

Structural Biology in the Multi-Omics Era

Caitlyn McCafferty, Eric J Verbeke, Edward M Marcotte, and David W Taylor

J. Chem. Inf. Model., **Just Accepted Manuscript** • DOI: 10.1021/acs.jcim.9b01164 • Publication Date (Web): 04 Mar 2020

Downloaded from pubs.acs.org on March 9, 2020

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

Structural Biology in the Multi-Omics Era

Caitlyn L. McCafferty^{1*}, Eric J. Verbeke¹, Edward M. Marcotte^{1,2,3}, David W. Taylor¹⁻⁴

¹Department of Molecular Biosciences, University of Texas at Austin, Austin, TX 78712, USA.

²Institute for Cellular and Molecular Biology, University of Texas at Austin, Austin, TX 78712, USA.

³Center for Systems and Synthetic Biology, University of Texas at Austin, Austin, TX 78712, USA

⁴LIVESTRONG Cancer Institutes, Dell Medical School, Austin, TX 78712, USA.

*Correspondence to: clmccafferty@utexas.edu

Abstract

Rapid developments in cryo-electron microscopy have opened new avenues to probe the structures of protein assemblies in their near native states. Recent studies have begun applying single particle analysis to heterogeneous mixtures, revealing the potential of structural-omics approaches that combine the power of mass spectrometry and electron microscopy. Here, we highlight advances and challenges in sample preparation, data processing, and molecular modeling for handling increasingly complex mixtures. Such advances will help structural-omics methods extend to cellular level models of structural biology.

1
2
3 With the sequencing of thousands of genomes, large biological data sets (-omics data)
4 have become pervasive in most fields of biology, including development^{1, 2}, the classification of
5 organisms^{3, 4}, and disease⁵⁻⁷, among many others. Disciplines embracing -omics strategies reach
6 well beyond the central dogma of biology—genomics, transcriptomics, and proteomics—into such
7 areas as metabolomics⁸, epigenomics⁹, pharmacogenomics¹⁰, and interactomics¹¹. And, as with
8 these other endeavors, structural biology too has expanded to embrace -omics approaches.
9
10
11

12
13
14 Major historic interactions of structural biology and -omics approaches have included, for
15 example, electron tomography¹² to provide cellular context and spatial information to complement
16 proteomics and interactomics data¹³⁻¹⁵, many efforts at proteome-scale modeling of 3D structures
17 and interactions¹⁶⁻¹⁸, and the entire field of structural genomics¹⁹⁻²². Structural genomics has
18 employed techniques such as X-ray crystallography, NMR spectroscopy, and electron
19 microscopy (EM) to solve structures of purified macromolecules in a high-throughput manner,
20 targeting new protein folds and entire proteomes, which have been supplemented by molecular
21 modeling and structure prediction to extend structural insights to new molecules.
22
23
24
25

26 27 **The potential of shotgun cryo-EM methods**

28

29
30 More recently, advances in single particle cryo-electron microscopy (cryo-EM) have
31 opened interesting new opportunities to connect -omics approaches and structural biology. In
32 particular, cryo-EM boasts several important features: it only requires small amounts of sample,
33 there is no requirement for crystal screening and optimization, and as a result, it is possible to
34 capture several states of a macromolecular machine of interest. Cryo-EM is also capable of
35 imaging a large field of individual macromolecular complexes in a single image. With the advent
36 of direct electron detectors, ultra-stable electron microscopes, automated data collection
37 strategies²³ and real-time data processing²⁴, the 'resolution revolution' in cryo-EM provides a
38 definite route forward for increasing the throughput of structural biology²⁵. We can anticipate that
39 structures from these methods, in combination with electron tomography, will produce
40 information-rich cell atlases capturing high-resolution structures of the proteome and its spatial
41 context that will synergize with other -omics approaches. Here, we focus specifically on efforts to
42 increase single particle cryo-EM applicability to increasingly complex and heterogenous samples,
43 approaching cell lysates in complexity (as in shotgun cryo-EM), thus furthering the transformation
44 of cryo-EM into a pipeline for structural-omics.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Mass spectrometry combined with electron microscopy has been shown to be well suited for characterizing the architecture of protein complexes without purifying a specific target molecule, as demonstrated in yeast¹⁶, *Desulfovibrio vulgaris*²⁶, macrophage cytoplasm²⁷, the nuclear pore complex²⁸⁻³⁰, and most recently, in *Plasmodium falciparum*³¹. Protein-protein interactions identified through mass spectrometry in conjunction with advances in 3D structure determination have been used to investigate the architecture of multiple distinct protein complexes from mixtures such as fractionated cell lysate or even single cells³²⁻³⁴. Thus far, such studies have largely been limited to the identification of protein complexes that were easily recognizable (e.g. the proteasome and ribosome) or of high enough resolution to identify the proteins by comparing contiguous stretches of highly resolved amino acids to a reference proteome³¹. Currently, the field lacks robust and systematic computational pipelines for sorting, identifying, and molecular modeling of the myriad of structures that can potentially be solved from mixtures. The question remains: how can we break through these barriers?

