



This is an electronic reprint of the original article. This reprint may differ from the original in pagination and typographic detail.

Hoffecker, Ian T.; Yang, Yunshi; Bernardinelli, Giulio; Orponen, Pekka; Högberg, Björn A computational framework for DNA sequencing microscopy

Published in: Proceedings of the National Academy of Sciences

DOI: 10.1073/pnas.1821178116

Published: 04/09/2019

Document Version Peer-reviewed accepted author manuscript, also known as Final accepted manuscript or Post-print

Please cite the original version:

Hoffecker, I. T., Yang, Y., Bernardinelli, G., Orponen, P., & Högberg, B. (2019). A computational framework for DNA sequencing microscopy. *Proceedings of the National Academy of Sciences*, *116*(39), 19282-19287. Article 1821178116. https://doi.org/10.1073/pnas.1821178116

This material is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of the repository collections is not permitted, except that material may be duplicated by you for your research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered, whether for sale or otherwise to anyone who is not an authorised user.

A Computational Framework for DNA Sequencing Microscopy

lan T. Hoffecker^a, Yunshi Yang^a, Giulio Bernardinelli^a, Pekka Orponen^b, and Björn Högberg^{a,2}

^aDepartment of Medical Biochemistry and Biophysics, Karolinska Institutet, S-17177 Stockholm, Sweden; ^bDepartment of Computer Science, Aalto University, FI-00076 Aalto, Finland

This manuscript was compiled on August 16, 2019

We describe a method whereby micro-scale spatial information such as the relative positions of biomolecules on a surface can be trans-2 ferred to a sequence-based format and reconstructed into images 3 without conventional optics. Barcoded DNA polony amplification techniques enable one to distinguish specific locations of a surface 5 by their sequence. Image formation is based on pairwise fusion of 6 uniquely tagged and spatially adjacent polonies. The network of 7 polonies connected by shared borders forms a graph whose topol-8 9 ogy can be reconstructed from pairs of barcodes fused during a polony crosslinking phase, the sequences of which are determined 10 by recovery from the surface and next-gen sequencing. We devel-11 oped a mathematical and computational framework for this principle 12 called Polony Adjacency Reconstruction for Spatial Inference and 13 Topology and show that Euclidean spatial data may be stored and 14 transmitted in the form of graph topology. Images are formed by 15 transferring molecular information from a surface of interest, which 16 we demonstrated in silico by reconstructing images formed from 17 stochastic transfer of hypothetical molecular markers. The theory 18 developed here could serve as a basis for an automated, multiplex-19 able, and potentially super resolution imaging method based purely 20 on molecular information.

21

next gen sequencing | DNA microscopy | polonies | DNA computing | graph theory |

icroscopic imaging has traditionally relied on optics to amplify signals derived from initially confined spatial 2 regions. Exceptions include atomic force microscopy which im-3 ages by utilizing a probe to interact with the sample. DNA has 4 a high information density, with storage levels of 5.5 petabits 5 per cubic millimeter achieved (1), making it an attractive 6 medium for encoding spatial information at microscales. In 7 this paper, we present a theoretical foundation for a spatial information encoding approach that utilizes DNA sequencing 9 and graph theory that could be used to generate whole images. 10 DNA-driven reactions can be coupled to optically-acquired

11 spatial information such as with proximity ligation assay 12 (PLA) (2), and DNA-PAINT (3) where molecular interactions 13 mediated by DNA are discovered using fluorescence. There is 14 also a family of techniques for connecting spatial locations with 15 single cell RNA sequencing data: using a priori knowledge of 16 spatial marker genes to associate unknown genes to approxi-17 mate locations, the *a priori* data being in most cases obtained 18 by microscopy such as with *in situ* hybridization or modelling 19 of spatial expression patterns to retrieve locations of associ-20 ated genes (4-9). Alternatively, direct microscopy-based in 21 situ sequencing methods achieve precise context-sensitive spa-22 tial transcriptomic information without needing to scramble 23 spatial data by dissociation prior to sequencing (10, 11). 24

Encoding spatial information in a way that is preserved 25 in the scrambling during isolation and recovery from *in situ* 26

contexts that can then be read and recovered with sequencing 27 is a major challenge. A few techniques achieve this by encoding 28 spatial information directly into a molecular format, e.g. in the 29 form of DNA read during sequencing along with transcriptomic 30 data. These methods are based on artificial generation of an 31 addressable surface using printing or lithography (12–14). 32

