



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Super-resolution fight club

Citation for published version:

Sage, D, Pham, T-A, Babcock, H, Lukes, T, Pengo, T, Chao, J, Velmurugan, R, Herbert, A, Agrawal, A, Colabrese, S, Wheeler, A, Archetti, A, Rieger, B, Ober, R, Hagen, GM, Sibarita, J-B, Ries, J, Henriques, R, Unser, M & Holden, S 2019, 'Super-resolution fight club: assessment of 2D and 3D single-molecule localization microscopy software', *Nature Methods*, vol. 16, no. 5, pp. 387-395.
<https://doi.org/10.1038/s41592-019-0364-4>

Digital Object Identifier (DOI):

[10.1038/s41592-019-0364-4](https://doi.org/10.1038/s41592-019-0364-4)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Methods

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 Super-resolution fight club: Assessment of 2D & 3D single- 2 molecule localization microscopy software

3 *Daniel Sage*^{*+1}, *Thanh-An Pham*⁺¹, *Hazen Babcock*², *Tomas Lukes*^{3,4}, *Thomas Pengo*⁵, *Jerry Chao*^{6,7}, *Ramraj*
4 *Velmuruga*^{7,8}, *Alex Herbert*⁹, *Anurag Agrawal*¹⁰, *Silvia Colabrese*^{1,11}, *Ann Wheeler*¹², *Anna Archetti*¹³, *Bernd*
5 *Rieger*¹⁴, *Raimund Ober*^{6,7,15}, *Guy M. Hagen*¹⁶, *Jean-Baptiste Sibarita*^{17,18}, *Jonas Ries*¹⁹, *Ricardo Henriques*²⁰,
6 *Michael Unser*¹, *Seamus Holden*^{*+21}

7 *Corresponding authors: daniel.sage@epfl.ch, seamus.holden@ncl.ac.uk.

8 +Equal contribution

9 1: Biomedical Imaging Group, School of Engineering, Ecole Polytechnique Fédérale de Lausanne
10 (EPFL), Switzerland
11 2: Harvard Center for Advanced Imaging, Harvard University, Cambridge, Massachusetts, USA
12 3: Laboratory of Nanoscale Biology & Laboratoire d'Optique Biomédicale, STI - IBI, EPFL, Lausanne,
13 Switzerland
14 4: Department of Radioelectronics, FEE, Czech Technical University, Prague, Czech Republic
15 5: University of Minnesota Informatics Institute, University of Minnesota Twin Cities, USA
16 6: Department of Biomedical Engineering, Texas A&M University, College Station, Texas, USA
17 7: Department of Molecular and Cellular Medicine, Texas A&M University Health Science Center,
18 College Station, Texas, USA
19 8: Department of Microbial Pathogenesis and Immunology, Texas A&M University Health Science
20 Center, Bryan, Texas, USA
21 9: MRC Genome Damage and Stability Centre, School of Life Sciences, University of Sussex, Brighton,
22 UK
23 10 : Double Helix LLC, Boulder, Colorado, USA
24 11 : Istituto Italiano di Tecnologia, Genova, Italy
25 12: Advanced Imaging Resource, Institute of Genetics and Molecular Medicine, University of
26 Edinburgh, Edinburgh, UK
27 13 : Laboratory of Experimental Biophysics, École Polytechnique Fédérale de Lausanne (EPFL),
28 Lausanne, Switzerland
29 14: Department of Imaging Physics, Delft University of Technology, The Netherlands
30 15: Centre for Cancer Immunology, University of Southampton, Southampton, UK
31 16: UCCS center for the Biofrontiers Institute, University of Colorado at Colorado Springs, Colorado,
32 USA
33 17: Interdisciplinary Institute for Neuroscience, University of Bordeaux, Bordeaux, France
34 18: Interdisciplinary Institute for Neuroscience, Centre National de la Recherche Scientifique (CNRS)
35 UMR 5297, Bordeaux, France
36 19: European Molecular Biology Laboratory (EMBL), Cell Biology and Biophysics Unit, Heidelberg,
37 Germany
38 20: Quantitative Imaging and Nanobiophysics Group, MRC Laboratory for Molecular Cell Biology,
39 University College London, UK
40 21: Centre for Bacterial Cell Biology, Institute for Cell and Molecular Biosciences, Newcastle
41 University, UK

42 **ABSTRACT**

43 With the widespread uptake of 2D and 3D single molecule localization microscopy, a large set of
44 different data analysis packages have been developed to generate super-resolution images. To guide
45 researchers on the optimal analytical software for their experiments, in a large community effort we
46 designed a competition to extensively characterise and rank these options. We generated realistic
47 simulated datasets for popular imaging modalities – 2D, astigmatic 3D, biplane 3D, and double helix
48 3D – and evaluated 36 participant packages against these data. This provides the first broad
49 assessment of 3D single molecule localization microscopy software, provides a holistic view of how
50 the latest 2D and 3D single molecule localization software perform in realistic conditions, and
51 ultimately provides insight into the current limits of the field.

52 INTRODUCTION

53 Image processing software is central to single molecule localization microscopy (SMLM¹⁻³). Efficient
54 and automated image processing is essential to extract the super-resolved positions of individual
55 molecules from thousands of raw microscope images, containing millions of blinking fluorescent
56 spots. Improvements in SMLM image processing have been crucial in maximizing spatial resolution
57 and reducing imaging time of SMLM for compatibly with live cell imaging⁴⁻⁶. If SMLM is to achieve a
58 resolving power approaching that of electron microscopy, the analysis software employed needs to
59 be robust, accurate, and performing at current algorithmic limits. This can only be achieved through
60 rigorous quantification of SMLM software performance.

61 The first localization microscopy software challenge was carried out in 2013 to benchmark 2D SMLM
62 software⁷. But biology is not just a 2D problem, and a key focus of localization microscopy is the
63 imaging of 3D imaging of nanoscale cellular processes^{8,9}. 3D localization microscopy is a more
64 difficult image processing problem than 2D SMLM. In addition to finding the center of diffraction
65 limited spots to super-resolve lateral position, 3D SMLM algorithms must also extract axial
66 information from the image, usually by measuring small changes in the shape of a point spread
67 function¹⁰ (PSF).

68 Despite the widespread use of 3D localization microscopy, and challenging nature of 3D SMLM
69 image processing, the performance of software for 3D single molecule localization microscopy has
70 previously only been assessed for 2-3 software packages at a time, and without standard test data or
71 metrics¹¹⁻¹⁴. In the absence of common reference datasets and reliable assessment, it is not possible
72 to objectively assess how different software affects final image quality, or which algorithmic
73 approaches are most successful. Crucially, end-users cannot determine which 3D SMLM software
74 package and imaging modality is optimal for their application.

75 We therefore ran the first 3D localization microscopy software challenge, to assess the performance
76 of 3D SMLM software. We assessed software performance on simulated datasets designed for
77 maximum realism, incorporating experimentally derived point spread functions, using biologically
78 inspired structures, signal to noise levels based closely on common experimental conditions, and
79 modelling fluorophore photophysics. We assessed software performance on synthetic datasets for
80 three popular 3D SMLM modalities: astigmatic imaging¹⁰, biplane imaging¹⁵ and double helix point
81 spread function microscopy¹⁶. We also assessed astigmatism software performance on two real
82 STORM datasets. Furthermore, we ran a second 2D localization microscopy software challenge to
83 assess performance of the latest 2D SMLM software.

84 RESULTS

85 Competition design

86 We established a broad committee from the SMLM community, including experimentalists and
87 software developers, to define the scope of the challenge, ensure realism of the datasets and define
88 analysis metrics. We opened this discussion to all interested parties in an online discussion forum¹⁷.

89 In 2016, we ran a first round of the 3D SMLM competition with explicit submission deadlines,
90 culminating in a special session at the 6th annual Single Molecule Localization Microscopy
91 Symposium (SMLMS 2016). Since then, the challenge has been opened to continuously accept new
92 entries. Thirty-six software packages have been entered in the competition thus far, including four
93 packages used in commercial software (**Table S1, Supplementary Note 1**). Participation in the
94 competition actually led at least eight teams to modify their software to support additional 3D
95 SMLM modalities, showing how competition can foster microscopy software development.

96 Realistic 3D simulations

97 Testing super-resolution software on experimental data lacks the ground truth information required
98 for rigorous quantification of software performance. Therefore, realistic simulated datasets are

99 required. A critical challenge to in simulating 3D SMLM data was to accurately model the
100 experimental microscope PSF for each 3D modality. 3D SMLM inherently involves addition of
101 aberrations to the microscope PSF to encode the Z-position of the molecule. For the PSF models
102 included in the competition: astigmatic (AS), double helix (DH), and biplane (BP), we observed that
103 the PSFs showed complex aberrations not well described by simple analytical models (**Fig. S1**). Even
104 experimental 2D PSFs showed significant aberrations away from the focal plane (**Fig. S1**).

105 We thus combined experimental 3D PSFs with simulated ground truth by performing simulations
106 using PSFs directly derived from experimental calibration data (**Fig. 1, Methods**). We generated
107 simulated datasets over a range of spot densities and signal to noise levels, for simulated
108 microtubule- and endoplasmic reticulum-like structures, using a 4-state model for photophysics¹⁸
109 (**Methods**).

110 **Quantitative performance assessment of 3D software**

111 We assessed software performance by 26 quality metrics (**Supplementary Note 2**). The complete set
112 of summary statistics, axially resolved performance and super-resolved images is available for each
113 competition software on the competition website. We built an interactive ranking and graphing
114 interface for ranking and plotting software performance by any metric, including new user defined
115 metrics (**Fig. S2**). Detailed individual software reports can also be accessed, along with a tool for
116 side-by-side comparison of software (**Fig. S2, S3**).

