



BIROn - Birkbeck Institutional Research Online

Patwardhan, A. and Brandt, R. and Butcher, S.J and Collinson, L. and Gault, D. and Grünewald, K. and Hecksel, C. and Huisken, J.T. and Iudin, A. and Jones, M.L. and Korir, P.K. and Koster, A.J. and Lagerstedt, I. and Lawson, C.L. and Mastrorade, D. and McCormick, M. and Parkinson, H. and Rosenthal, P.B. and Saalfeld, S. and Saibil, Helen R. and Sarntivijai, S. and Solanes Valero, I. and Subramaniam, S. and Swedlow, J.R. and Tudose, I. and Winn, M. and Kleywegt, G.J. (2017) Building bridges between cellular and molecular structural biology. eLife 6 , ISSN 2050-084X.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/19169/>

Usage Guidelines:

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html> or alternatively contact lib-eprints@bbk.ac.uk.

1 ***Building bridges between cellular and molecular structural***
2 ***biology***

3
4 Authors: Ardan Patwardhan^{1*}, Robert Brandt², Sarah J. Butcher³, Lucy Collinson⁴,
5 David Gault⁵, Kay Grünewald⁶, Corey Hecksel^{7,8}, Juha T. Huiskonen⁶, Andrii Iudin¹,
6 Martin L. Jones⁴, Paul K. Korir¹, Abraham J. Koster⁹, Ingvar Lagerstedt^{1,10}, Catherine
7 L. Lawson¹¹, David Mastronarde¹², Matthew McCormick¹³, Helen Parkinson¹, Peter
8 B. Rosenthal⁴, Stephan Saalfeld¹⁴, Helen R. Saibil¹⁵, Sirarat Sarntivijai¹, Irene
9 Solanes Valero^{1,16}, Sriram Subramaniam¹⁷, Jason R. Swedlow¹⁸, Ilinca Tudose¹,
10 Martyn Winn¹⁹, Gerard J. Kleywegt^{1*}

11
12
13 Affiliations:

14 ¹European Molecular Biology Laboratory, European Bioinformatics Institute,
15 Wellcome Genome Campus, Hinxton CB10 1SD, UK

16
17 ²FEI Visualization Sciences Group, Mérignac, France

18
19 ³Institute of Biotechnology and Department of Biosciences, University of Helsinki,
20 Helsinki, Finland

21
22 ⁴The Francis Crick Institute, London, UK

23
24 ⁵Centre for Gene Regulation and Expression, University of Dundee, Dundee, UK

25
26 ⁶Division of Structural Biology, Wellcome Trust Centre for Human Genetics, Oxford,
27 UK

28
29 ⁷National Center for Macromolecular Imaging, Verna and Marrs McLean Department
30 of Biochemistry and Molecular Biology, Baylor College of Medicine, Houston, Texas

31
32 ⁸Current affiliation: Electron Bio-Imaging Centre, Diamond Light Source Ltd, Didcot,
33 UK

34
35 ⁹Department of Molecular Cell Biology, Leiden University Medical Center, Leiden,
36 The Netherlands

37
38 ¹⁰Current affiliation: Computational Chemistry and Cheminformatics, Lilly UK,
39 Windlesham, UK

40
41 ¹¹Center for Integrative Proteomics Research and Research Collaboratory for
42 Structural Bioinformatics, Rutgers, The State University of New Jersey, New Jersey,
43 USA

44
45 ¹²Department of Molecular, Cellular, and Developmental Biology, University of
46 Colorado, Boulder, Colorado, USA

47

48 ¹³Kitware, Inc., Carrboro, NC, USA

49

50 ¹⁴Janelia Research Campus, Howard Hughes Medical Institute, Ashburn, VA, USA

51

52 ¹⁵Institute of Structural and Molecular Biology, Department of Crystallography,
53 Birkbeck College, London, UK

54

55 ¹⁶Current affiliation: University of Vic - Central University of Catalonia, Barcelona,
56 Spain

57

58 ¹⁷Center for Cancer Research, National Cancer Institute, Bethesda, Maryland, USA

59

60 ¹⁸Centre for Gene Regulation and Expression and Division of Computational Biology,
61 University of Dundee, Dundee, UK

62

63 ¹⁹Scientific Computing Department, Science and Technology Facilities Council,
64 Research Complex at Harwell, Didcot, United Kingdom

65

66 **Competing financial interests**

67 n/a

68

69 *Correspondence should be addressed to:

70 Patwardhan: Phone: + 44 1223 492649; Fax: + 44 1223 494 468; Email:

71 ardan@ebi.ac.uk

72

73 Kleywegt: Phone: + 44 1223 492663; Fax: + 44 1223 494 468; Email:

74 gerard@ebi.ac.uk

75

76

77 **Impact Statement**

78 The integration of structural data from different imaging scales requires the
79 development of standards and tools for representing the segmentation and
80 transformation of data, and for the annotation of biological structures.

81

82

83

84

85

86

87

88 **Abstract**

89

90 The integration of cellular and molecular structural data is key to understanding the
91 function of macromolecular assemblies and complexes in their *in vivo* context. Here
92 we report on the outcomes of a workshop that discussed how to integrate structural
93 data from a range of public archives. The workshop identified two main priorities: the
94 development of tools and file formats to support segmentation (that is, the
95 decomposition of a three-dimensional volume into regions that can be associated
96 with defined objects), and the development of tools to support the annotation of
97 biological structures.