Challenges in sample preparation of heterogeneous mixtures

In fact, even before the challenges of molecular modelling of mixtures of structures obtained from shotgun cryo-EM methods, several challenges exist for high-throughput cryo-EM data collection and processing of mixtures. Sample preparation is often a major bottleneck in structural studies. In our hands, finding suitable freezing conditions for heterogeneous mixtures has proven equally difficult as a single purified sample³⁵, with the addition of several new challenges. Notably, in the case of cell extracts, the presence of dominating, highly-abundant macromolecules can make screening difficult, especially when the size and shape of other, less abundant proteins are unfamiliar. Although multiple orthogonal chromatographic separations might help simplify mixtures, we find that sample preparation with similar size macromolecules improves chances of success. We have also found that different buffers in combination with different support substrates such as graphene oxide can produce an additional 'purification' step, ultimately determining which complexes are present on the grid. Furthermore, many 3D reconstructions are built from large datasets containing hundreds of thousands of particles per complex. Scaling this to samples containing tens to hundreds of complexes, which may be present in different quantities, could prove challenging simply from a data collection perspective. It will also be important to incorporate improved denoising and particle picking algorithms to assist users in picking difficult to recognize particles with multiple shapes and sizes³⁶⁻³⁸. Despite these challenges, several groups have already produced multiple structures to <5 Å resolution from fractionated lysates^{31, 32}.

1
2
3 While sample preparation methods are being worked on for investigating fractionated, or
4 whole cell lysates, there already exist many approaches which can be used to reduce the
5 complexity or target specific molecules from a mixture. Modified grid surfaces have been used for
6 capturing proteins by His-tag^{39,40}, biotin⁴¹, and antibody affinity⁴². These approaches can alleviate
7 the need for purification, target low abundance proteins, help with orientation bias, and can be
8 readily integrated in combination with clonal sets such as the ASKA library⁴³. Other approaches
9 include using microfluidic devices which can isolate and enrich target molecules⁴⁴. So far, many
10 of these studies have been limited to identifying only a few symmetric molecules from a mixture
11 and scaling these approaches for high-throughput has yet to be attempted.
12
13
14
15
16
17

18 **Advances in data processing of heterogeneous mixtures**

19
20
21 Apart from optimizing sample preparation and data collection, new data processing
22 schemes will also need to be introduced. Currently, most cryo-EM data processing software
23 operates under the assumption that samples contain one dominant structure which may contain
24 conformational or subunit heterogeneity. In order to adapt these software for use on highly
25 heterogenous samples, we developed an auxiliary algorithm based on the principles of the
26 projection-slice theorem to presort particles into homogenous subsets prior to conventional 3D
27 classification and therefore avoid the need to guess the number of underlying structures present
28 in the data³⁵. A subsequent challenge will be to identify the resulting models, which can range
29 from low- to high-resolution. Recently, the cryoID software package was introduced which uses a
30 unique approach to sequence by structure from highly-resolved, contiguous amino acids in a 3D
31 reconstruction³¹. However, the challenges from sample preparation suggest that it is more likely
32 that these studies will produce a number of low- to mid-resolution maps, and there still remains a
33 significant challenge for identifying and modelling low- to mid-resolution reconstructions from a
34 mixture when their identities are not known *a priori*.
35
36
37
38
39
40
41
42
43

44 **Approaches for docking atomic models into low- to mid-resolution reconstructions**

45
46 Due to the likelihood that lower abundance proteins in mixtures will only achieve low- to
47 mid- resolution 3D reconstructions, if simply as a function of fewer particles, there will continue to
48 be a need to better leverage other structural data. For this reason, an important focus remains
49 improving approaches for fitting both predicted and currently available atomic structures into these
50 lower resolution 3D reconstructions (**Figure 1**). These range from user-intensive to computation-
51 intensive approaches. Ideally, given the ambiguity of fitting numerous subunits into 3D
52 reconstructions of unknown identity, one would prefer a quick, efficient, and computationally
53
54
55
56
57
58
59
60

1
2
3 driven method. The challenge of fitting subunits into a 3D reconstruction becomes increasingly
4 difficult for multi-subunit complexes and may be additionally complicated by considerations of
5 symmetry. Techniques such as MBP and Fab labeling of individual subunits have been used to
6 identify specific subunits within multi-subunit complexes^{45, 46}. While this would prove cumbersome
7 for identifying proteins in multiple complexes within a cell lysate, it may be useful for targeting a
8 specific complex of interest.
9
10
11
12

13
14 One commonly employed user-driven approach for fitting atomic structures into 3D
15 reconstructions involves segmenting the maps either manually or using the Segger tool⁴⁷ followed
16 by rigid-body docking using Fit-in-Map into these segmented regions in UCSF Chimera^{48, 49}.
17 Scoring of this approach can be optimized using a flexible fitting tool^{50, 51} such as MDFF⁵⁰, which
18 applies forces proportional to the density gradient of the EM map, while conserving
19 stereochemistry, to fit atomic structures into EM maps with resolutions as low as 15 Å. While
20 these methods may work well if structural information is known *a priori*, any manual approach of
21 rigid docking faces the possibility of getting caught in a local minimum, suffering from user bias,
22 and requiring numerous user hours. Furthermore, fitting atomic models into complexes becomes
23 extremely challenging when their identities are incompletely known.
24
25
26
27
28
29

