

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

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

7 \*Corresponding authors: [daniel.sage@epfl.ch](mailto:daniel.sage@epfl.ch), [seamus.holden@ncl.ac.uk](mailto: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 (SMLM<sup>1-3</sup>). 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 spots.  
56 Improvements in SMLM image processing have been crucial in maximizing spatial resolution and  
57 reducing imaging time of SMLM for compatibly with live cell imaging<sup>4-6</sup>. 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 software<sup>7</sup>. 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<sup>8,9</sup>. 3D localization microscopy is a more difficult  
64 image processing problem than 2D SMLM. In addition to finding the center of diffraction limited spots  
65 to super-resolve lateral position, 3D SMLM algorithms must also extract axial information from the  
66 image, usually by measuring small changes in the shape of a point spread function<sup>10</sup> (PSF).

67 Despite the widespread use of 3D localization microscopy, and challenging nature of 3D SMLM image  
68 processing, the performance of software for 3D single molecule localization microscopy has previously  
69 only been assessed for 2-3 software packages at a time, and without standard test data or metrics<sup>11-</sup>  
70 <sup>14</sup>. In the absence of common reference datasets and reliable assessment, it is not possible to  
71 objectively assess how different software affects final image quality, or which algorithmic approaches  
72 are most successful. Crucially, end-users cannot determine which 3D SMLM software package and  
73 imaging modality is optimal for their application.

74 We therefore ran the first 3D localization microscopy software challenge, to assess the performance  
75 of 3D SMLM software. We assessed software performance on simulated datasets designed for  
76 maximum realism, incorporating experimentally derived point spread functions, using biologically  
77 inspired structures, signal to noise levels based closely on common experimental conditions, and  
78 modelling fluorophore photophysics. We assessed software performance on synthetic datasets for  
79 three popular 3D SMLM modalities: astigmatic imaging<sup>10</sup>, biplane imaging<sup>15</sup> and double helix point  
80 spread function microscopy<sup>16</sup>. We also assessed astigmatism software performance on two real  
81 STORM datasets. Furthermore, we ran a second 2D localization microscopy software challenge to  
82 assess performance of the latest 2D SMLM software.

## 83 RESULTS

### 84 Competition design

85 We established a broad committee from the SMLM community, including experimentalists and  
86 software developers, to define the scope of the challenge, ensure realism of the datasets and define  
87 analysis metrics. We opened this discussion to all interested parties in an online discussion forum<sup>17</sup>.

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

### 95 Realistic 3D simulations

96 Testing super-resolution software on experimental data lacks the ground truth information required  
97 for rigorous quantification of software performance. Therefore, realistic simulated datasets are  
98 required. A critical challenge to in simulating 3D SMLM data was to accurately model the experimental

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

104 We thus combined experimental 3D PSFs with simulated ground truth by performing simulations using  
105 PSFs directly derived from experimental calibration data (**Fig. 1, Methods**). We generated simulated  
106 datasets over a range of spot densities and signal to noise levels, for simulated microtubule- and  
107 endoplasmic reticulum-like structures, using a 4-state model for photophysics<sup>18</sup> (**Methods**).

## 108 **Quantitative performance assessment of 3D software**

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

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

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

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

## 125 **Performance of 3D software**

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

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

### 130 *Astigmatic localization microscopy*

131 Astigmatic localization microscopy is probably the most popular 3D SMLM modality, reflected by the  
132 highest number of software submissions in the 3D competition (**Fig. 2**). For astigmatism, we observed  
133 a large spread of software performance, even for the most straightforward high SNR, low spot density  
134 (LD) conditions (**Fig. 3, Supplementary Table 2**). The best-in-class software (SMAP-2018<sup>19</sup>) has  
135 significantly better localization error and Jaccard index performance than average (lateral RMSE 26 nm  
136 best vs 38 nm average, axial RMSE 29 nm best vs 66 nm average, Jaccard index 85 % best vs 74 %  
137 average). Clearly, the quality of the image reconstruction depends strongly on choice of 3D software.

138 To investigate the reasons for software variation, we inspected plots of software performance as a  
139 function of axial position in the low density, high SNR dataset for best-in-class and representative  
140 middle-range software (**Fig. S4A**). We observed that a key cause of the spread in software  
141 performance is variation in software performance away from the focal plane. Near the focal plane,  
142 most software packages perform well. However, the axial and lateral RMSE away from the plane of  
143 focus is significantly higher for the best in class software, and the Jaccard index is also slightly improved

144 (Fig. S4A). This is also visibly apparent in the super-resolved images (Fig. 4A). We observed that best-  
145 in-class software had a Z-range (the FWHM range of axially resolved software recall, Methods) of  
146 1170 nm, greater than two-thirds of the simulated range. Outside this range, the recall and Jaccard  
147 index dropped sharply, probably due the large increase in PSF size and decrease in effective SNR at  
148 large defocus (Fig. S1).

149 When we examined results for the low SNR, low density dataset (Fig. 2A, 3F), we found an expected  
150 two-fold degradation in best-in-class RMSE (lateral RMSE 39 nm, axial RMSE 60 nm), due to the  
151 decrease in image SNR. However, the best-in-class software (SMolPhot<sup>20</sup>) Jaccard index was  
152 effectively constant between the low and high SNR datasets (86 % vs 85 %), although the Z-range did  
153 drop at lower SNR (930 nm vs 1120 nm). The best astigmatism software packages were thus  
154 remarkably good at finding spots at low SNR, even away from the focal plane.