Herein, we describe a computational framework for a 33 method called Polony Adjacency Reconstruction for Spatial 34 Inference and Topology (PARSIFT), for the purpose of encod-35 ing images, for example of the positions of specific molecules 36 relative to others on a 2D plane, directly into a DNA-based 37 format without transduction of information through any other 38 medium without a priori surface addressing. PARSIFT uti-39 lizes the connectivity of vertices in a graph of paired DNA 40 sequences to infer Euclidean spatial adjacency and next-gen 41 sequencing to recover that information a posteriori. 42

Encoding of topological data in DNA sequence format is 43 possible by using DNA barcodes (unique molecular identifiers), 44 i.e. randomized stretches of bases within a sequence of synthetic 45 DNA. Barcodes associated with spatial patches can establish 46 an identity for those locations, each patch distinguishable from 47 another by sequence. A DNA barcode with 10 bases has over 48 a million possible sequences, and larger barcodes can be used 49 to create effectively unique labels in a system. The basic 50 unit of topological data is an edge or association between two 51 adjacent patches by physically linking between their barcodes. 52 Topological mapping with barcoding has been used to infer 53 neural connectomes by building a network from cells sharing 54 common barcodes left by cell-traversing viruses (15) as well 55 as features of DNA origami (16). 56

We can barcode surface patches using polony generation methods like bridge amplification (17), a 2-primer rolling 57

58

Significance Statement

Traditional microscopy is based on the propagation of interactions between light and small scale objects up to larger scales. Such information may be encoded in DNA and transmitted with next-gen sequencing to be later reconstructed and visualized computationally. We provide a mathematical framework and computational proof of concept for a form of DNA-sequencing based microscopy that may be used to construct whole images without the use of optics. Such an approach can be automated in a parallel and multiplexable way that current optical and scanning-based techniques are unable to achieve.

ITH, GB, and BH conceived project. ITH and YY implemented the in silico proof of concept. ITH, YY, PO, and BH developed the mathematical theory. ITH, YY, GB, PO, and BH wrote the manuscript.

The authors declare no conflict of interest.

¹To whom correspondence should be addressed. E-mail: bjorn.hogberg ki.se



Fig. 1. Encoding and recovering metrics through polony adjacency. (a) Seed molecules with unique barcode sequences land randomly on a surface of primers. (b) Local amplification of seed molecules produces sequence-distinct polonies. (c) Saturation of polonies occurs when polonies are blocked from further growth by encountering adjacent polonies, forming a tessellated surface. (d) Random crosslinking of adjacent strands leads to pairwise association of nearby barcodes. (e) Recovery and sequencing of barcode pairs enables reconstruction of a network with similar relative positions of polonies as the original surface.

circle amplification (18), template walking amplification (19), 59 or packing of barcoded beads (20). Unique "seed" strands are 60 captured by primer strands on the surface (Figure 1a) and 61 locally amplified in the immediate vicinity where they landed. 62 This generates numerous distinct patches, or "polonies", of 63 amplified DNA (Figure 1b). Within each, all DNA is derived 64 from a single seed molecule. Any of the above techniques 65 could be applied to our method, though we focus herein on 66 the polony-amplification-by-surface-primers approach. 67

By growing polonies on a surface of primers to saturation 68 (Figure 1c), i.e. when growing polonies encounter the bound-69 aries of other adjacent polonies, a tessellation of neighboring 70 71 polonies forms. Each polony has a limited number of immedi-72 ately adjacent neighboring polonies with their own respective barcodes. Though each patch is associated with a unique 73 sequence according to its parent seed molecule, isolation of 74 this DNA and subsequent sequencing would scramble informa-75 tion about the polony's position and its neighboring polonies. 76 Thus the critical step is to crosslink strands (SI Appendix Fig. 77 S1) from each polony to strands from adjacent polonies (Fig-78 79 ure 1d) in a way that enables both barcodes to be sequenced together in a single read. Recovery of the strands, i.e. strip-80 ping them from the surface followed by next-gen sequencing 81 (by any means including non-optical approaches such as Ox-82 ford Nanopore) thus preserves topological association between 83 neighboring polonies as pairs of barcodes — a complete set of 84 which constitutes the whole topological network of adjacent 85 polonies (Figure 1e). For random seed distributions we show 86 that topological information alone, constrained by being a 87