98

99

100 Introduction

101

102 To obtain an integrated view of how molecular machinery operates inside cells,
103 biologists are increasingly combining structural data at different length scales,
104 obtained using a range of techniques such as electron tomography, electron
105 microscopy, NMR spectroscopy and X-ray crystallography. Structural data is held in
106 public archives such as the Electron Microscopy Data Bank (EMDB; [emdb-](http://emdb-empiar.org)
107 empiar.org; Tagari et al., 2002), the Electron Microscopy Public Image Archive
108 (EMPIAR; empiar.org; Iudin et al., 2016), and the Protein Data Bank (PDB;
109 wwpdb.org; Bernstein et al., 1977)

110

111 Integration between PDB and EMDB data is based on atomic models in the PDB that
112 have been fitted to or built into EMDB volume maps. For purified biological
113 molecules or larger defined complexes this approach is done routinely. Sequence
114 information from the models can be used to link to other bioinformatics resources
115 such as the Universal Protein Resource (UniProt; uniprot.org/; UniProt Consortium,
116 2013). However, atomic models are not always available for a variety of reasons,
117 such as when molecular averaging fails to obtain high-resolution features or the
118 inherently lower resolution studies when molecules are imaged in more complex or
119 even cellular environments. In such cases, the identification of features often relies
120 on prior knowledge or correlation of structural data obtained at different scales.

121

122 Once features have been identified, segmentation defined here as the
123 decomposition of the 3D volume into regions that can be associated with defined
124 objects, can be employed to facilitate and visualise the interpretation of the map. For
125 example, in a recent study the segmentation of electron and soft X-ray tomography
126 reconstructions was used to study leakage and breakage of the membranes in
127 erythrocytes infected by *Plasmodium falciparum*, and documented the dramatic
128 changes in the morphology of cells during egress (Hale et al., 2016). The soft X-ray
129 tomograms provided overviews of the membrane compartments in intact, vitrified
130 cells (Figure 1). It should be noted that the word 'segmentation' may have different
131 interpretations: for example, in whole animal, pre-clinical and medical imaging,
132 segmentation includes a concept of a model that is used for fitting of the features. In
133 this manuscript we limit the definition to the separation of density into distinct sub-
134 domains.

135

136 In tomography, where multiple copies of nearly identical objects are found, 3D sub-
137 tomogram averaging and 3D classification may be employed to obtain higher
138 resolution reconstructions. This process often involves combining information from
139 multiple tomograms. Since the higher resolution afforded by sub-tomogram
140 averaging provides more structural detail, displaying sub-tomogram averages at the
141 original tomogram positions and orientations may reveal important information about
142 the organization and distribution of the object within a cellular and functional context.
143 If properly annotated such data can be further mined with other questions in mind by
144 other researchers. For example, researchers recently created composite maps of
145 Lassa virus particles by inserting the sub-tomogram average structure of the Lassa
146 virus glycoprotein spike back into the original tomographic reconstructions, revealing

147 the organisation and copy number of the spikes on the virus surface (Figure 2; Li et
148 al., 2016). Another example revealed the lateral clustering of viral membrane
149 proteins mediating membrane fusion (Maurer et al., 2013).

150

151 The archiving of segmentation data in EMDB entries was identified as an area
152 requiring urgent attention in previous workshops on “Data-Management Challenges
153 in 3D Electron Microscopy” in 2011 (Patwardhan et al., 2012) and “A 3D Cellular
154 Context for the Macromolecular World” in 2012 (Patwardhan et al., 2014), as was
155 the improved biological annotation of structural data to make it more accessible to
156 the wider biological audience and to enable integration with structural and other
157 bioinformatics resources. Crucially for data integration we need “structured biological
158 annotation” which is here defined as the association of data with identifiers (e.g.,
159 accession codes from UniProt) and ontologies taken from well established
160 bioinformatics resources. (Ontologies are formal collections of statements defining
161 concepts, relationships and constraints; for example, the mitochondrial large and
162 small ribosomal subunits are parts of the mitochondrial ribosome which, in turn, is a
163 part of the mitochondrion). To our knowledge, none of the segmentation formats
164 widely used in electron microscopy and related fields currently support structured
165 biological annotation. Furthermore, spatial transformations relating sub-tomograms
166 to their parent tomograms are not currently captured in EMDB. Moreover, wider
167 usage of both segmentation and transformation data by non-expert users is hindered
168 by a plurality of formats.

169

170 To discuss and address the challenges of representing and capturing segmentations
171 and transformation data, the Protein Data Bank in Europe ([PDBe](#)) organised an
172 expert workshop on “3D Segmentations and Transformations - Building Bridges
173 between Cellular and Molecular Structural Biology” in December 2015. The
174 objectives were:

- 175 • To identify data models and formats for representing segmentation and
176 transformation data that could provide support for structured biological
177 annotation, thus facilitating their use by EMDB and enabling data-exchange
178 between different software packages
- 179 • To gain a better understanding of the challenges involved in the annotation of
180 electron microscopy data and develop requirements in terms of tools and
181 strategies to facilitate annotation.

182 Here we report and discuss the main outcomes of the workshop, which was attended
183 by a range of participants including software developers, users of segmentation
184 software, ontology experts, and experts in structure and data archiving.

185

186 **Data models and file formats for segmentations and** 187 **transformations**

188 Prior to the workshop, PDBe developed a draft data model to support segmentations
189 and their annotations in EMDB that could accommodate segmentation descriptions
190 from a range of existing formats and software packages as well as structured
191 biological annotation. It supported the key features of major segmentation packages

192 such as Amira (www.fei.com/software/amira/), IMOD (Kremer et al., 1996) and
193 Chimera (Pettersen et al., 2004), and provided scope for extension and flexibility as
194 the field developed. However, the draft data model did not cover minor features (e.g.,
195 surface rendering parameters), especially those that are only relevant in the context
196 of a particular software package. The data model was implemented in an XML
197 schema with the following features:

198 a) **Support for hierarchical segmentation description.** This is important for
199 representing segmentations from (semi-)automatic approaches that naturally result
200 in a hierarchal segmentation, such as Segger (which iteratively groups the results of
201 the initial watershed segmentation into a hierarchy; Pintilie et al., 2010).