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

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

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

#### 174 *Double helix point spread function localization microscopy*

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

186 Double helix software performance decreased significantly for the low spot density low SNR condition  
187 (best software, SMAP-2018), particularly in terms of best-in-class Jaccard index (66 % low SNR vs 77 %  
188 high SNR, Fig. 3D-E, S8, S9A). DH Jaccard index was also significantly worse than astigmatism results  
189 at either high or low SNR (85 % high SNR, 86 % low SNR). This indicates that it was quite hard to  
190 successfully find localizations in the low SNR DH dataset, likely because the large size of the DH PSF

191 spreads emitted photons over a large area, lowering effective image SNR. DH PSF designs with reduced  
192 Z-range but more compact PSF would likely be less sensitive to this issue<sup>21</sup>.

193 Double helix software performed poorly on the high spot density datasets at high SNR (best software  
194 CSpline<sup>22</sup>), especially in terms of the Jaccard index (**Fig. 3D-E, S9A**, best lateral RMSE 67 nm, best axial  
195 RMSE 69 nm, best Jaccard index 46 %). The poor performance at high spot density is again probably  
196 because the large DH PSF size increases spot density and decreases SNR (**Fig. S1**). DHPSF performance  
197 at high spot density and low SNR was also not reliable (**Fig. 3D-F, S9A**, best software, SMAP-2018).

### 198 *Biplane localization microscopy*

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

210 At high spot density and high SNR, best-in-class biplane software (SMAP-2018) showed acceptable  
211 performance (**Fig. 3D-F, S7, S9B**, best lateral RMSE 43 nm, best axial RMSE 49 nm, best Jaccard index  
212 61 %). Uniquely among the 3D modalities, best-in-class biplane software also gave acceptable  
213 performance at high spot density and low SNR (**Fig. 3D-F, S7, S9B**, best lateral RMSE 55 nm, best axial  
214 RMSE 72 nm, best Jaccard index 61 %, best software SMAP-2018).

### 215 **Performance of 2D software**

216 We next assessed the performance of 2D SMLM software. For the pseudo-ER 2D dataset, at low  
217 density best-in-class software (ADCG<sup>23</sup>) performed substantially better than the class average  
218 (**Fig. S10, S11**, lateral RMSE 31 nm vs 36 nm average, Jaccard index 90 % best vs 72 %). Low density  
219 results for the brighter fluorophore microtubules dataset were similar to the dimmer pseudo-ER  
220 dataset (**Fig. S10, S12** best software SMolPhot). For the very high density 2D dataset, which had 25x  
221 higher spot density than the LD dataset, best-in-class software (ADCG) showed excellent performance  
222 (**Fig. S10**, lateral RMSE, 45.5 nm, Jaccard index 75%). Best-in-class performance (ADCG) on the dimmer  
223 fluorophore data at high spot density was also strong (**Fig. S10**, best lateral RMSE 51 nm, best Jaccard  
224 index 70 %).

### 225 **Algorithms**

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

227 1) *Non-iterative* software regroups pixels in the local neighborhood of the candidates, like  
228 interpolation, center of mass (QuickPALM<sup>24</sup>) or template matching (WTM<sup>25</sup>). These often older  
229 algorithms are fast but tend to achieve poor performance.

230 2) *Single emitter fitting* software is usually built on a multi-step strategy of detection, spot localization,  
231 and optional spot rejection. The detection step finds bright spots in noisy images on the pixel grid. The  
232 selection of candidates is usually performed by local maximum search after a denoising filter. Others  
233 rely on more complex algorithms like the wavelet transform (WaveTracer<sup>26</sup>). We did not observe  
234 software ranking to depend noticeably on the choice of optimization scheme: least-square, weighted  
235 least-square or maximum-likelihood estimator.

236 3) *Multi-emitter fitting* software groups clusters of overlapping spots, and simultaneously fits  
237 multiple model PSFs to the data. Typically, fitted spots are added to the cluster until a stopping  
238 condition is met<sup>4,5</sup>. This leads to improved localization performance at high spot density, at the cost  
239 of reduced speed. This class of software (e.g., 3D-DAOSTORM<sup>11</sup>, CSpline, PeakFit, ThunderSTORM<sup>27</sup>)  
240 was amongst the top performers in each 2D and 3D competition category.

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

245 4) *Compressed sensing algorithms*. One subset of these algorithms utilize deconvolution with sparsity  
246 constraints to reconstruct super-resolved images<sup>28-30</sup>. Although deconvolution approaches can give  
247 good results, they are limited by the necessary use of a sub-pixel grid; increased localization precision  
248 requires smaller grid resolution, which must be balanced against increased computational time.  
249 Recent approaches address this issue by localizing the point sources in a gridless manner under some  
250 sparsity constraint (ADCG, SMfit, SOLAR\_STORM, TVSTORM<sup>31</sup>). This software class consistently gave  
251 the overall best performance for 2D high-density (ADCG 1<sup>st</sup>, FALCON<sup>30</sup> 2<sup>nd</sup>, SMfit 3<sup>rd</sup>).

252 5) *Other approaches*. Of the alternative algorithmic approaches used, the annihilating filter-based  
253 method LEAP<sup>32</sup> gave good performance for biplane imaging. Recently, we received the first challenge  
254 submission from a deep learning SMLM software (DECODE); these promising preliminary results are  
255 available on the competition website.

#### 256 *Post-hoc temporal grouping*

257 Because molecule on-time is stochastically distributed across multiple frames, a common post-  
258 processing approach to improve localization precision is to group molecules detected multiple times  
259 in adjacent frames, and average their position<sup>33</sup> (**Supplementary Note 4**). Temporal grouping was used  
260 by the top performers (including SMolPhot, MIATool<sup>34</sup> and SMAP-2018), and is visibly apparent as a  
261 more punctate super-resolved image (**Fig. 4A**).