30
31 The development of integrative methods allows for a more hands-off approach, eliminating
32 some of these biases⁵²⁻⁵⁴. These approaches combine data retrieved from various experiments
33 such as yeast two-hybrid (Y2H) assays, mutagenesis, cross-linking, small angle X-ray scattering,
34 electron microscopy, and X-ray crystallography to build the multi-protein model^{55, 56}. Such
35 methodologies have been successful in building models for a number of multi-protein complexes
36 such as the nuclear pore complex⁵⁷, 16S rRNA complexed with methyltransferase A small
37 subunit⁵⁶, and the BBSome⁵⁸. Recently, several models predicted by integrative modeling were
38 validated against their experimentally determined high-resolution structures⁵². The results
39 showed that for all atom models the positions of subunit centers were within 5 Å of the true model,
40 demonstrating the power of this approach⁵⁹⁻⁶⁴. For those structures with resolution higher than 10
41 Å, secondary structure elements can not only be detected, but orientation and connectivity may
42 be predicted to validate the integrative models⁶⁵. While these methods are promising for building
43 a single multi-protein assembly with abundant data, they are computationally intensive, and
44 whether they will be equally applicable to mixtures of multiple complexes from structural-omics
45 data remains untested. Methods that could simplify model building by further constraining possible
46 orientations, interactions, or flexibility, may help moving forward.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Approaches for identifying molecular machines within complex mixtures

Due to the size and complexity of the data that describes extremely heterogeneous samples, corresponding mass spectrometry data becomes pivotal in identifying the proteins present, estimating their relative abundances, and identifying those that interact to form complexes in the sample. Previous studies have shown that machine learning combined with co-fractionation mass spectrometry can be used to detect proteins that interact to form complexes based on their elution profiles from multiple separation techniques⁶⁶. These predicted complexes can be prioritized by relative abundance for modeling. Additionally, identification of previously solved structures could reduce the number of 3D reconstructions which need to be considered for subsequent modeling. Pipelines such as GEM-PRO could accomplish this by streamlining rapid searches of the PDB by returning protein structures given a gene or protein sequence, while also evaluating quality of the structures and preparing sequences for comparative modeling for those that do not have a known structure⁶⁷.

Recently, improved shape-based searches for protein complexes have been developed to better accommodate the low- to mid-resolution EM data produced from tomography⁶⁸. Such shape search tools might prove useful to search 3D reconstructions in order to identify those known from prior structures. The 3D reconstructions that have been resolved and identified could then be used to revisit raw micrographs and pick specific particles with template matching approaches⁶⁹. The remaining 3D models would subsequently have to be built *de novo*, based on, e.g., protein identities from mass spectrometry performed on the same samples. Importantly, beyond the structures of proteins already solved and available in the Protein Data Bank⁷⁰, 3D structural models have now been computationally generated by many research groups at the proteome scale, a success of the Protein Structure Initiative (such as those indexed by the Uniprot⁷¹ database), using techniques of comparative modeling^{67, 72}, evolutionary couplings⁷³, or even *ab initio*⁷⁴ approaches.

Any structural modeling of native protein assemblies would most likely require prior knowledge of which specific protein-protein interactions were occurring^{75, 76}, as well as the stoichiometries of the interacting subunits. The latter, if unknown, might be obtainable using mass spectrometry^{57, 66, 77-79}. Other approaches to deciphering stoichiometry might include using volume constraints where volumes of different numbers of individual subunits are compared to the volume of a 3D reconstruction. Cross-linking mass spectrometry, where large numbers of pair-wise protein interactions may be identified, can help in elucidating protein interaction partners⁸⁰.

1
2
3 Additionally, other pair-wise restraints may be added, such as protein docking predictions, to
4 reveal new assemblies⁸¹⁻⁸³. However, protein docking becomes significantly more complex with
5 more than two proteins and no knowledge of interaction interfaces or order of assembly.
6
7

8 **Moving towards structural-omics**

9
10
11 Given knowledge of interacting subunits and their stoichiometries, the task becomes fitting
12 them into the correct map in the correct assembly. The problem resembles a jigsaw puzzle, where
13 subunits must fit into the molecular envelope while respecting mutual packing interfaces. In
14 general, such packing problems are known to be NP-complete⁸⁴ and cannot be solved
15 computationally in polynomial time. Nonetheless, additional restraints can be brought to bear to
16 reduce the search complexity. For example, like a puzzle, one might determine interacting
17 interfaces among the subunits, either by docking¹⁸ or more approximate approaches, ideally
18 algorithms that are rapid and partner-specific. In our own work, we have developed reduced
19 representations of protein surfaces to help predict complementary interaction interfaces, which
20 add a measure of robustness to minor structural deformations upon binding⁸⁵. Combinations of
21 such packing restraints could then be employed to help pack and refine 3D protein structures to
22 EM maps. In parallel, researchers have improved computational search algorithms for packing
23 problems by using reduction or backtracking^{86, 87}, and the potential exists to crowd-source the
24 problem, employing the visual acuity of humans to manually fit subunits into 3D reconstructions⁸⁸.
25
26
27
28
29
30
31
32
33

34
35 Structural-omics stands to benefit strongly from the cryo-EM resolution revolution, in turn
36 these approaches have the potential to greatly enhance our understanding of biology from a
37 systems perspective. Towards this end, it is already clear that various low- to high-resolution
38 complexes may be reconstructed from a cell lysate using single particle electron microscopy. The
39 development of new computational tools to efficiently sort and build atomic models into these low-
40 to mid-resolution reconstructions, or to solve the high-resolution structures from mixtures of
41 increasing complexity, will certainly help to further advance this field and put it on a path towards
42 even richer structural cell atlases.
43
44
45
46
47

48 **Acknowledgements**

49
50
51 This work was supported in part by Welch Foundation Research Grants F-1938 (to D.W.T.) and
52 F-1515 (to E.M.M.), Army Research Office Grant W911NF-15-1-0120 (to D.W.T.), a Robert J.
53 Kleberg, Jr. and Helen C. Kleberg Foundation Medical Research Award (to D.W.T.), and grants
54 from the National Institutes of Health (GM122480, DK110520, HD085901) to E.M.M.. C.L.M is an
55
56
57
58
59
60