Fig. 2. Encoding and recovering metrics via topology. (a) Nine seed molecule points distributed randomly on a plane, the induced Voronoi tessellation T (grav lines), its Delaunay diagram D (blue lines), and the untethered graph G. (b) The distribution of Euclidean distances associated with a given topological distance (path with the fewest edges between two points) sampled for random Poisson Delaunay triangulations (5000 samples per topological distance value). (c) Euclidean distances normalized to the average length of a typical Poisson Delaunay edge (Equation 9.9 (21)) plotted versus topological distance for different Poisson intensities, exhibiting linearity between topological and Euclidean distance. (d) The untethered graph: a set of nodes (black) and edges (red) that constitutes the information preserved after dissociation from spatial context. (e) Reconstructed planar embedding of the initially untethered graph (red lines) using the Tutte embedding approach and corresponding Voronoi tessellation (gray lines). (f) Alignment of reconstructed embedding from e with the original Delaunay diagram from a.

2D planar network with known boundary geometry, retains significant spatial metrics of the original distribution. By generating such a mappable surface, we propose that localization of molecules bound to the surface can be done by covalent association with polonies, enabling inference of molecular spatial distributions and construction of images with polonies as pixels.

1. Results and Discussion

A. Voronoi Tessellation as a Model of Polony Saturation. ${\rm The}$ spatial distribution of polonies on a surface, the *a priori* Euclidean information that is not explicitly accessible after isolation, can be preserved by associations between adjacent polony sequences and recovered with sequencing. Information that is 100 available after sequencing and subsequent transformations of 101 that data are then referred to as a posteriori. 102

Assume that seed molecule amplification on a bounded 2D 103 surface, say in the shape of a disk, takes the form of uniform 104

88

89

90

91

92

93

94

95

96

97

98

99

circular growth. At the point of saturation, polonies have 105 amplified to the extent that their expanding boundaries are 106 restricted from further growth, having encountered neighboring 107 polonies. The system of polonies then forms a planar Voronoi 108 109 tessellation T (SI Appendix A), appearing as a characteristic 110 mosaic of polygons with the property that every point within a given cell is closer to its parent seed point than any others. 111 T can also be represented by its plane dual *Delaunay diagram* 112 D = (P, L) whose vertices P are the seed points of T and 113 edges L are the line segments connecting the seed points of 114 adjacent cells (polonies). By the geometric characteristics of 115 T, all the faces of D are triangles (22)(Section 9). 116

We refer to the graph defined purely by its vertices and 117 edges without spatial considerations as the untethered graph. 118 Figure 2a presents a miniature Voronoi tessellation T formed 119 from 9 seed points within a square and its Delaunay diagram 120 D. The unterhered graph G = (V, E) (Figure 2d) is obtained 121 from D by omitting all geometric information, retaining only 122 topological characteristics of the Delaunay diagram D. This 123 includes a topological distance function t(i, j) defined as the 124 fewest number of edges that must be traversed to get from 125 one vertex i to another j, but no other information about 126 the spatial origins of G is explicitly stored (e.g. no Euclidean 127 coordinates of the original points). 128

B. Topological Metrics as a Proxy for Euclidean Metrics. Let $P = \{p_k \mid k = 1, ..., N\}$ be a planar placement of N seed points, resulting from a Poisson-distributed seeding with intensity (i.e. polony density) λ over an area A. Thus $N \approx \lambda A$, and an untethered graph representation G = (V, E) of the true Delaunay diagram D can be obtained by:

 $V = \{1, \dots, N\},\$ $E = \{\{i, j\} \mid \text{barcodes } w_i \text{ and } w_j \text{ co-occur, } i, j = 1, \dots, N\}$

Since sufficiently long barcodes are with high probability 129 unique (SI Appendix B), we treat pairs of barcodes as unique 130 markers of polony adjacency. We postulate that with a suffi-131 ciently dense Poisson-distributed placement P, the topological 132 metric on G (with an appropriate linear scaling) approxi-133 mates well the actual Euclidean metric of the points in P134 (SI Appendix D-E). Figure 2b shows the Euclidean distance 135 distributions for increasing topological distances from a refer-136 ence vertex, for a large collection of Delaunay triangulations 137 of Poisson random point sets. Figure 2c then plots the scaled 138 (Equation 9.9 (21)) average Euclidean distances as a func-139 tion of topological distances for Delaunay triangulations of 140 141 random point sets generated by Poisson processes of increasing intensity λ , showing crucially that the two variables are 142 proportional. 143

