

1 **Robust integrated intracellular organization of the human iPSC cell: where,** 2 **how much, and how variable**

3 Matheus P. Viana, Jianxu Chen¹, Theo A. Knijnenburg¹, Ritvik Vasan¹, Calysta Yan¹, Joy E. Arakaki,
4 Matte Bailey, Ben Berry, Antoine Borensztein, Jackson M. Brown, Sara Carlson, Julie A. Cass,
5 Basudev Chaudhuri, Kimberly R. Cordes Metzler, Mackenzie E. Coston, Zach J. Crabtree, Steve
6 Davidson, Colette M. DeLizo, Shailja Dhaka, Stephanie Q. Dinh, Thao P. Do, Justin Domingus, Rory
7 M. Donovan-Maiye, Tyler J. Foster, Christopher L. Frick, Griffin Fujioka, Margaret A. Fuqua, Jamie L.
8 Gehring, Kaytlyn A. Gerbin, Tanya Grancharova, Benjamin W. Gregor, Lisa J. Harrylock, Amanda
9 Haupt, Melissa C. Hendershott, Caroline Hookway, Alan R. Horwitz, Chris Hughes, Eric J. Isaac,
10 Gregory R. Johnson, Brian Kim, Andrew N. Leonard, Winnie W. Leung, Jordan J. Lucas, Susan A.
11 Ludmann, Blair M. Lyons, Haseeb Malik, Ryan McGregor, Gabe E. Medrash, Sean L. Meharry, Kevin
12 Mitcham, Irina A. Mueller, Timothy L. Murphy-Stevens, Aditya Nath, Angelique M. Nelson, Luana
13 Paleologu, T. Alexander Popiel, Megan M. Riel-Mehan, Brock Roberts, Lisa M. Schaeffbauer,
14 Magdalena Schwarzl, Jamie Sherman, Sylvain Slaton, M. Filip Sluzewski, Jacqueline E. Smith,
15 Youngmee Sul, Madison J. Swain-Bowden, W. Joyce Tang, Derek J. Thirstrup, Daniel M. Toloudis,
16 Andrew P. Tucker, Veronica Valencia, Winfried Wiegraebe, Thushara Wijeratna, Ruian Yang,
17 Rebecca J. Zaunbrecher, Allen Institute for Cell Science, Graham T. Johnson, Ruwanthi N.
18 Gunawardane, Nathalie Gaudreault, Julie A. Theriot², Susanne M. Rafelski³

19 ¹These authors contributed equally

20 ³For correspondence: susanner@alleninstitute.org

21
22 Affiliations:

23 All authors: Allen Institute for Cell Science, 615 Westlake Ave N, Seattle, WA, 98125, USA.

24 ²Department of Biology and Howard Hughes Medical Institute, University of Washington, Seattle, WA
25 98195

26 27 28 **Summary**

29 Despite the intimate link between cell organization and function, the principles underlying intracellular
30 organization and the relation between organization, gene expression and phenotype are not well
31 understood. We address this by creating a benchmark for mean cell organization and the natural
32 range of cell-to-cell variation. This benchmark can be used for comparison to other normal or
33 abnormal cell states. To do this, we developed a reproducible microscope imaging pipeline to
34 generate a high-quality dataset of 3D, high-resolution images of over 200,000 live cells from 25
35 isogenic human induced pluripotent stem cell (hiPSC) lines from the Allen Cell Collection. Each line
36 contains one fluorescently tagged protein, created via endogenous CRISPR/Cas9 gene editing,
37 representing a key cellular structure or organelle. We used these images to develop a new multi-part
38 and generalizable analysis approach of the locations, amounts, and variation of these 25 cellular
39 structures. Taking an integrated approach, we found that both the extent to which a structure's
40 individual location varied ("stereotypy") and the extent to which the structure localized relative to all

41 the other cellular structures (“concordance”) were robust to a wide range of cell shape variation, from
42 flatter to taller, smaller to larger, or less to more polarized cells. We also found that these cellular
43 structures varied greatly in how their volumes scaled with cell and nuclear size. These analyses
44 create a data-driven set of quantitative rules for the locations, amounts, and variation of 25 cellular
45 structures within the hiPSC as a normal baseline for cell organization.

46

47 **Introduction**

48 A living cell must organize all of its millions of subcellular components and processes in space
49 and time through as many as four orders of magnitude. At the nanometer scale, specific molecular
50 interactions permit the assembly of macromolecules and organelles to perform and regulate cell
51 function. More global cell behaviors, however, can occur over scales of tens of microns, such as the
52 coordinated protrusion of a cell front and retraction of a rear during cell migration (Lauffenburger and
53 Horwitz, 1996). Identifying the rules of cell organization and understanding how they facilitate global
54 behaviors across this broad span of spatial scales is an immensely complex and daunting task that
55 must also incorporate dynamic changes across a broad temporal spectrum. However, to understand
56 cell organization at the level of the major intracellular machinery and organelles (cellular structures),
57 requires the study of only ~25-50 of these structures. This enormously reduces the dimensional
58 complexity, making feasible the quest for an interpretable and testable set of rules that govern cell
59 organization and how this organization changes as cells transition to alternative normal or abnormal
60 cell behaviors. For example, measuring the locations of each of these cellular structures relative to all
61 the others, as well as the total volume occupied by each structure, creates a rich set of quantitative
62 rule-building constraints for generating and testing models of cell organization (Johnson et al., 2015;
63 Macklin et al., 2020).

64 A significant potential challenge, even for this approach, is that cells must behave robustly yet
65 respond sensitively to their ever-changing environments. As a result, a population of normal,
66 putatively identical cells might exhibit significant cell-to-cell variability. Thus, it is important to establish
67 a baseline with which different kinds of cells can be compared. This baseline should represent the
68 typical, or mean, cell within the population, as well as the full range of normal variation of the
69 population itself. An abnormal cell phenotype may exhibit not only a shift in the mean but also a shift
70 in the variation (Roggiani and Goulian, 2015). Therefore, a meaningful and useful description of
71 normal cell organization requires quantitative measurements, not just of the locations or amounts of
72 each of the cellular structures, but also how they vary within a large group of normal cells.

73 To establish a normal baseline for cell organization, we turned to human induced pluripotent
74 stem cells (hiPSCs), which represent an early embryonic cell state and an ideal human model system.
75 hiPSCs are naturally immortal, karyotypically normal, and can be induced to differentiate into other

76 cell types (Drubin and Hyman, 2017). We previously developed methods to generate a series of
77 isogenic clonal hiPSC lines expressing fluorescent protein tags for visualizing specific organelles and
78 cellular structures via endogenous CRISPR/Cas9 gene editing. We performed extensive quality
79 control on these lines to create the *Allen Cell Collection* (Roberts et al., 2017a). In this work, we
80 imaged 25 lines from this collection, each containing one fluorescently tagged protein representing a
81 key cellular structure or organelle.

82 Here we present the *hiPSC Single-Cell Image Dataset*, an unprecedented collection of high-
83 resolution, 3D images of over 200,000 live cells. To analyze this large-scale dataset, we develop
84 generalizable, quantitative methods that permit direct comparisons of the similarity of overall cell and
85 nuclear shapes for 3D cell image data, to build a simple and human-interpretable “shape space”. This
86 approach facilitates the robust identification of clusters of cells that are most similar to each other in
87 their overall shape. We also introduce a generalizable method to parameterize fluorescence intensity
88 distributions in 3D cell images; this method allows the actual distribution observed for a particular cell
89 to be robustly “morphed” into another similar cell shape, without losing substantial quantitative
90 information about fluorescence localization. We then apply these methods to determine which
91 structures are most highly stereotyped with respect to their cell-to-cell variation, and also which pairs
92 of structures are most similar to each other, throughout the normal hiPSC shape space. These
93 analyses create a data-driven set of quantitative rule-building constraints for the locations, amounts,
94 and variation of 25 cellular structures within the hiPSC as a normal baseline for cell organization and
95 a fundamental benchmark for comparison with future analyses of cell shape and cell organization for
96 cells in different states.

97

98 **Results**

99

100 **An hiPSC Single-Cell Image Dataset contains over 200,000 live, high-resolution, 3D cells** 101 **spanning 25 cellular structures**

102 We previously developed methods and quality control workflows to create the Allen Cell
103 Collection (Roberts et al., 2017a) and www.allencell.org of hiPSC lines, each expressing a single
104 endogenously tagged protein representing a particular organelle or cellular structure. Here we created
105 15 additional cell lines and used a total set of 25 cell lines permitting a holistic view of cells at the level
106 of their the major organelles, cellular structures, and compartments (**Table 1**).

107 hiPSCs grow in tightly packed, epithelial-like monolayer colonies (Roberts et al., 2017a),
108 requiring well-defined imaging assay guidelines for reproducible data collection. We grow these cells
109 on Matrigel-coated glass plates compatible with high-resolution imaging while preserving their normal
110 pluripotent state (Roberts et al., 2017a). We built an, automated microscopy imaging pipeline to

111 **Table 1:** Fluorescently tagged cellular structures used to create the hiPSC Single-Cell Image Dataset

Structure	Protein	Gene	Location	AICS line
nucleoli (DFC)	fibrillarin	<i>FBL</i>	nuclear	AICS-0014 cl. 6
nucleoli (GC)	nucleophosmin	<i>NPM1</i>	nuclear	AICS-0057 cl. 50
nuclear speckles	SON	<i>SON</i>	nuclear	AICS-0094 cl. 24
cohesins	SMC-1A	<i>SMC1A</i>	nuclear	AICS-0068 cl. 9
histones	H2B	<i>H2BC11</i>	nuclear	AICS-0061 cl. 36
nuclear envelope	lamin B1	<i>LMNB1</i>	nuclear periphery	AICS-0013 cl. 210
nuclear pores	Nup153	<i>NUP153</i>	nuclear periphery	AICS-0069 cl. 88
ER (Sec61 beta)	Sec61 beta	<i>SEC61B</i>	cytoplasm	AICS-0010 cl. 55
ER (SERCA2)	SERCA2	<i>ATP2A2</i>	cytoplasm	AICS-0046 cl. 51
mitochondria	Tom20	<i>TOMM20</i>	cytoplasm	AICS-0011 cl. 27
peroxisomes	PMP34	<i>SLC25A17</i>	cytoplasm	AICS-0033 cl. 115
endosomes	Rab-5A	<i>RAB5A</i>	cytoplasm	AICS-0040 cl. 35*
lysosomes	LAMP-1	<i>LAMP1</i>	cytoplasm	AICS-0022 cl. 37
Golgi	sialyltransferase 1	<i>ST6GAL1</i>	cytoplasm	AICS-0025 cl. 44*
centrioles	centrin-2	<i>CETN2</i>	cytoplasm	AICS-0032 cl. 19***
microtubules	alpha-tubulin	<i>TUBA1B</i>	cytoplasm	AICS-0012 cl. 105
plasma membrane	CAAX	<i>AAVS1</i>	cell periphery	AICS-0054 cl. 91***
actin filaments	beta-actin	<i>ACTB</i>	cell periphery	AICS-0016 cl. 184
actin bundles	alpha-actinin-1	<i>ACTN1</i>	cell periphery	AICS-0007 cl. 79**
actomyosin bundles	non-muscle myosin IIB	<i>MYH10</i>	cell periphery	AICS-0024 cl. 80
gap junctions	connexin-43	<i>GJA1</i>	cell periphery	AICS-0053 cl. 16
tight junctions	ZO-1	<i>TJP1</i>	cell periphery	AICS-0023 cl. 20
desmosomes	desmoplakin	<i>DSP</i>	cell periphery	AICS-0017 cl. 65
adherens junctions	beta-catenin	<i>CTNNB1</i>	cell periphery	AICS-0058 cl. 67
matrix adhesions	paxillin	<i>PXN</i>	cell periphery	AICS-0005 cl. 50

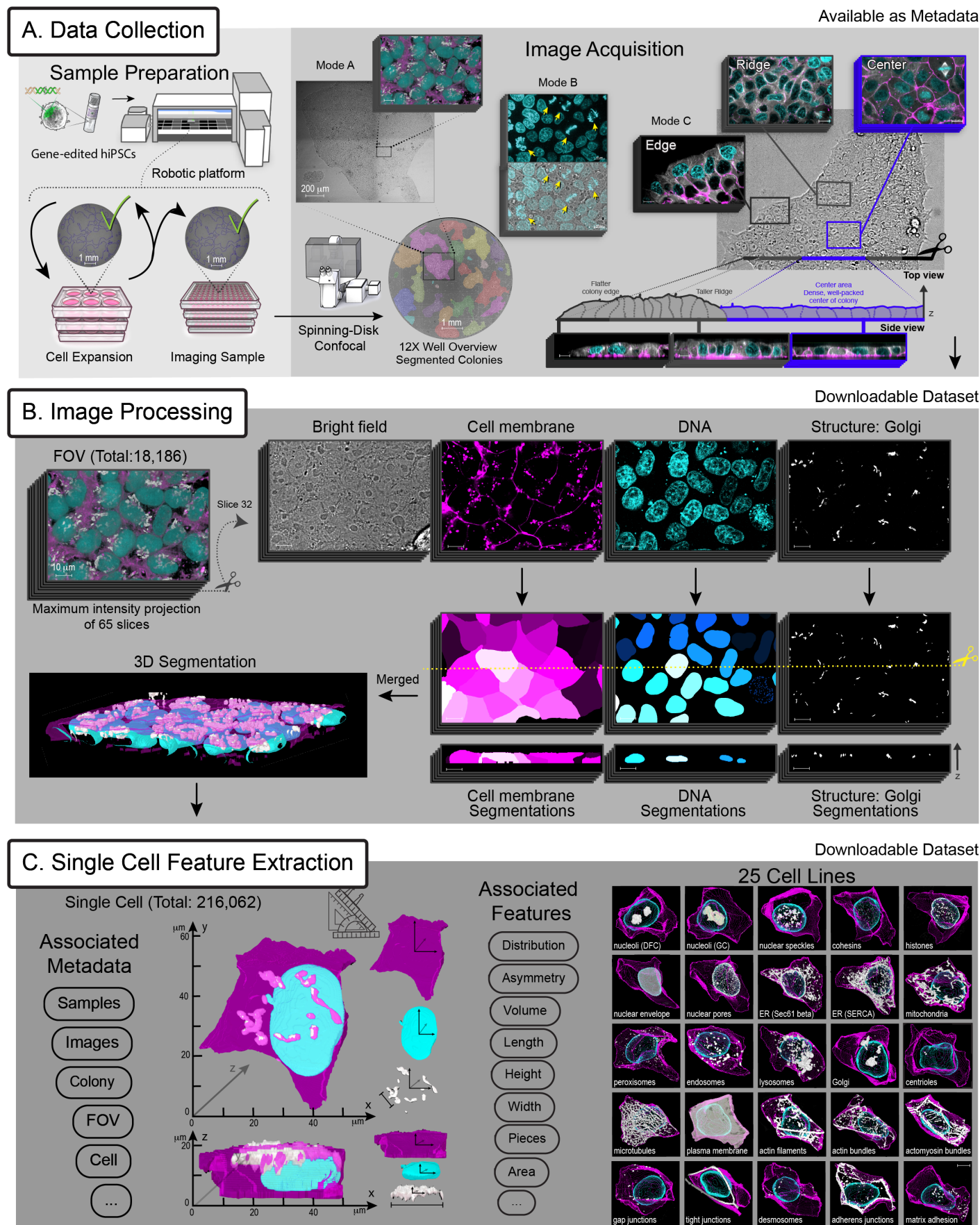
* biallelic edit, ** line not yet released, *** mTagRFP-T fluorophore tag instead of mEGFP fluorophore tag
 Bold AICS line number indicates cell lines published previously (Roberts, et al., 2017).