#### 262 *Choice of PSF model*

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

#### 269 *Multi-algorithm packages*

270 Several software packages take a Swiss army knife approach of integrating multiple optional  
271 localization algorithms into one program, to be flexible enough to suit various experimental  
272 conditions<sup>19,27</sup>. SMAP-2018 and ThunderSTORM achieved strong across-the-board performance  
273 supporting this rationale.

#### 274 *Software run time*

275 Software run time is important both for ease of use and real time analysis. We did not observe  
276 correlation between software localization performance (Efficiency) and software run time (**Fig. S13A**).  
277 We thus created an alternative ranking metric, *Efficiency-Runtime*, which gave 25 % weighting to run  
278 time (**Supplementary Note 2.7, Fig S13B**). Many good performers in the efficiency-only ranking were  
279 relatively fast and thus retained good ranking (SMAP-2018, SMolPhot, 3D-DAOSTORM). Interestingly,

280 two software packages highly optimized for speed gained top ranking in this analysis: pSMLM-3D<sup>35</sup>  
281 and QC-STORM.

### 282 *Diagnostic tools for software and algorithm performance*

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

### 290 **Assessment on real STORM data**

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

## 302 **DISCUSSION**

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

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

319 Of the different algorithm classes, well-tuned single-emitter and multi-emitter fitting algorithms (each  
320 capable of dealing well with occasional molecule overlap) gave good results for low density 3D SMLM.  
321 We also found that several software packages for astigmatic or biplane imaging gave adequate  
322 performance for the challenging case of high molecule densities, as long as the image SNR was high.  
323 Current software packages gave poor performance when molecule density was high and image SNR  
324 was low. These results indicate that with current algorithms high density 3D SMLM performance is  
325 mediocre at high SNR and poor at low SNR. Surprisingly, multi-emitter fitting did not show significant



326 improvement over well-tuned single emitter fitting for the 3D high-density datasets; this may indicate  
327 that significant potential for improvement remains in this category.

328 Many software packages did not apply temporal grouping<sup>33</sup>, resulting in reduced software  
329 performance. Since temporal grouping is a simple step for maximum precision, we urge all software  
330 developers to integrate this approach into their software as an optional final step in the localization  
331 process.

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

339 In future, we plan to extend the SMLM challenge into an open platform with a fully automated  
340 assessment process, and where new competition simulations and assessment metrics can easily be  
341 created and contributed by the community. It will be important to account for new technologies and  
342 developments in SMLM, such as scientific CMOS cameras<sup>6</sup>, in future simulations. It would also be  
343 exciting to adapt the tools developed in the SMLM challenge to other classes of super-resolution  
344 microscopy, such as fluorescence-fluctuation-based super-resolution microscopies (*e.g.*, 3B<sup>36</sup>, SOFI<sup>37</sup>,  
345 SRRF<sup>38</sup>) and structured illumination microscopy<sup>39</sup>.

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

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

357

## 358 ACKNOWLEDGEMENTS

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

## 386 AUTHOR CONTRIBUTIONS

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

## 394 REFERENCES

- 395 1. Betzig, E. *et al.* Imaging Intracellular Fluorescent Proteins at Nanometer Resolution. *Science*  
396 **313**, 1642–1645 (2006).
- 397 2. Hess, S. T., Girirajan, T. P. K. & Mason, M. D. Ultra-High Resolution Imaging by Fluorescence  
398 Photoactivation Localization Microscopy. *Biophys. J.* **91**, 4258–4272 (2006).
- 399 3. Rust, M. J., Bates, M. & Zhuang, X. Sub-diffraction-limit imaging by stochastic optical  
400 reconstruction microscopy (STORM). *Nat Methods* **3**, 793–795 (2006).
- 401 4. Holden, S. J., Uphoff, S. & Kapanidis, A. N. DAOSTORM: an algorithm for high- density super-  
402 resolution microscopy. *Nat Meth* **8**, 279–280 (2011).
- 403 5. Huang, F., Schwartz, S. L., Byars, J. M. & Lidke, K. A. Simultaneous multiple-emitter fitting for  
404 single molecule super-resolution imaging. *Biomed. Opt. Express* **2**, 1377–1393 (2011).

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

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

## 481 METHODS

### 482 1. CHALLENGE ORGANIZATION

483 We first ran the 3D SMLM software challenge as a time limited competition, with a results session  
484 hosted as a special session of the 6<sup>th</sup> Annual Single Molecule Localization Microscopy Symposium in  
485 August 2016. The competition has now been converted to a permanent software challenge accepting  
486 new submissions. Special thanks is due to the software SMAP and 3D-WTM<sup>25</sup> that participated in all  
487 eight categories (*density x modality*). The current list of participants is at:

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

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

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

### 493 2. LOCALIZATION MICROSCOPY SIMULATIONS

#### 494 2.1. Structure, noise levels and spot densities

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

505 The underlying sample structure is formalized in a continuous space which allows rendering of digital  
506 images at any scale, from very high resolution (up to 1 nm/pixel) to low resolution (camera resolution:  
507 100 nm/ pixel). The continuous-domain 3D curve is represented by means of a polynomial spline. The  
508 sample is imaged in a  $6.4 \times 6.4 \mu\text{m}^2$  field of view, and the center lines of the microtubules have limited  
509 variation along the *z* (vertical) axis, *i.e.*, less than 1.5  $\mu\text{m}$ . The fluorescent markers are uniform  
510 randomly distributed over the structure according to the required density. The photon emission rate  
511 of each fluorophore is controlled by a photo-activation model (see below). The exact locations of all  
512 fluorophores are stored at high precision floating-point numbers expressed in nanometers. This  
513 ground-truth file is used for conducting objective evaluations without human bias.

