



## King's Research Portal

DOI:

[10.1038/nmeth.3612](https://doi.org/10.1038/nmeth.3612)

*Document Version*

Peer reviewed version

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Rubin-Delanchy, P., Burn, G. L., Griffié, J., Williamson, D. J., Heard, N. A., Cope, A. P., & Owen, D. M. (2015). Bayesian cluster identification in single-molecule localization microscopy data. *NATURE METHODS*, 12(11), 1072-1076. <https://doi.org/10.1038/nmeth.3612>

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

# Bayesian Cluster Identification in Single-Molecule Localisation Microscopy Data

Patrick Rubin-Delanchy<sup>1</sup>, Garth L Burn<sup>2</sup>, Juliette Griffié<sup>2</sup>, David Williamson<sup>3</sup>, Nicholas A Heard<sup>4</sup>, Andrew Cope<sup>5</sup> & Dylan M Owen<sup>2</sup>

<sup>1</sup>Heilbronn Institute for Mathematical Research, School of Mathematics, University of Bristol, Bristol, UK, <sup>2</sup>Department of Physics and Randall Division of Cell and Molecular Biophysics, King's College London, London, UK, <sup>3</sup>Manchester Collaborative Centre for Inflammation Research, University of Manchester, Manchester, UK, <sup>4</sup>Department of Mathematics, Imperial College London, London, UK, <sup>5</sup>Academic Department of Rheumatology, Division of Immunology, Infection and Inflammatory disease, King's College London, London, UK.

Correspondence should be addressed to [patrick.rubin-delanchy@bristol.ac.uk](mailto:patrick.rubin-delanchy@bristol.ac.uk) and [dylan.owen@kcl.ac.uk](mailto:dylan.owen@kcl.ac.uk)

Single-molecule identification-based super-resolution microscopy techniques such as photo-activated localisation microscopy (PALM) and stochastic optical reconstruction microscopy (STORM) produce pointillist data sets of molecular coordinates. While many algorithms exist for the identification and localisation of molecules from the raw image data, methods for analysing the resulting point patterns for properties such as clustering have remained relatively under-studied. Here, we present the first model-based Bayesian approach to evaluate molecular cluster assignment proposals which, in this article, are generated by analysis based on Ripley's K-function. The method is also the first to take full account of the individual localisation precisions calculated for each emitter. The technique is validated using simulated and experimental data from which we characterise the clustering behaviour of CD3 $\zeta$ , an important subunit of the CD3-T cell receptor complex required for T cell function, in resting and activated primary human T cells.

Conventional fluorescence microscopes produce images of the distribution of fluorophores in the sample convolved with the microscope Point Spread Function (PSF). Due to diffraction, this PSF typically has a width of hundreds of nanometres meaning the resulting image has a resolution, as assessed by the Rayleigh criterion, of  $\sim 200$  nm. Several strategies now exist to circumvent this resolution limit<sup>1</sup>. Some of these, such as Stimulated Emission Depletion (STED) microscopy, rely on narrowing the excitation spot of a confocal microscope by means of a toroidal depletion beam and the process of stimulated emission<sup>2,3</sup>. Despite the increased resolution, these produce conventional fluorescence images, i.e., arrays of pixels with values representing the fluorescence intensity at those locations. Quantification can be performed in the same way as for conventional microscopes.

Another strategy is based on Single-Molecule Localisation Microscopy (SMLM)<sup>4-7</sup>. This relies on the temporal separation of the excitation of fluorophores in the sample whose PSFs would otherwise overlap at the detector. The position of each fluorophore can then be estimated from the centres of the PSFs. Many algorithms are available to extract the x-y coordinates of the molecules<sup>8-10</sup>. Each emitter can be localised to a precision between 10 and 30 nm. Common strategies for the temporal separation of molecules involve intra-molecular rearrangements to switch from dark to fluorescent states or the exploitation of non-emitting molecular radicals<sup>11, 12</sup>. These strategies are typically pursued using photoactivatable or photoconvertible fluorescent proteins or small molecule probes coupled with a reducing buffer and immunostaining protocols<sup>13</sup>. We refer to all such strategies as SMLM.

Unlike non-pointillist microscopy methods, SMLM imaging does not produce a conventional image. Instead, the raw data is a list of the x-y coordinates of all the fluorophores, each with an associated, estimated localisation precision. The analysis of spatial point patterns (SPPs) requires a different statistical toolkit to the analysis of pixel arrays, only now being explored in the context of SMLM.

Several techniques for analysing SPPs generated from SMLM have been proposed. To investigate and quantify clustering behaviour, widely used are Ripley's K-function<sup>14-16</sup> and pair-correlation (PC) analysis<sup>17, 18</sup>. Both rely on drawing a series of concentric shapes – circles in the case of the K-function and tori in the case of PC – around each localisation and counting the number of neighbours enclosed. These allow the degree of clustering at different spatial scales to be determined. In the case of the K-function, the values at each localisation can be interpolated to create cluster maps to which thresholds can then be applied<sup>19</sup>.