112 reproducibly generate the living colonies, imaged the cells in 3D using spinning-disk confocal
113 microscopes to collect standardized 3D information, and processed these images to create the hiPSC
114 Single-Cell Image Dataset (**Figure 1**). We mostly imaged cells halfway towards the centers of large,
115 well-packed colonies, as cells behaved most consistently in this region. We also captured variations in
116 colony area and locations within the colony, and enriched for images with mitotic cells when
117 necessary (**Figure 1A**). To keep track of the position of each image within each colony, we collected
118 low magnification overview images of the entire well prior to the high magnification imaging.

119 We included fluorescent cell membrane and DNA dyes to reference the locations of
120 intracellular fluorescent protein (FP)-tagged structures relative to the cell boundary and the nucleus or
121 mitotic chromosomes. Cells were imaged live and in 3D at high resolution (120x magnification, 1.25
122 numerical aperture, NA), generating 18,186 fields of view (FOVs) in four acquisition channels,
123 representing the fluorescently tagged protein, the cell membrane and DNA dyes, and the transmitted
124 light channel (**Figure 1A&B**). Measuring the locations of each of the 25 cellular structures within cells
125 required segmentation approaches that demarcate structure, as well as the cell and nuclear
126 boundaries within these 3D images. To do this, we used the *Allen Cell and Structure Segmenter* (the
127 Segmenter), a fully-accessible, Python-based 3D segmentation software package (Chen et al., 2018).
128 For each of the 25 cellular structures, we used the tagged protein to identify the location and
129 morphology of the structure, rather than the location of the FP-tagged protein, itself (**Figure S1**). The
130 tightly packed, epithelial-like nature of hiPSCs, as well as the need for highly-accurate 3D cell
131 boundaries to minimize cellular structure misassignment to neighboring cells required deep learning-
132 based segmentation approaches to create a robust, scalable, and highly accurate 3D cell and nuclear
133 segmentation algorithm ((Chen et al., 2018) and *3D Segmentation* in Methods) applicable to all FOVs
134 in this dataset (**Figure S1**).

135 From each FOV, individual cells were segmented using the plasma membrane dye, resulting
136 in a single-cell image dataset consisting of 216,062 cells (**Figure 1C**). Every individual cell is labeled
137 with a unique ID, permitting the persistent association of relevant metadata including the full set of
138 experimental parameters, position within the original FOV, and structure segmentations with
139 versioned software captured for future data provenance. Additional cell, nuclear, and structural
140 features were extracted and associated with each cell ID, e.g. cellular structure volume, generating a
141 rich single-cell image dataset for analysis. Both the FOV images and the single-cell dataset are
142 available for use by the community as downloadable files (allencell.org/data-downloading.html) and
143 through interactive online visual analysis tools that require no software installation or expertise
144 (cfe.allencell.org). For the analyses described below, we used the subset of 203,737 interphase cells,
145 excluding the 11,238 cells undergoing mitosis and a few outliers (**Table S1**).

Figure 1



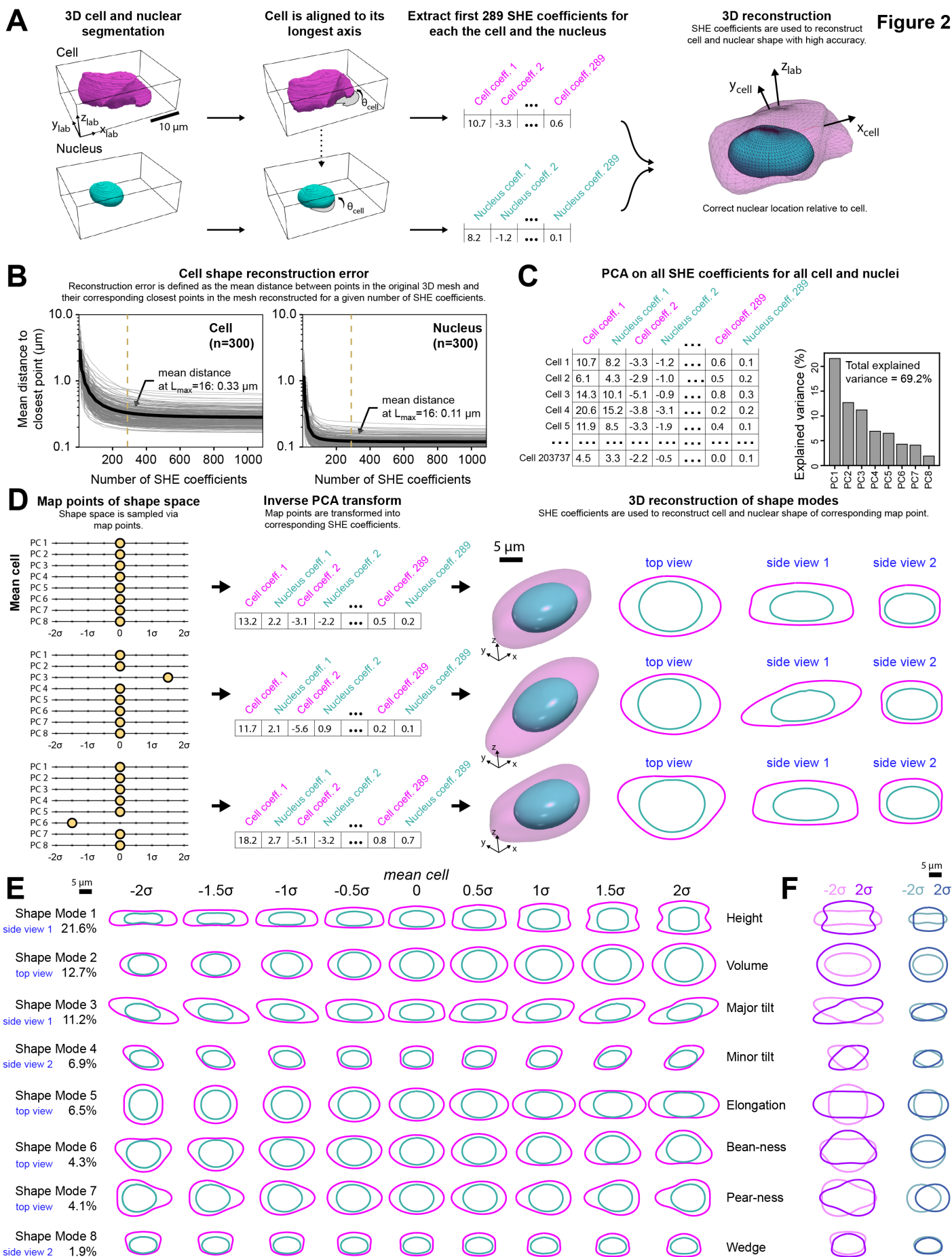
147 **Figure 1.** An hiPSC Single-Cell Image Dataset contains over 200,000 live, high-resolution, 3D cells
148 spanning 25 cellular structures. The dataset was generated by a microscopy pipeline composed of
149 three main parts; Data Collection, Image Processing and Single Cell Feature Extraction. **A)** Data
150 Collection: the sample preparation starts with a vial of frozen gene-edited hiPSCs from a line from the
151 Allen Cell Collection, expressing an endogenous, fluorescently tagged protein representing a
152 particular cellular structure. The cell cultures are expanded in 6-well plates on an automated cell
153 culture platform. At each passage cells are seeded into optical grade, glass bottom 96-well plates to
154 create imaging samples. Bright field overview images of each well are inspected and only wells
155 meeting pre-determined quality controls are passaged from the 6-well plates and imaged from the 96-
156 well plates. The image acquisition of live cells starts with a 12X overview image of each well on a
157 spinning-disk confocal microscope. Imaging sessions are conducted using three modes. In mode A,
158 the 12X overview images of colonies are segmented by an automated script to generate sets of
159 coordinates for positions within imageable colonies, located approximately halfway between the
160 colony edge and colony center. Imageable colonies are those that meet size, morphology, and
161 position-within-a-well criteria. In mode B, the microscope operator adjusts the location of the field of
162 view (FOV) to enrich for mitotic cells via appropriate cell and DNA morphology visible with live bright
163 field viewing and confirmed by DNA staining (yellow arrows). In mode C, three regions of colonies are
164 imaged, the edge, ridge (just inward from the edge), and center. Cells were labeled with fluorescent
165 DNA and membrane dyes and then imaged at each pre-selected colony position. Z-stacks were
166 acquired at 120X in four channels, representing the bright field, cell membrane dye (magenta), DNA
167 dye (cyan) and the fluorescently tagged cellular structure (grayscale), also shown in (B). Mode A and
168 C panels show Golgi (via sialyltransferase) and microtubules (via alpha-tubulin), respectively. **B)**
169 Image Processing: The hiPSC Single-Cell Image Dataset consists of a total of 18,186 curated FOVs,
170 which are available for download. An example z-stack is shown. On the left is the maximum intensity
171 projection of all 65 slices with all fluorescent channels combined, in the colors indicated in the panels
172 on the right. “Cutting” the z-stack in half exposes the view of a single slice (slice 32) in the middle of
173 the stack, shown for each individual channel, including the bright field channel. We applied 3D
174 segmentation algorithms to each of the fluorescent channels to identify boundaries in 3D of the cells
175 via the membrane dye (magenta), the nuclei and mitotic DNA via the DNA dye (cyan), and each of the
176 25 cellular structures via their fluorescent protein tag (grayscale; Golgi shown here). Resulting 3D
177 segmentations for cell membrane, DNA, and structure channels are also shown as a side view, the
178 xz-cross-section along the yellow dotted line. All segmentation algorithms were developed and
179 performed using the Allen Cell and Structure Segmenter. **C)** Single Cell Feature Extraction: A total of
180 216,062 single cells were segmented from the FOVs. Metadata related to the sample, experiment,
181 and microscopy was collected and associated with each individual cell. Appropriate features were
182 extracted for each cell from the cell, the nucleus or mitotic DNA, and the cellular structure
183 segmentations, including measurements such as the height and volume. These cells, including both
184 the images and the segmentations as well as the metadata and features are all available for
185 download. The hiPSC Single-Cell Image Dataset includes 25 cell lines representing key organelles
186 and cellular structures located throughout all of the major compartments of the cell. One
187 representative cell example per structure is shown as a 3D visualization in the 5x5 grid. Scale bars
188 are 10 μm unless otherwise noted.
189

190 **A PCA-based cell and nuclear shape space reveals interpretable modes of shape variation**

191 To embrace the great diversity of these 203,737 3D images of cells spanning 25 cellular
192 structures and directly compare cellular organization across this large population, we built a cell and
193 nuclear shape-based coordinate system (**Figure 2**), adapting a standard Principal Component
194 Analysis (PCA)-based dimensional reduction approach (Pincus and Theriot, 2007). First, we aligned
195 all cells to their centroids, preserving both, the biologically relevant, apical-basal axis (z-axis in lab
196 frame of reference) and the longest axis of the cell perpendicular to that axis (longest axis in x-y

197 plane). Next, we used a spherical harmonic expansion (SHE, (Marshall et al., 1996; Ruan and
198 Murphy, 2019)) to accurately parameterize each 3D cell and nuclear shape with a set of orthogonal
199 periodic basis set functions, defined on the surface of a sphere. For each cell and nuclear boundary,
200 we retained the first 16 degrees of the SHE, corresponding to 289 coefficients for each shape. The
201 joint vector of 578 SHE coefficients for each cell was sufficient for accurate reconstruction of the
202 original cell and nuclear shape with high spatial precision (**Figure 2A&B**). The joint vectors for all cells
203 (578 SHE coefficients) were then subjected to PCA. We found that the first eight principal components
204 represented about 70% of the total variation in cell and nuclear shape (**Figure 2C**). With this
205 dimensionality reduction, the cell and nuclear shapes for each individual cell can be approximately
206 reconstructed from a small vector with only eight components. This dimensionality reduction also
207 organizes the cells into a simple, intuitive 8-dimensional generative “shape space”. For example, the
208 origin (0,0,0,0,0,0,0,0) of the shape space can be reconstructed via the values of the SHE coefficients
209 representing this location in the 8-dimensional coordinate system, and then be visualized as an
210 idealized cell shape that statistically represents the average, or mean, shape of all of the cells in the
211 data set (**Figure 2D**). Similarly, idealized shapes can be reconstructed by traversing across each of
212 the eight orthogonal axes in the shape space.

213 To build a human-interpretable understanding of the modes of shape variation in our
214 population, we reconstructed cell and nuclear shapes at regular intervals separated by 0.5 standard
215 deviation units along every axis of this shape space (**Figure 2E, Movie S1**). These idealized cells
216 represent “map points” within the shape space, that can be used to identify and cluster individual real
217 cells that are similar in shape to each idealized map point and to each other. Intuitively, these
218 mathematically orthogonal modes of shape variation appear to describe expected variable cell shape
219 features that are independent of one another. For example, Shape Mode 1, the mode representing
220 the greatest amount of shape variation, appeared to largely reflect the height of the cell, and Shape
221 Mode 2 appeared to largely reflect the overall volume of the cell (**Figure S2A**). The fact that these two
222 biologically meaningful modes of shape variation correspond to the two top modes identified by the
223 PCA indicates that, within this dataset, the total height of the cell is largely independent of its overall
224 volume. Indeed, for the hiPSCs grown in self-organized colonies, the colony size and cell position
225 within the colony appear to be the primary determinants of cell height (see “Statistical Analysis of ...”
226 section in Methods). The remaining Shape Modes 3 to 8 represented other systematic ways the
227 shapes of these epithelial-like cells might be expected to vary, including tilting along the major or
228 minor xy-axes (Shape Modes 3 and 4) or elongation along the major axis (Shape Mode 5). In Shape
229 Modes 1, 2, and 5, nuclear shape changed concomitant with cell shape, while in the other shape
230 modes, nuclear shape changed very little as the shape mode axis was traversed. Instead, for these
231 modes it was the position and orientation of the nucleus within the cell that adjusted concomitant with
232 cell shape (**Figure 2E and MovieS1**). For completeness, we also independently calculated shape



234 **Figure 2.** A Principal component analysis (PCA)-based cell and nuclear shape space reveals
235 interpretable modes of shape variation in hiPSCs. **A)** Segmented 3D single-cell images of a cell and
236 its nucleus are used as the input for a 2D alignment algorithm. The cell image is rotated in the xy-
237 plane by θ_{cell} degrees around its centroid such that its longest axis becomes parallel to the x-axis. The
238 same rotation angle θ_{cell} is applied to the segmented nuclear image. The resulting aligned images of
239 the cell and nucleus are used as the input for spherical harmonics expansion (SHE) of degree $L_{\text{max}} =$
240 16 resulting in a total of 289 SHE coefficients for each the cell and the nucleus. These 578
241 coefficients, together, now can be used to reconstruct the cell and nuclear shape as two separate 3D
242 meshes with high fidelity. After reconstruction, the nuclear mesh is translated to the correct position
243 relative to the cell centroid. x_{lab} , y_{lab} and z_{lab} denote the lab frame of reference and x_{cell} and y_{cell}
244 represent the x and y coordinates in the rotated cell frame of reference. **B)** Mean distance between
245 points in the original meshes of cell (left) and nucleus (right) to their corresponding closest points in
246 the reconstructed meshes as the number of coefficients in the SHE increases. Each gray line is one
247 cell (left; n=300 randomly selected samples) or nucleus (right; n=300 randomly selected samples).
248 Black lines represent the mean. The dashed vertical lines indicate the number of coefficients for SHE
249 degree $L_{\text{max}} = 16$. **C)** SHE coefficients were calculated for all of the n=203,737 cells and nuclei in the
250 analysis dataset. PCA was used to reduce the dimensionality from 2x289 SHE coefficients into the
251 first eight principal components (PCs). **D)** Each PC was normalized into units of standard deviation,
252 generating eight shape modes, which together are referred to as the cell and nuclear shape space.
253 Each shape mode is sampled at nine map points. These map points are located at -2σ to 2σ in steps
254 of 0.5σ ($\sigma =$ standard deviation). Only one shape mode is permitted to vary at a time. The top
255 example shows the mean cell and nuclear shape, represented by the map point (0,0,0,0,0,0,0).
256 These nine map points for each of the eight shape modes are used as the input for an inverse PCA
257 transform to obtain the corresponding SHE coefficients and their corresponding 3D reconstructions at
258 these map points. The top view corresponds to an intersection between the 3D mesh of the cell and
259 nucleus reconstructions and the xy-plane. Side views 1 and 2 correspond to an intersection between
260 the 3D meshes and the xz- or yz-planes, respectively. **E)** 2D projections of 3D meshes obtained for
261 each of the nine map point bins of each of the eight shape modes. The center bin in all modes is the
262 identical mean cell shape. The most relevant of the three possible views is shown for each mode, as
263 indicated on the far left. All three views for each shape mode and map point can be seen in **Movie S1**.
264 Human-interpretable names for these shape modes are indicated on the right. **F)** Overlay of mesh
265 projections of the cell (magenta) and nucleus (cyan) for the two most extremes map points (at -2σ ,
266 lighter shade and $+2\sigma$, darker shade) of each shape mode.
267

268 spaces using only the overall cell shape and only the nuclear shape, which generally showed similar,
269 biologically meaningful, modes of variation (**Figure S2B&C**).