514 *Noise levels.* We generated data at three different signal-to-noise ratio (SNR) levels, based on real  
515 signal to noise levels encountered under common SMLM experimental scenarios: *N1*, fixed cells  
516 antibody labelled with organic dye<sup>10</sup>, high signal, medium background; *N2*, fluorescent protein  
517 labelling<sup>1</sup>, low signal, low background; and *N3*, live cell affinity dye labelling<sup>40,41</sup>, high signal, high  
518 background.

519 *Spot density.* As performance at different density of active emitters is a key challenge for SMLM  
520 software, we generated 3D competition datasets at both sparse emitter density  
521 ( $0.25 \text{ mol. [molecule]} \mu\text{m}^{-2}$ ), *3D LD* and high emitter density ( $2.5 \text{ mol.} \mu\text{m}^{-2}$ ), *3D HD*. For the 2D  
522 competition, we generated a sparse ( $0.5 \text{ mol.} \mu\text{m}^{-2}$ ), *2D LD*, and very high density dataset ( $5 \text{ mol.} \mu\text{m}^{-2}$ ),  
523 *2D HD*.

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

529

## 530 **2.2. Photophysics activation model**

531 We incorporated a 4-state model of fluorophore photophysics<sup>18</sup>, including a transient dark state (dye  
532 blinking) and a bleaching pathway (**Fig. S18C**). Given a list of source locations from the structure  
533 simulator, fluorophore blinking was simulated by a 4-states Markov chain model. The states are ON,  
534 OFF, BLEACH, DARK and the transitions are Poisson distributed (**Fig. S18C**), except for the OFF to ON  
535 transitions which follow a uniform random distribution to reflect that in typical experimental  
536 conditions, constant imaging density is maintained by tuning the photoactivation rate during the  
537 experiment. All switching is calculated at sub-frame resolution and then total fluorophore on-time  
538 was integrated over each frame.

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

$$540 \quad T_{LIFETIME} = \frac{1}{\frac{1}{T_{ON}} + \frac{1}{T_{BLEACH}}}$$

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

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

$$545 \quad F = \frac{\Phi P \sigma}{e}$$

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

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

## 553 **2.3. Experimental Point Spread Function**

554 Model PSFs, stored as high resolution look up tables, were derived from experimentally measured  
555 PSFs. Although the algorithmic approach is distinct, the concept of accurately modelling the  
556 experimental PSF based on calibration data bears relation to the PSF phase retrieval approach  
557 previously employed by Hanser and coworkers<sup>42</sup>.

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

561 To acquire experimental PSFs, we took 100 nm Tetraspek beads (Invitrogen) adsorbed to #1.5 (170  $\mu$ m  
562 thick) coverglass, imaged in water. The excitation wavelength was between 640 nm and 647 nm, and  
563 a Cy5 emission filter was used. Data acquisition parameters for each modality are listed in  
564 **Supplementary Table 5**.

565 The experimental PSFs used to generate the simulated data are available on the competition website.  
566 As the goal of this study was to compare software obtained on typical SMLM microscopes, we  
567 deliberately chose PSFs representative of common implementations of each 3D modality. However,  
568 additional PSF engineering should improve results of any specific modality, for example adaptive-  
569 optics corrected astigmatism<sup>43</sup>, or reduced Z-range, higher SNR DH-PSF designs<sup>21</sup>.

570 The experimental point spread functions used here were measured for fluorescent beads adsorbed to  
571 the microscope cover slip, and should be appropriate simulations of SMLM data acquired within a few  
572 microns of the cover slip. Performing SMLM imaging at greater depths, *e.g.*, in tissue or even deep  
573 within single cells, with oil immersion objectives will cause spherical aberration due to refractive index  
574 mismatch<sup>44</sup>. In order to accurately simulate SMLM data acquired at depth, the experimental PSFs  
575 could be acquired at a matching depth, by embedding fluorescent beads in agarose. Alternatively, the  
576 PSF for beads at the coverslip could be measured and explicitly calculated via phase retrieval, and then  
577 convolved with the appropriate degree of spherical aberration<sup>44</sup>.

578

## 579 **2.4. Simulation PSF construction**

580 For each modality, 3-6 beads were selected within a small ( $< 32 \mu\text{m}$ ) region, to minimize PSF variation  
581 due to spherical aberration. Images for each selected bead were interpolated in XY to a pixel size of  
582 10 nm. Beads were then coaligned by cross-correlation on the in-focus frame. Coaligned beads were  
583 averaged in XY to minimize pixel quantization artefacts and to increase SNR. Where necessary, Z-stacks  
584 were interpolated to a Z-step size of 10 nm. A central Z-range of  $1.5 \mu\text{m}$  was selected that represents  
585 151 optical planes with a Z-step of 10 nm. The Z-range covers  $-750 \text{ nm}$  to  $+750 \text{ nm}$ . The plane of best  
586 focus was chosen as the simulation 0 nm plane. Each model PSF was normalized such that the total  
587 intensity of the PSF in the in-focus frame within a diameter of 3 FWHM from the PSF center was equal  
588 to 1.

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

591 In biplane super-resolution microscopy, emitted fluorescence is split into two simultaneously imaged  
592 channels, with a small (500-1000 nm) defocus introduced between the two channels<sup>15</sup>. As the small  
593 defocus should introduce minimal additional aberration into an optical system, we semi-synthetically  
594 constructed a realistic biplane PSF from the experimental 2D PSF. The two defocused PSFs were  
595 constructed by duplicating the 2D PSF and offsetting it by  $-250 \text{ nm}$  and  $250 \text{ nm}$  for each Z-plane.

