

Sharing Data from Molecular Simulations

Mark James Abraham, Rossen Pavlov Apostolov, Jonathan Barnoud, Paul Bauer, Christian Blau, Alexandre M.J.J. Bonvin, Matthieu Chavent, John Damon Chodera, Karmen #ondi#-Jurki#, Lucie Delemotte, Helmut Grubmüller, Rebecca J. Howard, E. Joseph Jordan, Erik Lindal, O.H. Samuli Ollila, Jana Selent, Daniel G. A. Smith, Phill James Stansfeld, Johanna K. S. Tiemann, Mikael Trellet, Christopher J. Woods, and Artem Zhmurov

J. Chem. Inf. Model., **Just Accepted Manuscript** • DOI: 10.1021/acs.jcim.9b00665 • Publication Date (Web): 17 Sep 2019

Downloaded from pubs.acs.org on September 23, 2019

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.

SCHOLARONE™
Manuscripts

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

Sharing Data from Molecular Simulations

1
2
3
4 1
5
6
7
8
9
10 2 *Mark Abraham*¹, *Rossen Apostolov*², *Jonathan Barnoud*^{3,†}, *Paul Bauer*¹, *Christian Blau*¹,
11
12 3 *Alexandre M.J.J. Bonvin*⁴, *Matthieu Chavent*^{5,*}, *John Chodera*⁶, *Karmen Čondić-Jurkić*^{6,7},
13
14 4 *Lucie Delemotte*¹, *Helmut Grubmüller*⁸, *Rebecca J. Howard*⁹, *E. Joseph Jordan*⁹, *Erik*
15
16 5 *Lindahl*⁹, *O. H. Samuli Ollila*¹⁰, *Jana Selent*¹¹, *Daniel G. A. Smith*¹², *Phillip J. Stansfeld*¹³,
17
18 6 *Johanna K.S. Tiemann*¹⁴, *Mikael Trellet*⁴, *Christopher Woods*¹⁵, *Artem Zhmurov*¹
19
20
21 7

8 AUTHOR ADDRESS

- 22
23
24 9 *1- Science for Life Laboratory, Department of Applied Physics, KTH Royal Institute of Technology, Box*
25
26 10 *1031, SE-171 21 Solna*
27
28 11 *2- PDC Center for High Performance Computing, School of Electrical Engineering and Computer*
29
30 12 *Science, KTH Royal Institute of Technology, Stockholm, Sweden*
31
32 13 *3- University of Groningen, Netherlands*
33
34 14 *4- Utrecht University, Faculty of Science, Bijvoet Center, Utrecht, the Netherlands*
35
36 15 *5- IPBS, Université Paul Sabatier, Toulouse, France*
37
38 16 *6- Computational and Systems Biology Program, Sloan Kettering Institute, Memorial Sloan Kettering*
39
40 17 *Cancer Center, New York, USA*
41
42 18 *7- Open Force Field Consortium*
43
44 19 *8- Max Planck Institute for Biophysical Chemistry, Goettingen, Germany*
45
46 20 *9- Science for Life Laboratory, Department of Biophysics and Biochemistry, Stockholm University,*
47
48 21 *Box 1031, SE-171 21 Solna*
49
50 22 *10- Institute of Biotechnology, University of Helsinki, Finland*
51
52 23 *11- Research Programme on Biomedical Informatics, Hospital del Mar Medical Research Institute*
53
54 24 *(IMIM) & Department of Experimental and Health Sciences, Pompeu Fabra University, Barcelona,*
55
56 25 *Spain*
57
58 26 *12- The Molecular Sciences Software Institute, Blacksburg, USA*
59
60 27 *13- Department of Biochemistry, University of Oxford, Oxford, UK*
28
29 28 *14- Institute of Medical Physics and Biophysics, Faculty of Medicine, University Leipzig, Leipzig 04107,*
30
31 *Germany*
32
33 *15- University of Bristol, Bristol, UK*
34
35 *†- current address: Intangible Realities Laboratory, University of Bristol, UK*
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

1
2
3 33 **KEYWORDS** Data Sharing, Open Science, Reproducibility, File Standard, Molecular
4
5 34 Simulation
6
7

8
9 35 **ABSTRACT** Given the need for modern researchers to produce open, reproducible scientific
10
11 36 output, the lack of standards and best practices for sharing data and workflows used to
12
13 37 produce and analyze molecular dynamics (MD) simulations have become an important issue
14
15 38 in the field. There are now multiple well-established packages to perform molecular dynamics
16
17 39 simulations, often highly tuned for exploiting specific classes of hardware, and each with
18
19 40 strong communities surrounding them, but with very limited interoperability/transferability
20
21 41 options. Thus, the choice of the software package often dictates the workflow for both
22
23 42 simulation production and analysis. The level of detail in documenting the workflows and
24
25 43 analysis code varies greatly in published work, hindering reproducibility of the reported results
26
27 44 and the ability for other researchers to build on these studies. An increasing number of
28
29 45 researchers are motivated to make their data available, but many challenges remain in order
30
31 46 to effectively share and reuse simulation data. To discuss these and other issues related to
32
33 47 best practices in the field in general, we organized a workshop in November 2018 (
34
35 48 <https://bioexcel.eu/events/workshop-on-sharing-data-from-molecular-simulations/>). Here, we
36
37 49 present a brief overview of this workshop and topics discussed. We hope this effort will spark
38
39 50 further conversation in the MD community to pave the way towards more open, interoperable
40
41 51 and reproducible outputs coming from research studies using MD simulations.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

53 Introduction

54 Molecular simulations have become increasingly powerful and accessible in recent
55 years, due in part to the rise of HPC¹⁻³ and GPU-powered clusters and powerful desktop
56 computers⁴ as well as the development of user-friendly software to set-up simulations^{5,6}. The
57 underlying physical models and methods have also improved over the years to address ever
58 more complex biological and chemical questions^{7,8}. Finally, the number of users and available
59 tools is continuously increasing, as is the amount and complexity of workflows and produced
60 outputs^{9,10}. In this context, defining best practices related to documentation of protocols and
61 code used to generate and/or analyze Molecular Dynamics (MD) simulations is becoming
62 more important than ever¹¹. A set of guidelines for reporting results obtained using molecular
63 dynamics techniques and an opportunity to share data, similar to what structural biologists
64 have achieved with the world-wide Protein Data Bank¹² (wwPDB), should generally help to
65 improve the quality, reproducibility, statistics, and re-use of the published results.