270

271 **Building integrated average cells throughout the shape space via SHE coefficient-based** 272 **parameterization and 3D morphing**

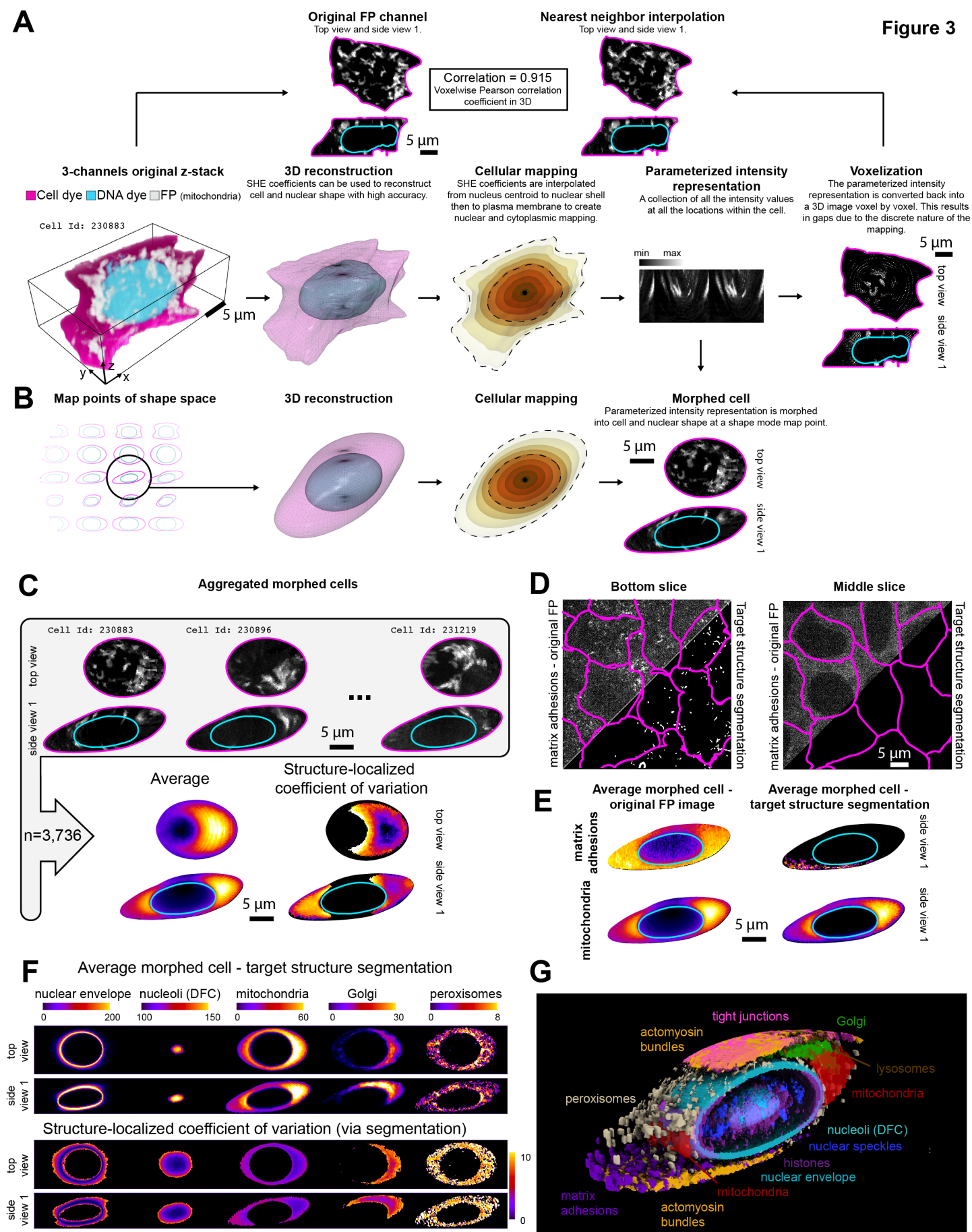
273 The human-interpretable understanding of each of the shape modes in this cell and nuclear
274 shape space now permits us to take advantage of the variation of cell and nuclear shape within this
275 dataset in two significant ways. First, this standardized shape space permits clustering of similarly
276 shaped cells, facilitating investigation of the location of cellular structures while keeping any chosen
277 3D spatial constraint constant. For example, we can measure how variable the locations of
278 mitochondria are within cells of similar height. Second, the ability to analyze cellular structure location
279 throughout this shape space permits us to consider each cell as its own “experiment” in intracellular
280 organization, representing a particular point in the overall cell shape space comprising this normal

281 population. Thus, we can ask how robust the location of a cellular structure may be when it is
282 subjected to systematic variation in cell and nuclear shape. For example, we can compare differences
283 in structure locations or their variations between flat and tall cells, small and large cells, or cells with
284 shapes that are less or more polarized. In brief, this shape space creates an opportunity to investigate
285 how the rules of cellular organization change in response to a set of naturally occurring shape
286 perturbations compared to the “mean” cell shape.

287 To directly and quantitatively compare similarly shaped individual cells and their contents near a
288 particular part of the shape space, we needed to map all of the possible locations of the contents of
289 these cells into one identically bounded cell shape. Therefore, we developed a method to “morph” all
290 of the locations of all of the points within a cell into the idealized reconstructed cell shape that best
291 represents that cell shape (**Figure 3**). We took advantage of the SH expansion describing the outer
292 cell boundary and the outer nuclear boundary and interpolated between the relevant SH coefficients.
293 This generates a “*cytoplasmic mapping*” of successive 3D concentric shells between the nuclear and
294 cellular boundaries at a specified spacing. Similarly, we generated a “*nuclear mapping*” from the
295 nuclear centroid to the nuclear boundary. We then created a “*parameterized intensity representation*”
296 of all of the intensity values at all of these mapped locations within the cell. This parameterized
297 intensity representation can then be transformed back into any cell or nuclear shape (**Figure 3A**). As
298 a proof of concept, we first performed this internal mapping and transformation of all of the fluorescent
299 signal within an individual cell back into that cell’s own original shape, permitting us to measure how
300 well spatial information is conserved using this approach. Since this internal mapping is discrete, the
301 resultant reconstructed intensity image will have gaps, which were filled using nearest neighbor
302 interpolation. We then calculated the voxel-wise Pearson correlation of the original and recreated
303 images of the same cell in 3D for individual cells representing each of the 25 cellular structures. We
304 found that for most structures this correlation was very high, above $r=0.8$ (**Figure S3A**). Only those
305 structures that localized to separate discrete spots displayed slight reductions in these correlation
306 values, likely due to the discrete nature of both the parameterized intensity representation and the
307 structures themselves.

308 We next applied this approach to morph the parameterized intensity representation of each
309 cell into the idealized cell and nuclear shape representing a nearby map point location in the shape
310 space, creating a ‘morphed cell’ (**Figure 3B**). Now we can choose any map point within the shape
311 space and identify a cluster of individual cells around that point, then create morphed versions of
312 these cells and their contents to fit within the exact same idealized shape. In this way the location of
313 the contents of each of these cells could be directly and quantitatively compared. The set of morphed
314 cells within a chosen region in the shape space could also be aggregated via their parameterized
315 intensity representations to generate an average of all of the intensities mapped within the cell and
316 nuclear shape, or to quantify the variation in intensities via the coefficient of variation (**Figure 3C**).

Figure 3



318 **Figure 3.** Building integrated average cells throughout the shape space via SHE coefficient-based
319 parameterization and 3D morphing. **A)** Bottom left image, labeled as *3-channel original z-stack*,
320 shows a 3D visualization of the original fluorescent protein (FP) intensities of tagged mitochondria (via
321 Tom20, grayscale) in a single cell and nucleus, visualized via cell membrane dye (magenta) and DNA
322 dye (cyan). Moving rightward along the bottom row are the steps to create the parametric intensity
323 representation of the mitochondria via the FP signal in this cell. The second image, labeled *3D*
324 *reconstruction*, shows the SHE-based 3D reconstruction meshes of the segmentations of this cell and
325 nucleus. The third image, labeled *cellular mapping*, shows the result of interpolating the SHE
326 coefficients to create a series of concentric mesh shells (indicated by different colors) from the
327 centroid of the nucleus (black dot) to the nuclear boundary (inner dashed contour) to create the
328 nuclear mapping and from that nuclear boundary to the cell boundary (outer dashed contour) to create
329 the cytoplasmic mapping. The intensity values in the FP channel are recorded at each mesh vertex
330 location, resulting in a *parameterized intensity representation* that is shown in a matrix format in the
331 fourth image. This parameterized intensity representation is then converted back into a 3D image,
332 voxel by voxel, into the same reconstructed cell and nuclear shape, shown in the fifth image, labeled
333 *voxelization*. Here the top view and side view 1 are shown with the intensity image in the FP channel
334 in grayscale and the cell and nuclear boundaries in magenta and cyan lines, respectively. The top left
335 image, labeled *original FP channel*, is the top view and side view of the same cell as in the 3-channel
336 original z-stack panel on the bottom row. The intensity image in the FP channel, in this case
337 mitochondria, is shown along with the cell and nuclear segmentations (magenta and cyan lines,
338 respectively). The top right image, labeled *nearest neighbor interpolation*, is the voxelized
339 parameterized intensity representation, now with gaps filled using nearest neighbor interpolation.
340 Voxel-wise Pearson correlation in 3D is used to compare the input image (original FP channel) with
341 the image reconstructed via the parametric intensity representation (nearest neighbor interpolation).
342 **B)** The same 3D reconstruction and cellular mapping procedure is now applied to a cell and nuclear
343 shape at any map point in the shape space, shown here to the Shape Mode 3 map point
344 $(0,0,1.5\sigma,0,0,0,0)$. In the fourth image, labeled *morphed cell*, the parameterized intensity
345 representation of the FP channel is morphed into this shape-space based cell and nuclear shape
346 creating a morphed cell with morphed structure location. **C)** Top panel shows images of three different
347 example cells with tagged mitochondria, each located near the Shape Mode 3 map point
348 $(0,0,1.5\sigma,0,0,0,0)$, morphed into the reconstructed cell and nuclear shape of that map point. These
349 three and all other morphed cells with tagged mitochondria within that map point bin in the shape
350 space can be aggregated voxel by voxel to create an average morphed cell representing the average
351 mitochondria locations (image labeled *average*) in that part of the shape space. Morphed cells can
352 also be aggregated by calculating the standard deviation at each voxel of the morphed cell shape
353 (**Figure S3**). The average and standard deviation morphed cells can be combined to calculate the
354 *structure-localized coefficient of variation*, representing a quantitative measure of how variable the
355 location of a structure is at any given voxel. **D)** FOV images of multiple cells (cell membrane indicated
356 by magenta lines) with labeled matrix adhesions (via paxillin) at two z positions in the z-stack. Top left
357 triangles in each image show the original FP image. Matrix adhesions are visible near the bottom of
358 the cells (left) but considerable FP-tagged paxillin signal is visible both at the bottom and center (right)
359 of cells. Bottom right triangles in each image show the result of the matrix adhesion specific
360 segmentation, in which only the matrix adhesions themselves remain near the bottom of the cells as
361 the target of the segmentation. **E)** Average morphed cells at the Shape Mode 3 map point
362 $(0,0,1.5\sigma,0,0,0,0)$, representing matrix adhesions (top row) and mitochondria (bottom row)
363 generated using either the original FP images (left column) or the target structure segmentations (right
364 column). **F)** Top view and side view 1 of average and structure-localized coefficient of variation
365 morphed cells at the Shape Mode 3 map point $(0,0,1.5\sigma,0,0,0,0)$, based on the target structure
366 segmentations, representing five distinct cellular structures. See **Figure S3** for examples for all 25
367 structures and **DataFile S1** for numbers of cells aggregated at each shape space bin. **G)** Eleven
368 structures are rendered simultaneously to illustrate their relative spatial relationships in this 3D
369 visualization (actomyosin bundles and mitochondria are labeled twice highlight their dual locations).
370 The average morphed cells representing each of these structures at the Shape Mode 3 map point