596 This yielded five high SNR model PSFs with an isotropic voxel size of  $10 \times 10 \times 10 \text{ nm}^3$ .

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

## 601 **2.5. Noise model**

602 A constant mean autofluorescent background was added to the noise-free simulated images, and  
603 these images were then fed through the noise model representing Poisson distributed fluorescence  
604 emission recorded on a high quantum efficiency back-illuminated EMCCD<sup>45,46</sup>.

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

- 606 •  $\sigma_S$ , the shot noise produced by the fluorescence background and signal and the spurious  
607 charge. Shot noise can be derived from the second moment of the Poisson distribution
- 608 •  $\sigma_R$ , the read noise of EMCCD camera, which is described by second moment of the Gaussian  
609 distribution

- $\sigma_{EM}$ , the electron multiplication noise introduced by the gain process, which is described by the second moment of the Gamma distribution<sup>46</sup>.

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

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

The final simulated photon electrons will thus be given by:

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

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

which leads to the final pixel counts:

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

## 2.6. Depth-dependent lateral distortion/ wobble

As the PSF models are experimentally derived, the 3D estimated localizations exhibit a depth-dependent lateral distortion, here called *wobble*. This optical distortion is due to a combination of a systematic offset (arbitrary definition of PSF center) and optical aberrations<sup>47</sup>. In order to compare estimated and true localizations, we correct this effect during the assessment (**Methods 3.1**).

## 2.7 Comparison of software results between different modalities.

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

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

## 3. SOFTWARE ASSESSMENT

### 3.1 Protocol

Each localization file submitted by the participants was manually checked for erroneous systematic errors in the definition of the dataset coordinate system, such as offsets, XY axis flips or clear scaling errors. Datasets were then programmatically standardized into a consistent output format. All



651 modifications are publicly available. If required, the modifications consisted of columns reordering,  
652 reversing axes, XY axis swap, and shifting the lateral positions by a half camera pixel.

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

661 The pairing between true and estimated localizations was performed frame by frame. For every frame,  
662 we identified the localizations that are close enough to a ground-truth position as true-positives (TP),  
663 the spurious localizations as false-positives (FP) and the undetected molecules as false-negatives (FN).  
664 The procedure matches two sets of localizations. We deployed the presorted nearest-neighbor search  
665 for its efficiency, with a linking threshold of 250 nm. The results are effectively similar to the  
666 computationally intensive Hungarian algorithm<sup>7</sup>.

#### 667 *Photon thresholding*

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

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

#### 678 *Wobble correction*

679 The centroid of experimental point spread functions shifts laterally by as much as 50 nm, as a function  
680 of axial position<sup>10,47</sup>. This is most often ignored by localization software, and instead corrected post-  
681 hoc by reference to a calibration curve<sup>37</sup>. Since our simulated PSF is experimentally derived, it was  
682 necessary to correct for these artefactual shifts between the observed localizations and ground truth,  
683 as part of the assessment process. This correction was performed using calibration data uploaded by  
684 competitors, similar to the correction typically performed on experimental data<sup>47</sup>.

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

690 In certain specific cases (identified on the competition website), at the request of authors, we did not  
691 apply this correction, for example because the software explicitly considered the whole 3D PSF during  
692 fitting and was thus immune to this lateral shift artefact. For accurate results, application of lateral  
693 shift correction is critical for analysis of localization microscopy simulations using experimentally  
694 derived PSFs, as can be seen by comparison of typical software results with and without wobble  
695 correction (**Fig. S19**).

### 696 3.2 Metrics

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

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

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

709 
$$RMSE \text{ lateral (RMSE Lateral) [nm]: } \sqrt{\frac{1}{TP} \sum_{i \in S \cap T} (x_i^s - x_i^t)^2 + (y_i^s - y_i^t)^2}.$$

710 
$$RMSE \text{ axial (RMSE Axial) [nm]: } \sqrt{\frac{1}{TP} \sum_{i \in S \cap T} (z_i^s - z_i^t)^2}.$$

- 711 2. *Jaccard index (JAC, %).* In addition to localization precision, SMLM image resolution depends  
712 critically on number of localized molecules<sup>48</sup>, so it is crucial for SMLM software to accurately  
713 detect a large fraction of molecules in a dataset, and minimize false localizations. For every  
714 frame, we identified the localizations that are close enough to a ground-truth position as  
715 true-positives (TP), the spurious localizations as false-positives (FP) and the undetected  
716 molecules as false-negatives (FN). We then computed the *Jaccard index* (JAC, %), which  
717 measures the fraction of correctly detected molecules in a dataset,

718 
$$JAC = 100 \frac{TP}{TP + FP + FN}$$

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

723 
$$E = 100 - \sqrt{(100 - JAC)^2 + \alpha^2 RMSE^2}$$

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

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

736 We also calculated localization based metric results as a function of axial position. We proceeded by  
737 considering a subset of activations lying within an interval of axial positions (*i.e.*, from the true  
738 localizations). Then, most of the metrics (*e.g.*, Recall) are locally computed. This yields a curve  
739 providing information on the depth performance of each software / modality.

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

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

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

### 759 **3.3 Baseline Localization Software**

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

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

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

## 772 **4 REAL DATA ASSESSMENT**

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

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

## 781 **5 DATA AVAILABILITY**

### 782 **5.1 Data availability statement**

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

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

### 790 **5.2 Code availability statement**

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

## 792 **REFERENCES, ONLINE METHODS**

793 40. Carlini, L. & Manley, S. Live Intracellular Super-Resolution Imaging Using Site-Specific Stains.  
794 *ACS Chem. Biol.* **8**, 2643–2648 (2013).

