cryoSPARC



Ali Punjani University of Toronto Structura Biotechnology Inc.

EMBO Birkbeck Sept 2019

CryoSPARC

Algorithm Development, Software, Performance

- Complete single-particle processing workflow
- Freely available for academia, supported by industry



- Under heavy development, regularly updated
- Methods and algorithms
 - Computational statistics, optimization theory, machine learning
- Professionally designed and built UI/software platform
- Ultra high-performance





- Thin film of liquid sample on grid
- Rapid plunge-freezing to vitrify and preserve sample
- Results in protein particles frozen in random orientations in thin film of ice



Vitrified Sample in Amorphous Ice

- Ice thickness/quality
- Particle concentration
- Sample heterogeneity
- Orientation distributions
- Denaturation/aggregation





- Pre-processing
- Locate particles in ice
- Exclude junk/broken particles
- 2D-to-3D reconstruction
- Resolving heterogeneity
- High-resolution refinement

Single-particle workflow



Tutorial Dataset

T20S Proteasome

- Data collected at NRAMM (Titan Krios)
- Publically depositited at EMPIAR (10025)
- 0.6575 A pixel size (super resolution)
- 8k x 8k, 38 frames
- 53 electrons/A²
- Large, stable, no heterogeneity, no flexibility
- Subset of 20 movies
- TIFF format compression



Start Motion Correction

T20S Proteasome

- Full Frame Motion Correction (Multi) job
- Use output from import job
- Default params

General Principle

- Maximize matching score over trajectories
- Global motion and local motion



Full frames or patches

Variant of *alignparts_Imbfgs* [Rubinstein & Brubaker 2015]

- Maximize cross-correlation between aligned average and frames
- Optimize over trajectories with gradient based quasi-Newton optimizer (LMBFGS)
- Apply a smoothness penalty to trajectory



Variant of alignparts

- Maximize cross-correlation between aligned average and frames
- Optimize over trajectories with gradient based quasi-Newton optimizer (LMBFGS)
- Apply a smoothness penalty to trajectory





Variant of alignparts

- Maximize cross-correlation between aligned average and frames
- Optimize over trajectories with gradient based quasi-Newton optimizer (LMBFGS)
- Apply a smoothness penalty to trajectory



How much smoothing?

- New cross-validation approach for data-driven optimal smoothing:
- Optimal Smoothing
 - No tuning parameters
 - 10 seconds per movie



Cross Validation Error

New patch-based Local motion correction

- Global rigid motion and local sample deformation together
- No need for particle positions
- Performs interpolation to create motion corrected micrographs
- Dose weighting



CTF Estimation

New Patch-based Local CTF Estimation

- Fast GPU Implementation
- Automatically estimates defocus variation for tilted, bent, deformed samples
- Accurate for all particle sizes and type including flexible and membrane proteins



Noble et al eLife 2018

CTF Estimation

New Patch-based Local CTF Estimation

- Fast GPU Implementation
- Automatically estimates defocus variation for tilted, bent, deformed samples
- Accurate for all particle sizes and type including flexible and membrane proteins



Exposure curation

Often tedious

- Based primarily on motion trajectories and CTF fit quality scores
- Looking for thinnest ice, least motion
- Exposure curation tool



Particle Picking

Template or blob based picking

- Correlation-based scoring
 - Template bias (Einstein-from noise)
- Local power score
 - Unbiased measure of signal content
- GPU Implementation:
 - < 1 second per micrograph
- Unbiased blob-based particle picker
 - Circular and elliptical blobs





Inspect and curate massive datasets

- Remove junk particles and improve homogeneity of particle stacks
- Scales well with dataset size
- Online-EM algorithm