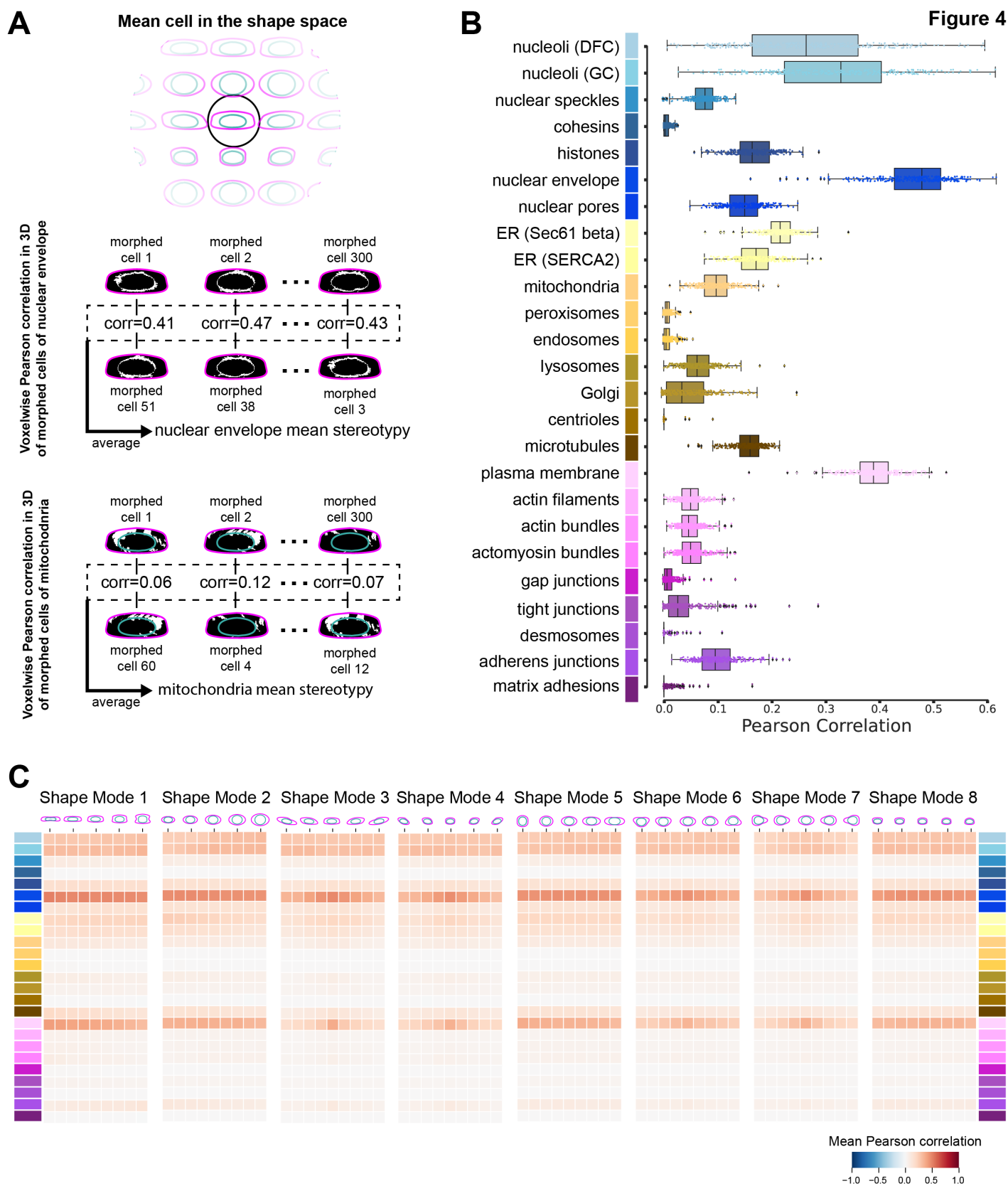
371 (0,0,1.5 σ ,0,0,0,0), based on the target structure segmentations were combined in this image. For
372 each of these, the average structure image was segmented using the default *Surface* option found in
373 the *Volume Viewer* window of ChimeraX. Thresholds for each channel were selected arbitrarily to
374 clarify dominant localization patterns observed in the voxel intensities. See **Movie S2** for rotating
375 image.
376

377 This parametric intensity representation takes all intensities in the image into account, including any
378 FP-tagged protein not localized to the target structure that the protein represents. For example,
379 EGFP-tagged paxillin localized to matrix adhesions at the bottom of the cell but also throughout the
380 cytoplasm. However, the segmentation target defined for this cell line included only the high intensity
381 regions representing the matrix adhesions (**Figure 3D**). Applying this same cell morphing approach to
382 the segmented versions of the cellular structure images (rather than to the FP images directly) permits
383 the creation of average morphed cells containing the locations of the cellular structures that each
384 tagged protein represents (**Figure 3E**). Our remaining analyses in this paper focus on the segmented
385 structure images; but conceptually the same approach could also be applied to the raw intensity
386 images.

387 We clustered all cells in the dataset arbitrarily into nine bins along each of the eight shape
388 modes at the standard deviation intervals shown in **Figure 2 (Figure S3C)** to create a total of 65 cell
389 shape map points (the center bin is the same in all modes), into which we morphed each of the 25
390 structures (**Figure 3F and Figure S3D**). By direct visual inspection, we found that the average
391 morphed cells accurately represented the location patterns of these structures in individual,
392 unmorphed cells (**Figure S3D**). We could then combine the average morphed locations of each of the
393 25 structures into the same cell shape (11 integrated structures visualized in **Figure 3G and Movie**
394 **S2**), creating integrated average morphed cells, which we did for each of the 65 map point cell
395 shapes.
396

397 **The location stereotypy of cellular structures depends on the structure but not the cell shape**

398 To measure how variable the location of each individual cellular structure is within the cell, we
399 used individual morphed cell images based on the structure segmentations for each structure at each
400 map point bin in the cell and nuclear shape space. We calculated the 3D voxel-wise Pearson
401 correlation between pairs of individual morphed cell images within a shape bin for each of the 25
402 cellular structures (**Figure 4A**) and averaged those correlation values to generate a measure of the
403 “*location stereotypy*” of each structure (**Figure 4B**). Structures with a high stereotypy value have little
404 cell-to-cell variation in their overall absolute positions for similarly shaped cells, while structures with a
405 low stereotypy value may be found in distinct locations even for two cells whose shapes are very
406 similar. Comparing the average stereotypy for each structure permitted us to rank structures that are
407 most to least stereotyped in their location within the mean cell and nuclear shape. The most



409 **Figure 4.** The location stereotypy of cellular structures depends on the structure but not the cell
410 shape. **A)** Overview of the process to calculate the stereotypy of cellular structures within the mean
411 cell shape, using the nuclear envelope (via lamin B1) and mitochondria (via Tom20) as examples.
412 Segmented images of each cellular structure within 300 cells located in the mean cell bin in the shape
413 space were each morphed into the mean cell shape, creating, for example, 300 nuclear envelope and
414 300 mitochondria morphed cells. The voxel-wise Pearson correlation was calculated for 300 unique
415 pairs of morphed cells of same cellular structure and the results were organized as a correlation list.
416 The mean value of the correlation list was defined as the mean location stereotypy for that structure.
417 **B)** Box plots corresponding the values in the correlation list (see panel A) for each of the 25 cellular
418 structures, represented by unique colors to the left of the y-axis. Dots represent the raw data (n=300),
419 vertical black lines represent first and third quartile, boxes represent the interquartile range and the
420 vertical black line inside the box is the mean. **C)** The process described in panel A was performed for
421 cells in each of the 72 shape space bins to calculate the average location stereotypy for all 25 cellular
422 structures throughout the shape space. Each heatmap value corresponds to the mean stereotypy of
423 all 25 cellular structures for a given shape mode. Each row in the heatmap represents a different
424 cellular structure, indicated by the same colors as in panel B. Columns in the heatmaps represent the
425 nine binned map points along each shape mode (see **Figure S3C**). N = 300 morphed cells for each
426 cellular structure and shape mode bin or the maximum number of cells available (see **DataFile S1**
427 and Methods).

429 stereotyped structures were the nuclear envelope (lamin B1) and the plasma membrane (CAAX
430 domain of K-Ras, “CAAX”). These observations are effectively positive controls, because these two
431 structures should be very similar to the cell and nuclear boundary shapes that were used as fixed
432 points in the SHE interpolation. In decreasing order of stereotypy, the next highest were two nucleolar
433 compartments, the Dense Fibrillar Component (DFC, via fibrillarin and the Granular Component (GC,
434 via nucleophosmin), followed by the ER (both Sec61 beta and SERCA). Structures with the least
435 location stereotypy within the mean cell included those with a low number of discrete separated
436 locations near the top or bottom of the cell such as centrioles (via centrin-2), desmosomes
437 (desmoplakin), and matrix adhesions (paxillin). They were followed by slightly increased stereotypy for
438 cohesins (SMC-1A), endosomes (Rab-5A) and peroxisomes (PMP34). To control for any effects of
439 variable spacing within locations of each of the 25 cellular structures, we performed a systematic
440 downsampling of the voxel size (**Figure S4A**) and found little change to the stereotypy order of the
441 structures (**Figure S4B**).

442 We next investigated how much the location stereotypy changed in response to the set of
443 naturally occurring cell shape perturbations represented by the systematic changes in cell shape
444 along each of the eight shape modes compared to the “mean” cell shape. Strikingly, we found very
445 little change in the magnitude or rank order of the location stereotypy throughout the entire shape
446 space, demonstrating that the stereotypy of all of these 25 structures was extremely robust to overall
447 cell shape variation (**Figure 4C** and **Figure S4C&D**).

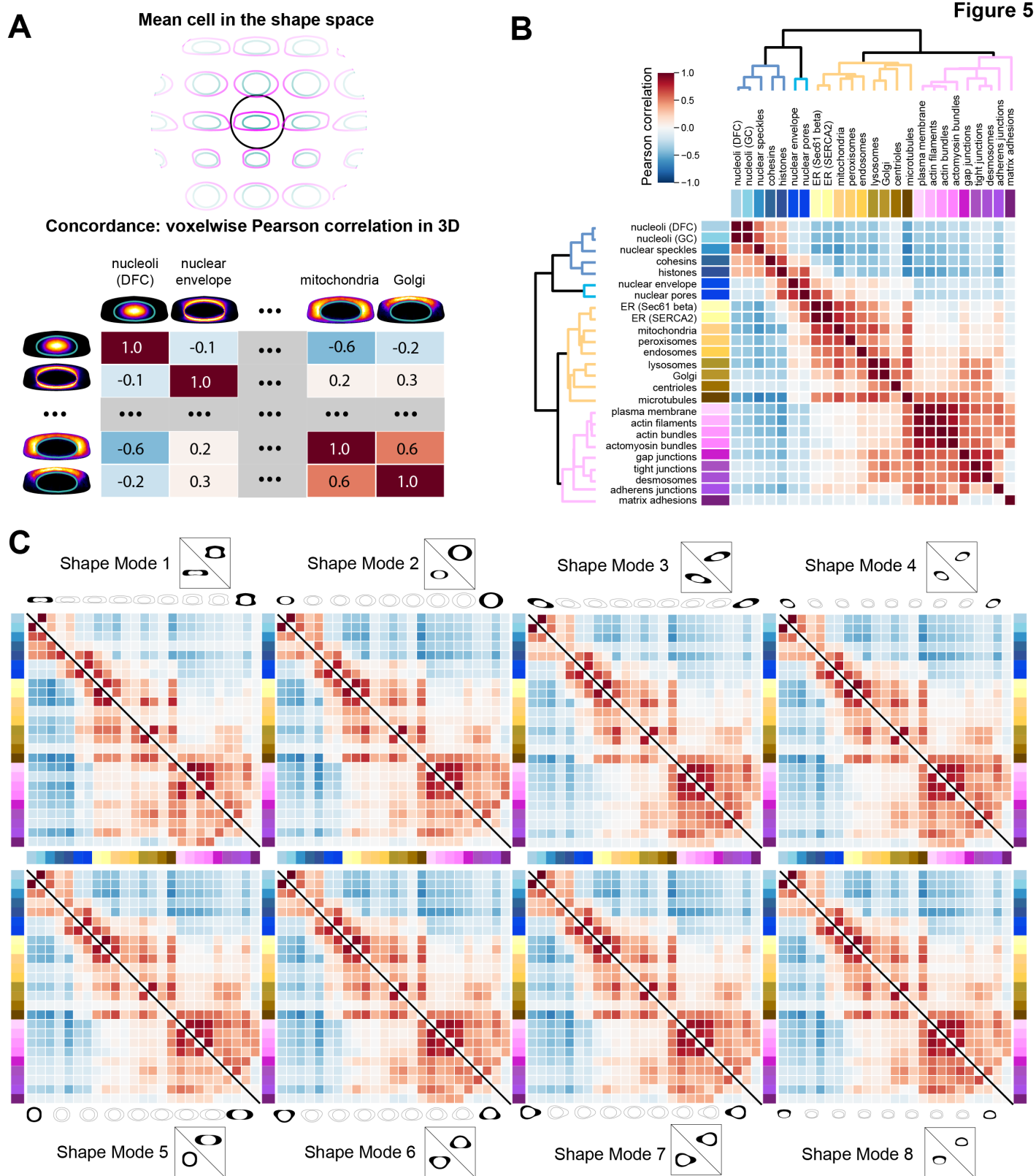
448
449 **The location concordance of all 25 cellular structures to each other suggests a robust, ordered**
450 **compartmentalization of the cell**

451 The analysis described above enables quantitative ranking of the cell-to-cell variation in
452 localization of each tagged cellular structure relative to that same structure in a different cell. Also of
453 interest is the relative similarity of the absolute localizations of each tagged structure as compared
454 with every other structure. To measure the relationships of the locations of each of the 25 cellular
455 structures relative to all the others, we calculated the 3D voxel-wise Pearson correlation between the
456 average morphed cell images for all pairs of structures (pairwise structure location “concordance”)
457 within the mean cell shape (**Figure 5A**). We then performed a hierarchical clustering analysis of the
458 concordance values. This clustering is purely data-driven based on the images alone. Importantly, we
459 found that the location concordance of these cellular structures clustered naturally into an ordered
460 compartmentalization of the cell, from the center of the nucleus outward (**Figure 5B** and colors in
461 **Table 1**). The four top-level clusters included structures localized to the nucleus, nuclear periphery,
462 cytoplasm, and cell periphery, respectively. A priori, we expected to find strong concordance for
463 several cellular structure sets, including the two nucleolar structures (DFC and GC), the two structures
464 at the nuclear periphery (nuclear envelope and nuclear pores), the two ER tags (Sec61 beta and
465 SERCA), and the three structures with primary localization to the apical cell-cell contacts (gap
466 junctions (connexin-43), tight junctions (ZO-1), and desmosomes (desmoplakin)). The concordance
467 hierarchy confirmed the expected strong concordances within each of these sets, validating this
468 analysis approach.

469 We also identified several other notably high relative concordances such as the tight concordance
470 between lysosomes (LAMP-1) and Golgi (sialyltransferase 1), consistent with their enrichment in
471 location in the cytoplasm near the top of the cells and the known role of Golgi in regulating lysosome
472 localization (Hao et al., 2018; Wang and Hong, 2002). Mitochondria (Tom20) and peroxisomes
473 (PMP34) were more tightly concordant with each other than either structure was with endosomes
474 (Rab-5A), even though direct visual examination of individual peroxisome-tagged and endosome-
475 tagged cells did not easily highlight this distinction. However, this observation is consistent with the
476 known association between mitochondria and peroxisomes (reviewed in (Fransen et al., 2017)).

477 Next, we investigated how much the concordance between all pairs of the 25 cellular
478 structures changed in response to changes in cell shape, as described above for the stereotypy
479 analysis. Overall, we found very little change in the hierarchical compartmentalization of these 25
480 structures throughout the shape space (**Figure 5C**, **Figure S5**). Some structures showed an overall
481 decrease in the magnitude of concordance with other structures in the shape mode bins furthest from
482 the mean. These structures also had greatly decreased numbers of cells available in these bins for
483 this calculation, for example actin filaments (beta-actin) or cohesins (SMC-1A) in the furthest bins of
484 Shape Mode 1, and so these decreases may not be biologically meaningful (**DataFile S1**).