1
2
3 NSF Graduate Research Fellow supported by the National Science Foundation (2019238253). D.W.T
4 is a CPRIT Scholar supported by the Cancer Prevention and Research Institute of Texas
5 (RR160088) and an Army Young Investigator supported by the Army Research Office (W911NF-
6 19-1-0021).
7
8
9
10
11
12

13 References:
14
15

- 16 1. Wang, Z.; Gerstein, M.; Snyder, M., RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Gen.* **2009**, 10, 57.
- 17 2. Kumar, P.; Tan, Y.; Cahan, P., Understanding development and stem cells using single
18 cell-based analyses of gene expression. *Development* **2017**, 144, 17-32.
- 19 3. Joyce, A. R.; Palsson, B. Ø., The model organism as a system: integrating 'omics' data
20 sets. *Nat. Rev. Mol. Cell Biol.* **2006**, 7, 198.
- 21 4. Raupach, M. J.; Amann, R.; Wheeler, Q. D.; Roos, C., The application of “-omics”
22 technologies for the classification and identification of animals. *Org. Divers. Evol.* **2016**, 16, 1-
23 12.
- 24 5. Karczewski, K. J.; Snyder, M. P., Integrative omics for health and disease. *Nat. Rev.*
25 *Gen.* **2018**, 19, 299.
- 26 6. Hasin, Y.; Seldin, M.; Lusi, A., Multi-omics approaches to disease. *Genome Biol.* **2017**,
27 18, 83.
- 28 7. Potter, S. S., Single-cell RNA sequencing for the study of development, physiology and
29 disease. *Nat. Rev. Nephrol.* **2018**, 14, 479.
- 30 8. Riekeberg, E.; Powers, R., New frontiers in metabolomics: from measurement to insight.
31 *F1000Research* **2017**, 6.
- 32 9. Jones, P. A.; Baylin, S. B., The epigenomics of cancer. *Cell* **2007**, 128, 683-692.
- 33 10. Daly, A. K., Pharmacogenetics: a general review on progress to date. *Br. Med. Bull.*
34 **2017**, 124, 65-79.
- 35 11. Luck, K.; Sheynkman, G. M.; Zhang, I.; Vidal, M., Proteome-scale human interactomics.
36 *Trends Biochem. Sci.* **2017**, 42, 342-354.
- 37 12. Lučić, V.; Förster, F.; Baumeister, W., Structural studies by electron tomography: from
38 cells to molecules. *Annu. Rev. Biochem.* **2005**, 74, 833-865.
- 39 13. Güell, M.; Van Noort, V.; Yus, E.; Chen, W.-H.; Leigh-Bell, J.; Michalodimitrakis, K.;
40 Yamada, T.; Arumugam, M.; Doerks, T.; Kühner, S., Transcriptome complexity in a genome-
41 reduced bacterium. *Science* **2009**, 326, 1268-1271.
- 42 14. Kühner, S.; van Noort, V.; Betts, M. J.; Leo-Macias, A.; Batisse, C.; Rode, M.; Yamada,
43 T.; Maier, T.; Bader, S.; Beltran-Alvarez, P., Proteome organization in a genome-reduced
44 bacterium. *Science* **2009**, 326, 1235-1240.
- 45 15. Yus, E.; Maier, T.; Michalodimitrakis, K.; van Noort, V.; Yamada, T.; Chen, W.-H.;
46 Wodke, J. A.; Güell, M.; Martínez, S.; Bourgeois, R., Impact of genome reduction on bacterial
47 metabolism and its regulation. *Science* **2009**, 326, 1263-1268.
- 48 16. Aloy, P.; Böttcher, B.; Ceulemans, H.; Leutwein, C.; Mellwig, C.; Fischer, S.; Gavin, A.-
49 C.; Bork, P.; Superti-Furga, G.; Serrano, L., Structure-based assembly of protein complexes in
50 yeast. *Science* **2004**, 303, 2026-2029.
- 51 17. Baker, D.; Sali, A., Protein structure prediction and structural genomics. *Science* **2001**,
52 294, 93-96.
53
54
55
56
57
58
59
60