66 Here, we would like to focus on the term reproducibility. The definition of reproducibility
67 and its distinction from replicability can vary between disciplines¹³⁻¹⁵, but in this context, we
68 will broadly define reproducibility as the ability to re-implement the workflows of published work
69 and obtain similar behavior for observables of interest as well as define the appropriate way
70 to measure/calculate and report these observables¹⁶. Reproducibility is a long-standing issue
71 for molecular modeling¹⁷ and a key step toward better reproducibility and improved
72 collaboration is making data more accessible and workflows interoperable. This can help
73 reduce the entry barrier for the newcomers, but it could also help the existing practitioners to
74 focus on answering scientific questions rather than wasting time in redeveloping existing sets
75 of parameters or translating files formats to pass from one software to another. To reach this
76 goal, it is now necessary to overcome several difficulties:

- 77 • First, there is now a multitude of package-specific file formats and object models.

78 This variety, although increasing the efficiency for each package, introduces limitations in the

1
2
3 79 interoperability and creates friction for users juggling with various software to generate and
4
5 80 analyze their data.

6
7
8 81 • Second, there is still a lack of exhaustive documentation related to new software
9
10 82 development. The proliferation of various libraries and toolkits definitely opens up new
11
12 83 avenues of research, but documenting the entire workflow from building a molecular model
13
14 84 and parameterization to data analysis and visualization has become more complex. The
15
16 85 method sections in publications often lack sufficient details to successfully re-implement the
17
18 86 protocol or repeat the study from scratch, and default parameters to run a simulation may vary
19
20 87 from one software version to another.

21
22
23 88 • Last but not least, there is no consensus to share data. The recent years have
24
25 89 seen developments of different open data platforms, but the (ever-increasing) size of the
26
27 90 generated trajectories makes it difficult to share simulation data efficiently. The absence of
28
29 91 appropriate infrastructure, guidelines, and incentives further complicate the situation^{18,19}.

30
31
32 92
33
34
35 93 In general, we are witnessing a growing effort to make science more open by
36
37 94 researchers themselves and increasingly so by funders and journals^{20,21}. Soon, it may be
38
39 95 mandatory to share data and deposit models obtained from hybrid/integrative approaches
40
41 96 combining molecular modeling and experimental results²². Finding a way to consistently share
42
43 97 data, workflows, and protocols will be thus necessary to ensure an efficient information
44
45 98 exchange. Defining best practices and coming up with solutions should be a community effort
46
47 99 to achieve the best outcome for everyone involved. In an effort to start a discussion around
48
49 100 these questions, we organized a BioExcel workshop on *Sharing Data from Molecular*
50
51 101 *Simulations* (SDMS) in Stockholm, November 2018. In this paper, we present a summary of
52
53 102 discussions broadly focused on 4 topics:
54
55
56
57
58
59
60

- 104 • Standardization of file formats
- 105 • Streamlining molecular simulations data
- 106 • Tools for trajectory file sharing
- 107 • Reproducibility of molecular simulations

108 Each topic was introduced by 2 researchers and then openly discussed by all participants. All
109 the presentations and the discussions were recorded and are accessible here:
110 <https://bioexcel.eu/sdms18-recordings/>. The slides for the majority of the talks can be found
111 here: <https://doi.org/10.5281/zenodo.2652703>.

113 **Standardization of file formats**

114 While in structural biology the established PDB file format was stable for decades¹²,
115 the MD simulations field has a tendency to produce a multitude of input/output formats each
116 related to one MD package^{1,23-27}. With the rapid growth in complexity, size, and number of
117 macromolecular structures led by advances in experimental techniques, even the canonical
118 PDB format is now evolving to allow rendering and analyzing larger files with a gain in
119 performance²⁸. This evolution may also encourage the MD community to update its file formats
120 to deal with larger and more heterogeneous data.

121 A new jointly developed format would need to be modular and flexible enough to not
122 only take into account current but also catch future needs. Here arises a first question: What
123 are the current and future needs of the MD community for such format? While particle
124 coordinates are the current main feature both for input and output standards, other features
125 need to be discussed such as physical/chemical descriptions of the model, experimental data
126 used to create the model, technical details related to the simulation (such as algorithms used,
127 sampling method, and forcefield). Different formats may be used as templates such as
128 MMTF²⁸, MMCIF²⁹, JSON (<http://www.json.org/>), TNF³⁰. At this workshop we all agreed that it

1
2
3 129 would be a great advance if this new standard can follow the FAIR principle³¹: Findable,
4
5 130 Accessible, Interoperable, and Reproducible/Reusable. Many details remain to be discussed
6
7
8 131 and the standardization question cannot be solved in one workshop with only a small sample
9
10 132 of the MD community but need to be discussed by all main software developers joined with
11
12 133 users to ensure usability. To do so another workshop will be held soon in New York to further
13
14 134 discuss the question of file format and MD packages interoperability:
15
16
17 135 <https://molssi.org/2019/07/29/molssi-workshop-molecular-dynamics-software-interoperability/>
18

19 136 .
20
21 137 For further details and discussions interested readers can watch associated videos from the
22
23 138 2018 workshop:

- 26 139 • [Introduction of the topic](https://youtu.be/2S3qjBIE6Y4) by Mark Abraham (<https://youtu.be/2S3qjBIE6Y4>)
- 27
28 140 • [Preliminary talk I](https://youtu.be/Hvy8-gyTmj8) by Erik Lindahl (<https://youtu.be/Hvy8-gyTmj8>)
- 29
30 141 • [Preliminary talk II](https://youtu.be/48Eb2MLHoYU) from Alexandre Bonvin (<https://youtu.be/48Eb2MLHoYU>)
- 31
32 142 • [Breakout discussions](https://youtu.be/4fnV5EFXDpc) presented by Phillip Stansfeld, Mikael Trellet, Daniel Smith
33
34 143 and Johanna Tiemann (<https://youtu.be/4fnV5EFXDpc>)
- 35
36
37 144

39 145 **Streamlining molecular simulations data**

40
41
42 146 The MD simulation is often not a means and an end in itself but instead is run as part of a
43
44 147 larger workflow. Such workflows involve joining together the output of many independent
45
46 148 programs, such as those used for parameterizing molecules, those for performing molecular
47
48 149 dynamics, and those for trajectory analysis. Managing the data movement between different
49
50 150 programs in this workflow is challenging for several reasons:

- 51
52
53 151 1. The file formats used by different programs in the workflow may be incompatible,
54
55 152 thereby preventing certain combinations of tools from being used together.
- 56
57
58
59
60

- 1
2
3 153 2. The features and forcefields supported by different programs in the workflow may be
4
5 154 incompatible, thereby forcing researchers to choose algorithms and forcefields based
6
7
8 155 on software compatibility rather than for good scientific reasons.
9
10 156 3. Different programs may implement features or forcefields in different ways, thereby
11
12 157 meaning that the results of running the workflow will depend on the exact combination
13
14 158 of programs (and possibly program versions) used. It is generally not possible to mix-
15
16
17 159 and-match different programs and get the same results.
18
19 160

21 161 These challenges have forced researchers to develop workflows using specific
22
23 162 software packages and specific forcefields. This creates divisions within the community and
24
25
26 163 makes it difficult to write workflows that function equally well across a number of forcefields
27
28 164 and a number of different software packages.

30 165 One of the solutions to this problem is the development of programs that
31
32 166 convert/handle molecular information between the different file formats such as VMD³²,
33
34 167 cpptraj³³, MDAnalysis^{34,35}, mdtraj³⁶, LOOS^{37,38} and many others for trajectory analysis and
35
36
37 168 TopoGromacs³⁹, CHARMM-GUI⁴⁰, CHAMBER⁴¹, ParmEd
38
39 169 (<http://parmed.github.io/ParmEd/html/index.html#>), InterMol⁴²
40
41 (https://github.com/shirtsgroup/InterMol), and others for topology generation and editing. The
42
43 170 aim of these programs is to translate as much information as possible from one molecular file
44
45 171 format into another. One recent example is BioSimSpace (<https://biosimspace.org/>), which
46
47 172 provides wrappers that simplify the generation of the command files that are used to control
48
49 173 the running of simulations. This allows researchers to write workflows that are independent of
50
51 174 the choice of the underlying packages used to perform the simulation. BioSimSpace aims to
52
53 175 run all stages of the workflow using the simulation software installed on the researcher's
54
55 176 computer that is compatible with the forcefield chosen for the specific calculation.
56
57
58
59
60

1
2
3 178 While translators and program wrappers like ParmEd and BioSimSpace solve some
4
5 179 of these problems, they are not a universal solution. They do not solve the issue that different
6
7
8 180 simulation programs use different algorithms (or interpretations of algorithms, for example,
9
10 181 different implementations of thermostats or integrators), or that different programs store and
11
12 182 represent molecular information in different ways (e.g. SHAKE information for constraining
13
14 183 bonds is represented in the molecular topology in GROMACS, while it is a simulation
15
16 184 command parameter in NAMD and AMBER). This means MD properties/observables
17
18 185 computed with one package will be systematically different by an often small but statistically
19
20 186 significant amount from those computed with a different package as shown for free energy
21
22 187 calculations⁴³. Thus, the version and name of the MD program used to produce a simulation
23
24 188 result will affect that result, and must be reported accordingly. Furthermore, MD simulations
25
26 189 outputs are mainly trajectories which (1) represent ensemble averages (2) are chaotic in that
27
28 190 small differences in initial conditions cause large differences in the subsequent dynamics
29
30 ('butterfly effect'). This adds another layer of complexity and needs also a consensus on how
31
32 191 to further analyze/process these trajectories to provide the final quantities of interest.
33
34
35 192

36
37 193 The recordings of this session can be found here:

- 38
39 194
- 40 • [Introduction to the topic](https://youtu.be/6xOfN0y_uoQ) by John Chodera (https://youtu.be/6xOfN0y_uoQ)
 - 41
 - 42 • [Preliminary talk I](https://youtu.be/YPYeujSD-6Y) by Philip Stansfeld (<https://youtu.be/YPYeujSD-6Y>)
 - 43
 - 44 • [Preliminary talk II](https://youtu.be/w1d1xtbGhHc) by Christopher Woods (<https://youtu.be/w1d1xtbGhHc>)
 - 45
 - 46 • [Breakout discussions](#) presented by Christian Blau, Christopher Woods, Jonathan
47
48 Barnoud and Mark Abraham (<https://youtu.be/Z-JfBU3Emug>)
- 49 198
50
51 199

52 53 200 **Tools for trajectory file sharing**

54
55 201 The benefits of sharing data together with the peer-reviewed publication, preprint or as a self-
56
57 202 standing research output seem to be many - from receiving additional credit for one's work to
58
59
60

1
2
3 203 improving reproducibility, reusability or offering potentially new avenues of research^{20,44}. Some
4
5 204 disciplines, such as protein crystallography or genomics, have open data practices well
6
7 205 integrated into their workflow, with metadata being collected throughout the workflow, and
8
9 206 those practices are a *de facto* standard in scholarly communication. However, data sharing in
10
11 207 the MD community still has not become widely adopted because best practice guidelines or
12
13 208 journal recommendations on how to share MD simulations are yet to be established and
14
15 209 adopted by the whole community. Making data sharing a standard practice in the field faces
16
17 210 both technical and cultural challenges, although these are currently being tackled by some
18
19 211 ongoing initiatives and solutions^{20,45,46}. Thus, the development of best practices and guidelines
20
21 212 for simulation data sharing will be of tremendous value, especially if created with the FAIR
22
23 213 principles in mind³¹. To do so, we need to address several important questions regarding *what*
24
25 214 *data* should be shared, *how* and *where*.

26
27 215 Answering to the *what data* question would need longer discussions not limited to a
28
29 216 small group of individuals but involving the whole community and especially all the MD
30
31 217 packages (another workshop will be held soon to help starting to answer to this question:
32
33 218 <https://molssi.org/2019/07/29/molssi-workshop-molecular-dynamics-software-interoperability/>
34
35 219). The emergence of dedicated tools is now helping to answer to the *how* question. Software
36
37 220 such as MDsrv⁴⁷, HTMoL⁴⁸, Mol* (<https://molstar.org>), Molmil⁴⁹ are now taking advantage of
38
39 221 the WebGL API for sharing trajectories through interactive visualization on the web⁵⁰.

