

# Generative modeling of biological shapes and images using a probabilistic $\alpha$ -shape sampler

Emily T. Winn-Núñez<sup>1,†</sup>, Hadley Witt<sup>2,3</sup>, Dhananjay Bhaskar<sup>4</sup>, Ryan Y. Huang<sup>5</sup>, Jonathan S. Reichner<sup>2,3</sup>, Ian Y. Wong<sup>2,6</sup>, and Lorin Crawford<sup>7-9†</sup>

1 Division of Applied Mathematics, Brown University, Providence, RI, USA

2 Graduate Program in Pathobiology, Brown University, Providence, RI, USA

3 Division of Surgical Research, Department of Surgery, Rhode Island Hospital, Providence, RI, USA

4 Department of Genetics, Yale School of Medicine, New Haven, CT USA

5 Department of Computer Science, Brown University, Providence, RI USA

6 School of Engineering, Legoretta Cancer Center, Brown University, Providence, RI USA

7 Microsoft Research, Cambridge, MA, USA

8 Department of Biostatistics, Brown University, Providence, RI, USA

9 Center for Computational Molecular Biology, Brown University, Providence, RI, USA

† Corresponding E-mails: emily\_winn-nunez@brown.edu; lcrawford@microsoft.com

## Abstract

Understanding morphological variation is an important task in many areas of computational biology. Recent studies have focused on developing computational tools for the task of sub-image selection which aims at identifying structural features that best describe the variation between classes of shapes. A major part in assessing the utility of these approaches is to demonstrate their performance on both simulated and real datasets. However, when creating a model for shape statistics, real data can be difficult to access and the sample sizes for these data are often small due to them being expensive to collect. Meanwhile, the current landscape of generative models for shapes has been mostly limited to approaches that use black-box inference—making it difficult to systematically assess the power and calibration of sub-image models. In this paper, we introduce the  $\alpha$ -shape sampler: a probabilistic framework for generating realistic 2D and 3D shapes based on probability distributions which can be learned from real data. We

demonstrate our framework using proof-of-concept examples and in two real applications in biology where we generate (i) 2D images of healthy and septic neutrophils and (ii) 3D computed tomography (CT) scans of primate mandibular molars. The  $\alpha$ -shape sampler R package is open-source and can be downloaded at <https://github.com/lcrawlab/ashapesampler>.

## Author Summary

Using shapes and images to understand genotypic and phenotypic variation has proven to be an effective strategy in many biological applications. Unfortunately, shape data can be expensive to collect and, as a result, sample sizes for analyses are often small. Despite methodological advancements in shape statistics and machine learning, benchmarking standards for evaluating new computational tools via data simulation is still underdeveloped. In this paper, we present a probability-based pipeline called the  $\alpha$ -shape sampler which has the flexibility to generate new and unobserved shapes based on an input set of data. We extensively evaluate the generative capabilities of our pipeline using 2D cellular images of neutrophils and 3D mandibular molars from two different suborders of primates.

## Introduction

Shape statistics has become an integral component of several applications within computational biology including medical imaging<sup>1</sup>, geometric morphometrics<sup>2–4</sup>, and cell biology<sup>5,6</sup>. Recently, there has been a focus to develop computational tools that address the subimage analysis problem: given a collection of images or shapes, find the features that best explain the variation between them with respect to a response variable<sup>7</sup>. One example of this type of analysis is identifying the biologically-relevant atomic and residue-level differences between two protein structural ensembles<sup>8</sup>. To date, several approaches have been proposed with the aim to quantify the global variation between images and shapes including some in topological data analysis<sup>9–12</sup>, methods leveraging landmark-based<sup>13–15</sup> or diffeomorphic-based representations<sup>2,16–18</sup>, and tools that use “functional maps” to identify similarities and differences between shapes via a learned latent space<sup>19</sup>.

Despite the many methodological advances being made for the subimage selection problem in shape analysis, there has yet to be a principled framework to assess the power and limitations of these new

tools. Traditionally, there are two common strategies for benchmarking feature selection methods in computational biology: (i) by analyzing real data where there is a “ground truth” about which features are associated with a given phenotype of interest, or (ii) by using simulations where synthetic data is generated such that we know the causal relationship between features and the response variable. Both of these strategies have well-established statistical practices for tabular data (e.g., gene expression in genomics) but they become increasingly difficult to implement when working with shapes. Using data from real biomedical studies for methodological benchmarking is a challenge because shape-based modalities can be hard to collect. On the other hand, when data is able to be collected, sample sizes within studies are usually small, which both compromises the statistical power of the methods being assessed and inhibits the ability to study algorithmic robustness to variance between observations. Lastly, the relationship between shape and phenotype is largely speculative for many biological applications. For example, there have been radiomic studies which have proposed an association between tumor morphology and survival prognostics for patients with glioblastoma, but the exact biological mechanisms connecting the two remains unknown<sup>1,20</sup>.

Simulation studies are an alternative way to evaluate newly developed computational tools in shape analyses. The key to performing these studies is to have an interpretable generative model such that the process for creating synthetic (yet realistic) shapes is well understood. This facilitates the ability to assess how powered a tool is at identifying causal features driving the morphological variation across samples. In general, algorithmic frameworks for generating synthetic shapes consists of two steps: (i) a procedure to generate a point set and (ii) a set of rules for reconstructing a shape from those points. Multiple end-to-end shape generation pipelines have been introduced in the literature but each have their own sets of limitations. For example, to sample random points from a probability distribution over a manifold, one theoretically needs to know the manifold itself which can be impractical to estimate for many applications<sup>21–24</sup>. Recently, there have been machine learning algorithms that have been developed for generating point clouds and reconstructing shapes using dual generators<sup>25</sup>, diffusion-based methods<sup>26</sup>, encoders<sup>27</sup>, and generative adversarial networks<sup>28,29</sup>; but each of these frameworks lack transparency into the generative process for creating new synthetic shapes<sup>30</sup>. From a more mathematical perspective, several methods have been proposed to infer shapes from randomly generated point clouds. Many of these approaches use Čech and Vietoris Rips complexes<sup>31,32</sup>; however, unfortunately, they require tens (and sometimes hundreds) of simplicial complexes to be constructed for each point set resulting in long

84 runtimes. There are 2D shape reconstruction methods based on contours<sup>33</sup> and curves<sup>34</sup>, but their theory  
85 does not directly translate to higher dimensional objects<sup>35</sup>. Lastly, probability-based shape generative  
86 pipelines are still in their infancy and have thus far relied on component vector analysis where parts of 2D  
87 and 3D objects are broken into smaller components and the assembly/connectivity between components  
88 are hidden variables learned by a pre-specified model<sup>36,37</sup>.

89 In this work, we present the  $\alpha$ -shape sampler: a probabilistic framework which takes in a collection of  
90 real shapes or images as input and generates new synthetic ones with features that both quantitatively  
91 and qualitatively resemble data in the input set. Methodologically,  $\alpha$ -shapes require a single numerical  
92 parameter  $\alpha$  for reconstruction which can be interpreted as a measure of shape detail or granularity  
93 (Fig 1). They can also be generated in  $O(P \log P)$  time where  $P$  is the number of points in the point  
94 cloud that is input into the algorithm<sup>38</sup>. As part of our contributions, we introduce a scalable naïve,  
95 data-driven algorithm to estimate the *reach*<sup>39</sup> for a given set of shapes and theoretically relate it to  
96 the numerical  $\alpha$  parameter. Altogether these properties allow our proposed framework to scale and  
97 accommodate the growing sizes of emerging imaging and shape-based databases. It is worth mentioning  
98 that, while the mathematical concept of reach has been used extensively in topological data analysis to  
99 reconstruct shapes and sample point clouds<sup>40,41</sup>, to our knowledge, we are the first to tie it  $\alpha$ -shapes  
100 parameter for an end-to-end generative modeling pipeline. Shape generation using  $\alpha$ -shapes has been  
101 previously studied in two-dimensions where the underlying manifold is known<sup>42</sup> and to learn about  
102 shape boundaries<sup>43,44</sup>; while shape reconstruction with  $\alpha$ -shapes has primarily been studied in three  
103 dimensions<sup>45–47</sup>. They have also been previously used structural biology application in ecology<sup>48,49</sup> but,  
104 overall, the focus of these studies was to understand the interpretation of the parameter of  $\alpha$  itself rather  
105 than attempting to use  $\alpha$ -shapes as a basis to create a framework for generating new data.

106 Throughout the rest of the paper, we will describe the  $\alpha$ -shape sampler using a combination of  
107 probability theory, topology, and tools from differential geometry. We then translate the theoretical  
108 components of the pipeline into a series of algorithmic steps for practical implementation. Finally, we  
109 illustrate the utility of our approach on small proof-of-concept examples (annuli in two dimensions and  
110 tori in three dimensions) and real datasets (neutrophils in two dimensions and primate mandibular molars  
111 in three dimensions). We find that the  $\alpha$ -shape sampler is effective at generating new shapes which honor  
112 major local and global characteristics of realistic data, while also maintaining algorithmic transparency  
113 so that the pipeline can be used for a wide-range of biological applications.

# Results

## Algorithmic overview of the $\alpha$ -shape sampler

Statistically,  $\alpha$ -shapes are convenient because they require a single numerical parameter  $\alpha$  to encode all connectivity information for a point set. For example, in Fig 1a-c, we see that all points are  $\alpha$ -extreme (i.e., on the border); while in Fig 1d, we see that  $\alpha$  becomes large enough such that one point is not  $\alpha$ -extreme and is therefore an *interior* point of the shape. Finally, in Fig 1e, there are three interior points and the rest are boundary or  $\alpha$ -extreme points. An extension of this figure showing different  $\alpha$ -shapes being formed as a function of the number of points sampled from a unit square and the parameter  $\alpha$  can be found in Fig S1. With this theory in mind, a probability distribution on  $\alpha$ -shapes can be explicitly estimated via uniform point sampling on a given (approximate) manifold and then shapes can be constructed from that point set using  $\alpha$  (see Supporting Information). Recent work has investigated using the  $\alpha$  parameter as a shape characteristic<sup>48,49</sup> but, to our knowledge, it has yet to be used for shape generation. This is likely due to the requirement that point sets need to be in general position, a characteristic often not seen in nature. However, we work within the confines of this assumption in return for theoretical soundness, statistical simplicity, and algorithmic transparency.

We will detail our probabilistic generative framework while assuming that we are working with shapes that are  $d = 2$  or 3-dimensions. The  $\alpha$ -shape sampler involves five key steps (see Fig 2a). To begin, the pipeline receives real shapes; throughout the rest of this paper, we will refer to these input data as “reference” shapes. Note that we depict these reference shapes as binary masks in Fig 2, but the  $\alpha$ -shape sampler software can take shape data in any format as input. In the second step, the reference shapes are aligned, scaled (if applicable), and converted to triangular meshes which we treat as simplicial complexes. In the third step, the reference meshes are used in a generative algorithm which, in the fourth step, outputs newly generated shapes in the form of new  $\alpha$ -complexes. In the fifth and final step, these newly generated  $\alpha$ -complexes are converted back into binary masks (or any other data representation), to match the same format as the original input reference data.

We assume that all reference shapes from a phenotypic class (e.g., healthy cells or molars from a given species of primate) have vertices sampled from the same underlying manifold and that the variation observed across shapes within the class stems from a finite sampling of points. With this in mind, the generative algorithm proportion of the  $\alpha$ -shape sampler is comprised of four main steps (see Fig 2b).

143 First, the  $N$  collection reference meshes are input into the algorithm. We represent the  $i$ -th reference  
144 mesh as  $K_i = \{V_i, E_i, F_i, T_i\}$  which is collection of vertices  $V_i$ , edges  $E_i$ , faces  $F_i$ , and tetrahedra  $T_i$  (if  
145 applicable). In the second step, we estimate the reach  $\tau_i$  for every  $i$ -th reference mesh by computing the  
146 distance to edge neighbors and the circumcenter distance to neighboring faces (and tetrahedra for 3D  
147 objects) for each boundary vertex in the complex  $p \in \partial K_i$ . After completing this for all  $N$  reference  
148 shapes, we have a vector of shape-specific reach estimates  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_N)$ . In the third step, we select  
149  $2 \leq J \leq N$  reference shapes from the input dataset which we use as a basis to generate new shapes.  
150 Here, we combine the point clouds from the  $J$  shapes into a joint partial point and take the minimum  
151 between their corresponding values in  $\boldsymbol{\tau}$  to be the reach estimate  $\hat{\tau}_J$ . Next, we sample candidate points  
152 for the newly generated shapes from balls of radius  $\hat{\tau}_J/8$  around vertices in the joint partial point cloud.  
153 A radius of  $\hat{\tau}_J/8$  is chosen to force newly sampled points to remain relatively close to the boundary of the  
154 reference point cloud. Each new candidate point is accepted or rejected according to a probability-based  
155 rule with parameter  $\theta$  (see Materials and Methods). The  $\theta$  parameter is the minimum number of points  
156 in the joint partial point cloud that need to neighbor the new candidate point in order to accept it. It  
157 effectively determines the level of confidence needed to believe that a randomly sampled point is from  
158 the same underlying manifold as the reference data. Once we have the newly sampled point cloud, in  
159 the fourth step of the algorithm, we set  $\alpha = \hat{\tau}_J - \epsilon$ , where  $\epsilon > 0$  is arbitrarily small, and generate the  
160  $\alpha$ -complexes for new shapes. By default, the  $\alpha$ -shape sampler software sets  $J = 2$ ,  $\theta = d$  (i.e., same  
161 dimensions as the input data), and  $\epsilon = 0.001$  (see URLs). Unless otherwise stated, these are the values  
162 that we also use to generate all of the results presented throughout the rest of the paper.

163 There are two important components to the implementation of our pipeline. First, the  $\alpha$ -shape  
164 sampler uses a function to compute the reach for each shape that is completely separate from the shape  
165 generation function (again see Fig 2b). This serves two purposes: (i) it increases computational speed by  
166 avoiding redundant calculations, and (ii) it provides an informal check for potential outlier shapes before  
167 using those shapes as reference inputs for the generative part of algorithm (e.g., this can be done by  
168 empirically assessing the tails of the distribution for  $\boldsymbol{\tau}$ ). Second, setting  $\alpha$  to be just under  $\hat{\tau}_J$  for some  
169 subset of reference shapes guarantees that we will preserve the original homology and most of the local  
170 geometry that is present in the reference dataset without losing any features or generating any atypical  
171 ones. Theoretical details of our implementation are fully detailed in the Materials and Methods and  
172 Supporting Information.