18. Vakser, I. A., Protein-protein docking: From interaction to interactome. *Biophys. J.* **2014**, 107, 1785-1793.
19. Kim, S.-H., Shining a light on structural genomics. *Nat. Struct. Mol. Biol.* **1998**, 5, 643.
20. Skolnick, J.; Fetrow, J. S.; Kolinski, A., Structural genomics and its importance for gene function analysis. *Nat. Biotechnol.* **2000**, 18, 283.
21. Stevens, R. C.; Yokoyama, S.; Wilson, I. A., Global efforts in structural genomics. *Science* **2001**, 294, 89-92.
22. Chandonia, J.-M.; Brenner, S. E., The impact of structural genomics: expectations and outcomes. *Science* **2006**, 311, 347-351.
23. Li, Y.; Cash, J. N.; Tesmer, J. J.; Cianfrocco, M. A., High-throughput cryo-EM enabled by user-free preprocessing routines. *bioRxiv* **2019**.
24. Tegunov, D.; Cramer, P., Real-time cryo-electron microscopy data preprocessing with Warp. *Nat. Methods* **2019**, 1-7.
25. Kühlbrandt, W., The resolution revolution. *Science* **2014**, 343, 1443-1444.
26. Han, B.-G.; Dong, M.; Liu, H.; Camp, L.; Geller, J.; Singer, M.; Hazen, T. C.; Choi, M.; Witkowska, H. E.; Ball, D. A., Survey of large protein complexes in *D. vulgaris* reveals great structural diversity. *PNAS* **2009**, 106, 16580-16585.
27. Maco, B.; Ross, I. L.; Landsberg, M. J.; Mouradov, D.; Saunders, N. F.; Hankamer, B.; Kobe, B., Proteomic and electron microscopy survey of large assemblies in macrophage cytoplasm. *Mol. Cell. Proteomics* **2011**, 10, M111. 008763.
28. Alber, F.; Dokudovskaya, S.; Veenhoff, L. M.; Zhang, W.; Kipper, J.; Devos, D.; Suprpto, A.; Karni-Schmidt, O.; Williams, R.; Chait, B. T., The molecular architecture of the nuclear pore complex. *Nature* **2007**, 450, 695.
29. Alber, F.; Dokudovskaya, S.; Veenhoff, L. M.; Zhang, W.; Kipper, J.; Devos, D.; Suprpto, A.; Karni-Schmidt, O.; Williams, R.; Chait, B. T., Determining the architectures of macromolecular assemblies. *Nature* **2007**, 450, 683.
30. Beck, M.; Hurt, E., The nuclear pore complex: understanding its function through structural insight. *Nat. Rev. Mol. Cell Biol.* **2017**, 18, 73.
31. Ho, C.-M.; Li, X.; Lai, M.; Terwilliger, T. C.; Beck, J. R.; Wohlschlegel, J.; Goldberg, D. E.; Fitzpatrick, A. W.; Zhou, Z. H., Bottom-up structural proteomics: cryoEM of protein complexes enriched from the cellular milieu. *Nat. Methods* **2019**, 1-7.
32. Kastritis, P. L.; O'Reilly, F. J.; Bock, T.; Li, Y.; Rogon, M. Z.; Buczak, K.; Romanov, N.; Betts, M. J.; Bui, K. H.; Hagen, W. J., Capturing protein communities by structural proteomics in a thermophilic eukaryote. *Mol. Syst. Biol.* **2017**, 13, 936.
33. Verbeke, E. J.; Mallam, A. L.; Drew, K.; Marcotte, E. M.; Taylor, D. W., Classification of single particles from human cell extract reveals distinct structures. *Cell Rep.* **2018**, 24, 259-268. e3.
34. Yi, X.; Verbeke, E. J.; Chang, Y.; Dickinson, D. J.; Taylor, D. W., Electron microscopy snapshots of single particles from single cells. *J. Biol. Chem.* **2019**, 294, 1602-1608.
35. Eric J. Verbeke, Y. Z., Andrew P. Horton, Anna L. Mallam, David W. Taylor, Edward M. Marcotte, Separating distinct structures of multiple macromolecular assemblies from cryo-EM projections. *J. Struct. Biol.* **2019**, 107416.
36. Bepler, T.; Noble, A. J.; Berger, B., Topaz-Denoise: general deep denoising models for cryoEM. *bioRxiv* **2019**, 838920.
37. Wagner, T.; Merino, F.; Stabrin, M.; Moriya, T.; Antoni, C.; Apelbaum, A.; Hagel, P.; Sitsel, O.; Raisch, T.; Prumbaum, D., SPHIRE-crYOLO is a fast and accurate fully automated particle picker for cryo-EM. *Commun. Biol.* **2019**, 2, 1-13.
38. Bepler, T.; Morin, A.; Noble, A. J.; Brasch, J.; Shapiro, L.; Berger, B. Positive-unlabeled convolutional neural networks for particle picking in cryo-electron micrographs. In Research in computational molecular biology:... Annual International Conference, RECOMB...: proceedings. RECOMB (Conference: 2005-), 2018; NIH Public Access: 2018; Vol. 10812; pp 245-247.