117 We focused our primary analysis on metrics directly assessing performance in detecting individual
118 molecules. This was based on three key metrics (**Methods**):

- 119 1. *Root mean squared localization error* (RMSE) between measured molecule position and the
120 ground truth.
- 121 2. *Jaccard index* (JAC). This quantifies the fraction of correctly detected molecules in a dataset.
- 122 3. *Efficiency* (*E*). For ranking purposes, we developed a single summary statistic for overall
123 evaluation of software performance combining RMSE and Jaccard index, which we term the
124 *efficiency* (**Methods**).

125 Choice of ranking metric is discussed in **Supplementary Note 2**, where several alternative ranking
126 metrics are also presented.

127 **Performance of 3D software**

128 Complete rankings for each imaging modality and spot density are presented (**Fig. 2**), together with
129 summary information on all competition software (**Supplementary Table 1, Supplementary Note 1**).

130 After assembling an overall summary of best performers for each competition category, we
131 investigated the performance of software within each imaging modality.

132 *Astigmatic localization microscopy*

133 Astigmatic localization microscopy is probably the most popular 3D SMLM modality, reflected by the
134 highest number of software submissions in the 3D competition (**Fig. 2**). For astigmatism, we
135 observed a large spread of software performance, even for the most straightforward high SNR, low
136 spot density (LD) conditions (**Fig. 3, Supplementary Table 2**). The best-in-class software (SMAP-
137 2018¹⁹) has significantly better localization error and Jaccard index performance than average
138 (lateral RMSE 26 nm best vs 38 nm average, axial RMSE 29 nm best vs 66 nm average, Jaccard index
139 85 % best vs 74 % average). Clearly, the quality of the image reconstruction depends strongly on
140 choice of 3D software.

141 To investigate the reasons for software variation, we inspected plots of software performance as a
142 function of axial position in the low density, high SNR dataset for best-in-class and representative
143 middle-range software (**Fig. S4A**). We observed that a key cause of the spread in software

144 performance is variation in software performance away from the focal plane. Near the focal plane,
145 most software packages perform well. However, the axial and lateral RMSE away from the plane of
146 focus is significantly higher for the best in class software, and the Jaccard index is also slightly
147 improved (**Fig. S4A**). This is also visibly apparent in the super-resolved images (**Fig. 4A**). We observed
148 that best-in-class software had a Z-range (the FWHM range of axially resolved software recall,
149 **Methods**) of 1170 nm, greater than two-thirds of the simulated range. Outside this range, the recall
150 and Jaccard index dropped sharply, probably due the large increase in PSF size and decrease in
151 effective SNR at large defocus (**Fig. S1**).

152 When we examined results for the low SNR, low density dataset (**Fig. 2A, 3F**), we found an expected
153 two-fold degradation in best-in-class RMSE (lateral RMSE 39 nm, axial RMSE 60 nm), due to the
154 decrease in image SNR. However, the best-in-class software (SMolPhot²⁰) Jaccard index was
155 effectively constant between the low and high SNR datasets (86 % vs 85 %), although the Z-range did
156 drop at lower SNR (930 nm vs 1120 nm). The best astigmatism software packages were thus
157 remarkably good at finding spots at low SNR, even away from the focal plane.

158 We compared best-in-class software performance to Cramér-Rao lower bound (CRLB) theoretical
159 limits (**Fig. S5, S6, Supplementary Note 3**). Close to the focus, best-in-class software was near the
160 CRLB (within 25 %), but significant deviations from the CRLB occurred > 200 nm (**Fig. S6**). This could
161 be due to difficulty in distinguishing signal from false positives away from focus.

162 Astigmatic software performance dropped for the challenging high spot density datasets (**Fig. 2A, 3**).
163 For the high SNR high spot density dataset (best software, SMolPhot), localization error increased
164 and Jaccard index decreased significantly compared to the low density condition (lateral RMSE best
165 HD 51 nm vs best LD 27 nm, axial RMSE best HD 66 nm vs best LD 29 nm, Jaccard index best HD 66 %
166 vs best LD 85 %). Inspection of the super-resolved images (**Fig. S7**) nevertheless shows qualitatively
167 acceptable results for the HD dataset, particularly in the lateral dimension. In some circumstances,
168 the performance reduction at 10x higher spot density could be acceptable for 10x faster, potentially
169 live-cell-compatible, imaging speed. We also observed a large spread of software performance for
170 the high density datasets, probably because a significant fraction of the software packages were
171 primarily designed for low density conditions.

172 We observed poor performance for the most challenging low SNR high spot density astigmatism
173 dataset (**Fig. 2A, 3, S8**, best software SMolPhot). Best-in-class localization precision and Jaccard
174 index decreased significantly (lateral RMSE 76 nm, axial RMSE 101 nm, Jaccard index 58 %). These
175 data suggest that low SNR high density 3D astigmatic localization microscopy entails significant
176 reduction in image resolution.

177 *Double helix point spread function localization microscopy*

178 We next analyzed the performance of the double helix software (**Fig. 3D-F, S9A**). For the software in
179 the high SNR low spot density condition, double helix software showed more uniform performance
180 than astigmatism. Best-in-class software (SMAP-2018) showed only a limited improvement
181 compared with average software (**Fig. 3D-F**, lateral RMSE, 27 nm best vs 37 nm average; axial RMSE
182 21 nm best vs 34 nm average; Jaccard index 77 % best vs 73 % average). In general software
183 localization performance was close to the CRLB (**Fig. S6**). We observed that performance of the
184 software away from the focal plane is relatively uniform (**Fig. 4A, S4A**), and best-in-class Z-range at
185 high SNR was large at 1180 nm (**Fig. S4A, Supplementary Table 2**). Double helix imaging may show
186 less software-to-software variation and larger Z-range at low spot density than astigmatic imaging
187 because the PSF shape and intensity are fairly constant as a function of Z; unlike astigmatic imaging,
188 where spot size, shape and intensity vary greatly as a function of Z (**Fig. S1**).

189 Double helix software performance decreased significantly for the low spot density low SNR
190 condition (best software, SMAP-2018), particularly in terms of best-in-class Jaccard index (66 % low

191 SNR vs 77 % high SNR, **Fig. 3D-E, S8, S9A**). DH Jaccard index was also significantly worse than
192 astigmatism results at either high or low SNR (85 % high SNR, 86 % low SNR). This indicates that it
193 was quite hard to successfully find localizations in the low SNR DH dataset, likely because the large
194 size of the DH PSF spreads emitted photons over a large area, lowering effective image SNR. DH PSF
195 designs with reduced Z-range but more compact PSF would likely be less sensitive to this issue²¹.

196 Double helix software performed poorly on the high spot density datasets at high SNR (best software
197 CSpline²²), especially in terms of the Jaccard index (**Fig. 3D-E, S9A**, best lateral RMSE 67 nm, best
198 axial RMSE 69 nm, best Jaccard index 46 %). The poor performance at high spot density is again
199 probably because the large DH PSF size increases spot density and decreases SNR (**Fig. S1**). DHPSF
200 performance at high spot density and low SNR was also not reliable (**Fig. 3D-F, S9A**, best software,
201 SMAP-2018).

202 *Biplane localization microscopy*

203 Best-in-class biplane software (SMAP-2018), at low spot density and for both high and low SNR,
204 delivered the best performance in any modality (high SNR: lateral RMSE 12.3 nm, axial RMSE 21.7
205 nm, Jaccard 87 %), despite a slightly decreased image SNR for the biplane simulations (**Methods**).
206 We observed a large spread in software performance in terms of lateral RMSE and Jaccard index,
207 with the best-in-class software significantly outperforming the other competitors (**Fig. S9B, 2D**). At
208 low spot density, best-in-class biplane software (SMAP-2018) showed good performance as a
209 function of Z, with high Jaccard index over almost the entire Z-range of the simulations, and with a Z-
210 range of 1200 nm at high SNR (**Fig. S4AC, Supplementary Table 2**). The axial RMSE was relatively
211 uniform as a function of Z and close to the CRLB limit (**Fig. S6**). As axial and lateral RMSE are both
212 averaged over the entire Z-range, the strong biplane results arise from good performance across a
213 large Z-range (**Fig. S4**).

214 At high spot density and high SNR, best-in-class biplane software (SMAP-2018) showed acceptable
215 performance (**Fig. 3D-F, S7, S9B**, best lateral RMSE 43 nm, best axial RMSE 49 nm, best Jaccard index
216 61 %). Uniquely among the 3D modalities, best-in-class biplane software also gave acceptable
217 performance at high spot density and low SNR (**Fig. 3D-F, S7, S9B**, best lateral RMSE 55 nm, best
218 axial RMSE 72 nm, best Jaccard index 61 %, best software SMAP-2018).

219 **Performance of 2D software**

220 We next assessed the performance of 2D SMLM software. For the pseudo-ER 2D dataset, at low
221 density best-in-class software (ADCG²³) performed substantially better than the class average
222 (**Fig. S10, S11**, lateral RMSE 31 nm vs 36 nm average, Jaccard index 90 % best vs 72 %). Low density
223 results for the brighter fluorophore microtubules dataset were similar to the dimmer pseudo-ER
224 dataset (**Fig. S10, S12** best software SMolPhot). For the very high density 2D dataset, which had 25x
225 higher spot density than the LD dataset, best-in-class software (ADCG) showed excellent
226 performance (**Fig. S10**, lateral RMSE, 45.5 nm, Jaccard index 75%). Best-in-class performance (ADCG)
227 on the dimmer fluorophore data at high spot density was also strong (**Fig. S10**, best lateral RMSE 51
228 nm, best Jaccard index 70 %).

229 **Algorithms**

230 We identified several classes of algorithm participant software (**Supplementary Table 1**):

231 1) *Non-iterative* software regroups pixels in the local neighborhood of the candidates, like
232 interpolation, center of mass (QuickPALM²⁴) or template matching (WTM²⁵). These often older
233 algorithms are fast but tend to achieve poor performance.