On this basis we propose that by finding a proper straight-144 line planar embedding of the untethered graph G we approxi-145 mate also the metric properties of the underlying Delaunay 146 diagram D and the corresponding Voronoi tessellation T. A 147 straight-line embedding of G in a plane is determined by the 148 placement P' of its vertices, from which the line segments 149 L' corresponding to the edges can be deduced, thus denoted 150 as $\langle G, P' \rangle$. Our hidden a priori embedding is the Delaunay 151 diagram $D = \langle G, P \rangle$, and the goal is to approximate this with 152 a good a posteriori embedding $\langle G, P' \rangle$. 153

One constraint on our candidate $\langle G, P' \rangle$ is that it must be planar, i.e. no two edges may cross each other. This is due to the physical assumption that the barcode-pairings correspond 156 to polony adjacencies and thus cannot bridge non-neighboring 157 polonies. There are several efficient algorithms for finding a 158 plane embedding of a planar graph, one of which is the *Tutte* 159 or *barycentric embedding* (23), applicable to Delaunay-diagram 160 type graphs. Another quality constraint is that an average 161 spatial density of the *a posteriori* vertex positions λ' should 162 be obtained from the final distribution with no systematic 163 variation across the reconstructed area. Finally, if we were to 164 generate a new Delaunav triangulation from the reconstructed 165 points (as can be done from any arbitrary set of points), this 166 should produce a similar set of edges as the original untethered 167 graph that was the basis for reconstruction. 168

Our reconstruction approach (flow diagram SI Appendix 169 Fig. S2) starts by determining the outer or boundary face of 170 the Delaunay diagram D underlying the unterhered graph G. 171 This can practically be done by finding the face in any planar 172 embedding of G that has the most vertices with an intermediate 173 planar embedding, because in D all faces except the boundary 174 face are few-vertex triangles. Fixing the placement of the 175 vertices on the boundary face, we then compute positions for 176 the other vertices of G by Tutte's algorithm, which simply 177 places each vertex at the average (barycenter) of its neighbors' 178 positions. In the case of a Delaunay-diagram type graph with 179 the boundary face a convex polygon, this system is guaranteed 180 to be non-degenerate (23), and the result will be a crossing-free 181 straight-line embedding of G. 182

If spatial characteristics of the original Euclidean boundary 183 are known — for instance if we specify that all boundary points 184 must lie on a circle of known radius — then the embedding 185 may also be scaled to match the original Euclidean metrics. 186 Figure 2e shows the Tutte embedding of the unterthered graph 187 (Figure 2d) with boundary points arranged uniformly around 188 the unit disk. For comparison, we have aligned the recon-189 structed graph with the original Delaunay diagram (Figure 2f) 190 by linearly transforming the planar graph to minimize the 191 distance between corresponding vertices. We can see that 192 relative positions are preserved albeit with local distortion 193 that leads to slight displacement of each reconstructed vertex 194 relative to its original seed counterpart. The algorithm thus 195 returns approximate relative spatial positions of polonies from 196 an input of paired polonies. 197

C. Simulation and Reconstruction by Embedding. We simu-198 late the primer lawn as a hexagonally packed disk of area 199 A with M primer sites as the region of interest (ROI) (Fig-200 ure 3a). We simulate a random seeding at a polony density 201 λ by selecting $N = \lambda A$ random sites, followed by pairing of 202 adjacent polony primer sites and scrambling of edge data prior 203 to reconstruction. Figure 3b shows how crosslinking leads 204 to random pairing of adjacent sites, some of which are self-205 pairing events (providing no additional pairing information) 206 and some of which are cross-polony sites that can be used to 207 deduce the presence of a spatial boundary, with the fraction 208 of information-bearing cross-pairs diminishing with the rela-209

tive site density $\rho \stackrel{\text{def}}{=} M/(A\lambda)$ or average number of sites per polony (SI Appendix Fig. S3). The probabilistic nature of the pairing opens up the possibility to miss an existing boundary, particularly when the boundary is small or when ρ is low. A 2000 polony-simulated surface is shown in Figure 3c, and SI Appendix Fig. S4 shows site-linking and the corresponding 210





Fig. 3. Simulation of polony adjacency reconstruction. (a) Lattice diagram of primer lawn and polonies denoted with color and Voronoi cell boundaries. Filled circles indicate seed locations. (b) Illustration of random site pairing between adjacent primer sites. (c) Alignment between *a priori* and *a posteriori* points from a. (d) Larger simulated surface with a polony density $\lambda = 2000$ polonies/unit area and a relative site density $\rho = 50$ sites per polony on average. (e) Reconstructed graph (red lines) and corresponding Voronoi tessellation (gray lines) computed using the Tutte embedding approach from scrambled edges derived from the simulated surface in d.