The methods presented above have several key shortcomings. They often require calibration data or user-selected analysis parameters that strongly influence the output. This problem is exacerbated by batch-processing, meaning that regions are often analysed with the same sub-optimal parameters. The methods also do not take any account of the individual localisation precisions for each point. Finally, these are model-free methods, which makes it inherently difficult to judge performance and interpret results.

Here, we present a model-based, Bayesian approach to cluster analysis of SPPs generated by SMLM. The quality of a given assignment of molecules to clusters is evaluated against its (marginal) posterior probability, computed on the basis of a fully-specified model for the data, including the localisation precisions. This provides a principled mechanism for choosing between clustering proposals generated by different algorithms and settings. In this article, clustering proposals are generated using a strategy based on the K-function<sup>14</sup>, with variable spatial scale and threshold. We therefore generate several thousand candidate clustering proposals per ROI, from which the optimum is selected according to the Bayesian model. Code is available in the Supplementary Material.

We demonstrate using simulated SPP data that we can accurately evaluate molecular clustering in a variety of conditions. The technique is then used to compare the clustering behaviour of CD3 $\zeta$ -mEos3 in resting T cells *versus* at the T cell immunological synapse. Here, it is accepted that proteins, including the CD3 $\zeta$  subunit, arrange

into microclusters upon synapse formation. While many other biological processes involve the clustering of proteins at the cell surface, this application is especially informative because both the K-function strategy and PC have been applied previously<sup>14, 15, 20, 21</sup>. For experimental data, it is important that artefacts caused by multiple blinking of individual fluorophores and overlapped PSFs, inherent to the methodology of SMLM, are removed (or accounted for). Our algorithm does not attempt to correct for, or be robust to, multiple blinking. Therefore, our method generates quantitatively reliable results only when multiple blinking has been corrected, as is possible with PALM data. Here, this was achieved using ThunderSTORM<sup>22</sup> localisation software which includes blink correction based on the method of Annibale *et al*<sup>23</sup> previously validated using mEos, and is able to fit multiple emitters to overlapping PSFs. Our algorithm is applicable to data from other SMLM implementations, however, because of the difficulties of correcting multiple blinking, results must be interpreted appropriately.

## RESULTS

We begin by assuming a single coordinate for each molecule in the region of interest (ROI), generated by the localization software, in our case ThunderSTORM<sup>24</sup>. The 2D molecular positions are modelled as a set of Gaussian distributed clusters overlaid on a completely spatially random (CSR) background. These molecular coordinates are then disturbed by Gaussian distributed errors as a result of the localisation process. The errors have different standard deviations, which are treated as known. In fact, they are estimated from the raw microscopy data based on the number of collected photons, PSF width, local background noise and camera pixel size<sup>25</sup>.

The cluster centres themselves are assumed to be uniformly distributed over the ROI and their radii (standard deviation) are drawn from a user-supplied prior distribution. Localisations are assigned independently to the CSR background with a fixed prior probability and the remaining localisations are clustered according to the Dirichlet process<sup>26</sup>. We compute the posterior probability of any given assignment of localisations to clusters (a clustering proposal) with respect to the above model. The calculation is deterministic, unlike with many Bayesian models, requiring only numerical integration over one dimension (see Supplementary Methods).

To generate clustering proposals we use a method based on Ripley's K-function<sup>16</sup>. Every localisation is allocated a clustering score,  $L$ , as proposed by Getis<sup>27</sup>.  $L$  is a function of the number of localisations within a distance,  $r$ , of that point, normalised by the mean molecular density of the ROI. Localisations with a value of  $L$  below a certain threshold,  $T$ , are assigned to the background.  $T$  can be interpreted as the minimum local density required for a point to be assigned to a cluster. Any two remaining localisations within a distance  $2r$  of each other are then connected and the connected components form clusters (**Fig. 1a**). By scanning  $r$  and  $T$  we generate of the order of 10,000 cluster proposals which are then assigned a posterior probability. The highest scoring proposal is retained, key descriptors are extracted. Although other proposal mechanisms are possible, e.g. K-means<sup>28</sup>, KDE clustering<sup>29</sup>, agglomerative clustering<sup>30</sup> or Density-based Spatial Clustering of Applications with Noise (DBSCAN)<sup>31</sup>, this approach is attractive because it has a straightforward geometrical interpretation and can be rapidly computed.

In a representative simulated data set (**Fig. 2b**), the posterior probability is calculated for a range of values of  $r$  and thresholds (**Fig. 2c**). The dashed line indicates positions where  $L(r)$  is thresholded at  $r$ , i.e., the line  $T=r$ .  $L(r)$  being greater (smaller) than  $r$  indicates that points are more (less) clustered at that scale than would be expected under CSR. It is intuitive that a clustering model should favour thresholding  $L(r)$  above  $r$ . Four  $r$ - $T$  combinations are selected and the clustering proposals generated by each are shown (**Fig. 2d**). The highest scoring is proposal 2. The others illustrate three different manifestations of a sub-optimal selection of  $r$  and  $T$ . In proposal 1, several small, spurious clusters are identified largely due to the small value of  $r$  used. In proposal 3, the threshold is too stringent and localisations at the cluster extrema are assigned to the background. Finally, in proposal 4, clusters are merged due to a large value of  $r$ .