795 41. Shim, S.-H. *et al.* Super-resolution fluorescence imaging of organelles in live cells with  
796 photoswitchable membrane probes. *Proc. Natl. Acad. Sci.* **109**, 13978–13983 (2012).

797 42. Hanser B. M., Gustafsson M. G. L., Agard D. A. & Sedat J. W. Phase-retrieved pupil functions in  
798 wide-field fluorescence microscopy. *J. Microsc.* **216**, 32–48 (2004).

799 43. Izeddin, I. *et al.* PSF shaping using adaptive optics for three-dimensional single-molecule  
800 super-resolution imaging and tracking. *Opt. Express* **20**, 4957–4967 (2012).

801 44. McGorty, R., Schnitzbauer, J., Zhang, W. & Huang, B. Correction of depth-dependent  
802 aberrations in 3D single-molecule localization and super-resolution microscopy. *Opt. Lett.* **39**, 275–  
803 278 (2014).

804 45. Hirsch, M., Wareham, R. J., Martin-Fernandez, M. L., Hobson, M. P. & Rolfe, D. J. A Stochastic  
805 Model for Electron Multiplication Charge-Coupled Devices – From Theory to Practice. *PLOS ONE* **8**,  
806 e53671 (2013).

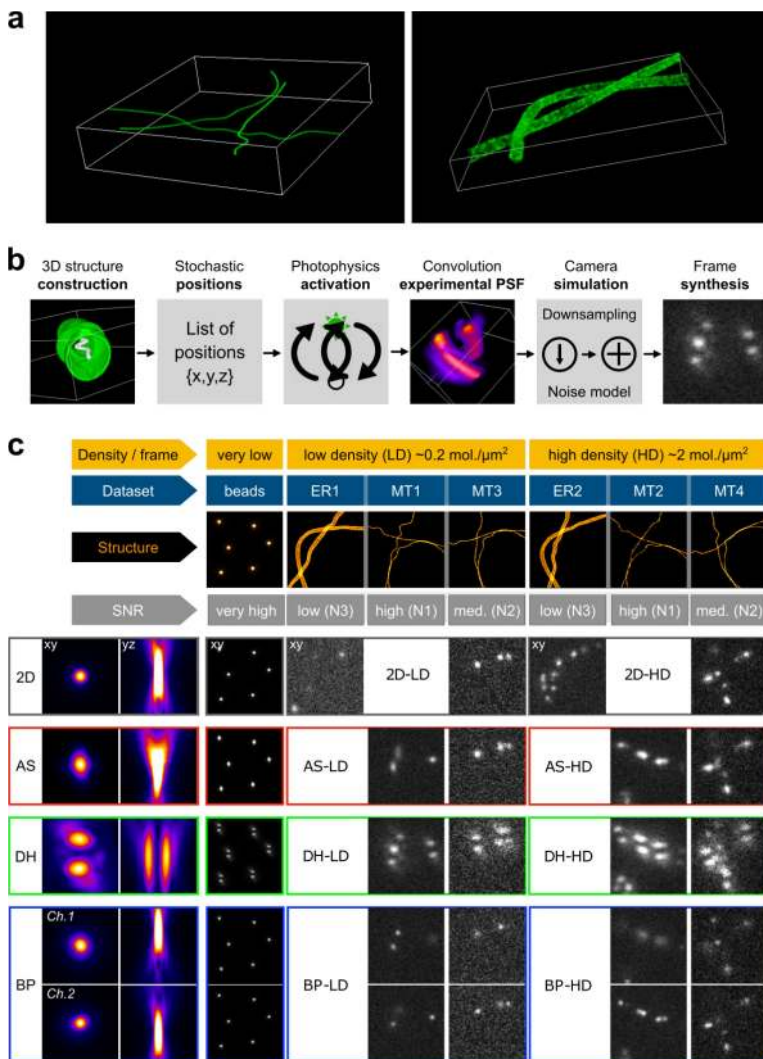
807 46. Basden, A. G., Haniff, C. A. & Mackay, C. D. Photon counting strategies with low-light-level  
808 CCDs. *Mon. Not. R. Astron. Soc.* **345**, 985–991 (2003).

809 47. Carlini, L., Holden, S. J., Douglass, K. M. & Manley, S. Correction of a Depth-Dependent Lateral  
810 Distortion in 3D Super-Resolution Imaging. *PLoS ONE* **10**, e0142949 (2015).

811 48. Baddeley, D. & Bewersdorf, J. Biological Insight from Super-Resolution Microscopy: What We  
812 Can Learn from Localization-Based Images. *Annu. Rev. Biochem.* **87**, 965–989 (2018).