234 2) *Single emitter fitting* software is usually built on a multi-step strategy of detection, spot
235 localization, and optional spot rejection. The detection step finds bright spots in noisy images on the
236 pixel grid. The selection of candidates is usually performed by local maximum search after a

237 denoising filter. Others rely on more complex algorithms like the wavelet transform (WaveTracer²⁶).
238 We did not observe software ranking to depend noticeably on the choice of optimization scheme:
239 least-square, weighted least-square or maximum-likelihood estimator.

240 3) *Multi-emitter fitting* software groups clusters of overlapping spots, and simultaneously fits
241 multiple model PSFs to the data. Typically, fitted spots are added to the cluster until a stopping
242 condition is met^{4,5}. This leads to improved localization performance at high spot density, at the cost
243 of reduced speed. This class of software (e.g., 3D-DAOSTORM¹¹, CSpline, PeakFit, ThunderSTORM²⁷)
244 was amongst the top performers in each 2D and 3D competition category.

245 As expected, single- and multiple-emitter fitting methods both performed well on low density data.
246 For the 2D challenge, multi-emitter fitting showed a clear advantage over single emitter fitting at
247 high density. Surprisingly however, well-tuned single-emitter fitting algorithms (SMolPhot, SMAP-
248 2018) outperformed multi-emitter algorithms for the 3D high density conditions.

249 4) *Compressed sensing algorithms*. One subset of these algorithms utilize deconvolution with
250 sparsity constraints to reconstruct super-resolved images²⁸⁻³⁰. Although deconvolution approaches
251 can give good results, they are limited by the necessary use of a sub-pixel grid; increased localization
252 precision requires smaller grid resolution, which must be balanced against increased computational
253 time. Recent approaches address this issue by localizing the point sources in a gridless manner under
254 some sparsity constraint (ADCG, SMfit, SOLAR_STORM, TVSTORM³¹). This software class consistently
255 gave the overall best performance for 2D high-density (ADCG 1st, FALCON³⁰ 2nd, SMfit 3rd).

256 5) *Other approaches*. Of the alternative algorithmic approaches used, the annihilating filter-based
257 method LEAP³² gave good performance for biplane imaging. Recently, we received the first challenge
258 submission from a deep learning SMLM software (DECODE); these promising preliminary results are
259 available on the competition website.

260 *Post-hoc temporal grouping*

261 Because molecule on-time is stochastically distributed across multiple frames, a common post-
262 processing approach to improve localization precision is to group molecules detected multiple times
263 in adjacent frames, and average their position³³ (**Supplementary Note 4**). Temporal grouping was
264 used by the top performers (including SMolPhot, MIATool³⁴ and SMAP-2018), and is visibly apparent
265 as a more punctate super-resolved image (**Fig. 4A**).

266 *Choice of PSF model*

267 Most software used a variant of Gaussian PSF model. A few participants designed more accurate PSF
268 models. Either diffraction theory was used (MIATool, LEAP) or spline fitting of an analytical function
269 to the experimental PSF was adopted (CSpline, SMAP-2018). Although simple Gaussian model PSFs
270 were sufficient to obtain best-in-class performance for the 2D and astigmatic modalities (ADCG,
271 PeakFit, SMolPhot), top results for the more optically complex biplane and double helix modalities
272 were exclusively software using non-Gaussian PSF models (SMAP-2018, CSpline, MIATool, LEAP).

273 *Multi-algorithm packages*

274 Several software packages take a Swiss army knife approach of integrating multiple optional
275 localization algorithms into one program, to be flexible enough to suit various experimental
276 conditions^{19,27}. SMAP-2018 and ThunderSTORM achieved strong across-the-board performance
277 supporting this rationale.

278 *Software run time*

279 Software run time is important both for ease of use and real time analysis. We did not observe
280 correlation between software localization performance (Efficiency) and software run time (**Fig.**

281 **S13A**). We thus created an alternative ranking metric, *Efficiency-Runtime*, which gave 25 % weighting
282 to run time (**Supplementary Note 2.7, Fig S13B**). Many good performers in the efficiency-only
283 ranking were relatively fast and thus retained good ranking (SMAP-2018, SMolPhot, 3D-
284 DAOSTORM). Interestingly, two software packages highly optimized for speed gained top ranking in
285 this analysis: pSMLM-3D³⁵ and QC-STORM.

286 *Diagnostic tools for software and algorithm performance*

287 During our analysis, we frequently noticed common types of deviation between software results and
288 ground truth which were easily diagnosed by visual inspection (**Fig. S14, S15**). This included not only
289 obvious issues of poor localization precision or spot averaging at high density, but also more subtle
290 problems such as a common error of structural warping which significantly reduced software
291 performance. On the competition website, we provide detailed diagnostic software reports including
292 multiple examples of software performance on individual frames to help developers to identify
293 algorithm and software limitations and maximize software performance (**Fig. S3, S16**).

294 **Assessment on real STORM data**

295 We investigated the performance of a representative subset of astigmatism software on real STORM
296 datasets of well characterized test structures, microtubules and nuclear pore complex, NPC (**Fig. 4B,**
297 **S17**). This qualitative assessment was consistent with findings for simulated data. No performance
298 difference between single and multi-emitter fitters was observed, which is not surprising since spot
299 density in these datasets was low. Relatively poor software performance was immediately obvious
300 from visual inspection (QuickPALM). Temporal grouping noticeably improved resolution (3D-
301 DAOSTORM, CSpline, MIATool, SMAP-2018). Gaussian fitting software . Interestingly, although
302 Gaussian/ Bessel PSF modelling software (3D-DAOSTORM, MIATool, ThunderSTORM) gave high
303 resolution images, software which modelled the experimental PSF via spline fitting (CSpline, SMAP-
304 2018) gave noticeably improved resolution of fine structural features such as the top and bottom of
305 the NPC (**Fig. 4B**) or the hollow core of antibody-labelled microtubules (**Fig. S17**).

306 **DISCUSSION**

307 The strongest conclusion we draw from the 3D localization microscopy challenge is that choice of
308 localization software greatly affects the quality of final super-resolution data, even at “easy” high
309 SNR, low spot density conditions. Biplane performance was particularly dependent on software
310 choice, with only one software (SMAP-2018) achieving near-Cramér-Rao lower bound performance.
311 Double helix SMLM showed less sensitivity to choice of software than biplane, with astigmatic SMLM
312 intermediate between the two. The best software in each modality performed close to the Cramér-
313 Rao lower bounds over a wide focal range and successfully detected most molecules, even at low
314 signal to noise. Average software in all three modalities was significantly worse, with the obtained
315 axial resolution being particularly sensitive to software choice.

316 The second major conclusion is that localization software that explicitly includes the experimental
317 PSF in the fitting model gives a significant performance increase for 3D SMLM. For the more optically
318 complex biplane and double helix modalities in particular, the best results were from software which
319 incorporated non-Gaussian PSF models (SMAP-2018, CSpline, MIATool). This result also highlights
320 the importance of accurate PSF modelling in 3D SMLM simulations. The performance advantage of
321 experimental PSF fitting software would not have been observable had simulations been generated
322 with a simple Gaussian PSF.

323 Of the different algorithm classes, well-tuned single-emitter and multi-emitter fitting algorithms
324 (each capable of dealing well with occasional molecule overlap) gave good results for low density 3D
325 SMLM. We also found that several software packages for astigmatic or biplane imaging gave
326 adequate performance for the challenging case of high molecule densities, as long as the image SNR
327 was high. Current software packages gave poor performance when molecule density was high and

328 image SNR was low. These results indicate that with current algorithms high density 3D SMLM
329 performance is mediocre at high SNR and poor at low SNR. Surprisingly, multi-emitter fitting did not
330 show significant improvement over well-tuned single emitter fitting for the 3D high-density datasets;
331 this may indicate that significant potential for improvement remains in this category.

332 Many software packages did not apply temporal grouping³³, resulting in reduced software
333 performance. Since temporal grouping is a simple step for maximum precision, we urge all software
334 developers to integrate this approach into their software as an optional final step in the localization
335 process.

336 The second 2D localization microscopy challenge provided the opportunity to reassess the state of
337 the field. The performance of best-in-class 2D software over a range of conditions, at both high and
338 low spot density, was very strong. Interestingly, the top three performers in the 2D high density
339 condition were all compressed sensing algorithms (ADCG, FALCON, SMfit). In low density 2D
340 conditions, the best single-emitter, multi-emitter and compressed sensing algorithms all gave
341 comparable, excellent, performance. We speculate that performance in the low spot density 2D
342 category might now be near optimal levels.

343 In future, we plan to extend the SMLM challenge into an open platform with a fully automated
344 assessment process, and where new competition simulations and assessment metrics can easily be
345 created and contributed by the community. It will be important to account for new technologies and
346 developments in SMLM, such as scientific CMOS cameras⁶, in future simulations. It would also be
347 exciting to adapt the tools developed in the SMLM challenge to other classes of super-resolution
348 microscopy, such as fluorescence-fluctuation-based super-resolution microscopies (*e.g.*, 3B³⁶, SOFI³⁷,
349 SRRF³⁸) and structured illumination microscopy³⁹.

350 The results of this competition show that the best 2D and 3D localization microscopy software have
351 formidable algorithmic performance. However, a problem that often hinders adoption of new SMLM
352 algorithms is that only a small subset of algorithms is packaged in, or compatible with fast, well-
353 maintained, user-friendly software packages, which include all stages of the SMLM data analysis
354 pipeline – analysis, visualization and quantification. This remains a key outstanding challenge for the
355 field.

356 Both the 3D and 2D localization microscopy software challenges remain open and continuously
357 updated on the competition website. This continuously evolving analysis of SMLM software
358 performance provides software developers with a robust means of benchmarking new algorithms,
359 and helps to ensure that super-resolution microscopists use software that gets the best out of their
360 hard-won data.

