# Shape-to-graph Mapping Method for Efficient Characterization and Classification of Complex Geometries in Biological Images

William Pilcher[1], Xingyu Yang[2], Anastasia Zhurikhina[1], Olga Chernaya[1], Yinghan Xu[1], Peng Qiu[1], Denis Tsygankov[1,*]

[1] Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology and Emory University School of Medicine, Atlanta, GA, 30332, USA

[2] School of Biological Sciences, Georgia Institute of Technology, Atlanta, GA, 30332, USA

* Corresponding author.

E-mail address: denis.tsygankov@bme.gatech.edu (D. Tsygankov)

Short Title: Morphometric analysis of cell formations

24 **Abstract**

25

26 With the ever-increasing quality and quantity of imaging data in biomedical research comes the

27 demand for computational methodologies that enable efficient and reliable automated extraction

28 of the quantitative information contained within these images. One of the challenges in providing

29 such methodology is the need for tailoring algorithms to the specifics of the data, limiting their

30 areas of application. Here we present a broadly applicable approach to quantification and

31 classification of complex shapes and patterns in biological or other multi-component formations.

32 This approach integrates the mapping of all shape boundaries within an image onto a global

33 information-rich graph and machine learning on the multidimensional measures of the graph. We

34 demonstrated the power of this method by (1) extracting subtle structural differences from visually

35 indistinguishable images in our phenotype rescue experiments using the endothelial tube

36 formations assay, (2) training the algorithm to identify biophysical parameters underlying the

37 formation of different multicellular networks in our simulation model of collective cell behavior,

38 and (3) analyzing the response of U2OS cell cultures to a broad array of small molecule

39 perturbations.

40

41

42    Author Summary

43

44    In this paper, we present a methodology that is based on mapping an arbitrary set of outlines onto

45    a complete, strictly defined structure, in which every point representing the shape becomes a

46    terminal point of a global graph. Because this mapping preserves the whole complexity of the

47    shape, it allows for extracting the full scope of geometric features of any scale. Importantly, an

48    extensive set of graph-based metrics in each image makes integration with machine learning

49    routines highly efficient even for a small data sets and provide an opportunity to backtrack the

50    subtle morphological features responsible for the automated distinction into image classes. The

51    resulting tool provides efficient, versatile, and robust quantification of complex shapes and patterns

52    in experimental images.

53

54    **Introduction**

55

56    Quantitative characterization of cell shapes and their organization within multicellular formations

57    is critically important for many biomedical applications, including tissue engineering (Gupta et al.

58    2009), phenotypic cell-based screening (Conrad et al. 2004, Viros et al. 2008), and testing

59    platforms for drug discovery (Murphy et al. 2010, Zanella et al. 2010). However, broadly

60    applicable and comprehensive morphometric analysis of complex geometries in imaging data

61    remains a challenging task. Here we present an approach that allows for an efficient and precise

62    extraction and classification of structural features in arbitrarily complex cellular patterns, including

63    subtle variations that are difficult to decipher using visual inspection or a set of standard geometric

64    measures.

65

66    Currently, a number of methods have been developed for the analysis of morphological changes

67    among individual cells (Carpenter et al. 2006, Selinummi et al. 2005, Tsygankov et al. 2014).

68    Some targeted approaches for extracting structural features in specific applications have been also

69    reported (Guidolin et al. 2004, Khoo et al. 2011, Lin et al. 2005, Nguyen et al. 1994), but there is

70    still a need for a *general* methodology allowing for automated comparative analysis of complex

71    multicellular formations. In particular, it is difficult to study the effects of a small perturbation in

72    the extracellular environment on the collective behavior of many cells and the patterns resulting

73    from their complex interactions (Chernaya et al. 2018). This problem is exacerbated when working

74    with experimental systems that allow for a precise control of different physical conditions

75    generating large and diverse sets of imaging data. To address this issue, we have developed a

76    general approach, which automatically generates a rich set of interpretable features from images

77   of cellular structures. These features are computed using a mathematically precise mapping of the

78   boundaries outlining all shapes in an image onto a global graphical structure. This graphical

79   structure captures multiple features relating to the width of the cellular objects, the shapes and

80   roughness of the boundaries, as well as the connectivity and density of the cell clusters across the

81   image. Using these features, we can identify images with similar structures, cluster images into

82   groups based on structural patterns, and use the image-level characteristics for regression tasks.

83   With this approach, one can cluster and visualize the differences between multicellular patterns

84   based on high-level features, while still retaining the ability to interpret and understand the features

85   defining each image type.

86

87   Unlike other graphical approaches which utilize morphological thinning (Boizeau et al. 2013,

88   Carpentier et al. 2012, Guidolin et al. 2004) or rely on a heavily pruned skeleton (Grélard et al.

89   2017, Ogniewicz and Kübler 1995, Rohde et al. 2008, Styner et al. 2003, Wearne et al. 2005,

90   Xiong et al. 2010), ours exploits the exhaustive image-scale graph to capture both fine features on

91   the boundary of the structures and coarse features of the objects' shapes. Furthermore, this

92   approach is not limited to only work on networked structures. One can use this method to

93   characterize changes in patterns of isolated cells and cell clusters, dense cellular networks, or any

94   mixture of such formations.

95

96   As a testing system for our methodology, we first used an endothelial tube formation assay along

97   with a computational model that simulates the formation of cellular patterns under controlled

98   perturbations of the biomechanical properties of the cells. The tube formation assay is a useful *in-*

99   *vitro* tool to screen for treatments that affect early stages of vasculogenesis. Healthy vascular

5

100    endothelial cells cultured on Matrigel form dense cellular networks across the dish. Environmental

101    or genetic perturbations can alter the resulting structure, leading to more irregular networks or

102    completely isolated cell clusters. The standard approach to quantify these assays is to count the

103    number of tubules (connections between cell clusters) or measure the percent coverage of a cellular

104    network within a certain field of view (Arnaoutova and Kleinman 2010). While these approaches

105    can be used to screen for treatments that are strongly pro- or anti-angiogenic, they are not precise

106    enough to distinguish between more similar patterns.

107

108    For experimental perturbation of collective cell behavior, we used knockdowns of the three CCM

109    proteins, with and without treatment by a Rho-associated protein kinase (ROCK) inhibitor, H1152

110    (Chernaya et al. 2018). These knockdowns all negatively affect tube formation and lead to either

111    small isolated cell clusters or sparse patterns with large tubules depending on the targeted protein.

112    Inhibition of ROCK partially rescues tube formation, increasing both tubule count and coverage,

113    although the resulting cellular networks appear much more disorganized compared to wild-type.

114    Here we show that features from the shape-to-graph mapping can differentiate images from these

115    experiments, including the cases when images do not seem to be distinguishable and explain the

116    differences between these visually similar groups using the features extracted from the mapping.

117

118    In addition to in-vitro assays, we utilized a simulated model that allowed us to generate a range of

119    different multicellular patterns depending on two biomechanical characteristics: the stability of

120    cell-cell contacts and the strength of cell-matrix adhesion (Chernaya et al. 2018). Altering these

121    properties can create structures ranging from completely isolated cellular clusters to interconnected

122    networks, all with varying densities. We apply our approach to predict the model parameters used

6

123    to generate each *in-silico* image, demonstrating that these features can capture the trends in the

124    way cellular structures progressively change due to the controlled modulation of the biomechanical

125    properties of the system.

126

127    Finally, to show that our methodology is not limited to mesh-like cell formations characteristic to

128    specific cell types, we applied it to completely different type of data from a large imaging set

129    publicly available at the Broad Bioimage Benchmark Collection [BBBC022v1] (Gustafsdottir et

130    al. 2013, Ljosa et al. 2012). Specifically, we analyzed confluent cultures of U2OS cells subjected

131    to an extensive set of small molecule treatments. The global (image-scale) nature of our graph

132    structure, which captures both the shapes of all individual cells and their relative spatial positioning

133    in the field of view, allowed us to outperform the conventional shape metrics in terms of precision

134    and sensitivity of the phenotypic classification.

135

136    Collectively, the performed data analysis illustrated the power of our approach for both single cell

137    and multicellular pattern characterization, capturing apparent and subtle geometric variations using

138    a small set of images or a large high throughput scans, while providing a way to backtrack and

139    interpret geometric features responsible for the classification outcome.
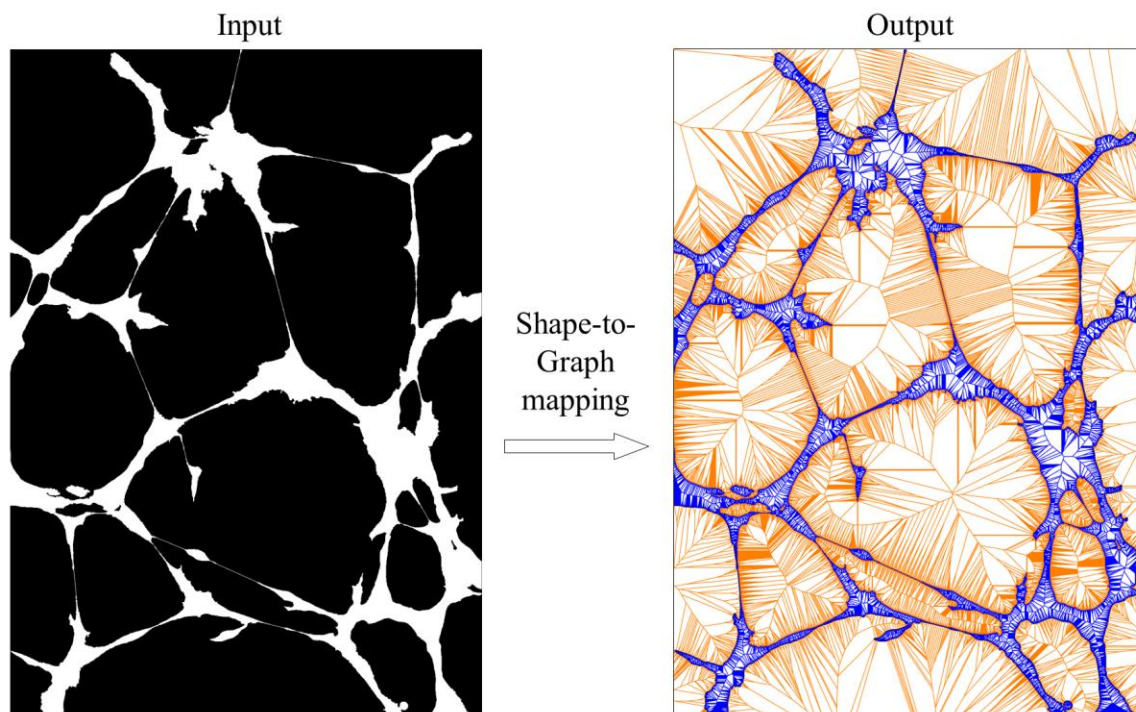
140

141 **Results**

142

143 *Shape-to-graph Mapping*

144

145 Our shape-to-graph mapping is a generalization of the Voronoi Diagram to accept the edges

146 outlining a shape as inputs. The traditional Voronoi Diagrams only operate discrete sets of points

147 such as in the default MATLAB algorithm (Aurenhammer 1991). Our algorithm is based on a

148 sweep-circle method (Xin et al. 2013) modified to work with line inputs. The algorithm has

149 $O(n \log n)$ complexity, where $n$ is the number of inputs, which scales linearly with image

150 resolution provided the same image content. Thus, the first step in the processing pipeline is to

151 take any binary images as an input, and output a graphical structure, which maps all piecewise

152 linear boundaries in the image to a unique image-scale graph spanning both the foreground and

153 background of the image **(Fig. 1)**.