40
41 222 Other fields of research can help us to answer to the *where* question. Existing
42
43 223 databanks, such as wwPDB⁵¹ and Galaxy (<https://usegalaxy.org>), have been recognized by
44
45 224 the scientific community. However, the establishment of an analogous, specialized platform
46
47 225 for MD data, poses a great challenge, given the current lack of long-term support for the
48
49 226 infrastructure projects of this kind. It is not clear yet who should be responsible for building
50
51 227 such platform and how this infrastructure could be funded in a sustainable way, preferably
52
53 228 without relying on short-term research grants, to cover the costs of development, maintenance
54
55
56
57
58
59
60

1
2
3 229 and data hosting. In the meantime, community-driven, special-purpose platforms like the
4
5 230 GPCRmd (<http://www.gpcrmd.org>), Lipidbook⁵² and NMRlipids⁴⁵
6
7 231 (<http://nmrlipids.blogspot.com>), Ligandbook⁵³, MoDEL⁵⁴ and BIGNASim⁵⁵ lead the way,
8
9 232 providing specialized platforms for deposition and analysis of G protein-coupled receptors
10
11 233 (GPCR), lipids, small molecules, proteins, and nucleic acids, respectively. General data
12
13 234 sharing resources like Zenodo (<https://zenodo.org>), FigShare (<https://figshare.com>), Open
14
15 235 Science Framework (<https://osf.io>) and others, also provide an opportunity for every
16
17 236 researcher to deposit their simulation files and trajectories. Nevertheless, those resources may
18
19 237 not provide sometimes enough space to sustainably store MD simulations outputs (with file
20
21 238 size limits ranging between 5 GB and 50 GB).

22
23
24 239 To establish an efficient sharing culture, a systematic approach to developing tools
25
26 240 and sharing guidelines is necessary, with the participation of the entire community in such
27
28 241 activities and efforts. An open and inclusive discussion about best practices in data sharing,
29
30 242 identification of short-term solutions based on the currently available frameworks and tools,
31
32 243 as well as developing a strategy and requirements for future solutions bespoke to MD
33
34 244 community and their needs is necessary. More details about the discussions taking place at
35
36 245 the workshop can be found in the following videos:

- 37
38
39 246
 - [Introduction to the topic](https://youtu.be/mvesL9Y_9xU) by Daniel Smith (https://youtu.be/mvesL9Y_9xU)
 - [Preliminary talk I](https://youtu.be/VOT6fEc7Iuc) by Johanna Tiemann (<https://youtu.be/VOT6fEc7Iuc>)
 - [Preliminary talk II](https://youtu.be/TVS75j48mQ8) by Jana Selent (<https://youtu.be/TVS75j48mQ8>)
 - [Breakout discussions](https://youtu.be/Uls1isntUPY) presented by John Chodera, Karmen Čondić-Jurkić, Samuli
44
45 249 Ollila and Lucie Delemotte (<https://youtu.be/Uls1isntUPY>)

46
47 250
48
49 251

252 **Reproducibility of molecular simulations**

253 MD simulations are chaotic and as such, the definition of reproducible results is non-
254 trivial. First, the distinction between repeatability (by the same team and the same
255 computational setup), replicability (by a different team and the same computational setup) and

1
2
3 256 reproducibility (by a different team, and with a different experimental setup) should be made
4
5 257 ¹⁴. Differences in outputs from these three perspectives may indicate different types of errors
6
7
8 258 (bugs in software, human errors, or different choices along the workflow - choice of code, force
9
10 259 field, system setup and more). The variability of parameters and dependence of the final
11
12 260 results on both software and hardware makes it complicated (but also often unnecessary) to
13
14 261 achieve the exact replication/repetition of any given setup, and untangling all the effects would
15
16
17 262 be a difficult task. Focusing on a set of observables that can be calculated and preferably
18
19 263 validated against experiments might be a better way of approaching reproducibility in this
20
21 264 particular field. Similarly, focusing at observables which, despite the underlying chaoticity of
22
23 265 the detailed dynamics, are reproducible without too large variation might be beneficial.
24
25
26 266 Reaching an agreement on which observables we should aim to reproduce and how to
27
28 267 properly calculate and report these values is thus desirable. For this, educational efforts are
29
30 268 needed: best practice dissemination in terms of calculating statistical properties, for example,
31
32
33 269 are crucial¹⁶. Coming up with standard benchmarks would also help, where the performance
34
35 270 of different software/forcefield combinations for selected tasks could be compared.

36
37 271 In practice, data sharing would help with replicability and reproducibility. Practical
38
39
40 272 challenges come from the size of data sets. However, one can envision sharing at least
41
42 273 minimal data sets to improve

43
44 274 ● methods reproducibility: provide sufficient details to replicate the study; this is
45
46
47 275 in principle already done in publications, but authors, reviewers, and editors
48
49 276 should pay special attention to the question, and sharing directly all input files
50
51 277 should be mandatory,

52
53 278 ● raw data reproducibility: share a minimum amount of data in the form of MD
54
55
56 279 simulation snapshots, or even better whole trajectories, on existing data
57
58 280 sharing repositories - Zenodo, Figshare, OSF, and

281 • results and inferential reproducibility: share among other analysis code,
282 pipeline/workflow and example used.

283 Inspiration can be found in other research fields (e.g. Genomics⁵⁶ or Proteomics⁵⁷) and existing
284 dedicated initiatives, like MemProtMD⁵⁸ (<http://memprotmd.bioch.ox.ac.uk>), NMRlipids project
285 (www.nmr lipids.blogspot.fi) and GPCRmd (<http://www.gpcrmd.org>), show that small groups of
286 people focused on a narrow topic can create the necessary structure to share even large
287 datasets in an efficient way. For further details and discussions interested readers can watch
288 associated videos:

- 289 • [Introduction to the topic](#) by Karmen Čondić-Jurkić
290 (<https://youtu.be/IUTQgOXDEP8>)
- 291 • [Preliminary talk I](#) by Helmut Grubmüller (<https://youtu.be/cliVmGlrKag>)
- 292 • [Preliminary talk II](#) by Samuli Ollila (<https://youtu.be/46s33SonsiU>)
- 293 • [Breakout discussions](#) presented by Mikael Trellet, Alexandre Bonvin, Mark
294 Abraham and Christopher Woods (https://youtu.be/ex0_bqmJwE8)