## 173 2D proof-of-concept study with simulated annuli

174 To demonstrate the  $\alpha$ -shape sampler, we begin with a two-dimensional (2D) toy example where we  
 175 simulate  $N = 50$  “real” (i.e., reference) annuli with inner radius  $r = 0.25$ , outer radius  $R = 0.75$ ,  
 176 and thickness equal to  $R - r = 0.5$ . Each reference annulus is constructed by sampling  $P = 500$   
 177 points uniformly from the annulus and then connecting them using true  $\alpha = 0.15$ . The reach value for  
 178 these reference annuli is given by the inner radius of the hole such that  $\tau = 0.25$ . We consider these  
 179 measurements to be the “ground truth” during evaluation.

180 Using the real 2D annuli as input data, we generate another  $N^* = 10$  annuli using the  $\alpha$ -shape sampler.  
 181 Figs 3a and 3b show that the generated annuli preserve the homology of the original reference shapes (i.e.,  
 182 each generated shape is singular connected component and has exactly one hole). To further evaluate  
 183 how “realistic” the geometric characteristics were for the newly generated annuli, we first identified their  
 184  $\alpha$ -extreme points and separated them into two categories: (i) radii less than 0.5 and (ii) radii greater  
 185 than 0.5. The averages of both categories were used to numerically define each generated shape’s inner  
 186 and outer radii, respectively. The thickness of each generated shape was then found by subtracting the  
 187 inner radius from the outer radius. Table S1 gives the root mean square error (RMSE) for each of these  
 188 characteristics when comparing the generated annuli to the real reference annuli. Overall, we see relatively  
 189 low RMSEs (values below 0.01 for each category) which aligns with the aesthetic similarity between the  
 190 shapes seen in Fig 3. It is important to note that the mean estimated reach for the generated annuli  
 191 produced by the  $\alpha$ -shape sampler was  $\hat{\tau} = 0.1749 \pm 0.009$ . While less than the true value of  $\tau = 0.25$ ,  
 192 this is unexpected given that we are estimating the reach from the data directly (rather than estimating  
 193 it using the true radius). Indeed, we would rather our estimate of the reach be smaller than the truth  
 194 and, consequently, have to sample more points rather than our estimate of  $\tau$  be too large and we lose  
 195 geometric information about the shapes.

## 196 3D proof-of-concept study with simulated tori

197 Next, we extend our demonstration of the  $\alpha$ -shape sampler to a three-dimensional (3D) toy example  
 198 where we simulate  $N = 50$  “real” (i.e., reference) tori with major radius  $R = 0.75$  and minor radius  
 199  $r = 0.25$ . Each reference torus was constructed by sampling  $P = 5000$  points uniformly using the  
 200 `alphashape3d` R package<sup>50</sup> with the Computational Geometry Algorithms Library (CGAL)<sup>51</sup> where the

points were connected using a true  $\alpha = 0.25$ . The reach value for these reference tori is  $\tau = 0.5$  which corresponds to the radius of the hole (or tube) of the tori. As with the previous proof-of-concept study using the 2D annuli, we again consider the geometric measurements for the real reference tori to be the “ground truth” during our assessment. Examples of the real reference tori can be found in Fig 3c.

Using the real 3D tori as input data, we generate another  $N^* = 10$  tori using the  $\alpha$ -shape sampler. An example of these generated shapes can be found in Fig 3d. Here, we see that the generated tori qualitatively preserve the homology of the original data where each have one connected component and one hole. To get estimates of the major and minor radii for the generated torus, we start by examining their boundary points. The following relates the major and minor radii for a torus centered at the origin

$$r^2 = \left( \sqrt{x^2 + y^2} - R \right)^2 + z^2$$

where  $(x, y, z)$  are the Cartesian coordinates of the boundary points for the torus. Rearranging the above equation then yields the following relationship

$$(x^2 + y^2 + z^2) = 2R\sqrt{x^2 + y^2} + (r^2 - R^2).$$

By treating  $Y = x^2 + y^2 + z^2$  to be a response variable,  $X = 2\sqrt{x^2 + y^2}$  to be a covariate,  $\beta = R$  to be a coefficient, and  $\varepsilon = (r^2 - R^2)$  to be a residual, the above rewritten equation mirrors a linear model. As a result, we can use ordinary least squares to estimate the appropriate values for  $(\beta, \varepsilon)$ . This then allows us to infer corresponding estimates for the major  $R$  and minor  $r$  radii for each generated torus, respectively.

Table S2 compares the major and minor radii estimates for the tori generated by the  $\alpha$ -shape sampler to same characteristics in the original reference shapes. We see that while the major radius  $R$  is well preserved (RMSE = 0.002), the minor radius  $r$  is slightly larger for the generated shapes (RMSE = 0.02). This result translated to a slightly larger thickness for the generated tori (again see Fig 3d). While generally still close, it does demonstrate a potential shortcoming in our data-driven approach for shape generation where our random sampling algorithm can be prone to accept points outside of the reference boundary, particularly for shapes with smooth surfaces. While this issue may be corrected via some post-processing step to assure that the generated shapes are on a desired scale, we still caution that the probabilistic nature of the  $\alpha$ -shape sampler is not perfect and may lead to a slight distortion of shape



geometry. It is still worth noting that, despite the slightly larger thickness, the mean reach estimate for the generated tori produced by our algorithm was  $\hat{\tau} = 0.402 \pm 0.006$ , lower than the true value of  $\tau = 0.5$ . Again, the lower reach estimate helps us to preserve the majority of geometric and topological characteristics in the generated shapes even if the overall scale is slightly misrepresented.

## Comparison of real and generated shapes based on primary human neutrophils from healthy and septic patients

Human cells display diverse and dynamic morphologies, driven by the rich interplay between the intracellular cytoskeleton and matrix adhesion during cell migration<sup>52,53</sup>. For example, neutrophils are versatile “first responders” of the innate immune system that are rapidly recruited to tissue sites of injury and infection<sup>54</sup>. Neutrophils become adherent and polarized after “activation” by proinflammatory mediators<sup>55</sup>, exhibiting a leading edge with protrusive pseudopods as well as a trailing edge with a contractile uropod<sup>56</sup>. Indeed, such polarized morphologies appear to be correlated with faster neutrophil motility, but can be considerably more heterogeneous for slower moving cells<sup>57</sup>. Further, neutrophils exhibit profound defects in migration and antimicrobial function during sepsis, an aberrant host response to infection that can result in multi-organ failure and death<sup>58</sup>. An unresolved problem is to meaningfully classify neutrophils, since they plastically transition through distinct phenotypic states but also occur as distinct subsets defined by biomarkers and gene expression<sup>59</sup>.

As a first case study, we applied the  $\alpha$ -shape sampler to two-dimensional cell shapes acquired from phase microscopy images of primary human neutrophils. Briefly, neutrophils were isolated from consented healthy donors and septic patients at Rhode Island Hospital (with approval from the Institutional Review Board), then plated at compliant polyacrylamide hydrogel substrates functionalized with fibronectin (see Materials and methods and Witt et al.<sup>60</sup> for more details). Representative cell morphologies were manually traced, converted to binary masks, and then turned into simplicial complexes (similar to what was shown in Fig 2a). The  $\alpha$ -shape sampler was used to synthetically generate additional cells using the default parameters  $J = 2$ ,  $\theta = 2$ , and  $\epsilon = 0.001$ . The training set consisted of approximately  $N = 20$  neutrophil shapes each from the healthy donors and septic patients, which were then used to generate  $N^* = 25$  new neutrophil shapes from each class. Qualitatively, real healthy neutrophils exhibited relatively rounded and compact morphologies (including a uropod<sup>56</sup>) with a typical diameter of  $\sim 10$  microns ( $\mu\text{m}$ ), which were visually similar to the generated healthy neutrophils (Fig 4a). In comparison, real sep-

tic neutrophils exhibited greater ruffling and elongated protrusions relative to real healthy neutrophils, which was visually recapitulated in the generated septic neutrophils (again see Fig 4a). Further, septic neutrophils showed greater spread areas than healthy neutrophils, with diameters approaching 15-20  $\mu\text{m}$ . These differences between healthy and septic neutrophil shapes were captured in the reach estimates produced by the  $\alpha$ -shape sampler, with the healthy neutrophils having mean  $\hat{\tau} = 3.3689 \times 10^{-3} \pm 1.0152 \times 10^{-3}$  compared to the septic neutrophils having mean  $\hat{\tau} = 5.3409 \times 10^{-3} \pm 4.6246 \times 10^{-3}$ . The larger mean  $\tau$  can be explained by the larger variation along the border of the septic neutrophils, while the larger standard deviation reflects the greater single cell heterogeneity in shape.

To further quantify the differences between the real healthy and septic neutrophils and the similarities between real and generated neutrophils, 33 shape characteristics were calculated including area, perimeter length, compactness, and number of protrusions (see Materials and methods and Bhaskar et al.<sup>61</sup> for more details). These vectors were then projected onto a two-dimensional space using a manifold regularized autoencoder (MRAE)<sup>62</sup> as applied to Potential of Heat-diffusion for Affinity-based Transition Embedding (PHATE) coordinates (Fig 4b). In this lower dimensional representation, real healthy neutrophils are roughly grouped together for larger MRAE1, while real septic neutrophils are roughly grouped together for smaller MRAE1; although, there is not a large separation of these two groupings. Moreover, generated healthy neutrophils also group together with real healthy neutrophils for larger MRAE1, while generated septic neutrophils group together with real septic neutrophils for smaller MRAE1. These general trends were confirmed to be independent of the choice of dimension reduction method, including Uniform Manifold Approximation (UMAP)<sup>63</sup>, PHATE<sup>64</sup>, Principal Component Analysis (PCA), and a generalized autoencoder with an Adam optimizer and mean square error loss (Fig S3). For PCA, in particular, the top two principal components were most heavily weighted by area and perimeter in the loadings. Although MRAE is more difficult to interpret due to the nonlinear representation, the components were similarly weighted by area and perimeter but also solidity and circularity (based on inspection of cell shapes for varying MRAE1 and MRAE2).

Additional examination of these shape metrics revealed statistically significant quantitative differences between healthy and septic neutrophils (Fig 4c and Table S3). Notably, healthy real and generated neutrophils had comparable median area of  $\sim 125 \mu\text{m}^2$  ( $P$ -value = 0.066). Moreover, septic real neutrophils had a median area of  $246 \mu\text{m}^2$ , but septic generated neutrophils had a significantly larger median area of  $332 \mu\text{m}^2$  ( $P$ -value =  $1.71 \times 10^{-4}$ ). Similarly, healthy real and generated neutrophils had median

perimeters of  $\sim 47 \mu\text{m}$  ( $P$ -value = 0.189), while septic real neutrophils had a median perimeter of  $75 \mu\text{m}$  and septic generated neutrophils had a median perimeter of  $89 \mu\text{m}$  ( $P$ -value = 0.0032). In comparison, circularity (expressed as a ratio between 0 and 1 describing similarity to a circle, with 1 denoting a perfect circle), solidity (the fraction of the area of the cell over the area of the convex hull), convexity (the ratio of the convex hull perimeter to the cell perimeter), and compactness (the ratio of the diameter of the circle with the same area of the cell to the major axis of rectangular fit) showed statistically significant differences between the real healthy and septic neutrophils that were maintained by the generated healthy and septic neutrophils, but no statistically significant differences between the real and generated neutrophils.

In order to elucidate this discrepancy between real and generated septic neutrophil shapes, we re-examined how the  $\alpha$ -shape generator was sampling from the training set to define a “manifold” based on the union of point clouds from  $J = 2$  reference shapes (see Material and methods). Without perfect alignment, in this setting, the corresponding combined manifold will retain the outermost protruding points associated with both reference shapes, which will bias the generated shape towards larger areas and perimeters. Since septic real neutrophils exhibit pronounced single cell heterogeneity, the inclusion of a few unusually large cells with this pairwise sampling skewed the shape distribution of septic generated neutrophils towards larger areas and perimeters. In comparison, healthy real and generated neutrophils exhibited no statistical difference in any of the measured shape features, likely since they were more homogeneous in shape. It should be noted that the septic neutrophils could include some subsets that are more dysregulated (perhaps prematurely released from the bone marrow) and others that are phenotypically more similar to healthy neutrophils. If so, the presence of this latter subset could obfuscate the separation of healthy and septic neutrophils by morphology.

## Comparison of real and generated shapes based on primate mandibular molars

As a final case study with three-dimensional shapes, we applied the  $\alpha$ -shape sampler to a dataset consisting of  $N = 15$  computed tomography (CT) scans of mandibular molars from two suborders of primates: 8 of these real teeth came from the genus *Microcebus* of the Strepsirrhine suborder and the remaining 7 came from the *Tarsius* of the Haplorhini suborder<sup>2,65,66</sup>. The  $\alpha$ -shape sampler was used to synthetically generate an additional  $N^* = 10$  teeth from each genus using the parameters  $J = 2$ ,  $\theta = 0$ , and  $\epsilon = 0.001$ . In this analysis, we had to set  $\theta = 0$  because the CT scans for each molar came in the form of boundary meshes, which are technically a “hollowed” representation of fully dense 3D objects (see Materials and

methods). This effectively meant that each reference tooth had volumes equal to 0. As a result, we had to avoid setting  $\theta > 0$  to keep the acceptance probability of new candidate points from being nearly 0 (i.e., we would reject nearly 100% of new candidate points).

It is worth briefly noting that the original dataset started with  $N = 10$  *Microcebus* teeth and  $N = 18$  *Tarsius* teeth, respectively. Some of these references were removed from the analysis after we estimated their reach values (see, again, the second step in Fig 2b) and observed some distinct outliers which would affect our ability to generate new and realistic shapes downstream. For the *Microcebus* genus, the teeth we used in our analysis had estimated reach values in the range  $\hat{\tau} \in [0.0242, 0.0793]$ , while the unused teeth had values  $\hat{\tau} = \{0.253, 0.459, 1.597\}$  (somewhere  $10\times$  to  $100\times$  larger than the rest of the data). Similarly, for the *Tarsius* genus, data for our analysis was restricted to teeth with estimated reach values which fell in the range of  $\hat{\tau} \in [0.0241, 0.1124]$ , while the omitted teeth had reaches between  $\hat{\tau} \in [0.9358, 6.6698]$ . When using all teeth, even with proper alignment and scaling, we generated unrealistic shapes (e.g., synthetic teeth with six or eight roots, which do not occur in either species). A key feature of the  $\alpha$ -shape sampler is that it allows users to use the estimated  $\hat{\tau}$  to identify reference shapes that are outliers relative to the rest of input dataset. This can be used to proactively prune reference shapes or use the  $\hat{\tau}$  values post hoc to diagnose why the algorithm produced a shape that does not fit with the original set.

