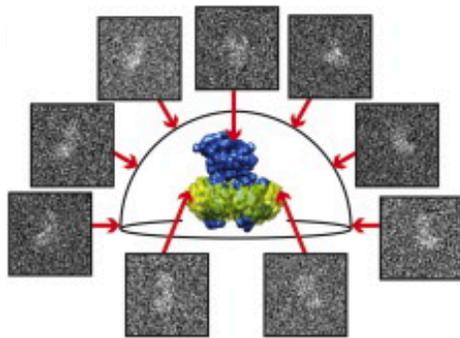
~13300 particles

- 128x128 pixels
  - pixel size 2.2 Å
  - Coarse-grained residues
- With uniformly distributed:
- random orientations
  - random defocus [0.5,3] micro-m
  - random signal-to-noise ratio (SNR) [0.001, 0.1]

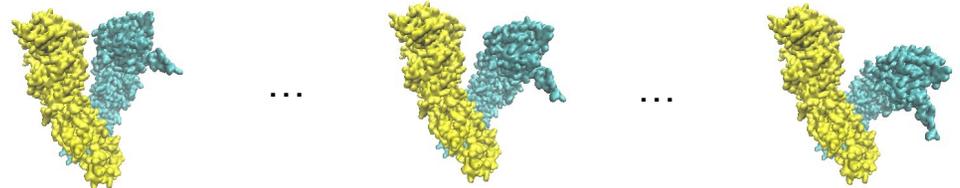
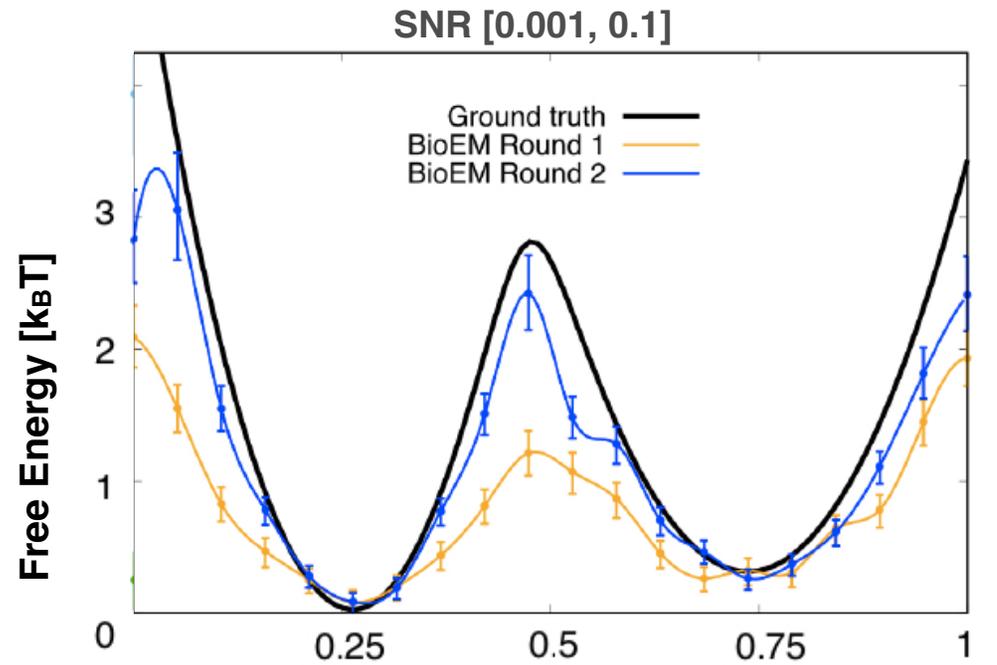


# Projection direction accuracy is important

The projection direction is unknown.