361

362 ACKNOWLEDGEMENTS

363 *Authors acknowledge the following funding sources: a Newcastle University Research Fellowship and*
364 *a Wellcome Trust & Royal Society Sir Henry Dale Fellowship grant number 206670/Z/17/Z to SH; an*
365 *European Research Council (ERC) under the European Union's Horizon 2020 research and innovation*
366 *programme, Grant Agreement no. 692726 to DS, TAP, MU; UK BBSRC grants BB/M022374/1,*
367 *BB/P027431/1, BB/R000697/1 grant and MRC grants MC-UU-12018/2, MR/K015826/1 to RH;*
368 *European Research Council (ERC) grant CoG-724489, CellStructure to JR; FranceBioImaging*
369 *infrastructure ANR-10-INBS-04 to J.-B.S; National Institutes of Health grant 1R15GM128166-01 to*
370 *GMH; and NSF SBIR grants 1353638, 1534745 to Double Helix LLC. We thank R. Piestun at University*
371 *of Colorado for providing DH-PSF phase mask designs to Double Helix LLC. We thank all the*
372 *localization microscopy challenge participants for their contribution: Hazen Babcock (3D-DAOSTORM,*
373 *Cspline, L1H), Fabian Hauser (3D-STORM Tools), Shigeo Watanabe (3D-WTM,WTM), Nicholas Boyd*
374 *(ADCG), Junhong Min, Kyong Jin and Jong Chul Ye (ALOHA, FALCON), Hervé Rouault (B-recs),*
375 *Emmanuel Soubies (CELO-STORM), Artur Speiser, Srinivas Turagas and Jakob Macke (DECODE), Alex*
376 *von Diezmann, Camille Bayas and W. E. Moerner (Easy-DHPSF), Thomas Vomhof and Jochen Reichel*
377 *(FIRESTORM), Hanjie Pan (LEAP), Ann Wheeler (Localizer), Zhen-li Huang and Yujie Wang*
378 *(MaLiang), J. Chao, R. Velmurugan, A. V. Abraham and R. J. Ober (MIATool), Hendrik Deschout*
379 *(mlePALM), Thomas Pengo (Octane, PeakSelector), Yi-na Wang (PALMER), Alex Herbert*
380 *(PeakFit), Koen Martens and Johannes Hohlbein (pSMLM-3D), Luchang Li (QC-STORM), Ricardo*
381 *Henriques (QuickPALM), G. Tamas and J. Sinko (RainSTORM), Steve Wolter and Markus Sauer*
382 *(RapidSTORM), Manfred Kirchgessner and Frederik Gruell (SFP Estimator), Yiming Li and Jonas Ries*
383 *(SMAP), Hayato Ikoma (SMfit), A. Loot, A. Valdmann, M. Eltermann, M. Kree and M. Pärs*
384 *(SMolPhot), Yoon J. Jung, Anthony Barsic Rafael Pietsun, and Nikta Fakhri (SOLAR_STORM), Anna*
385 *Archetti (STORMChaser), Martin Ovesny, Guy Hagen and Pavel Krizek (ThunderSTORM), Jiaqing*
386 *Huang (TVSTORM), Adel Kechkar, Corey Butler and Jean-Baptiste Sibarita (WaveTracer) and Benoît*
387 *Lelandais (ZOLA-3D). We thank the SMLMS 2016 organizers (S. Manley and A. Radenovic, EPFL) for*
388 *hosting a localization microscopy challenge special session. We also thank Double Helix LLC and*
389 *Molecular Devices LLC for sponsoring the SMLMS 2016 special session. The sponsors had no input or*
390 *influence on the research.*

391 AUTHOR CONTRIBUTIONS

392 DS and SH conceived and coordinated the study. DS, SH, TAP, AAr, HB, SC, AW, GMH, RH, TL, TP, JBS
393 designed the study. SH, AAg, RH, JBS collected experimental PSFs. DS, TAP, SH, TL wrote simulation
394 code. BR shared unpublished software. DS generated simulated datasets. JR shared experimental
395 STORM data. AH, JR, JC, RV provided feedback and quality control on simulations and analysis
396 methods. TAP carried out the assessment of software performance. TAP, DS, SH analysed
397 and interpreted the results. DS, HB, RO, BR, GMH, JBS, JR, RH, MU, SH directed research. SH, DS, TAP
398 wrote the manuscript with feedback from all authors.

399 REFERENCES

- 400 1. Betzig, E. *et al.* Imaging Intracellular Fluorescent Proteins at Nanometer Resolution. *Science*
401 **313**, 1642–1645 (2006).
- 402 2. Hess, S. T., Girirajan, T. P. K. & Mason, M. D. Ultra-High Resolution Imaging by Fluorescence
403 Photoactivation Localization Microscopy. *Biophys. J.* **91**, 4258–4272 (2006).
- 404 3. Rust, M. J., Bates, M. & Zhuang, X. Sub-diffraction-limit imaging by stochastic optical
405 reconstruction microscopy (STORM). *Nat Methods* **3**, 793–795 (2006).
- 406 4. Holden, S. J., Uphoff, S. & Kapanidis, A. N. DAOSTORM: an algorithm for high- density super-
407 resolution microscopy. *Nat Meth* **8**, 279–280 (2011).
- 408 5. Huang, F., Schwartz, S. L., Byars, J. M. & Lidke, K. A. Simultaneous multiple-emitter fitting for
409 single molecule super-resolution imaging. *Biomed. Opt. Express* **2**, 1377–1393 (2011).

- 410 6. Huang, F. *et al.* Video-rate nanoscopy using sCMOS camera-specific single-molecule
411 localization algorithms. *Nat. Methods* **10**, 653–658 (2013).
- 412 7. Sage, D. *et al.* Quantitative evaluation of software packages for single-molecule localization
413 microscopy. *Nat. Methods* **12**, 717–724 (2015).
- 414 8. Huang, B., Jones, S. A., Brandenburg, B. & Zhuang, X. Whole-cell 3D STORM reveals
415 interactions between cellular structures with nanometer-scale resolution. *Nat Meth* **5**, 1047–1052
416 (2008).
- 417 9. Shtengel, G. *et al.* Interferometric fluorescent super-resolution microscopy resolves 3D
418 cellular ultrastructure. *Proc. Natl. Acad. Sci.* **106**, 3125–3130 (2009).
- 419 10. Huang, B., Wang, W., Bates, M. & Zhuang, X. Three-Dimensional Super-Resolution Imaging
420 by Stochastic Optical Reconstruction Microscopy. *Science* **319**, 810–813 (2008).
- 421 11. Babcock, H., Sigal, Y. M. & Zhuang, X. A high-density 3D localization algorithm for stochastic
422 optical reconstruction microscopy. *Opt. Nanoscopy* **1**, 1–10 (2012).
- 423 12. Ovesný, M., Křížek, P., Švindrych, Z. & Hagen, G. M. High density 3D localization microscopy
424 using sparse support recovery. *Opt. Express* **22**, 31263–31276 (2014).
- 425 13. Min, J. *et al.* 3D high-density localization microscopy using hybrid astigmatic/ biplane
426 imaging and sparse image reconstruction. *Biomed. Opt. Express* **5**, 3935–3948 (2014).
- 427 14. Zhang, S., Chen, D. & Niu, H. 3D localization of high particle density images using sparse
428 recovery. *Appl. Opt.* **54**, 7859–7864 (2015).
- 429 15. Juette, M. F. *et al.* Three-dimensional sub-100 nm resolution fluorescence microscopy of
430 thick samples. *Nat. Methods* **5**, 527–529 (2008).
- 431 16. Pavani, S. R. P. *et al.* Three-dimensional, single-molecule fluorescence imaging beyond the
432 diffraction limit by using a double-helix point spread function. *Proc. Natl. Acad. Sci.* **106**, 2995–2999
433 (2009).
- 434 17. Collaboration through competition. *Nat. Methods* **11**, 695 (2014).
- 435 18. Annibale, P., Vanni, S., Scarselli, M., Rothlisberger, U. & Radenovic, A. Quantitative Photo
436 Activated Localization Microscopy: Unraveling the Effects of Photoblinking. *PLOS ONE* **6**, e22678
437 (2011).
- 438 19. Li, Y. *et al.* Real-time 3D single-molecule localization using experimental point spread
439 functions. *Nat. Methods* (2018). doi:10.1038/nmeth.4661
- 440 20. Loot A. , Valdmann A., Eltermann M., Kree M., Pärs M. SMolPhot Software. Available at:
441 <https://bitbucket.org/ardiloot/>. (Accessed: 28th January 2019)
- 442 21. Grover, G., DeLuca, K., Quirin, S., DeLuca, J. & Piestun, R. Super-resolution photon-efficient
443 imaging by nanometric double-helix point spread function localization of emitters (SPINDLE). *Opt.*
444 *Express* **20**, 26681–26695 (2012).
- 445 22. Babcock, H. P. & Zhuang, X. Analyzing Single Molecule Localization Microscopy Data Using
446 Cubic Splines. *Sci. Rep.* **7**, 552 (2017).
- 447 23. Boyd, N., Schiebinger, G. & Recht, B. The Alternating Descent Conditional Gradient Method
448 for Sparse Inverse Problems. *SIAM J. Optim.* **27**, 616–639 (2017).
- 449 24. Henriques, R. *et al.* QuickPALM: 3D real-time photoactivation nanoscopy image processing in
450 ImageJ. *Nat Meth* **7**, 339–340 (2010).
- 451 25. Takeshima, T., Takahashi, T., Yamashita, J., Okada, Y. & Watanabe, S. A multi-emitter fitting
452 algorithm for potential live cell super-resolution imaging over a wide range of molecular densities. *J.*
453 *Microsc.* **271**, 266–281 (2018).
- 454 26. Kechkar, A., Nair, D., Heilemann, M., Choquet, D. & Sibarita, J.-B. Real-Time Analysis and
455 Visualization for Single-Molecule Based Super-Resolution Microscopy. *PLOS ONE* **8**, e62918 (2013).
- 456 27. Ovesný, M., Křížek, P., Borkovec, J., Švindrych, Z. & Hagen, G. M. ThunderSTORM: a
457 comprehensive ImageJ plug-in for PALM and STORM data analysis and super-resolution imaging.
458 *Bioinformatics* **30**, 2389–2390 (2014).
- 459 28. Soubies, E., Blanc-Féraud, L. & Aubert, G. A Continuous Exact l0 Penalty (CELO) for Least
460 Squares Regularized Problem. *SIAM J. Imaging Sci.* **8**, 1607–1639 (2015).

