

**SUPPLEMENTARY MATERIALS FOR  
“SAFARI: SHAPE ANALYSIS FOR AI-SEGMENTED IMAGES”**

BY ESTEBAN FERNÁNDEZ MORALES<sup>1</sup>, SHENGJIE YANG<sup>2</sup>, SY HAN CHIOU<sup>1</sup>, CHUL  
MOON<sup>3</sup>, CONG ZHANG<sup>1</sup>, BO YAO<sup>2</sup>, GUANGHUA XIAO<sup>2,†</sup>, AND QIWEI LI<sup>1,\*</sup>

<sup>1</sup>*Department of Mathematical Sciences, The University of Texas at Dallas, \* [qiwei.li@utdallas.edu](mailto:qiwei.li@utdallas.edu)*

<sup>2</sup>*Quantitative Biology Research Center, Department of Population and Data Sciences, The University of  
Texas Southwestern Medical Center, † [guanghua.xiao@utsouthwestern.edu](mailto:guanghua.xiao@utsouthwestern.edu)*

<sup>3</sup>*Department of Statistical Science, Southern Methodist University*

## S1. Supplement to Shape Representations.

S1.1. *Binary Matrix.* Let  $\mathbf{M}_{W \times L}$  be a binary matrix representing an arbitrary  $W$ -by- $L$  image, containing a 4-connected region, where the foreground and background are composed of ones and zeros, respectively. Additionally, we can represent each pixel in the image as a point in a 2-dimensional discrete plane, that is, each entry  $M_{wl} \in \mathbf{M}$  can be denoted as a point  $(l, w) \in \mathbb{N}^2$ . Furthermore, to differentiate between foreground and background points, let  $I_R: \mathbb{N}^2 \rightarrow \{0, 1\}$  be the indicator function for an image matrix given by

$$I_R(l, w) = \begin{cases} 1 & \text{if } (l, w) \text{ is a foreground pixel,} \\ 0 & \text{if } (l, w) \text{ is a background pixel.} \end{cases}$$

The indicator function  $I_R$  and distribution of points  $(l, w)$ 's will be used to recreate the region's contour in a two dimensional Cartesian plane, known as the polygonal chain.

S1.2. *Polygonal Chain.* The entries of the binary matrix  $\mathbf{M}_{W \times L}$  that make up the contour of the region can be extracted by the Moore-Neighbor tracing algorithm, modified by Jacob's stopping criteria, with the 1) starting boundary point, 2) direction to traverse the boundary (clockwise or counter-clockwise), and 3) pixel connectivity (Gonzalez, Woods and Eddins, 2020). EBIImage, and as a result also SAFARI, uses a 4-connectivity. Therefore, the first argument is trivial. For the starting boundary point, the point in the lowest left-most location is chosen. Specifically, let

$$S = \left\{ (l_i, w_i) \mid I_R(l_i, w_i) = 1 \wedge w_i = \min_{1 \leq k \leq W} w_k \right\}$$

be the collection of points that make up the region and are located in the lowest  $y$ -coordinate such that

$$(l_1, w_1) = \begin{cases} S_1 & |S| = 1 \\ \min_{l_1} S & \text{otherwise} \end{cases}$$

Applying the Moore-Neighbor tracing algorithm to the region matrix, results in the points that make up the boundary of the region, specifically, where the boundary begins at the lowest left-most area of the region and traverse through the boundary in a clockwise direction. From the boundary points, we can create a sequence of points, known as the closed polygonal chain, that represents the boundary of the region by creating a closed and simple polygon. We can see the binary matrix and its corresponding polygonal chain in Figure S2. Let  $N = \sum_{l,w} I_R(l, w)$  be the total number of points that make up the region and  $\mathbf{P}_{(n+1) \times 2}$  be the collection of points that make up the closed polygonal chain of the region such that 1)  $n \leq N$  is the number of boundary points, 2)  $P_i = (x_i, y_i) \in \mathbb{N}^2$ , for  $i = 1, \dots, n + 1$ , and 3)  $(x_1, y_1) = (x_{n+1}, y_{n+1})$ . Through the polygonal chain, we can derive further two- and one-dimensional shape representations that can be used to compute specific descriptors.

S1.3. *Chain Code.* The slope of a shape's contour can be approximated by the directional changes between two consecutive boundary points. These directional changes can be encoded to, essentially, assign a number (from 0 to 7) to each possible relative direction resulting in an encoding list, each element known as the chain code, that provides a compact representation of the shape's contour (Wirth, 2004; Agu, 2014). Let  $\mathbf{c}$  be a  $1 \times n$  vector representing the chain codes of the polygonal chain where each

entry  $c_i \in \mathbb{N} \cap [0, 7]$  corresponds to a direction in the 8-way split of the unit circle and determined by a series of steps.