485
486



488 **Figure 5.** The location concordance of all 25 cellular structures to each other suggests a robust,
489 ordered compartmentalization of the cell. **A)** Overview of the process to calculate the location
490 concordance between all pairwise-combinations of the 25 cellular structures within the mean cell
491 shape. The voxel-wise Pearson correlation was calculated between pairs of average morphed cells,
492 based on the structure segmentations, morphed into the mean cell shape. This created a correlation
493 matrix including each of the 25 cellular structures with elements of this matrix representing the
494 location concordance between two cellular structures. **B)** The heatmap representing the location
495 concordance for every pair of 25 cellular structures in the mean cell shape. Each heatmap value
496 corresponds to the Pearson correlation value between the two indicated structures. The correlation
497 matrix is used as input for a clustering algorithm to produce the dendrogram shown alongside the
498 heatmap. Dendrogram branches are color coded according to major cell compartments (nucleus in
499 blue, nuclear periphery in cyan, cytoplasm in orange and the cell periphery in magenta). Lengths of
500 dendrogram branches represent the distance between clusters. **C)** The process to create the
501 concordance heatmap for the mean cell shape in B was repeated for the reconstructed cell and
502 nuclear shapes at the -2σ and 2σ shape space map points for each of the eight shape modes. Each
503 heatmap represents one shape mode. The lower triangle represents shape space map point -2σ and
504 the upper triangle represents shape space map point 2σ . For sake of clarity, diagonals are colored in
505 white and black lines are used to separate the lower and upper triangles. The number of cells
506 analyzed for each cellular structure and shape mode bin can be found in **DataFileS1**.
507

508 **The impact of cell and nuclear size on the variation in cellular structure size is structure-** 509 **dependent**

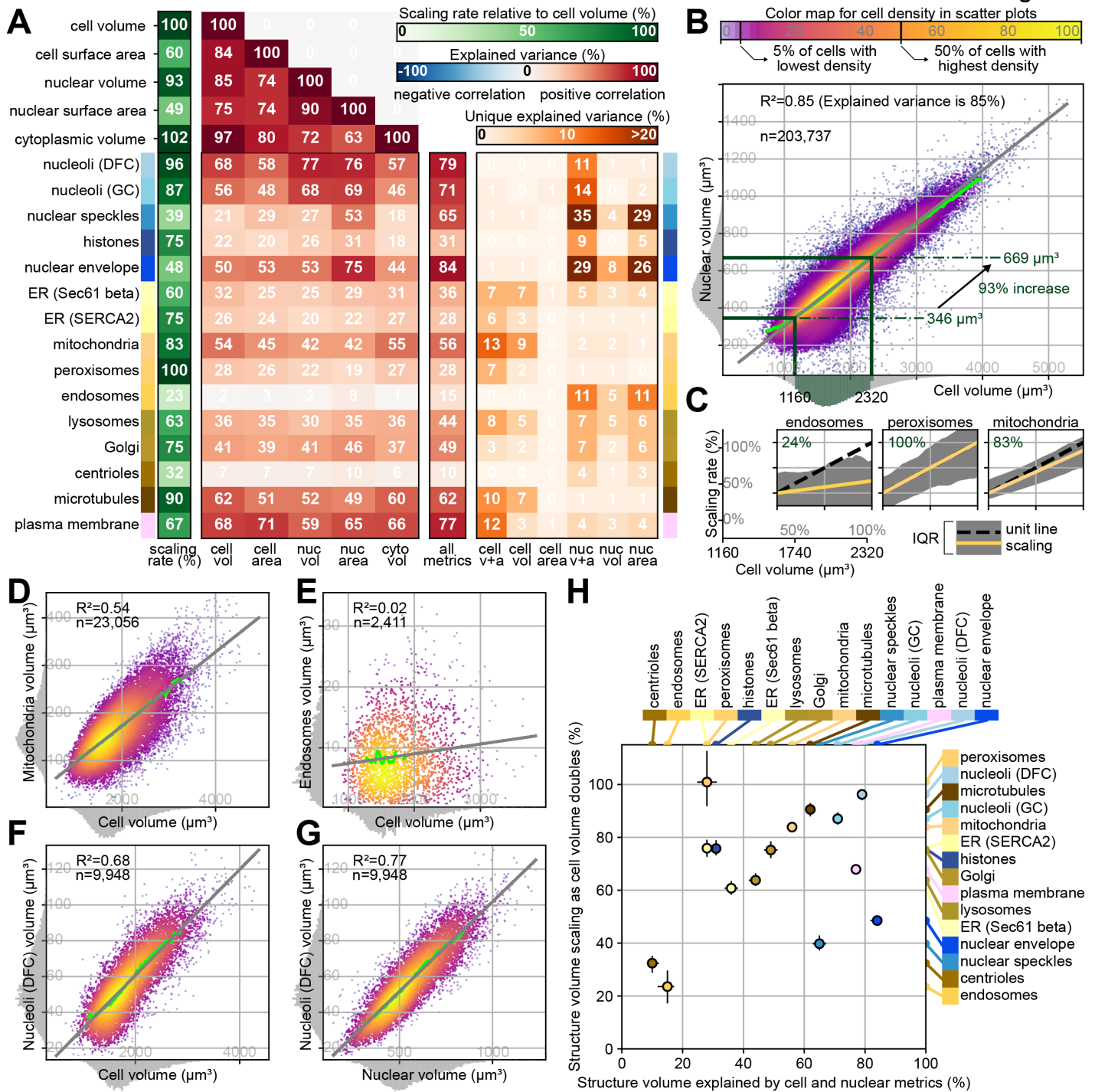
510 Intracellular structures exhibit cell-to-cell variation not only in their locations but also in how
511 much of a given structure is present in the cell. So far, we have found that neither the variability in
512 each location of each structure in the cell nor their relative locations to each other changed much with
513 cell volume (Shape Mode 2; **Figure 5C**). However, it has previously been shown that the volume of
514 several cellular structures in the cell does correlate with overall cell volume, including the nucleus and
515 mitochondria (reviewed in (Marshall, 2020)). We therefore used our large dataset to perform a
516 systematic and comparative analysis of the relationship between cellular structure volume and five
517 relevant size metrics (cell volume, cell surface area, nuclear volume, nuclear surface area, and
518 cytoplasmic volume) for the 15 cellular structures in this dataset validated for structural volume
519 analysis (**Figure 6** and **Figure S1**). We used simple linear regression to fit the data and calculated the
520 percent of the variation in cellular structure volumes that can be explained by each of the four cell and
521 nuclear size metrics (“percent explained variance”; **Figure 6A**). The rolling average, a non-linear
522 model fit, and analyses in which we considered the geometrical relationship between the volume and
523 surface area of the roundest nuclei all showed similar results, validating the simple linear regression
524 approach (**Figure 6B**, **Figure S6A&B**, and Methods). We found that the percent explained variance
525 attributable to these overall cell size metrics was substantially greater for some structures, such as
526 mitochondria (Tom20; 54%) than for other structures, such as endosomes (Rab-5A; 2%, **Figure**
527 **6D&E**). We also found that for nuclear structures like the nucleolar DFC (fibrillarin), more of the
528 variance in their volumes could be explained by nuclear volume than by cell volume (77% vs. 68%,
529 respectively; **Figure 6F&G**).

530 Each of the five cell and nuclear size metrics themselves, of course, also correlate with each
531 other (**Figure 6A**), thus obscuring their potential to independently explain the variance in the volumes
532 of the 15 cellular structures. To disentangle these correlations, we applied a multivariate model and
533 calculated the total percentage of the variance explained for each of these structures by the
534 combination of all four direct cell and nuclear size metrics (“total explained variance”; see “all metrics”
535 column in **Figure 6A**). For all but two of the cellular structures, the total explained variance was at
536 least 28%; but this percentage varied widely depending on the structure (x-axis in **Figure 6H**). At the
537 lowest end were the centrioles (centrin-2), which we expected to be very small as they are discrete
538 structures that should not get bigger as cells grow, and thus invariant with all size metrics. At the
539 highest end were the nuclear envelope (laminB1) and the plasma membrane (CAAX), which we
540 expected would correlate well with nuclear and cell surface areas, respectively. Notably, the volumes
541 of all three nuclear body structures (nucleolar DFC, GC, and speckles) were the next-most tightly
542 correlated to the optimal linear combination of cell and nuclear size metrics.
543 We then used the multivariate model to calculate the *unique* contributions of both cell size metrics vs.
544 both nuclear size metrics vs. the unique contributions of each of the four metrics individually (**Figure**
545 **6A**). For all five nucleus-related structures, the variance in structure volume was better explained by
546 nuclear size metrics than by cellular size metrics. For the nuclear envelope, more of the variance was
547 uniquely attributable to the nuclear surface area than nuclear volume; this anticipated result confirmed
548 the validity of this approach. Unexpectedly, the variance in nuclear speckle (SON) volumes was most
549 uniquely attributable to the nuclear surface area and not the nuclear volume, although speckles
550 localize throughout the nucleoplasm.

551 Of the cytoplasmic structures, microtubules (alpha-tubulin, see **Figure S1** for target
552 segmentation of microtubule bundles), which localize throughout the cytoplasm (**Figure S5D**), had the
553 highest percent variance explained by the optimal combination of the four size metrics, followed next
554 by mitochondria, Golgi and lysosomes (x-axis in **Figure 6H**). Endosomes (Rab-5A) had one of the
555 lowest percent explained variance values, almost as low as centrioles, even though they are spread
556 out throughout the cytoplasm. For the cytoplasmic structures, some variation in their structure
557 volumes was uniquely attributable to either cell or nuclear metrics; but in all cases the unique
558 contribution of cell surface area on its own was negligible. While nuclear structures seem to be most
559 tightly coupled to nuclear size metrics, cellular structures range more widely in how well the variance
560 in their volumes was uniquely attributable to cell versus nuclear size metrics. We explored whether
561 cell and nuclear shape might explain some of the variation in cellular structure volumes but found
562 contributions from other shape modes to be negligible (**Figure S6C**). Overall, these results show that
563 how well cell and nuclear size metrics account for the variation in cytoplasmic structure volumes is
564 structure-dependent, consistent with the wide range of cell functions that these structures regulate.

565

Figure 6



567 **Figure 6.** The impact of cell and nuclear size on the variation in cellular structure size is structure-
568 dependent. **A)** Heatmap in four parts summarizing the results of a systematic, comparative analysis of
569 the relationship between the volumes of 15 cellular structures and four cell and nuclear size metrics
570 (cell and nuclear volume and surface area, also referred to as *Cell vol*, *Cell area*, *Nuc vol*, *Nuc area* in
571 the heat map columns). The number of cells for each cellular structure can be found in **Table S1**. The
572 very left of the heatmap shows the same compartmentalized cellular structure coloring scheme as in
573 other figures. Heatmap part 1: the leftmost column, labeled *scaling rate* and colored in green indicates
574 the percent scaling of cellular structure volumes relative to a doubling in cell volume. For example, a
575 value of 83% for mitochondria indicates that mitochondrial volume is increased by 83% when the cell
576 doubles in size (100%) (B and C). Heatmap part 2: the next five columns, each labeled with the four
577 cell and nuclear size metrics plus cytoplasmic volume (the difference between cell and nuclear
578 volume, labeled *Cyto vol*) are colored with the blue-red heatmap color range, labeled *explained*
579 *variance*. These columns indicate the percent of variation in each cellular structure volume (and
580 surface area for the cell and nucleus; rows) that can be statistically explained using the metrics
581 indicated in each column. Negative values would represent a negative correlation relationship
582 between the two variables (row and column), but are not present in this heatmap. These percent of
583 explained variance values are a measure of the tightness of the coupling between cellular structure
584 volume and specific cell and nuclear size metrics. Heatmap part 3: the center single column, labeled
585 *all metrics*, uses a multivariate model that includes the four cell and nuclear size metrics. The values
586 in this heat map column represent the total percent of variation in each cellular structure volume that
587 can be statistically explained using a combination of all four metrics. Heatmap part 4: The last six
588 columns colored with the orange heatmap color range, labeled *unique explained variance*, show the
589 percentage of variation in each cellular structure volume that can be uniquely attributed to a single
590 metric (each of the four cell and nuclear size metrics) or a pair of metrics (cell volume plus surface
591 area as the cell size metrics – labeled *cell v+a*, nuclear volume plus surface area as the nuclear size
592 metrics, labeled *nuc v+a*). This number is computed as the difference between the total explained
593 variance (the all metrics column) and the variance explained by a model using all four cell and nuclear
594 size metrics except for the metric (or pairs of metrics) indicated in that column. For example, the
595 second orange heat map column, labeled *Cell vol*, indicates the percentage of explained variance that
596 is lost when cell volume is removed from the multivariate model. Thus, this is the percent of explained
597 variance that can be uniquely attributed to cell volume. **B)** Scatterplot comparing cell volume (x-axis)
598 and nuclear volume (y-axis) across all cells (n=203,737). Cells are colored based on an empirical
599 density estimate. The green line is a running average. The gray line depicts a linear regression model
600 where variation in the nuclear volume (y-axis) is explained as a linear function of the cell volume (x-
601 axis). The explained variance (R^2) in nuclear volume is 85% as stated in the top-left of the plot. The
602 linear regression model is also used to calculate the scaling rate, i.e. how much larger (in %) nuclear
603 volume is when cell volume doubles. Specifically, the regression model is evaluated for the cell
604 volume interval from 1,160 to 2,320 μm^3 (where the cell volume doubles) to determine to scaling
605 percentage for nuclear volume. **C)** Line plots showing the relative volume scaling rate for three cellular
606 structures (endosomes, peroxisomes and mitochondria) over the same cell volume doubling range as
607 in **B**, from 1,160 to 2,320 μm^3 . The yellow lines represent the scaling rate, also indicated by the
608 numbers in the top left corner of each of these plots. The regions filled in gray represent the
609 interquartile range (IQR) measured across cells that were binned in 10 cell volume bins (y-axis). The
610 xy-axes to the far left are used to indicate the values of the tick marks in each of the three plots. **D-G)**
611 Similar scatterplots as in (B), correlating the volumes of mitochondria (D), endosomes (E), nucleoli
612 (DFC, F and G) with either cell or nuclear volume (x-axes) along with statistical measures. **H)**
613 Scatterplot comparing the total percent explained variance (x-axis contains the values in the all
614 metrics column of the heatmap in (A) and the relative volume scaling rate (y-axis contains the values
615 in the scaling rate column in the green heatmap in (A) across all of the 15 cellular structures. The
616 error bars depict the 5-95% confidence intervals using a bootstrap analysis. The markers along the
617 top and right side of the plot indicate the ranked order of the structures for that metric.
618 Compartmentalized cellular structure coloring scheme is used to help identify specific structures.