- 461 29. Babcock, H. P., Moffitt, J. R., Cao, Y. & Zhuang, X. Fast compressed sensing analysis for super-
462 resolution imaging using L1-homotopy. *Opt. Express* **21**, 28583–28596 (2013).
- 463 30. Min, J. *et al.* FALCON: fast and unbiased reconstruction of high-density super-resolution
464 microscopy data. *Sci. Rep.* **4**, 4577 (2014).
- 465 31. Huang, J., Sun, M., Ma, J. & Chi, Y. Super-Resolution Image Reconstruction for High-Density
466 Three-Dimensional Single-Molecule Microscopy. *IEEE Trans. Comput. Imaging* **3**, 763–773 (2017).
- 467 32. Pan, H., Simeoni, M., Hurley, P., Blu, T. & Vetterli, M. LEAP: Looking beyond pixels with
468 continuous-space Estimation of Point sources. *Astron. Astrophys.* **608**, A136 (2017).
- 469 33. Durisic, N., Laparra-Cuervo, L., Sandoval-Álvarez, Á., Borbely, J. S. & Lakadamyali, M. Single-
470 molecule evaluation of fluorescent protein photoactivation efficiency using an in vivo nanotemplate.
471 *Nat. Methods* **11**, 156–162 (2014).
- 472 34. Chao, J., Ward, E. S. & Ober, R. J. A software framework for the analysis of complex
473 microscopy image data. *IEEE Trans. Inf. Technol. Biomed. Publ. IEEE Eng. Med. Biol. Soc.* **14**, 1075–
474 1087 (2010).
- 475 35. Martens, K. J. A., Bader, A. N., Baas, S., Rieger, B. & Hohlbein, J. Phasor based single-
476 molecule localization microscopy in 3D (pSMLM-3D): An algorithm for MHz localization rates using
477 standard CPUs. *J. Chem. Phys.* **148**, 123311 (2017).
- 478 36. Cox, S. *et al.* Bayesian localization microscopy reveals nanoscale podosome dynamics. *Nat.*
479 *Methods* **9**, 195–200 (2012).
- 480 37. Dertinger, T., Colyer, R., Iyer, G., Weiss, S. & Enderlein, J. Fast, background-free, 3D super-
481 resolution optical fluctuation imaging (SOFI). *Proc. Natl. Acad. Sci.* **106**, 22287–22292 (2009).
- 482 38. Gustafsson, N. *et al.* Fast live-cell conventional fluorophore nanoscopy with ImageJ through
483 super-resolution radial fluctuations. *Nat. Commun.* **7**, (2016).
- 484 39. Gustafsson, M. G. L. Surpassing the lateral resolution limit by a factor of two using structured
485 illumination microscopy. SHORT COMMUNICATION. *J. Microsc.* **198**, 82–87 (2000).
- 486

487 METHODS

488 1. CHALLENGE ORGANIZATION

489 We first ran the 3D SMLM software challenge as a time limited competition, with a results session
490 hosted as a special session of the 6th Annual Single Molecule Localization Microscopy Symposium in
491 August 2016. The competition has now been converted to a permanent software challenge
492 accepting new submissions. Special thanks is due to the software SMAP and 3D-WTM²⁵ that
493 participated in all eight categories (*density x modality*). The current list of participants is at:

494 <http://bigwww.epfl.ch/smlm/challenge2016/index.html?p=participants>

495 All datasets, methods, participations, and results of the challenge 2016 made available at
496 <http://bigwww.epfl.ch/smlm/challenge2016/>. Software for simulation and analysis is hosted on the
497 competition GitHub repository: <https://github.com/SMLM-Challenge/Challenge2016/>

498 A Life Sciences Reporting Summary is associated with this manuscript on the Nature Methods
499 website.

500 2. LOCALIZATION MICROSCOPY SIMULATIONS

501 2.1. Structure, noise levels and spot densities

502 *Structure.* The synthetic datasets were designed to be similar to images derived from real cellular
503 structures. We defined mathematical models for cellular structures that imitate cytoskeletal
504 filaments such as microtubules and larger tubular structures such as the endoplasmic reticulum or
505 mitochondria (**Fig. S18A**). These structures have a tubular shape in the 3D space. For the 3D
506 competition, we simulated synthetic 25 nm diameter microtubules (**Fig. 1**). Pseudo-microtubules are
507 defined with their central axis elongating in a 3D space having an average outer diameter of 25 nm
508 with an inner, hollow tube of 15 nm diameter. For the 2D competition, in addition to synthetic
509 microtubules (MT), we simulated larger diameter 150 nm cylinders, called pseudo-endoplasmic
510 reticulum (pseudo-ER), designed to approximate larger cellular structures such as mitochondria and
511 the endoplasmic reticulum (ER) (**Fig. 1**).

512 The underlying sample structure is formalized in a continuous space which allows rendering of digital
513 images at any scale, from very high resolution (up to 1 nm/pixel) to low resolution (camera
514 resolution: 100 nm/ pixel). The continuous-domain 3D curve is represented by means of a
515 polynomial spline. The sample is imaged in a $6.4 \times 6.4 \mu\text{m}^2$ field of view, and the center lines of the
516 microtubules have limited variation along the z (vertical) axis, *i.e.*, less than 1.5 μm . The fluorescent
517 markers are uniform randomly distributed over the structure according to the required density. The
518 photon emission rate of each fluorophore is controlled by a photo-activation model (see below). The
519 exact locations of all fluorophores are stored at high precision floating-point numbers expressed in
520 nanometers. This ground-truth file is used for conducting objective evaluations without human bias.

521 *Noise levels.* We generated data at three different signal-to-noise ratio (SNR) levels, based on real
522 signal to noise levels encountered under common SMLM experimental scenarios: *N1*, fixed cells
523 antibody labelled with organic dye¹⁰, high signal, medium background; *N2*, fluorescent protein
524 labelling¹, low signal, low background; and *N3*, live cell affinity dye labelling^{40,41}, high signal, high
525 background.

526 *Spot density.* As performance at different density of active emitters is a key challenge for SMLM
527 software, we generated 3D competition datasets at both sparse emitter density
528 ($0.25 \text{ mol. [molecule] } \mu\text{m}^{-2}$), *3D LD* and high emitter density ($2.5 \text{ mol. } \mu\text{m}^{-2}$), *3D HD*. For the 2D
529 competition, we generated a sparse ($0.5 \text{ mol. } \mu\text{m}^{-2}$), *2D LD*, and very high density dataset
530 ($5 \text{ mol. } \mu\text{m}^{-2}$), *2D HD*.

531 Together, these simulated conditions closely resemble experimental 3D and 2D data under a range
532 of challenging conditions of SNR, spot density, axial thickness and structure summarized in
533 **Supplementary Table 3**. In addition, we provide simulated z-stacks of bright beads for software
534 calibration. The competition datasets (**Supplementary Table 4**) are available online on the
535 competition website.

536

537 **2.2. Photophysics activation model**

538 We incorporated a 4-state model of fluorophore photophysics¹⁸, including a transient dark state (dye
539 blinking) and a bleaching pathway (**Fig. S18C**). Given a list of source locations from the structure
540 simulator, fluorophore blinking was simulated by a 4-states Markov chain model. The states are ON,
541 OFF, BLEACH, DARK and the transitions are Poisson distributed (**Fig. S18C**), except for the OFF to ON
542 transitions which follow a uniform random distribution to reflect that in typical experimental
543 conditions, constant imaging density is maintained by tuning the photoactivation rate during the
544 experiment. All switching is calculated at sub-frame resolution and then total fluorophore on-time
545 was integrated over each frame.

546 Due to two decay paths, the actual mean lifetime of the state ON is

$$T_{LIFETIME} = \frac{1}{\frac{1}{T_{ON}} + \frac{1}{T_{BLEACH}}}$$

547 Switching rates were chosen to approximate photoactivatable fluorescent proteins $T_{ON}=3$ frames,
548 $T_{DARK}=2.5$ frames, and $T_{BLEACH}=1.5$ frames.

549 Fractional fluorophore ON-times per frame (between 0 and 1) were multiplied by the mean flux of
550 photon emission. The flux of photons expressed in photons/seconds was given by the relation

$$F = \frac{\Phi P \sigma}{e}$$

551 Φ is the quantum yield of the dye, P is power of the laser in W/cm^2 , $e = h c / \lambda$ is the energy of one
552 photon, $\sigma = 1000 \ln(10) \epsilon / N_A$ is the absorption cross section in cm^2 and ϵ is the molar extinction
553 coefficient (EC) or absorptivity in cm^2/mol which is a characteristic of a given fluorophore. The laser
554 power was Gaussian distributed over the field of view. At the end of this process a list of XY
555 positions, on-frames and (noise-free) intensities for all activated fluorophores was obtained.

556 Analysis of the resulting simulated photon counting distribution is presented in **Supplementary**
557 **Note 5** and **Figure S23**.

558 **2.3. Experimental Point Spread Function**

559 Model PSFs, stored as high resolution look up tables, were derived from experimentally measured
560 PSFs. Although the algorithmic approach is distinct, the concept of accurately modelling the
561 experimental PSF based on calibration data bears relation to the PSF phase retrieval approach
562 previously employed by Hanser and coworkers⁴².