154

8

155    *Figure 1.* *An illustration of the shape-to-graph mapping. Algorithm input is a binary image with the*

156    *foreground (value 1) shown in white and the background (value 0) shown in black. Algorithm output is an*

157    *image-scale graph structure. The part of the graph in the foreground (defined later in the text as in-graph)*

158    *is shown in blue, while the part in the background (out-graph) is shown in orange.*
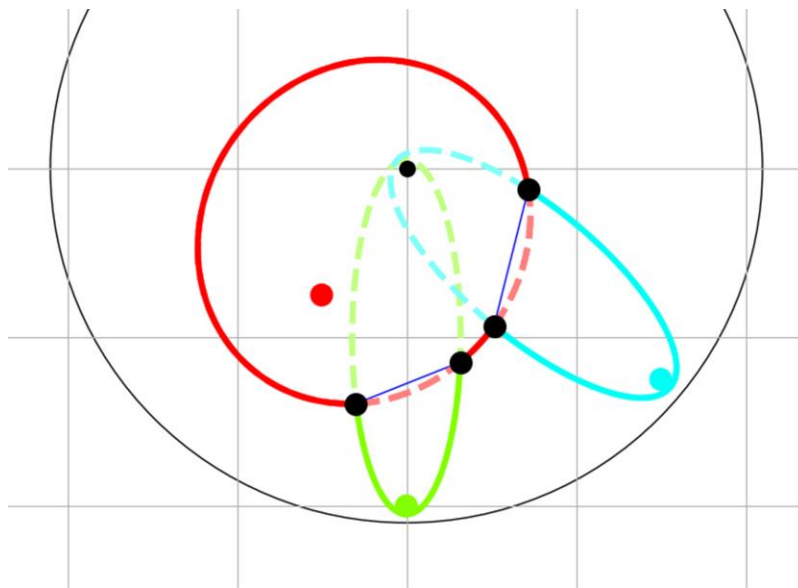
159

160    ***Graph construction process***

161

162    A Voronoi diagram consists of vertices, which are the centers of the largest circles that can be

163    packed within a given set of inputs, such that no input element lies within the circles. Thus each

164    graph vertex is the center of a circle tangent to three or more input elements, while the graph edges

165    are bisectors between two inputs. Our graph satisfies these definitions but presents a generalized

166    version of the Voronoi Diagram, which is derived from inputs that can include both a set of points

167    and a set of line segments. However, our main interest is an input of pixel-scale line segments

168    forming the boundaries in a binary image.

169

170    This graph can be constructed by searching through all circles tangent to any combination of three

171    inputs, and removing circles which contain an input within it. However, this approach would have

172    $O(n^3)$ complexity, where $n$ is the input size. Instead, we use a sweep-circle method, in which we

173    compute the Voronoi diagram within an expanding circle centered at the origin. Each input

174    generates a bisector with the sweep circle (Xin et al. 2013). Such bisector can be an ellipse for a

175    point input or a parabola for a line input. When a new input enters the sweep circle, it's bisector

176    will intercept with another bisector within the sweep circle. The set of all bisector segments that

177    are not contained within another bisector is referred to as the *beachfront* (**Fig. 2**). The interceptions

178    between two arcs of the beachfront always lie on bisectors between the inputs, which trace  out
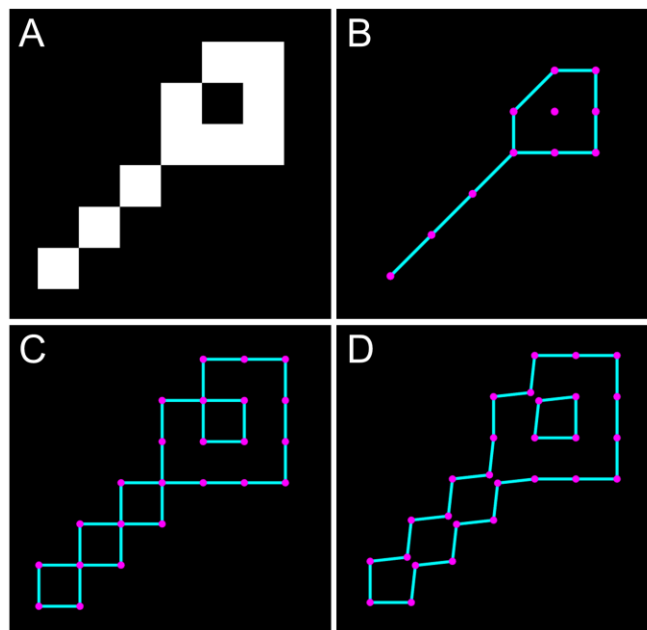
9

179    edges in the Voronoi diagram. To find the graph vertices, we only need to test inputs which have

180    adjacent arcs on the beachfront. The ordering of arcs on the beachfront are stored within a red-

181    black balanced binary tree (Xin et al. 2013), therefore the position of a new point within the

182    beachfront can be found with a binary search. Thus, the complexity with this approach scales as

183    $O(n \log n)$ with the number of inputs. For a more detailed, formal description of the sweep-circle

184    algorithm, see (Xin et al. 2013). Constricted this way, each Voronoi vertex has three Voronoi edges

185    Even in cases when the Voronoi vertex is equidistant to four or more inputs, such as the center of

186    a regular polygon, multiple Voronoi vertices are created at the same position, each with a degree

187    of three and a zero-length edge connecting them.



188

189    **Figure 2.** *Sweep-circle Voronoi algorithm for the graph construction. In this algorithm, a sweep circle*

190    *(grey circle) expands from the center of the image (grey dot). Each input point (colored dots) forms a*

191    *bisector (colored ellipses) with the expanding sweep circle. The beachfront is a set of all outer most portions*

192    *(solid elliptical arcs) of these bisectors. The intersections between the ellipses (black dots) trace out*

193    *Voronoi edges (blue lines). When two intersection points merge, pinching out a beachfront arc, a new*

194    *Voronoi vertex is formed.*

195

196    To extract the algorithm input from a binary image, we trace the boundaries along the half-pixel

197    border separating the background and foreground pixels. This is different from the conventional

198    tracing of boundaries along the pixel centers but ensures that a horizontal or vertical line of pixels

199    will have the width of one, rather than zero, which allows us to include pixel-size features to the
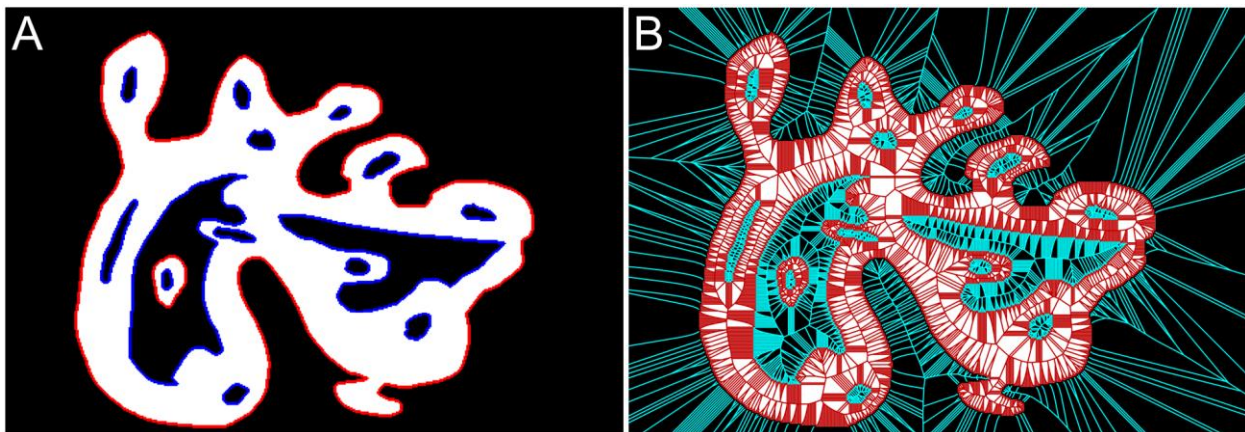
200    image analysis. (**Fig. 3A-C**)



201

*Figure 3. Boundary tracing.* ***A.*** *A simplified example of an input image.* ***B.*** *The conventional tracing of the boundary (implemented in MATLAB) along the centers of the pixels at the edge of a foreground object.* ***C.*** *Our algorithm traces the boundary directly along the lines separating the foreground and background pixels.* ***D.*** *An illustration of how the algorithm eliminates all boundary self-crossings by a small non-disruptive off-diagonal shift (here the shift was exaggerated for the illustration purposes).*

207

208    If two boundary points overlap, such as when two foreground pixels are connected diagonally,

209    these points are separated in the off-diagonal direction by a very small distance (we used 1/20 of

210    the pixel size) to ensure that boundaries in the image never intersect or self-cross but

211    unambiguously enclose the corresponding objects and holes (**Fig. 3D**).

212     *Graph annotation*

213

214     All connected components in the foreground (objects) and background (holes) of the binary image

215     are identified and assigned a unique numerical label. Boundaries are additionally categorized into

216     two types: *exterior boundaries* that completely enclose a foreground object and *interior*

217     *boundaries* that enclose a hole and, in turn, are enclosed by an object (**Fig. 4A**). Once the complete

218     graph is contracted, we will refer to the part of the graph situated in the image foreground as *in-*

219     *graph* and the part in the image background as *out-graph* (**Fig. 4B**).

220



221     *Figure 4. **A.** Boundary annotation: exterior boundaries are shown in red, while interior boundaries are*

222     *should blue. **B.** Overall graph annotation: in-graph is shown in red, while the out-graph is shown in cyan.*

223

224     Graph vertices that are equidistant to exactly two different boundaries form a sequence of vertices

225     that we call *bridges*. Different bridges come together at graph vertices that are equidistant to three

226     or more different boundaries and identified here as *hubs* (**Fig. 5A**). Additionally, a sequence of

227     vertices that connect two looped bridges associated with the same boundary, which may occur

228     when there is a hole within an extended protrusion of an object, is referred here as a *connector*.

