(S9716)

Visualizing ATP-Dependent Substrate-Processing Dynamics of the Human 26S Proteasome at Near-Atomic Resolution

Youdong Mao Peking University School of Physics Dana-Farber Cancer Institute, Harvard Medical School

GTC 2019, San Jose

History of Electron Microscopy

The Nobel Prize in Chemistry 2017



Timeline of Cryo-EM for Complex Dynamics



Cryo-EM for Complex Dynamics of Human Proteasome



Single-Particle Reconstruction of Cryo-EM



Ubiquitin-Proteasome Pathway



For the discovery of ubiquitin-mediated protein degradation

The Nobel Prize in Chemistry 2004



Photo: D. Porges

Aaron Ciechanover Prize share: 1/3



Photo: D. Porges

Avram Hershko

Prize share: 1/3



Photo from the Nobel Foundation archive.

Irwin Rose Prize share: 1/3

The largest molecular machine for degradation

28 subunits in 20S core particle (CP)
19 subunits in 19S regulatory particle (RP)



An assembly of multiple molecular sub-machines



- Ubiquitin receptor: capture ubiquitintagged substrates
- Deubiquitinase: remove ubiquitin tags from the substrate
- AAA-ATPase Unfoldase: unfold the substrate
- Degradation
 chamber: degrade
 the substrate

Substrate-engaged proteasome



Y. Dong, S. Zhang, ..., Y. Mao. *Nature* 565: 49-55 (2018).

Substrate-engaged human proteasome



Y. Dong, S. Zhang, ..., Y. Mao. Nature 565: 49-55 (2018).

Local resolution



High-resolution density maps



Y. Dong, S. Zhang, ..., Y. Mao. *Nature* doi: 10.1038/s41586-018-0736-4 (2018).

Observation of single magnesium ions



Y. Dong, S. Zhang, ..., Y. Mao. Nature doi: 10.1038/s41586-018-0736-4 (2018).

Substrate-processing dynamics



ATP hydrolysis powers the molecular motor



Dynamics of the motor unfolding substrate

□ Coordination of three adjacent ATPases



Three modes of ATP hydrolysis



Workflow of Single-Particle Image Analysis



Hierarchical focused 3D classification



Comparison of approaches in analyzing protein dynamics

	NMR	Single-molecule FRET	Small-angle X- ray scattering	Cryo-EM			
Time resolution/scal e	fs	ms	ms	ms to seconds			
Real-time	Yes	Yes	Yes	Pseudo			
Spatial resolution	angstrom	>1 angstrom	nanometer	> 2.5 angstrom			
Full atomic modelling	Yes	No	No	Yes			
Protein size	<300kDa	Any	Any	>100 kDa			
Labeling Yes		Yes	No	No			
Energy landscaping	Complete	Partial	Low-resolution and partial	Complete			

Summary of GPU acceleration for cryo-EM

	Software	GPU Acceleration	Based on Deep Learning					
Motion Correction	MotionCor2	Yes	ms					
CTF parameters	GCtf	Yes	Yes					
Particle picking	DeepEM	Yes	nanometer					
Deep 2D classification	ROME	Νο	No					
3D classification	RELION	Yes	Any					
Deep 3D classification	ROME	Νο	No					
Cryo-EM refinement	RELION CryoSPARC	Yes	Low-resolution and partial					

Convolutional Neural Network



The 6C-3S-12C-2S-12C-2S CNN we designed for KLH dataset.

Y. Zhu, Q. Ouyang, Y. Mao. BMC Bioinformatics 18, 348 (2017).

The learned feature maps of an example input image



Layer C1 Feature Maps





Layer S2 Feature Maps





Layer C3 Feature Maps

0400

Layer C5 Feature Maps



Layer S4 Feature Maps



Layer S6 Feature Maps



Input Image

DeepEM Algorithm Workflow



Keyhole Limpet Hemocyanim (KLH)



Zhu, Y., et al. IEEE Trans. Med. Imaging 22, 1053-1062 (2003)

Testing on some challenging datasets



Effect of training dataset sizes



DeepEM exhibits good performance on low SNR



DeepEM on GPU

Geforce GTX 970, running Matlab 2016a and CUDA
 8.0: 40-190 seconds per micrograph

□ Tesla K20c GPU with CUDA 8.0:

	Pdb1f07	KLH	10,005	10,184	10,075
Micrograph Size	1024*1024	2048*2048	3710*3710	3838*3710	4096*4096
Particle Size	100*100	272*272	180*180	256*256	300*300
Preprocessing	0.17	0.57	3.08	3.72	3.75
Segmentation	0.55	3.64	9.92	9.58	11.23
Classification	2.26	3.12	8.63	6.84	2.75
Postprocessing	0.34	4.59	20.03	6.99	30.38
Total Time of PIXER	3.32	11.94	41.67	47.17	48.07
DeepPicker	10.47	23.75	80.76	81.34	95.43
DeepEM	40.56	80.54	65.47	39.75	54.38

Y. Zhu, Q. Ouyang, Y. Mao. *BMC Bioinformatics* 18, 348 (2017). J. Zhang, ..., F. Sun, F. Zhang. *BMC Bioinformatics* 20, 41 (2019).

Summary

- DeepEM outperforms existing programs in particle recognition, allowing particle picking, selection and verification in an integrated fashion.
- With better training dataset or iterative training, one can improve the accuracy of particle recognition. In other words, the program can be trained to be "smarter".
- □ DeepEM can be easily applied to batch processing.
- We expect that DeepEM marks the inception of applications of modern AI technology in expediting cryo-EM structure determination.

Latent variable space Data space A(s; W)mapping S22 $p(\mathbf{s}) = \frac{1}{K} \sum_{k=1}^{K} \delta(\mathbf{s} - \mathbf{s}_{k}) \qquad T_{\mathbf{r}}(X_{ij}) = \operatorname{CTF}_{ij} \left[A(\mathbf{s}; \mathbf{W}) \right]_{j} + N_{ij}$ $\mathbf{t}_{1} \quad p(\mathbf{Y}_{i} | \mathbf{s}, \mathbf{W}, \boldsymbol{\beta}) = \prod_{i=1}^{J} \left(\frac{\beta_{ij}}{2} \right)^{\frac{1}{2}} \exp\{-\frac{\beta_{ij}}{2} \left(Y_{ij} - \operatorname{CTF}_{ij} \left[\mathbf{A}(\mathbf{s}; \mathbf{W}) \right]_{j} \right)^{2} \}$ S. $\left[A(s;\mathbf{W})\right]_{j} = \sum_{m=1}^{M} \varphi_{m}(s) W_{mj} \qquad \varphi_{m}(s) = \begin{cases} \exp\{-\frac{\|s - \mu_{m}\|^{2}}{2\sigma^{2}}\}, m \le M_{NL} \\ 1, m = M_{NI} + 1 \end{cases}$

Statistical Manifold Learning: Generative Tomographic Mapping

- □ The shortest line between two points s_i and s_j in the latent space is mapped to geodesic line on a manifold between A(s_i) and A(s_j) in the data space.
- Latent space: the structural difference resulting from changes in the image orientation and/or the molecular conformations
- Data space: the Fourier transform of particle images

J. Wu, Y. Ma, C. Congdon, B. Brett, S. Chen, Q. Ouyang, Y. Mao. PLoS ONE 12, e0182130 (2017).

SML optimizes a MAP objective function \Box The maximum-likelihood estimator $\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(\mathbf{Y}_i | \Theta)$

□ The maximum-a-posterior (MAP) estimate

$$\hat{\boldsymbol{\Theta}} = \operatorname{argmax}\left[\sum_{i=1}^{N} \log p(\boldsymbol{Y}_{i} | \boldsymbol{\Theta}) + \log p(\boldsymbol{\Theta})\right]$$
$$p(\mathbf{W}) = \left(\frac{\alpha}{2\pi}\right)^{\frac{MJ}{2}} \exp\left(-\frac{\alpha}{2}\sum_{m=1}^{M}\sum_{j=1}^{J}W_{mj}^{2}\right)$$

- Maximization of the objective allows a numeric solutions to the model parameters that define the mapping of latent variables to the manifold in the data space.
- We developed Expectation-Maximization algorithm with CTF embedded for numerical solution.