A comparison of the quality controlled real teeth and the generated teeth from the  $\alpha$ -shape sampler can be found in Fig 5a-5d. Overall, we chose this specific collection of molars for our analysis because of the phylogenetic relationship between the *Microcebus* and the *Tarsius* (Fig 5e)<sup>67</sup>. Morphologists and evolutionary anthropologists have previously used this data to understand variations of the paraconid, the cusp of a primitive lower molar. The paraconids do not appear in other genera<sup>68,69</sup> and are only retained by *Tarsius* which allows this genus of primate to eat a wider range of foods<sup>70</sup>. When using these teeth as reference data in our shape generation pipeline, we see that the  $\alpha$ -shape sampler is indeed able to produce newly generated teeth that qualitatively preserve key features shared between both species (e.g., the four roots) as well as recapitulate species-specific variation that is driven by the presence of the paraconids in the *Tarsius*. More specifically, the generated *Microcebus* teeth are missing the distinguished paraconid that is captured in the generated *Tarsius* teeth (again see Fig 5a-5d), repeating the patterns we see in the real data.

To further assess the quality of the shapes produced by the  $\alpha$ -shape sampler, we follow Turner et al.<sup>12</sup> and used Procrustes analysis<sup>71,72</sup> to assign 400 landmarks onto each reference and newly generated

tooth (Materials and methods). The  $(400 \times 3)$ -dimensional matrix of landmark points for each shape was reshaped to a scalar vector of length 1200. This was then projected onto a two-dimensional space using the manifold regularized autoencoder (MRAE) on PHATE coordinates (Fig 5f). As expected, we see the real *Microcebus* and the real *Tarsius* teeth form distinctly separate groups along both MRAE1 and MRAE2. We also see the generated *Microcebus* teeth group together with the real *Microcebus* teeth, while the generated *Tarsius* teeth group together with the real *Tarsius* teeth. These general trends were again confirmed to be independent of the choice of dimension reduction method (Fig S4). For a more quantitative analysis, we also computed the average pairwise Euclidean distance between each tooth group (e.g., Table S4). Here, we observe that the generated *Microcebus* and generated *Tarsius* teeth are nearly twice as close to their respective real groups than to any other group. We attribute the nonzero distance between the generated and real teeth to the fact that we end up accepting all randomly sampled points during our shape generation algorithm (see Materials and methods).

## Discussion

In this paper, we introduced the  $\alpha$ -shape sampler: a probability-based generative model for two-dimensional and three-dimensional shapes. The underlying theoretical innovation of connecting the mathematical concept “reach” with the  $\alpha$  parameter in  $\alpha$ -shapes allows us to implement a data-driven algorithm with the scalability to accommodate the growing sizes of emerging imaging and shape-based databases. We applied our generative pipeline to both 2D and 3D datasets and demonstrated its ability to successfully capture important geometric, morphometric, and topological characteristics of complex objects. In the main text, we focus on demonstrating our generative model when reference shapes are available. This is meant to approximate the reality that the underlying manifold for shapes observed in many biological applications is often unknown. In the Supporting Information, we derive theory and discuss how to generate new shapes when the true manifold is indeed known and available (Fig S5-S12). This includes detailing how one might sample new shapes directly from probability distributions (code for this “exact” approach is also included in our open-source R package; see URLs).

The current implementation of the  $\alpha$ -shape sampler framework offers many directions for future development. For example, there are a few considerations to be made when choosing the  $J$  number of reference shapes and the  $\theta$  threshold for accepting new candidate points in the  $\alpha$ -shape sampler

pipeline. Almost counter-intuitively, the smaller we select  $J$  to be, the more variation there will be in the generated shapes. This is because the joint point cloud starts to converge as the number of  $J$  shapes that are included grows. Additionally, the number of  $J$  reference shapes limits the number of new shapes that can be produced. Combinatorially, we can only generate  $\binom{N}{J}$  new shapes. While this may be seen as a limitation, it also prevents us from augmenting a study with generated shapes that are too far outside of what has been observed in real data. Similarly, when selecting  $\theta$ , our suggestion is to choose  $\theta = d$  (the dimension of the shape space) so that one avoids noisy points and edges around the boundary. The exception to this rule is when the reference shapes are in the form of boundary meshes which are technically a lower dimensional representation of the full shape data. For example, the primate teeth meshes analyzed in the main text are two-dimensional simplices in three dimensions. In this case, we recommend  $\theta = 0$  such that all points are accepted. While this removes the possibility for noise and variation between iterative runs of the  $\alpha$ -shape sampler, even choosing  $\theta = 1$  will result in such a strict threshold of acceptance that the new shape will be a few isolated points scattered in space. We believe this happens because the volume of intersection of a two-dimensional surface mesh with a three-dimensional ball is 0 due to the mesh having Lebesgue measure 0. While the generated shapes may end up being thicker meshes, this can be fixed via post-processing of the data. To avoid this issue, it is best to use shapes that are “filled” in (such as the neutrophil example), but sometimes this is not feasible or practical for the given dataset.

In its current form, the  $\alpha$ -shape sampler performs considerably better when the reference shapes in the input dataset are well aligned. Indeed, alignment was performed with the simulated annuli and tori (Fig 3), as well as with the mandibular molars which included landmarks amenable to unsupervised learning methods (Fig 5). In comparison, neutrophil morphologies lacked such landmarks and so shapes were only centered on their centroids (Materials and methods). Nevertheless, real and generated shapes for healthy neutrophils were statistically similar, since the real morphologies exhibited comparable areas and were relatively compact (Fig 4). However, some generated shapes for septic neutrophils considerably exceeded the corresponding real shapes in area and perimeter, since the  $\alpha$ -shape sampler generates manifolds that retains the outermost protruding points associated with both shapes (Table S3). To address this artifact, we attempted to rescale shapes after generation to match areas and perimeters, which distorted circularity and convexity. Alternatively, aligning neutrophils along their long axis tended to bias towards the generation of more elongated morphologies. It is conceivable that septic neutrophils

with very different morphologies belong to different subsets, and so the generated cell is a chimera based on different subsets without a plausible biological basis. These issues could be addressed in highly heterogeneous populations by sampling a larger number of single cells to limit the biasing effect of outliers, and to discard any generated cells that deviate excessively from the real shape distribution. Future work could also utilize additional information based on cell migration or tractions<sup>60,73,74</sup>, along with single-cell genomics<sup>75</sup> to gain additional insight into septic cell phenotype. Finally, this approach could be effective for other cell types, such as analyzing the epithelial-mesenchymal transition, since the associated spindle-like morphology displays more consistent landmarks for shape alignment<sup>76–78</sup>.

From a statistical perspective, the assumption that all points in the input data point clouds are uniformly distributed over the same underlying manifold may not be suitable for all applications. When points are not uniformly distributed, the calculation of reach becomes less precise because there is too much variance between boundary points. As a result, the  $\tau$  estimate ends up too big in some parts of the point cloud and too small in others, leading to the loss of local geometric information and the possible addition of global topological information, both of which hinder the ability to generate new realistic shapes that properly fit in the same class as the input dataset. Where points are not uniformly distributed, it may be the case that  $\alpha$ -shapes are the appropriate tool for modeling shapes, as was studied in Gerritsen<sup>79</sup>. This is particularly true when points have additional contextual meaning (e.g., molecular structures such as proteins or strands of DNA) or in cases where meshes are very detailed in some areas and less so in others. An immediate future avenue of work is to extend our pipeline to work for weighted  $\alpha$ -shapes<sup>80</sup>, coupled  $\alpha$ -shapes<sup>81</sup>, and  $\beta$ -shapes<sup>82</sup> to fit a broader range of applications.

## URLs

Code for the  $\alpha$ -shape sampler and data simulations is available at <https://www.github.com/lcrawlab/ashapesampler>. Slicer auto3dgm paradigm is available at <https://toothandclaw.github.io/>. Binary masks of the healthy and septic neutrophils and 3D meshes of the primate mandibular molars are available on the Harvard Dataverse at <https://doi.org/10.7910/DVN/K9A0EG>. Scripts to reproduce the results in this paper are also publicly available and can be found at [https://github.com/lcrawlab/ashapesampler\\_paper\\_results](https://github.com/lcrawlab/ashapesampler_paper_results).

# Materials and methods

## Introduction on $\alpha$ -shapes

In this work, we consider a shape to be the simplicial complex approximation of a compact Riemannian manifold embedded in Euclidean space. We use the same definitions for simplices and simplicial complexes as presented in Edelsbrunner and Harer<sup>83</sup>. We also assume that all shapes considered in a given phenotypic class (e.g., healthy septic cells or molars from a given species of primate) have vertices sampled from the same underlying manifold and that the variation observed across shapes within the class stems from a finite sampling of points. When we know the true underlying manifold, we can generate shapes using hierarchical probability distributions (see Supporting Information). The demonstration of the  $\alpha$ -shape sampler in the main text (and what we detail throughout this section) demonstrates how we can generate new shapes when we have data instead of the underlying manifold. Given our applications in the main text, we will derive the details of our probabilistic generative framework while assuming that we are working with shapes that are  $d = 2$  or 3 dimensions; however, also note that the theory we present is generally applicable to larger finite dimensions as well.

We define  $\alpha$ -shapes using Voronoi cells and the Deluanay triangulation. The main motivation behind this choice is that it mirrors how we compute  $\alpha$ -shapes in practice and we believe that this construction provides a more intuitive framing for understanding the parameters in our sampling algorithm. For a more rigorous definition, we refer the reader to Edelsbrunner et al.<sup>38</sup>. To begin, we assume that all points are in general position. That is, in the  $d$ -th dimension<sup>84</sup>, we assume the following:

- No  $d + 1$  points are colinear or coplanar;
- No  $d + 2$  points are cocircular or cospherical;
- No points form a smallest circle or cicumsphere of radius  $\alpha$ ;
- No points lie on the smallest circumsphere of  $d + 1$  other points.

In practice, this assumption is relatively strict and rarely occurs naturally; however, in the Supporting Information, we prove that this assumption holds true in our generative algorithm so long as points are sampled uniformly. In real data applications, users can either ignore the points during the estimation of reach  $\tau$  (e.g., as we do with the primate mandibular molars) or perturb the points slightly to correct for this assumption (e.g., as we do with the segmented images of the neutrophils).



Let  $\mathcal{S}$  denote a set of  $P$  points in  $\mathbb{R}^d$  in general position. The *Voronoi cell* of a point  $p \in \mathcal{S}$  is the set of points in  $\mathbb{R}^d$  for which  $p$  is the closest. We denote the Voronoi cell as the following

$$\mathcal{V}(p) = \{x \in \mathbb{R}^d \mid \|x - p\| \leq \|x - p'\|, \forall p' \in \mathcal{S} - p\}. \quad (1)$$

The *Voronoi diagram* of  $\mathcal{S}$  is then the union of all Voronoi cells and takes up the space of  $\mathbb{R}^d$ . The *Delaunay complex* of  $\mathcal{S}$  is isomorphic to the nerve of the Voronoi diagram. As long as the points of  $\mathcal{S}$  are in general position, the Delaunay complex of  $\mathcal{S}$  is well-defined and forms the convex hull of the points  $\mathcal{S}$  in  $\mathbb{R}^d$ . This is often referred to as the *Delaunay triangulation* of  $\mathcal{S}$  and is denoted by

$$DT(\mathcal{S}) = \left\{ \mathcal{S}^* \subset \mathcal{S} \mid \bigcap_{p \in \mathcal{S}^*} \mathcal{V}(p) \neq \emptyset \right\}, \quad (2)$$

where  $\mathcal{S}^*$  is a subset of points in  $\mathcal{S}$  and  $\emptyset$  represents the empty set. The example in Fig 1 depicts the Delaunay triangulation and the convex hull for a point set. Instead of Voronoi cells which together take up the entire space, we can look at subsets of those cells. Let  $\mathcal{B}_\alpha(p)$  denote a ball of radius  $\alpha$  centered at point  $p$ . Furthermore, let  $\mathcal{R}_p(\alpha) = \mathcal{B}_\alpha(p) \cap \mathcal{V}(p)$  denote the intersection of the Voronoi cell of  $p$  and the ball of radius  $\alpha$  centered at  $p$  (e.g., see the gray shapes in Fig 1). The union of  $\mathcal{R}_p(\alpha)$  for all points  $p \in \mathcal{S}$  form a cover of  $\mathcal{S}$ , the nerve of which forms the  $\alpha$ -complex which we will denote as  $\mathcal{S}_\alpha$ . The boundary of  $\mathcal{S}_\alpha$  defines the  $\alpha$ -shape. Formally, the border is defined by  $\alpha$ -extreme points, which are the points  $p^* \in \mathcal{S}$  such that there exists a ball of radius  $\alpha$  with  $p^*$  on the border where the complement of the disc contains all other points in  $\mathcal{S}$ . In Fig 1a-c, we see that all points are  $\alpha$ -extreme; while in Fig 1d, we see that  $\alpha$  becomes large enough such that one point is not  $\alpha$ -extreme and is therefore an *interior* point of the shape. Finally, in Fig 1e, there are three interior points and the rest are boundary or  $\alpha$ -extreme points.

## Estimating the reach parameter $\tau$

Assume that we have a dataset with  $N$  shapes or images. We will refer to these samples as “reference shapes” from which we will generate new shapes. Let  $K_i = \{V_i, E_i, F_i, T_i\}$  denote the mesh for the  $i$ -th observation in the reference set comprised of a collection of vertices  $V_i$ , edges  $E_i$ , faces  $F_i$ , and tetrahedra  $T_i$  (if applicable). Recall that (i) we assume that all vertices for reference shapes in the same phenotypic class come from the same underlying manifold, and (ii) most real shape and imaging data do not readily

come in the form of  $\alpha$ -shapes or  $\alpha$ -complexes. In order to generate new shapes, we must derive an appropriate point set from the reference shapes (both in terms of location in space and in the total number of vertices) and we must find an appropriate value of  $\alpha$ . To do so, we use the concept of *reach* (denoted by  $\tau$ ) as presented in Aamari et al.<sup>39</sup>, which can also be related to the inverse of the condition number as introduced in Niyogi et al.<sup>85</sup> (see Supporting Information for a formal definition). In practice,  $\tau$  is the minimum distance from the boundary of a shape to its medial axis and can be approximated as either the minimum distance between connected components or the minimum radius of any holes (or voids) in a shape.