202 b) **Different representations of segmentations.** Contours and simple geometric
203 primitives such as spheres and lines are often used to delineate regions of interest
204 (ROIs) when segmentation is performed manually. In automatic segmentation the
205 segments are typically represented as surface meshes and/or 3D volume masks. In
206 the latter case, run-length encoding and limited bit-depth are commonly used
207 techniques to minimise memory requirements. It could be argued that it would be
208 useful to have only one canonical representation and convert all the individual
209 representations to it. However, representing geometric primitives such as spheres as
210 surface meshes could lead to substantial increases in storage size and decreases in
211 accuracy of the descriptions.

212 c) **Support for externally defined (i.e., as separate files) 3D volume masks.** It
213 may be useful to allow separation between the metadata (annotations) and the
214 actual segmentations (e.g., to lessen the burden on tools and web-services that only
215 require the metadata). The data model accommodates links to external files (and
216 locations within these files) for representing segments.

217 d) **Segment colours.** In some application areas, colour is used to identify objects of
218 the same kind, so it is important that such information is not lost.

219
220 The draft data model was intended primarily for internal use in EMDB. However, the
221 meeting participants strongly favoured a broader scope so that the format could
222 serve the entire biological segmentation field. This would also make it easier to
223 support the development of translators between different formats and possibly
224 contribute to a reduction of the number of formats (or at least prevent further
225 proliferation of formats). Representatives for several major software packages used
226 for segmentation including IMOD (D.M.), Amira (R.B.) and Chimera (Tom Goddard,
227 personal communication) have expressed a commitment to providing read/write
228 capabilities for the developed format if standard libraries are made available.

229
230 The draft data model included support for various colour models including RGB, HSV
231 and colour names. Participants argued that it would be sufficient to support only the
232 most commonly used one, namely the RGB model, as the other models can be
233 converted to it.

234
235 Participants also noted that it might be useful to allow quantification of the estimated
236 certainty of a biological annotation, for example a score for the agreement between a
237 sub-tomogram average and a corresponding region from an originating tomogram.
238 There may also be alternative biological annotations in various combinations (logical

239 OR, XOR, AND, etc.). The quantification of alternative annotations could become
240 very complex to represent and use, and the participants agreed to initially limit the
241 scope to a single annotation per segment and to let the need for more complex
242 representations be driven by actual use cases.

243

244 Concerning the transformations between sub-tomogram averages and tomograms,
245 the participants agreed that this information should be incorporated into the
246 segmentation data model; it simply requires adding support for multiple
247 transformations of the same 3D volume representation. It was agreed that the
248 convention to define affine transformations should be well-defined in terms of the
249 transformation, the order in which they are applied, the direction of the
250 transformation, and the orientations and origins of coordinate systems.

251

252 With respect to correlative multi-modal imaging it was recognized that there would
253 eventually be a need to go beyond affine transformations, for example to represent
254 distortions and deformations of slice data, but the participants did not come to a
255 conclusion about a coherent extensible format. Often, a segment consists of multiple
256 spatially transformed copies of the same primitive. This is also relevant for sub-
257 tomogram averages as the same volume is to be spatially transformed into multiple
258 locations within a tomogram. To accommodate these situations, every segment can
259 be associated with a list of transformations. This representation will also be useful in
260 the context of template matching for describing the transformations between the
261 template and the 3D volume.

262

263 The draft data model was developed in XSD ([XML Schema Definition](#)). The definition
264 of data models is greatly facilitated by tools that enable GUI-based development of
265 schemas such as [Oxygen](#) and [XMLSpy](#). Code generators such as [generateDS](#)
266 create object-model wrappers from schemas that enable reading, writing and
267 manipulation of XML files, thus allowing for rapid prototyping. Various XML validators
268 also allow the correctness of a file relative to a schema to be tested. However,
269 concerns were raised about the verbosity of the XML format and the efficiency with
270 which it can be used. Participants proposed that while XML may be the natural
271 format for a schema defined in XSD, it would be useful to consider other more
272 compact and efficient formats such as JSON and [HDF5](#) (a binary format that allows
273 for efficient representation of hierarchal metadata and data in a single container).
274 Both JSON and HDF5 are now widely supported with libraries in most major
275 programming languages, including Python and C/C++, to facilitate reading and
276 writing. To this end, utilities to convert between the XML, JSON and HDF5
277 representations of the segmentation data model are currently in development at
278 PDBe.

279

280 Future format development will be an iterative process involving extensive
281 consultation with relevant stakeholders to obtain consensus in and support from the
282 community of developers, yielding a format that they will support. A "[Segmentation
283 and transformation file format working group](#)" has been established by a subset of the
284 workshop participants, and other developers working on segmentation who are
285 interested in joining the group are asked to contact AP.

286
287 PDBE has already modified the data model based on the feedback from the meeting,
288 and this will continue in several rounds of consultation with the working group. The
289 schema is versioned to keep track of changes. To facilitate adoption of the format,
290 dubbed EMDB-SFF (SFF=Segmentation File Format), PDBE is developing
291 translators to/from other commonly used formats. The code for these translators is
292 provided as free open source and distributed via the [CCP-EM SVN repository](#).
293 Comments on the [schema](#) should be sent to AP.

294 **Structured biological annotation**

295 As previously explained, structured biological annotation is the association of data
296 with identifiers and ontologies taken from well-established bioinformatics resources.
297 The use of structured biological annotation is not common practice in the electron
298 microscopy or structural biology communities. Therefore, ontology experts were
299 invited to the workshop to explain why these are useful and what resources and tools
300 are available for assigning annotations. Use-cases such as mouse imaging data
301 helped to explain the principles and practice of structured biological annotation. By
302 the end of the meeting there was a clearer appreciation of the importance of
303 structured biological annotation for searching and linking imaging data across
304 different scales, between different imaging and structural databases and with other
305 bioinformatics resources.

306
307 Structured annotation would enable the seamless integration of structural, imaging
308 and bioinformatics data from different resources, thus making it possible to provide
309 problem-centric views of biology that incorporate structural and imaging data and are
310 easily accessible by the broader biological community (and in contrast to the highly
311 specialised structure-centric resources that are available today and mainly serve
312 domain-specific communities). However, there were concerns that many in the
313 electron microscopy community would find navigating the landscape of ontologies
314 challenging and that this approach would only gain traction in the community if tools
315 were developed to simplify the biological annotation process.