First, we determine the angle between the vector composed of the difference between the two consecutive points and the  $x$ -axis, that is, let  $\theta_i = \arctan2(\mathbf{d}_i)$  be the resulting angle where  $\mathbf{d}_i = P_i - P_{i+1}$  is the difference. Since  $\theta_i \in [-\pi, \pi]$ , we have to transform the angle to

$$\hat{\theta}_i = \begin{cases} \theta_i & \theta_i \geq 0, \\ \theta_i + 2\pi & \theta_i < 0, \end{cases}$$

such that  $\hat{\theta}_i \in [0, 2\pi)$ . As a result, we can now determine the corresponding chain code

$$c_i = \left\lfloor \frac{\hat{\theta}_i}{\pi/4} \right\rfloor.$$

Clearly, we can see that this procedure splits the unit circle into eight equal parts. Additionally, if the directional change does not exactly align within the eight splits, then the rounding operator  $\lfloor \cdot \rfloor$  will approximate the chain code to the nearest integer.

**S1.4. Curvature Chain Code.** Let  $\Delta \mathbf{c}$  be a  $1 \times n$  vector representing the curvature chain code such that each entry is formed from a transformation of the difference between two consecutive chain codes, that is, let  $\Delta d_i = c_i - c_{i+1}$  and

$$\Delta c_i = \begin{cases} \Delta d_i - 7 & \Delta d_i > 2, \\ \Delta d_i + 7 & \Delta d_i < -2, \\ \Delta d_i & \text{otherwise.} \end{cases}$$

This simple chain code derivation estimates the curvature and contains information on the convexity of the shape (Wirth, 2004).

**S1.5. Radial Lengths.** Let  $\mathbf{r}$  be a  $1 \times n$  vector of radial lengths, that is, each entry  $r_i = \|P_i - P_c\|$ , for  $i = 1, \dots, n$ , is the Euclidean distance from the boundary point to the shape's centroid. We define the centroid of the polygonal chain  $P_c = (x_c, y_c)$  as

$$x_c = \frac{1}{6A} \sum_{i=1}^n (x_i + x_{i+1})(x_i y_{i+1} - x_{i+1} y_i)$$

$$y_c = \frac{1}{6A} \sum_{i=1}^n (y_i + y_{i+1})(x_i y_{i+1} - x_{i+1} y_i)$$

where

$$A = \frac{1}{2} \sum_{i=1}^n (x_i y_{i+1} - x_{i+1} y_i)$$

is the signed area of the shape, obtained using Gauss's area formula. Clearly, the radial lengths are not scale-invariant (as the Euclidean distance is not). Therefore, to properly analyze the structure of  $\mathbf{r}$ , the individual radial lengths must be normalized.

**S1.6. Normalized Radial Lengths.** Let  $r_{(n)}$  be the maximum radial length in  $\mathbf{r}$  such that we can introduce a  $1 \times n$  vector of normalized radial lengths, denoted as  $\tilde{\mathbf{r}}$ , where each entry  $\tilde{r}_i = r_i / r_{(n)}$ ,  $i = 1, \dots, n$ . By normalizing the radial lengths, we have obtain a 1-dimensional signal that is scale-invariant and which we can use to analyze the fine details of the shape's contour (Wirth, 2004).

TABLE S1  
*Components of the output list resulting from the SAFARI procedure.*

Component	Description
<code>desc</code>	A <code>data.frame</code> object of the shape features corresponding to each segmented ROI.
<code>holes</code>	An integer matrix containing the holes within each ROI, labeled according to the regions.
<code>id</code>	A character vector that is identical to the <code>id</code> argument.
<code>k</code>	A specified factor to enlarge the polygonal chain by with the default being 3.
<code>n</code>	Number of resulting segmented regions.
<code>plg.chains</code>	A <code>list</code> object where each component is the polygonal chain of a segmented ROI.
<code>regions</code>	An integer matrix containing the segmented ROI, labeled from largest to smallest.

\*It takes about 0.3 seconds for SAFARI to run a moderate size binary image ( $300 \times 300$  pixels).

TABLE S2  
*Data notation for shape representations.*

Name	Data	Support
Binary Matrix	$M_{W \times L}$	$M_{wl} \in \{0, 1\}$
Polygonal Chain	$P_{(n+1) \times 2}$	$P_i = (x_i, y_i) \in \mathbb{N}^2$
Convex Hull Chain	$P_{CH}$	$P_{CH_i} \in \mathcal{P}$
Minimum Bounding Box Chain	$P_{MBB}$	$P_{MBB_i} \in \mathbb{R}_+^2$
Chain Code	$c = [c_i]_{1 \times n}$	$c_i \in \mathbb{N} \cap [0, 7]$
Curvature Chain Code	$\Delta c = [\Delta c_i]_{1 \times n}$	$\Delta c_i \in \mathbb{Z} \cap [-2, 2]$
Radial Lengths	$r = [r_i]_{1 \times n}$	$r_i \in \mathbb{R}_+$
Normalized Radial Lengths	$\tilde{r} = [\tilde{r}_i]_{1 \times n}$	$\tilde{r}_i \in [0, 1]$

Table S3: Overview of the shape features available in the SAFARI package.