At a high level, we estimate the reach  $\tau_i$  for the  $i$ -th reference shape by using the boundary points of its simplicial complex  $p \in \partial K_i$  (i.e., the  $\alpha$ -extreme points in an  $\alpha$ -shape). We do this because the boundary information is all that is relevant to estimating reach. The collection of  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_N)$  values from the  $N$  reference shapes are then used to estimate an appropriate value of  $\alpha$  for the newly generated shapes. Other theoretical methods for estimating reach using an underlying manifold have been proposed<sup>39,86,87</sup>, but we use this approximate estimate to optimize computational speed. By connecting  $\alpha$  to  $\tau$ , we ensure the preservation of major topological and geometric characteristics for the simplicial complex derived from the  $\alpha$  parameter over a point set. The reach estimates  $\tau$  can also be used to sample a point set for the new shapes, both in point set size (i.e., how many vertices we need to sample from the underlying manifold) and in point density. We substitute the minimum number of points needed to preserve the homology of the underlying manifold with an  $\alpha$ -dense cover using the main result in Niyogi et al.<sup>85</sup> (Supporting Information). Once  $\tau$  is derived from the input reference dataset, the appropriate  $\alpha$  can be selected and a new point set can be sampled—the combination of which will allow use to generate new shapes.

Algorithmically, the process of estimating the reach  $\tau_i$  for the  $i$ -th reference shape is done using the following procedure.

- Examining a boundary vertex  $p \in \partial K_i$ , we first learn its distance to neighboring sets of vertices  $q \in \mathcal{N}_i(p)$  by studying the corresponding edges  $E_i$  that are present in the mesh. We save the largest of these distances using the Euclidean distance,  $d_E = \max_{q \in \mathcal{N}_i(p)} \|p - q\|$ .
- Next, we define  $\mathcal{C}_p$  to be the set of circumcenters of all faces in  $F_i$  and tetrahedra  $T_i$  containing  $p$ . These circumcenters are the points at which any three or four points would meet in the Voronoi

513 diagram and, hence, where faces and tetrahedra would form in the resulting  $\alpha$ -complex. We also  
 514 save the largest of these distances  $d_C = 2 \max_{c \in \mathcal{C}_p} \|p - c\|$ . Here, we take twice the value of the  
 515 circumcenter distance in an effort to preserve consistency across dimensions. Recall that for  $d_E$ , we  
 516 consider the entire lengths of edges, not just the midpoints. The circumcenter can be interpreted  
 517 as a rough estimate of a “midpoint” for faces and tetrahedra; as a result, we multiply that value  
 518 by 2 to capture the full “distance”  $d_C$ .

- 519 • Once we have these two distances corresponding to edges and circumcenters involving point  $p$ , we  
 520 take the maximum which we denote as  $d_p = \max(d_E, d_C)$ . Each value  $d_p$  indicates how large  $\alpha$   
 521 needs to be in order to recover the geometric properties in a localized region of the reference mesh.
- 522 • In practice, we find the next furthest point outside of the minimum  $d_p$  range because it establishes  
 523 the largest that  $\alpha$  can be without us losing any geometric information. To do so, we consider the  
 524 set of vertices in  $V_i$  that do not share an edge with  $p$  but are more than  $d_p$  distance away. Formally,  
 525 this set is  $V_p^* = \{v \in V_i \mid \|v - p\| > d_p\}$ . The  $\tau$  value for a given point is computed as

$$526 \quad \tau_p = \min_{s \in V_p^*} \|s - p\|. \quad (3)$$

527 In the event that  $V_p^*$  is empty (e.g., when  $p$  shares an edge, face, or tetrahedra with all other vertices  
 528 in  $V_i$ ), we take  $\tau_p = d_p$ .

- 529 • The reach for the  $i$ -th mesh shape is approximated by

$$530 \quad \tau_i \approx \frac{1}{|\partial K_i|} \sum_{p \in \partial K_i} \tau_p, \quad (4)$$

531 which is the mean  $\tau$  value for all boundary points in the shape where  $|\partial K_i|$  denotes the cardinality  
 532 of the set.

533 Note that other summary statistics could be used in the final step, such as taking the minimum  $\tau_p$  across  
 534 all points, but empirically we find that taking the mean gives robust estimates and keeps outliers from  
 535 artificially deflating the value of  $\tau_i$ . For example, in theory, the true reach estimate would take the  
 536 minimum of  $\tau_p$  over all boundary points; however, a small outlier  $\tau_p$  value would lead to a small  $\tau_i$  when  
 537 we take the minimum and that would result in computational bottlenecks when we later generate shapes.

Therefore, we choose to trade the precise theoretical implementation for computational scale without compromising major shape information. Repeating this procedure for all  $N$  meshes in the dataset yields a collection of estimated reach parameters  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_N)$  which we will use to generate new point clouds and shapes.

## Algorithm for generating new shapes

When generating new shapes, the first task is to create a corresponding point cloud. This step requires developing a method for sampling points from some underlying manifold  $\mathcal{M}$ . Ideally, one could fit a function to each reference shape from a given dataset, average the functions to approximate the true manifold, and then sample new points directly from that manifold via rejection sampling to simulate uniformity. This strategy is similar to what Diaconis et al.<sup>21</sup> illustrates on the torus; however, this same approach is computationally infeasible for modern datasets with tens to hundreds of shapes. One could use techniques from manifold learning to generate point clouds, but the available techniques involve black-box methods such as dual generators<sup>25</sup> and autoencoders<sup>27</sup>. While these approaches have been shown to be useful for assessing predictive models, these do not provide enough interpretability to learn much about the underlying functional representation of the manifold. We could recover a function for each shape using Gaussian processes, as what is done in Albrecht et al.<sup>13</sup>, but to practically implement this strategy, we need to have access to landmarks for each shape. Once we have our point set, we need to find an  $\alpha$  parameter for the shape to dictate how to reconstruct the shape. Most imaging and shape datasets will not be in the form of  $\alpha$ -complexes as the points in many applications are not in general position. As a result, we need an algorithm that can give us both an accurate point cloud from the underlying sub-manifold and the correct parameter for constructing the  $\alpha$ -shape.

Sampling uniformly from balls with radius  $\hat{\tau}/8$  around points in a given reference point cloud allows us stay close to what we assume to be the true manifold without directly calculating the manifold itself. Additionally, while this procedure is not exactly the same as uniform sampling (i.e., points that are closer together will have balls with greater overlap), we conjecture that the overall sampling ends up matching the true density of the point set. Adding a rejection-like step to the sampling scheme then gives the algorithm robustness to outlying points or atypical features that are present in shapes from the reference dataset. We will work with the “approximate manifold” given by the union of balls around the corresponding reference point clouds of radius  $\hat{\tau}/8$ ; call this manifold  $\widehat{\mathcal{M}} \approx \mathcal{M}$ . In practice, we

avoid estimating or calculating the underlying manifold, but we stay true to the given reference data by implementing a “rejection sampling-like” algorithm via the following five step procedure:

1. Choose  $2 \leq J \leq N$  number of reference shapes from the input dataset to serve as references and combine their corresponding point clouds into a joint set denoted as  $\mathcal{Q}$ .
2. Determine the number of candidate points  $\mathbf{y}$  to sample based on a ball of radius  $\hat{\tau}_J/8$  centered around reference points  $\mathbf{x} \in \mathcal{Q}$ . Here,  $\hat{\tau}_J$  is the minimum value in  $\boldsymbol{\tau}$  corresponding to the subset of  $J$  selected reference shapes. Note that this  $\hat{\tau}_J$  value will change depending on the subset of  $J$  reference shapes chosen for the generation of new shapes. The variation of  $\hat{\tau}_J$  across different subsets of reference shapes contributes to the variation observed in newly generated shapes.
3. Given a reference point  $x \in \mathcal{Q}$  in the joint point cloud of the  $J$  reference shapes, sample random candidate points  $\mathbf{y}$  from  $\mathcal{B}_{\hat{\tau}_J/8}(x)$ — that is, sample random points  $\mathbf{y}$  from a small ball of radius  $\hat{\tau}_J/8$  centered at point  $x$ .
4. Calculate the number of additional points in the joint point cloud  $z \in \mathcal{Q}$  that lie within a ball centered at each candidate point  $y$  which we define as  $p_Q(y) = \#\{z \in \mathcal{Q} \mid z \in \mathcal{B}_{\hat{\tau}_J/4}(y)\}$ . This number does not include the original reference point  $x$  from the previous step. Next, choose  $\theta \leq p_Q(y)$  to be the minimum number of points needed to accept each new candidate point  $y$ . This sets up the following accept-reject decision rules for the generation of new shapes where:
  - If  $p_Q(y) \geq J\theta$ , accept point  $y$ .
  - If  $p_Q(y) < J\theta$ , accept point  $y$  with rate  $1 - \exp\{-2(p_Q(y) - \theta)/J\theta\}$ .
  - If  $p_Q(y) < \theta$ , reject point  $y$ .

We detail the logic behind this rejection rule below.

5. Repeat these steps for all points in the combined point cloud  $x \in \mathcal{Q}$ .

There are a few key takeaways in the procedure specified above. First, we sample new points uniformly from one ball at a time rather than from the union of balls. This means that the new point cloud will reflect the density of the combined point cloud  $\mathcal{Q}$  from the subset of  $J$  reference shapes. Second, to add some variance to the sampled point cloud and to ensure confidence in the newly sampled points, we

593 implement the following rejection-like rule:

$$594 \quad f(y) = \begin{cases} 0 & p_Q(y) < \theta \\ 1 - \exp\{-2(p_Q(y) - \theta)/J\theta\} & \theta \leq p_Q(y) < J\theta \\ 1 & p_Q(y) \geq J\theta \end{cases} \quad (5)$$

595 where, again,  $p_Q(y)$  is the number of points in the joint point cloud  $\mathcal{Q}$  that are within a  $\hat{\tau}_J/4$  radius of the  
 596 candidate point  $y$ ;  $\theta$  is the minimum number of points we require from the reference point cloud  $\mathcal{Q} - x$  to  
 597 be within a  $\hat{\tau}_J/4$  radius of the candidate point  $y$  in order to accept  $y$  (as the reference point  $x$  is already  
 598 within that radius by definition); and  $J$  is again the number of reference shapes. Note that in Eq (5),  
 599 the choice of  $J$  will affect the rate of acceptance and will approach 1 as  $p_Q(y) \rightarrow J\theta$ . The three-part  
 600 rule in Eq (5) is designed to accommodate three scenarios when we consider to accept a newly sampled  
 601 point  $y$ . If  $p_Q(y) < \theta$ , then there are fewer neighboring reference points than desired and indicates that  
 602 the candidate point  $y$  is likely to be far away from the boundary of the point cloud. We have little  
 603 confidence that these points are from the manifold that we wish to approximate  $\widehat{\mathcal{M}}$  and so, consequently,  
 604 we reject these points. In the scenario where  $p_Q(y) \geq J\theta$ , the candidate point  $y$  is near more than  $\theta$  real  
 605 points (on average) from the  $J$  reference shapes. In this case, we have high confidence that  $y$  is from the  
 606 approximated manifold  $\widehat{\mathcal{M}}$  and automatically accept it as a newly sampled point.

607 In the middle scenario, where  $\theta \leq p_Q(y) < J\theta$ , we want a rule that allows for some uncertainty in  
 608  $y$  as a function of the number of nearby points  $p_Q(y)$  from the  $J$  reference shapes. Here, we choose  
 609  $1 - \exp\{-2(p_Q(y) - \theta)/J\theta\}$ , which is the cumulative distribution function (CDF) for an exponential  
 610 random variable with rate  $J\theta/2$  that is shifted to be 0 when  $p_Q(y) = \theta$  (i.e., the threshold for the  
 611 minimum number of points needed to accept each new candidate point  $y$ ). The exponential distribution  
 612 is typically used to model the amount of time until some specific event occurs—where there are fewer  
 613 large values and more small values. The main motivation behind this choice is to reward candidate points  
 614  $y$  that with higher values of  $p_Q(y)$ . When we have  $J = 2$  reference shapes, the rate of the distribution will  
 615 be  $\theta$ ; as we add more reference shapes to the algorithm, the rate at which we find more neighboring points  
 616 for any candidate point  $y$  will increase. In practice, using our proposed rejection-like rule, the acceptance  
 617 rate will be roughly 100% for randomly drawn candidate points that are near the interior of the point  
 618 cloud (particularly in regions where the  $J$  reference shapes being used all overlap). Intuitively, the rate of

acceptance will decrease for new candidate points that are sampled near the boundaries of the  $J$  reference shapes. The range of the overall acceptance probability will depend on the intraclass heterogeneity of the reference dataset and the quality of alignment of the point clouds during preprocessing.

## **Patient blood sample collection and primary neutrophil isolation**

Blood was drawn from healthy donors or septic patients with written informed consent at Rhode Island Hospital, in accordance with the guidelines and approval of the Institutional Review Board. Briefly, healthy donors had no known acute infection or chronic systemic disease within one month prior to the blood draw. We did not collect blood from minors, pregnant women, prisoners, mentally retarded or mentally disabled patients or volunteers. Septic patients from the surgical intensive care unit (ICU) and the trauma ICU displayed at least two systemic inflammatory response syndrome criteria with a source of infection, and enrolled within 48 hours of their diagnosis or admission. Patients also had to be at least 18 years of age without a massive blood transfusion. Further details on study design are documented elsewhere<sup>60</sup>.

For both healthy donors and septic patients, 10-30 milliliters (mL) of blood was collected in EDTA-containing Vacutainer tubes. Buffy coat was separated by centrifugation with Histopaque-1077 with an additional sedimentation step for neutrophils using 3% Dextran (400-500 kDa). Any contaminating erythrocytes were eliminated by hypotonic lysis, and neutrophils were then resuspended in cation-free HBSS media.