563 Images of fluorescent beads were recorded for each modality (**Supplementary Table 5**). Signal to
564 noise ratio of recorded PSFs was maximized in all cases by maximizing exposure time and averaging
565 over several frames to increase dynamic range.

566 To acquire experimental PSFs, we took 100 nm Tetraspek beads (Invitrogen) adsorbed to #1.5
567 (170 μm thick) coverglass, imaged in water. The excitation wavelength was between 640 nm and 647
568 nm, and a Cy5 emission filter was used. Data acquisition parameters for each modality are listed in
569 **Supplementary Table 5**.

570 The experimental PSFs used to generate the simulated data are available on the competition
571 website. As the goal of this study was to compare software obtained on typical SMLM microscopes,
572 we deliberately chose PSFs representative of common implementations of each 3D modality.
573 However, additional PSF engineering should improve results of any specific modality, for example
574 adaptive-optics corrected astigmatism⁴³, or reduced Z-range, higher SNR DH-PSF designs²¹.

575 The experimental point spread functions used here were measured for fluorescent beads adsorbed
576 to the microscope cover slip, and should be appropriate simulations of SMLM data acquired within a
577 few microns of the cover slip. Performing SMLM imaging at greater depths, *e.g.*, in tissue or even
578 deep within single cells, with oil immersion objectives will cause spherical aberration due to
579 refractive index mismatch⁴⁴. In order to accurately simulate SMLM data acquired at depth, the
580 experimental PSFs could be acquired at a matching depth, by embedding fluorescent beads in
581 agarose. Alternatively, the PSF for beads at the coverslip could be measured and explicitly calculated
582 via phase retrieval, and then convolved with the appropriate degree of spherical aberration⁴⁴.

583

584 **2.4. Simulation PSF construction**

585 For each modality, 3-6 beads were selected within a small (< 32 μm) region, to minimize PSF
586 variation due to spherical aberration. Images for each selected bead were interpolated in XY to a
587 pixel size of 10 nm. Beads were then coaligned by cross-correlation on the in-focus frame. Coaligned
588 beads were averaged in XY to minimize pixel quantization artefacts and to increase SNR. Where
589 necessary, Z-stacks were interpolated to a Z-step size of 10 nm. A central Z-range of 1.5 μm was
590 selected that represents 151 optical planes with a Z-step of 10 nm. The Z-range covers -750 nm to
591 +750 nm. The plane of best focus was chosen as the simulation 0 nm plane. Each model PSF was
592 normalized such that the total intensity of the PSF in the in-focus frame within a diameter of 3
593 FWHM from the PSF center was equal to 1.

594 For the DH PSF, the transmission of the combined phase mask system was measured as 96 %, which
595 was approximated as 100 % brightness relative to the 2D and astigmatic PSFs.

596 In biplane super-resolution microscopy, emitted fluorescence is split into two simultaneously imaged
597 channels, with a small (500-1000 nm) defocus introduced between the two channels¹⁵. As the small
598 defocus should introduce minimal additional aberration into an optical system, we semi-
599 synthetically constructed a realistic biplane PSF from the experimental 2D PSF. The two defocused
600 PSFs were constructed by duplicating the 2D PSF and offsetting it by -250 nm and 250 nm for each Z-
601 plane.

602 This yielded five high SNR model PSFs with an isotropic voxel size of 10x10x10 nm³.

603 The ground truth XY=0 was defined as the image center of mass of the in-focus frame of the model
604 PSF, and Z=0 was defined as the in-focus frame. Accounts for shifts in the fitted XY center of the
605 model PSF by localization software due to systematic offsets and Z-dependent variation of the model
606 PSF center of mass are dealt with below (wobble correction).

607 **2.5. Noise model**

608 A constant mean autofluorescent background was added to the noise-free simulated images, and
609 these images were then fed through the noise model representing Poisson distributed fluorescence
610 emission recorded on a high quantum efficiency back-illuminated EMCCD^{45,46}.

611 The proposed noise model assumed as main contributions to the stochastic noise:

- 612
- 613 • σ_S , the shot noise produced by the fluorescence background and signal and the spurious charge. Shot noise can be derived from the second moment of the Poisson distribution

- 614 • σ_R , the read noise of EMCCD camera, which is described by second moment of the Gaussian
615 distribution
616 • σ_{EM} , the electron multiplication noise introduced by the gain process, which is described by
617 the second moment of the Gamma distribution⁴⁶.
618

619 We assumed as camera parameters the ones specified for the Photometrics Evolve Delta 512 EMCCD
620 camera (values for other manufacturer's EMCCDs are similar):

- 621 • QE = 0.9, Evolve quantum efficiency at 700 nm absorption wavelength.
622 • $\sigma_R = 74.4$ electrons, manufacturer measured root mean square noise for Evolve 512 camera
623 • $c = 0.002$ electrons, manufacturer quoted spurious charge (clock induced charge only, dark
624 counts negligible)
625 • $EM_{gain} = 300$
626 • $e_{adu} = 45$ electron per analog to digital unit (ADU), analog to digital conversion factor
627 • $G = 0.9 \cdot 300 / 45 = 6$, total system gain
628 • BL = 100 ADU

629 The final simulated photon electrons will thus be given by:

$$n_{ie} = \mathcal{P}(QE \cdot n_{photn} + c)$$

$$n_{oe} = \Gamma(n_{ie}, EM_{gain}) + \mathcal{G}(0, \sigma_R)$$

630 which leads to the final pixel counts:

$$ADU_{out} = \min\left(\frac{n_{oe} - n_{oe} \bmod e_{ADU}}{e_{peradu}} + BL, 65535\right)$$

631 **2.6. Depth-dependent lateral distortion/ wobble**

632 As the PSF models are experimentally derived, the 3D estimated localizations exhibit a depth-
633 dependent lateral distortion, here called *wobble*. This optical distortion is due to a combination of a
634 systematic offset (arbitrary definition of PSF center) and optical aberrations⁴⁷. In order to compare
635 estimated and true localizations, we correct this effect during the assessment (**Methods 3.1**).

636 **2.7 Comparison of software results between different modalities.**

637 The intensities of the PSF in each imaging modality were normalized to facilitate comparison of
638 results between different modalities. Software results between 2D, 3D AS and 3D DH modalities are
639 expected to be directly comparable.

640 For the biplane model PSF, as the emitted fluorescence is split into two channels, the intensity in
641 each of the two simulated biplane channels was additionally reduced by 50 %. We note that a
642 simulation bug meant that the fluorescence background was not reduced by 50 % as intended,
643 leading to artificially high background for the biplane simulation. *I.e.*, the background in each of the
644 two biplane channels is the same as in the single channel of the other modalities. However, due to
645 the low background level in the 3D simulations, the effect on image SNR and thus localization error
646 is small (see **Fig. S5, S6**), less than 5 nm near the plane of focus. Therefore, as long as the small drop
647 in image SNR is taken into account, approximate comparisons of the biplane data to the other
648 modalities can still be made.

649 **3. SOFTWARE ASSESSMENT**

650 **3.1 Protocol**

651 Each localization file submitted by the participants was manually checked for erroneous systematic
652 errors in the definition of the dataset coordinate system, such as offsets, XY axis flips or clear scaling

653 errors. Datasets were then programmatically standardized into a consistent output format. All
654 modifications are publicly available. If required, the modifications consisted of columns reordering,
655 reversing axes, XY axis swap, and shifting the lateral positions by a half camera pixel.

656 The assessment pipeline includes three main parts: localization processing, the pairing between true
657 and estimated localization and the metrics calculations. The first one depends on the assessment
658 settings. There are two switchable properties: photon thresholding and wobble correction. Their
659 combinations yield four different assessment settings. Up to 64 assessment runs per software were
660 possible (*i.e.*, 4 modalities, 4 datasets per modality). For any setting, we excluded the fluorophores
661 within a lateral distance of 450 nm from the border. This value corresponds to the radius of the
662 largest PSF, *i.e.*, Double Helix. The activations too close from the border are more difficult to localize
663 and could bias the results.

664 The pairing between true and estimated localizations was performed frame by frame. For every
665 frame, we identified the localizations that are close enough to a ground-truth position as true-
666 positives (TP), the spurious localizations as false-positives (FP) and the undetected molecules as
667 false-negatives (FN). The procedure matches two sets of localizations. We deployed the presorted
668 nearest-neighbor search for its efficiency, with a linking threshold of 250 nm. The results are
669 effectively similar to the computationally intensive Hungarian algorithm⁷.

670 *Photon thresholding*

671 A photon threshold was required primarily due to the use of a realistic fluorophore blinking model.
672 Since a fluorophore could activate/ bleach at any point in a simulated frame, this led to many frames
673 containing very dim, undetectable localizations, *e.g.*, where a molecule had been active for one or
674 more frames previously, and then bleached during the first 5 % of a frame. These fractional
675 localizations should also be present but practically undetectable in an experimental dataset.

676 We decided to focus the software analysis on the localizations where the molecule was active for the
677 majority of a frame, to be consistent with experimental expectations. Therefore, we implemented a
678 photon threshold means where we kept the 75% brightest ground truth fluorophore activations.
679 Because this was performed *after* the pairing step, observed localizations that were paired to
680 discarded ground truth activations were also removed from the metric calculations.

681 *Wobble correction*

682 The centroid of experimental point spread functions shifts laterally by as much as 50 nm, as a
683 function of axial position^{10,47}. This is most often ignored by localization software, and instead
684 corrected post-hoc by reference to a calibration curve³⁷. Since our simulated PSF is experimentally
685 derived, it was necessary to correct for these artefactual shifts between the observed localizations
686 and ground truth, as part of the assessment process. This correction was performed using calibration
687 data uploaded by competitors, similar to the correction typically performed on experimental data⁴⁷.