Contrast Transfer Function (CTF) and Aberration

- □ CTF is the Fourier transform of the point spread function of transmission electron microscope
- □ CTF characterize how the images are distorted by the aberration of objective lens





Combined Effect of CTF and Noise

Noiseless and no CTF

Noiseless with CTF



Heavy noise with CTF

Expectation-maximization algorithm

□ In the E-step, we evaluate the posterior probability

$$R_{ki}(\boldsymbol{\Theta}^{[n]}) = p(\boldsymbol{s}_{k} | \boldsymbol{Y}_{i}, \boldsymbol{W}^{[n]}, \boldsymbol{\beta}^{[n]}) = \frac{p(\boldsymbol{Y}_{i} | \boldsymbol{s}_{k}, \boldsymbol{W}^{[n]}, \boldsymbol{\beta}^{[n]}) p(\boldsymbol{s}_{k})}{\sum_{k'=1}^{K} p(\boldsymbol{Y}_{i} | \boldsymbol{s}_{k'}, \boldsymbol{W}^{[n]}, \boldsymbol{\beta}^{[n]}) p(\boldsymbol{s}_{k'})}$$

□ In the M-step, we maximize the MAP estimator

$$Q(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{[n]}) = \sum_{i=1}^{N} \sum_{k=1}^{K} R_{ki}(\boldsymbol{\Theta}^{[n]}) \ln p(\boldsymbol{Y}_{i} | \boldsymbol{\Theta}) + \ln p(\boldsymbol{\Theta})$$

$$\sum_{i=1}^{N} \sum_{k=1}^{K} R_{ki}(\Theta^{[n]}) \beta_{ij}^{[n]} \varphi_{m}(s_{k}) CTF_{ij}(Y_{ij} - CTF_{ij} \sum_{m'=1}^{M} \varphi_{m'}(s_{k}) W_{m'j}^{[n+1]}) - \alpha W_{mj}^{[n+1]} = 0$$
$$\frac{1}{\beta_{ij}^{[n+1]}} = \sum_{k=1}^{K} R_{ki}(\Theta^{[n]}) \left(Y_{ij} - CTF_{ij} \sum_{m=1}^{M} \varphi_{m}(s_{k}) W_{mj}^{[n+1]}\right)^{2}$$

Deep classification by SML



- There is a probability calculated for each image assigned to a given class
- □ The class average is CTF corrected and is a probability-weighted average.

SML vs. MAP: ~40% higher angular accuracy



Effect of SNR on angular error







SML vs. MAP classification SML MAP (maximum a posteriori)



17,103 inflammasome particles are classified into 300 reference-free classes. Only classes whose particle numbers were larger than 9 are exhibited.

Deeper classification identifies hidden heterogeneity



(free RP) (free RP) (RP - CP) (RP - CP) (RP - CP)