### Performance and sensitivity analysis

SMLM localisation data were simulated under four different clustering scenarios. In the first, the Standard Conditions, a 3000 x 3000 nm area contains 2000 localisations. These comprise 10 Gaussian clusters with radius 50 nm containing 100 localisations each and 1000 localisations (50%) in the background. Each localisation is then disturbed by Gaussian noise with variance drawn from a Gamma distribution with mean 30 nm and standard deviation 13 nm (emulating the localisation error of the microscope). These parameters were chosen to approximately reflect typical clustering behaviour of proteins at the immunological synapse<sup>14, 15, 20, 21</sup>. The three remaining scenarios have the same parameters as the Standard Conditions except where stated otherwise. The second scenario is a sparse data set containing only 200 localisations with 10 per cluster and 100 in the background. In the third, the cluster radii are 100 nm. Finally, the fourth scenario has 10 localisations per cluster with 900 (90%) in the background. 100 ROIs were simulated for each scenario.

Representative example of each of the four simulated scenarios are shown (**Fig. 2a**) with corresponding heat-maps displaying the log-posterior probability (**Fig. 2b**). The highest scoring  $r$ - $T$  combinations are encircled and the generated proposals displayed (**Fig. 2c**). Histograms of three key cluster descriptors – cluster radii (empirical standard deviation of the localisations), number of localisations per cluster and percentage of localisations in clusters were generated (**Supplementary Fig. 1**). A thorough characterisation of our algorithm tested on varying simulation parameters is also provided (**Supplementary Fig. 2**).

Our algorithm substantially outperforms currently available cluster analysis methods, e.g. Getis and Franklin's Local Point Pattern analysis and DBSCAN and offers definite advantages over other approaches, e.g. Ripley's K-function and pair correlation (**Supplementary Figs. 3- 5**, and Supplementary Materials and Methods). There, we also test our algorithm against more challenging conditions, including an uneven background (**Supplementary Figs. 6 and 7**), very small clusters (multimers) (**Supplementary Fig. 8**) or clusters with variable size (**Supplementary Fig. 9**). A side by side comparison of our algorithm, Getis's method and DBSCAN on three

example conditions clearly demonstrates the superiority of our approach (**Supplementary Fig. 10**). A sensitivity analysis to prior settings is also provided (**Supplementary Fig. 11**).

### **Analysis of protein clustering in primary human T cells**

We analysed SMLM data of the CD3 $\zeta$  subunit of the TCR-CD3 complex, fused to the photoswitchable fluorescent protein mEos3 at the plasma membrane of CD4<sup>+</sup> primary human T cells that had formed an immunological synapse on anti-CD3/28 coated glass coverslips. Non-activating poly-L-lysine coated coverslips were used as a control. After 4 minutes of incubation on the coverslips at 37°C, cells were pH-shift fixed and imaged. Photoswitching of mEos3 was achieved using 405 nm laser light and switched proteins imaged using 564 nm excitation. Details of the sample preparation method and imaging can be found in the Supplementary Materials and Methods. Multiple blinking of fluorophores and overlapped PSFs were compensated for using ThunderSTORM localisation software<sup>24</sup> and optimal settings estimated using the method of Annibale *et al*<sup>23</sup> (**Supplementary Fig. 12** and Supplementary Materials and Methods). Note that while mEos3 displays multiple-blinking during PALM data acquisition, this effect can be effectively corrected (due to the different time-scales of photo-switching and photo-blinking), thus rendering the input data appropriate for our algorithm. The localisation precisions were calculated using the method of Quan *et al*<sup>25</sup> and representative histograms of these values are shown (**Supplementary Fig. 12**).

From SMLM images of resting and activated T cells (**Fig. 3a**), 3000 x 3000 nm regions ( $n = 30$  per condition) were selected and the localisations plotted (**Fig. 3b**). For each, the Log Posterior Probability heat map is shown (**Fig. 3c**) with the highest scoring proposal (**Fig. 3d**). Beeswarm plots of the percentage of localisations in clusters, number of clusters per region, cluster radii and relative density of localisations inside and outside clusters are shown (**Fig. 3e**). T-tests were used and their p-values were computed by permutation (see Supplementary Materials and Methods). In good agreement with previous reports<sup>21, 32</sup>, CD3 $\zeta$  was clustered in stimulated and non-stimulated cells, and cluster parameters were significantly altered. Despite no large scale changes in the percentage of localisations found in clusters ( $29 \pm 2\%$  in PLL to  $33 \pm 1\%$  in activated cells,  $P > 0.05$ ), a significant increase in the number of clusters and a significant decrease in the size of clusters was observed from resting to activated cells ( $8 \pm 1$  clusters per region on PLL versus  $20 \pm 3$  in activated cells;  $P \leq 0.005$  and  $82 \pm 4$  nm radius versus  $48 \pm 2$  nm after activation;  $P \leq 0.0005$ ; **Fig. 3e**). In addition, a significant increase in the density of localisations in clusters relative to non-clustered regions was observed ( $7 \pm 1$  versus  $14 \pm 1$ ,  $P \leq 0.0005$ ). Results are consistent when we divide localisations into two equally sized data sets (**Supplementary Fig. 13**). As with the simulations, we analysed all experimental data using three well established cluster analysis methods; Getis and Franklin's cluster maps, Ripley's K-function and Pair Correlation (**Supplementary Fig. 14**).