296 This article summarizes the discussions started during the workshop held in Stockholm
297 in November 2018. As may be noted by the reader, these discussions have not solved the
298 issues about sharing data that our field is facing. Of course, this has never been the goal of
299 such a small workshop. This workshop was intended to start asking relevant questions. Thus,
300 this document (and the videos associated) can be seen as a road map for future
301 developments. It is now crucial to build a community responsible for transforming these ideas
302 into actions. This community needs to represent a diversity of perspectives by including both
303 MD users and developers, newcomers and more seasoned practitioners, PhD students and
304 postdocs, who are performing MD simulations on a daily basis, and PIs, who may hold the
305 bigger picture views. As a community building effort, we are planning to regularly organize

1
2
3 306 more specific workshops aiming to address some of the issues raised in this article or to
4
5 307 expand the scope of newly recognized problems. Of course, the structure of the workshops
6
7 308 limits the number of participants, but care will be taken to ensure the aforementioned diversity
8
9
10 309 of perspectives and roles in the field. In an effort to include as many users as possible in this
11
12 310 discussion, the best practices guidelines that will emerge from these workshops will be
13
14 311 submitted to the Living Journal of Computational Molecular Science
15
16 312 (<http://www.livecomsjournal.org/>). This journal "... *provides a venue where authors can submit*
17
18 313 *living documents that are updated on an ongoing basis as websites or Wikipedia articles could*
19
20 314 *be, but which still have clear authorship and provide a mechanism for authors to get publication*
21
22 315 *credit for their work.*"⁵⁹ Hence, researchers interested to help us shape new practices to share
23
24 316 data will be able to provide their feedback or directly contribute to the forthcoming document
25
26 317 (as per the general idea laid out here: https://livecomsjournal.github.io/about/paper_code/).
27
28 318 We hope that our work will act as a first step in a community-driven process of defining best
29
30 319 practices for tool development and application in the molecular dynamics field.
31
32
33
34
35
36
37
38
39
40
41
42

43 320

44 321

45 322 AUTHOR INFORMATION

46 323 **Corresponding Author**

47 324 *Correspondence: matthieu.chavent@ipbs.fr , @Matth_Chavent

48 325

49 326 **Author Contributions**

50
51 327 The manuscript was written through contributions of all authors. All authors have given
52
53 328 approval to the final version of the manuscript.
54
55
56
57
58
59 329
60

1
2
3 330 **Funding Sources**
4

5
6 331 The workshop was supported by BioExcel Centre of Excellence (www.bioexcel.eu).
7
8

9 332 **Acknowledgement**
10

11 333 This work was supported by BioExcel Centre of Excellence (www.bioexcel.eu) funded by the
12

13 334 European Union contracts H2020-INFRAEDI-02-2018-823830 and H2020-EINFRA-2015-1-
14

15 335 675728. MC acknowledges support from CNRS-MITI grants PEPS MPI 2018 and
16

17 336 “Modélisation du vivant” 2019. This work was supported by grants from the Gustafsson
18

19 337 Foundation and Science for Life Laboratory to LD. HG has been supported by Max Planck
20

21 338 Society and the German Research Foundation (DFG), Cluster of excellence Multiscale
22

23 339 Imaging and the DFG priority programmes 1648, 755, and 803. OHSO acknowledges financial
24

25 340 support from Academy of Finland (315596). DGAS thanks the National Science Foundation
26

27 341 for support under Grant No. ACI-1547580. PJS would like to thank Wellcome [208361/Z/17/Z],
28

29 342 the BBSRC [BB/P01948X/1, BB/R002517/1, BB/S003339/1 and BB/I019855/1] and MRC
30

31 343 [MR/S009213/1] for funding. JKST acknowledges support from the Deutsche
32

33 344 Forschungsgemeinschaft (DFG) HI1502/1-2. CW acknowledges support from the EPSRC
34

35 345 (EP/N018591/1). We thank Oliver Beckstein and David Mobley for their careful reading and
36

37 346 their comments.
38
39
40
41
42
43

44 347 **Link to the SDMS18 recordings:** <https://bioexcel.eu/sdms18-recordings/>

45 348 **Discussions from Twitter:** can be retrieved/extended by using the hashtag #SDMS18

46 349 **Several participants from this workshop can be contacted/followed on Twitter:**

47 350 @the_mabraham, @jbarnoud, @amjjbonvin, @Matth_Chavent, @jchodera, @karmecon,

48 351 @DelemotteLab, @CompBioPhys, @eriklindahl, @NMRlipids, @dga_smith, @pstansfeld,