Category	Data	Name	Formulae	Range	Interpretation	Invariance			Reference
						Rotation	Scale	Translation	
$M$	Net Area		$A_{net}(M) = \sum_{i \in \mathcal{R}} A_{net} = A(M)$	N	Number of pixels that make up the region.	•	•	$\mathcal{O}(LW)$	Agar (2014)
	Thickness <sup>1,2</sup>		$\text{Thickness}(M) = \nu(\tilde{R})$	N	Number of erosion steps that can be applied to the region before the area equals zero.	•	•	$\mathcal{O}(L^2W^2)$	Wirth (2004)
	Elongation <sup>3</sup>		$\text{Elongated}(M) = \frac{A(M)}{2 \sqrt{\text{Elongated}(M)^2}}$	$\mathbb{R}_+$	Relationship between the area of the region and the square of its thickness.	•	•	$\mathcal{O}(L^2W^2)$	Wirth (2004)
$P$	Filled Area <sup>4,5</sup>		$A_{\text{filled}}(P) = \frac{1}{2} \sum_{i=1}^n  x_{i+1} - x_i   y_i  \approx A(M)$	$\mathbb{R}_+$	Approximate area of the region.	•	•	$\mathcal{O}(n)$	Agar (2014)
	Perimeter		$\text{Perimeter}(P) = \sum_{i=1}^n  F_{i+1} - F_i _2 \approx P(M)$	$\mathbb{R}_+$	Approximate length of the boundary that makes up the region.	•	•	$\mathcal{O}(n)$	Agar (2014)
	Circularity		$\text{Circularity}(P) = 4\pi \cdot A_{\text{filled}}(P) / \text{Perimeter}(P)^2$	[0, 1]	Compactness normalized against a filled circle.	•	•	$\mathcal{O}(n)$	Agar (2014)
	Flare Length		$\text{Flare}_{\text{sub}}(P) = \frac{\text{Per}(P) - \sqrt{4\pi(P^2 - 16 \cdot A_{\text{filled}}(P))}}{4}$	$\mathbb{R}_+$	Measurement of the flare of a region.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Flare Width		$\text{Flare}_{\text{sub}}(P) = A_{\text{filled}}(P) / \text{Flare}_{\text{sub}}(P)$	$\mathbb{R}_+$	Measurement of the flare of a region.	•	•	$\mathcal{O}(n)$	Wirth (2004)
Geometric	Convex Area		Measurement can be obtained by applying the same formula as in Filled Area to $P_{\text{CH}}$ .						
	Convex Perimeter		Measurement can be obtained by applying the same formula as in Perimeter to $P_{\text{CH}}$ .						
	$P_{\text{CH}}$ Roundness <sup>6</sup>		$\text{Roundness}(P_{\text{CH}}) = 4\pi \cdot A_{\text{filled}}(P) / \text{Perimeter}(P_{\text{CH}})^2$	[0, 1]	Relationship between the area of the region and the area of a circle with the same convex perimeter i.e. area-to-perimeter ratio.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Convexity <sup>7</sup>		$\text{Convexity}(P_{\text{CH}}) = \text{Perimeter}(P_{\text{CH}}) / \text{Perimeter}(P)$	[0, 1]	Measures how much the region differs from a convex shape.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Solidity <sup>8</sup>		$\text{Solidity}(P_{\text{CH}}) = A_{\text{filled}}(P) / A_{\text{filled}}(P_{\text{CH}})$	[0, 1]	Measures the density of the region.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Major Axis Length		$A_{\text{convex}}(P_{\text{MAB}}) = \frac{1}{2} \sqrt{4A_{\text{filled}}(P) + \frac{1}{4} (b_{\text{MAB}} - a_{\text{MAB}})^2}$	$\mathbb{R}_+$	Length of the region.	•	•	$\mathcal{O}(1)$	
	Major Axis Angle		$A_{\text{convex}}(P_{\text{MAB}}) = \arctan\left(\frac{b_{\text{MAB}}}{a_{\text{MAB}}}\right) \ni j = \arg\max_{\theta \in [0, 2\pi)} \ b_{\theta}\ _2$	$[-\frac{\pi}{2}, \frac{\pi}{2}]$	Orientation of the region.	•	•	$\mathcal{O}(1)$	
	Minor Axis Length		$A_{\text{convex}}(P_{\text{MAB}}) = \frac{1}{2} \sqrt{4A_{\text{filled}}(P) + \frac{1}{4} (b_{\text{MAB}} - a_{\text{MAB}})^2}$	$\mathbb{R}_+$	Width of the region.	•	•	$\mathcal{O}(1)$	
	Bounding Box Area		Measurement can be obtained by computing $A_{\text{convex}}(P_{\text{MAB}}) \cdot A_{\text{convex}}(P_{\text{MAB}})$ .						
	Eccentricity		$\text{Eccentricity}(P_{\text{MAB}}) = A_{\text{convex}}(P_{\text{MAB}}) / A_{\text{convex}}(P_{\text{MAB}})$	[0, 1]	Measures the ellipticity of the region.	•	•	$\mathcal{O}(1)$	Wirth (2004)
Boundary	Curv		$\text{Curv}(P_{\text{MAB}}) = A_{\text{convex}}(P_{\text{MAB}}) / \text{Flare}_{\text{sub}}(P)$	$\mathbb{R}_+$	Measures the degree to which a region is 'curved up'.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Bending Energy		$E_b(\Delta) = \frac{1}{2} \sum_{i=1}^n \Delta x_i^2$	$\mathbb{R}_+$	Energy necessary for a rod to be bent like the region.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Total Abs. Curvature <sup>10</sup>		$\kappa_{\text{total}}(\Delta) = \frac{1}{2} \sum_{i=1}^n  \Delta x_i $	$\mathbb{R}_+$	Another measurement for curvature.	•	•	$\mathcal{O}(n)$	Wirth (2004)
	Radial Mean		$\bar{r}(P) = \frac{1}{n} \sum_{i=1}^n r_i = \bar{r}$	[0, 1]	Measure macroscopic changes in the boundary of the region.	•	•	$\mathcal{O}(n)$	Kilbay, Palmeri and Fox (1993)
	Radial S.D.		$\sigma(r) = \sqrt{\frac{1}{n} \sum_{i=1}^n (r_i - \bar{r})^2} = s_r$	[0, 1]	Similar to Radial Mean but can also indicate fine boundary changes.	•	•	$\mathcal{O}(n)$	Kilbay, Palmeri and Fox (1993)
	Entropy <sup>11,12</sup>		$E_r(P) = -\sum_{i=1}^n p_i \log(p_i) \ni p_i = P(r_i < r < r_i + \Delta)$	$\mathbb{R}_+$	Probabilistic measure of how well the radial lengths can be estimated.	•	•	$\mathcal{O}(n)$	Kilbay, Palmeri and Fox (1993)
	Area Ratio <sup>13</sup>		$A_{\text{ratio}}(P) = \frac{1}{n} \sum_{i=1}^n p_i \log(p_i) \ni p_i = P(r_i < r < r_i + \Delta)$	[0, 1]	Measures the macroscopic characteristics of the region.	•	•	$\mathcal{O}(n)$	Kilbay, Palmeri and Fox (1993)
	Zero Crossing Count <sup>14</sup>		$Z_{\text{zero-crossing}}(P) = \sum_{i=1}^n \left\{ (r_i - r) \cdot (r_{i+1} - r) < 0 \right\}$	N	Captures the fine detail of the boundary.	•	•	$\mathcal{O}(n)$	Kilbay, Palmeri and Fox (1993)
	Normalized Moment Classifier <sup>15</sup>		$\text{NMF}(P) = \left\{ \frac{1}{n} \sum_{i=1}^n (r_i - \bar{r})^2 \right\}^{1/2} - \left[ \frac{1}{n} \sum_{i=1}^n (r_i - \bar{r})^4 \right]^{1/4} / \bar{r}$	[0, 1]	Measures the roughness of the region, based on the moments of inertia of $\tilde{r}$ .	•	•	$\mathcal{O}(n)$	Pohlman et al. (1996)
	Topological	Number of Holes <sup>16</sup>		$N_H(M) =  S(H) _{\text{max}} \ni H = I_H(M) - M$	N		•	•	$\mathcal{O}(L^2W^2)$
Number of Punctures			$N_P(M) = A(\tilde{R}) \ni \tilde{R} = \tilde{R} - \tilde{R} = \tilde{R} - (I_H(M) \oplus E) \circ E$	N		•	•	$\mathcal{O}(L^2W^2)$	