229     Identifying all bridges, hubs, and connectors allows us to partition the whole graph into non-

230    overlapping *subgraphs* uniquely associated with each interior or exterior boundary (**Fig. 5B**).

231    Extracting features from their subgraphs is central to our methodology.
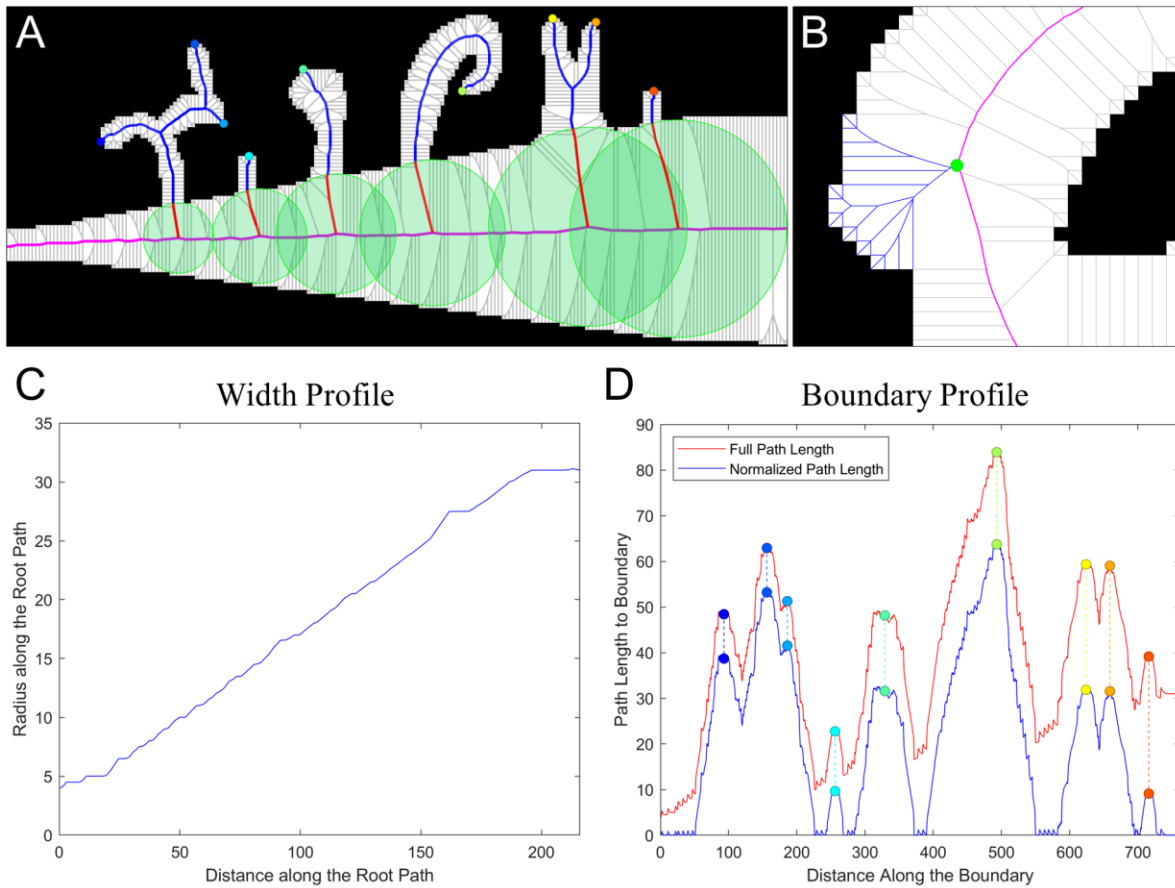


232

233    ***Figure 5.*** *The key elements of the graph. **A.** All bridges (red), hubs (green), and connectors (blue) of the*

234    *in-graph. **B.** Partitioning of the in-graph into subgraphs (shown with unique colors). Each non-overlapping*

235    *subgraph is associated with exactly one interior or exterior boundary.*

236    ***Graph-based Feature Extraction***

237

238    Each vertex in the constructed graph represents the center of a circle inscribed within the object.

239    A subgraph with no bridges, such as the graph within a single-boundary object with no holes or a

240    hole with no objects inside, is a single tree with the root node being the center of the largest

241    inscribed circle. Otherwise, a subset of vertices located on the graph bridges and connectors of the

242    associated subgraph acts as a set of the roots, from which graph edges branch out towards the

243    corresponding boundary (**Fig. 6A,B**). Constructed this way, each subgraph is outlined by the

244    boundary on one side and by a continuous sequence of bridges and connectors on the other side.

245    We will call this sequence of bridges and connectors the *root path*. Again, in case of objects with

246    no holes, there are no bridges, and the root path is defined as the longest path to the boundary

13

247     which passes through the single root node. Based on this construction, we derive two primary

248     metrics for each subgraph, which we call the *width profile* and the *boundary profile.*



249

**Figure 6.** *The primary graph metrics.* ***A.*** *An example of paths along the graph edges from the root path*

*(magenta) to the tips of object protrusions. The inscribed circles (green) provide a measure for the width*

*profile. The parts of the paths (blue) outside the circles provide a measure for the normalized boundary*

*profile.* ***B.*** *An illustration of path (blue) branching from a root (green) to the boundary, so that each*

*boundary point has an associated root node and a shortest path to this node along the graph edges.* ***C.*** *The*

*resulting width profile showing the inscribed circle radii for every node on the root path.* ***D.*** *The resulting*

*boundary profile before subtracting the radii of the corresponding root nodes (red) and after subtracting*

*(blue). The colored points at the local maxima of the boundary profile correspond to the protrusion tips in*

***A.***

14

259    The *width profile* describes coarse variations in the subgraph's width defined as the radii of the

260    inscribed circles with the centers located at the vertices of the root path (**Fig. 6C**). When computed

261    in background regions, this captures local variations in density. The *boundary profile* captures the

262    size of any protrusion or bump which lies along the boundary. The boundary profile is computed

263    by measuring the shortest distance along the subgraph edges from all points along the boundary to

264    the corresponding root nodes. By using distances along the subgraph edges, we accurately

265    characterize the size of these features even if the boundary is highly curved. To ensure that the

266    boundary profile is not sensitive to the same variations in object size as the width profile, the

267    boundary profile is normalized at each point by subtracting the radius on the inscribed circle with

268    the center at the root node where the path to that boundary point begins (**Fig. 6D**).

269

270    Because each boundary has a corresponding subgraph in both the in-graph and out-graph parts of

271    the full graph, each boundary has a foreground and background width profile along with a

272    foreground and background boundary profile. The only exception would be the most outward

273    boundaries, for which out-graphs extend to infinity. To resolve this issue, we constrain the graph

274    within the image by using the image boundary as the most outward boundary.

275

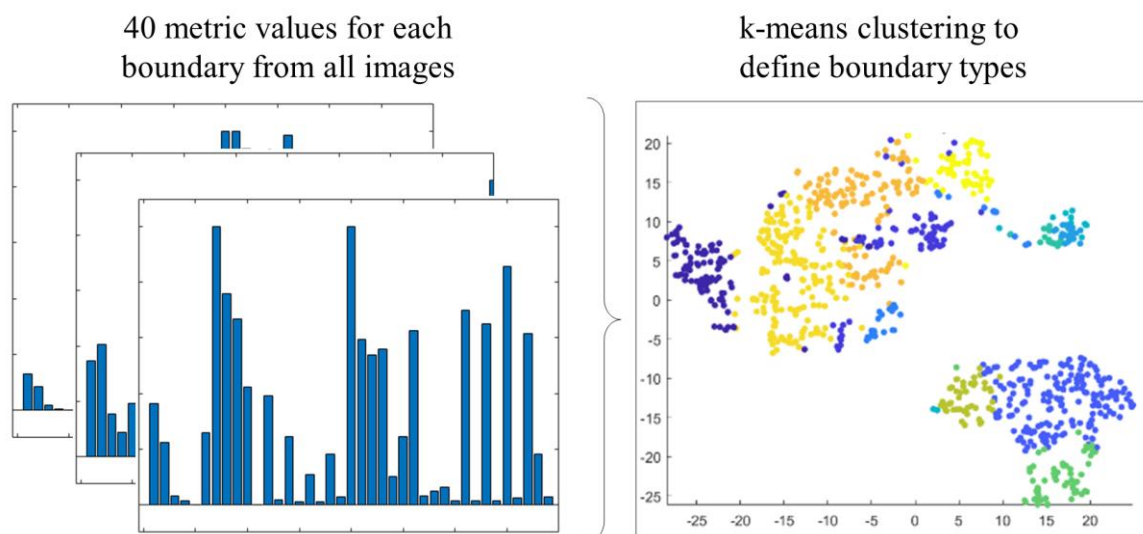276    *Per-Image Structural Features*

277

278    In order to characterize or compare complex geometric structures such as multicellular patterns,

279    per-boundary classification would be insufficient as we must consider the features of all

280    boundaries to account for the overall structure of a pattern in an image. Thus we construct a set of

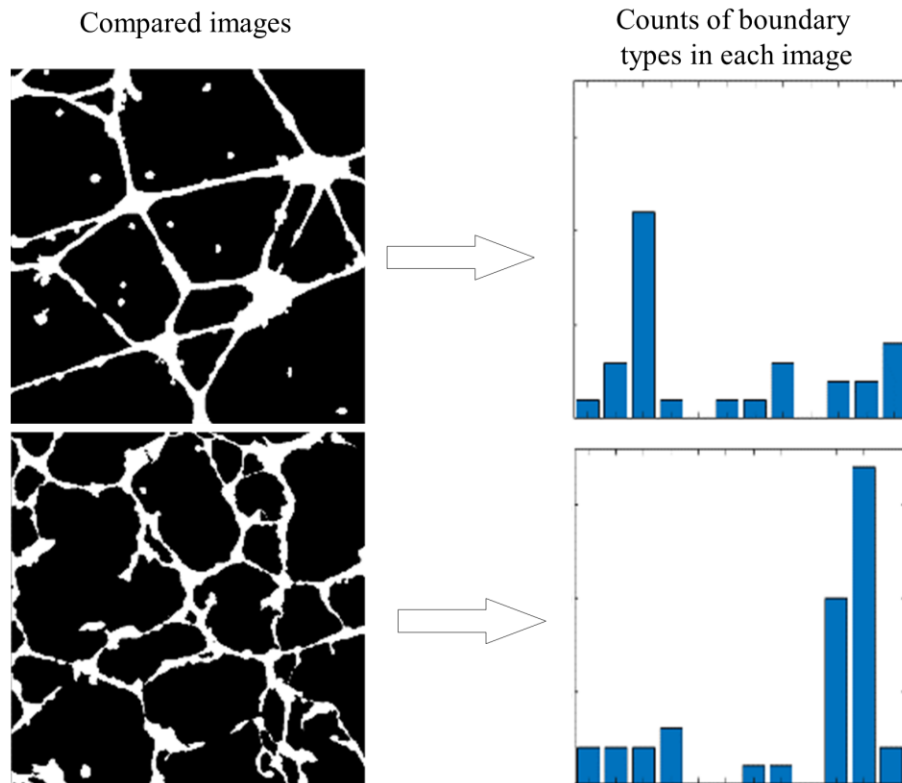281    per-image features derived from our graph-based per-boundary features.