## DISCUSSION

Super-resolution imaging based on localising individual molecules is becoming increasingly widespread. While methods to localise molecules from the raw fluorescence data have been extensively analysed<sup>8, 10, 33</sup>, the subsequent interrogation of the point pattern data has been relatively under-studied. We have demonstrated a new, Bayesian cluster analysis algorithm for SMLM data.

Unlike previously demonstrated methods based on the generation of cluster maps which involve an interpolation algorithm to generate the surface<sup>19</sup>, the new method is not prone to artefacts in sparse data sets, e.g. from low copy-number proteins. The method also has the possibility to allow faster imaging as less localisations are required to accurately identify and characterise clustering. Increasing the speed of SMLM data acquisition and processing has been one of the major goals to move the technique into the domain of live cell imaging<sup>33</sup>.

The algorithm is only weakly sensitive to the prior settings and this is a major advantage over previous methods, where the initial choice of spatial scale and threshold has a large effect on the final results. In addition, here, all ROIs are analysed with the parameters which are estimated to be optimal for that specific region, rather than diverse regions being treated equally. The method is the first to take full account of the localisation precisions rather than treating all localisations as exact. We have stress-tested the algorithm under challenging conditions, for example, finding the detectability limit to be around 6 localisations per cluster. We hypothesise that a Bayesian model that explicitly targets small features would be more successful in detecting small multimers. Indeed, an interesting avenue of future research would be to develop a number of different models to capture the diversity of point patterns observed in SMLM data, including fibres, meshes, areas of exclusion and so on.

It is well known that raw SMLM data can exhibit artefacts, due to the photophysical nature of the process<sup>22, 23, 34</sup>, whereby individual molecules can re-excite and thus generate multiple localisations. In addition, due to the stochastic nature of the activation process, it is possible for several PSFs to overlap at the detector, causing errors in the extracted coordinates. Our algorithm does not attempt to correct or to be robust to multiple-blinking effects. If there is suspicion that these have not been adequately addressed by the localisation software, then the output of our algorithm should be interpreted with caution. In our case, we acquired experimental data using PALM, for which multiple blinking can be corrected. This is because of the different timescales of molecular photo-conversion and photo-blinking. In other experimental conditions, for example, when using small molecule dyes, such corrections may not be possible. Therefore, the outputs of the algorithm may contain artefacts, in particular, spurious clusters. Our algorithm remains a valuable exploratory tool for such data.

Our method therefore allows the accurate and principled quantification of clustering behaviour in SMLM data in a manner that is more automatic, robust and objective than previously possible. In this initial case, we focused on a model consisting of circular, Gaussian distributed clusters overlaid on a CSR background. In future, it will be possible to create generative models with different clustering characteristics. Evaluation of SMLM data against such models may allow a better understanding of the biophysical principles underlying protein clustering.

## METHODS

Methods and any associated references are available in the online version of the paper at

<http://www.nature.com/naturemethods/>.

Note: Supplementary information is available on the Nature Methods website.

## ACKNOWLEDGEMENTS

DMO Acknowledges funding from the ERC (FP7 Starter Grant 337187).

## AUTHOR CONTRIBUTIONS

PRD, NAH and DMO conceived the method. PRD, JG and DMO performed the analysis. PRD and DMO wrote the manuscript. GB acquired cell data. GB, DW and AC provided new materials.

## COMPETING FINANCIAL INTERESTS

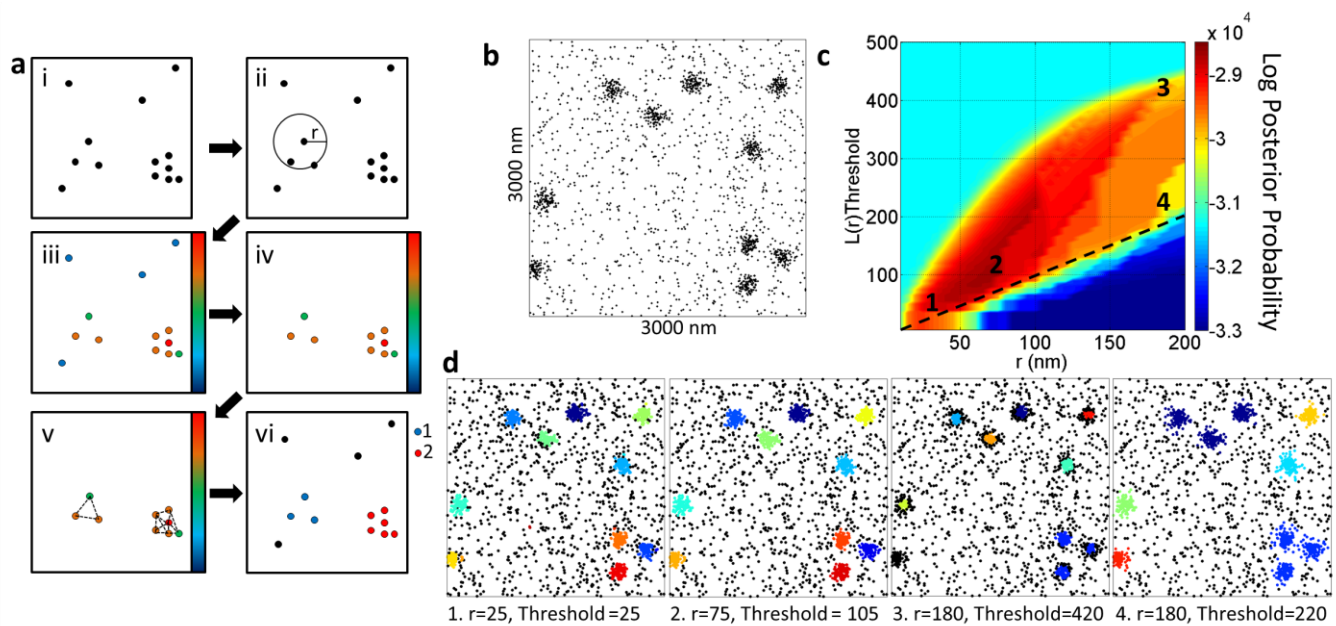
The authors declare no competing financial interests.