²¹⁶ Delaunay triangulation of a 500 polony example.

We reconstructed the topological network from the scram-217 bled edges and performed intermediate embedding, boundary 218 face determination, and Tutte embedding (Figure 3c). For 219 this larger reconstruction, spatial uniformity is more appar-220 ent, we see that Voronoi cells take on the approximate size of 221 polonies in the *a priori* surface, observe no obvious systematic 222 changes in mesh density across the length of the ROI, and note 223 the absence of crossed edges. Besides the Tutte embedding 224 strategy, we developed 2 additional approaches for approxi-225 mating Euclidean metrics from the unterhered graph. One 226 is a non-deterministic spring relaxation (24). This approach 227 does not strictly require a crossing-free planar embedding, and 228 can thus lead to provably false positions involving non-planar 229 adjacency, however this feature could also be advantageous 230 231 if natural interpenetration of adjacent polonies leads to such topology. The last approach (SI Appendix F) is based on the 232 notion of topological distance t(i, j) and its role as a proxy 233 for Euclidean distance. We extend the principle of geometric 234 triangulation, whereby the set of distances of a point to other 235 points in a plane can be converted to Cartesian coordinates, to 236 incorporate t(i, j) as a surrogate for Euclidean distance. In one 237 variant of this method, a total topological distance matrix is 238 239 reduced to two principal component vectors approximating the x and y coordinate vectors. In the alternative variant, t(i, j)240 of each vertex are only measured out to peripheral vertices, re-241 ducing systematic distortions. A comparison of reconstructed 242 meshes from the different approaches is shown in SI Appendix 243 244 Fig. S5-S6.

D. Stamping and Image Formation. Knowledge of polony loca tions could be exploited to provide spatial information about

Fig. 4. Voronoi image formation. (a) An image is overlaid on a surface of primer sites. (b) Molecular markers representing different targets (R, G, and B) contact-transferred to the polony surface and each covalently linked to a polony barcode. (c) Monte Carlo sampling to determine if a marker is associated with a given site and if so which target by taking the probability from the RGB value normalized to 1 at the corresponding position in the image. (d) Tallying of markers and empty sites within a polony/Voronoi cell determines the color and brightness of that "pixel". A subsequent image (lower pane) is formed by coloring each cell accordingly. (e) Larger scale reconstruction from scrambled edge data using the Tutte embedding approach with 30,000 polonies. (f) Closeup of e revealing individual Voronoi "bixels"

objects of interest. We devised a basic model of image re-247 construction from the principle of contact or diffusion-based 248 transfer of molecules of interest to the mapped surface, i.e. 249 a kind of molecular stamp. As proof of concept, we use an 250 image (Figure 4a) as a representation of a hypothetical proba-251 bility distribution of 3 types of molecular markers labeled with 252 identifying sequences called "red", "green", and "blue". The 253 image represents a surface of interest that we would like to 254 sample from, for example a cell surface covered in oligo-tagged 255 antibodies, each of which would be coupled enzymatically to 256 a given polony upon contact (Figure 4b) or diffusing RNA 257 molecules like in (14). The color of the image corresponds to 258 the density of such markers and thus the probability that a 259 marker of a particular color is placed on the polony surface. 260 To simulate molecule transfer, the overlaid lattice of primer 261 sites denotes points where a Monte Carlo sampling will occur 262 in the corresponding position in the image. If the image pixel 263 at a given primer site location has an RGB value dominated by 264 red and green for example, then there is a higher probability 265 of that site being occupied by either a green or red marker 266 (Figure 4c). Realistically, molecular transfer introduces distor-267 tion, e.g. from curvature of cell membranes or lateral diffusion 268 of mRNAs. 269