<sup>1</sup> We note that  $\tilde{r} = |M| \cdot (1 - \frac{\omega}{M} - \frac{M}{\omega} - \frac{E}{M} - \frac{M}{\omega}) = 0$  is the set containing the repeated erosion steps applied to the region, where  $E$  is a structuring element.  
<sup>2</sup> The holes within the region must be filled before computing the measurement.  
<sup>3</sup> This measurement applies to curved regions, unlike eccentricity.  
<sup>4</sup> Estimator measurement using Gauss' area formula for polygons (Bradon, 1986).  
<sup>5</sup> Approximation accounts for irregularities, such as an irregular boundary.  
<sup>6</sup> Approximation accounts for irregularities, such as an irregular boundary.  
<sup>7</sup> This measurement accounts for irregularities, such as an irregular boundary.  
<sup>8</sup> A measure less than one accounts for a region with an irregular boundary or containing holes.  
<sup>9</sup> The measurement decreases as the degree to which the region is 'curved up' increases.  
<sup>10</sup> The measurement decreases as the degree to which the region is 'curved up' increases.  
<sup>11</sup> Relies on the 100-bin histogram of the radial lengths where  $P_R$  is the  $k^{\text{th}}$  entry and  $\omega = k / 100 \cdot \Delta = 0.01$ .  
<sup>12</sup> Also incorporates the roundness and roughness of the region.  
<sup>13</sup> This measurement can be simulated as how much of the region is outside the circle with a fine defined by  $\tilde{r}$ .  
<sup>14</sup> This measurement can be simulated as how much of the region is outside the circle with a radius defined by  $\tilde{r}$ .  
<sup>15</sup> A rough boundary will have a greater normalized moment classifier than a smooth boundary.  
<sup>16</sup> We define  $H: X \rightarrow Y$  and  $I_S: X \rightarrow Y$ , where  $X, Y$  are binary matrices, as the procedures to fill the holes within a region and segment the region within an image, respectively (Gomales, Woods and Edlins, 2002).