316
317 It was also discussed whether annotation should be performed by the depositor or by
318 EMDB curators. While curators could be trained to a high level of expertise in the
319 use of ontologies, they would not necessarily have enough knowledge about the
320 sample and the specifics of the biological system underlying the study. It was
321 concluded that depositors should perform the annotation, with curators overseeing
322 and checking annotations.

323

324 **Tools for structured biological annotation**

325 Structured biological annotation for electron microscopy will rely on a range of
326 established ontologies such as Gene Ontology (GO; Gene Ontology, 2008),
327 Experimental Factor Ontology (EFO; Malone et al., 2010), Protein Ontology (PRO;
328 Natale et al., 2014), Cellular Microscopy Phenotype Ontology (CMPO; Jupp et al.,
329 2016), NCBI organismal classification ([NCBITaxon](#)), integrated cross-species for
330 anatomical structures (UBERON; Mungall et al., 2012, imaging modality and sample

331 preparation from Fbbi (Orloff et al., 2013), Foundational Model of Anatomy ([FMA](#))
332 and Cell Ontology (CL; Diehl et al., 2016). It may also include identifiers from
333 resources such as UniProt and the [Complex Portal](#), which in turn contain cross-
334 reference information to other useful standardised vocabularies and common
335 terminology identifiers, such as the [OMIM](#) and KEGG (Kanehisa et al., 2016)
336 pathways. This cross-reference information is useful when linking data coded with
337 these terminologies to the ontologies.

338
339 Several of these resources provide application programming interfaces (APIs) that
340 can be used to access the information programmatically and provide search
341 functionality. The Samples, Phenotypes and Ontologies Team (SPOT) at EMBL-EBI
342 has developed tools such as [Zooma](#) and the Ontology Lookup Service ([OLS](#); Jupp et
343 al., 2015), which aggregate information from a wide range of ontologies and provide
344 APIs to access these tools. These APIs can be used when building tools for
345 segmentation annotation to provide simplified views and search facilities for
346 ontological terms.

347
348 At the workshop, PDBe presented mock-ups of a web-based segmentation
349 annotation tool (SAT; Figure 3). This tool would allow a user to add structured
350 biological annotation to segmentations obtained from a variety of different software
351 packages and then output an annotated segmentation file in EMDB-SFF that could
352 be deposited to EMDB or EMPIAR. Annotation could either be done during
353 deposition, in which case the biological annotation from the segmentation file could
354 be harvested by the deposition system to facilitate the deposition process, or it could
355 be done post deposition. The workflow would consist of: (i) the user uploading
356 segmentation files (there could be several if the segments have been saved as
357 separate files) and the corresponding map (unless it is already released in EMDB or
358 EMPIAR); (ii) conversion to an EMDB-SFF file; (iii) use of a GUI-based interface to
359 view the segmentations overlaid on the map and to select segmentations and add
360 annotation; (iv) output of a fully annotated EMDB-SFF file that could be uploaded to
361 EMDB (Figure 4).

362
363 Two different options were presented for how annotation could take place (Figure 3).
364 Many macromolecular systems for which data are deposited in EMDB fall into broad
365 categories such as ribosomes, proteasomes, chaperonins and so on: for each of
366 these categories, and with the added information about taxonomy, lists of likely
367 components could be generated to facilitate annotation (Figure 3A). Similarly for
368 cellular level annotation, lists of cellular components could be used. As it would not
369 be possible to cover every potential scenario with pre-defined lists, the other option
370 is to provide a search facility that offers potentially applicable terms from available
371 ontologies (Figure 3B).

372
373 The workshop participants expressed strong support for the development of the SAT
374 and the functionality depicted in the mock-ups but raised concerns about a number
375 of issues: the upload of data to a web server – some users may find it challenging to
376 upload large maps and segmentation; the need to annotate segmentations twice –
377 users would typically add free text annotations in the software used for the

378 segmentation and would the need to re-annotate in the SAT; finding the 'right'
379 metadata terms (particularly in cases where a search yields more than one term, and
380 it is not clear which is the most relevant term); annotating a hierarchical
381 segmentation. (The SAT mock-up accommodates annotation on only one level of
382 hierarchy: this might be sufficient in many cases, but it could become problematic as
383 more automated segmentation techniques are developed and their usage expands.)
384

385 A desktop version of the SAT would help users concerned about the upload of large
386 amounts of data to a web-server. Another option would be to integrate the
387 functionality for structured biological annotation into existing packages such as
388 IMOD, Chimera and Amira; this would also avoid the problem of users having to
389 annotate the segmentations twice. This alternative would require the development of
390 libraries and widgets that facilitate the use of ontologies and the EMDB-SFF by third
391 parties. For example, the program for segmentation in IMOD already has a 'Name
392 Wizard' plugin that helps the user to choose standardized object names from a CSV
393 file: however, additional development would be need to provide access to on-line
394 ontologies.
395

396 Participants agreed that PDBe should start by developing the web-based SAT
397 because it could reuse a number of components that are already being used in other
398 electron microscopy-related web services (such as the Volume slicer; Salavert-
399 Torres et al., 2016), followed by the desktop version. Once the SAT reaches a
400 certain level of maturity PDBe could work with third-party developers to integrate the
401 annotation functionality into their packages.
402

403 By far the greatest challenge is developing the functionality to find the appropriate
404 biological metadata (Malone et al., 2016) and tools such as Zooma and OLS will be
405 useful for this purpose. A "[Segmentation annotation working group](#)" has been
406 established by a subset of the workshop participants to provide data sets and use
407 cases to aid the design of the SAT and to help with its testing. Members of the
408 electron microscopy community and related communities who are interested in
409 joining the group are asked to contact AP.
410