According to the reconstruction procedure, a Voronoi tes-270 sellation is produced from the final set of vertex positions -271 each cell of which constitutes a pixel that can be used to form 272 an image. The final RGB value of the cell can be determined 273 by tallying the markers that have associated with the primer 274 sites in the polony as well as the number of un-associated sites 275 (Figure 4d). The Voronoi-images shown in Figure 4e and f 276 were generated with the scrambling step that removes any 277

spatial information of the original image and reconstructed 278 using our algorithm. Note that global rotation and chirality 279 are not explicitly preserved from the original image. To place 280 this 30,000 pixel image in experimentally relevant terms, we 281 282 point to a recent spatial transcriptomics manuscript (20)283 where circular discs of barcoded 10 μm beads (in their case sequenced optically in situ to obtain sequence addresses) are 284 used to capture transcriptomic data from tissue slices. Ro-285 drigues et al report a typical size of 70,000 10 μm beads per 286 3 mm disk and obtain approximate single-cell resolution (see 287 also SI Appendix H). Image reconstructions from the four 288 approximation approaches are compared in SI Appendix Fig. 289 S5-S6. 290

E. Assessment of Distortion and Precision. We may charac-291 terize reconstruction quality by defining a distortion metric. 292 The *a priori* seed distribution points have a 1-to-1 corre-293 spondence with points in the *a posteriori* reconstruction, and 294 since we generated the *a priori* points ourselves, we can di-295 rectly compare corresponding original and inferred positions 296 by applying a linear transform Tx(P) (rotation, mirroring, 297 scaling, and translation) to the set of reconstructed points 298 that minimizes net displacement between the two distributions. 299 Distortion is thus defined as the set of displacements: Df =300 $(Df_i,...Df_N) \stackrel{\text{def}}{=} d(P,Tx(P')) \mid \min(\sum_{i=1}^N d(p_i,Tx(p'_i))).$ Averaged over multiple runs, we obtain 2D histograms (Figure 5a 301 302 and SI Appendix Fig. S7-S8) of distortion as a function of 303 position in the ROI. Increasing the polony density (λ) re-304 duces average distortion $\overline{Df} = \frac{1}{N} \sum_{i=1}^{N} Df_i$ (Figure 5d and SI Appendix Fig. S10) whereas changes in the site density ρ (Fig-305 306 ure 5f and SI Appendix Fig. S11) has a negligible effect on \overline{Df} 307 except at $\rho < 100$ sites per polony near the point of network 308 disconnection from absent edges. Examining a single simula-309 310 tion (Figure 5b) we can visualize typical distortions, persistent over limited local scales and occurring with greater probability 311 near the boundaries. Analysis of the radial distribution of 312 this instance (Figure 5c) reveals this as a mild systematic 313 worsening near the boundary, an artifact introduced by the al-314 gorithm's treatment of vertices on the boundary. SI Appendix 315 Fig. S9 compares single instance distortions for the different 316 317 reconstruction approaches.

We also characterize reconstruction quality with Leven-318 shtein distance $(lev_{G,G'})$, the number of edits needed to make 319 two graphs identical, between the unterhered graph and set 320 321 of edges derived from a Delaunay triangulation D' generated 322 from the final reconstructed coordinates. Importantly, this metric is based only on a *posteriori* information, so it can 323 be used in an experimental context where knowledge of the 324 underlying distribution is unavailable. It weakly but positively 325 correlates with distortion for a given λ (SI Appendix Fig. S13). 326 $lev_{G,G'}$ grows linearly with λ (Figure 5e and SI Appendix Fig. 327 S10), and like distortion is relatively constant as a function 328 329 of ρ with a transient catastrophic breakdown at low ρ (Figure 5g SI Appendix Fig. S11). We also measured a classical 330 resolution, the full width half maximum (FWHM) of a point 331 spread function (Figure 5h), by sampling the inferred posi-332 tion of a single site (taking its position to be the centroid of 333 whatever Voronoi cell it lands in). Like distortion, FWHM is 334 approximately $\propto 1/\sqrt{\lambda}$ (Figure 5i) indicating that to halve the 335 minimum size of distinguishable features, one should quadru-336 ple λ (SI Appendix Fig. S12). In experimental terms, polonies 337

generated from techniques like template walking amplification, which forms polonies from sites that must be near the packing limit of oligo surface immobilization, can be on the order of nanometers (19) (SI Appendix G).