TABLE S4

*Patient characteristics of the National Lung Screening Trial (NLST) and The Cancer Genome Atlas (TCGA) datasets. Values are either mean  $\pm$  standard deviation, or number (percentage). In the case of the survival time, we use the median instead of the mean.*

	NLST	TCGA
Number of Patients	143	61
Age (in years)	64.01 $\pm$ 5.19	58.26 $\pm$ 12.48
Survival Time (in days)	1517 $\pm$ 730.04	403.93 $\pm$ 291.52
Karnofsky Score (0-100)	–	81.15 $\pm$ 12.92
<b>Status</b>		
Alive	98 (68.53%)	21 (34.42%)
Dead	45 (31.47%)	40 (65.57%)
<b>Gender</b>		
Male	80 (55.94%)	42 (68.85%)
Female	63 (44.06%)	19 (31.15%)
<b>Smoking Status</b>		
Yes	75 (52.45%)	–
No	68 (47.55%)	–
<b>Cancer Stage</b>		
Stage I	95 (66.43%)	–
Stage II	15 (10.49%)	–
Stage III	23 (16.08%)	–
Stage IV	10 (6.99%)	–

TABLE S5

*Univariate analysis of individual shape features in the National Lung Screening Trial (NLST) dataset. A Cox proportional-hazards (CoxPH) model was fitted to each centered and scaled feature, clustered to adjust for patients with multiple samples.*

	Coefficients	Hazard Ratio (HR)	Standard Error (SE)	Robust Standard Error	<i>p</i> -value*
Net Area	0.2679	1.3072	0.0864	0.1173	<b>0.0224</b>
Thickness	0.2733	1.3143	0.0943	0.1119	<b>0.0146</b>
Elongation	-0.1467	0.8636	0.1153	0.1035	0.1563
Area Filled	0.2675	1.3066	0.0869	0.1126	<b>0.0175</b>
Perimeter	0.2900	1.3364	0.1026	0.1280	<b>0.0235</b>
Circularity	0.1623	1.1762	0.1066	0.1163	0.1627
Convex Area	0.2853	1.3301	0.0888	0.1182	<b>0.0158</b>
Convex Perimeter	0.3467	1.4144	0.1017	0.1331	<b>0.0092</b>
Roundness	0.1249	1.1331	0.1140	0.1268	0.3243
Convexity	0.1165	1.1236	0.1154	0.1359	0.3913
Solidity	0.2796	1.3226	0.1210	0.1383	<b>0.0433</b>
Major Axis Length	0.3934	1.4820	0.1046	0.1359	<b>0.0038</b>
Major Axis Angle	-0.0213	0.9790	0.1114	0.1166	0.8553
Minor Axis Length	0.2857	1.3307	0.1020	0.1206	<b>0.0179</b>
Bounding Box Area	0.3092	1.3624	0.0908	0.1175	<b>0.0085</b>
Eccentricity	-0.1456	0.8645	0.1101	0.1239	0.2398
Fibre Length	0.3093	1.3625	0.0952	0.1132	<b>0.0063</b>
Fibre Width	0.2860	1.3310	0.1028	0.1281	<b>0.0255</b>
Curl	-0.0705	0.9319	0.1160	0.1222	0.5640
Bending Energy	0.0054	1.0054	0.1074	0.1010	0.9572
Total Abs. Curvature	0.0074	1.0074	0.1072	0.0991	0.9405
Radial Mean	0.0682	1.0706	0.1107	0.1271	0.5916
Radial S.D.	-0.0586	0.9431	0.1079	0.1052	0.5776
Entropy	-0.0624	0.9395	0.1060	0.1037	0.5475
Area Ratio	-0.0504	0.9509	0.1094	0.1131	0.6560
Zero Crossing	-0.0354	0.9652	0.1087	0.1074	0.7418
Normalized Moment	-0.1409	0.8686	0.1152	0.1383	0.3082
Number of Holes	0.2289	1.2572	0.0844	0.0911	<b>0.0120</b>
Number of Protrusions	0.2877	1.3334	0.1029	0.1282	<b>0.0248</b>

\*Bolding signifies features with *p*-value  $\leq 0.05$

TABLE S6

Comparison of the results obtained in our univariate study, compared to those in Wang et al. (2018). The p-value in Wang et al. (2018) corresponds to either the sum of the shape feature for all regions or the shape feature for the main region, whichever was more significant. Additionally, bolding signifies features with p-value  $\leq 0.05$  and we show features not present in either study.