15

282

283    To this end, we start with associating each boundary with 40 features, including distribution

284    metrics for the width profile and boundary profile, along with the area and perimeter of each

285    boundary. Half of the features computed for each boundary come from the corresponding in-graph

286    and half from the out-graph. The full list of features is provided in the **Supplemental Table 1**.

287    Next, we perform k-means clustering on the list of all boundaries across all provided images (**Fig.**

288    **7**). This process creates a histogram of $N$ boundary types within each image. The goal of this

289    clustering is to automatically differentiate boundaries based on a combination of their roughness,

290    the size and shape of the enclosed objects and holes, and the relative separation of these objects

291    and holes. This means that holes or objects with the same shape may lie in different clusters if the

292    cellular structure around the hole is thicker or thinner, or if the object lies in a more or less dense

293    region. The count or frequency of the boundary types in each image then serves as a per-image

294    feature (**Fig. 8**). The specific interpretation of each boundary type depends on the nature of data

295    presented in the images under investigation, but this is what ultimately allows us to understand

296    differences in the structural organization of the patterns in imaging data sets, as we show in the

297    next section.



298

299    ***Figure 7.*** *Boundary type identification. We use 40 metrics extracted for each boundary from all the images*

300    *in a given set and use k-means to associate each boundary with one of the N classes.*



301

302    ***Figure 8.*** *Per-image characterization. For each image, we extract the counts of boundaries that belong to*

303    *each of N boundary types, which were determined using k-means clustering on the 40 boundary features.*

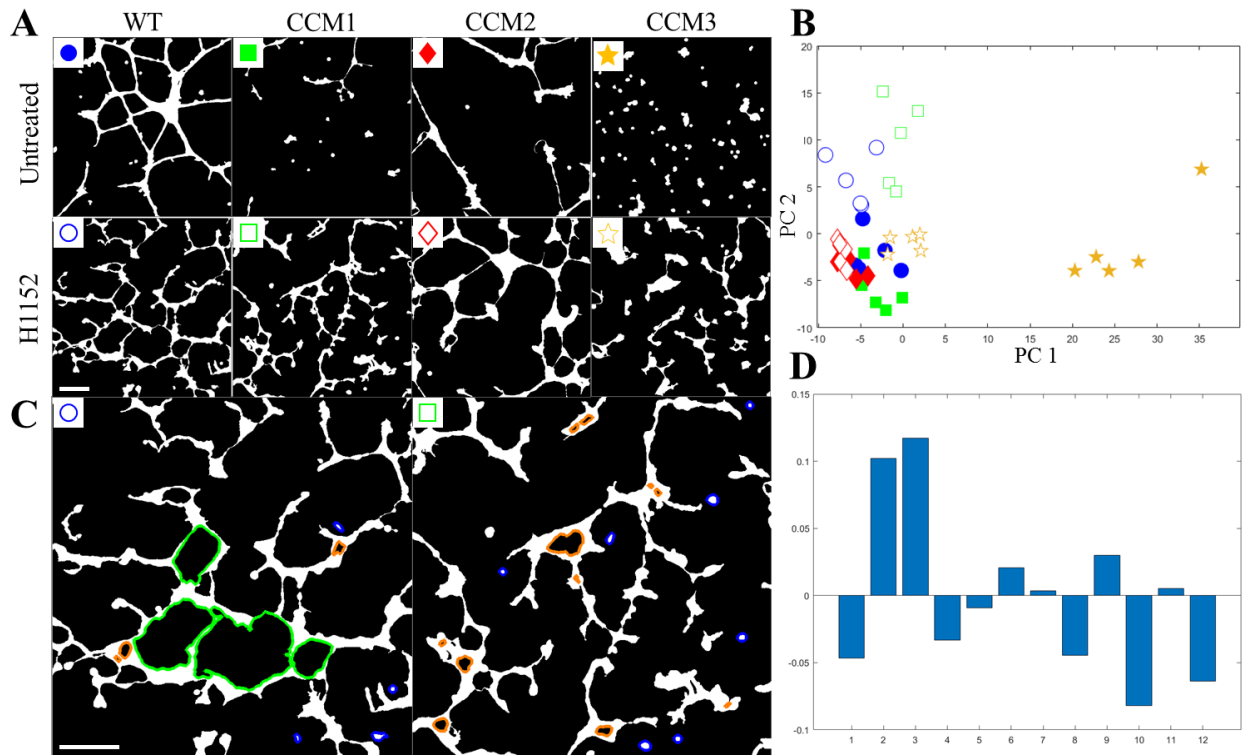304    ***Analysis of In-Vitro Tube Formations***

305

306    In this section we test the ability of our method to identify subtle structural difference in a small

307    set of images from an in vitro endothelial tube formation assay (the experimental data has been

308    previously published in (Chernaya et al. 2018)). The set includes images of the control cell (wild-

309    type HUVEC) and cells with knockdown (KD) of the three Cerebral Cavernous Malformation

310    (CCM) proteins, CCM1 (or KRIT1), CCM2, and CCM3 (or PDCD10), which disrupts the integrity

311    of multicellular mesh. In addition, the control and KD cells were treated with an inhibitor of Rho-

17

312    associated protein kinase (ROCK), which was shown to be over-activated in CCM KD cultures

313    (Chernaya et al. 2018). The treatment with the ROCK inhibitor H1152 partially rescues the wild-

314    type (WT) phenotype, although the resulting cellular patterns in the tube formation assay do not

315    closely match the WT patterns. Previously, we showed that although the diseased and the H1152

316    treated phenotypes are clearly different from the *untreated* WT phenotype, some treated cultures

317    are indistinguishable from the *treated* WT cells both visually and based on the traditional

318    geometric measures(Chernaya et al. 2018). Here we show that our shape-to-graph approach allows

319    us to identify the distinguishing features in all the phenotypes, including the ones with subtle

320    disparities that are not apparent upon visual inspection. The latter are of the main interest from the

321    methodology testing perspective.

322

323    For each of eight phenotypes (WT, CCM1, CCM2, CCM3, WT$^{H1152}$, CCM1$^{H1152}$, CCM2$^{H1152}$,

324    CCM3$^{H1152}$), we used five representative fields of view **(Fig. 9A)**. The boundaries were clustered

325    into 12 boundary types using k-means clustering. The optimal number of boundary types was

326    selected by performing 3-nearest neighbor classification on each image, where the class of each

327    image was determined by the class corresponding to the three most similar boundary type

328    histograms in the image set. Twelve clusters had a 90% classification accuracy (**Supplemental**

329    **Fig. S1**).

18

*Figure 9. Comparison of in-vitro tube formation assay structures with eight different phenotypes. A. The eight phenotypes resulted from WT and the knockdown of three CCM proteins, all with and without treatment by the ROCK inhibitor. Knockdown of the CCM proteins is associated with the disruption of the otherwise connected mesh. ROCK inhibitor leads to a more connected but still noticeably disorganized network. The scale bar is 200 μm. B. The first two principal components of each image's boundary type histogram. Images of a similar type and appearance tend to have similar histograms. Here, the markers indicate the corresponding images in A. C. Two images from $WT^{H1152}$ and $CCM1^{H1152}$ that appear visually similar but have significantly different boundary type counts. Boundaries that are responsible for the difference are highlighted in blue and cyan. The scale bar is 200 μm. D. The difference of the boundary type frequency histograms for $CCM1^{H1152}$ and $WT^{H1152}$. Boundary types 2 and 3 (Blue, orange) corresponding to small, isolated objects and small holes in wider locations in the network, appear significantly more often in $CCM1^{H1152}$ formations as compared to otherwise similar $WT^{H1152}$ structures. $WT^{H1152}$ structures tend to have more of boundary type 10 (Green), which are medium sized holes with more bumps and protrusions extending into the hole.*

19

345

346     Principal component analysis (PCA) was performed on the matrix of per-image boundary

347     histograms. Generally, images of the same class group together and exist in space near images

348     with similar structural features (**Fig. 9B**). Groups that are visually distinct, such as CCM3 cultures,

349     which have several small cellular clusters, appear far from H1152-treated cultures with fully

350     connected cell networks. Similarly, images with thicker structures, such as in CCM2$^{H1152}$ cultures,

351     appear further in principal component space from images with thinner structures, such as in

352     CCM1$^{H1152}$ and WT$^{H1152}$ cultures. Visually similar structures of CCM1$^{H1152}$ and WT$^{H1152}$ (**Fig. 9C**),

353     which both have many thin, disorganized connections, appear nearer to each other in principal

354     component space. Significantly different boundary types between sets of images can be identified

355     from the average boundary frequency histograms (**Fig. 9D).** This difference corresponds to an

356     increased frequency of three boundary types: type 2 consists of the *small isolated objects in regions*

357     *of high density* which appear more often in CCM1$^{H1152}$ cultures (blue boundaries in **Fig. 9C**); type

358     3 includes small holes in thick regions of the cellular structure, which also occur more frequently

359     in  CCM1$^{H1152}$ (orange boundaries in **Fig. 9C**); type 10 includes medium size holes, typically with

360     more bumps or protrusions from the cellular network extending into the hole, which occurs more

361     frequently in WT$^{H1152}$ samples (green boundaries in **Fig. 9C**). Descriptions of the boundary types

362     can be determined by analyzing the distribution of the original boundary metrics within each type

363     (**Supplemental Fig. S2**).