## REFERENCES

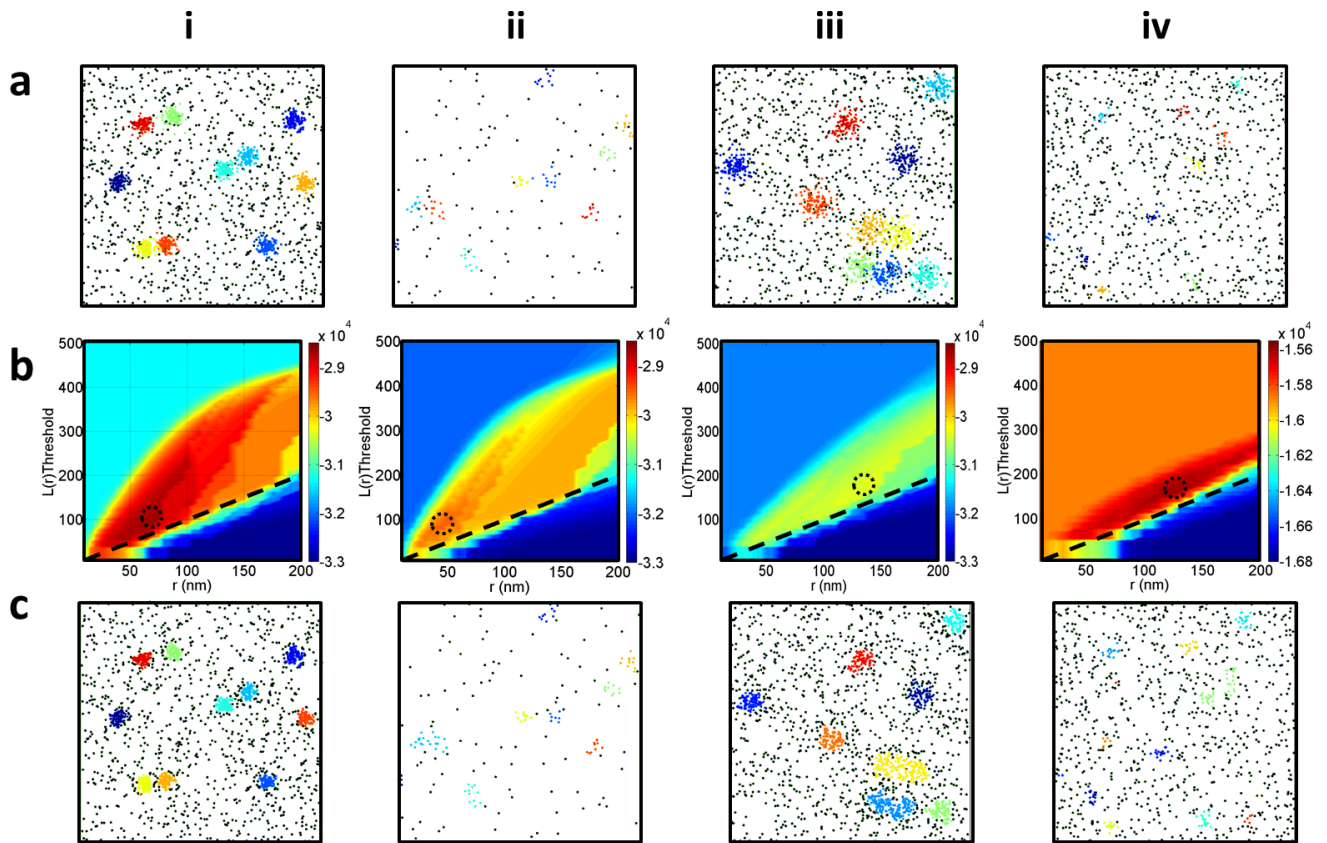
1. Huang, B. Super-resolution optical microscopy: multiple choices. *Curr. Opin. Chem. Biol* **14**, 10-14 (2010).
2. Hell, S.W. & Wichmann, J. Breaking the diffraction resolution limit by stimulated emission: Stimulated-emission-depletion fluorescence microscopy. *Opt. Lett.* **19**, 780-782 (1994).
3. Chmyrov, A. et al. Nanoscopy with more than 100,000 'doughnuts'. *Nat. Meth.* **10**, 737-740 (2013).
4. Betzig, E. et al. Imaging Intracellular Fluorescent Proteins at Nanometer Resolution. *Science* **313**, 1642-1645 (2006).
5. Rust, M.J., Bates, M. & Zhuang, X. Sub-diffraction-limit imaging by stochastic optical reconstruction microscopy (STORM). *Nat. Meth.* **3**, 793-796 (2006).
6. Heilemann, M. et al. Subdiffraction-Resolution Fluorescence Imaging with Conventional Fluorescent Probes. *Angew. Chem. Int. Ed.* **47**, 6172-6176 (2008).
7. Hess, S.T., Girirajan, T.P.K. & Mason, M.D. Ultra-High Resolution Imaging by Fluorescence Photoactivation Localization Microscopy. *Biophys. J.* **91**, 4258-4272 (2006).
8. Wolter, S. et al. rapidSTORM: accurate, fast open-source software for localization microscopy. *Nat. Meth.* **9**, 1040-1041 (2012).
9. Holden, S.J., Uphoff, S. & Kapanidis, A.N. DAOSTORM: an algorithm for high- density super-resolution microscopy. *Nat. Meth.* **8**, 279-280 (2011).
10. Henriques, R. et al. QuickPALM: 3D real-time photoactivation nanoscopy image processing in ImageJ. *Nat. Meth.* **7**, 339-340 (2010).
11. van de Linde, S. et al. Direct stochastic optical reconstruction microscopy with standard fluorescent probes. *Nat. Protocols* **6**, 991-1009 (2011).
12. Heilemann, M., van de Linde, S., Mukherjee, A. & Sauer, M. Super-Resolution Imaging with Small Organic Fluorophores. *Angew. Chem. Int. Ed.* **48**, 6903-6908 (2009).
13. Dempsey, G.T. et al. Photoswitching Mechanism of Cyanine Dyes. *JACS* **131**, 18192-18193 (2009).