	Our Study	Wang's Paper
Average Tumor Probability	–	0.78
Net Area	<b>0.0224</b>	<b>0.0033</b>
Thickness	<b>0.0146</b>	–
Elongation	0.1563	–
Area Filled	<b>0.0175</b>	<b>0.0029</b>
Perimeter	<b>0.0235</b>	<b>0.0034</b>
Circularity	0.1627	<b>0.019*</b>
Convex Area	<b>0.0158</b>	<b>0.0047</b>
Convex Perimeter	<b>0.0092</b>	–
Roundness	0.3243	–
Convexity	0.3913	–
Solidity	<b>0.0433</b>	0.16
Major Axis Length	<b>0.0038</b>	<b>0.0099</b>
Major Axis Angle	0.8553	0.92
Minor Axis Length	<b>0.0179</b>	<b>0.030</b>
Bounding Box Area	<b>0.0085</b>	–
Eccentricity	0.2398	0.13
Extent	–	0.34
Fibre Length	<b>0.0063</b>	–
Fibre Width	<b>0.0255</b>	–
Curl	0.5640	–
Bending Energy	0.9572	–
Total Abs. Curvature	0.9405	–
Radial Mean	0.5916	–
Radial S.D.	0.5776	–
Entropy	0.5475	–
Area Ratio	0.6560	–
Zero Crossing	0.7418	–
Normalized Moment	0.3082	–
Number of Holes	<b>0.0120</b>	<b>0.0033</b>
Number of Protrusions	<b>0.0248</b>	–
Number of Regions	–	0.48

\*Measure corresponds to a variation of the formula we used.



TABLE S7

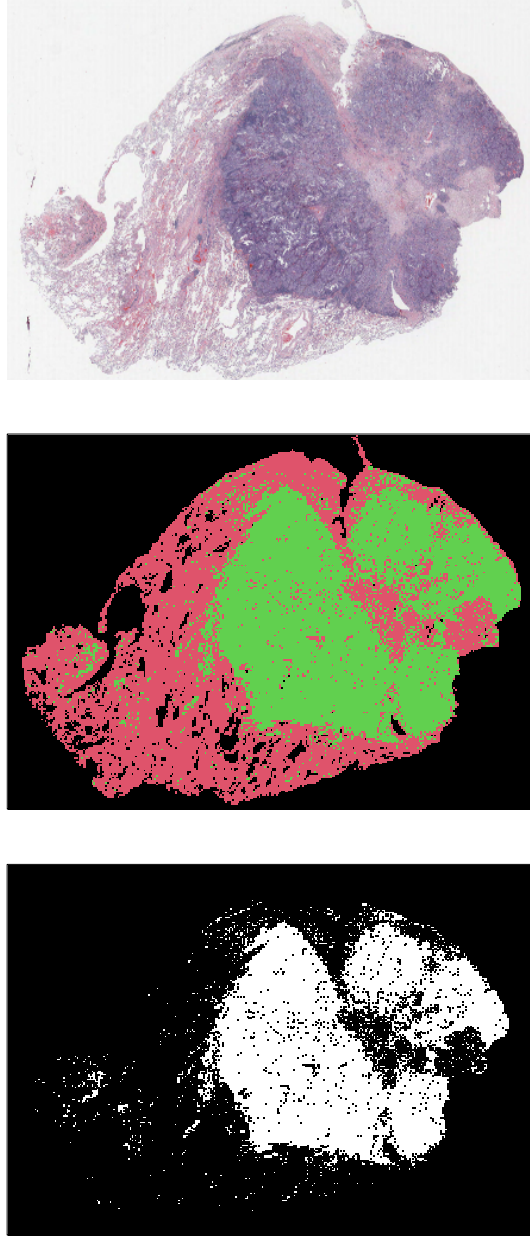
*Univariate analysis of individual shape features in The Cancer Genome Atlas (TCGA) dataset. A Cox proportional-hazards (CoxPH) model was fitted to each centered and scaled feature.*

	Coefficients	Hazard Ratio (HR)	Standard Error (SE)	<i>p</i> -value
Net Area	0.5784	1.7832	0.1706	<b>&lt;0.001</b>
Thickness	0.327	1.3868	0.1821	0.0726
Elongation	-0.1706	0.8431	0.287	0.5522
Area Filled	0.4361	1.5467	0.1843	<b>0.0180</b>
Perimeter	0.4182	1.5192	0.1975	<b>0.0342</b>
Circularity	0.0962	1.101	0.1631	0.5555
Convex Area	0.555	1.742	0.1936	<b>0.0041</b>
Convex Perimeter	0.4917	1.6351	0.2001	<b>0.0140</b>
Roundness	0.2225	1.2492	0.1839	0.2264
Convexity	0.099	1.1041	0.1743	0.5700
Solidity	0.1957	1.2161	0.1807	0.2790
Major Axis Length	0.4776	1.6121	0.2006	<b>0.0173</b>
Major Axis Angle	0.5612	1.7528	0.2104	<b>0.0076</b>
Minor Axis Length	0.5108	1.6666	0.198	<b>0.0099</b>
Bounding Box Area	0.5371	1.7111	0.1913	<b>0.005</b>
Eccentricity	0.0695	1.072	0.1501	0.6435
Fibre Length	0.2615	1.2989	0.1776	0.1410
Fibre Width	0.3954	1.485	0.1966	<b>0.0443</b>
Curl	-0.0652	0.9368	0.2386	0.7845
Bending Energy	0.14	1.1502	0.165	0.3964
Total Abs Curvature	0.1501	1.1619	0.1655	0.3646
Radial Mean	0.2061	1.2289	0.171	0.2281
Radial Sd	0.0786	1.0818	0.1649	0.6336
Entropy	0.146	1.1572	0.1624	0.3689
Area Ratio	0.0869	1.0908	0.1696	0.6084
Zero Crossing	0.2271	1.255	0.1638	0.1656
Normalized Moment	0.2168	1.2421	0.1712	0.2053
Number Holes	0.4667	1.5947	0.1963	<b>0.0174</b>
Number Protrusions	0.3618	1.4359	0.1926	0.0603