364

365

366

367
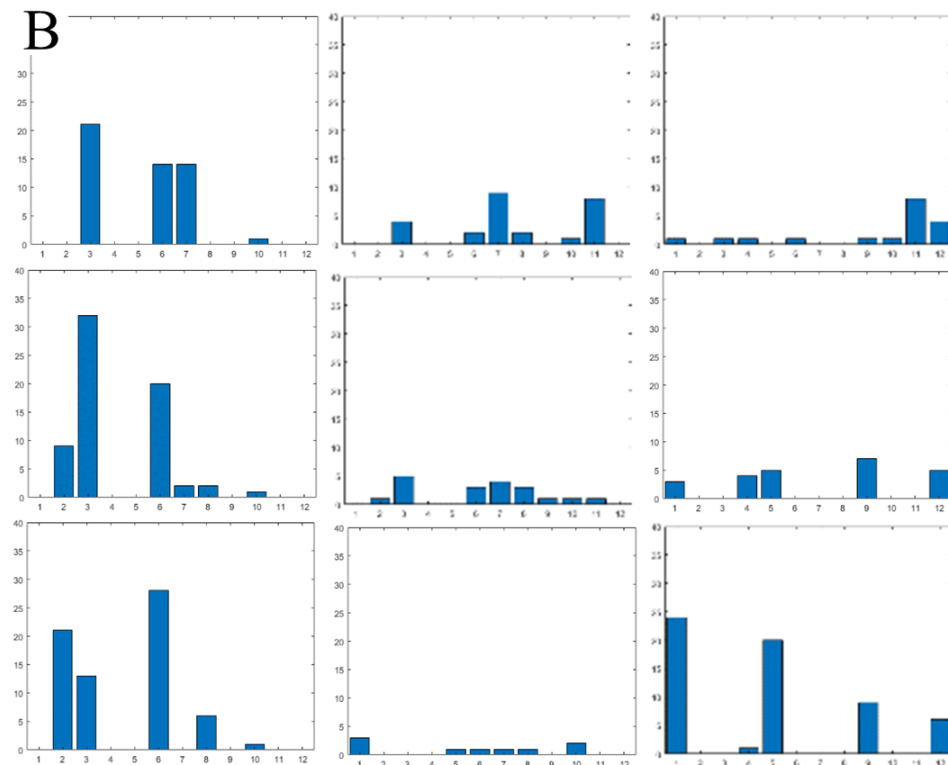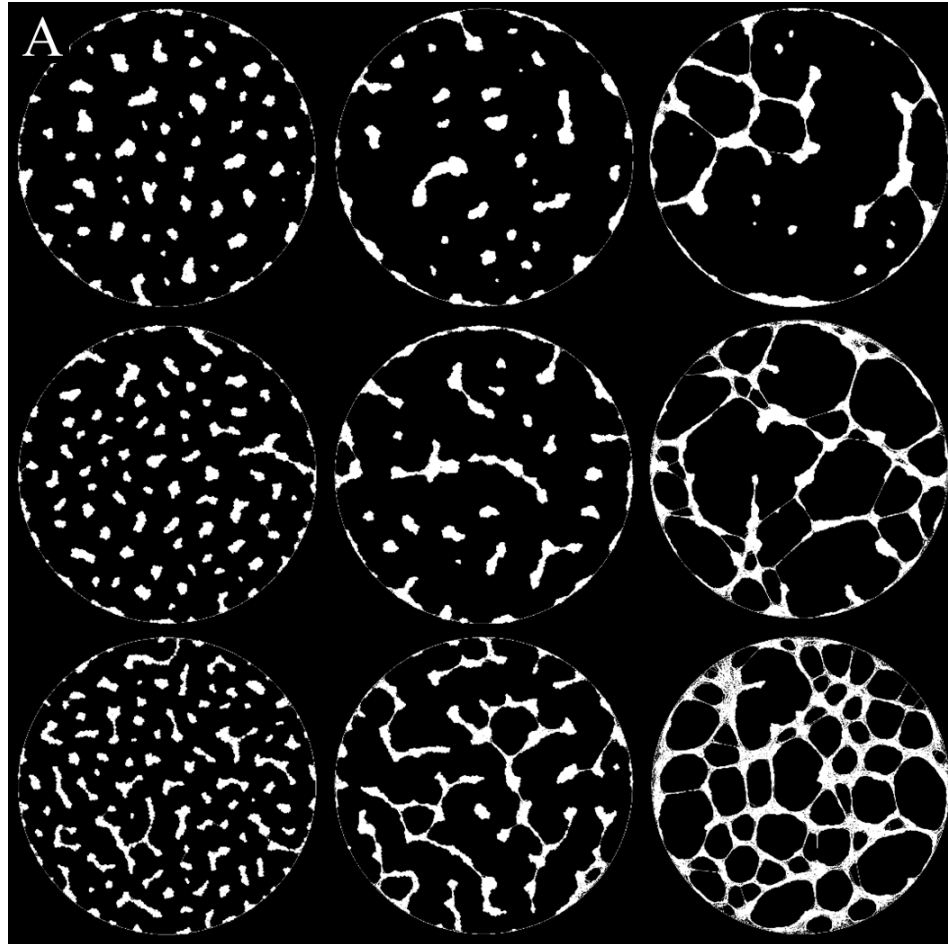
20

368    *Analysis of Simulated Data*

369

370    We used a previously developed computational model of endothelial tube formation (Chernaya et

371    al. 2018) to simulate 100 images of different cellular patterns corresponding to changes in two

372    biomechanical characteristics of cell interaction.

373

374    In this simulation model, each individual cell from a large group (hundreds to thousands) of cells

375    sparsely distributed over the substrate surface is represented as an extendable half-ellipsoid with

376    stochastically extending and retracting protrusions. Protrusions that extend downwards are

377    responsible for cell-substrate interactions, while protrusions that extend sideways along the surface

378    are responsible for cell-cell interactions. Cells form attachments when protrusions either reach

379    deep enough into the substrate, or when it reaches another cell. Retraction of the attached

380    protrusions leads to the cell movement, changes in cell shapes, and the buildup of the mechanical

381    stress that can lead to the contact breakage. Ultimately, because of these cell-cell and cell-substrate

382    interactions, the multicellular system evolves to form different patterns depending on the model

383    parameters at the cell level. Two key parameters of interest here are the stability of cell-cell and

384    cell-ECM adhesions. With properly selected values of the parameters, the model produces a dense

385    cellular network closely resembling wild-type endothelial cells in our in-vitro tube formation assay.

386    Reducing the values of each parameter leads to either a more sparse network or a number of

387    isolated cell clusters, similar to the behavior of cell with the knockdown of CCM1 and CCM3. It

388    is important to note here that even with a fixed set of parameters, the stochastic nature of protrusion

389    dynamics and a random initial distribution of cells make the structures resulted in simulations vary;

390    so that multiple patterns can be generated for the same phenotype similar to the experimental data.
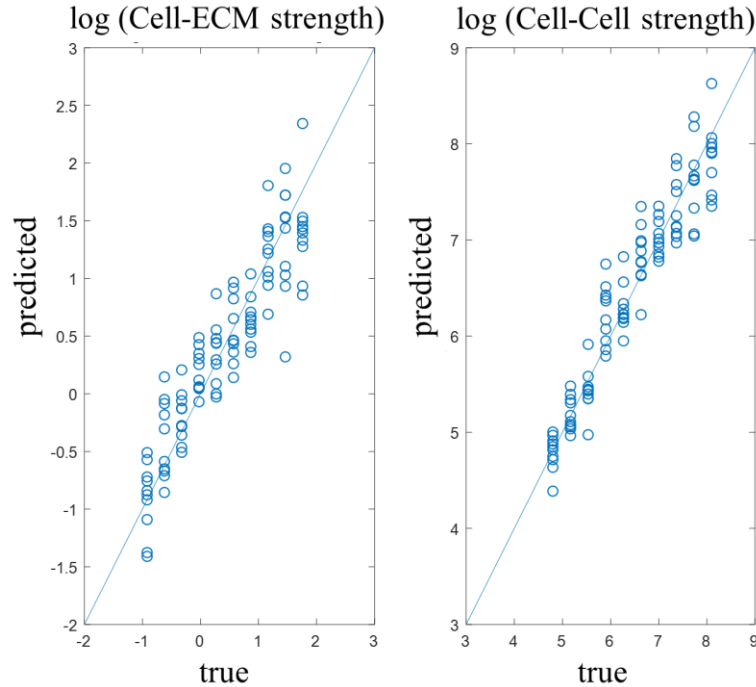
21

391    As we vary the two parameters representing the stability of cell contacts, our simulations allow us

392    to generate a sequence of cell formations with progressively changing structures (**Fig. 10A**).

393    Variation in the stability of cell-cell contact, the parameter $\kappa_{lat}$ in the probability of contact

394    breakage $P_{cell-cell} = 1 - exp(-l^2/\kappa_{lat}^2)$, where $l$ is the extension of the contact spring in the

395    model, has a strong impact on the boundary metrics. As this parameter is increased, cells go from

396    forming completely isolated cell clusters to a completely interconnected network. This leads to an

397    overall reduction in boundary types corresponding to isolated cell clusters, and a shift towards

398    networked structures with medium to large sized holes. The other parameter, $\kappa_{bott}$ in the

399    probability of cell-substrate contact breakage $P_{cell-ECM} = 1 - exp(-l^2/\kappa_{bott}^2)$, primarily affects

400    the velocity of cell movement and the resulting density of the cell clusters. The way this parameter

401    impacts the resulting structure depends on the network connectivity in the multicellular pattern,

402    but generally controls the density of the structure, with low values causing cells to form larger and

403    more sparse clusters.

22

404

405     ***Figure 10. A.*** *Nine representative images of multicellular formations out of 100 that were generated by*

406     *varying two parameters: the strength of cell-ECM adhesion (vertical axis) and the stability of cell-cell*

407     *contacts (horizontal axis).* ***B.*** *Variations of the two parameters result in visible changes in the boundary*

408     *type histograms.*

409

410     We applied our shape-to-graph mapping to the 100 generated images, extracted the boundary

411     features, and clustered boundaries to create a histogram of boundary types for each image. By

412     plotting the boundary type histograms, we can see the trends in the boundary type distribution

413     when the two parameters are varied (**Fig. 10B**) as described above. A multi-regression model was

414     used to predict the log-transformed values of the two model parameters based on the count of each

415     boundary type in each image (**Fig. 11**). If these parameters have a predictable impact on the

416     resulting multicellular pattern, and if the shape-to-graph mapping captures features that properly

417     reflect these changes, then this multi-regression model should be able to reproduce trends in the

418     two model parameters purely from the structural aspects of the cell patterns in the resulting images.

419     Indeed, our approach allowed us to predict the parameter values with high accuracy: log-

420     transformed cell-cell adhesion had a mean average error of 0.2392 with values ranging from 5 to

421     8 and a correlation coefficient of 0.9977, while log-transformed cell-ECM adhesion had a mean

422     average error of 0.2782 and a correlation coefficient of 0.86695. Twelve boundary clusters were

423     used based on cross validation performance.