411 **Discussion**

412 The EMDB-SFF data model has undergone a round of updates based on the
413 feedback from the meeting. The development of the file format and the segmentation
414 annotation tools will be iterative, with user-testing and feedback from the working
415 groups being integral parts of the process. The file format and tools are expected to
416 be ready by late 2017, although they might not offer all the features discussed
417 above.
418

419 Wide acceptance and support of the EMDB-SFF format by software developers
420 working on segmentations and transformations will be crucial. Providing well-
421 documented open-source tools for working with the format will help in this regard,
422 and the Collaborative Computational Project for Electron cryo-Microscopy ([CCP-EM](#))

423 has committed to distributing these tools, and including them in training events for
424 users and developers. However, the scope of the format is not limited to the cryo-EM
425 field. For example, segmentation is an essential element of the workflow for
426 interpreting data in 3D scanning electron microscopy (3D-SEM; Patwardhan et al.,
427 2014). It will also be possible to provide support for segmentations for other imaging
428 modalities (and also for imaging on other length scales), although the range of
429 biological ontologies and vocabularies will need to be expanded. It should also be
430 possible to support techniques that combine imaging modalities (such as correlative
431 light and electron microscopy), but this will involve extra work on the transformation
432 model.

433
434 It was clear from the discussions regarding the annotation of segmentations that
435 there are significant language barriers between the fields. Overcoming these barriers
436 is a prerequisite for progress, as is the development of new tools that will facilitate
437 annotation.

438
439 This workshop was an important milestone in that it defined concrete actionable
440 outcomes to address the challenges involved in the integration of cellular and
441 molecular structural data in the public archives. This integration will provide
442 researchers with "problem-centric views" of data from many different sources, and
443 will also help the wider biological and medical communities by making make
444 structural data more accessible.

445 **Acknowledgements**

446 The workshop and the work on EMDB-SFF at PDBe were supported by the UK
447 Medical Research Council (MRC) with co-funding from the UK Biotechnology and
448 Biological Sciences Research Council (grant MR/L007835 to G.J.K.) and the
449 European Commission Framework 7 Programme (284209). Work on EMDB and
450 EMPIAR at PDBe is further supported by the Wellcome Trust (104948), the US
451 National Institutes of Health National Institute of General Medical Sciences (R01
452 GM079429), and EMBL-EBI. CCP-EM is supported by the MRC (MR/N009614/1).
453

454 **References**

- 455
456 1. Tagari, M., Newman, R., Chagoyen, M., Carazo, J.M. & Henrick, K. *Trends*
457 *Biochem. Sci.* **27**:589 (2002).
- 458 2. Iudin, A., Korir, P.K., Salavert-Torres, J., Kleywegt, G.J. & Patwardhan, A.
459 *Nat. Methods* **13**:387-388 (2016).
- 460 3. Bernstein, F.C., Koetzle, T.F., Williams, G.J., Meyer, E.F., Jr., Brice, M.D.,
461 Rodgers, J.R., Kennard, O., Shimanouchi, T. & Tasumi, M. *J Mol Biol* **112**:535-542
462 (1977).
- 463 4. UniProt Consortium. *Nucleic Acids Res* **41**:D43-47 (2013).