- 1
2
3 39. Kelly, D. F.; Abeyrathne, P. D.; Dukovski, D.; Walz, T., The Affinity Grid: a pre-fabricated
4 EM grid for monolayer purification. *J. Mol. Biol.* **2008**, 382, 423-433.
- 5 40. Benjamin, C. J.; Wright, K. J.; Bolton, S. C.; Hyun, S.-H.; Krynski, K.; Grover, M.; Yu, G.;
6 Guo, F.; Kinzer-Ursem, T. L.; Jiang, W., Selective capture of histidine-tagged proteins from cell
7 lysates using TEM grids modified with NTA-graphene oxide. *Sci. Rep.* **2016**, 6, 1-11.
- 8 41. Han, B.-G.; Walton, R. W.; Song, A.; Hwu, P.; Stubbs, M. T.; Yannone, S. M.; Arbeláez,
9 P.; Dong, M.; Glaeser, R. M., Electron microscopy of biotinylated protein complexes bound to
10 streptavidin monolayer crystals. *J. Struct. Biol.* **2012**, 180, 249-253.
- 11 42. Yu, G.; Li, K.; Huang, P.; Jiang, X.; Jiang, W., Antibody-based affinity cryoelectron
12 microscopy at 2.6-Å resolution. *Structure* **2016**, 24, 1984-1990.
- 13 43. Kitagawa, M.; Ara, T.; Arifuzzaman, M.; Ioka-Nakamichi, T.; Inamoto, E.; Toyonaga, H.;
14 Mori, H., Complete set of ORF clones of Escherichia coli ASKA library (A Complete Set of E.
15 coli K-12 ORF Archive): Unique Resources for Biological Research. *DNA Res.* **2005**, 12, 291-
16 299.
- 17 44. Schmidli, C.; Albiez, S.; Rima, L.; Righetto, R.; Mohammed, I.; Oliva, P.; Kovacic, L.;
18 Stahlberg, H.; Braun, T., Microfluidic protein isolation and sample preparation for high-resolution
19 cryo-EM. *PNAS* **2019**, 116, 15007-15012.
- 20 45. Lander, G. C.; Estrin, E.; Matyskiela, M. E.; Bashore, C.; Nogales, E.; Martin, A.,
21 Complete subunit architecture of the proteasome regulatory particle. *Nature* **2012**, 482, 186-
22 191.
- 23 46. Wang, Y.; Ding, Z.; Liu, X.; Bao, Y.; Huang, M.; Wong, C. C.; Hong, X.; Cong, Y.,
24 Architecture and subunit arrangement of the complete *Saccharomyces cerevisiae* COMPASS
25 complex. *Sci. Rep.* **2018**, 8, 1-10.
- 26 47. Pintilie, G.; Chiu, W., Comparison of Segger and other methods for segmentation and
27 rigid-body docking of molecular components in Cryo-EM density maps. *Biopolymers* **2012**, 97,
28 742-760.
- 29 48. Pettersen, E. F.; Goddard, T. D.; Huang, C. C.; Couch, G. S.; Greenblatt, D. M.; Meng,
30 E. C.; Ferrin, T. E., UCSF Chimera—a visualization system for exploratory research and
31 analysis. *J. Comput. Chem.* **2004**, 25, 1605-1612.
- 32 49. Goddard, T. D.; Huang, C. C.; Ferrin, T. E., Visualizing density maps with UCSF
33 Chimera. *J. Struct. Biol.* **2007**, 157, 281-287.
- 34 50. Trabuco, L. G.; Villa, E.; Mitra, K.; Frank, J.; Schulten, K., Flexible fitting of atomic
35 structures into electron microscopy maps using molecular dynamics. *Structure* **2008**, 16, 673-
36 683.
- 37 51. Kovacs, J. A.; Galkin, V. E.; Wriggers, W., Accurate flexible refinement of atomic models
38 against medium-resolution cryo-EM maps using damped dynamics. *BMC Struct. Biol.* **2018**, 18,
39 12.
- 40 52. Braitbard, M.; Schneidman-Duhovny, D.; Kalisman, N., Integrative structure modeling:
41 overview and assessment. *Annu. Rev. Biochem.* **2019**, 88.
- 42 53. Topf, M.; Lasker, K.; Webb, B.; Wolfson, H.; Chiu, W.; Sali, A., Protein structure fitting
43 and refinement guided by cryo-EM density. *Structure* **2008**, 16, 295-307.
- 44 54. Russel, D.; Lasker, K.; Webb, B.; Velázquez-Muriel, J.; Tjioe, E.; Schneidman-Duhovny,
45 D.; Peterson, B.; Sali, A., Putting the pieces together: integrative modeling platform software for
46 structure determination of macromolecular assemblies. *PLoS Biol.* **2012**, 10, e1001244.
- 47 55. Webb, B.; Viswanath, S.; Bonomi, M.; Pellarin, R.; Greenberg, C. H.; Saltzberg, D.; Sali,
48 A., Integrative structure modeling with the integrative modeling platform. *Protein Sci.* **2018**, 27,
49 245-258.
- 50 56. van Zundert, G. C.; Melquiond, A. S.; Bonvin, A. M., Integrative modeling of biomolecular
51 complexes: HADDOCKing with cryo-electron microscopy data. *Structure* **2015**, 23, 949-960.
- 52
53
54
55
56
57
58
59
60