Secondary MAP classification

Secondary SML classification

1000 reference-free classes computed in a few hours

*	2	-	\$	Ð	22	3	\$#	2	4	雜	4	-	10	all all	13	*	18	柳	感	-		and the second	200	ALC: N	A.	14			Contraction of the second	10.00						
	*	3	-	1	-	(5)	3	3	3	金	(0)	3	弦	63	S.	05	-	1	É.	1	CH-	S.D.	11	Sec.	100				No.							
·	3	40	÷	物	de la	3		24	24	1	5	*	5	3	100									and and		100	ie.		100							
1	赤	書		-	3	2	23	9	THE .	1	-	4	3	-	各	12	1. Sta	100		2	St.	ALL NO	111	100												
-	-	泰	2	*	4	2		14		197	2	(W)		All of the	31		tic	報	1	20	Ser.	191		100	10			and the second								
2	南	*	-	4	Se la		4	告	Se la	1	2	3	新		35	朝	3	dia dia	朝	24	100	Sale		13	ALC: NO				Cler.							
2	to	8	-	赤	3	132	-	恋	影	3	-	影	3	W.	-	(Er	100	2	-22	and the	all a		15		200				1611							
dit .	剱		-	3	No.	3	Ş.	*	B	3	No.	di.	帮	3	-	11	S.	(4)	金	2				1210			100	12.25								
de la	南	1	$\hat{\varphi}$	8	18	-	5	あ	1	3	37	the second	20	-	3	-33	-	12	A.R.	1	in the	3	ALC: N	1	N. N.			N. L.	No.							
秘	2	Ť.	部	-	-	12	3	-	0	-	-	12	14	Res and	27	1	No.	5	-45	The second			- 10					No.	Re Fa							
勢	题	-	2	-	3	萄	3	-	5	11	10		161	物	1	御	and the	125	419	22	in the	and the	1	2			Car	1.1		N.						
ē?	8	2	-	物	物	1	4	-27	制	133	12	12	-	·Re	3	德		30	14	No.			1	NAME OF TAXABLE		100	and the second	35	New York							
-	(3)	1	1	5	1	3	1	3	影	3	1	10	命	10	51	-	金	3			12	100	E.B.	124		and the	22	1								
\$	榆	d'	6	à	30	(1)	ŝ.	1	-	-	-	1	0	alle a	10-	始	12	13	111	(0)	歌	1	田		and the second	10				(and	10		A STATE	Sec. 1		
$\bar{\Delta}$	٢	-	de la	-	*	1	3	-	-	3	-	102	3	32	the.	帝	192	-	S.S.	S.M.	54	-	50	-	12		Res I		3		16 B			No.		
¢	Ð	赤	4	4	ŝ	魏	3	勢	-	10	3	14	1	100	15	4	3	家	39	3	6.	1	SY.	1	1402	K	No.	1			ALC: NO			The second		
4	楷	*	*	÷	2	1	1	翁		10	-	3	-	-	\$	10	1	10	in the	-	创	and the	133	100	Ser.		and the second s	S.	and the second			100		States		
合	\$	-	8	÷		13	3	儆	22	3	(6)	8	2	1	See.	S.	6	13		郡	The second	33	3	EN		-	26							and		
\$		-	$\widehat{\mathcal{A}}_{k}^{(i)}$	Å	3	2	Ŵ	3	金	墩	\$3)	R	物	5	感	in.	-		命	125	商	和	The second	1		9	SE I	A CON								
4	-	*	*	\$	物	ŝ	奔	-	1	197		-	3	-	-	ち	a.	-	SAN		14	122	14	5	-	4	AN AN							No.		
索	\tilde{k}	4	2	*	幸	ŝ	3	34	4	Q.	-	5	1	-	10	3	海	4	1	60	131	3	He was	6	13	14 A	124		E SAN							
\$	赤	10	4	\$	34	靈	¢,	-	-5	39	剑	語	2	3	2	22	100	100	12	122		14	SF.	-	13	11C	352	15	No.		121		No.			
\$	纺	物	3	٩	Q.	100	30	de la	嘲	12	1	3	語	D	3	-	影	and the second	部	5	AN AN					all a		-								
÷	\$	*	\$	\$	3	S.	ð	-	34	恭	1	-	の	3	S)		ST.	3	-	No.	-20	1		100	12		100							No.		
	\$	-	-	3	1	1	ĝ?	*	會	部	2	-	素	80	44	御	論	010	3	50	Ser.	-	No.	201	action 1	11an	10		100				(inter	5153		
4	Ş.	4	\$	\$	14	3	24	1	20	-	3	-	3	30	北	die .	-	5	1	-	-TE	10	24	Set.	A.S.	-	199		(internet	Contraction of the second	1	ALC: NO	151	Ser.		

Deep classification improves initial reconstruction



ROME-based initial model: 234 best class averages selected from 1000 reference-free classes of a cryo-EM dataset (117,471 images) within 3 hours on a cluster of 512 CPU cores (32 nodes)

- EMAN2-based initial model: all 128 ref-free class averages calculated by EMAN2 for >8 hours using the same dataset.
- The resolution of the ROME-based initial model is 10-Å (50%) higher, show in the FSC plots on the left.

Memory/cache optimization-code blocking



Access pattern

Access pattern after blocking

Memory/cache optimization-code blocking for adaptive

search



Access pattern

ROME vs. RELION: ~20-fold faster on the same hardware

 When GTM and ML algorithms in ROME are used for unsupervised deep classification, it outperforms RELION by 10-20 times



Performance of deep classification by ROME

- ROME can generate 1000 reference-free classes of real cryo-EM data within a few hours on a cluster of 500 CPU cores (32 nodes)
- □ All other programs tested have crashed on the same dataset in our test clusters when doing this

task



Speedup on Intel Xeon Phi Knights Landing processors



SOURCE: https://software.intel.com/en-us/articles/recipe-rome-1-0-sml-for-the-intel-xeon-phi-processor-7250

Features in our ROME v1.0 system

- □ ROME (<u>Refinement and Optimization with Machine-IEarning</u>) implements SML in HPC.
- Fully compatible with the I/O file format used in other software in structural biology and cryo-EM imaging, including RELION and SPIDER.
- Fully modernized code, designed and parallelized for Intel Xeon CPUs
- □ Designed to use both OpenMP and MPI.
- It can perform image classification at a scale of thousands of classes in a single run with improved accuracy, which is about 1-2 orders of magnitude greater than existing software in this area.