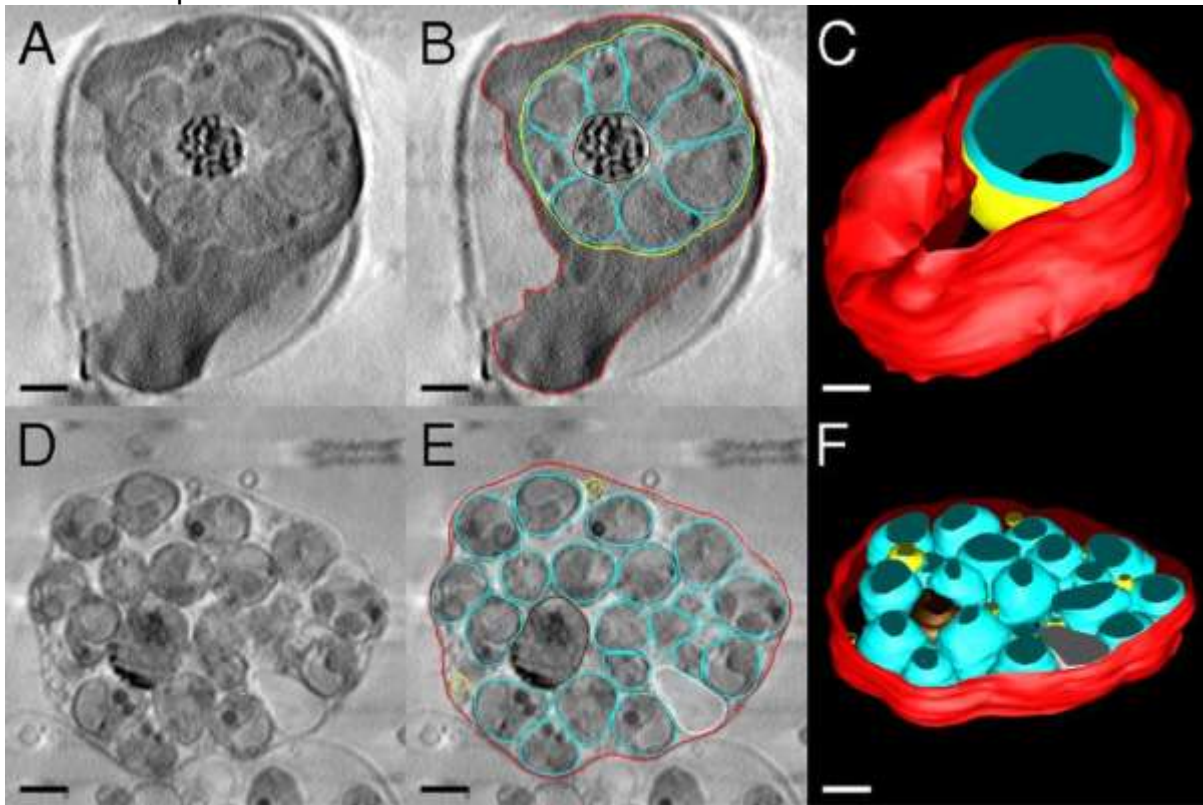
- 464 5. Hale, V.L., Watermeyer, J.M., Hackett, F., Vizcay-Barrena, G., van Ooij, C.,
465 Thomas, J.A., Spink, M.C., Harkiolaki, M., Duke, E., Fleck, R.A., Blackman, M.J. &
466 Saibil, H.R. *PNAS* **114**:3439-3444 (2017).
- 467 6. Li, S., Sun, Z., Pryce, R., Parsy, M.L., Fehling, S.K., Schlie, K., Siebert, C.A.,
468 Garten, W., Bowden, T.A., Strecker, T. & Huiskonen, J.T. *PLoS Pathogens*
469 **12**:e1005418 (2016).
- 470 7. Maurer, U.E., Zeev-Ben-Mordehai, T., Pandurangan, A.P., Cairns, T.M.,
471 Hannah, B.P., Whitbeck, J.C., Eisenberg, R.J., Cohen, G.H., Topf, M., Huiskonen,
472 J.T. & Grunewald, K. *Structure* **21**:1396-1405 (2013).
- 473 8. Patwardhan, A., Carazo, J.M., Carragher, B., Henderson, R., Heymann, J.B.,
474 Hill, E., Jensen, G.J., Lagerstedt, I., Lawson, C.L., Ludtke, S.J., Mastronarde, D.,
475 Moore, W.J., Roseman, A., Rosenthal, P., Sorzano, C.O., Sanz-Garcia, E., Scheres,
476 S.H., Subramaniam, S., Westbrook, J., Winn, M., Swedlow, J.R. & Kleywegt, G.J.
477 *Nat Struct Mol Biol* **19**:1203-1207 (2012).
- 478 9. Patwardhan, A., Ashton, A., Brandt, R., Butcher, S., Carzaniga, R., Chiu, W.,
479 Collinson, L., Doux, P., Duke, E., Ellisman, M.H., Franken, E., Grunewald, K.,
480 Heriche, J.K., Koster, A., Kuhlbrandt, W., Lagerstedt, I., Larabell, C., Lawson, C.L.,
481 Saibil, H.R., Sanz-Garcia, E., Subramaniam, S., Verkade, P., Swedlow, J.R. &
482 Kleywegt, G.J. *Nat Struct Mol Biol* **21**:841-845 (2014).
- 483 10. Kremer, J.R., Mastronarde, D.N. & McIntosh, J.R. *J. Struct. Biol.* **116**:71-76
484 (1996).
- 485 11. Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M.,
486 Meng, E.C. & Ferrin, T.E. *J Comput Chem* **25**:1605-1612 (2004).
- 487 12. Pintilie, G.D., Zhang, J., Goddard, T.D., Chiu, W. & Gossard, D.C. *J. Struct.*
488 *Biol.* **170**:427-438 (2010).
- 489 13. Gene Ontology, C. *Nucleic Acids Res* **36**:D440-444 (2008).
- 490 14. Malone, J., Holloway, E., Adamusiak, T., Kapushesky, M., Zheng, J.,
491 Kolesnikov, N., Zhukova, A., Brazma, A. & Parkinson, H. *Bioinformatics* **26**:1112-
492 1118 (2010).
- 493 15. Natale, D.A., Arighi, C.N., Blake, J.A., Bult, C.J., Christie, K.R., Cowart, J.,
494 D'Eustachio, P., Diehl, A.D., Drabkin, H.J., Helfer, O., Huang, H., Masci, A.M., Ren,
495 J., Roberts, N.V., Ross, K., Ruttenberg, A., Shamovsky, V., Smith, B., Yerramalla,
496 M.S., Zhang, J., AlJanahi, A., Celen, I., Gan, C., Lv, M., Schuster-Lezell, E. & Wu,
497 C.H. *Nucleic Acids Res* **42**:D415-421 (2014).
- 498 16. Jupp, S., Malone, J., Burdett, T., Heriche, J.K., Williams, E., Ellenberg, J.,
499 Parkinson, H. & Rustici, G. *J Biomed Semantics* **7**:28 (2016).

- 500 17. Mungall, C.J., Torniai, C., Gkoutos, G.V., Lewis, S.E. & Haendel, M.A.
501 *Genome Biol* **13**:R5 (2012).
- 502 18. Orloff, D.N., Iwasa, J.H., Martone, M.E., Ellisman, M.H. & Kane, C.M. *Nucleic*
503 *Acids Res* **41**:D1241-1250 (2013).
- 504 19. Diehl, A.D., Meehan, T.F., Bradford, Y.M., Brush, M.H., Dahdul, W.M.,
505 Dougall, D.S., He, Y., Osumi-Sutherland, D., Rutenber, A., Sarntivijai, S., Van
506 Slyke, C.E., Vasilevsky, N.A., Haendel, M.A., Blake, J.A. & Mungall, C.J. *J Biomed*
507 *Semantics* **7**:44 (2016).
- 508 20. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. *Nucleic*
509 *Acids Res* **44**:D457-462 (2016).
- 510 21. Jupp, S., Burdett, T., Leroy, C. & Parkinson, H. 2015. Collaborative ontology
511 development using the webulous architecture and Google App. *SWAT4LS*
512 *International Conference: Semantic Web Applications and Tools for Life Sciences*.
513 Cambridge. Available at: [http://www.swat4ls.org/wp-](http://www.swat4ls.org/wp-content/uploads/2015/10/SWAT4LS_2015_paper_32.pdf)
514 [content/uploads/2015/10/SWAT4LS_2015_paper_32.pdf](http://www.swat4ls.org/wp-content/uploads/2015/10/SWAT4LS_2015_paper_32.pdf)
- 515 22. Salavert-Torres, J., Iudin, A., Lagerstedt, I., Sanz-Garcia, E., Kleywegt, G.J. &
516 Patwardhan, A. *J Struct Biol* **194**:164-170 (2016).
- 517 23. Malone, J., Stevens, R., Jupp, S., Hancocks, T., Parkinson, H. & Brooksbank,
518 C. *PLoS Comput Biol* **12**:e1004743 (2016).
- 519 24. Muller, A., Beeby, M., McDowall, A.W., Chow, J., Jensen, G.J. & Clemons,
520 W.M., Jr. *Microbiologyopen* **3**:702-710 (2014).
- 521 25. Santarella-Mellwig, R., Pruggnaller, S., Roos, N., Mattaj, I.W. & Devos, D.P.
522 *PLoS Biology* **11**:e1001565 (2013).
- 523 26. Bennett, A., Liu, J., Van Ryk, D., Bliss, D., Arthos, J., Henderson, R.M. &
524 Subramaniam, S. *J Biol Chem* **282**:27754-27759 (2007).
- 525 27. Bennett, A.E., Narayan, K., Shi, D., Hartnell, L.M., Gousset, K., He, H.,
526 Lowekamp, B.C., Yoo, T.S., Bliss, D., Freed, E.O. & Subramaniam, S. *PLoS*
527 *Pathogens* **5**:e1000591 (2009).
- 528 28. Liu, J., Bartesaghi, A., Borgnia, M.J., Sapiro, G. & Subramaniam, S. *Nature*
529 **455**:109-113 (2008).
530
531
532