**424**

*Figure 11.* *A linear regression model was trained to predict log-transformed model parameters from the boundary type histograms. The mean average error in predicting cell-cell adhesion was 0.2392, while predicting the strength of cell-ECM adhesion had the mean average error of 0.2782.*

**425**

**426**

**427**

**428**

**429** *Analysis of individual cells in a high throughput assay profiling small-molecules-induced cell*
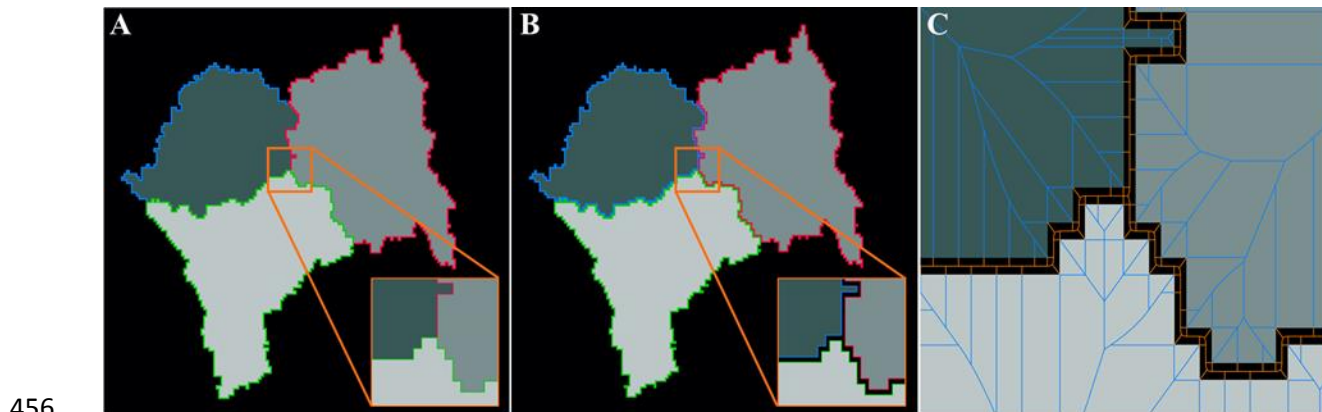
**430** *cultures.*

**431**

**432** In the previous sections we have focused on the analysis of complex multicellular formation with

**433** a mesh-like structures. However, our methodology is not limited to that particular type of data and

**434** can be adapted for the analysis of any images that can be segmented into the object(s) of interest

**435** and the background. To illustrate this statement, we applied our method to analyze *individually*

**436** *segmented* cells in a large publicly available image set with cell cultures subjected to phenotype

**437** perturbations by a variety of small molecules. We used image set BBBC022v1 , available from

**438** the Broad Bioimage Benchmark Collection (Ljosa et al. 2013). The original dataset consists of

25

439    fluorescent microscopy images of U2OS cells treated with one of over 1600 compounds. Five

440    fluorescent channels were captured for each field of view. The dyes used for visualization included

441    Hoechst 33342 (nuclei), concanavalin A (endoplasmic reticulum), SYTO 14 (nucleoli), phalloidin

442    (actin), and WGA (Golgi complex). A CellProfiler (Carpenter et al. 2006) pipeline provided with

443    the dataset was used to segment individual cells in each field of view via the watershed algorithm.

444    The samples were split into 20 plates with 384 wells each. Nine fields of view were obtained for

445    each well.

446

447    In the previous sections, our analysis relied on the input images for the shape-to-graph algorithm

448    in the form of binary masks, in which the extracted boundaries separated the cellular structure from

449    the background. However, in the imaging data we use here, each cell is treated as an individual

450    object, and therefore may share a boundary with either the background or other cells. This can

451    cause some cell boundaries to overlap **(Fig. 12A).** To ensure cell boundaries do not overlap, we

452    added a subpixel separation of the boundaries by shifting boundary points half-way from the

453    previously defined half-pixel boundaries towards the corresponding pixel center **(Fig. 12B)**. This

454    means one-pixel wide objects are thinned to have a width of half a pixel, and a half-pixel size gap

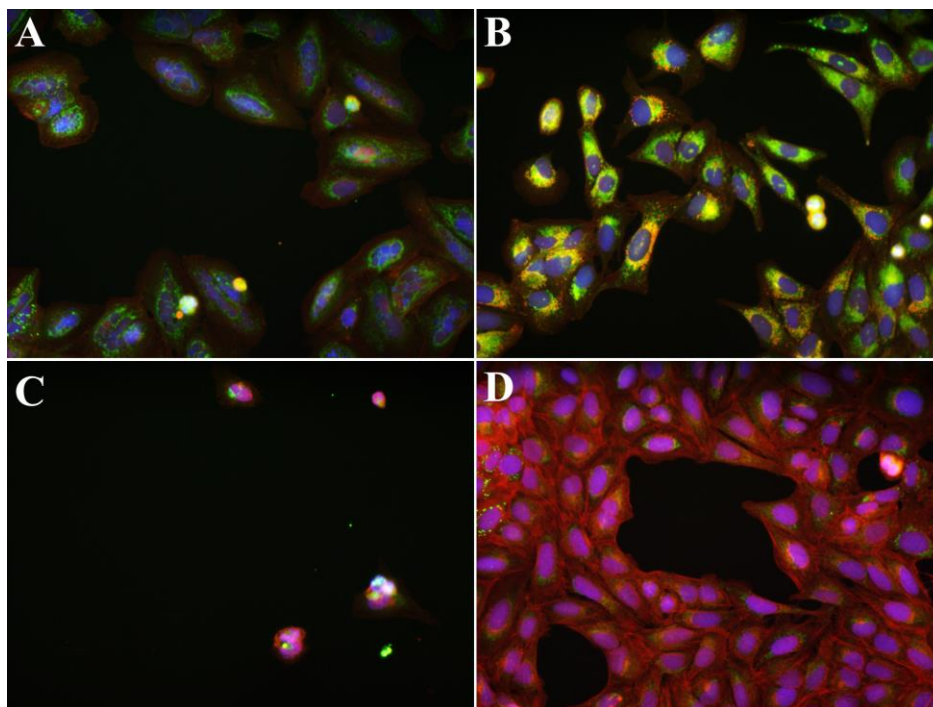455    is enforced to appear between two touching objects.

456

457     *Figure 12. A modified boundary tracing for individual cells in a tight cluster. **A.** With the previously*

458     *described boundary tracing, boundaries of contacting cells will overlap. **B.** The tracing routine is modified*

459     *to place boundary points halfway between the pixel center and our original half-pixel type tracing. This*

460     *creates a half-pixel gap between bordering cells. **C.** Parts of the out-graph for each cell (orange) lies within*

461     *this gap. Thus, the image out-graphs will include the out-graph nodes between all the contacting cells,*

462     *effectively encoding the spatial distribution of the cells in the image.*

463

464     With this processing approach, the cells are presented as individual objects embedded in an image-

465     scale mesh-like background **(Fig. 12C),** so that the graph representation of the background (out-

466     graph) encodes the information about the positional organization of all the cells and degree of

467     confluency of the whole cell culture.

468

469     For our analysis, we selected 11 compounds which the authors identified as forming strong clusters

470     based on their known mechanism of action and the 824 textural and morphological features they

471     extracted for each image. These compound clusters include tubulin modulators (fenbendazole,

472     oxibendazole, taxol) **(Fig. 13A)**, modulators of neuronal receptors (fluphenazine, metoclopramide,

473     procaine) **(Fig. 13B)**, and structurally related cardenolide glycosides (digoxin, lanatoside C,

474     peruvoside, neriifolin, digitoxin) **(Fig. 13C)**. We also included control samples from the same

475     assays **(Fig. 13D)**. We investigated if we could predict these mechanisms of action utilizing the

476     shape metrics derived from our shape-to-graph mapping. To this end, we extract the previously

477     described set of measures for each object in each image. The mean and standard deviation of these

478     per-cell metrics are computed across each well. To account for variance between plates, we

479     subtracted the feature vector of each well by the median feature vector of the control wells in the
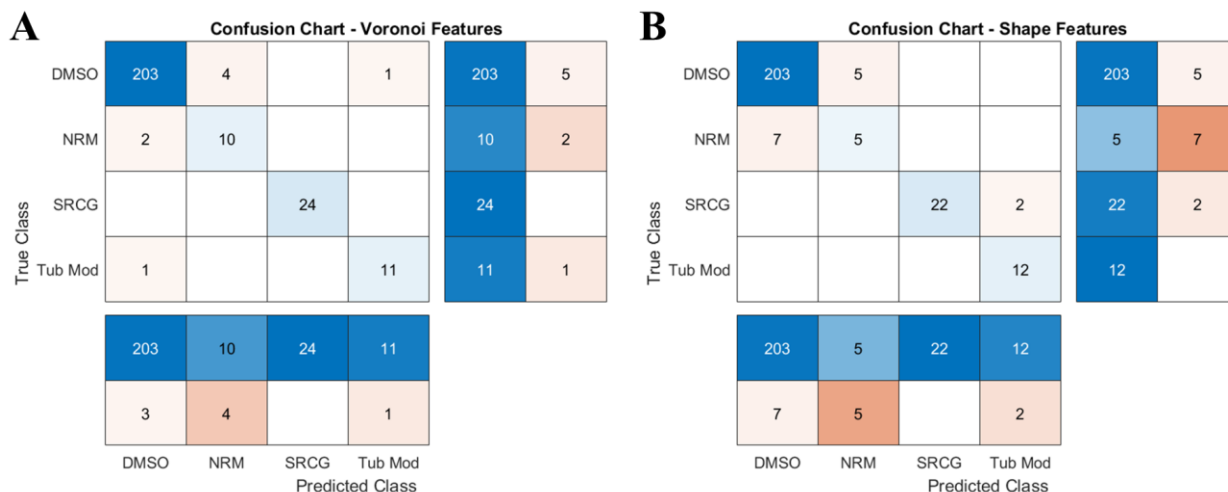
480    same plate. In the end, this resulted in 208 control wells, 12 samples of tubulin modulators, 12

481    samples of neuronal receptor modulators, and 24 samples of structurally related cardenolides.

482



483    **Figure 13.** *Images from the U2OS dataset. Red channel is phalloidin, blue is Hoechst 33342, and green is*

484    *WGA. A) Example image from the untreated group. B) Image of cells treated with taxol from the tubulin*

485    *modulators group. C) Image of cells treated with metoclopramide from the modulator of neuronal receptors*

486    *group. D) Image of cells treated with digoxin from the structurally related cardenolide glycosides group.*

487

488    Once the metrics were extracted, each plate was individually held-out, and a decision tree trained

489    on the wells in the remaining 19 plates were used to predict the held-out well labels. Shape-to-

490    graph features had a mean $F_1$ score of 0.916 (defined as $2 * \frac{precision * recall}{precision + recall}$, where $recall =$

491    $\frac{true\ positive}{true\ positive + false\ negative}$ , and $precision = \frac{true\ positive}{true\ positive + false\ positive}$ ), while the original

492    published shape features (Gustafsdottir et al. 2013) had an $F_1$ score of 0.826. Notably, the shape-

493    to-graph mapping had much better performance on the 'Modulators of Neuronal Receptors'

28

494    category, with a class $F_1$ score of 0.769 versus 0.455 for the original shape features (**Fig. 14)** and

495    each class appears to form tighter, more distinct clusters with the new features (**Supplemental Fig.**

496    **S3**). This treatment is the one which most strongly resembles the control dataset, but the cells tend

497    to be much less dense relative to the control wells. This reduced density is captured in the out-

498    graph radius metrics for each cell (**Supplemental Fig. S4**).