Nogales & Scheres, Mol Cell 2015



Path collective variable

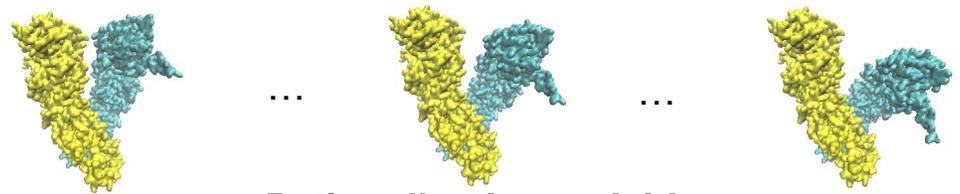
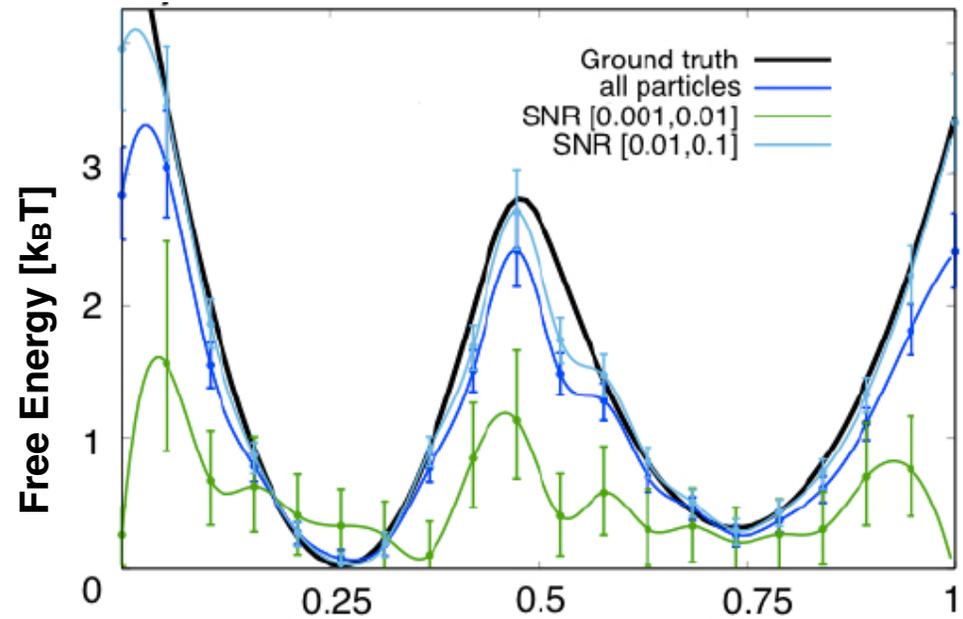
# Signal-to-noise ratio (SNR)

A poor FE profile is recovered from the low SNR group

High SNR provides a slightly better profile than using all particles.

Adding “bad” particles does not hinder the FE recovery (if there are some “good” ones in the group).

Orient. Round 2



Path collective variable

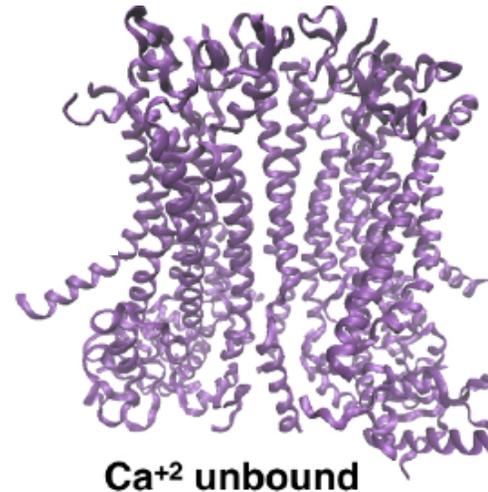
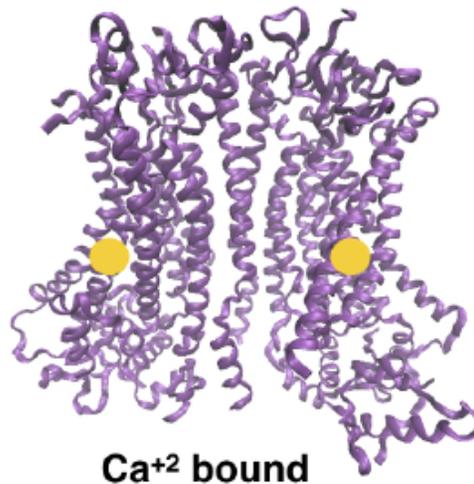
# Real data: TMEM16F a calcium-activated ion channel