## **Polyacrylamide gel preparation and neutrophil imaging**

Briefly, polyacrylamide gel substrates were polymerized on a 25 millimeters (mm) glass coverslip, using 3% acrylamide and 0.2% bisacrylamide for a Young's modulus of  $E = 1.5$  kPa, along with fluorescent red 0.5  $\mu\text{m}$  carboxylate-modified polystyrene beads. Gel substrates were then coated with human fibronectin (Gibco 33016015) using the photoactivatable crosslinker sulfo-SANPAH (Sigma 803332) and rinsed extensively. Further experimental details are documented elsewhere in Oakes et al.<sup>73</sup> and Witt et al.<sup>60</sup>, respectively. The polyacrylamide gel and coverslip were mounted in a coverslip holder, then covered with 1 mL of Leibovitz L-15 media. About 50,000 neutrophils were plated and allowed to adhere for 15 minutes. Approximately 20-60 adherent cells were imaged in phase microscopy using a Nikon TI-2 epifluorescent microscope using a 40X air objective with a 0.6 numerical aperture. An Okolab enclosure

around the TI-2 maintained the apparatus at 37° and 5% CO<sub>2</sub> for the duration of the experiments. Only adherent cells were selected for imaging. The  $N$  represents the number of individual neutrophils imaged and analyzed, with an  $n > 3$  for individual septic or healthy donors.

## Converting segmented neutrophil images to 2D simplicial complexes

To convert tif files into two-dimensional simplicial complexes, we used a multi-step procedure. For the healthy neutrophils, each image was first cropped to include only the middle 50%. Septic neutrophil images were already cropped. Next, the centroid of each shape was found using the median row and column; cells were centered by placing this centroid at the center of the new matrix. The black-and-white cell images were converted into a binary matrix representing black-and-white pixels. This matrix was then searched to find all the black pixels, which were used as vertices for the complex. To add randomness to the pixel points, all vertices were also perturbed within their pixel areas. Next, edges were formed by finding pairs of vertices that were either orthogonally or diagonally adjacent according to the matrix. However, in order to avoid overlapping edges, the upper left and downward right diagonals of each vertex were removed except when upper right and downward left diagonals could not exist (such that the overlap would be impossible). Finally, every three edges that could form a triangle were listed as a face to construct a group of adjacent faces, which was plotted to generate a 2D simplicial complex for the image.

## Evaluation of generated neutrophils

Representative cell morphologies were manually traced, converted to binary masks, and then turned into simplicial complexes (Fig 2a). The  $\alpha$ -shape sampler was used to synthetically generate additional cells with parameters  $J = 2$  and  $\theta = 2$ . These newly generated neutrophils were then converted to binary masks. We computed 33 geometric characteristics using the masks of the original and the generated shapes, respectively, including: area, perimeter length, number of protrusions, compactness, and others as described in Bhaskar et al.<sup>61</sup>. The vectors of these characteristics were projected onto a two-dimensional latent space using a manifold regularized autoencoder (MRAE)<sup>62</sup> where the loss function is the combination of a mean square error loss on the autoencoder itself and the “Potential of Heat-diffusion for Affinity-based Transition Embedding” (PHATE) coordinates in latent space. This combined loss function



is formally defined as the following

$$\mathcal{L}(\cdot) = \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 + \sum_{i=1}^N \sum_{i'=1}^N | \|z_i - z_{i'}\| - \|\phi(x_i) - \phi(x_{i'})\| |, \quad (6)$$

where  $\mathcal{L}(\cdot)$  denotes the loss function associated with the autoencoder;  $N$  is the number of shapes in the dataset;  $x_i$  is the input data for the  $i$ -th shape;  $\|\cdot\|$  is the  $L^2$ -norm;  $\hat{x}_i$  is the reconstructed version of the  $i$ -th shape as determined by the decoder portion of the MRAE;  $z_i$  is the two-dimensional latent embedding for the data associated with the  $i$ -th shape; and  $\phi(\cdot)$  is the PHATE function mapping the shape data to  $\mathbb{R}^2$ . The idea behind the loss function is to train an autoencoder to not only minimize the difference between the input and reconstructed data, but also force the latent space to behave as similarly as possible to the PHATE function  $\phi$ . Since PHATE is a dimensionality reduction method designed to honor the original local and global structure of high-dimensional data<sup>64</sup>, adding the extra loss component based on the PHATE coordinates in the latent space forces the autoencoder to also honor the original structure of the data as well.

In addition to the MRAE, we also assess the new shapes generated by the  $\alpha$ -shape sampler using other dimensionality reduction approaches including: the uniform manifold approximation projection (UMAP)<sup>63</sup>, PHATE, principal component analysis (PCA), and a generic autoencoder. Each of these analyses were used to demonstrate that our conclusions about the shapes produced by the  $\alpha$ -shape sampler are robust regardless of the unsupervised dimension reduction method that we choose. Briefly, UMAP was implemented with 5 nearest neighbors, 2 connected components, Euclidean distance, and a minimum distance set to 0.1. PHATE was implemented with 5 nearest neighbors, 2 connected components, a Von Neumann Entropy diffusion operator, log potential, Euclidean distance, and we used stochastic gradient descent for the multi-dimensional scaling method. Both the autoencoder and the MRAE were trained with 500 epochs.

## Evaluation of generated primate manibular molars

To generate synthetic primate manibular molars, we used parameters  $J = 2$  and  $\theta = 0$  in the  $\alpha$ -shape sampler software, which meant an automatic 100% acceptance rate of sampled points. Since the reference teeth data were given as two-dimensional surface meshes in three-dimensional space, they had volumes equal to 0. In this case, setting  $\theta > 0$  would send the acceptance probability of new candidate points to

nearly 0 (i.e., we would reject nearly 100% of new candidate points). Our evaluation for the generated shapes with this dataset were similar to the landmarking and subsequent dimensionality reduction analyses used in Turner et al.<sup>12</sup>. First, the reference teeth were aligned using the software package `auto3dgm`<sup>88</sup>. We then generated 10 new synthetic teeth each from the *Microcebus* and the *Tarismus* genera, respectively. We used Procrustes analysis<sup>71,72</sup> to assign 400 landmarks to each newly generated tooth so that these could also be aligned and scaled. The (400×3)-dimensional matrices of landmark points for both the newly generated and real reference teeth were reshaped to scalar vectors of length 1200. These were then projected onto a two-dimensional space using the same manifold regularized autoencoder (MRAE) and other dimensionality reduction techniques (UMAP, PHATE, PCA, and an autoencoder) as was done the neutrophils. UMAP was implemented with 5 nearest neighbors, 2 connected components, Euclidean distance, and a minimum distance set to 0.1. PHATE was implemented with 5 nearest neighbors, 2 connected components, a Von Neumann Entropy diffusion operator, log potential, Euclidean distance, and we used stochastic gradient descent for the multi-dimensional scaling method. Both the autoencoder and the MRAE were trained with 500 epochs. For quantitative results, we calculate Euclidean distances between the length 1200 scalar vectors representing each tooth and gather the pairwise distances to reaffirm that the generated teeth are appropriately spaced from the original reference datasets.

## Acknowledgements

This research was conducted using computational resources and services at the Center for Computation and Visualization (CCV), Brown University.

## Funding

This work was supported by the National Science Foundation Graduate Research Program (1644760 to ETWN), the National Institutes of Health (T32HL134625 and F31DE02874 to HW, R01AI116629 to JSR), Department of Surgery in the Rhode Island Hospital (to JSR), a Yale-Boehringer Ingelheim Biomedical Data Science Fellowship (to DB), the Army Research Office (W911NF2310385 to IYW), and a David & Lucile Packard Fellowship for Science and Engineering (to LC). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

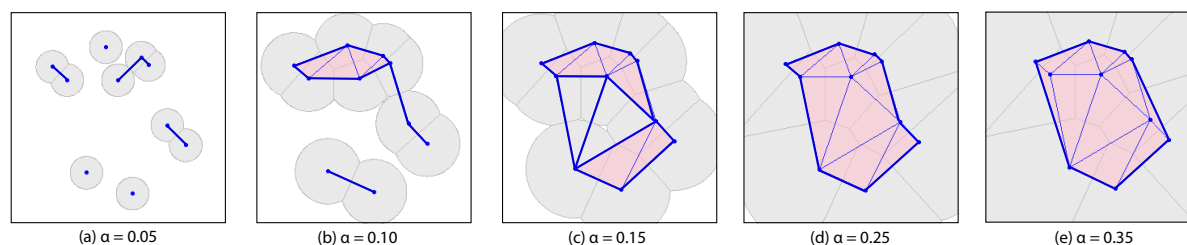
## 727 Competing Interests

728 DB is currently supported by a Boehringer Ingelheim Fellowship at Yale University. All other authors  
729 have declared that no competing interests exist.

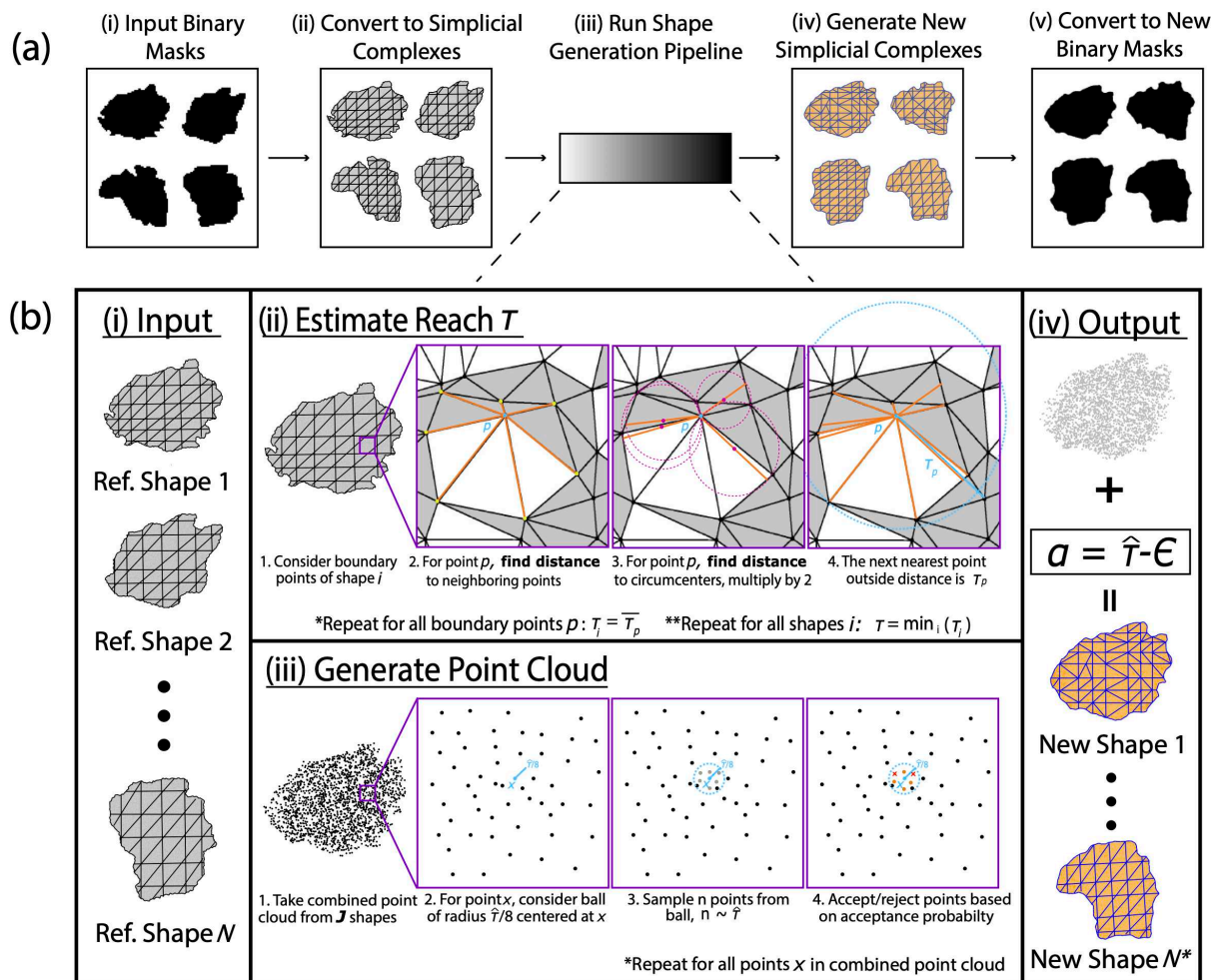
## 730 Author Contributions

- 731 • **Conceptualization:** Emily T. Winn-Núñez, Ian Y. Wong, Lorin Crawford
- 732 • **Data Curation:** Hadley Witt, Jonathan S. Reichner, Lorin Crawford
- 733 • **Formal Analysis:** Emily T. Winn-Núñez
- 734 • **Funding Acquisition:** Emily T. Winn-Núñez, Hadley Witt, Dhananjay Bhaskar, Jonathan S.  
735 Reichner, Ian Y. Wong, Lorin Crawford
- 736 • **Investigation:** Emily T. Winn-Núñez, Lorin Crawford
- 737 • **Methodology:** Emily T. Winn-Núñez, Dhananjay Bhaskar, Lorin Crawford
- 738 • **Project Administration:** Lorin Crawford
- 739 • **Resources:** Jonathan S. Reichner, Ian Y. Wong, Lorin Crawford
- 740 • **Software:** Emily T. Winn-Núñez
- 741 • **Supervision:** Jonathan S. Reichner, Ian Y. Wong, Lorin Crawford
- 742 • **Validation:** Emily T. Winn-Núñez, Ryan Y. Huang, Jonathan S. Reichner, Ian Y. Wong, Lorin  
743 Crawford
- 744 • **Visualization:** Emily T. Winn-Núñez, Ryan Y. Huang, Ian Y. Wong, Lorin Crawford
- 745 • **Writing - Original Draft:** Emily T. Winn-Núñez, Ian Y. Wong, Lorin Crawford
- 746 • **Writing - Final Draft:** Emily T. Winn-Núñez, Hadley Witt, Dhananjay Bhaskar, Ryan Y. Huang,  
747 Jonathan S. Reichner, Ian Y. Wong, Lorin Crawford