499

500    *Figure 14. Held-out plates were classified with a decision tree trained on the remainder of the dataset. The*

501    *new metrics derived with our approach tends to have better classification accuracies, especially for the*

502    *control class and the modulator of neuronal receptors (NRM). Mean $F_1$ score is 0.916 with the graph*

503    *derived metrics, and 0.826 with the CellProfiler shape metrics.*

504

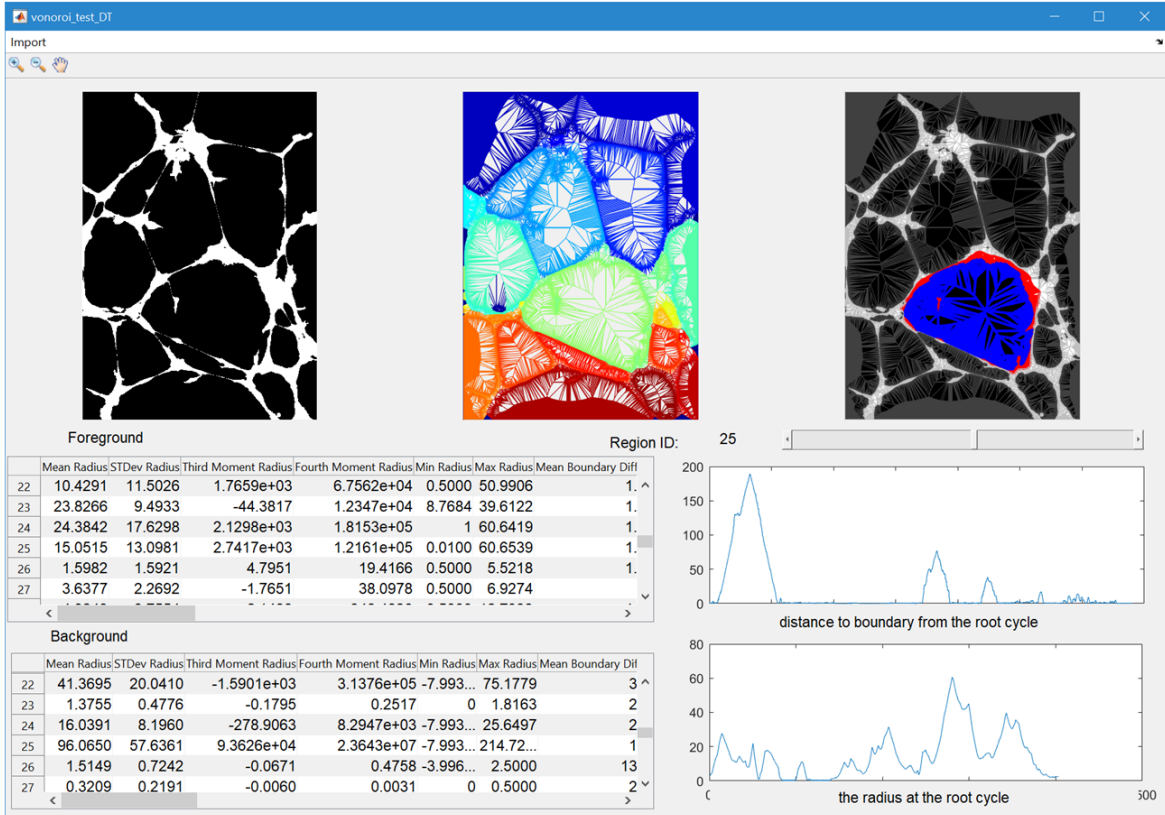505    *Graphical User Interface*

506

507    We have created a graphical user interface (GUI) to provide readers with a quick and easy way to

508    try our shape-to-graph mapping on their own data (**Fig. 15A**). The GUI can be used to generate

509    and display the shape-to-graph mapping for individual images. The user can cycle through all the

29

510    boundaries in the image and visualize their width and boundary profiles. A table of values of the

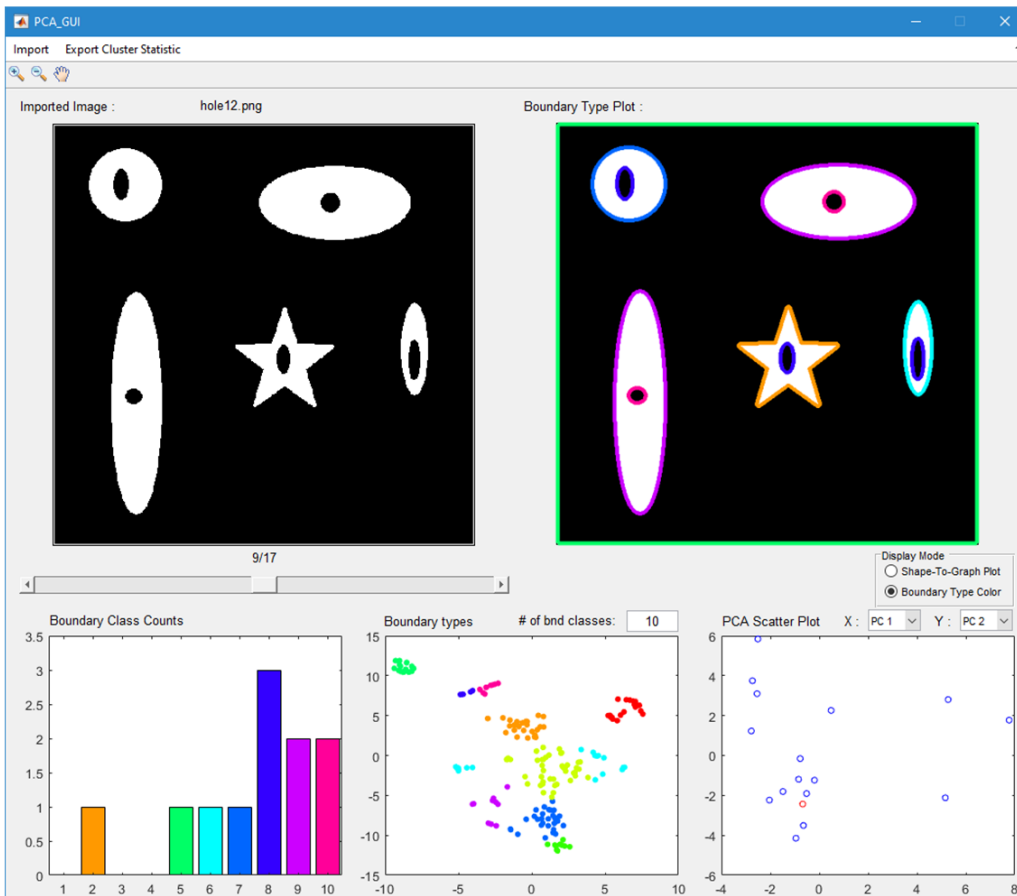511    forty measures for each boundary is also displayed.

512

513    Additionally, a graphical user interface is provided to generate boundary types from multiple

514    images **(Fig. 15B)**. The user can choose a number of boundary classes and inspect each image

515    from the imported set with its boundaries colored according to the class they were automatically

516    assigned based on the features from the shape-to-graph mapping (which can also be displayed).

517    These visualizations are accompanied with (1) a color-coded histogram showing the boundary type

518    distribution in the current image, (2) a t-SNE plot of the boundaries across all images, and (3) a

519    plot of two user-selected principal components calculated based on the boundary type histograms

520    across all the images. The point corresponding to the current image is highlighted in the PCA plot.

A


B


521

522    *Figure 15. Two Graphical User Interfaces for demonstrating the graph construction and analysis. **A.** GUI*

523    *for illustrating the shape-to-graph approach and the key concepts such as subgraph, in- and out-graphs,*

524    *and the width and boundary profiles. The user can cycle through the boundaries and see the 40 metrics*

525    *extracted for each boundary. **B.** GUI for processing multiple images. Boundaries are automatically*

526    *clustered and colored according to a user-specified number of boundary types. The bottom graphs are the*

527    *frequency of boundary types in the current image, a t-SNE of all the boundaries calculated by their features*

528    *and colored by their resulting class, and a PCA plot of all the images derived from their boundary type*

529    *histograms.*

530

## Discussion

531

532

533    In this paper we introduced a methodology for extracting, quantifying, and classifying structural

534    features of an arbitrarily complex pattern in a segmented image. The methodology is based on a

535    mathematically defined mapping of all boundaries in the binary image onto a global graph. The

536    graph preserves all the information specified by the boundaries but also provides an efficient and

537    precise way of defining meaningful metrics for further processing. We illustrated the power of this

538    approach by analyzing experimental images of human umbilical vein endothelial cells forming

539    multicellular patterns with different levels of connectivity depending on genetic (*ccm1*, *ccm2*,

540    *ccm3* knockdowns) and biochemical (Rho kinase inhibition) perturbations. We showed that all the

541    visually distinguishable patterns could be reliably grouped in different classes using principal

542    component analyses of boundary types that were defined based on a large set of graph measures.

543    We also showed that our method is sensitive enough to identify subtle differences in visually

544    similar patterns. More importantly, after classification, the geometric features that made such

545    differentiation possible can be backtracked for further analysis or verification. Thus, our method

546 allows not only for statistical quantification of pattern characteristics but also for the discovery of

547 structural features that are not apparent from visual inspection. This is particularly important for

548 research projects that aim to determine not only 'which' class of patterns a particular image

549 belongs to, but also 'why' it is so in term of intuitively understandable geometric features.

550

551 As another illustration of the strength of our method, we analyzed a set of images generated with

552 a simulation model with two control parameters responsible for the structural organization of the

553 multicellular patterns. We showed that after training the algorithm with a subset of images, it could

554 accurately predict the parameters used for the image generation. It is important to notice that the

555 stochastic nature of cell-cell interactions in the model creates a variability of patterns in different

556 simulations even with the same parameters, which can be interpreted as a noise in the data. Despite

557 this variability, we achieved the correlation coefficients between the predicted and the actual

558 values of the two control parameters as high as 0.9977 and 0.86695. This result shows that a

559 biological characteristic influencing the geometry of an observed structure or pattern can be

560 accurately quantified/predicted directly from the images once the algorithm is trained with a few

561 images for which this characteristic was measured. One of the applications of such quantification

562 would be an investigation of the transition dynamics between the known biological states (e.g.