49 352 @jOkaso, @chryswoods
50
51
52
53
54
55
56
57
58
59
60

353
354

355 REFERENCES

- 356 (1) Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E.
357 GROMACS: High Performance Molecular Simulations Through Multi-Level
358 Parallelism From Laptops to Supercomputers. *SoftwareX* **2015**, 1-2, 19–25.
- 359 (2) Lagardère, L.; Jolly, L.-H.; Lipparini, F.; Aviat, F.; Stamm, B.; Jing, Z. F.; Harger, M.;
360 Torabifard, H.; Cisneros, G. A.; Schnieders, M. J.; Gresh, N.; Maday, Y.; Ren, P. Y.;
361 Ponder, J. W.; Piquemal, J.-P. Tinker-HP: a Massively Parallel Molecular Dynamics
362 Package for Multiscale Simulations of Large Complex Systems with Advanced Point
363 Dipole Polarizable Force Fields. *Chem. Sci.* **2018**, 9, 956–972.
- 364 (3) Jung, J.; Nishima, W.; Daniels, M.; Bascom, G.; Kobayashi, C.; Adedoyin, A.; Wall,
365 M.; Lappala, A.; Phillips, D.; Fischer, W.; Tung, C. S.; Schlick, T.; Sugita, Y.;
366 Sanbonmatsu, K. Y. Scaling Molecular Dynamics Beyond 100,000 Processor Cores
367 for Large-Scale Biophysical Simulations. *J Comput Chem* **2019**.
- 368 (4) Stone, J. E.; Hardy, D. J.; Ufimtsev, I. S.; Schulten, K. GPU-Accelerated Molecular
369 Modeling Coming of Age. *J. Mol. Graph. Model.* **2010**, 29, 116–125.
- 370 (5) Doerr, S.; Harvey, M. J.; Noé, F.; De Fabritiis, G. HTMD: High-Throughput Molecular
371 Dynamics for Molecular Discovery. *J. Chem. Theory Comput.* **2016**, 12, 1845–1852.
- 372 (6) Jo, S.; Cheng, X.; Lee, J.; Kim, S.; Park, S. J.; Patel, D. S.; Beaven, A. H.; Lee, K. I.;
373 Rui, H.; Park, S.; Lee, H. S.; Roux, B.; MacKerell, A. D.; Klauda, J. B.; Qi, Y.; Im, W.
374 CHARMM-GUI 10 Years for Biomolecular Modeling and Simulation. *J Comput Chem*
375 **2017**, 38, 1114–1124.
- 376 (7) Marrink, S. J.; Corradi, V.; Souza, P. C. T.; Ingólfsson, H. I.; Tieleman, D. P.; Sansom,
377 M. S. P. Computational Modeling of Realistic Cell Membranes. *Chem. Rev.* **2019**,
378 *acs.chemrev.8b00460*.
- 379 (8) Bottaro, S.; Lindorff-Larsen, K. Biophysical Experiments and Biomolecular
380 Simulations: a Perfect Match? *Science* **2018**, 361, 355–360.
- 381 (9) Im, W.; Liang, J.; Olson, A.; Zhou, H.-X.; Vajda, S.; Vakser, I. A. Challenges in
382 Structural Approaches to Cell Modeling. *J. Mol. Biol.* **2016**, 428, 2943–2964.
- 383 (10) Chavent, M.; Duncan, A. L.; Sansom, M. S. Molecular Dynamics Simulations of
384 Membrane Proteins and Their Interactions: From Nanoscale to Mesoscale. *Curr.*
385 *Opin. Struct. Biol.* **2016**, 40, 8–16.
- 386 (11) Elofsson, A.; Hess, B.; Lindahl, E.; Onufriev, A.; Van Der Spoel, D.; Wallqvist, A. Ten
387 Simple Rules on How to Create Open Access and Reproducible Molecular
388 Simulations of Biological Systems. *PLoS Comput Biol* **2019**, 15, e1006649.
- 389 (12) Berman, H. M.; Kleywegt, G. J.; Nakamura, H.; Markley, J. L. The Protein Data Bank
390 at 40: Reflecting on the Past to Prepare for the Future. *Structure (London, England :
391 1993)* **2012**, 20, 391–396.
- 392 (13) Plesser, H. E. Reproducibility vs. Replicability: a Brief History of a Confused
393 Terminology. *Front Neuroinform* **2017**, 11, 76.
- 394 (14) Hinsén, K. ActivePapers: a Platform for Publishing and Archiving Computer-
395 Aided Research. *F1000Res* **2014**, 3, 289.
- 396 (15) Barba, L. A. Terminologies for Reproducible Research. *CoRR* **2018**.
- 397 (16) Grossfield, A.; Patrone, P. N.; Roe, D. R.; Schultz, A. J.; Siderius, D. W.; Zuckerman,
398 D. M. Best Practices for Quantification of Uncertainty and Sampling Quality in
399 Molecular Simulations [Article v1.0]. *Living Journal of Computational Molecular
400 Science* **2018**, 1.
- 401 (17) Walters, W. P. Modeling, Informatics, and the Quest for Reproducibility. *J Chem Inf
402 Model* **2013**, 53, 1529–1530.
- 403 (18) Graham, S. C.; Nagar, B.; Privé, G. G.; Deane, J. E. Molecular Models Should Not Be
404 Published Without the Corresponding Atomic Coordinates. *Proc Natl Acad Sci USA*
405 **2019**, 116, 11099–11100.
- 406 (19) Romero, R.; Yuen, T.; New, M. I.; Zaidi, M.; Haider, S. Reply to Graham Et Al.: in
407 Silico Atomistic Coordinates and Molecular Dynamics Simulation Trajectories of the