## Cryo-EM Studies of TMEM16F Calcium-Activated Ion Channel Suggest Features Important for Lipid Scrambling

(2019)

Shengjie Feng,<sup>1,4</sup> Shangyu Dang,<sup>2,4</sup> Tina Wei Han,<sup>1</sup> Wenlei Ye,<sup>1</sup> Peng Jin,<sup>1</sup> Tong Cheng,<sup>1</sup> Junrui Li,<sup>2</sup> Yuh Nung Jan,<sup>1,2,3</sup>  
Lily Yeh Jan,<sup>1,2,3,\*</sup> and Yifan Cheng<sup>2,3,5,\*</sup>  
<sup>1</sup>Department of Physiology, University of California, San Francisco, San Francisco, CA 94158, USA  
<sup>2</sup>Department of Biochemistry and Biophysics, University of California, San Francisco, San Francisco, CA 94158, USA  
<sup>3</sup>Howard Hughes Medical Institute, University of California, San Francisco, San Francisco, CA 94158, USA

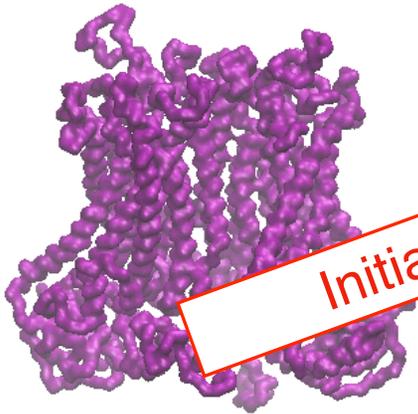
Generated from two different data sets: with and without  $\text{Ca}^{+2}$



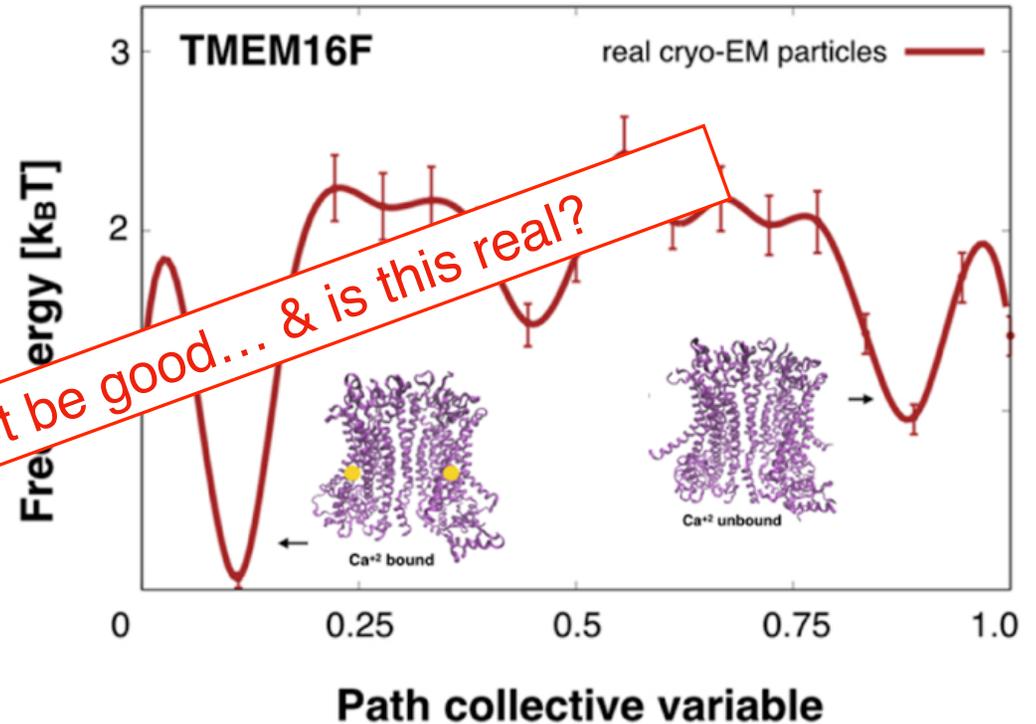
All particles  
available:  
EMPIAR 10278