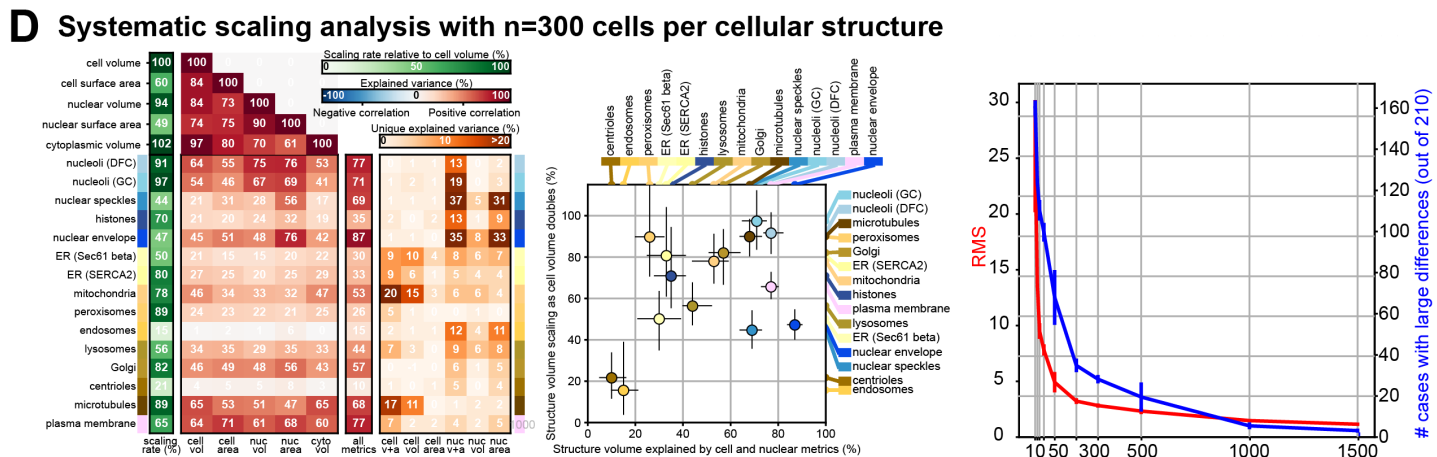
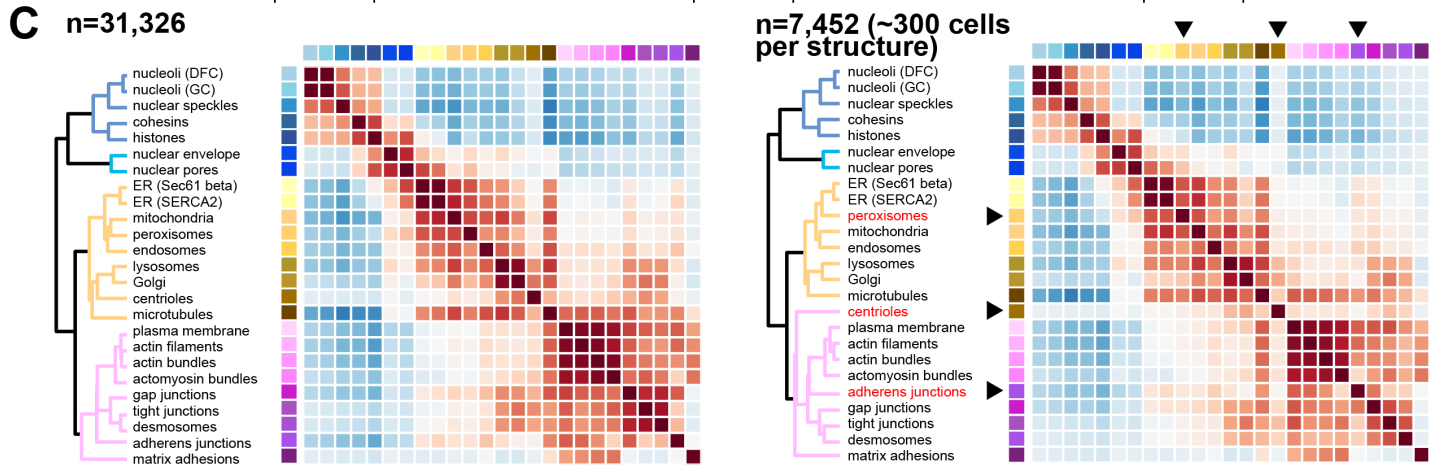
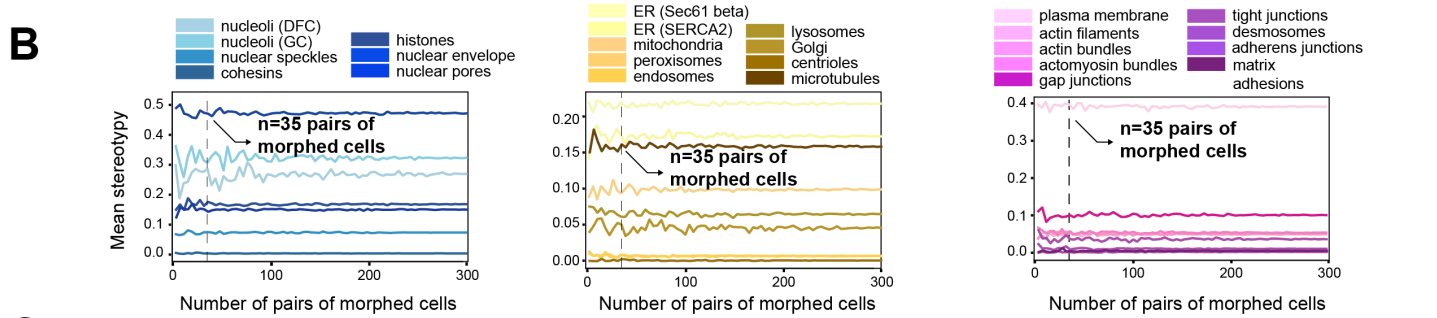
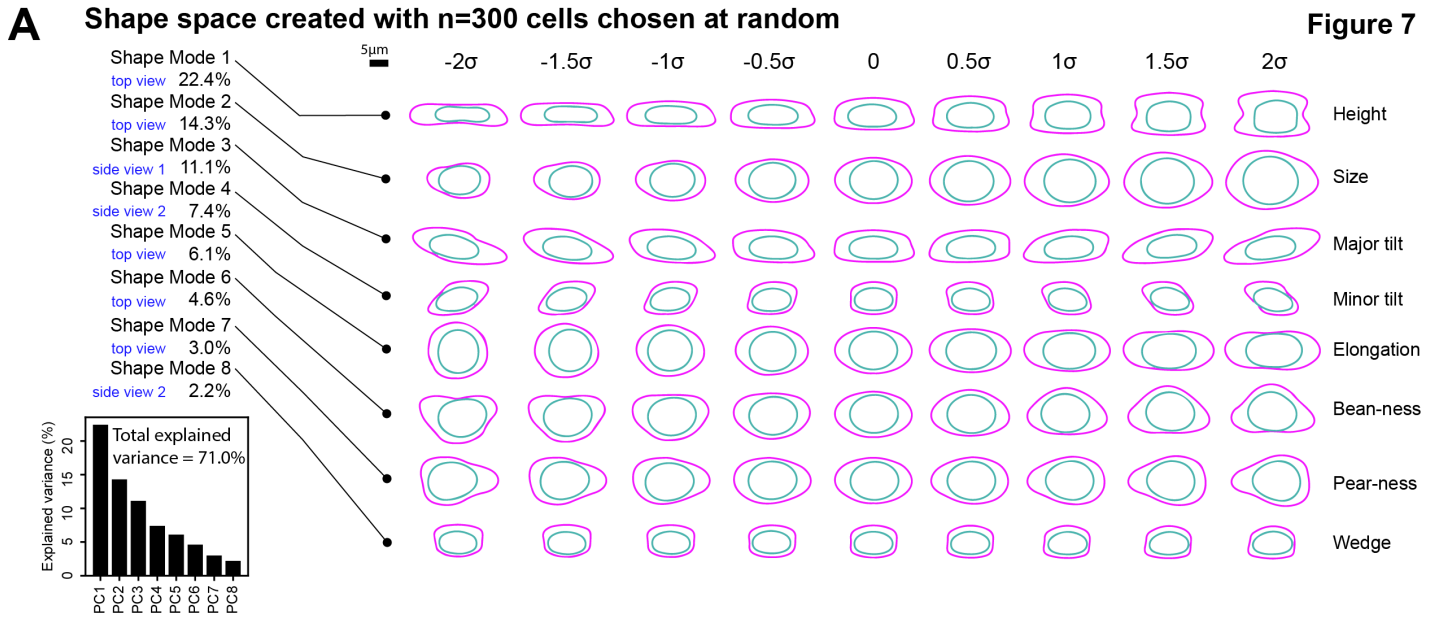
619 We also measured the relative volume “*scaling rates*” for each of these 15 cellular structures
620 as the percentage increase in structure volume given one doubling in cell volume over a range that
621 was well represented in our cell population, specifically from 1160 μm^3 to 2320 μm^3 (**Figure 6A-C**).
622 For example, the volume of mitochondria increased by an average 83% (from 108 to 198 μm^3) for this
623 doubling in cell volume (an increase of 100%). The structures with the greatest relative scaling rates
624 were the peroxisomes (via PMP34), followed closely by both nucleolar structures and then
625 microtubules (y-axis in **Figure 6H**), all of which nearly doubled in structure volume with the doubling of
626 cell volume. The structures with the lowest relative volume scaling rates were also the structures
627 identified as having the lowest explained variance, that is the endosomes and centrioles. For most
628 structures, however, we observed relative scaling rates of at least 60%, consistent with the simple
629 expectation that larger cells typically would also have larger organelles. We observed lower scaling
630 rates for the two structures whose volumes correlated most strongly with nuclear surface area, the
631 nuclear envelope and nuclear speckles. This is consistent with surface area generally scaling less
632 quickly than volume, for example, doubling the size of a perfect sphere leads to only a 59% increase
633 in its surface area. The peroxisomes stood out as exhibiting an unusual pattern of both a high relative
634 volume scaling rate and great variability in peroxisome volume from cell to cell. This systematic
635 analysis of the relationship between cellular structure volume and cell and nuclear size metrics
636 creates a rich set of quantitative constraints for modeling intracellular organization.

637

638 **Downsampling the hiPSC Single-Cell Image Dataset demonstrates generalizability of this** 639 **multi-part analysis approach**

640 The over 200,000 individual cells spanning 25 cellular structures within the hiPSC Single-Cell
641 Image Dataset permitted us to develop this multi-part, data-driven, computational approach that
642 generates a collection of quantitative rule-building constraints for the locations, amounts, and degree
643 of variation of a set of cellular structures within a population of 3D cell images. The four outputs of this
644 approach are 1) a 3D cell and nuclear shape space with human-interpretable orthogonal shape
645 modes; 2) measurements of the location stereotypy throughout the shape space; 3) measurements of
646 the location concordance between all pairs of cellular structures throughout the shape space; and 4)
647 measurements of the variation in the volumes of the cellular structure relative to cell and nuclear size.
648 With this approach we have created a fundamental benchmark for comparison of these analyses of
649 healthy, normal, undifferentiated hiPSCs with future studies of other populations of cells in different
650 cell states, including differentiated cells, or cells in pathological states generated by pharmacological
651 or genetic perturbations. However, broader adoption of this quantitative approach will depend on how
652 generalizable these analyses are to other kinds of data sets, particularly those that have a
653 substantially smaller number of cells.

654



656 **Figure 7.** Downsampling the hiPSC Single-Cell Image Dataset demonstrates generalizability of this
657 multi-part analysis approach. **A)** Cell and nuclear shape space generated with 300 randomly selected
658 cells from the analysis dataset. The figure shows 2D projections of 3D meshes obtained for each of
659 the nine map point bins of each of the eight shape modes. The center bin in all modes is the identical
660 mean cell shape. The most relevant of the three possible views is shown for each mode, as indicated
661 on the far left. Humanly interpretable names for these shape modes are indicated on the right. PCA
662 was used to reduce the dimensionality from 2x289 SHE coefficients into the first 8 PCs. Total
663 explained variance by each component is shown in the bottom left plot. **B)** Mean stereotypy as a
664 function the number of pairs of morphed cells used to compute the voxel-wise Pearson correlations.
665 Cellular structures are grouped into three clusters: nucleus and nuclear periphery, cytoplasm and cell
666 periphery. Dashed vertical lines indicate the minimum number of pairs of cells necessary to recover
667 the stereotypy ranking shown in **Figure S4**. **C)** The left heatmap represents the location concordance
668 for every pair of 25 cellular structures in the mean cell shape as in **Figure 5B** computed on all
669 $n=31,326$ cells within the mean cell shape bin. The right heatmap represents the location
670 concordance for every pair of 25 cellular structures in the mean cell shape calculated on a dataset
671 composed by 300 randomly chosen cells per structure (except $n=252$ for nuclear speckles, see
672 **DataFile S1**) for a total number of $n=7,452$ cells within this mean cell shape bin. Each heatmap value
673 corresponds to the Pearson correlation value between the two indicated structures. The correlation
674 matrix is used as input for a clustering algorithm to produce the dendrogram shown alongside the
675 heatmap. Dendrogram branches are color coded according to major cell compartments (nucleus in
676 blue, nuclear periphery in cyan, cytoplasm in orange and the cell periphery in magenta). Lengths of
677 dendrogram branches represent the distance between clusters. Dendrogram on the left and right side
678 only differs by the cluster assignment of three cellular structures, the peroxisomes, centrioles, and
679 adherens junctions, highlighted in red and marked with arrowheads. **D)** Systematic scaling analysis
680 with $n=300$ randomly chosen cells per cellular structure. The Heatmap on the left is the equivalent to
681 FigureScalingA and the scatterplot in the center is the equivalent of **Figure 6H**, but with the
682 downsampled number of cells per structures. The plot on the right shows the effect of downsampling
683 the dataset on the complete set of statistical measurements calculated and shown in the heatmap on
684 the left. Here, the original set of measurements were compared to a new set of measurements
685 calculated on a series of downsampled versions of the dataset with n cells per structure randomly
686 selected ($n=10, 20, 30, 50, 100, 200, 300, 500, 1,000, 1,500$; x-axis), repeated three times. The root-
687 mean-square (RMS) difference between the two sets of measurements is shown in red and the
688 number of cases, out of 300, where the absolute difference between the measurements was larger
689 than 5% is shown in blue. A dataset size of 300 cells lies near the inflection point of both metrics.
690

691 We therefore assessed the minimal dataset size required to maintain the scientific conclusions
692 of each of these four analyses (**Figure 7**). For the cell and nuclear shape space, seven of the first
693 eight shape modes were clearly recapitulated with just 300 randomly chosen segmented cells and
694 nuclei (**Figure 7A**). The location stereotypy was greatly invariant to sample size (**Figure 7B**); the
695 correct location stereotypy rank order of all 25 structures could be entirely recapitulated with just 35
696 pairs of cells per structure per shape space bin. Next, we calculated the cellular structure pairwise
697 concordance based on a random sampling of 300 cells per structure within the mean cell bin and
698 found that only three structures changed their location in the hierarchical clustering dendrogram
699 (**Figure 7C**). For these two analyses, sufficient cells would need to be imaged to ensure the required
700 number of cells per desired shape space bin. Similarly, we repeated the full set of statistical analysis
701 of cellular structure volume variation with a systematic reduction in numbers of cells per structure and
702 found that we could recapitulate all of the biological observations reported above with 300 cells per

703 structure (**Figure 7D**). Overall, each of the four analyses could be successfully performed based on
704 numbers of cells that could be reasonably collected by a single investigator imaging on a standard
705 laboratory microscopy over the course of a few days.

706

707 **Discussion**

708 A major goal of cell biology is to determine how a subset of expressed genes dictate cellular
709 phenotypes. To address this enormous challenge, we must develop approaches that can reduce the
710 amount and complexity of the information contained in cell behaviors. While many others are
711 approaching this using genomics and proteomics, our strategy is to approach this question from the
712 perspective of cellular organization, because it is both a key readout and driver of cell behavior. We
713 choose a dimensionally-reduced approach by focusing on the level of the major cellular structures. To
714 do this, we developed new ways to convert raw image data of cells and their structures into
715 dimensionally reduced information in a form that both summarizes the raw data and embraces the
716 vast cell-to-cell variability observed even within a population of putatively identical cells. Our work
717 toward this goal has included creating a collection of isogenic cell lines with FP-tagged cellular
718 structures (the Allen Cell Collection), building and standardizing an automated microscopy imaging
719 pipeline, creating new deep learning image processing algorithms, generating a large high-quality
720 dataset, inventing a new multi-part analysis approach that generates a collection of quantitative rule-
721 building constraints (rules), and building tools to facilitate the democratization of these results such as
722 an online 3D viewer to access the data. With this approach we determined where, how much and
723 how variable the various cellular structures are in an integrated and holistic manner. The initial
724 objective is to identify quantitative relationships that can become rules of organization, and then to
725 use these data to eventually develop and improve biological models and ultimately laws for
726 understanding and predicting cellular organization in a wide variety of biological contexts.