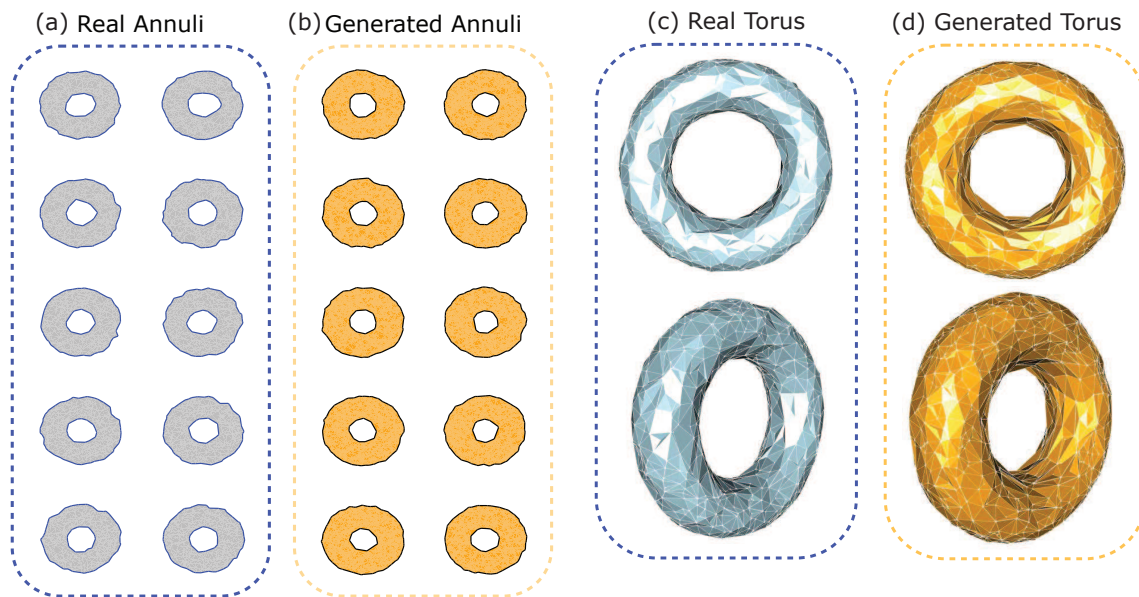
# Figures and Tables



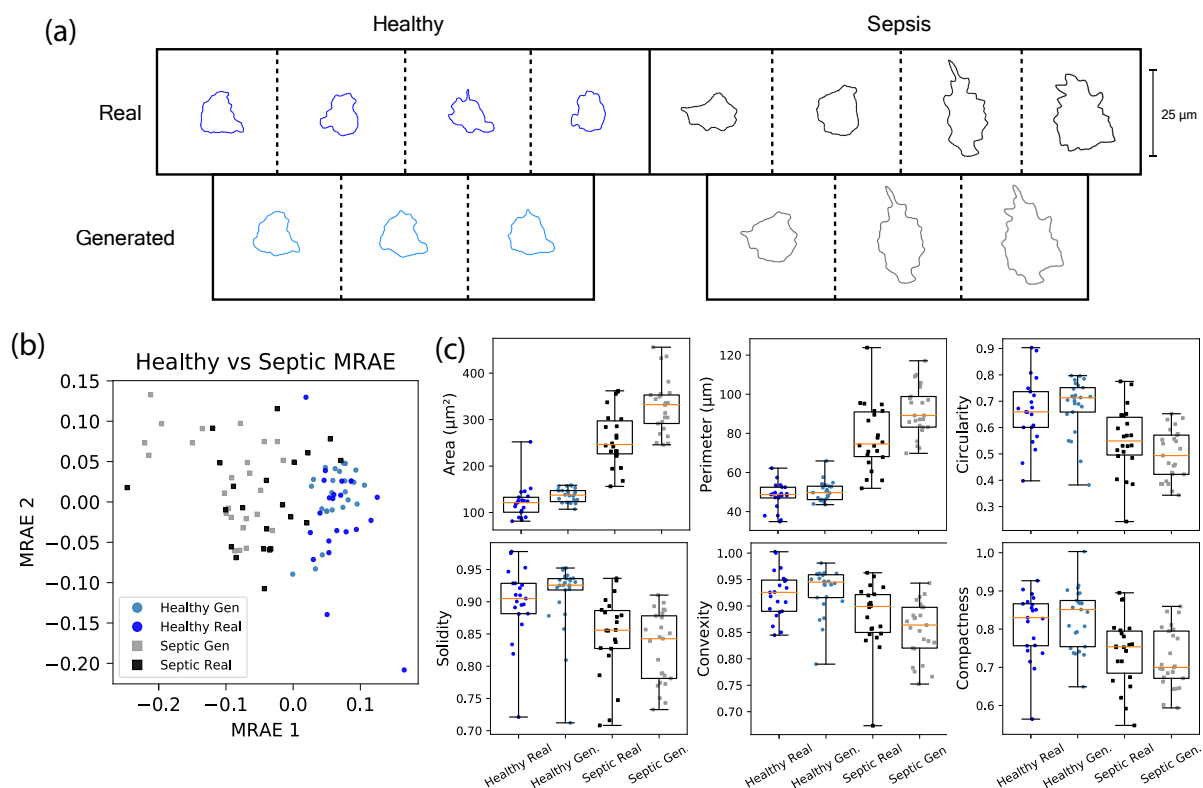
**Figure 1. An example of various  $\alpha$ -shapes for the same set of points under different choices for the numerical parameter  $\alpha$ .** Here, we consider different parameter values (a)  $\alpha = 0.05$ , (b),  $\alpha = 0.10$ , (c)  $\alpha = 0.15$ , (d)  $\alpha = 0.2$ , and (e)  $\alpha = 0.35$ . In each panel, the gray shapes are the intersection of balls of radius  $\alpha$  and the Voronoi cells at each point. The pink triangles are then faces representing the collective interior, and the blue lines are edges of the  $\alpha$ -complex. The bold blue edges are known as the “boundary edges” and denote the  $\alpha$ -shape for each panel. In (a) and (b), where  $\alpha$  is smaller, we have disconnected components. In (c), we see an instance where edges may form the boundary of a face, but the face is not quite yet filled in since the three Voronoi cells have not collectively met. In (d), the faces are filled in and one of the points becomes an interior point while the rest remain  $\alpha$ -extreme points. In (e),  $\alpha$  is large enough such that the given  $\alpha$ -complex is the Delaunay triangulation and convex hull of the point set. When determining how to generate a new shape from an existing dataset, we use information within the given simplicial complex to determine how many points are needed, where the points should be sampled, and the appropriate  $\alpha$  parameter to connect the points. For a more detailed overview and theoretical discussion of concepts surrounding  $\alpha$ -shapes, see Materials and Methods and Supporting Information.



**Figure 2. Schematic overview of the  $\alpha$ -shape sampler: a probabilistic framework for simulating realistic 2D and 3D images and shapes.** (a) A general illustration of the pre- and post-processing workflow in the  $\alpha$ -shape sampler software. In step (i), the user inputs data of real shapes in some format—in this case, binary masks for illustration. We refer to these data as “reference” shapes. In step (ii), the reference masks are converted to triangular meshes which are treated as simplicial complexes. In step (iii), the reference meshes are input into the shape generation pipeline which, in step (iv), outputs newly generated shapes in the form of  $\alpha$ -complexes. Finally, in step (v), these generated  $\alpha$ -complexes are converted back to match the same format as the original input data (again, here, binary masks). (b) Details underlying the algorithm for generating new shapes via the  $\alpha$ -shape sampler. (i) A collection meshes from  $N$  reference shapes are given to the software. For simplicity, we assume that these shapes are from the same phenotypic class and, thus, their points are from the same manifold. (ii) Next, we estimate the reach  $\tau_i$  for each reference shape by computing the distance to edge neighbors for each point (i.e., vertex in the mesh) and the circumcenters to neighboring faces (note that we also evaluate tetrahedra for 3D objects). The next closest vertex is the value  $\tau_p$  for point  $p$ , and the smallest  $\tau_p$  among all points is the value of  $\tau_i$  for the  $i$ -th reference shape. We then take the minimum  $\tau = (\tau_1, \dots, \tau_N)$  to be the representative estimate of the reach  $\hat{\tau}$  for all reference shapes. (iii) We create a partial point cloud by combining points from  $J$  reference shapes in our input dataset, where  $2 \leq J \leq N$ . Next, we sample new points from a ball of radius  $\hat{\tau}/8$  around vertices in the partial point cloud. Each new point is accepted or rejected according to a probability-based rule. (iv) Once we have the newly sampled point cloud, we set  $\alpha = \hat{\tau} - \epsilon$ , where  $\epsilon > 0$  is arbitrarily small, and generate the  $\alpha$ -complexes for new shapes.

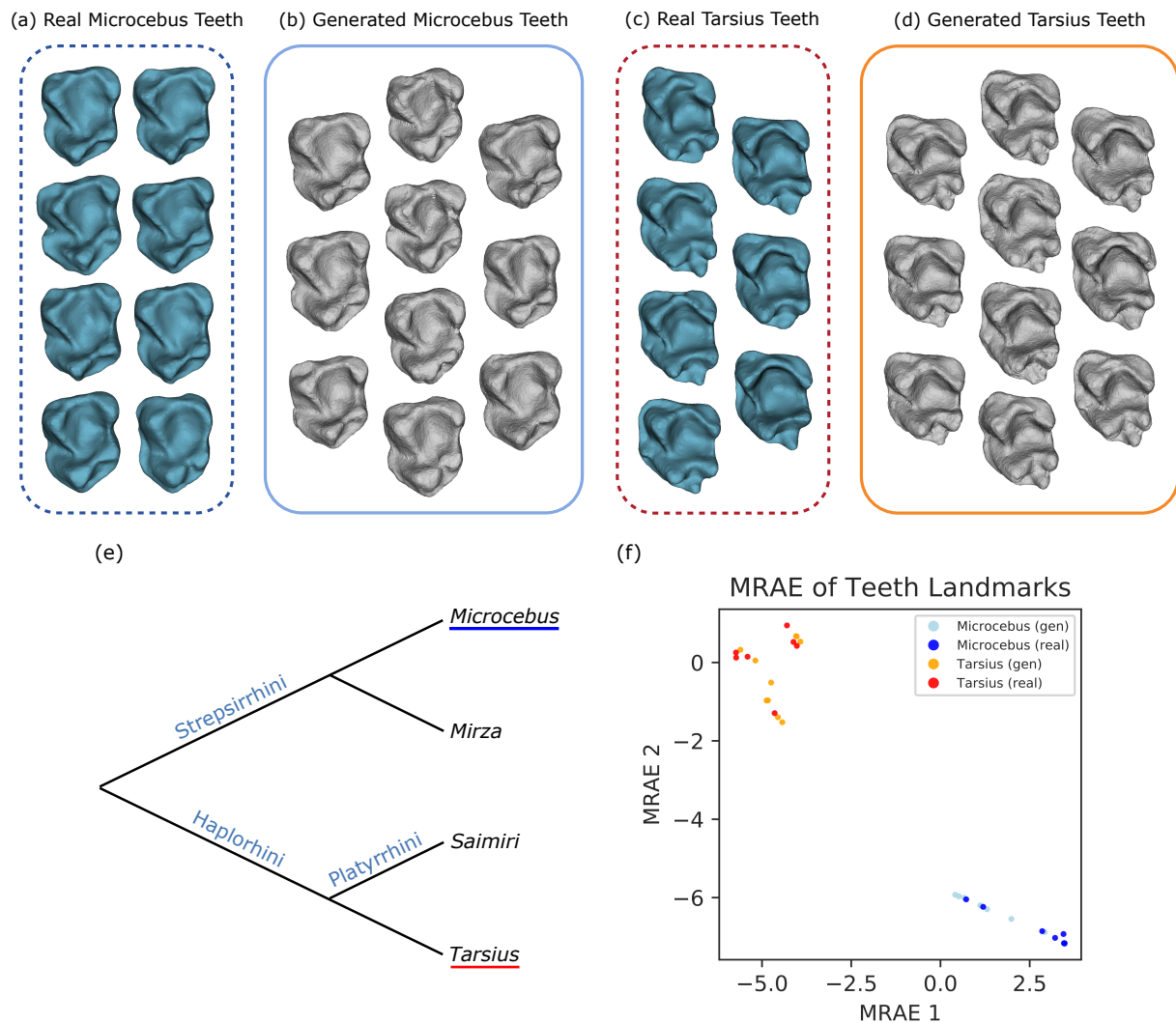


**Figure 3. Qualitative comparisons of real and generated 2D annuli and 3D tori using the  $\alpha$ -shape sampler.** Panels (a) and (b) show real (gray) and generated (orange) annuli. Similarly, in panels (c) and (d), we show real (gray) and generated (orange) tori. Overall, we see that the  $\alpha$ -shape sampler generates slightly thicker shapes than the examples in the original dataset (see Tables S1 and S2 for a quantitative evaluation). Nonetheless, the generated shapes preserve the most important topological property in that they all have exactly one connected component and exactly one hole.



**Figure 4. Application of the  $\alpha$ -shape sampler to generate synthetic 2D images of healthy and septic neutrophils.** (a) Examples of real healthy (blue), generated healthy (light blue), real septic (black), and generated septic (gray) neutrophils in gels with stiffness 1.5 kilopascals (kPa). Each synthetic neutrophil in the second row was generated using the two shapes it sits in between in the row above. Variation in the newly generate cells can be most seen along the boundary, which is a function of the sampling process in the  $\alpha$ -shape pipeline. When comparing the generated and real cells, perhaps most noticeable are (i) the differences in area and (ii) the number of protrusions in the healthy versus septic cells. (b) We use a manifold regularized autoencoder (MRAE) to show that the generated shapes cluster and intermix with real cells in their respective categories. This provides evidence that the images being generated by the  $\alpha$ -shape sampler are realistic. (c) We compute the area, perimeter, circularity, solidity, convexity, and compactness of each real and generated cell. Next, we compare the distribution of these measurements for the healthy and septic groups, respectively. Here, if the  $\alpha$ -shape is able to preserve geometric and morphological characteristics while generating new data, then we would expect the distributions of these measurements to line up within a group. Note that due to the high heterogeneity and difficulty aligning shapes, the generated septic neutrophils are slightly larger in area and perimeter than the real ones. However, the generated neutrophils with the  $\alpha$ -shape sampler still capture other key shape characteristics.





**Figure 5. Application of the  $\alpha$ -shape sampler to generate synthetic 3D primate mandibular molars.** Here, we qualitatively compare meshes of (a) real *Microcebus*, (b) generated *Microcebus*, (c) real *Tarsius*, and (d) generated *Tarsius* teeth. Morphologically, we know that tarsier teeth have an additional high cusp (highlighted in red) which allows this genus of primate to eat a wider range of foods<sup>70</sup>. Here, we see that the generated *Tarsius* teeth from the  $\alpha$ -shape sampler preserve the unique paraconids. In panel (e), we show the phylogenetic relationship between the *Microcebus* and *Tarsius* genus. It has been estimated that the divergence dates of the *Microcebus* and *Mirza* from *Tarsius* happened around five million years before the branching of *Tarsius* from *Saimiri*<sup>67</sup>. (f) We use a manifold regularized autoencoder (MRAE) to show that the generated teeth cluster and intermix with the real *Microcebus* and *Tarsius* teeth, respectively. Figure S4 shows that the same results hold regardless of the dimensionality reduction technique that is used.