533 **Figure 1**

534 **Segmentation of *Plasmodium falciparum*-infected erythrocytes**

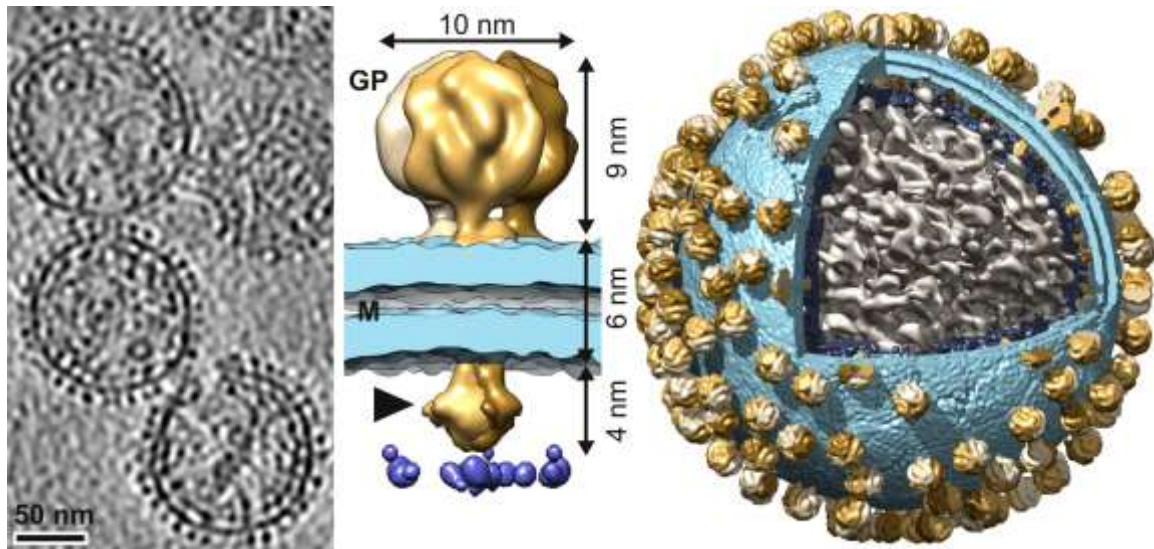
535 Soft X-ray tomography shows loss of mechanical integrity of the red cell membrane
536 in the final stages of egress. Panels A-C depict schizonts treated with a selective
537 malarial cGMP-dependent protein kinase G inhibitor (C2), and panels D-F depict
538 schizonts treated with a broad-spectrum cysteine protease inhibitor, E64, which
539 allows parasitophorous vacuole membrane (PVM) rupture but prevents erythrocyte
540 membrane rupture, resulting in merozoites trapped in the blood cell. (A) Slice from
541 tomogram of C2-arrested schizont. (B) Outlines of erythrocyte membrane (red), PVM
542 (yellow), and parasites (cyan) in the tomogram slice in A. (C) 3D rendering of the
543 schizont. The vacuole (yellow) is densely packed with merozoites (cyan) that have
544 been collectively rather than individually rendered, for clarity. The overall height of
545 the cell is $\sim 5 \mu\text{m}$. (D) Tomogram slice from an E64-arrested schizont, shown with
546 outlining of membranes in E. Remnants of the PVM are visible. (F) 3D rendering of
547 the schizont. Figure and legend adapted with permission from Hale et al., 2017.
548 Scale bar $1 \mu\text{m}$.



549
550

551 **Figure 2**
552 **Arrangement of Lassa virus glycoprotein spikes on the virion surface**

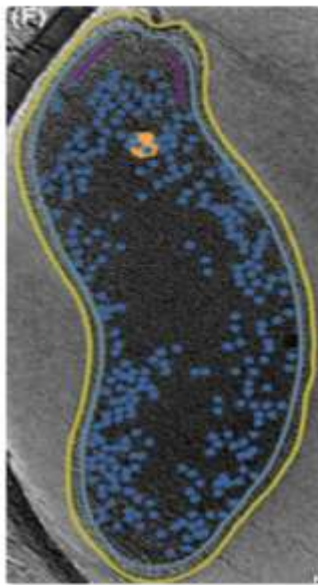
553
554 Left to right: A slice from a tomographic volume of Lassa viruses, a sub-tomogram
555 average of the glycoprotein spike, and the sub-tomogram average inserted back
556 onto a virus reconstruction. Images adapted from Li et al., 2016 (under a [CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)
557 license).