342

2. Discussion and Conclusion

The three reconstruction methods (Tutte embedding, spring 343 relaxation, and topological distance matrix) succeed in pro-344 ducing approximations of the original seed distributions that 345 can be used to generate images. Tutte embedding exhibited 346 the best estimated algorithmic complexity (based on run time 347 scaling with λ , SI Appendix Fig. S14) making it the fastest 348 technique which becomes significant for large reconstruction 349 problems ($\lambda > 10,000$ polonies/unit area). Both Tutte embed-350 ding and spring relaxation had the lowest distortion levels, with 351 Tutte embedding exhibiting slightly better Df and $lev_{G,G'}$ 352 scaling with λ . Tutte embedding was sensitive to catastrophic 353 failure at low ρ , with singly-connected edges crashing the 354 reconstruction, and all four approaches were sensitive to dis-355 joint subgraphs - making noisy and unconnected graph data 356 a likely challenge for experimental scenarios. SI Appendix 357 Fig. S13 and SI Appendix I discuss our attempts to move 358 towards an algorithm that optimally exploits the available 359 information, and future research should seek to establish a 360 provably maximum-entropy reconstruction that is efficient and 361 deterministic. 362

Along these lines, utilizing information such as the number 363 of self-pairing events could be useful to extract more informa-364 tion and weight edges according to estimated polony size and 365 better control point placement. Alternatively, low-information 366 content self-pairing events could be prohibited through a bipar-367 tite network approach whereby only pairings between A-type 368 and B-type polonies would be allowed (SI Appendix Fig. S15). 369 The bridge amplification approach to polony generation leaves 370 the possibility of doing this with two species of independent 371 primers on the surface and two interpenetrating/overlapping 372 and independently saturated polony surfaces. Another possi-373 ble approach is series growth of polonies. In the basic concept 374 presented in previous sections, a primer of uniform sequence is 375 assumed, however generation of a saturated layer of polonies 376 that could then be used as primers for a subsequent polony 377 generation step would then result in an overlapping of every 378 2nd-layer polony with multiple 1st-layer polonies. This would 379 result in efficient pairing of barcodes without the need for 380 subsequent crosslinking steps. 381

At the time of publication, we are aware of immediately 382 prior works whose contributions are complementary to ours 383 on development of DNA-sequencing based microscopy (25.384 26). The former work experimentally demonstrates DNA 385 microscopy with images of mRNA in cells using locally confined 386 cDNA amplifications and polymerase extension-based fusion 387 of barcodes to connect spatial patches. Their approach differs 388 from ours through the fact that fusion events are used as 389 a direct distance metric, whereas our data instead relies on 390 topology as a proxy for Euclidean metrics. The latter work uses 391 series proximity ligation to associate planar spatial patches 392 and form a network, utilizing a spring relaxation approach for 393 reconstruction. 394

A. Conclusion. PARSIFT is a concept for microscopic image reconstruction using spatial information encoding in DNA base 396

format. We showed an *in silico* proof of concept by construct-397 ing a pipeline for taking decoupled edge data, generated from 398 simulated polony distributions, that are then reassembled into 399 a topological network and embedded in a Euclidean plane, re-400 401 suming spatial characteristics of the original seed distribution. 402 We saw that global distortions are low enough to resolve whole images. We hold that this framework and pipeline for recon-403 struction could be exploited for image acquisition of micro-404 and nano-scale surfaces with molecular libraries of potentially 405 very high multiplicity and with throughput automated in a 406 way that would not be possible with most optical approaches. 407

408 **Supporting Information (SI).** The code is available at 409 https://github.com/Intertangler/parsift

ACKNOWLEDGMENTS. The authors thank Ferenc Fördős for
insightful discussions. This work was supported by Åke Wiberg
Stiftelsen grant for medical research M17-0214 to ITH., the Knut
and Alice Wallenberg project grant 2017.0114, the Knut and Alice
Wallenberg Academy Fellow grant KAW2014.0241 to BH, and by
Academy of Finland grant 311639 to PO.