- 1
2
3 57. Kim, S. J.; Fernandez-Martinez, J.; Nudelman, I.; Shi, Y.; Zhang, W.; Raveh, B.;
4 Herricks, T.; Slaughter, B. D.; Hogan, J. A.; Upla, P., Integrative structure and functional
5 anatomy of a nuclear pore complex. *Nature* **2018**, 555, 475.
- 6 58. Chou, H.-T.; Apelt, L.; Farrell, D. P.; White, S. R.; Woodsmith, J.; Svetlov, V.; Goldstein,
7 J. S.; Nager, A. R.; Li, Z.; Muller, J., The molecular architecture of native BBSome obtained by
8 an integrated structural approach. *Structure* **2019**, 27, 1384-1394. e4.
- 9 59. Miled, N.; Yan, Y.; Hon, W.-C.; Perisic, O.; Zvelebil, M.; Inbar, Y.; Schneidman-Duhovny,
10 D.; Wolfson, H. J.; Backer, J. M.; Williams, R. L., Mechanism of two classes of cancer mutations
11 in the phosphoinositide 3-kinase catalytic subunit. *Science* **2007**, 317, 239-242.
- 12 60. Schweppe, D. K.; Chavez, J. D.; Lee, C. F.; Caudal, A.; Kruse, S. E.; Stuppard, R.;
13 Marcinek, D. J.; Shadel, G. S.; Tian, R.; Bruce, J. E., Mitochondrial protein interactome
14 elucidated by chemical cross-linking mass spectrometry. *PNAS* **2017**, 114, 1732-1737.
- 15 61. Murakami, K.; Elmlund, H.; Kalisman, N.; Bushnell, D. A.; Adams, C. M.; Azubel, M.;
16 Elmlund, D.; Levi-Kalisman, Y.; Liu, X.; Gibbons, B. J., Architecture of an RNA polymerase II
17 transcription pre-initiation complex. *Science* **2013**, 342, 1238724.
- 18 62. Leitner, A.; Joachimiak, L. A.; Bracher, A.; Mönkemeyer, L.; Walzthoeni, T.; Chen, B.;
19 Pechmann, S.; Holmes, S.; Cong, Y.; Ma, B., The molecular architecture of the eukaryotic
20 chaperonin TRiC/CCT. *Structure* **2012**, 20, 814-825.
- 21 63. Luo, J.; Cimermancic, P.; Viswanath, S.; Ebmeier, C. C.; Kim, B.; Dehecq, M.; Raman,
22 V.; Greenberg, C. H.; Pellarin, R.; Sali, A., Architecture of the human and yeast general
23 transcription and DNA repair factor TFIIH. *Mol. Cell* **2015**, 59, 794-806.
- 24 64. Shi, Y.; Fernandez-Martinez, J.; Tjioe, E.; Pellarin, R.; Kim, S. J.; Williams, R.;
25 Schneidman-Duhovny, D.; Sali, A.; Rout, M. P.; Chait, B. T., Structural characterization by
26 cross-linking reveals the detailed architecture of a coatomer-related heptameric module from the
27 nuclear pore complex. *Mol. Cell. Proteomics* **2014**, 13, 2927-2943.
- 28 65. Lindert, S.; Alexander, N.; Wötzel, N.; Karakaş, M.; Stewart, P. L.; Meiler, J., EM-fold: de
29 novo atomic-detail protein structure determination from medium-resolution density maps.
30 *Structure* **2012**, 20, 464-478.
- 31 66. Havugimana, P. C.; Hart, G. T.; Nepusz, T.; Yang, H.; Turinsky, A. L.; Li, Z.; Wang, P. I.;
32 Boutz, D. R.; Fong, V.; Phanse, S., A census of human soluble protein complexes. *Cell* **2012**,
33 150, 1068-1081.
- 34 67. Brunk, E.; Mih, N.; Monk, J.; Zhang, Z.; O'Brien, E. J.; Bliven, S. E.; Chen, K.; Chang, R.
35 L.; Bourne, P. E.; Palsson, B. O., Systems biology of the structural proteome. *BMC Syst. Biol.*
36 **2016**, 10, 26.
- 37 68. Han, X.; Sit, A.; Christoffer, C.; Chen, S.; Kihara, D., A global map of the protein shape
38 universe. *PLoS Comput. Biol.* **2019**, 15, e1006969.
- 39 69. Rickgauer, J. P.; Grigorieff, N.; Denk, W., Single-protein detection in crowded molecular
40 environments in cryo-EM images. *Elife* **2017**, 6, e25648.
- 41 70. Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.;
42 Shindyalov, I. N.; Bourne, P. E., The protein data bank. *Nucleic Acids Res.* **2000**, 28, 235-242.
- 43 71. Apweiler, R.; Bairoch, A.; Wu, C. H.; Barker, W. C.; Boeckmann, B.; Ferro, S.; Gasteiger,
44 E.; Huang, H.; Lopez, R.; Magrane, M., UniProt: the universal protein knowledgebase. *Nucleic*
45 *Acids Res.* **2004**, 32, D115-D119.
- 46 72. Lam, S. D.; Das, S.; Sillitoe, I.; Orengo, C., An overview of comparative modelling and
47 resources dedicated to large-scale modelling of genome sequences. *Acta Crystallographica*
48 *Section D: Structural Biology* **2017**, 73, 628-640.
- 49 73. Marks, D. S.; Hopf, T. A.; Sander, C., Protein structure prediction from sequence
50 variation. *Nat. Biotechnol.* **2012**, 30, 1072.
- 51 74. Simons, K. T.; Bonneau, R.; Ruczinski, I.; Baker, D., Ab initio protein structure prediction
52 of CASP III targets using ROSETTA. *Proteins: Struct., Funct., Bioinf.* **1999**, 37, 171-176.
- 53
54
55
56
57
58
59
60