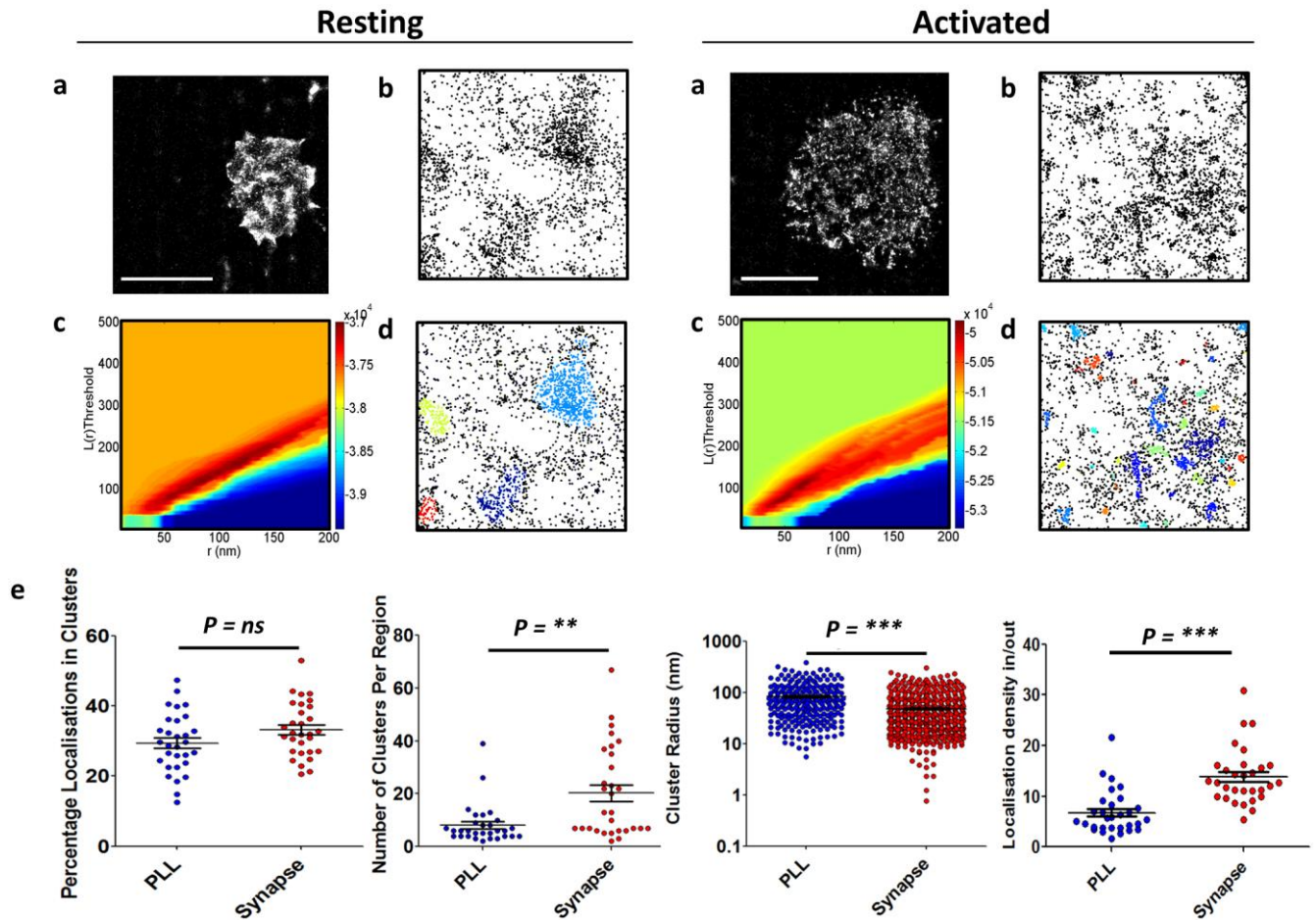
14. Williamson, D.J. et al. Pre-existing clusters of the adaptor Lat do not participate in early T cell signaling events. *Nat. Immunol.* **12**, 655-662 (2011).
15. Rossy, J., Owen, D.M., Williamson, D.J., Yang, Z. & Gaus, K. Conformational states of the kinase Lck regulate clustering in early T cell signaling. *Nat. Immunol.* **14**, 82-89 (2013).
16. Ripley, B.D. Modelling spatial patterns. *J. R. Stat. Soc. Series B Stat. Methodol.* **39**, 172-192 (1977).
17. Sengupta, P. et al. Probing protein heterogeneity in the plasma membrane using PALM and pair correlation analysis. *Nat. Meth.* **8**, 969-975 (2011).
18. Veatch, S.L. et al. Correlation Functions Quantify Super-Resolution Images and Estimate Apparent Clustering Due to Over-Counting. *PLoS ONE* **7**, e31457 (2012).
19. Owen, D.M. et al. PALM imaging and cluster analysis of protein heterogeneity at the cell surface. *J. Biophoton.* **3**, 446-454 (2010).
20. Sherman, E. et al. Functional Nanoscale Organization of Signaling Molecules Downstream of the T Cell Antigen Receptor. *Immunity* **35**, 705-720 (2011).
21. Lillemeier, B.F. et al. TCR and Lat are expressed on separate protein islands on T cell membranes and concatenate during activation. *Nat. Immunol.* **11**, 90-96 (2010).
22. Annibale, P., Vanni, S., Scarselli, M., Rothlisberger, U. & Radenovic, A. Identification of clustering artifacts in photoactivated localization microscopy. *Nat. Meth.* **8**, 527-528 (2011).
23. Annibale, P., Vanni, S., Scarselli, M., Rothlisberger, U. & Radenovic, A. Quantitative Photo Activated Localization Microscopy: Unraveling the Effects of Photoblinking. *PLoS ONE* **6**, e22678 (2011).