813

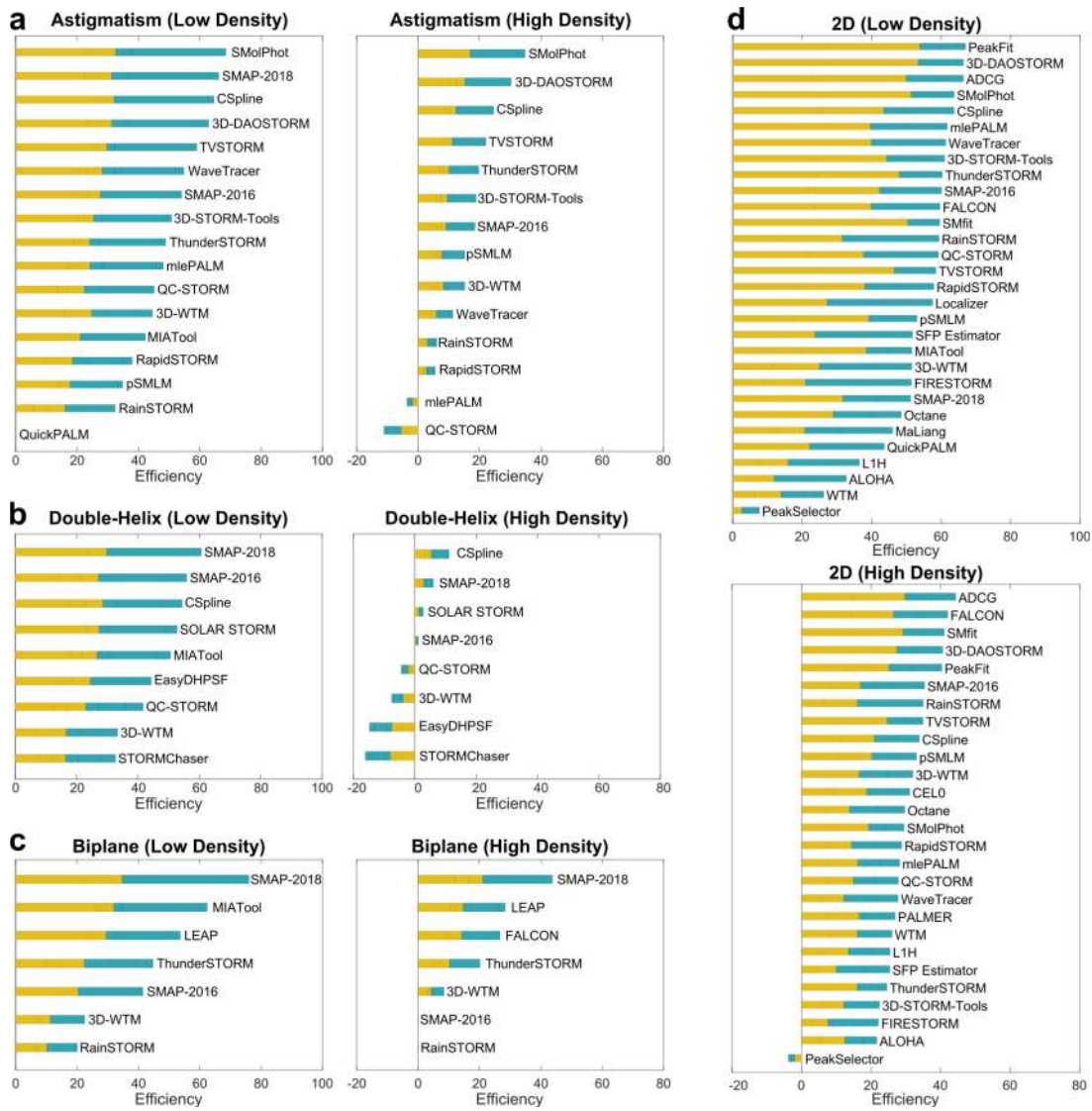
814 **FIGURES**



815

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

827

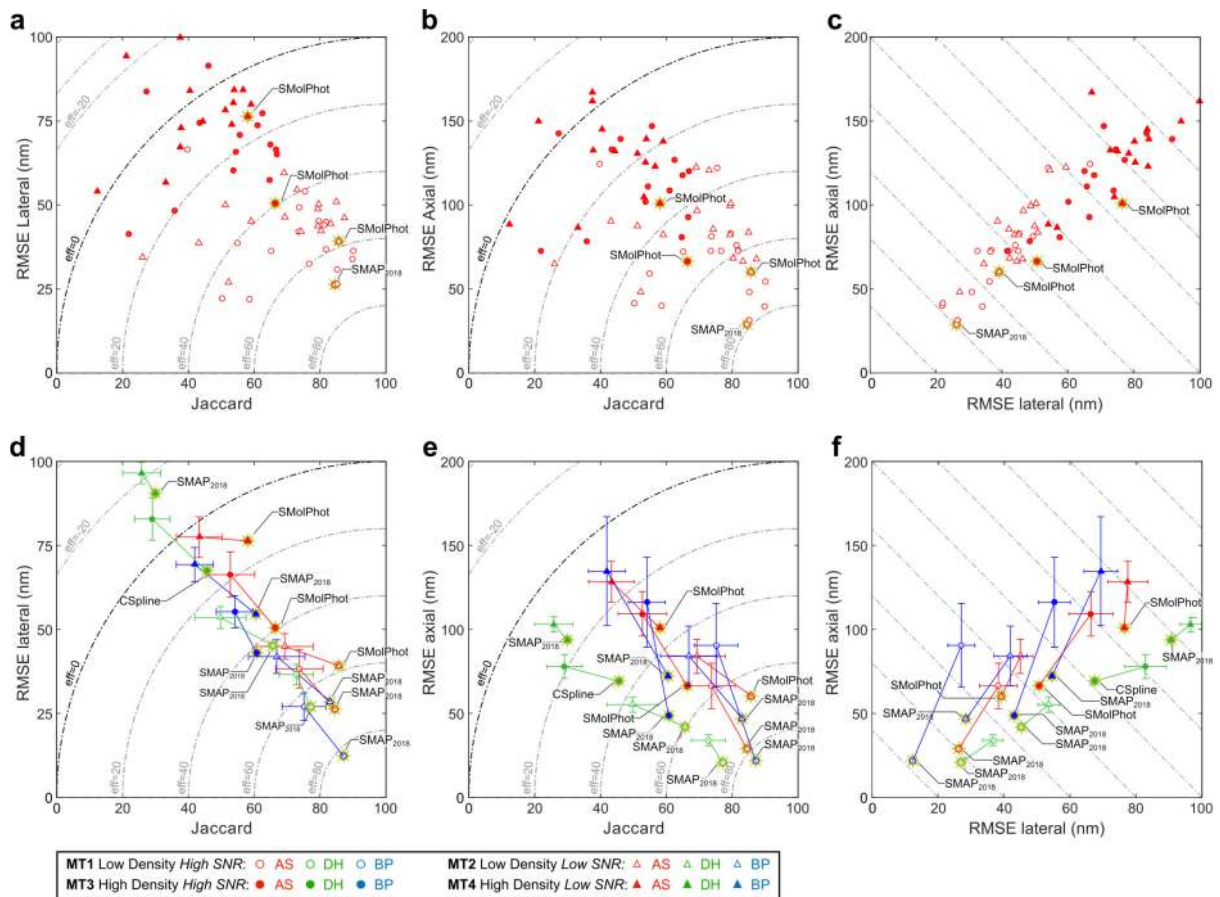


828

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

833

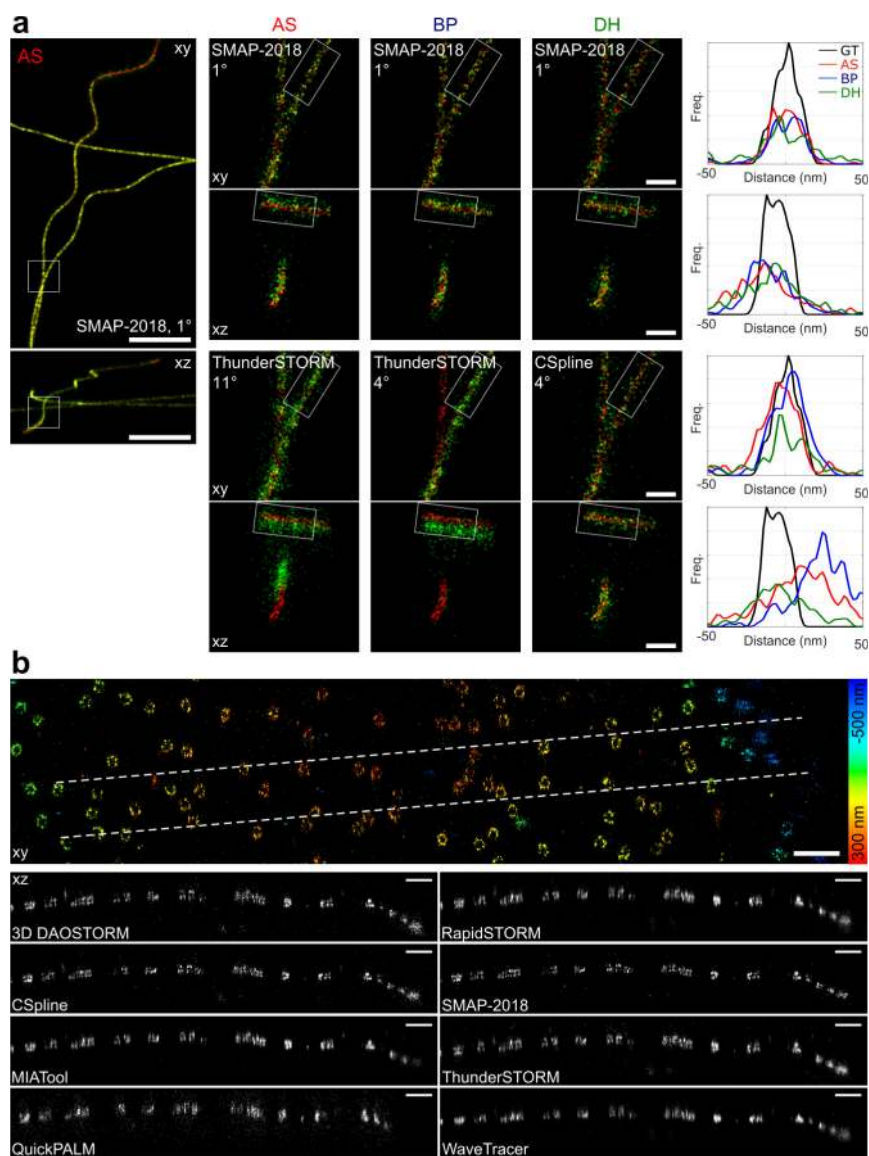




834

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

841



842

843 **Figure 4: Super-resolved images of software results for simulated and real competition datasets. A.** *Xy*  
 844 *and xz projection images of 3D competition datasets for representative software. Top: best-in-class*  
 845 *software in each modality, for high SNR low density dataset. Bottom: representative average software.*  
 846 *Left: xy and xz overview images for winning AS software. Middle: xy and xz zoom images of boxed*  
 847 *regions in left panel, for winning and mid-range software, each modality. Right: xy and xz line profiles*  
 848 *of winning and mid-range software for each modality, for boxed regions in middle panel. Image colors:*  
 849 *red, ground truth; green, software results. Line profiles: GT, ground truth, black; AS, astigmatism, red;*  
 850 *BP, biplane, blue; DH, double helix, green. Panel key: Software-name Dataset-ranking°. Scale bar: full*  
 851 *image, 1  $\mu$ m, magnified regions, 100 nm. B. Astigmatism software results for real nuclear pore complex*  
 852 *3D STORM data. Top: Super-resolved overview image in xy for 3D-DAOSTORM software, color coded*  
 853 *for depth. Bottom: xz orthoslices along 600 nm wide dashed region indicated in top panel for 8*  
 854 *astigmatism software packages. Scale bars, 500 nm.*