125,000 particles with 50 classes in 12 minutes on 1x GPU

-	1425 ptcls	Diff peck	793, pre5	758 ptcb	698 ptcls	668 ptcis 9.6 A 1 ess	612 piets	602 picts	S61 ptcls
536 ptcis	534 ptcis 13.5 A 1 ess	508 ptcis	471 ptcis	466 ptcls	462 ptcls	437 ptcls	421 ptcls	407 ptcls	401 ptcls
382 ptels	354 ptcls	336 picis 13.4 A 1 ess	330 ptcls	313 ptels 13.7 A 1 ess	299 ptris 16.6 A 1 ess	297 ptcls 21.7 A 1 ess	276 ptcls 972,5 A 1 ess	260 ptcls	259 ptcls
252 ptcls 21.8 A 1 ess	240 ptcls 23.1 A 1 ess	227 ptcls 21.7 A 1 ess	217 ptcls 21.9 A 1 ess	212 ptcls 22.7 A 1 ess	204 ptcls 22.1 A 1 ess	203 ptcls 25.8 A 1 ess	104 ptcls 23.7 A 1 ess	169 ptcls 22.6 A1 ess	160 ptcls 23.5 A 1 ess
157 ptcls 22.3 A 1 ess	152 ptcls 24.7 A 1 ess	140 ptcls 23.7 A 1 ess	139 ptcls 23.8 A 1 ess	129 ptels 24.1 A 1 ess	127 ptdis 972,5 A 1 ess	117 ptcls 24.3 A Less	110 ptck 25.1 A 1 ess	96 ptcls 25.7.4 1 ess	91 ptcls 26.3 A-1 ess

125,000 particles with 50 classes in 12 minutes on 1x GPU 1,000,000 particles with 50 classes in 2 hours on 1x GPU

5588 ptcts 5.2 A 1 ess	5379 ptcls	5119 ptcts	4713 ptcls	4487 ptcls	4164 ptc/s	4087 ptcls	3857 ptcts	3825 ptcls	3779 ptcis
3760 ptcls	3732 ptcls 6.0 A 1 ess	3711 ptcls	3531 ptcls 5.9 A 1 ess	3529 ptcls 2 5.9 A 1 ess	3384 ptcls	3349 ptcls	3321 ptcls 6.0 A 1 ess	3291 ptcls	3160 ptcls
3041 ptcls 5.9 A 1 ess	2743 ptcls	2676 ptcls	2584 ptcls 5.9 A 1 ess	2531 ptcls 6.0 A 1 ess	2513 ptcls	2501 ptcls	2304 ptcls	2263 ptcls	2206 ptcls 6.0 A 1 ess
1922 ptcls	1598 ptcls	1596 ptcls 7.3 A 1 ess	1588 ptcls 16.1 A 1 ess	1568 ptcls	1255 pixes 652.3 A 1 ess	1220 ptcls 9.5 A 1 ess	1161 ptcls 12.6 A 1 ess	1133 ptcls 14.8 A 1 ess	1021 ptcls 7.6 A 1 ess
937 ptcls	741 ptcls 19.2 A 1 ess	733 ptcls 22.1 A 1 ess	711 ptcls 17.1 A 1 ess	641 ptcis 17.5 A 1 ess	559 ptcls 18.3 A 1 ess	256 ptcls 31.7 A 2 ess	254 ptcls 28,5 A 2 ess	109 ptcls 24.4 A 1 ess	1 ptcls

125,000 particles with 50 classes in <mark>12 minutes</mark> on 1x GPU 1,000,000 particles with 50 classes in <mark>2 hours</mark> on 1x GPU 200,000 particles with 200 classes in <mark>2 hours</mark> on 1x GPU



SPA 3D Reconstruction

What makes the problem difficult

- Unknown pose of each particle
 - 3D Orientation + 2D Shift
- High noise level
 - Irreducible due to beam damage
- Many particles (100,000+)
- Corruption by microscope contrast transfer function
- Multiple conformational states or distinct particles
- Flexibility and disorder





Probabilistic graphical model

• Fully "Bayesian" treatment:

 $p(V|X_1, X_2, \dots X_N)$

- Meaning we want to know all the possible 3D structures that could explain the images
- Arrow from V to X is the image formation model



Making inference computationally feasible

- Maximum probability estimate: $\max_V p(V|X_1,X_2,...X_N)$
- Now only looking for the single estimate of V that best explains the data
- Different choices of V yield all SPA techniques



Various choices of V

• Homogeneous refinement:

V = V

• Discrete Heterogeneity:

 $V = V_k | k = C_i$

• Local refinement/multiple rigid bodies:

 $V = V_1 + R_{rel}(V_2)$



• In all cases, need to solve an optimization problem: $\max_V p(V|X_1,X_2,...X_N)$