Is there a population of the  $\text{Ca}^{+2}$ -unbound state in the  $\text{Ca}^{+2}$ -bound set?

# Real data: TMEM16F a calcium-activated ion channel



Initial path might not be good... & is this real?



Particles from the entire set:

- EMPIAR 10278
- 15000 images
- 256x256 pixels
- pixel size 1 Å

# Opening Pandora's box

- ***Extension to many dimensions***
- Experimental validation system
- How to optimize the molecular path?
- Comparing many models vs many images optimization.



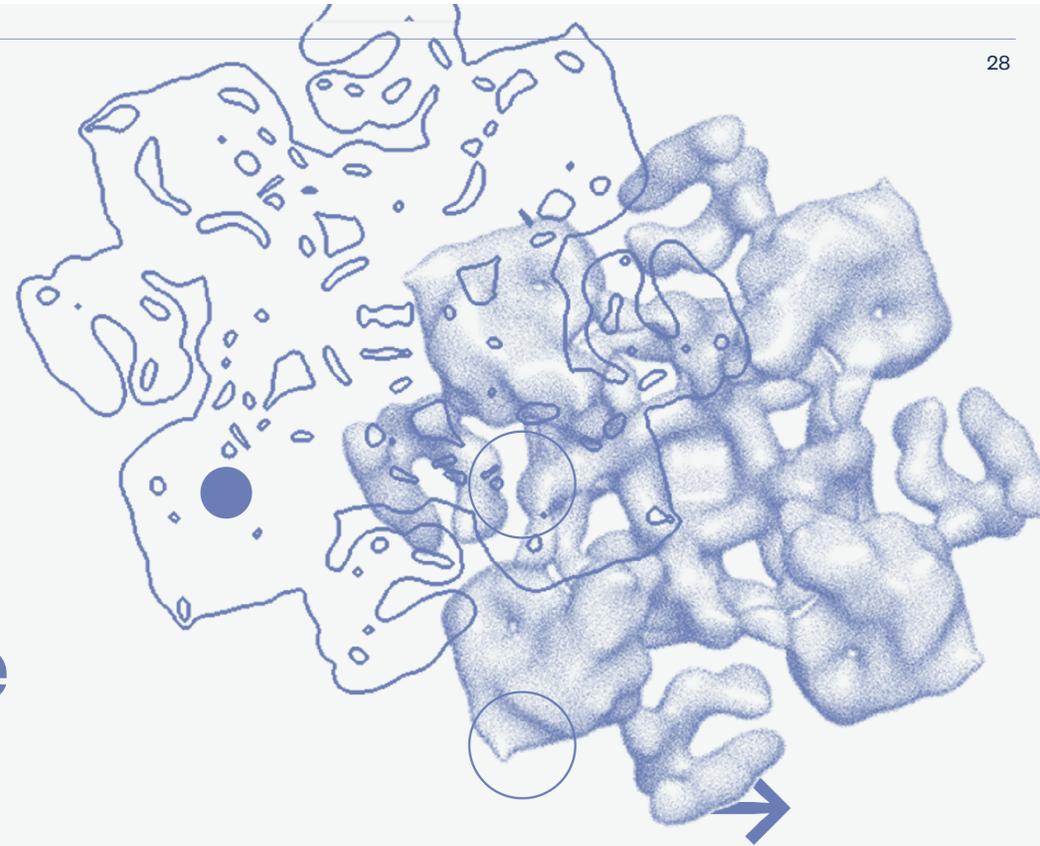


WS Tang



EH Thiede

# Cryo-EM ensemble reweighting\*



*Extension to multiple dimensions: working directly in configuration space.*



\*Tang et al. 2022 <https://arxiv.org/abs/2212.05320>



# Theory

## Bayes' theorem

$$p(H|Y) \propto p(H)p(Y|H)$$

we would like to know

For a set of *iid* cryo-EM particles  
(not an averaged observable)

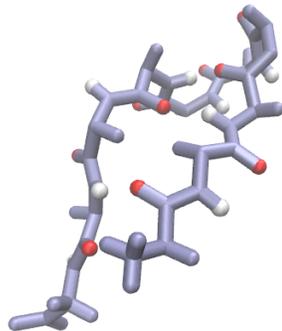