560
561

562 **Figure 3**
563 **Mock-up of a possible Segmentation-Annotation Tool (SAT)**
564

565 Image slices are shown with the segmentations overlaid. (A) The top right panel
566 presents a tree that enables the user to select the segment to be annotated, and
567 existing annotations are shown in the middle right panel. The bottom right panel
568 provides pre-defined lists of annotation terms for frequently studied assemblies and
569 complexes. The image in the left panel is adapted from Müller et al., 2014 (under a
570 [CC BY 3.0](#) license). (B) The top right and middle right panels are similar to those in
571 A. The bottom right panel provides a search option to find relevant terms. The image
572 in the left panel is adapted from Santarella-Mellwig et al., 2013 (under a [CC BY 4.0](#)
573 license).
574



▼ *Campylobacter jejuni*

- Outer Membrane (OM)
- Inner Membrane (IM)
- Chemoreceptors
- Ribosomes
- Storage Granules

Descriptors	Category	Database	ID

Edit annotations

Search for terms Frequently used terms Delete annotations

Assemblies and complexes

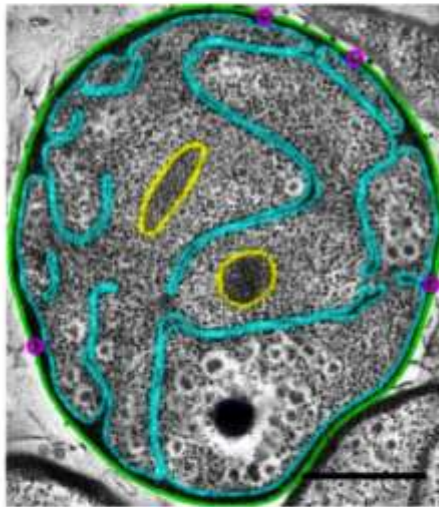
Descriptors	Category	Database	ID	Select
Ribosome	Cellular component	GO	GO_0005840	<input checked="" type="checkbox"/>
GroEL-GroES	Cellular component	GO	GO_1990220	<input type="checkbox"/>
RNA Polymerase Complex	Cellular component	GO	GO_0030880	<input type="checkbox"/>

Organelles

Previously used terms

Add annotations

(A)



▼ *Gemmata obscuriglobus*

- Outer Membrane (OM)
- Inner Membrane (IM)
- DNA
- OM invaginations

Descriptors	Category	Database	ID
membrane invagination	biological process	GO	GO_0010324

Edit annotations

Search for terms Frequently used terms Delete annotations

outer membrane

Descriptors	Category	Database	ID	Select
outer membrane	cellular component	GO	GO_0019867	<input checked="" type="checkbox"/>
outer membrane lipoprotein b1c	material entity	PRO	PR_000022228	<input type="checkbox"/>

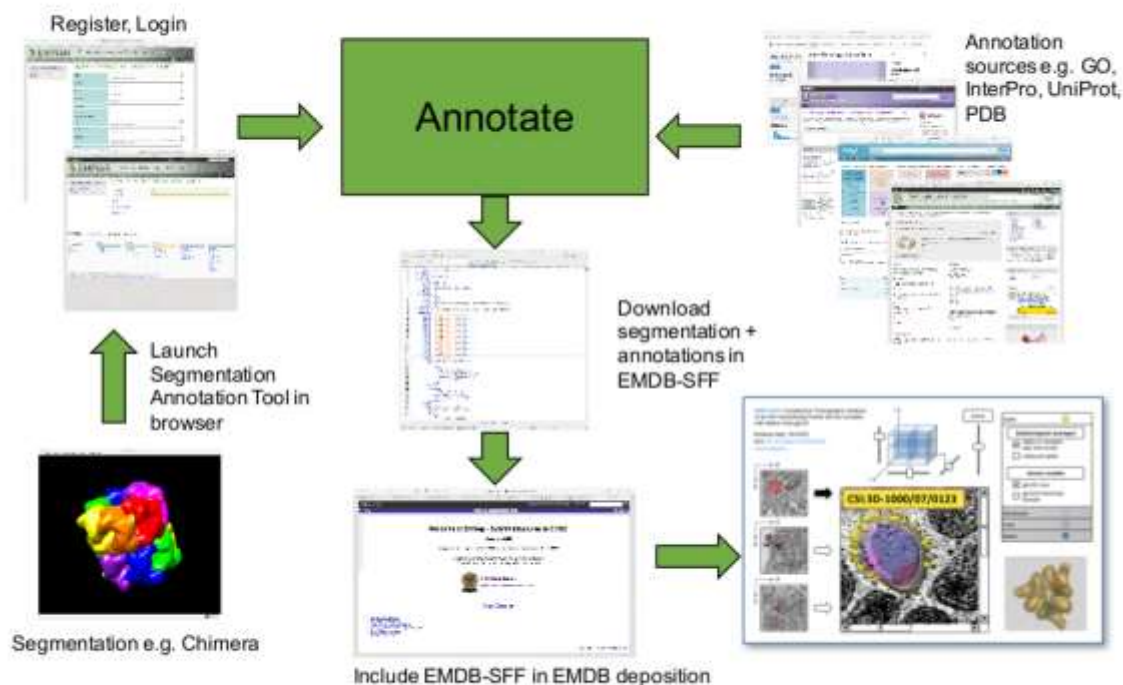
Add annotations

(B)

575
576
577
578

579 **Figure 4**
580 **Segmentation-annotation workflow**
581

582 A user launches the Segmentation-Annotation Tool and uploads segmentations
583 obtained with third-party software. After the segmentation has been annotated with
584 biologically meaningful terms, a segmentation file is written in EMDB-SFF format;
585 this file can be uploaded to the Electron Microscopy Data Bank when the structure is
586 deposited. Once released, the EMDB-SFF file can be used for the integration of
587 structural data between different imaging scales and across resources. The Volume
588 browser mock-up (bottom right) contains images adapted from Bennett et al., 2007
589 and Bennett et al., 2009 (under a [CC0 1.0](#) license). The 3D rendering was generated
590 from EMDB entry EMD-5020 and PDB entry 3dno (Liu et al., 2008).



591