688 Three scenarios were proposed to the participants: no correction was applied during the
689 assessment; the correction was based on a file provided by the participant itself or the correction
690 was calculated by ourselves. The latter nevertheless requires the participant to localize a stack of
691 beads we provided. Since the true positions of the beads are known, the difference between the
692 estimated and true positions could be calculated and averaged. It thus yields the values for wobble
693 correction.

694 In certain specific cases (identified on the competition website), at the request of authors, we did
695 not apply this correction, for example because the software explicitly considered the whole 3D PSF
696 during fitting and was thus immune to this lateral shift artefact. For accurate results, application of
697 lateral shift correction is critical for analysis of localization microscopy simulations using

698 experimentally derived PSFs, as can be seen by comparison of typical software results with and
699 without wobble correction (**Fig. S19**).

700 **3.2 Metrics**

701 We calculated a large number of analysis metrics to quantify the performance of software relative to
702 ground truth. These are discussed in detail in **Supplementary Note 2**. The metrics are split into two
703 categories: localization based and image based metrics.

704 *Localization based metrics.* This directly relies on the localizations positions and notably includes the
705 Recall, the Precision, the Jaccard Index, the RMSE (axial and lateral) and the consolidated Z-range.
706 For the calculation of average software performance (**Fig. 3D-F, S10**) outlier software with an
707 efficiency less than $eff=0$ ($eff=-30$ for 3D high density dataset) were excluded from the
708 measurement. The key metrics of assessment were:

709 1. *Root mean squared localization error (RMSE).* The foremost consideration for localization
710 software is how accurately it finds the position of labelled molecules. This was quantified as
711 the root mean squared difference between the measured molecule position, x_i^s , and the
712 ground truth position, x_i^t , in both the lateral (XY) and axial (Z) dimensions.

713
$$RMSE \text{ lateral (RMSE Lateral) [nm]: } \sqrt{\frac{1}{TP} \sum_{i \in SN_T} (x_i^s - x_i^t)^2 + (y_i^s - y_i^t)^2}.$$

714
$$RMSE \text{ axial (RMSE Axial) [nm]: } \sqrt{\frac{1}{TP} \sum_{i \in SN_T} (z_i^s - z_i^t)^2}.$$

715 2. *Jaccard index (JAC, %).* In addition to localization precision, SMLM image resolution depends
716 critically on number of localized molecules⁴⁸, so it is crucial for SMLM software to accurately
717 detect a large fraction of molecules in a dataset, and minimize false localizations. For every
718 frame, we identified the localizations that are close enough to a ground-truth position as
719 true-positives (TP), the spurious localizations as false-positives (FP) and the undetected
720 molecules as false-negatives (FN). We then computed the *Jaccard index* (JAC, %), which
721 measures the fraction of correctly detected molecules in a dataset,

$$JAC = 100 \frac{TP}{TP + FP + FN}$$

722 3. *Efficiency (E).* For ranking purposes, we developed a single summary statistic for overall
723 evaluation of software performance, which we term the *efficiency (E)*, encapsulating both
724 the software's ability to find molecules, measured by the Jaccard index, and the software's
725 ability to precisely localize molecules.

$$E = 100 - \sqrt{(100 - JAC)^2 + \alpha^2 RMSE^2}$$

726 The trade-off between these two metrics is controlled by a parameter α . In a retrospective
727 analysis, we chose $\alpha = 1 \text{ nm}^{-1}$ for the lateral efficiency E_{lat} , $\alpha = 0.5 \text{ nm}^{-1}$ for the axial efficiency
728 E_{ax} , based on the linear regression slope between the localization errors and Jaccard index
729 (**Fig. S20J-K**). Using this definition, an average software performance has an efficiency in the
730 range 25-75, a perfect software would have the maximum efficiency of 100. Overall 3D
731 efficiency was calculated as the average of lateral and axial efficiencies. Overall software
732 rankings (**Fig. 2**) were calculated as the sum of rankings for high and low SNR datasets.

733 *Image based metrics.* The image based metrics are computed from a rendered image and includes
734 the Signal-to-Noise Ratio (SNR) and the Fourier Ring / Shell Correlation (FRC/FSC). To render the
735 image, we added the contribution of each localized molecule at the corresponding pixels. A
736 contribution takes the form of a 3D additive Gaussian with a Full-Width Half Maximum (FWHM) of
737 20 nm. A complete list of all computed metrics is presented in the **Supplementary Note 2**.

738 We also calculated localization based metric results as a function of axial position. We proceeded by
739 considering a subset of activations lying within an interval of axial positions (*i.e.*, from the true

740 localizations). Then, most of the metrics (*e.g.*, Recall) are locally computed. This yields a curve
741 providing information on the depth performance of each software / modality.

742 In order to summarize software axial performance, we analyzed how the recall varied as a function
743 of Z. A typical recall versus axial position curve (**Fig. S4**) will drop at positions far from the focal
744 plane, *i.e.*, where software can no longer detect spots to defocus. We first smoothed the curve using
745 a sliding window. Then we computed the software Z-range, defined as the full width half maximal
746 Recall of the smoothed curve (**Fig. S21**). This quantity is visually intuitive and useful for discussion of
747 the recall performance if considered alongside a plot of recall vs axial position. However, because
748 FWHM recall depends on the maximal recall, ranking based on this procedure would promote a
749 software which poorly performed everywhere (*i.e.*, flat curve), whereas a software which performed
750 well in the focal plane but less well outside would obtain a worse FWHM recall. This observation
751 leads us to produce a so-called consolidated Z-range, by multiplying the Z-range value by the
752 maximal Recall, which should provide a robust metric that avoids the previous case scenario.

753 *Principal component analysis.* In order to analyse the relationship between analysis metrics we
754 computed the covariance matrix between each metric (**Fig. S22A**) and the principal component
755 analysis (PCA) on the metrics (**Fig. S22B-D**). Each metric was standardized before applying the
756 covariance and the PCA. For convenience, we took the additive inverse of the metrics for which
757 lower values are best (*i.e.*, FP, FN, RMSE, FRC, FSC).

758 Summary statistics and detailed results for each software are available on the competition website
759 (<http://bigwww.epfl.ch/smlm/challenge2016/index.html?p=results>), which also includes a tool for
760 side-by-side comparison of the results of multiple software packages

761 **3.3 Baseline Localization Software**

762 We developed a minimalist Java tool software that performs localizations of bright emitters on the 4
763 modalities of the challenge 2016: 2D, Astigmatism, Double-Helix, and Biplane. This
764 SMLM_BaselineLocalization software is only designed to establish the performance baseline for the
765 SMLM challenge. It has intentionally limited lines of code and relies only on few threshold
766 parameters to localize particles. It has basic calibration tool that has to run on a z-stack of beads to
767 find the linear $f(x)$ relation between the axial position Z and the shape of the bead.

- 768 • Astigmatism: $Z = f(W_x - W_y)$, where W_x and W_y are respectively an estimation of the size in X
769 and Y.
- 770 • Double-Helix: $Z = f(\theta)$, where θ is the angle formed the pairing of two close points.
- 771 • Biplane: $Z = f(W_{\text{left}} - W_{\text{right}})$, where W_{left} and W_{right} are respectively an estimation of the size of
772 the spots in left and the right plane.

773 The Java code is available: <https://github.com/SMLM-Challenge/Challenge2016>

774 **4 REAL DATA ASSESSMENT**

775 Astigmatism software was tested on previously published real 3D STORM datasets of microtubules
776 and nuclear pore complex¹⁹. The tubulin dataset corresponds to the raw data for **Fig. S6** in Ref ¹⁹,
777 and the nuclear pore complex dataset corresponds to raw data for **Fig. S9** in Ref ¹⁹. Key acquisition
778 parameters for data analysis are summarized on the competition website.

779 Data were analyzed by software authors or expert users, and submitted via the competition website.
780 All data were drift corrected via cross-correlation. STORM images were rendered with a constant
781 Gaussian blur with 3 nm standard deviation and saturated by 0.1 – 0.5 %. The complete scripts used
782 for assessment and image rendering are available on the competition GitHub page.

783 5 DATA AVAILABILITY

784 5.1 Data availability statement

785 Simulated competition datasets are available at <http://bigwww.epfl.ch/smlm/challenge2016/>,
786 together with the parameters used to generate the data. The ground truth list of simulated molecule
787 positions for each competition dataset remains secret in order to allow the software challenge to
788 remain continuously open to new submissions. However, ground truth data are available for the
789 simulated training datasets.

790 Raw data for this study are uploaded on the Nature Methods website. The data corresponding to
791 specific figures are listed with the Supplementary information.