563 predicting the onset of a diseased phenotype).

564

565 Our methodology works for any binary images. Because we construct the graph for both

566 foreground and background, the extracted features characterize the geometry of individual objects,

567 connectivity in networked structures, as well as the relative organization of isolated objects. This

568 fact makes our method highly versatile and generally applicable. We illustrated this statement

33

569    reanalyzing a subset of previously published data set from a high throughput assay profiling small-

570    molecule-induced U2OS cell cultures (Gustafsdottir et al. 2013). We used the same processing

571    pipeline as in the original study but apply the geometric features from our shape-to-graph mapping.

572    By comparing a combined metric of precision and sensitivity, the $F_1$ score, we showed that our

573    graph representation of the image content provides an improvement in classification performance

574    of 10% for the three major mechanisms-of-action clusters and 40% for the cluster that differs the

575    least from the wild type cultures. Saying that, it is important to notice that the initial, pre-processing

576    step of segmentation is critical and the presented method can be only as accurate as allowed by the

577    quality of microscopy and the segmentation routine.

578

579

580    **Materials and Methods**

581

582    *Cell culture*

583

584    Human umbilical cord endothelial cells HUVEC (Lonza, Walkersville, MD) were maintained in

585    EGM-2 medium (Lonza) at 37°C/5% CO2 and passaged every 3 to 4 days for up to 6 passages at

586    a 1:5 sub-culturing ratio. For tube formation experiments, 4.5-5x10$^3$ cells were plated into each

587    well of angiogenesis μ-slides (ibidi, Fitchburg, WI) coated with 10 μl of growth factor reduced

588    phenol red-free Matrigel (Corning, Corning, NY), and incubated for up to 18 hrs.

589

590    *Microscopy*

591

592   For endothelial tubule formation imaging, cells plated on Matrigel were incubated with

593   CellMask™ Green Plasma Membrane Stain (Invitrogen, Carlsbad, CA) for 15 min at 37°C. The

594   media was changed to phenol-free EGM-2 supplemented with 2% FBS and growth factors

595   (PromoCell GmbH). Images were acquired using PerkinElmer UltraVIEW VoX spinning disk

596   confocal microscope (PerkinElmer, Waltham, MA). Image processing and analysis were

597   performed using ImageJ software (NIH). Images in Figure 9 represent a 1.2 mm by 1.2 mm areas.

598   With the plating density of ~ 400 cells per mm$^2$, there is ~600 cells in each image.

599

600   *Gene expression knockdown*

601

602   To achieve knockdown of CCM protein expression, cells were infected with PLKO.1 vector based

603   lentiviruses carrying shRNAs for human krit1 (RHS4533-EG889), ccm2 (RMM4534-EG216527),

604   and pdcd10 (RHS4533-EG11235) genes (Dharmacon, Lafayette, CO). Lentiviral particles,

605   prepared and purified by VectorBuilder technical service group (VectorBuilder, Santa Clara, CA)

606   were added to EGM-2 media supplemented with 8μ/mL polybrene for 48 hrs. Transduced cells

607   were selected through their resistance to puromycin added to the growth media in the concentration

608   of 2.5 μg/ml. Expression knockdown was measured by real-time PCR with TaqMan gene

609   expression assays. Phenotypic experiments were conducted between 6 and 10 days after infection.

610

611   *Image Preprocessing*

612

613   Simulated images in vector format were rendered at 1024x1024 resolution. By design, the model

614   generates binary images with all interacting cells and their protrusions being the foreground of the

615    image. All holes smaller than 100 pixels were automatically filled. Multiple fields of view were

616    sampled from experimental images of tube formation at a fixed resolution of 690x690 pixels. The

617    images were segmented with a simple threshold followed by manual corrections to under

618    segmented tubules. Cellular debris below 50 pixels in size were automatically removed.

619

620    Boundaries were extracted from each binary image. Linear pixel-size segments that connect

621    boundary points serve as the input to the shape-to-graph mapping algorithm. Rather than defining

622    boundary points at the center of each pixel at the edge of an object, points on the boundary were

623    placed on the half-pixel border between an object and the background. This ensures that any object

624    within the boundary has a non-zero area and any protruding part of an object has a non-zero width.

625    When operating on label images, boundaries are extracted from the largest four-connected

626    components for each label. Boundary points are placed half-way between the center of the pixel

627    and the half-pixel edge used for binary images. This creates a half-pixel sized gap between objects

628    which share a boundary, and any objects which are one pixel wide will have a width in the Voronoi

629    diagram of 0.5px.

630

631    **Acknowledgements**

632

638     Research Alliance, and the Georgia Tech Foundation through their support of the Marcus Center

639     for Therapeutic Cell Characterization and Manufacturing (MC3M) at Georgia Tech.

640

641

642     **References**

643

644     Arnaoutova I, Kleinman HK. 2010. In vitro angiogenesis: endothelial cell tube formation on gelled

645     basement membrane extract. Nat Protoc 5:628-635.

646     Aurenhammer F. 1991. Voronoi Diagrams - a Survey of a Fundamental Geometric Data Structure.

647     Computing Surveys 23:345-405.

648     Boizeau M-L, Fons P, Cousseins L, Desjobert J, Sibrac D, Michaux C, Nestor A-L, Gautret B, Neil K,

649     Herbert C. 2013. Automated image analysis of in vitro angiogenesis assay. Journal of laboratory automation

650     18:411-415.

651     Carpenter AE, et al. 2006. CellProfiler: image analysis software for identifying and quantifying cell

652     phenotypes. Genome Biol 7:R100.

653     Carpentier G, Martinelli M, Courty J, Cascone I. 2012. Angiogenesis analyzer for ImageJ. Pages 198-201.

654     4th ImageJ User and Developer Conference proceedings.

655     Chernaya O, Zhurikhina A, Hladyshau S, Pilcher W, Young KM, Ortner J, Andra V, Sulchek TA,

656     Tsygankov D. 2018. Biomechanics of Endothelial Tubule Formation Differentially Modulated by Cerebral

657     Cavernous Malformation Proteins. iScience 9:347-358.

658     Conrad C, Erfle H, Warnat P, Daigle N, Lorch T, Ellenberg J, Pepperkok R, Eils R. 2004. Automatic

659     identification of subcellular phenotypes on human cell arrays. Genome Res 14:1130-1136.

660     Grélard F, Baldacci F, Vialard A, Domenger J-P. 2017. New methods for the geometrical analysis of tubular

661     organs. Medical image analysis 42:89-101.

662   Guidolin D, Vacca A, Nussdorfer GG, Ribatti D. 2004. A new image analysis method based on topological

663   and fractal parameters to evaluate the angiostatic activity of docetaxel by using the Matrigel assay in vitro.

664   Microvasc Res 67:117-124.

665   Gupta D, Venugopal J, Prabhakaran MP, Dev VR, Low S, Choon AT, Ramakrishna S. 2009. Aligned and

666   random nanofibrous substrate for the in vitro culture of Schwann cells for neural tissue engineering. Acta

667   Biomater 5:2560-2569.

668   Gustafsdottir SM, et al. 2013. Multiplex Cytological Profiling Assay to Measure Diverse Cellular States.

669   Plos One 8.

670   Khoo CP, Micklem K, Watt SM. 2011. A comparison of methods for quantifying angiogenesis in the

671   Matrigel assay in vitro. Tissue Eng Part C Methods 17:895-906.

672   Lin G, Bjornsson CS, Smith KL, Abdul-Karim MA, Turner JN, Shain W, Roysam B. 2005. Automated

673   image analysis methods for 3-D quantification of the neurovascular unit from multichannel confocal

674   microscope images. Cytometry Part A: The Journal of the International Society for Analytical Cytology

675   66:9-23.

676   Ljosa V, Sokolnicki KL, Carpenter AE. 2012. Annotated high-throughput microscopy image sets for

677   validation. Nature Methods 9:637-637.

678   ---. 2013. Annotated high-throughput microscopy image sets for validation (vol 9, pg 637, 2012). Nature

679   Methods 10:445-445.

680   Murphy EA, et al. 2010. Disruption of angiogenesis and tumor growth with an orally active drug that

681   stabilizes the inactive state of PDGFRbeta/B-RAF. Proc Natl Acad Sci U S A 107:4299-4304.

682   Nguyen M, Shing Y, Folkman J. 1994. Quantitation of angiogenesis and antiangiogenesis in the chick

683   embryo chorioallantoic membrane. Microvasc Res 47:31-40.

684   Ogniewicz RL, Kübler O. 1995. Hierarchic voronoi skeletons. Pattern recognition 28:343-359.

685   Rohde GK, Ribeiro AJ, Dahl KN, Murphy RF. 2008. Deformation-based nuclear morphometry: Capturing

686   nuclear shape variation in HeLa cells. Cytometry Part A: The Journal of the International Society for

687   Analytical Cytology 73:341-350.

688    Selinummi J, Seppala J, Yli-Harja O, Puhakka JA. 2005. Software for quantification of labeled bacteria

689    from digital microscope images by automated image analysis. Biotechniques 39:859-863.

690    Styner M, Gerig G, Lieberman J, Jones D, Weinberger D. 2003. Statistical shape analysis of

691    neuroanatomical structures based on medial models. Medical image analysis 7:207-220.

692    Tsygankov D, Bilancia CG, Vitriol EA, Hahn KM, Peifer M, Elston TC. 2014. CellGeo: a computational

693    platform for the analysis of shape changes in cells with complex geometries. J Cell Biol 204:443-460.

694    Viros A, Fridlyand J, Bauer J, Lasithiotakis K, Garbe C, Pinkel D, Bastian BC. 2008. Improving melanoma

695    classification by integrating genetic and morphologic features. PLoS Med 5:e120.

696    Wearne S, Rodriguez A, Ehlenberger D, Rocher A, Henderson S, Hof P. 2005. New techniques for imaging,

697    digitization and analysis of three-dimensional neural morphology on multiple scales. Neuroscience

698    136:661-680.

699    Xin SQ, Wang XN, Xia JZ, Mueller-Wittig W, Wang GJ, He Y. 2013. Parallel computing 2D Voronoi

700    diagrams using untransformed sweepcircles. Computer-Aided Design 45:483-493.

701    Xiong Y, Kabacoff C, Franca-Koh J, Devreotes PN, Robinson DN, Iglesias PA. 2010. Automated

702    characterization of cell shape changes during amoeboid motility by skeletonization. BMC systems biology

703    4:33.

704    Zanella F, Lorens JB, Link W. 2010. High content screening: seeing is believing. Trends Biotechnol

705    28:237-245.

706

## 707    **Supporting information**

708    1. Supplemental Information with Figures and Tables in a single PDF file.

709    2. All Scripts and GUIs in a single ZIP file.

710