$$p(H|Y) \propto p(H) \prod_i \int p(y_i|x_i)p(x_i|H)dx_i$$

data



CryoEM  
images  
 $Y=\{y_i\}$

Conformation  
 $x \in \mathbb{R}^{3N}$



Boltzmann distribution

$$p(x|H) = \frac{1}{Z} e^{-\beta H(x)}$$

Hamiltonian

## Approximation:

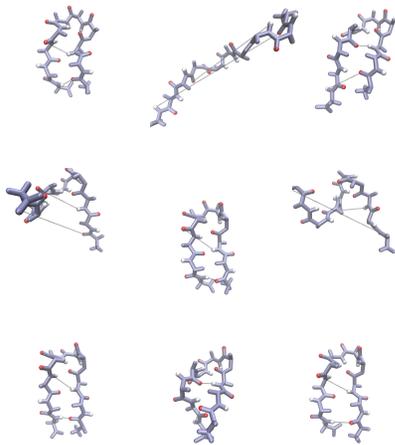
we parametrize the density by a set of  $\{x_m\}$  with weights  $\{\alpha_m\}$ .

data



CryoEM  
images  
 $Y = \{y_i\}$

Set of  $\{x_m\} \in \mathbb{R}^{3N}$



Approximate probability density

$$p(x|\{\alpha_m\}) = \sum_m \delta(x - x_m) \alpha_m$$

↑  
Goal: extract the  
weights

↑  
Sum to one

The posterior reduces to

$$p(\{\alpha_m\}|Y) \propto p(\{\alpha_m\}) \prod_i \sum_m p(y_i|x_m) \alpha_m$$

\*Same expression as in cryoBIFE (Sci Rep 2021), so we can use STAN!

# Test system Chignolin: Synthetic data

Sufficiently small to have a converged ensemble with MD.

But complex enough with three metastable states:

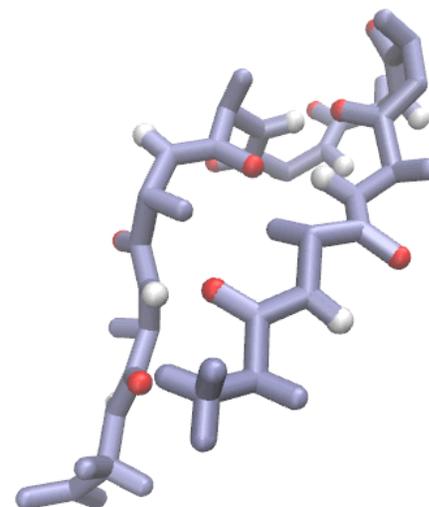
Folded



Unfolded

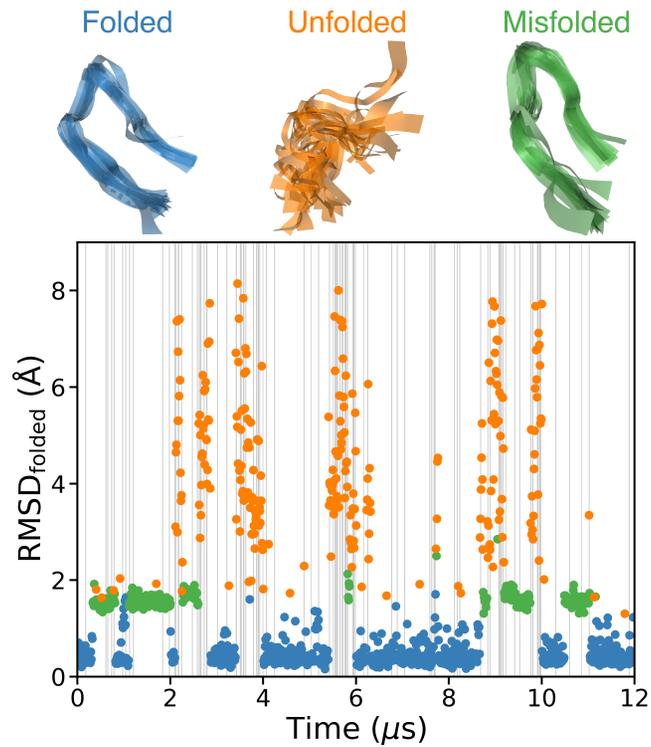


Misfolded



## Two independent ensembles:

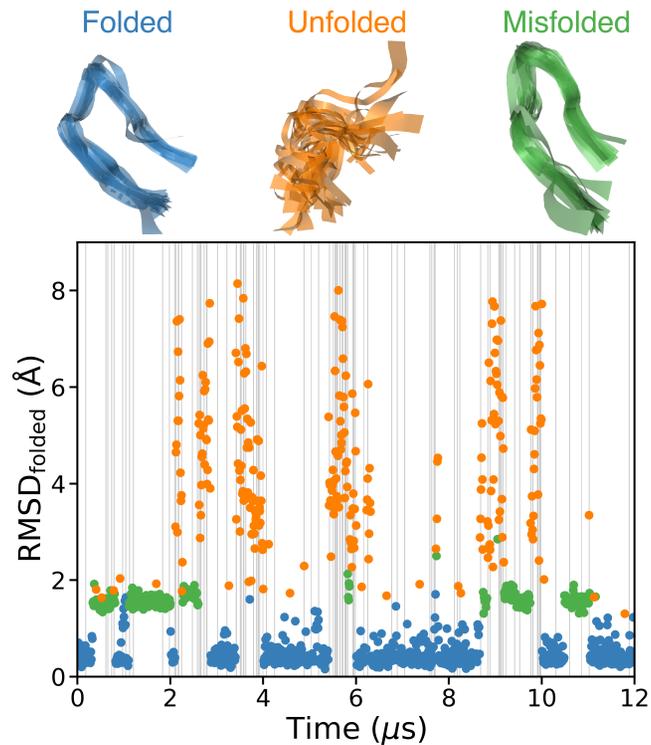
Structure-generating “in house”  
MD trajectory



Cluster conformations with K-medoids  
RMSD (50 centers) =  $\{x_m\}$

## Two independent ensembles:

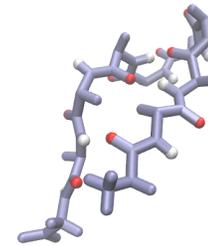
Structure-generating “in house”  
MD trajectory



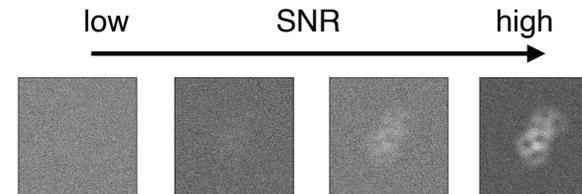
Cluster conformations with K-medoids  
RMSD (50 centers) =  $\{x_m\}$

Image-generating DE Shaw MD  
trajectory

- > Different force field
- > Different ensemble
- > Structures to generate images are not in the ensemble for refinement.