792 5.2 Code availability statement

793 All software is available at <https://github.com/SMLM-Challenge/Challenge2016>

794 REFERENCES, ONLINE METHODS

- 795 40. Carlini, L. & Manley, S. Live Intracellular Super-Resolution Imaging Using Site-Specific Stains.
796 *ACS Chem. Biol.* **8**, 2643–2648 (2013).
- 797 41. Shim, S.-H. *et al.* Super-resolution fluorescence imaging of organelles in live cells with
798 photoswitchable membrane probes. *Proc. Natl. Acad. Sci.* **109**, 13978–13983 (2012).
- 799 42. Hanser B. M., Gustafsson M. G. L., Agard D. A. & Sedat J. W. Phase-retrieved pupil functions
800 in wide-field fluorescence microscopy. *J. Microsc.* **216**, 32–48 (2004).
- 801 43. Izeddin, I. *et al.* PSF shaping using adaptive optics for three-dimensional single-molecule
802 super-resolution imaging and tracking. *Opt. Express* **20**, 4957–4967 (2012).
- 803 44. McGorty, R., Schnitzbauer, J., Zhang, W. & Huang, B. Correction of depth-dependent
804 aberrations in 3D single-molecule localization and super-resolution microscopy. *Opt. Lett.* **39**, 275–
805 278 (2014).
- 806 45. Hirsch, M., Wareham, R. J., Martin-Fernandez, M. L., Hobson, M. P. & Rolfe, D. J. A Stochastic
807 Model for Electron Multiplication Charge-Coupled Devices – From Theory to Practice. *PLOS ONE* **8**,
808 e53671 (2013).
- 809 46. Basden, A. G., Haniff, C. A. & Mackay, C. D. Photon counting strategies with low-light-level
810 CCDs. *Mon. Not. R. Astron. Soc.* **345**, 985–991 (2003).
- 811 47. Carlini, L., Holden, S. J., Douglass, K. M. & Manley, S. Correction of a Depth-Dependent
812 Lateral Distortion in 3D Super-Resolution Imaging. *PLoS ONE* **10**, e0142949 (2015).
- 813 48. Baddeley, D. & Bewersdorf, J. Biological Insight from Super-Resolution Microscopy: What We
814 Can Learn from Localization-Based Images. *Annu. Rev. Biochem.* **87**, 965–989 (2018).
- 815

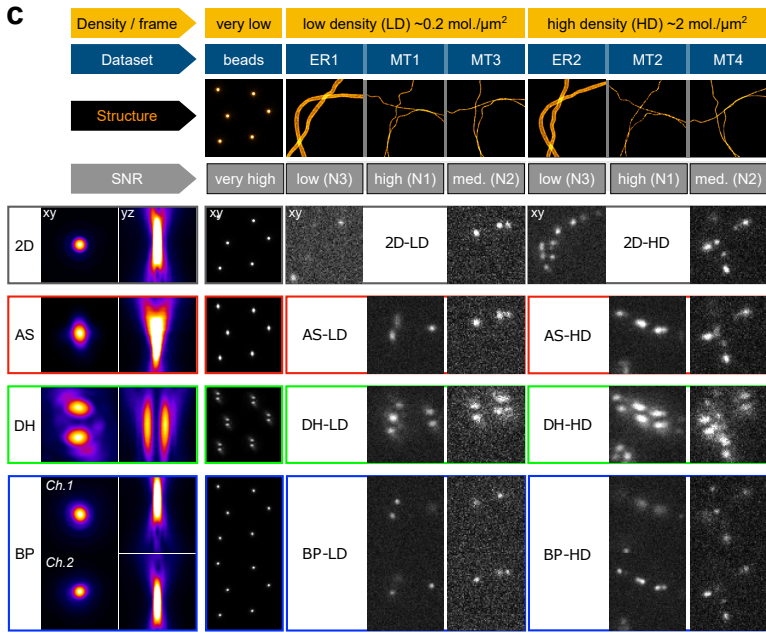
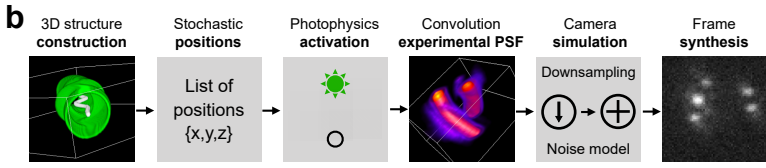
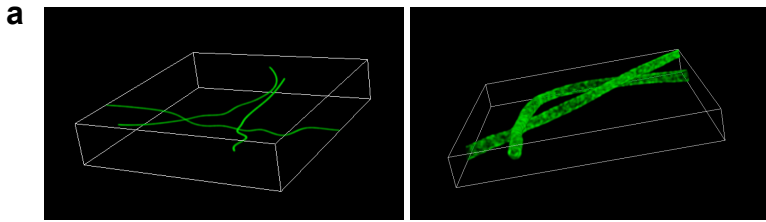
816 **FIGURE LEGENDS**

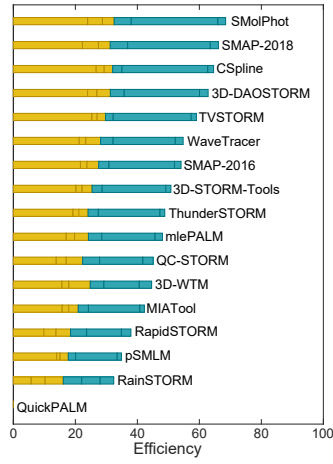
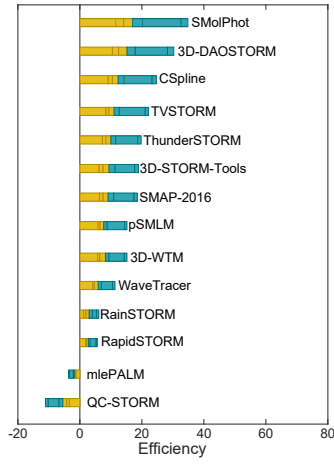
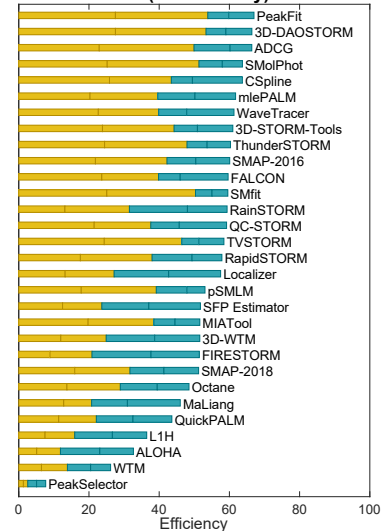
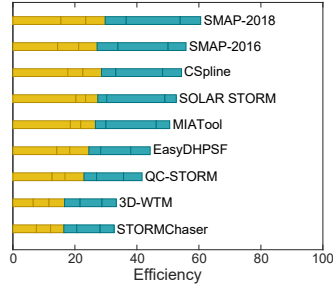
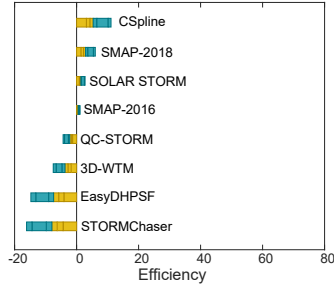
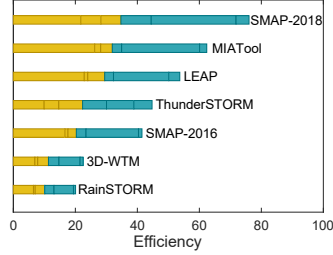
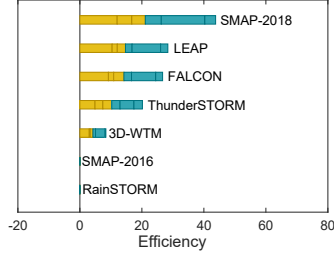
817 **Figure 1: Summary of SMLM challenge simulations.** **A.** 3D rendering of simulated microtubules and
818 endoplasmic reticulum samples. **B. Key simulation steps.** The structure is constructed from 3D tubes
819 continuously defined by three B-spline functions in the volume of interest. Membranes of the tubes
820 are densely populated with possible positions. Fluorophores follow a 4-state photophysics model.
821 Activations of a given frame are convolved with the experimental PSF and shot & camera noise is
822 added. **C.** Summary of all 16 challenge datasets, calibration data and experimental PSFs. Left column:
823 orthogonal projections of the experimentally-derived PSF. Right column: exemplar frame for each
824 competition dataset, characterized by structure (endoplasmic reticulum, E; microtubules, MT),
825 modality (2D; astigmatism, AS; double helix, DH; biplane, BP), density (low density, LD; high density,
826 HD) and SNR (noise level N1, N2, N3). *BP Ch. 1,2*, indicates two biplane channels with a relative focal
827 shift of 500 nm.

828 **Figure 2: Leaderboards for each competition modality, at low and high spot density.** Ranking is based
829 on software Efficiency, which combines Jaccard index (fraction of successfully detected molecules)
830 and localization precision (RMSE, root mean square error, lateral & axial). Orange, contribution of
831 high SNR dataset; blue, contribution of low SNR dataset.

832 **Figure 3: Comparison of 3D software performance.** Gold stars indicate top performers for each
833 dataset. Dashed lines in top, middle panels indicate overall efficiency (higher is better). **A-C.**
834 Localization error and spot detection performance of all astigmatic SMLM software. **D-E.** Average
835 (colored marker with *s.d.* error bars, sample sizes for each category indicated in **Supplementary**
836 **Table 2**) and best-in-class (colored marker with gold star) software performance for all competition
837 modalities. *AS, astigmatism; DH, double helix; BP, biplane.*

838 **Figure 4: Super-resolved images of software results for simulated and real competition datasets.** **A.**
839 *Xy and xz projection images of 3D competition datasets for representative software. Top: best-in-*
840 *class software in each modality, for high SNR low density dataset. Bottom: representative average*
841 *software. Left: xy and xz overview images for winning AS software. Middle: xy and xz zoom images of*
842 *boxed regions in left panel, for winning and mid-range software, each modality. Right: xy and xz line*
843 *profiles of winning and mid-range software for each modality, for boxed regions in middle panel.*
844 *Image colors: red, ground truth; green, software results. Line profiles: GT, ground truth, black; AS,*
845 *astigmatism, red; BP, biplane, blue; DH, double helix, green. Panel key: Software-name Dataset-*
846 *ranking°. Scale bar: full image, 1 μ m, magnified regions, 100 nm. **B. Astigmatism software results for**
847 *real nuclear pore complex 3D STORM data. Top: Super-resolved overview image in xy for 3D-*
848 *DAOSTORM software, color coded for depth. Bottom: xz orthoslices along 600 nm wide dashed*
849 *region indicated in top panel for 8 astigmatism software packages. Scale bars, 500 nm.**



a Astigmatism (Low Density)**Astigmatism (High Density)****d 2D (Low Density)****b Double-Helix (Low Density)****Double-Helix (High Density)****c Biplane (Low Density)****Biplane (High Density)****2D (High Density)**