727 This grand plan requires high-quality image data of cells and their structures. We introduce
728 such a dataset here: the hiPSC Single Cell Image Dataset with over 200,000 live cells in 3D and
729 spanning 25 major cellular structures. The scale and quality of this dataset permitted us to create a
730 multi-part, data-driven, generalizable approach that transforms the image data into a dimensionally
731 reduced and human-interpretable set of quantitative rule-building constraints for the locations,
732 amounts, and degree of variation of each of these 25 cellular structures. These constraints comprise a
733 quantitative benchmark for comparison to other cell types and states, e.g., those observed during
734 hiPSC differentiation as they transition from their epithelial-like state into more mesenchymal-like
735 cells, with consequent differences in intracellular organization. Second, these constraints are
736 inherently quantitative and therefore can be used to both produce and test simulations and models of
737 cell organization. For example, models that predict the organization of sets of cellular structures
738 based on conceptual or mechanistic constraints would need to recreate the location, variability,

739 amount, and relative positional dependencies of these structures. Comparing these quantitative rule-
740 building constraints with those of cells in other cell states and testing how underlying mechanisms can
741 generate these rules via models and simulations, will permit us to move towards a deeper
742 understanding of cell organization, aspiring to find general principles.

743 The analysis approaches presented here reveal some initial glimpses into new biologically
744 interesting phenomena. The cell and nuclear shape space that represents the hiPSC Single-Cell
745 Image Dataset revealed that we can reduce the vast complexity of 3D cell and nuclear shape into a
746 human-interpretable understanding of the mean cell shape and the variation around it. It also creates
747 a coordinate system by which we can now cluster groups of similar cells for further analysis and
748 identify outliers. This shape space exposes a relationship between the behavior of the overall cell
749 shape and the nuclear shape, which deserves deeper investigation. Next, the methods we developed
750 to morph cells and their cellular structures together within clusters of similarly shaped cells permit both
751 the analysis of cellular structure stereotypy and concordance. In principle, the overall concordance
752 among structures can span a range. At one extreme, all structures could be coupled, e.g. every
753 structure depending on every other structure; whereas at the other extreme, every structure could be
754 independent from every other structure. We found that the location concordance of all the cellular
755 structures clustered naturally into an ordered compartmentalization of the cell, from the center of the
756 nucleus outward. Then we tested whether this result was valid for only a particular cell shape, e.g.,
757 the mean cell shape, or changed with systematic changes in that shape. We were surprised at how
758 robust both the stereotypy and the concordance proved to be across all of the cell shape variation in
759 our population.

760 Our systematic analysis of how cellular structure volume relates to cell and nuclear size also
761 raises interesting questions. The apparent lack of correlation between endosome volumes and cell
762 and nuclear size could arise in several ways. First, endosome size may just depend on the size of
763 other cellular structures, highlighting a need for experimental data that includes the size of other
764 structures, going beyond just the cell and nuclear size. Second, we use Rab-5A as the marker for
765 endosomes; but it marks only a subset of endosomes, the early endosomes, and perhaps different
766 subsets of endosomes relate differently to cell and nuclear size. Interesting questions also follow from
767 our observation that the variation in the volumes of all nuclear structures was better explained by
768 nuclear size than cell size metrics; and furthermore, the nuclear surface area was more tightly
769 coupled to these structures than to nuclear volume. Nuclear speckles, for example, exhibited a
770 surprisingly strong relationship with nuclear surface area. This is intriguing in light of the possible
771 connection between transcript splicing (which occurs at nuclear speckles) and increased rates of
772 nuclear export (Valencia et al., 2008). This is the first time that the relationship between cellular
773 structure size and cell and nuclear size has been compared among so many different cellular
774 structures all in the same consistent experimental system.

775 Transcriptomics and proteomics are having a great impact in understanding how the “building
776 blocks” of the cell generate and regulate cell behavior and disease. Recent single-cell versions of
777 these studies, particularly for gene expression, together with dimension reduction approaches to
778 statistically identify and separate groups of similar individual cells, have permitted new insights into
779 different cell types and states. Our studies add a new dimension to these analyses by incorporating
780 the spatial organization of cell structures, that is, where and when these parts come together in space
781 and time to drive that function. Our approach relies on the well-established tight linkage between 3D
782 cell organization and cell function, aspiring to use this to identify cell types and states from images,
783 and relate them to single cell gene expression profiles (Gerbin et al., 2020). These approaches were
784 built with live cell imaging in mind and thus are poised to incorporate dynamics. Recently, studies
785 combining quantitative measures of sarcomere organization with gene expression in the same
786 individual cardiomyocytes demonstrates the importance of incorporating the spatial cell organization
787 metrics for a more complete classification of cell states (Gerbin et al., 2020).

788 Other recent systematic image-based approaches have catalogued the location of human
789 proteins in several cell types and used protein and structure locations within cells to identify
790 differences in intracellular spatial patterns among cells in distinct states (Caicedo et al., 2017; Gut et
791 al., 2018; Thul et al., 2017). Our work complements these approaches with its focus on analyses of
792 3D cell organization at the level of cellular structures, and on the generation of quantitative
793 measurements in a human-interpretable manner. Taken together, these studies bring a critical
794 missing dimension, i.e., the spatio-temporal component, to the single cell revolution (Aldridge and
795 Teichmann, 2020). Our study furthers this community goal by adding critical tools, data, and analyses
796 that show the importance of studying large populations of cells and embracing their variations to
797 further our understanding of the underlying rules that organize cells. As part of our mission, we aspire
798 to democratize this emerging area of research; the full image dataset and analysis algorithms
799 introduced here, as well as all the reagents, methods, and tools needed to generate them, are shared
800 in an easily accessible way (Allencell.org). This data is available to all for further biological analyses
801 and as a benchmark for new development of tools and approaches moving towards a holistic
802 understanding of cell behavior.

803 **Acknowledgments:**

804 We thank Joan Brugge, Gaudenz Danuser, Michael Elowitz, Tom Goddard, Quincey Justman,
805 Jennifer Lippincott-Schwartz, Wallace Marshall, Sean Palacek, Zach Pincus, and Tom Pollard for
806 helpful scientific discussions. The WTC line that we used to create our gene-edited cell lines was
807 provided by the Bruce R. Conklin Laboratory at the Gladstone Institute and UCS. We wish to
808 acknowledge the Allen Institute for Cell Science founder, Paul G. Allen, for his vision, encouragement
809 and support.

810

811

812 **Author Contributions**

813 J.M.B., J.A.C., J.C., T.P.D., M.A.F., N.G., K.A.G., B.W.G., R.N.G., A.H., M.C.H., C.H., A.R.H., G.T.J., J.J.L.,
814 A.N., A.M.N., L.P., S.M.R., M.M.R., B.R., L.M.S., M.S., J.S., J.E.S., J.A.T., D.J.T., D.M.T., A.P.T., V.V., M.P.V.,
815 W.W., C.Y. contributed to the development and/or design of the methods used in this paper. M.B., J.M.B.,
816 J.A.C., B.C., J.C., Z.J.C., S.D., S.D., T.P.D., R.M.D., T.J.F., G.F., T.G., L.J.H., H.C.H., E.J.I., G.T.J., G.R.J., B.K.,
817 J.J.L., G.E.M., S.L.M., K.M., A.N., L.P., T.A.P., M.M.R., L.M.S., M.S., J.S., S.S., M.F.S., M.J.S., D.J.T., D.M.T.,
818 R.V., M.P.V., W.W., T.W., C.Y., R.Y. contributed to the software developed to create, assess, and analyze the
819 dataset. J.E.A., A.B., J.A.C., J.C., M.E.C., S.Q.D., M.A.F., N.G., K.A.G., T.G., B.W.G., A.H., C.H., J.J.L., H.M.,
820 I.A.M., A.N., A.M.N., L.P., S.M.R., M.M.R., B.R., L.M.S., J.E.S., D.J.T., D.M.T., A.P.T., V.V., M.P.V., C.Y., R.J.Z.
821 contributed to validation of the cell lines used, images taken and data that lies herein to ensure reproducibility.
822 J.E.A., J.M.B., J.A.C., S.Q.D., R.M.D., M.A.F., K.A.G., T.G., R.N.G., A.H., G.R.J., T.A.K., J.J.L., H.M., L.P.,
823 S.M.R., M.M.R., B.R., L.M.S., J.S., Y.S., D.M.T., A.P.T., V.V., R.V., M.P.V., C.Y., R.J.Z. contributed to the
824 formal analysis that is shown in main figures and supplementary materials. J.E.A., A.B., S.C., M.E.C., C.M.D.,
825 S.Q.D., M.A.F., N.G., J.L.G., K.A.G., T.G., B.W.G., A.H., M.C.H., C.H., W.W.L., S.A.L., H.M., I.A.M., A.N.,
826 A.M.N., L.P., S.M.R., B.R., J.E.S., W.J.T., J.A.T., D.J.T., A.P.T., V.V., M.P.V., C.Y., R.Y., R.J.Z. contributed to
827 collecting the data, performing the experiments or analyzing the imaging data. B.B., J.M.B., J.A.C., J.D., M.A.F.,
828 B.W.G., A.H., A.N.L., R.M., T.L.M., M.M.R., B.R., J.S., J.E.S., D.M.T., W.W. contributed to providing the
829 reagents, materials or compute resources needed to accomplish this study. J.E.A., M.B., B.B., A.B., J.M.B.,
830 S.C., J.C., M.E.C., C.M.D., S.Q.D., R.M.D., T.J.F., M.A.F., N.G., J.L.G., K.A.G., T.G., B.W.G., A.H., M.C.H.,
831 C.H., T.A.K., W.W.L., J.J.L., S.A.L., K.M., I.A.M., A.N., M.M.R., B.R., J.S., S.S., J.E.S., M.J.S., W.J.T., D.J.T.,
832 D.M.T., A.P.T., C.Y. contributed to management activities to annotate, scrub data, and maintain research data.
833 J.C., M.E.C., T.P.D., N.G., C.H., A.R.H., G.T.J., T.A.K., A.N., S.M.R., J.A.T., R.V., M.P.V., C.Y. helped prepare,
834 create and present the original draft of this work. A.B., J.C., K.R.C.M., T.P.D., C.L.F., N.G., C.H., A.R.H., G.T.J.,
835 T.A.K., I.A.M., S.M.R., J.A.T., R.V., M.P.V., C.Y. helped provide critical review and editing of this work. J.A.C.,
836 K.R.C.M., T.P.D., N.G., G.T.J., T.A.K., M.M.R., D.M.T., R.V., M.P.V. helped prepare the visualization of the work
837 described, specifically via visual tools, resources, and data presentation. J.E.A., K.R.C.M., M.A.F., N.G., K.A.G.,
838 R.N.G., A.H., A.R.H., G.T.J., T.A.K., I.A.M., A.M.N., S.M.R., B.R., D.M.T., W.W. provided oversight and
839 leadership, including planning and execution of the cross-team project management. K.R.C.M., R.N.G., A.H.,
840 G.T.J., I.A.M., S.M.R., S.S., D.M.T. provided management, coordination of the project and administration of
841 resources.

842 MATERIALS AND METHODS

843

844 RESOURCE AVAILABILITY

845

846 Lead Contact

847 Further information and requests for resources and reagents should be directed to and will be
848 fulfilled by the Lead Contact, Susanne Rafelski (susanner@alleninstitute.org).

849

850 Material Availability

851 Using the Wild Type WTC-11 hiPSC line background (Kreitzer et al., 2013), we previously
852 generated the Allen Cell Collection of hiPSC lines in which each gene-edited cell line harbors a
853 fluorescent protein endogenously tagged to a protein representing a distinct cellular structure of the
854 cell (Roberts et al., 2017b). Fifteen additional Allen Cell Collection lines were generated using the
855 same methods in this study. The cell lines are described at <https://www.allencell.org> and are available
856 through Coriell at <https://www.coriell.org/1/AllenCellCollection>. For all non-profit institutions, detailed
857 MTAs for each cell line are listed on the Coriell website. Please contact Coriell regarding for-profit use
858 of the cell lines as some commercial restrictions may apply.

859

860 Data and Code Availability

861 The Datasets generated during this study are available at Quilt as packages:

- 862 • Full dataset: https://open.quiltdata.com/b/allencell/packages/aics/hipsc_single_cell_image_dataset
- 863 • 12X colony dataset:
864 https://open.quiltdata.com/b/allencell/packages/aics/hipsc_12x_overview_image_dataset
- 865 • Supplementary MYH10 repeat dataset:
866 https://open.quiltdata.com/b/allencell/packages/aics/hipsc_single_cell_image_dataset_supp_myh10
867 [0](#)
- 868 • Tutorials and demo for how to access the data for different purposes:
869 <https://github.com/AllenCell/quilt-data-access-tutorials>
- 870 • Original/source data for figures in the paper are available in Github: [https://github.com/aics-](https://github.com/aics-int/cvapipe_figure_notebooks)
871 [int/cvapipe_figure_notebooks](https://github.com/aics-int/cvapipe_figure_notebooks)

872

873 The code supporting the current study has been deposited in Github repositories and the released
874 code repositories and packages use the following packages in parts, including Numpy (Harris et al.,
875 2020), Scipy (Virtanen, 2020), Napari (Nicholas Sofroniew et al., 2019), Seaborn

- 876 <https://seaborn.pydata.org/citing.html>, Py Torch (Paszke et al., 2019), ITK (McCormick et al., 2014),
877 pandas (McKinney, 2011), matplotlib (Hunter, 2007), and label free (Ounkomol et al., 2018):
- 878 • Image segmentation code, trained models, and demo Jupyter notebooks have been released at
879 https://github.com/AllenCell/segmenter_model_zoo.
 - 880 • Segmentation code used to reproduce structure segmentations from a set of algorithms to choose
881 from, each with restricted numbers of parameters to tune are available at
882 <https://github.com/AllenCell/aics-segmentation>.
 - 883 • Mitotic image classifier code (Falcon and Cho, 2020; Paszke et al., 2019), (for both training and
884 testing) and all trained models are available at https://github.com/AllenCell/image_classifier_3d.
 - 885 • Code used to generate contact sheets for quality control single-cell visualizations of all segmented
886 cells is available at <https://github.com/AllenCellModeling/actk>
 - 887 • Code used for feature calculation:
 - 888 ○ aicsfeature (<https://github.com/AllenCell/aicsfeature>)
 - 889 ○ spherical harmonics parameterization (<https://github.com/AllenCell/aics-shparam>)
 - 890 ○ cytoplasmic parameterization (<https://github.com/AllenCell/aics-cytoparam>)
 - 891 • Code used to perform organelle size-scaling analysis
892 (<https://github.com/AllenCell/stemcellorganellesizescaling>)
 - 893 • Code used to perform morphing, compute shape modes, and calculate multi-resolution Pearson
894 correlation analysis on 3D single cell images (Rocklin, 2015)
895 (https://github.com/AllenCell/cvapipe_analysis)
 - 896 • Code to create 12X colony dataset and to perform cell height regression
897 (<https://github.com/AllenCell/colony-processing>)
 - 898 • Software will be shared under the Allen Institute Software License and Contribution Agreement,
899 subject to any applicable third-party licensing restrictions.
 - 900 • Datasets will be shared under the Allen Institute Terms of Use:
901 <https://alleninstitute.org/legal/terms-use/>.