Performance of deep classification by ROME

- ROME can generate 1000 reference-free classes of real cryo-EM data within a few hours on a cluster of 500 CPU cores (32 nodes)
- All other programs tested have crashed on the same dataset in our test clusters when doing this task



Speedup on Intel Xeon Phi Knights Landing (KNL) processors



SOURCE: https://software.intel.com/en-us/articles/recipe-rome-1-0-sml-for-the-intel-xeon-phi-processor-7250

MAP3D benchmark dataset 1 : ribosomes wi/wo EFG

□ 5 nodes Intel(R) Xeon(R) CP E5-2697 v4 @ 2.30GHz □ RELION 1.4

ROME:

RELION

MPI_Rank	Threads	Running_time(minutes)
10	18	222.95
10	36	166.03
20	9	187.03
20	18	139.37
45	4	160.9
45	8	125.68
90	2	151.08
90	4	120.0

ROME 1.1.0

MPI_Rank	Threads	Running_time(minutes)
10	36	26.1
5	72	21.98

MAP3D benchmark dataset 2 : Plasmodium ribosome



RELION 1.4

MPI_Rank	Threads	Running_time(minutes)
10	18	857.6
10	36	644.02
20	9	769.15
20	18	624.23
45	4	835.37
45	8	658.73
90	2	crash
90	4	crash

ROME 1.1.0

MPI_Rank	Threads	Running_time(minutes)
10	36	180.55
5	72	166.2

MAP3D for KNL benchmark (3D classification/refinement)



platform	CPU	Recommended Customer Price
KNL	Intel® Xeon Phi™ Processor 7230 16GB, 1.30 GHz, 64 core	\$1,992.00
Broadwell	Intel [®] Xeon [®] Processor E5-2697 v4 45M Cache, 2.30 GHz	\$2,702.00

Features in ROME v1.1.0 system

- □ ROME (<u>Refinement and Optimization with Machine-IEarning</u>) implements SML and MAP3D in HPC.
- Fully compatible with the I/O file format used in other software in structural biology and cryo-EM imaging, including RELION and SPIDER.
- Fully modernized code, designed and parallelized for Intel Xeon CPUs
- □ Designed to use both OpenMP and MPI.
- It can perform image classification at a scale of thousands of classes in a single run with improved accuracy, which is about 1-2 orders of magnitude greater than existing software in this area.

Summary

- □ SML enables unsupervised deep classification at a greater scale.
- ROME, the first HPC software implementing SML algorithm, is optimized based on the modern standard of parallelization.
- Despite SML classification is mathematically more complicated than MAP classification, it is more suited for parallelization in the latest CPU architecture with a higher throughput of vector processing.

ROME Official Website http://ipccsb.dfci.harvard.edu/rome

Join Us

for STRUCTURAL BIOLOGY

Research V People V Resources V About Intel® PCCSB V

Home I Resources I ROME System

Resources

ROME System

Sullivan Supercomputer

Tecnai Arctica Imaging Platform

Automated CryoEM Imaging System

ROME System

The ROME (Refinement and Optimization via Machine IEarning for cryo-EM) project is one of the major research efforts at the Intel® PCCSB. The project aims to develop a parallel computing software system for high-resolution cryo-EM structure determination and data analysis, which implements advanced machine learning approaches in modern computer sciences and runs natively in an HPC environment. The ROME project emphasizes the development and improvement of cutting-edge machine-learning algorithms that can tap into increasingly complicated cryo-EM data with low SNRs, high structural



heterogeneity, and increasing volume of big data. The ROME system is expected to be optimized on both Intel® Xeon multi-core CPUs and Intel® Xeon Phi many-core coprocessors. The areas of investigation in this critical program include but are not limited to statistical pattern recognition, nonlinear dimensionality reduction and big-data processing.

Enter the official ROME software release website

ABOUT INTEL® PCCSB

SITE NAVIGATION

WHAT'S UP

The Intel® Parallel Computing Center for Structural Biology (Intel® PCCSB) is an innovation-oriented research center hosted at Dana-Farber Cancer Institute and Harvard Medical School. ...**Read more** Core Research Mission | The ROME Project Publications | Software | Members News and Views | Press Release Training Program | Workshop | Symposium

Copyright © 2014 | Intel® Parallel Computing Center for Structural Biolog

Solution to sample heterogeneity



L. Zhang et al. Science 350: 404-409 (2015).

Acknowledgement



Yuanchen Dong



Shuwen Zhang



Zhaolong Wu

- Early work on substrate-free proteasome was done with collaboration with Prof. Marc Kirschner's group at Harvard.
- Recent data were collected at Electron Microscopy Laboratory and Cryo-EM platform at PKU.
- □ Most data processing was done at PKU HPC platform.
- Our work was partly funded by NSFC, Young 1000-Talent Plan, Intel Corporation, CLS@PKU and NIH.

Thank you for your attention!