Example images with different SNR (random pose & defocus)



For simplicity, we  
assume that the  
pose and  
defocus known.

## Can we recover the state populations?

The population is defined as the sum of the weights of the members (i.e., cluster centers) of each state.

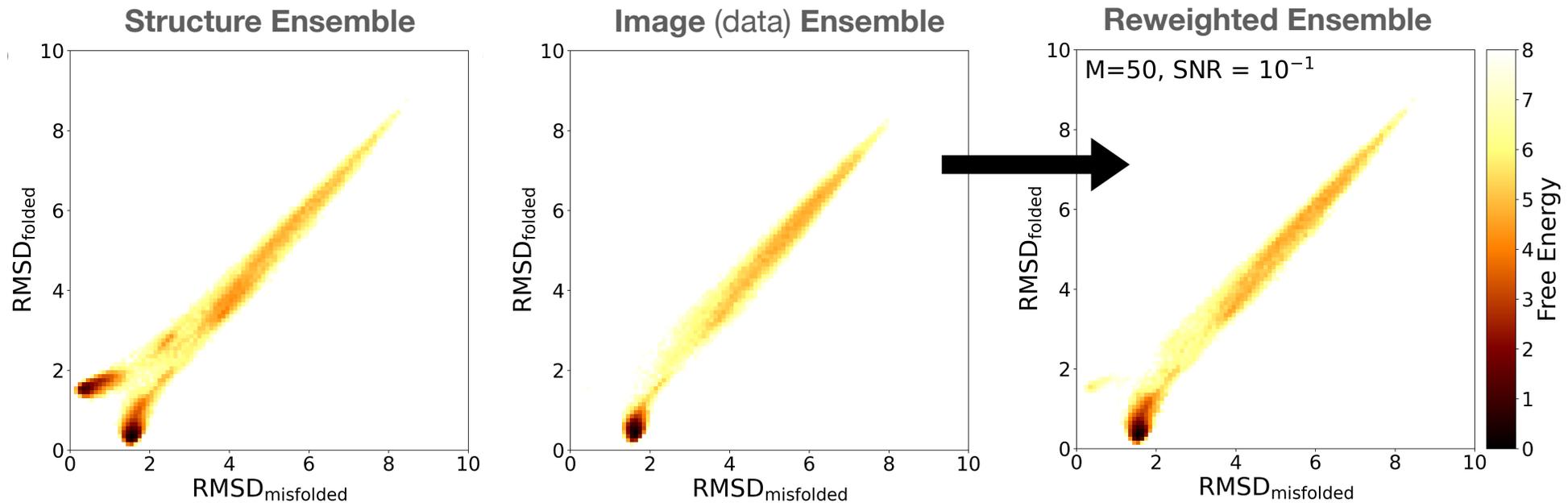
Population recovery\*:

	<b>%fld</b>	<b>%msfd</b>	<b>%unfd</b>
<b>Ground Truth</b>	<b>0.7707</b>	<b>0.0004</b>	<b>0.2289</b>
<b>No noise</b>	0.7784	0.0058	0.2152
<b>SNR = 1.0</b>	0.7786	0.0055	0.2160
<b>SNR = 0.1</b>	0.7787	0.0058	0.2155
<b>SNR = 0.01</b>	0.7693	0.0129	0.2178

Yes, even at low SNRs\*

\*If we have to optimize the pose then it becomes more challenging

● ○ *Post-processing of the refined ensemble:  
Free-energy over CVs of choice\*.*