Likelihood function

- Likelihood of observing an image given V depends on noise, CTF, etc
- Marginalization or maximization over poses *R*



$$\max_{V} p(V|X_{i=1..N}) = \max_{V} \prod_{i}^{N} p(X_{i}|V) \cdot p(V)$$
$$= \max_{V} \left(\prod_{i}^{N} \int_{R} p(X_{i}|V,R) dR \right) \cdot p(V)$$
$$\approx \max_{V} \max_{R_{i}} \prod_{i}^{N} p(X_{i}|V,R_{i}) \cdot p(V)$$

Optimization problems

Mathematical properties

$$\max_{x} f(x)$$

- Optimization is a very well studied topic in mathematics
- Almost all other computational problems in EM are optimization problems:
 - Motion correction
 - CTF estimation
 - Model building
- Major characteristics of optimization problems are key

Optimization problems

Property: Convexity



- Single global optimum
- Guaranteed unique solution
- Multiple local optima
- Need to "explore"
- Generally inexpensive to solve Generally exponential cost

Optimization problems

Property: Conditioning



- Least coupling of variables
- Simple to solve



- Highly coupled variables
- Difficult to solve

Optimization problem for SPA

Homogeneous case

- Globally Non-convex:
 - multiple local optima
- Poor-conditioning:
 - small change in angle changes the structure, and change in structure changes angles
- Locally convex and wellconditioned
 - Once near enough to a local optimum





Optimization algorithms

Coordinate Ascent

$$\max_{V} \max_{R_i} \prod_{i}^{N} p(X|V, R_i) \cdot p(V)$$

- Also known as iterative refinement, expectationmaximization
 - Hold first variable fixed, optimize the second
 - Hold second variable fixed, optimize the first
- Provably guaranteed to converge:
 - <u>Globally</u> for a convex problem
 - To the <u>nearest local optimum</u> for a non-convex problem
- Convergence rate:
 - Very fast for well conditioned problems
 - Very, very slow for poor condition

Optimization algorithms

Coordinate Ascent



- Hold first variable fixed, optimize the second
- Hold second variable fixed, optimize the first
- Each subproblem becomes simple





Iterative refinement (Expectation Maximization)



Optimization algorithms

Coordinate Ascent

- Well suited for optimization of our SPA problem <u>once near the optimal solution</u>
- How to solve without "knowing the answer"?
- How does this change in the more complex case of heterogeneity?
- Consider a different optimization algorithm:
 - Stochastic gradient descent
Optimization algorithms

Gradient Descent

- Can take large steps, in multiple dimensions at once
- Gradient is direction of steepest descent:

$$G = \frac{d}{dV} \left(\prod_{i}^{N} p(X|V, R_{i}) \cdot p(V) \right)$$

Optimization algorithms

Stochastic gradient descent

 $\max_{V} p(V|X_{1}, X_{2}, X_{3}, X_{4}, X_{5}, X_{6}, X_{7}, X_{8}, \dots X_{N})$ $\text{Iter 1} \quad \max_{V} p(V| \quad X_{2}, \quad X_{4}, X_{5}, \quad)$ $\text{Iter 2} \quad \max_{V} p(V|X_{1}, \quad X_{3}, \quad X_{7}, \quad)$ $\text{Iter 3} \quad \max_{V} p(V| \quad X_{6}, \quad X_{8}, \quad X_{N})$

Randomly select small subsets of images at each iteration



Many **noisy** incremental changes

Ab-initio reconstruction: SGD

- Stochastic optimization
 - Class of modern statistical methods
- Stochastic Gradient Descent (SGD)
 - Very successful variant
 - Powers modern deep learning



Ab-initio reconstruction: 80S Ribosome





Data from Wong et al. eLife 2014

Ab-initio reconstruction: SGD



Heterogeneity

Explore vs. Exploit

- Additional classes add more variables and more non-convexity (local optima)
- Plain SGD alone is no longer effective
- Exploration vs. Exploitation
- In SGD controlled by "step size" and "minibatch size"
- In cryoSPARC: "class similarity score"



Heterogeneous samples: ATPase



Heterogeneous samples: Holiday Junction



Heterogeneous samples: AAA+ Unfoldase



Important considerations

- Cannot fix missing views
- Too few classes: average of structures
- Too many classes: 3D classes start to just become 2D and have only a single view
- Reproducibility: can run multiple times each time will use a different random seed
- Symmetry can cause issues if enforced

Multi-reference Refinement (3D classification)

• One or more starting models (from ab-initio reconstruction)



High resolution refinement



- Iterative refinement can proceed from coarse structure
- Usually very computationally expensive
- Many refinements in a complex workflow