# References

1. Lorin Crawford, Anthea Monod, Andrew X. Chen, Sayan Mukerhjee, and Raúl Rabadán. Predicting clinical outcomes in glioblastoma: An application of topological and functional data analysis. *Journal of the American Statistical Association*, 115:1139–1150, 2020. doi: <https://doi.org/10.1080/01621459.2019.1671198>.
2. Doug M. Boyer, Yaron Lipman, Elizabeth St. Clair, Jesus Puente, Biren A. Patel, Thomas Funkhouser, Jukka Jernvall, and Ingrid Daubechies. Algorithms to automatically quantify the geometric similarity of anatomical surfaces. *Proceedings of the National Academy of Sciences*, 108(45):18221–18226, 2011. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1112822108. URL <https://pnas.org/doi/full/10.1073/pnas.1112822108>.
3. Tobias Theska, Bogdan Sieriebriennikov, Sara S Wighard, Michael S Werner, and Ralf J Sommer. Geometric morphometrics of microscopic animals as exemplified by model nematodes. *Nature Protocols*, pages 1–34, 2020.
4. Kory M. Evans, Olivier Larouche, Sara-Jane Watson, Stacy Farina, María Laura Habegger, and Matt Friedman. Integration drives rapid phenotypic evolution in flatfishes. *Proceedings of the National Academy of Sciences*, 118(18):e2101330118, 2021. doi: 10.1073/pnas.2101330118. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2101330118>.
5. Zachary Pincus and JA Theriot. Comparison of quantitative methods for cell-shape analysis. *Journal of Microscopy*, 227(2):140–156, 2007.
6. JC Caicedo, S Cooper, F Heigwer, S Warchal, P Qiu, C Molnar, AS Vasilevich, JD Barry, HS Bansal, O Kraus, M Wawer, L Paavolainen, MD Herrmann, M Rohban, J Hung, H Hennig, J Concannon, I Smith, PA Clemons, S Singh, P Rees, P Horvath, RG Linington, and AE Carpenter. Data-analysis strategies for image-based cell profiling. *Nat Methods*, 14(9):849–863, 2017.
7. Bruce Wang, Timothy Sudijono, Henry Kirveslahti, Tingran Gao, Douglas M. Boyer, Sayan Mukherjee, and Lorin Crawford. A statistical pipeline for identifying physical features that differentiate classes of 3d shapes. *The Annals of Applied Statistics*, 15, 2021. doi: 10.1214/20-AOAS1430.

- 775 8. Wai Shing Tang, Gabriel Monteiro da Silva, Henry Kirveslahti, Erin Skeens, Bibo Feng, Timothy  
776 Sudijono, Kevin K Yang, Sayan Mukherjee, Brenda Rubenstein, and Lorin Crawford. A topo-  
777 logical data analytic approach for discovering biophysical signatures in protein dynamics. *PLoS*  
778 *computational biology*, 18(5):e1010045, 2022.
- 779 9. Kun Meng, Jinyu Wang, Lorin Crawford, and Ani Eloyan. Randomness and statistical inference  
780 of shapes via the smooth euler characteristic transform. *arXiv*, 2023. URL [http://arxiv.org/](http://arxiv.org/abs/2204.12699)  
781 [abs/2204.12699](http://arxiv.org/abs/2204.12699).
- 782 10. Qotung Jiang, Sebastian Kurtek, and Tom Needham. The weighted euler curve transform for  
783 shape and image analysis. In *2020 IEEE/CVF Conference on Computer Vision and Pattern*  
784 *Recognition Workshops (CVPRW)*, pages 3685–3694, 2020. ISBN 978-1-72819-360-1. doi: [https:](https://doi.ieeecomputersociety.org/10.1109/CVPRW50498.2020.00430)  
785 [//doi.ieeecomputersociety.org/10.1109/CVPRW50498.2020.00430](https://doi.ieeecomputersociety.org/10.1109/CVPRW50498.2020.00430).
- 786 11. Peter Bubenick and Paweł Dłotko. A persistence landscape toolbox for topological summary statis-  
787 tics. *Journal of Symbolic Computation*, 78:91–114, 2017. doi: [https://doi.org/10.1016/j.jsc.2016.](https://doi.org/10.1016/j.jsc.2016.03.009)  
788 [03.009](https://doi.org/10.1016/j.jsc.2016.03.009).
- 789 12. K. Turner, S. Mukherjee, and D. M. Boyer. Persistent homology transform for modeling shapes  
790 and surfaces. *Information and Inference*, 3(4):310–344, 2014. ISSN 2049-8764, 2049-8772. doi: 10.  
791 1093/imaia/iau011. URL [https://academic.oup.com/imaia/article-lookup/doi/10.1093/](https://academic.oup.com/imaia/article-lookup/doi/10.1093/imaia/iau011)  
792 [imaia/iau011](https://academic.oup.com/imaia/article-lookup/doi/10.1093/imaia/iau011).
- 793 13. Thomas Albrecht, Marcel Lüthi, Thomas Gerig, and Thomas Vetter. Posterior shape models.  
794 *Medical Image Analysis*, 17(8):959–973, 2013. ISSN 13618415. doi: 10.1016/j.media.2013.05.010.  
795 URL <https://linkinghub.elsevier.com/retrieve/pii/S1361841513000844>.
- 796 14. Joshua Cates, Shireen Elhabian, and Ross Whitaker. Shapeworks: particle-based shape correspon-  
797 dence and visualization software. In *Statistical shape and deformation analysis*, pages 257–298.  
798 Elsevier, 2017.
- 799 15. Dennis Madsen, Andreas Morel-Forster, Patrick Kahr, dana Rahbani, Thomas Vetter, and Marcel  
800 Lüthi. A closest point proposal for mcmc-based probabilistic surface registration. In *European*  
801 *Conference on Computer Vision - ECCV*, pages 281–296, 2020. doi: [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-030-58520-4_17)  
802 [978-3-030-58520-4\\_17](https://doi.org/10.1007/978-3-030-58520-4_17).

- 803 16. Tingran Gao, Shahar Z Kovalsky, and Ingrid Daubechies. Gaussian process landmarking on man-  
804 ifolds. *SIAM Journal on Mathematics of Data Science*, 1(1):208–236, 2019.
- 805 17. Tingran Gao, Shahar Z Kovalsky, Doug M Boyer, and Ingrid Daubechies. Gaussian process land-  
806 marking for three-dimensional geometric morphometrics. *SIAM Journal on Mathematics of Data*  
807 *Science*, 1(1):237–267, 2019.
- 808 18. Yi Hong, Polina Golland, and Miaomiao Zhang. Fast geodesic regression for population-based  
809 image analysis. In *Medical Image Computing and Computer Assisted Intervention- MICCAI 2017:*  
810 *20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings,*  
811 *Part I 20*, pages 317–325. Springer, 2017.
- 812 19. Ruqi Huang, Panos Achlioptas, Leonidas Guibas, and Maks Ovsjanikov. Limit shapes—a tool for  
813 understanding shape differences and variability in 3d model collections. In *Computer Graphics*  
814 *Forum*, volume 38, pages 187–202. Wiley Online Library, 2019.
- 815 20. Lee Curtin, Paula Whitmire, Haylye White, Kamila M Bond, Maciej M Mrugala, Leland S Hu, and  
816 Kristin R Swanson. Shape matters: morphological metrics of glioblastoma imaging abnormalities  
817 as biomarkers of prognosis. *Scientific reports*, 11(1):23202, 2021.
- 818 21. Persi Diaconis, Susan Holmes, and Mehrdad Shahshahani. Sampling from a manifold. In *Institute*  
819 *of Mathematical Statistics Collections*, pages 102–125. Institute of Mathematical Statistics, 2013.  
820 ISBN 978-0-940600-84-3. doi: 10.1214/12-IMSCOLL1006. URL [http://projecteuclid.org/](http://projecteuclid.org/euclid.imsc/1379942050)  
821 [euclid.imsc/1379942050](http://projecteuclid.org/euclid.imsc/1379942050).
- 822 22. Paul Breiding and Orlando Marigliano. Random points on an algebraic manifold. *SIAM Journal*  
823 *on Mathematics of Data Science*, 2(3):683–704, 2020. ISSN 2577-0187. doi: 10.1137/19M1271178.  
824 URL <https://epubs.siam.org/doi/10.1137/19M1271178>.
- 825 23. Ferenc Fodor, Dániel I. Papvári, and Viktor Vígh. ON RANDOM APPROXIMATIONS BY  
826 GENERALIZED DISC-POLYGONS. *Mathematika*, 66(2):498–513, 2020. ISSN 0025-5793, 2041-  
827 7942. doi: 10.1112/mtk.12027. URL [https://onlinelibrary.wiley.com/doi/abs/10.1112/](https://onlinelibrary.wiley.com/doi/abs/10.1112/mtk.12027)  
828 [mtk.12027](https://onlinelibrary.wiley.com/doi/abs/10.1112/mtk.12027).

- 829 24. Kun Meng and Ani Eloyan. Principal manifold estimation via model complexity selection. *Journal*  
830 *of the Royal Statistical Society Series B: Statistical Methodology*, 83(2):369–394, 2021.
- 831 25. Cheng Wen, Baosheng Yu, and Dacheng Tao. Learning progressive point embeddings for 3d point  
832 cloud generation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
833 *(CVPR)*, pages 10261–10270. IEEE, 2021. ISBN 978-1-66544-509-2. doi: 10.1109/CVPR46437.  
834 2021.01013. URL <https://ieeexplore.ieee.org/document/9578874/>.
- 835 26. Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel  
836 diffusion. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 5806–5815.  
837 IEEE, 2021. ISBN 978-1-66542-812-5. doi: 10.1109/ICCV48922.2021.00577.
- 838 27. Lyne P. Tchapmi, Vineet Kosaraju, Hamid Rezatofighi, Ian Reid, and Silvio Savarese. TopNet:  
839 Structural point cloud decoder. In *2019 IEEE/CVF Conference on Computer Vision and Pattern*  
840 *Recognition (CVPR)*, pages 383–392. IEEE, 2019. ISBN 978-1-72813-293-8. doi: 10.1109/CVPR.  
841 2019.00047. URL <https://ieeexplore.ieee.org/document/8953650/>.
- 842 28. Ruihui Li, Xianzhi Li, Ka-Hei Hui, and Chi-Wing Fu. Sp-gan: Sphere-guided 3d shape generation  
843 and manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021.
- 844 29. Haoxu Zhang, Chenchen Qiu, Chao Wang, Bin Wei, Zhibin Yu, Haiyong Zheng, and Juan Li.  
845 Learning spectral normalized adversarial systems with stacked structure for high-quality 3d object  
846 generation. *Concurrency and Computation: Practice and Experience*, 33, 2021. doi: [https://](https://doi-org.revproxy.brown.edu/10.1002/cpe.5430)  
847 [doi-org.revproxy.brown.edu/10.1002/cpe.5430](https://doi-org.revproxy.brown.edu/10.1002/cpe.5430).
- 848 30. Ashley Mae Conard, Alan DenAdel, and Lorin Crawford. A spectrum of explainable and in-  
849 terpretable machine learning approaches for genomic studies. *Wiley Interdisciplinary Reviews:*  
850 *Computational Statistics*, page e1617, 2023.
- 851 31. Brittany Terese Fasy, Rafal Komendarczyk, Sushovan Majhi, and Carola Wenk. On the recon-  
852 struction of geodesic subspaces of  $\mathbb{R}^n$ . *International Journal of Computational Geometry & Appli-*  
853 *cations*, 32(1):91–117, 2022. ISSN 0218-1959, 1793-6357. doi: 10.1142/S0218195922500066. URL  
854 <https://www.worldscientific.com/doi/10.1142/S0218195922500066>.

- 855 32. Tyrus Berry, Timothy Sauer, and ,Department of Mathematical Sciences, Fairfax, VA 22030, USA.  
856 Consistent manifold representation for topological data analysis. *Foundations of Data Science*, 0  
857 (0):0–0, 2019. ISSN 2639-8001. doi: 10.3934/fods.2019001. URL [http://aimsciences.org/](http://aimsciences.org/article/doi/10.3934/fods.2019001)  
858 [/article/doi/10.3934/fods.2019001](http://aimsciences.org/article/doi/10.3934/fods.2019001).
- 859 33. Giovanni Bellettini, Valentina Beorchia, Maurizio Paolini, and Franco Pasquarelli. *Shape Recon-*  
860 *struction from Apparent Contours: Theory and Algorithms*. Computational Imaging and Vision.  
861 Springer Berlin Heidelberg, 1 edition, 2015. ISBN 978-3-662-45191-5.
- 862 34. Hao Guo, Feng Ju, Dongming Bai, Xiaoyong Wei, Lingyu Wang, and Bai Chen. Shape recon-  
863 struction for continuum robot based on pythagorean hodograph-bézier curve with imu and vision  
864 sensors. *IEEE Sensors Journal*, 23:8535–8344, 2023. doi: 10.1109/JSEN.2023.3248781.
- 865 35. Chengjie Niu, Yang Yu, Zhenwei Bian, Jun Li, and Kai Xu. Weakly supervised part-wise 3d shape  
866 reconstruction from single view rgb images. *Computer Graphics Forum*, 39:447–457, 2020. doi:  
867 10.1111/cgf.14158.
- 868 36. Henrik Lieng. A probabalistic framework for component-based vector graphics. *Computer Graphics*  
869 *Forum*, 36:195–205, 2017. doi: 10.1111/cgf.13285.
- 870 37. Evangelos Kalogerakis, Siddhartha Chaudhuri, Daphne Koller, and Vladlen Koltun. A probabal-  
871 istic model for component-based shape synthesis. *ACM Transactions on Graphics*, 31:1–11, 2012.  
872 doi: 10.1111/cgf.13285.
- 873 38. H. Edelsbrunner, D. Kirkpatrick, and R. Seidel. On the shape of a set of points in the plane. *IEEE*  
874 *Transactions on Information Theory*, 29(4):551–559, 1983. ISSN 0018-9448. doi: 10.1109/TIT.  
875 1983.1056714. URL <http://ieeexplore.ieee.org/document/1056714/>.
- 876 39. Eddie Aamari, Jisu Kim, Frédéric Chazal, Bertrand Michel, Alessandro Rinaldo, and Larry Wasser-  
877 man. Estimating the reach of a manifold. *Electronic Journal of Statistics*, 13(1), January 2019.  
878 ISSN 1935-7524. doi: 10.1214/19-ejs1551. URL <http://dx.doi.org/10.1214/19-EJS1551>.
- 879 40. Larry Wasserman. Topological data analysis. *Annual Review of Statistics and Its Application*, 5:  
880 501–532, 2018.