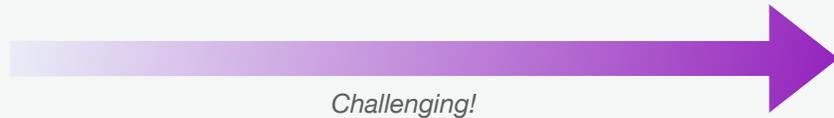
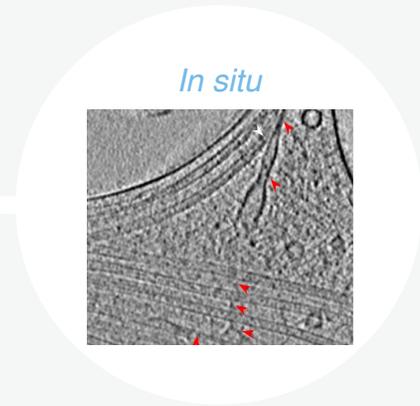
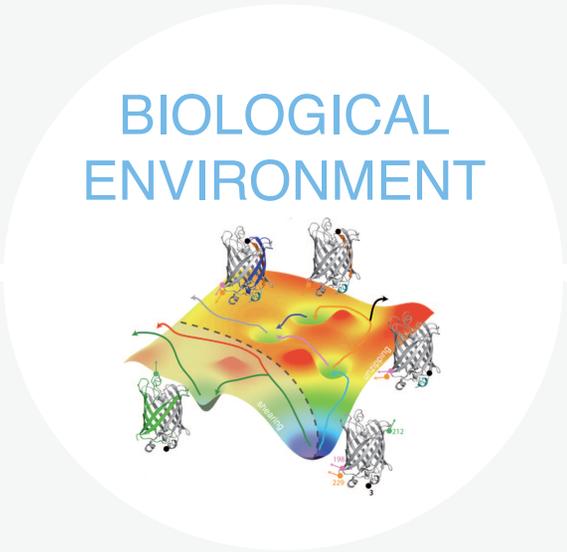
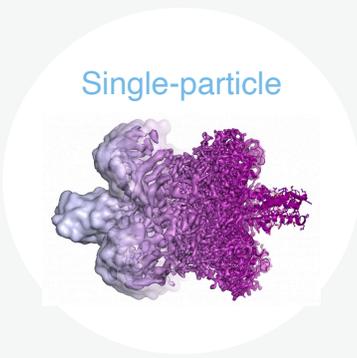


## Open questions:

- Validation system for conformational variability?
- How to generate a 'good' ensemble → bias MD



A stepping stone toward conformational landscapes *in situ*.



Challenging!

# Acknowledgments



*Structural and Molecular Biophysics, CCB/CCM*



## CCM



Leslie Greengard



Marylou Gabrie



Erik Thiede



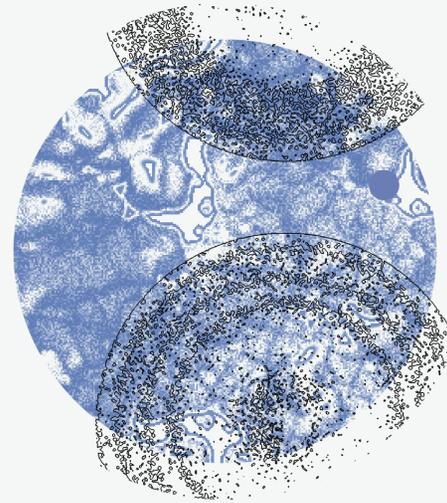
Alex Barnett



Bob Carpenter

Collaborators:  
*Steve Bonilla, Roberto Covino, Attila Szabo, Fabio Pietrucci, Karen Palacio-Rodriguez, Hadrien Vroylandt,....*

**THANKS FOR  
YOUR ATTENTION**



 FLATIRON  
INSTITUTE