- 1
2
3 75. Drew, K.; Lee, C.; Huizar, R. L.; Tu, F.; Borgeson, B.; McWhite, C. D.; Ma, Y.;
4 Wallingford, J. B.; Marcotte, E. M., Integration of over 9,000 mass spectrometry experiments
5 builds a global map of human protein complexes. *Mol. Syst. Biol.* **2017**, *13*, 932.
- 6 76. Giurgiu, M.; Reinhard, J.; Brauner, B.; Dunger-Kaltenbach, I.; Fobo, G.; Frishman, G.;
7 Montrone, C.; Ruepp, A., CORUM: the comprehensive resource of mammalian protein
8 complexes—2019. *Nucleic Acids Res.* **2018**, *47*, D559-D563.
- 9 77. Hernández, H.; Robinson, C. V., Determining the stoichiometry and interactions of
10 macromolecular assemblies from mass spectrometry. *Nat. Protoc.* **2007**, *2*, 715.
- 11 78. Skinner, O. S.; Schachner, L. F.; Kelleher, N. L., The Search Engine for Multi-Proteoform
12 Complexes: An Online Tool for the Identification and Stoichiometry Determination of Protein
13 Complexes. *Curr. Protoc. Bioinf.* **2016**, *56*, 13.30. 1-13.30. 11.
- 14 79. Smits, A. H.; Vermeulen, M., Characterizing protein–protein interactions using mass
15 spectrometry: challenges and opportunities. *Trends Biotechnol.* **2016**, *34*, 825-834.
- 16 80. Liu, F.; Rijkers, D. T.; Post, H.; Heck, A. J., Proteome-wide profiling of protein
17 assemblies by cross-linking mass spectrometry. *Nat. Methods* **2015**, *12*, 1179.
- 18 81. Dominguez, C.; Boelens, R.; Bonvin, A. M., HADDOCK: a protein– protein docking
19 approach based on biochemical or biophysical information. *J. Am. Chem. Soc.* **2003**, *125*, 1731-
20 1737.
- 21 82. Comeau, S. R.; Gatchell, D. W.; Vajda, S.; Camacho, C. J., ClusPro: an automated
22 docking and discrimination method for the prediction of protein complexes. *Bioinformatics* **2004**,
23 *20*, 45-50.
- 24 83. Gong, X.; Wang, P.; Yang, F.; Chang, S.; Liu, B.; He, H.; Cao, L.; Xu, X.; Li, C.; Chen,
25 W., Protein–protein docking with binding site patch prediction and network-based terms
26 enhanced combinatorial scoring. *Proteins: Struct., Funct., Bioinf.* **2010**, *78*, 3150-3155.
- 27 84. Demaine, E. D.; Demaine, M. L., Jigsaw puzzles, edge matching, and polyomino
28 packing: Connections and complexity. *Graphs Combin.* **2007**, *23*, 195-208.
- 29 85. McCafferty, C. L.; Marcotte, E. M.; Taylor, D. W., Simplified geometric representations of
30 protein structures identify complementary interaction interfaces. *bioRxiv* **2019**.
- 31 86. Knuth, D. E., Dancing links. *arXiv preprint cs/0011047* **2000**.
- 32 87. Lodi, A.; Martello, S.; Vigo, D., Heuristic algorithms for the three-dimensional bin packing
33 problem. *Eur. J. Oper. Res.* **2002**, *141*, 410-420.
- 34 88. Khatib, F.; Desfosses, A.; Players, F.; Koepnick, B.; Flatten, J.; Popović, Z.; Baker, D.;
35 Cooper, S.; Gutsche, I.; Horowitz, S., Building de novo cryo-electron microscopy structures
36 collaboratively with citizen scientists. *PLoS Biol.* **2019**, *17*, e3000472.
- 37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

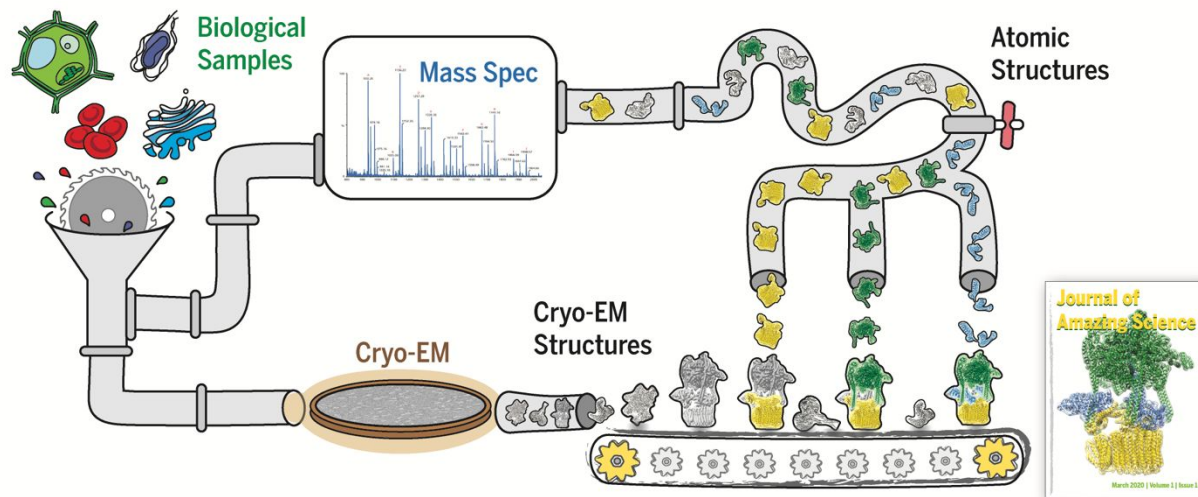
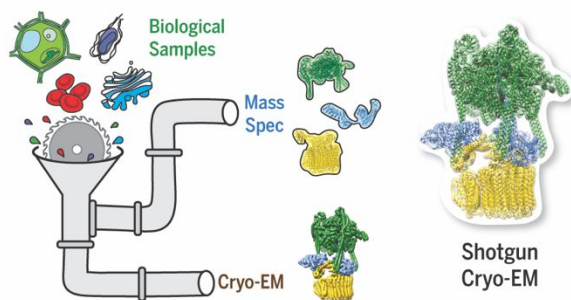


Figure 1: A structural-omics pipeline. A broad goal in the field is for a high-throughput, structural-omics approach for reconstructing complexes from a heterogeneous mixture. For example, whole cell lysates, organelle lysates, and heterogeneous mixtures might be analyzed by both cryo-EM and mass spectrometry. Cryo-EM produces multiple 3D reconstructions of protein complexes, while mass spectrometry provides identity information for the proteins present in the sample. To merge the two, even more efficient computational pipelines are needed to build or retrieve individual structures of proteins, organize them by interactions, assemble them into complexes, and match them to their 3D reconstructions obtained from a sample. Illustration by Angel Syrett.



TOC Figure

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