- 1
2
3 408 Glucocerebrosidase-Saposin C Complex. *Proc Natl Acad Sci USA* **2019**, *116*, 11101–
4 409 11102.
- 5 410 (20) Data Sharing and the Future of Science. *Nat Commun* **2018**, *9*.
- 6 411 (21) Introducing eLife's First Computationally Reproducible Article. **2019**.
- 7 412 (22) Burley, S. K.; Kurisu, G.; Markley, J. L.; Nakamura, H.; Velankar, S.; Berman, H. M.;
8 413 Sali, A.; Schwede, T.; TrewHELLa, J. PDB-Dev: a Prototype System for Depositing
9 414 Integrative/Hybrid Structural Models. *Structure (London, England : 1993)* **2017**, *25*,
10 415 1317–1318.
- 11 416 (23) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.;
12 417 Skeel, R. D.; Kalé, L.; Schulten, K. Scalable Molecular Dynamics with NAMD. *J*
13 418 *Comput Chem* **2005**, *26*, 1781–1802.
- 14 419 (24) Brooks, B. R.; Brooks, C. L.; Mackerell, A. D.; Nilsson, L.; Petrella, R. J.; Roux, B.;
15 420 Won, Y.; Archontis, G.; Bartels, C.; Boresch, S.; Caflisch, A.; Caves, L.; Cui, Q.;
16 421 Dinner, A. R.; Feig, M.; Fischer, S.; Gao, J.; Hodoscek, M.; Im, W.; Kuczera, K.;
17 422 Lazaridis, T.; Ma, J.; Ovchinnikov, V.; Paci, E.; Pastor, R. W.; Post, C. B.; Pu, J. Z.;
18 423 Schaefer, M.; Tidor, B.; Venable, R. M.; Woodcock, H. L.; Wu, X.; Yang, W.; York, D.
19 424 M.; Karplus, M. CHARMM: the Biomolecular Simulation Program. *J Comput Chem*
20 425 **2009**, *30*, 1545–1614.
- 21 426 (25) Salomon Ferrer, R.; Case, D. A.; Walker, R. C. An Overview of the Amber
22 427 Biomolecular Simulation Package. *Wiley Interdisciplinary Reviews: Computational*
23 428 *Molecular Science* **2012**, *3*, 198–210.
- 24 429 (26) Rackers, J. A.; Wang, Z.; Lu, C.; Laury, M. L.; Lagardère, L.; Schnieders, M. J.;
25 430 Piquemal, J.-P.; Ren, P.; Ponder, J. W. Tinker 8: Software Tools for Molecular Design.
26 431 *J. Chem. Theory Comput.* **2018**, *14*, 5273–5289.
- 27 432 (27) Eastman, P.; Friedrichs, M. S.; Chodera, J. D.; Radmer, R. J.; Bruns, C. M.; Ku, J. P.;
28 433 Beauchamp, K. A.; Lane, T. J.; Wang, L.-P.; Shukla, D.; Tye, T.; Houston, M.; Stich,
29 434 T.; Klein, C.; Shirts, M. R.; Pande, V. S. OpenMM 4: a Reusable, Extensible, Hardware
30 435 Independent Library for High Performance Molecular Simulation. *J. Chem. Theory*
31 436 *Comput.* **2013**, *9*, 461–469.
- 32 437 (28) Bradley, A. R.; Rose, A. S.; Pavelka, A.; Valasatava, Y.; Duarte, J. M.; Prlić, A.; Rose,
33 438 P. W. MMTF—an Efficient File Format for the Transmission, Visualization, and
34 439 Analysis of Macromolecular Structures. *PLoS Comput Biol* **2017**, *13*, e1005575.
- 35 440 (29) Bourne, P. E.; Berman, H. M.; McMahon, B.; Watenpaugh, K. D.; Westbrook, J. D.;
36 441 Fitzgerald, P. Macromolecular Crystallographic Information File. *Meth. Enzymol.*
37 442 **1997**, *277*, 571–590.
- 38 443 (30) Lundborg, M.; Apostolov, R.; Spangberg, D.; Gardenas, A.; Van Der Spoel, D.;
39 444 Lindahl, E. An Efficient and Extensible Format, Library, and API for Binary Trajectory
40 445 Data From Molecular Simulations. *J Comput Chem* **2014**, *35*, 260–269.
- 41 446 (31) Wilkinson, M. D.; Dumontier, M.; Aalbersberg, I. J. J.; Appleton, G.; Axton, M.; Baak,
42 447 A.; Blomberg, N.; Boiten, J.-W.; da Silva Santos, L. B.; Bourne, P. E.; Bouwman, J.;
43 448 Brookes, A. J.; Clark, T.; Crosas, M.; Dillo, I.; Dumon, O.; Edmunds, S.; Evelo, C. T.;
44 449 Finkers, R.; Gonzalez-Beltran, A.; Gray, A. J. G.; Groth, P.; Goble, C.; Grethe, J. S.;
45 450 Heringa, J.; 't Hoen, P. A. C.; Hooff, R.; Kuhn, T.; Kok, R.; Kok, J.; Lusher, S. J.;
46 451 Martone, M. E.; Mons, A.; Packer, A. L.; Persson, B.; Rocca-Serra, P.; Roos, M.; van
47 452 Schaik, R.; Sansone, S.-A.; Schultes, E.; Sengstag, T.; Slater, T.; Strawn, G.; Swertz,
48 453 M. A.; Thompson, M.; van der Lei, J.; van Mulligen, E.; Velterop, J.; Waagmeester,
49 454 A.; Wittenburg, P.; Wolstencroft, K.; Zhao, J.; Mons, B. The FAIR Guiding Principles
50 455 for Scientific Data Management and Stewardship. *Sci Data* **2016**, *3*, 160018.
- 51 456 (32) Humphrey, W.; Dalke, A.; Schulten, K. VMD: Visual Molecular Dynamics. *J Mol Graph*
52 457 **1996**, *14*, 33–38, 27–28.
- 53 458 (33) Roe, D. R.; Cheatham, T. E. PTRAJ and CPPTRAJ: Software for Processing and
54 459 Analysis of Molecular Dynamics Trajectory Data. *J. Chem. Theory Comput.* **2013**, *9*,
55 460 3084–3095.
- 56 461 (34) Michaud-Agrawal, N.; Denning, E. J.; Woolf, T. B.; Beckstein, O. MDAAnalysis: a
57 462 Toolkit for the Analysis of Molecular Dynamics Simulations. *J Comput Chem* **2011**,