- Church GM, Gao Y, Kosuri S (2012) Next-generation digital information storage in DNA. Science p. 1226355.
- Söderberg O, et al. (2006) Direct observation of individual endogenous protein complexes in situ by proximity ligation. *Nature Methods* 3(12):995.
- Jungmann R, et al. (2010) Single-molecule kinetics and super-resolution microscopy by fluorescence imaging of transient binding on DNA origami. *Nano letters* 10(11):4756–4761.
- Wang G, Moffitt JR, Zhuang X (2018) Multiplexed imaging of high-density libraries of RNAs with MERFISH and expansion microscopy. *Scientific Reports* 8(1):4847.
- Karaiskos N, et al. (2017) The Drosophila embryo at single-cell transcriptome resolution. Science 358(6360):194–199.
- Satija R, Farrell JA, Gennert D, Schier AF, Regev A (2015) Spatial reconstruction of single-cell gene expression data. *Nature Biotechnology* 33(5):495.
- Achim K, et al. (2015) High-throughput spatial mapping of single-cell RNA-seq data to tissue of origin. *Nature Biotechnology* 33(5):503.
- Halpern KB, et al. (2017) Single-cell spatial reconstruction reveals global division of labour in the mammalian liver. *Nature* 542(7641):352.
- Lein E, Borm LE, Linnarsson S (2017) The promise of spatial transcriptomics for neuroscience in the era of molecular cell typing. *Science* 358(6359):64–69.
- Wang X, et al. (2018) Three-dimensional intact-tissue sequencing of single-cell transcriptional states. Science p. eaat5691.
- Ke R, et al. (2013) In situ sequencing for rna analysis in preserved tissue and cells. Nature Methods 10(9):857.
- Lee JH, et al. (2015) Fluorescent in situ sequencing (fisseq) of rna for gene expression profiling in intact cells and tissues. *Nature Protocols* 10(3):442.
- Crosetto N, Bienko M, Van Oudenaarden A (2015) Spatially resolved transcriptomics and beyond. *Nature Reviews Genetics* 16(1):57.
- 442 14. Ståhl PL, et al. (2016) Visualization and analysis of gene expression in tissue sections by
 443 spatial transcriptomics. *Science* 353(6294):78–82.
- Peikon ID, et al. (2017) Using high-throughput barcode sequencing to efficiently map connectomes. Nucleic Acids Research 45(12):e115–e115.
- Schaus TE, Woo S, Xuan F, Chen X, Yin P (2017) A DNA nanoscope via auto-cycling proximity recording. *Nature communications* 8(1):696.
- Adessi C, et al. (2000) Solid phase DNA amplification: characterisation of primer attachment and amplification mechanisms. *Nucleic Acids Research* 28(20):e87–e87.
- Korfhage C, et al. (2017) Clonal rolling circle amplification for on-chip DNA cluster generation Biology Methods and Protocols 2(1).
- Ma Z, et al. (2013) Isothermal amplification method for next-generation sequencing. Proceedings of the National Academy of Sciences 110(35):14320–14323.
- Rodriques SG, et al. (2019) Slide-seq: A scalable technology for measuring genome-wide
 expression at high spatial resolution. *Science* 363(6434):1463–1467.
- Miles RE (1970) On the homogeneous planar Poisson point process. Mathematical Biosciences 6:85–127.
- de Berg M, van Krefeld M, Overmars M, Cheong O (2008) Computational Geometry: Algorithms and Applications, 3rd Rev. Ed. (Springer-Verlag).
- Tutte WT (1963) How to draw a graph. Proceedings of the London Mathematical Society 3(1):743–767.
- Kamada T, Kawai S, , et al. (1989) An algorithm for drawing general undirected graphs. Information Processing Letters 31(1):7–15.
- Weinstein JA, Regev A, Zhang F (2019) DNA microscopy: Optics-free spatio-genetic imaging by a stand-alone chemical reaction. *Cell.*
- Boulgakov A, Xiong E, Bhadra S, Ellington AD, Marcotte EM (2018) From space to sequence
 and back again: Iterative DNA proximity ligation and its applications to DNA-based imaging.
 bioRxiv.



Fig. 5. Reconstruction quality. (a) 2D histograms of average displacement values binned by relative position in the unit disk ($n = 5000/\lambda$ simulations per histogram) for varied parameters (λ and ρ). (b) Distortion in a single 2000 polony Tutte embedding with lines connecting *a priori* and *a posteriori* vertex locations. Color map indicates line length (max = unit disk diameter 2.0). (c) Radial profile of distortion in b and 5 point moving average (red line). (d) Log-log plot of average displacement versus λ (points single individual simulations reconstructions) and fixed $\rho = 500$ sites per polony showing displacement approximately $\propto 1/\sqrt{\lambda}$. (e) Linear plot of Levenshtein distance ($lev_{G,G'}$) between untethered and *a posteriori* Delaunay graphs as function of polony. (d and e: n = 25 simulations per λ value) (f) Plot of average displacement and (g) plot of $lev_{G,G'}$ each as a function of ρ for two values of λ , error bars represent standard deviation (f and g: n = 25 simulations per point, error bars: standard dev.) (h) Single instance of full width half maximum (FWHM) of *a posteriori* point spread function of a single site. (i) Log-log plot of FWHM versus λ , scaling approximately according to the negative square root of polony density.