24. Ovesný, M., Křížek, P., Borkovec, J., Švindrych, Z. & Hagen, G.M. ThunderSTORM: a comprehensive ImageJ plug-in for PALM and STORM data analysis and super-resolution imaging. *Bioinformatics* **30**, 2389-2390 (2014).
25. Quan, T., Zeng, S. & Huang, Z.-L. Localization capability and limitation of electron-multiplying charge-coupled, scientific complementary metal-oxide semiconductor, and charge-coupled devices for superresolution imaging. *J. Biomed. Opt.* **15**, 066005-066005-066006 (2010).
26. Ferguson, T.S. A Bayesian Analysis of Some Nonparametric Problems. *Ann. Stat.* **1**, 209-230 (1973).
27. Getis, A. & Franklin, J. Second-Order Neighborhood Analysis of Mapped Point Patterns. *Ecology* **68**, 473-477 (1987).
28. Lloyd, S. Least squares quantization in PCM. *IEEE Trans. Inf. Theor.* **28**, 129-137 (2006).
29. Hinneburg, A. & Gabriel, H.-H. in *Advances in Intelligent Data Analysis VII*, Vol. 4723. (eds. M. R. Berthold, J. Shawe-Taylor & N. Lavrač) 70-80 (Springer Berlin Heidelberg, 2007).
30. Johnson, S.C. Hierarchical clustering schemes. *Psychometrika* **32**, 241-254 (1967).
31. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. *KDD-96 Proceedings*, 226-231 (1996).
32. Neve-Oz, Y., Razvag, Y., Sajman, J. & Sherman, E. Mechanisms of localized activation of the T cell antigen receptor inside clusters. *Biochim. Biophys. Acta.* **1853**, 810-821 (2015).
33. Cox, S. et al. Bayesian localization microscopy reveals nanoscale podosome dynamics. *Nat. Meth.* **9**, 195-200 (2012).
34. Lee, S.-H., Shin, J.Y., Lee, A. & Bustamante, C. Counting single photoactivatable fluorescent molecules by photoactivated localization microscopy (PALM). *Proc. Natl. Acad. Sci.* **109**, 17436-17441 (2012).