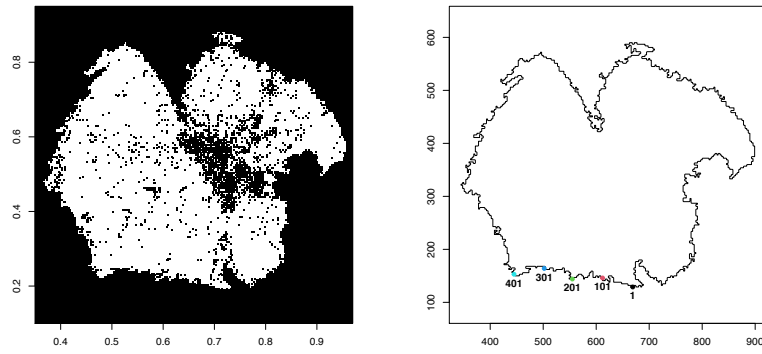
\*Bolding signifies features with *p*-value  $\leq 0.05$

**Table S8:** A summary of SAFARI and other related shape analysis tools in R.

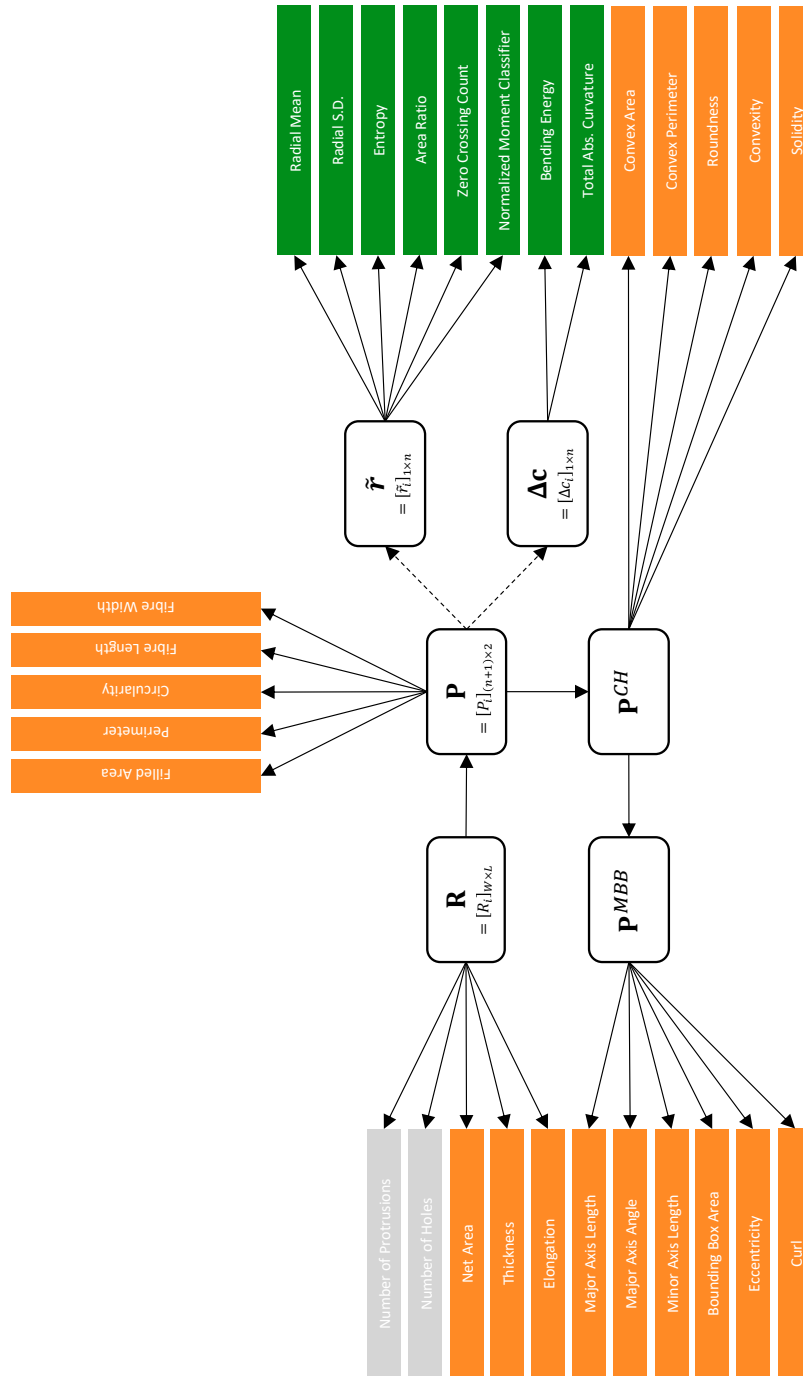
Tool	Description/Objectives	Data	# of Shape Features	Applications	Latest Version	Availability			Citation
						Article	Package	Web Tool	
<b>EImage</b>	This package provides general purpose functionality for image processing and analysis, by offering tools to segment cells and extract quantitative cellular descriptors and automating image processing tasks.	$M$	9	Microscopy-based cellular assays	4.36.0 (2021-12-19)	✓	✓		[1]
<b>SAFARI</b>	This package provides functionality for image processing and shape analysis in the context of segmented medical images generated by deep learning-based methods or standard image processing algorithms and produced from different medical imaging types. Specifically, offers tools to segment regions of interest and extract quantitative shape descriptors.	$M, P, c, r$	29	AI-segmented images	0.1.0 (2021-02-25)	✓	✓	✓	[2]
<b>shapes</b>	This package offers routines for the statistical analysis of landmark shapes.	$P$	0	Landmark shapes	1.2.6 (2021-03-30)	✓	✓		[3]
<b>wrtool</b>	This package intends to facilitate preprocessing and analyzing wood images toward automated recognition.	$M$	0	Wood images	1.0.0 (2016-11-08)		✓		[4]