2D-to-3D image alignment



- 5D pose search is expensive, for every image
- Existing techniques search exhaustively or locally
- Branch-and-bound reduces computational expense













High-speed high-resolution refinement

- Iterative refinement is very expensive – multiple rounds of alignment and reconstruction required
- Branch and bound (BnB) algorithm drastically reduces computation required without any loss in quality
- High-resolution structures in minutes
- Refinements are repeated dozens of times in a real cryoEM workflow

80S ribosome dataset: 105,000 particles refined to 3.2Å

11 min on 1x NVIDIA V100



Punjani et al. Nature Methods 2017

Measuring Resolution

- Dataset split in half
- Maps reconstructed from each half
- Compare corresponding shells of Fourier components
- Check where Fourier Shell Correlation drops below a threshold
- Mask is critical!



Iterative refinement (Expectation Maximization)



Preventing overfitting: Gold-Standard FSC [Scheres and Chen, 2012]



Orientation distribution



Local resolution and filtering

• Fast GPU implementations



Conventional refinement:

- Grounded in Fourier basis
 - Fourier slice theorem
 - CTF correction
 - 2D-3D Backprojection
 - Fourier Shell Correlation



Key intuition:

- Regularization in Fourier basis inappropriate
 - Disordered regions
 - Detergent, lipid nanodisc, etc
 - Solvent-protein boundary
 - Fractional occupancy
 - Flexible or highly dynamic regions
- Non-uniform refinement:

Use a basis that is localized in space and frequency





Oliver Clarke & Filippo Mancia, Columbia University

STRA6 receptor - 180 kDa membrane protein



Oliver Clarke & Filippo Mancia, Columbia University

STRA6 receptor - 180 kDa membrane protein



Non-uniform Refinement

Local and focused refinement

Masking, signal subtraction and non-uniform refinement





Local and focused refinement

Masking, signal subtraction and non-uniform refinement



Conformational Heterogeneity

"3D classification" is clustering

- Significant biological insight in conformational landscapes
- Existing algorithms can deal with discrete heterogeneity



- How many classes?
- Continuous heterogeneity?







• New algorithm!

3D Variability Analysis (New! cryoSPARC v2.9)





- New algorithm to solve for top K eigenvectors
- Eigenvectors correspond to molecular motions
- Accounting for CTF, viewing directions, all particles simultaneously, high resolution •





3D Variability Analysis – discrete heterogeneity

Directly determining 3D conformations

- Can measure "reaction coordinates" of individual particles
- Low-dimensional space directly shows clusters for discrete states
- Can run hierarchically



130K 50S Ribosome Particles (EMPIAR 10076) 3 variability dimensions



Non-trivial workflows

- Many tools needed for dealing with complex molecules
- Multiple iterations through processing pipeline required for near-atomic resolutions and identification of ligand binding sites
- Typical dataset sizes now multiple millions +


Robust software features and intuitive user experience



- ✓ Directly input and decompress raw data
- Project and workspace organization
- ✓ Tree view
- Drag and drop job builder
- ✓ Auto-tuning of hundreds of parameters
- ✓ Automated one-click workflows
- ✓ Interactive jobs
- ✓ Real-time experiment details and plots
- ✓ Direct downloads
- ✓ Computational resource manager
- ✓ Smart queuing/cluster integration
- ✓ User management
- ✓ Interoperability with other EM packages

cryoSPARC v2 software system



Multi-node and cluster support

Modularized

• Master and worker nodes



Scheduler

Smarter queueing

- Multiple independent lanes
- Each lane has workers on which jobs can be launched
- Within a lane, first in first out
- Jobs wait for their dependencies to complete before running
 - Easy to chain and queue jobs and let run autonomously



CryoSPARC Team



Ali Punjani



Suhail Dawood



Saara Punjani









Jay Yoo



David Fleet, PhD

University of Toronto



John Rubinstein, PhD

SickKids Hospital University of Toronto



Marcus Brubaker, PhD York University



Projects, Workspaces, Jobs

Tree structure with subtrees



- Jobs connect via results
- Complex workflows

Job builder

Drag and drop to create new jobs

- Connect inputs
- Set parameters
- Queue job



Results and Groups

High level and low level containers



Projects, Workspaces, Jobs

Tree structure with subtrees



• Workspaces cut up workflows into manageable chunks within a project