- 1
2
3 463 32, 2319–2327.
- 4 464 (35) Gowers, R.; Linke, M.; Barnoud, J.; Reddy, T.; Melo, M.; Seyler, S.; Domański, J.;
5 465 Dotson, D.; Buchoux, S.; Kenney, I.; Beckstein, O. MDAnalysis: a Python Package for
6 466 the Rapid Analysis of Molecular Dynamics Simulations; SciPy, 2016; pp 98–105.
- 7 467 (36) McGibbon, R. T.; Beauchamp, K. A.; Harrigan, M. P.; Klein, C.; Swails, J. M.;
8 468 Hernández, C. X.; Schwantes, C. R.; Wang, L.-P.; Lane, T. J.; Pande, V. S. MDTraj:
9 469 a Modern Open Library for the Analysis of Molecular Dynamics Trajectories. *Biophys.*
10 470 *J.* **2015**, *109*, 1528–1532.
- 11 471 (37) Romo, T. D.; Grossfield, A. LOOS: an Extensible Platform for the Structural Analysis
12 472 of Simulations. *Conf Proc IEEE Eng Med Biol Soc* **2009**, *2009*, 2332–2335.
- 13 473 (38) Romo, T. D.; Leioatts, N.; Grossfield, A. Lightweight Object Oriented Structure
14 474 Analysis: Tools for Building Tools to Analyze Molecular Dynamics Simulations. *J*
15 475 *Comput Chem* **2014**, *35*, 2305–2318.
- 16 476 (39) Vermaas, J. V.; Hardy, D. J.; Stone, J. E.; Tajkhorshid, E.; Kohlmeyer, A.
17 477 TopoGromacs: Automated Topology Conversion From CHARMM to GROMACS
18 478 Within VMD. *J Chem Inf Model* **2016**, *56*, 1112–1116.
- 19 479 (40) Lee, J.; Cheng, X.; Swails, J. M.; Yeom, M. S.; Eastman, P. K.; Lemkul, J. A.; Wei, S.;
20 480 Buckner, J.; Jeong, J. C.; Qi, Y.; Jo, S.; Pande, V. S.; Case, D. A.; Brooks, C. L., III;
21 481 MacKerell, A. D., Jr.; Klauda, J. B.; Im, W. CHARMM-GUI Input Generator for NAMD,
22 482 GROMACS, AMBER, OpenMM, and CHARMM/OpenMM Simulations Using the
23 483 CHARMM36 Additive Force Field. *J. Chem. Theory Comput.* **2015**, *12*, 405–413.
- 24 484 (41) Crowley, M. F.; Williamson, M. J.; Walker, R. C. CHAMBER: Comprehensive Support
25 485 for CHARMM Force Fields Within the AMBER Software. *International Journal of*
26 486 *Quantum Chemistry* **2009**, *109*, 3767–3772.
- 27 487 (42) Shirts, M. R.; Klein, C.; Swails, J. M.; Yin, J.; Gilson, M. K.; Mobley, D. L.; Case, D.
28 488 A.; Zhong, E. D. Lessons Learned From Comparing Molecular Dynamics Engines on
29 489 the SAMPL5 Dataset. *J. Comput. Aided Mol. Des.* **2016**, *31*, 147–161.
- 30 490 (43) Loeffler, H. H.; Bosisio, S.; Duarte Ramos Matos, G.; Suh, D.; Roux, B.; Mobley, D.
31 491 L.; Michel, J. Reproducibility of Free Energy Calculations Across Different Molecular
32 492 Simulation Software Packages. *J. Chem. Theory Comput.* **2018**, *14* (11), 5567–5582.
- 33 493 (44) Woelfle, M.; Olliaro, P.; Todd, M. H. Open Science Is a Research Accelerator. *Nature*
34 494 *Chem* **2011**, *3*, 745–748.
- 35 495 (45) Botan, A.; Favela-Rosales, F.; Fuchs, P. F. J.; Javanainen, M.; Kanduč, M.; Kulig, W.;
36 496 Lamberg, A.; Loison, C.; Lyubartsev, A.; Miettinen, M. S.; Monticelli, L.; Määttä, J.;
37 497 Ollila, O. H. S.; Retegan, M.; Róg, T.; Santuz, H.; Tynkkynen, J. Toward Atomistic
38 498 Resolution Structure of Phosphatidylcholine Headgroup and Glycerol Backbone at
39 499 Different Ambient Conditions. *J Phys Chem B* **2015**, *119*, 15075–15088.
- 40 500 (46) The PLUMED consortium. Promoting Transparency and Reproducibility in Enhanced
41 501 Molecular Simulations. *Nat. Methods* **2019**, *16*, 670–673.
- 42 502 (47) Tiemann, J. K. S.; Guixà-González, R.; Hildebrand, P. W.; Rose, A. S. MDsrv: Viewing
43 503 and Sharing Molecular Dynamics Simulations on the Web. *Nat. Methods* **2017**, *14*,
44 504 1123–1124.
- 45 505 (48) Carrillo-Tripp, M.; Alvarez-Rivera, L.; Lara-Ramírez, O. I.; Becerra-Toledo, F. J.;
46 506 Vega-Ramírez, A.; Quijas-Valades, E.; González-Zavala, E.; González-Vázquez, J.
47 507 C.; García-Vieyra, J.; Santoyo-Rivera, N. B.; Chapa-Vergara, S. V.; Meneses-Viveros,
48 508 A. HTMoL: Full-Stack Solution for Remote Access, Visualization, and Analysis of
49 509 Molecular Dynamics Trajectory Data. *J. Comput. Aided Mol. Des.* **2018**, *32*, 869–876.
- 50 510 (49) Bekker, G.-J.; Nakamura, H.; Kinjo, A. R. Molmil: a Molecular Viewer for the PDB and
51 511 Beyond. *J Cheminform* **2016**, *8*, 42.
- 52 512 (50) Hildebrand, P. W.; Rose, A. S.; Tiemann, J. K. S. Bringing Molecular Dynamics
53 513 Simulation Data Into View. *Trends Biochem. Sci.* **2019**.
- 54 514 (51) Berman, H.; Henrick, K.; Nakamura, H. Announcing the Worldwide Protein Data Bank.
55 515 *Nat. Struct. Biol.* **2003**, *10*, 980.
- 56 516 (52) Domański, J.; Stansfeld, P. J.; Sansom, M. S. P.; Beckstein, O. Lipidbook: a Public
57 517 Repository for Force-Field Parameters Used in Membrane Simulations. *J. Membr.*

- 1
2
3 518 *Biol.* **2010**, 236, 255–258.
- 4 519 (53) Domański, J.; Beckstein, O.; Iorga, B. I. Ligandbook — an Online Repository for Small
5 520 and Drug-Like Molecule Force Field Parameters. *Bioinformatics* **2017**, btx037.
- 6 521 (54) Meyer, T.; D'Abramo, M.; Hospital, A.; Rueda, M.; Ferrer-Costa, C.; Pérez, A.; Carrillo,
7 522 O.; Camps, J.; Fenollosa, C.; Repchevsky, D.; Gelpí, J. L.; Orozco, M. MoDEL
8 523 (Molecular Dynamics Extended Library): a Database of Atomistic Molecular Dynamics
9 524 Trajectories. *Structure (London, England : 1993)* **2010**, 18, 1399–1409.
- 10 525 (55) Hospital, A.; Andrio, P.; Cugnasco, C.; Codo, L.; Becerra, Y.; Dans, P. D.; Battistini,
11 526 F.; Torres, J.; Goñi, R.; Orozco, M.; Gelpí, J. L. BIGNASim: a NoSQL Database
12 527 Structure and Analysis Portal for Nucleic Acids Simulation Data. *Nucleic Acids Res.*
13 528 **2016**, 44, D272–D278.
- 14 529 (56) Kaye, J.; Heeney, C.; Hawkins, N.; de Vries, J.; Boddington, P. Data Sharing in
15 530 Genomics — Re-Shaping Scientific Practice. *Nature Reviews Genetics* 2009 10:5
16 531 **2009**, 10, 331–335.
- 17 532 (57) Martens, L.; Vizcaíno, J. A. A Golden Age for Working with Public Proteomics Data.
18 533 *Trends Biochem. Sci.* **2017**, 42, 333–341.
- 19 534 (58) Stansfeld, P. J.; Goose, J. E.; Caffrey, M.; Carpenter, E. P.; Parker, J. L.; Newstead,
20 535 S.; Sansom, M. S. P. MemProtMD: Automated Insertion of Membrane Protein
21 536 Structures Into Explicit Lipid Membranes. *Structure (London, England : 1993)* **2015**,
22 537 23, 1350–1361.
- 23 538 (59) Mobley, D. L.; Shirts, M. R.; Zuckerman, D. M. Why We Need the Living Journal of
24 539 Computational Molecular Science. *Living Journal of Computational Molecular Science*
25 540 **2017**, 1, 2031.
- 26 541
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

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