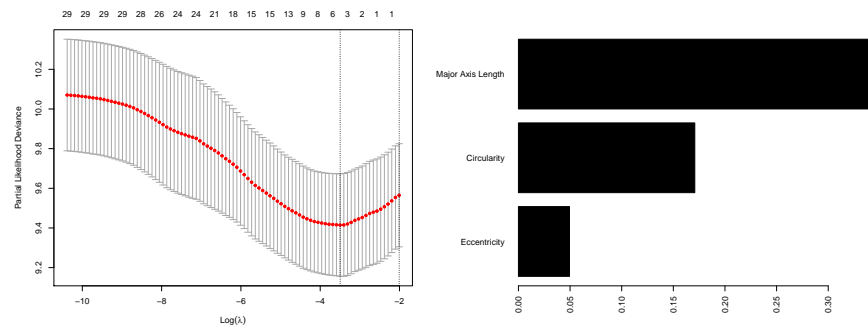
**Fig S1:** An example of a whole-slide image from the National Lung Screening Trial (NLST) cohort processed by an Automated Tumor Recognition System [Wang et al. \(2018\)](#) and then converted into a binary format. The images are whole-slide pathology image (top), segmented three-class image (middle), and segmented two-class or binary image (bottom).



**Fig S2:** An example of the binary matrix (left) and its corresponding polygonal chain (right). The polygonal chain also shows the starting point and four sample points to demonstrated the contour's clockwise direction.



**Fig S3:** Dependencies of shape features and representations. Colored boxes refer to the shape feature categories. Orange denotes geometric, green boundary, and grey topological.



**Fig S4:** Results from the regularized Cox proportional-hazards (CoxPH) model in **Downstream Analysis II: Predictive Performance**. The left figure shows the mean cross-validated errors, based on the Partial Likelihood Deviance, for each tuning parameter. The right figure shows the importance of each feature kept by the regularized Cox model, based on the magnitude of each coefficient.

## REFERENCES

- AGU, E. (2014). Digital Image Processing (CS/ECE 545) Lecture 8: Regions in Binary Images (Part 2) and Color (Part 1). Accessed: May 13, 2021.
- BRADEN, B. (1986). The Surveyor's Area Formula. *The College Mathematics Journal* **17** 326–337.
- GONZALEZ, R. C., WOODS, R. E. and EDDINS, S. L. (2020). *Digital image processing using MATLAB*, Third edition ed. Gatesmark Publishing, Knoxville.
- KILDAY, J., PALMIERI, F. and FOX, M. D. (1993). Classifying mammographic lesions using computerized image analysis. *IEEE Transactions on Medical Imaging* **12** 664–669.
- POHLMAN, S., POWELL, K. A., OBUCHOWSKI, N. A., CHILCOTE, W. A. and GRUNDFEST-BRONIATOWSKI, S. (1996). Quantitative classification of breast tumors in digitized mammograms. *Medical Physics* **23** 1337–1345.
- WANG, S., CHEN, A., YANG, L., CAI, L., XIE, Y., FUJIMOTO, J., GAZDAR, A. and XIAO, G. (2018). Comprehensive analysis of lung cancer pathology images to discover tumor shape and boundary features that predict survival outcome. *Scientific Reports* **8** 10393.
- WIRTH, M. A. (2004). Shape Analysis & Measurement. Accessed: May 13, 2021.