- 881 41. Eddie Aamari and Clément Levrard. Stability and minimax optimality of tangential delaunay com-  
882 plexes for manifold reconstruction. *Discrete & Computational Geometry*, 59(4):923–971, February  
883 2018. ISSN 1432-0444. doi: 10.1007/s00454-017-9962-z. URL <http://dx.doi.org/10.1007/s00454-017-9962-z>.  
884
- 885 42. Fausto Bernardini and Chandrajit L. Bajaj. Sampling and reconstructing manifolds using alpha-  
886 shapes. *Department of Computer Science Technical Reports*, 1997. URL <https://docs.lib.purdue.edu/cstech/1350>. Paper 1350.  
887
- 888 43. Ery Arias-Castro and Alberto Rodríguez-Casal. On estimating the perimeter using the alpha-shape.  
889 *Annales de l'Institut Henri Poincaré - Probabilités et Statistiques*, 53(3):1051–1068, 2017.
- 890 44. Harish Chintakunta and Hamid Krim. Distributed boundary tracking using alpha and delaunay-  
891 ~cech shapes. *arXiv preprint arXiv:1302.3982*, 2013.
- 892 45. Peer Stelldinger. Topologically correct surface reconstruction using alpha shapes and relations to  
893 ball-pivoting. In *2008 19th International Conference on Pattern Recognition*, pages 1–4. IEEE,  
894 2008.
- 895 46. Leonid Tcherniavski and Peer Stelldinger. A thinning algorithm for topologically correct 3d surface  
896 reconstruction. In *Proc. 8th IASTED Int. Conf. Vis., Imag., Image Process.(VIIP)*, page 119, 2008.
- 897 47. Frédéric Cazals and David Cohen-Steiner. Reconstructing 3d compact sets. *Computational Geom-*  
898 *etry*, 45(1-2):1–13, 2012.
- 899 48. James D. Gardiner, Julia Behnsen, and Charlotte A. Brassey. Alpha shapes: determining 3d shape  
900 complexity across morphologically diverse structures. *BMC Evolutionary Biology*, 18(1):184, 2018.  
901 ISSN 1471-2148. doi: 10.1186/s12862-018-1305-z. URL <https://bmcevolbiol.biomedcentral.com/articles/10.1186/s12862-018-1305-z>.  
902
- 903 49. Electa Cleveland, Angela Zhu, Bjorn Sandstede, and Alexandria Volkening. Quantifying different  
904 modeling frameworks using topological data analysis: a case study with zebrafish patterns. *arXiv*,  
905 2022. URL <http://arxiv.org/abs/2212.12034>.
- 906 50. Thomas Lafarge and Beatriz Pateiro-Lopez. alphashape3d: Implementation of the 3d alpha-shape

- 907 for the reconstruction of 3d sets from a point cloud. *CRAN*, 2023. URL <https://CRAN.R-project.org/package=alphashape3d>.
- 908
- 909 51. The CGAL Project. *CGAL User and Reference Manual*. CGAL Editorial Board, 5.6 edition, 2023.
- 910 URL <https://doc.cgal.org/5.6/Manual/packages.html>.
- 911 52. Alex Mogilner and Kinneret Keren. The shape of motile cells. *Curr Biol*, 19(17):R762–R771, 2009.
- 912 53. William R Holmes and Leah Edelstein-Keshet. A comparison of computational models for eukary-
- 913 otic cell shape and motility. *PLoS Comp Biol*, 8(12):e1002793, 2012.
- 914 54. Pei Xiong Liew and Paul Kubes. The neutrophil’s role during health and disease. *Phys Rev*, 99
- 915 (2):1223–1248, 2019.
- 916 55. Andrew E Ekpenyong, Nicole Toepfner, Edwin R Chilvers, and Jochen Guck. Mechanotransduction
- 917 in neutrophil activation and deactivation. *Biochim Biophys Acta Mol Cell Res*, 1853(11):3105–3116,
- 918 2015.
- 919 56. LE Hind, WJ Vincent, and A Huttenlocher. Leading from the back: The role of the uropod in
- 920 neutrophil polarization and migration. *Dev Cell*, 38(2):161–169, 2016.
- 921 57. Caleb K Chan, Amalia Hadjitheodorou, Tony Y-C Tsai, and Julie A Theriot. Quantitative compar-
- 922 ison of principal component analysis and unsupervised deep learning using variational autoencoders
- 923 for shape analysis of motile cells. *bioRxiv*, 2020. doi: 10.1101/2020.06.26.174474.
- 924 58. Yibing Wei, Jiyoun Kim, Harri Ernits, and Daniel Remick. The septic neutrophil—friend or foe.
- 925 *Shock*, 55(2):147–155, 2021.
- 926 59. M Palomino-Segura, J Sicilia, I Ballesteros, and A Hidalgo. Strategies of neutrophil diversification.
- 927 *Nat Immunol*, 24(4):575–584, 2023.
- 928 60. Hadley Witt, Zicheng Yan, David Henann, Christian Franck, and Jonathan Reichner. Mechanosen-
- 929 sitive traction force generation is regulated by the neutrophil activation state. *Sci Rep*, 13(1):11098,
- 930 2023.
- 931 61. Dhananjay Bhaskar, Darrick Lee, Knútsdóttir, Cindy Tan, MoHan Zhang, Pamela Dean, Calvin
- 932 Roskelley, and Leah Edelstein-Keshet. A methodology for morphological feature extraction and

- unsupervised cell classification. *bioRxiv*, 2019. doi: <https://doi.org/10.1101/623793>. URL <https://www.biorxiv.org/content/10.1101/623793v1.full.pdf>.
62. Sicong Lu, Huaping Liu, and Chunwen Li. Manifold regularized stacked autoencoder for feature learning. In *2015 IEEE International Conference on Systems, Man, and Cybernetics*, pages 2950–2955, 2015. doi: 10.1109/SMC.2015.513.
63. Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. UMAP: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018. ISSN 2475-9066. doi: 10.21105/joss.00861. URL <http://joss.theoj.org/papers/10.21105/joss.00861>.
64. Kevin R. Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B. Burkhardt, William S. Chen, Kristina Yim, Antonia van Elzen, Matthew J. Hirn, Ronald R. Coifman, and et al. Visualizing structure and transitions in high-dimensional biological data. *Nature Biotechnology*, 37(12):1482–1492, Dec 2019. doi: 10.1038/s41587-019-0336-3.
65. Tingrin Gao. *Hypoelliptic diffusion maps and their applications in automated geometric morphometrics*. Duke University, 2015.
66. Tingran Gao. The diffusion geometry of fibre bundles: Horizontal diffusion maps. *Applied and Computational Harmonic Analysis*, 50:147–215, 2021. doi: <https://doi.org/10.1016/j.acha.2019.08.001>.
67. Luca Pozzi, Jason A Hodgson, Andrew S Burrell, Kirstin N Sterner, Ryan L Raaum, and Todd R Disotell. Primate phylogenetic relationships and divergence dates inferred from complete mitochondrial genomes. *Molecular phylogenetics and evolution*, 75:165–183, 2014.
68. Debbie Guatelli-Steinberg. Primate dentition: An introduction to the teeth of non-human primates. *American Journal of Physical Anthropology*, 121:189–189, 2003. doi: <https://doi.org/10.1002/ajpa.10194>.
69. Elizabeth M. St. Clair and Doug M. Boyer. Lower molar shape and size in prosimian and platyrrhine primates. *American Journal of Physical Anthropology*, 161:237–258, 2016. doi: <https://doi.org/10.1002/ajpa.23021>.



- 959 70. Robin Huw Crompton, Russell Savage, and Iain R Spears. The mechanics of food reduction in  
960 tarsius bancanus. *Folia Primatologica*, 69(7):41–59, 1998.
- 961 71. Reema Al-Aifari, Ingrid Daubechies, and Yaron Lipman. Continuous procrustes distance between  
962 two surfaces. *Communications on Pure and Applied Mathematics*, 66(6):934–964, 2013. ISSN  
963 00103640. doi: 10.1002/cpa.21444. URL [https://onlinelibrary.wiley.com/doi/10.1002/cpa.](https://onlinelibrary.wiley.com/doi/10.1002/cpa.21444)  
964 21444.
- 965 72. J. C. Gower. Generalized procrustes analysis. *Psychometrika*, 40(1):33–51, 1975. ISSN 0033-3123,  
966 1860-0980. doi: 10.1007/BF02291478. URL <http://link.springer.com/10.1007/BF02291478>.
- 967 73. Patrick W Oakes, Dipan C Patel, Nicole A Morin, Daniel P Zitterbart, Ben Fabry, Jonathan S Re-  
968 ichner, and Jay X Tang. Neutrophil morphology and migration are affected by substrate elasticity.  
969 *Blood*, 114(7):1387–1395, 2009.
- 970 74. Jennet Toyjanova, Estefany Flores-Cortez, Jonathan S Reichner, and Christian Franck. Matrix  
971 confinement plays a pivotal role in regulating neutrophil-generated tractions, speed, and integrin  
972 utilization. *J Biol Chem*, 290(6):3752–3763, 2015.
- 973 75. G Crainiciuc, M Palomino-Segura, M Molina-Moreno, J Sicilia, DG Aragonés, JLY Li, R Madurga,  
974 JM Adrover, A Aroca-Crevillén, S Martin-Salamanca, AS Del Valle, SD Castillo, HCE Welch,  
975 O Soehnlein, M Graupera, F Sánchez-Cabo, A Zarbock, TE Smithgall, M Di Pilato, TR Mempel,  
976 PL Tharaux, SF González, A Ayuso-Sacido, LG Ng, GF Calvo, I González-Díaz, F Díaz-de María,  
977 and A Hidalgo. Behavioural immune landscapes of inflammation. *Nature*, 601(7893):415–421, 2022.
- 978 76. Susan E Leggett, Jea Yun Sim, Jonathan E Rubins, Zachary J Neronha, Evelyn Kendall Williams,  
979 and Ian Y Wong. Morphological single cell profiling of the epithelial–mesenchymal transition.  
980 *Integrative Biology*, 8(11):1133–1144, 2016.
- 981 77. Amanda S Khoo, Thomas M Valentin, Susan E Leggett, Dhananjay Bhaskar, Elisa M Bye, Shoham  
982 Benmelech, Blanche C Ip, and Ian Y Wong. Breast cancer cells transition from mesenchymal to  
983 amoeboid migration in tunable three-dimensional silk–collagen hydrogels. *ACS Biomater Sci Eng*, 5  
984 (9):4341–4354, 2019.

- 985 78. Susan E Leggett, Mohak Patel, Thomas M Valentin, Lena Gamboa, Amanda S Khoo, Eve-  
986 lyn Kendall Williams, Christian Franck, and Ian Y Wong. Mechanophenotyping of 3d multicellular  
987 clusters using displacement arrays of rendered tractions. *Proc Natl Acad Sci USA*, 117(11):5655–  
988 5663, 2020.
- 989 79. Barholomeus H. M. Gerritsen. *Using weighted alpha complexes in subsurface modelling: Recon-*  
990 *structing the shape of observed natural objects*. IOS Press, 2001. ISBN 978-90-407-2247-9.
- 991 80. Herbert Edelsbrunner. Weighted alpha shapes. Technical report, University of Illinois at Urbana-  
992 Champaign, 1992.
- 993 81. Yohai Reani and Omer Bobrowski. A coupled alpha complex. *arXiv*, 2021. doi: 10.48550/ARXIV.  
994 2105.08113. URL <https://arxiv.org/abs/2105.08113>. Publisher: arXiv Version Number: 1.
- 995 82. Donguk Kim, Mokwon Lee, Youngsong Cho, and Deok-Soo Kim. Beta-complex vs. alpha-complex:  
996 Similarities and dissimilarities. *IEEE Transactions on Visualization and Computer Graphics*, pages  
997 1–1, 2019. ISSN 1077-2626, 1941-0506, 2160-9306. doi: 10.1109/TVCG.2018.2873633. URL <https://ieeexplore.ieee.org/document/8496780/>.
- 999 83. Herbert Edelsbrunner and J. Harer. *Computational topology: an introduction*. American Mathe-  
1000 matical Society, 2010. ISBN 978-0-8218-4925-5. OCLC: ocn427757156.
- 1001 84. Herbert Edelsbrunner and Ernst Peter Mücke. Simulation of simplicity: a technique to cope with  
1002 degenerate cases in geometric algorithms. *ACM Transactions on Graphics*, 9(1):66–104, 1990.  
1003 ISSN 0730-0301, 1557-7368. doi: 10.1145/77635.77639. URL [https://dl.acm.org/doi/10.1145/](https://dl.acm.org/doi/10.1145/77635.77639)  
1004 [77635.77639](https://dl.acm.org/doi/10.1145/77635.77639).
- 1005 85. Partha Niyogi, Stephen Smale, and Shmuel Weinberger. Finding the homology of submanifolds  
1006 with high confidence from random samples. *Discrete & Computational Geometry*, 39(1):419–441,  
1007 2008. ISSN 0179-5376, 1432-0444. doi: 10.1007/s00454-008-9053-2. URL [http://link.springer.](http://link.springer.com/10.1007/s00454-008-9053-2)  
1008 [com/10.1007/s00454-008-9053-2](http://link.springer.com/10.1007/s00454-008-9053-2).
- 1009 86. Clément Berenfeld, John Harvey, Marc Hoffmann, and Krishnan Shankar. Estimating the reach of  
1010 a manifold via its convexity defect function. *Discrete & Computational Geometry*, 67(2):403–438,

- 1011 June 2021. ISSN 1432-0444. doi: 10.1007/s00454-021-00290-8. URL <http://dx.doi.org/10.1007/s00454-021-00290-8>.
- 1012
- 1013 87. Jean-Daniel Boissonnat, André Lieutier, and Mathijs Wintraecken. The reach, metric distortion,
- 1014 geodesic convexity and the variation of tangent spaces. *Journal of Applied and Computational*
- 1015 *Topology*, 3(1-2):29–58, June 2019. ISSN 2367-1734. doi: 10.1007/s41468-019-00029-8. URL
- 1016 <http://dx.doi.org/10.1007/s41468-019-00029-8>.
- 1017 88. Jesus Puente. *Distances and algorithms to compare sets of shapes for automated biological*
- 1018 *morphometrics*. Princeton University, 2013. URL <http://arks.princeton.edu/ark:/88435/dsp01sq87bt73n>.
- 1019