**Figure 1:** Workflow of the algorithm. a) Schematic of the cluster proposal generating mechanism. i) A raw data set consisting of molecular localisations within a square ROI. ii) Getis's method counts the number of localisations within a distance  $r$  of each point. iii) This allows each localisation to be assigned a score,  $L(r)$ , in pseudo-colour. iv) These values are then thresholded and v) those falling above are grouped into clusters by connecting any pair whose circles intersect. vi) Finally, this allows all localisations to be given a cluster label or to be assigned to the background, culminating in an overall cluster proposal. The algorithm searches through many combinations of  $r$  and  $T$  to generate many thousands of cluster proposals. b) A representative simulated data set from the Standard Conditions. c) A pseudo-coloured heat-map showing the Log Posterior Probability for a range of values of  $r$  and threshold  $T$ . Red represents the most probable combinations according to the model. The dashed line represents  $T = r$ . d) From this map, four  $r$  and  $T$  combinations are selected and the corresponding proposals shown. The highest scoring combination generates Proposal 2.



**Figure 2:** Four different clustering scenarios. i) Standard Conditions, ii) a sparse data set with only 10% as many localisations, iii) clusters which are twice the size and iv) only 10 localisations per cluster and 90% of localisations in the background. a) Representative simulated data. b) Log Posterior Probability heat-maps together with  $T = r$  line and highest scoring combination encircled. c) Highest scoring cluster proposal.



**Figure 3:** Comparison of the clustering behaviour of CD3 $\zeta$ -mEos3 in resting and activated primary human T cells. a) Representative SMLM images, scale bar 5  $\mu$ m. b) Example 3000 x 3000 nm region showing localisation coordinates. c) Log Posterior Probability heat-maps showing the highest scoring  $r$ - $T$  combination for a representative 3000 x 3000 nm ROI. d) Beeswarm plots presenting the percentage of localisations found in clusters (one point per ROI), the number of clusters per region (one point per ROI), the cluster radii (one point per cluster) and the relative density of localisations in clusters compared to the surrounding region (one point per ROI). \* =  $P \leq 0.05$ , \*\* =  $P \leq 0.005$ , \*\*\* =  $P \leq 0.0005$ .

## Online only methods

### Sample preparation

Primary human T cells were isolated from peripheral blood using Lymphoprep (Stemcell) followed by a naïve CD4 T cell negative selection kit (Miltenyi).  $1 \times 10^7$  naïve human CD4 T cells were transfected with 2  $\mu\text{g}$  CD3 $\zeta$ -mEOS3 using an Amaxa system (Lonza). Immune synapses were formed against activating coverslips coated with anti-CD3 (2  $\mu\text{g}/\text{mL}$ ) and anti-CD28 (5  $\mu\text{g}/\text{mL}$ ). Immune synapses were allowed to form for 5 minutes and were then pH-shift fixed (3% paraformaldehyde (PFA)-kPIPES at 80 mM, pH 6.8 for 5 minutes followed by 3%-PFA Borax at 100 mM for 10 minutes). Cells were imaged in phosphate buffered saline (PBS).

### SMLM imaging

SMLM imaging was performed on a Nikon N-STORM microscope using a 100X, 1.49N.A. oil-immersion total internal reflection fluorescence (TIRF) objective. Cells were imaged under TIRF illumination with a 563 nm laser with photo-activation at 405 nm. Fluorescence was collected in the wavelength range 575-625 nm on an Andor iXon DU897U EM-CCD camera. Camera integration time was 30 ms and a total of 20,000 frames were typically recorded. Molecular coordinates were localised using ThunderSTORM software. ThunderSTORM is able to separate multiple overlapped PSFs and compensates for the multiple blinking of individual fluorophores given a user supplied merge time and distance. Following software recommendations, up to four overlapped Gaussian PSFs were allowed. The optimal merge time was computed following the analysis method of Annibale *et al* as described for mEos (see Supplementary Figure 12). This has been shown to generate reliable molecular localisation coordinates free from multiple blinking effects and therefore appropriate for input into our algorithm. The distance is determined by the camera pixel size, in this case 100 nm. The merge time was determined to be three frames (30 ms). The photon threshold for single molecule identification was set at 500. ThunderSTORM corrects for sample drift during the acquisition using an autocorrelation approach. To calculate the localisation precision for each emitter, we selected the method of Quan *et al*, which accounts for the specific noise statistics of EM-CCD cameras.

### Permutation test

Significance levels for the changes in clustering properties of CD3 $\zeta$  in resting and activated primary human T cells were calculated based on a permutation test. This assumes only that the data are independent (in fact, exchangeable) under the null hypothesis.

The test statistic is  $T = |\bar{X}_1 - \bar{X}_2|$ , where  $\bar{X}_1$  and  $\bar{X}_2$  are the means of the two groups of values under analysis,  $X_1$  and  $X_2$ . The p-value of this test is the frequency that  $T^* \geq T$  where  $T^*$  is a simulated test statistic under the null hypothesis. Specifically,  $T^* = |\bar{X}_1^* - \bar{X}_2^*|$ , where  $\bar{X}_1^*$  and  $\bar{X}_2^*$  are two simulated groups constructed by sampling from the pooled values without replacement.

Through the procedure described by Gandy<sup>29, 30</sup>, we were able to bound the probability, due to simulation error, of reporting a p-value to be on the wrong side of the threshold (0.05, 0.005 or 0.0005) to  $1/10^6$ .