902

903 **METHOD DETAILS**

904

905 **Cell lines and cell culturing**

906 The gene-edited cell lines used in this study were created using the parental WTC-11 hiPSC
907 line, derived from a healthy, male donor (Kreitzer et al., 2013). Each gene-edited cell line harbors a
908 fluorescent protein endogenously tagged to a protein representing a distinct cellular structure (**Table**
909 **1**). The complete list of cell lines can be found in the Resource Table. The CRISPR/Cas9-mediated
910 genome editing methodology used to generate these cell lines was previously described in (Roberts

911 et al., 2017b). The tagging strategy for AAVS1 safe harbor targeting was altered to additionally
912 introduce a strong exogenous promoter for expression of CAAX-mTagRFP-T as described previously
913 (Hockemeyer et al., 2009; Ocegüera-Yanez et al., 2016). The identity of the unedited parental line
914 was confirmed with short tandem repeat (STR) profiling testing (29 allelic polymorphisms across 15
915 STR loci compared to donor fibroblasts (<https://www.coriell.org/1/AllenCellCollection>). Since WTC-11
916 is the only cell line used by the Allen Institute for Cell Science, edited WTC-11 cells were not re-tested
917 because they did not come into contact with any other cell lines.

918 The culture and handling protocols for all used hiPSC lines was internally approved by an
919 oversight committee and all procedures performed in accordance with the National Institutes of
920 Health, National Academy of Sciences, and Internal Society for Stem Cell Research guidelines. All
921 cell lines were expanded and grown on an automated cell culture platform developed on a Hamilton
922 Microlab STAR Liquid Handling System (Hamilton Company). This platform is summarized in part in
923 **Figure 1A**. Three daily workflows were performed on this platform, (1) plate maintenance, (2)
924 passaging, and (3) Matrigel coating of plates. Plate maintenance included the replacement of old
925 media with fresh media for both 6- and 96-well plates. Cells were cultured in a Cytomat 24 (Thermo
926 Fischer Scientific) at 37°C and 5% CO₂ in mTeSR1 medium with and without phenol red (STEMCELL
927 Technologies), supplemented with 1% penicillin-streptomycin (Thermo Fischer Scientific).

928 Cells were passaged every 4 days for up to 10 passages post-thaw. Cells were dissociated
929 into a single cell suspension with 37°C StemPro Accutase cell dissociation reagent (Thermo Fisher
930 Scientific) and counted with a Vi-CELL XR Series cell viability analyzer and associated Vi-CELL XR
931 sample vials (Beckman CoulterA). Cells were re-plated in mTeSR1 medium supplemented with 1%
932 penicillin-streptomycin (Thermo Fischer Scientific) and 10 mM Rho-associated protein kinase (ROCK)
933 inhibitor (Stemolecule Y-27632, STEMCELL Technologies) for 24 hr. Cell culture plates used for cell
934 expansion were clear-skirt, sterile, plastic 6-well plates with lid with condensation rings (Greiner Bio-
935 One). For imaging, samples were plated on glass-bottom, black-skirt, 96-well plates with #1.5 optical
936 grade cover glass (Cellvis). Cells were seeded at a density of 1.3x10³ to 3.0x10³ in 96-well plates
937 and at 80x10³ to 175x10³ in 6-well plates.

938 For most of the dataset, cell culture plates were coated with growth factor reduced (GFR)
939 Matrigel basement membrane matrix, phenol red-free (Lot # 5292003, Corning) diluted with Dulbecco's
940 modified eagle medium (DMEM)/F-12 (Thermo Fischer Scientific) for a final protein concentration of
941 0.337 mg/mL. Matrigel coating was performed at 4°C with 100 µl and 1,500 µl added to each 96-well
942 and 6-well, respectively. Plates were incubated at room temperature (RT) for 2hr and Matrigel
943 removed before cell seeding. For the last two cell lines (cohesins and nuclear speckles) imaged on
944 the pipeline, the Matrigel coating protocol was adjusted for improved cell plating. These cells were
945 also plated on a new lot of Matrigel basement membrane matrix, phenol red-free (Lot # 9021357) at a

946 final protein concentration of 0.185 mg/mL. For these samples, glass bottom 96-well plates coated
947 with Matrigel were incubated overnight at 4°C and for an additional 2 hr at 37°C before removing
948 Matrigel at RT. Further details for cell culture reagents and consumables can be found in the
949 Resource Table and standard protocols can be found at www.allencell.org.

950

951

952 **Cell culture and imaging sample quality control**

953 Rigorous and standardized quality control (QC) workflows for cell culture health were
954 performed at each passage, before imaging the cells at high resolution, and following the completion
955 of imaging a cell line. These QC workflows included cell and morphology assessment via microscopy,
956 cell stemness marker expression with flow cytometry, and outsourced cytogenetic analysis.

957 hiPSC morphology was evaluated by expert scientists for both plastic 6-well and glass-bottom
958 96-well plates and was examined 4 days post-passaging. Individual wells of plastic 6-well plates were
959 deemed ready for passage when cells reached ~85% confluency and displayed typical morphologies
960 associated with hiPSCs that have preserved the expression of stemness markers (Roberts et al.,
961 2017a). Exclusion criteria included, but were not limited to, under- or over-confluency, presence of
962 morphology associated with differentiating cells, and over 5% of cell death. The morphology of cells
963 and colonies grown on glass bottom 96-well plates was also examined prior to 3D field of view (FOV)
964 image acquisition. 12X well overview images were used to exclude wells that did not meet the
965 following four morphology criteria requirements: less than three occurrences of 1) colonies presenting
966 atypical crater-like morphology, 2) lifted colonies (ball-like morphology), 3) partially lifted colonies
967 (edges lifting) and 4) morphology associated with differentiation (Roberts et al., 2017a).

968 Following the completion of all 3D FOV image acquisition for a given cell line, two types of QC
969 were performed to ensure hiPSCs had retained stemness marker expression and normal G-band
970 karyotyping throughout the imaging period as previously described (Coston et al., 2020; Roberts et
971 al., 2017a). All cell lines imaged during the three years of data acquisition and included in the hiPSC
972 Single-Cell Image Dataset passed these QC requirements.

973

974

975 **Image acquisition**

976 The following methods are described in the order they were performed for a given image
977 acquisition workflow on the imaging pipeline. The image acquisition workflow and experimental setup
978 evolved over the three years of dataset collection and was versioned as such. Below is the list of all
979 pipeline image acquisition workflows and a description of each update and modification. The list of
980 pipeline workflow versions used to acquire each cell lines can be found in **Table 1**.

981

Workflow version	Description of changes and updates
Pipeline 4.0	Single camera system for interwoven acquisition of four channels. Original transmitted light = white light LED. Specimen exposed to dual peak emissions with highest at 460 nm and range of 400-700 nm. Collecting light with 525/50 bandpass filter, from 500-550 nm. Mode C acquisition performed without photoprotective cocktail.
Pipeline 4.1	Single camera system for interwoven acquisition of four channels. Same original transmitted light source as 4.0. Added photoprotective cocktail with mode C acquisition.
Pipeline 4.2	Single camera system for interwoven acquisition of four channels. Same original transmitted light source as 4.0. Photoprotective cocktail used with mode C acquisition. New emission filter added for acquisition of mTagRFP-T with a 600/50 nm band pass filter.
Pipeline 4.4	Second camera added to all systems. Using dual camera system to image four channels (bright-405-488-638 nm) with interwoven acquisition of 2X two channels simultaneously. New Transmitted Light: Red 740 nm LED Transmitted Light. Specimen exposed to peak emission 740 with narrow range. Collecting light with 706/95 bandpass filter, so from 660-750 nm. Using piezo z stage for fast movement in z.

982

983 **Microscopy**

984 Imaging was performed on Zeiss spinning-disk confocal microscopes with 10X/0.45 NA Plan-
985 Apochromat or 100X/1.25 W C-Apochromat Korr UV Vis IR objectives (Zeiss) and Zen 2.3 software
986 (blue edition; Zeiss) unless otherwise specified. The spinning-disk confocal microscopes were
987 equipped with a 1.2X tube lens adapter for a final magnification of 12X or 120X, respectively, a CSU-
988 X1 spinning-disk scan head (Yokogawa) and two Orca Flash 4.0 cameras (Hamamatsu). Standard
989 laser lines were used at the following laser powers measured with the 10X the objective; 405 nm at
990 0.28 mW, 488 nm at 2.3 mW, 561 nm at 2.4 mW and 640 nm at 2.4 mW unless otherwise specified.
991 The following Band Pass (BP) filter sets (Chroma) were used to collect emission from the specified
992 fluorophore; 450/50 nm for detection of DNA dye, 525/50 nm for detection of mEGFP tag, 600/50 nm
993 for detection of mTagRFP-T tag and 706/95 nm for detection of cell membrane dye. Images were
994 acquired with 200 ms exposure time unless otherwise specified. The microscope setup allowed us to
995 collect either all channels with a single camera (Pipeline 4.0-4.2, see description above) or two
996 channels simultaneously, either the bright field and mEGFP or the cell membrane and DNA dyes

997 (Pipeline 4.4, see description above). Cells were imaged in phenol red-free mTeSR1 media on the
998 stage of microscopes outfitted with a humidified environmental chamber to maintain cells at 37°C with
999 5% CO₂ during imaging. Transmitted light (bright field) images were acquired using a white LED light
1000 source with broad emission spectrum (pipeline 4.0-4.2) or a red LED light source with peak emission
1001 of 740 nm with narrow range and a BP filter 706/95 nm for bright field light collection (Pipeline 4.4
1002 only). A Prior NanoScan Z 100 mm piezo z stage (Zeiss) was used for fast acquisition in z (Pipeline
1003 4.4 only).

1004

1005 ***Well overview and manual well position selection***

1006 Typical imaging sessions started with a bright field overview image acquisition of wells from
1007 selected rows of a 96-well plate as 2D, 12X tiled images before cell membrane and DNA dye staining.
1008 These well overview images were used for final evaluation of cell morphology (see above) and
1009 manual or automated position selection for 3D FOV acquisition at 120X in wells satisfying QC criteria
1010 requirements (see description above). Manual selection of positions to be imaged at 120X was
1011 performed using the 12X overview images and stage function in Zen software. Manual position
1012 adjustments were also made at 120X using streaming bright field imaging to satisfy the requirement of
1013 each mode of imaging.

1014

1015 ***Imaging modes***

1016 Colony position selection was performed manually using the stage function in Zen software or
1017 as described below using an automated method (for the last six cell lines imaged) for imaging mode
1018 A. One position per colony was selected approximately half-way between the colony edge and center
1019 such that the imaged FOV did not fall at the edge nor at the center of the colony. In mode B, positions
1020 were also selected as per mode A followed by manual adjustment of the FOV position using
1021 transmitted light and streaming bright field imaging to navigate to a region enriched in mitotic cells.
1022 This mode was used to substantially increase the number of mitotic cells imaged in an FOV by 3-fold.
1023 Operators were trained on how to identify mitotic morphology from just bright field images using
1024 merged DNA dye and bright field images (**Figure 1A**, mode B). In mode C, three positions per colony
1025 were selected; a mid-center position area (as in mode A), a position right at the colony edge and a
1026 position just inward from the edge in an area referred as the ridge due to the tendency of these cells
1027 to grow taller until they flatten into the center area of the colony (**Figure 1A**; **Figure S7**). Due to the
1028 increased photosensitivity of the cells located at the edge of the colony, a photoprotective cocktail
1029 (see “Dye Staining” section below) was used when imaging in mode C to prevent premature cell
1030 retraction, blebbing and death. Mode C positions were selected manually for all cell lines.

1031

1032

1033 ***Automated well position selection***

1034 We developed an automated method that segments the colonies from a 12X well overview
1035 image and automatically suggests positions for 3D FOV acquisition based on the distance from the
1036 edge of a colony satisfying mode A criteria (see description above). Tiles from the well overview
1037 images were acquired with a 10% overlap and stitched using a processing function in Zen software.
1038 The automated position selection method segmented the colonies in the image with the following
1039 image processing steps developed in Python; 1) rescale intensity to increase contrast of colony edges
1040 from background, 2) apply Sobel filter (Scikit-Image) to identify colony edges and fill the holes to
1041 segment entire colony, 3) correct for segmentation artifacts with erosion, dilation on segmentation and
1042 removal of small objects, 4) identify centers of individual colonies with a distance map on binary
1043 segmentation, 5) separate connecting colonies using the center coordinates of colonies and binary
1044 segmentation of colony areas with watershed method, and finally 6) compute a distance map for each
1045 colony. Our criteria for position selection were as follows, 1) one position was selected per colony, 2)
1046 no more than 10 positions were selected per well, 3) the position had to be from a colony greater than
1047 34992 μm^2 (corresponding to colony size with a uniform flattened and well-packed central area), 4)
1048 the position had to be imaged approximately half-way between the edge and center of the colony
1049 (mid-center). To fulfill these criteria, the algorithm first filtered colonies based on their size, and
1050 selected mid-center colony positions (x-y coordinates). If more than 10 positions per well were
1051 automatically identified, the method gave preference to positions selected in larger colonies and in
1052 colonies closer to the center of the well. A graphical user interface was developed to assist users in
1053 viewing and confirming the proposed algorithm-generated positions for 3D FOV imaging. The user
1054 had the flexibility of moving, adding or deleting positions to finalize the list of FOV to be imaged at
1055 higher magnification for that imaging session. The list of positions was then saved as a text file with
1056 the stage coordinates and position number in the Zen software readable format (.czsh) and integrated
1057 into the experiment xml file for 3D FOV imaging at 120X.

1058
1059 ***DNA and cell membrane dye staining***

1060 Following well overview acquisition, the cell membrane and DNA of cells from selected wells
1061 were stained with fluorescent dyes. Wells were first incubated at 37°C and 5% CO₂ for 20 min with a
1062 DNA dye, NucBlue Live (Thermo Fisher Scientific, 1:16.66) diluted in phenol red-free mTeSR1
1063 medium. A cell membrane dye, CellMask Deep Red (CMDR, Thermo Fisher Scientific) was then
1064 added to the well (in the continued presence of NucBlue Live) at a final concentration of 5X (earlier
1065 lot) or 3X (last 2 lots, adjusted to provide equivalent contrast to noise ratio within a 2.5 hr imaging
1066 session) and the 96-well plate was incubated for an additional 10 min at 37°C and 5% CO₂. Each well
1067 was washed once with phenol red-free mTeSR1 medium before a final 200 μl of phenol red-free

1068 mTeSR1 medium was added per well and the plate returned to the stage of the microscope. In mode
1069 C acquisition, a photoprotective cocktail (1mM ascorbic acid, 0.3 U/ml OxyFluor and 10 mM lactate)
1070 was mixed into the phenol red-free mTeSR1 media before it was added to the well. For consistency,
1071 we limited the cell staining to a single row, or 10 wells, per plate at a time and imaged for a maximum
1072 of 2.5 hr post completion of the staining protocol. We limited the imaging time to 2.5 hours since we
1073 saw no adverse effects of the dyes on cell cycle (evaluated as % mitotic cells in cell colonies) or cell
1074 viability (evaluated as increased presence of dead cells on top of colonies) within that time frame. We
1075 imaged halfway between the edge and center of a colony to avoid imaging FOVs with reduced dye
1076 penetration at the center of large, tightly packed colonies.

1077

1078 ***3D FOV image acquisition***

1079 After the final wash with phenol red-free mTeSR1 media, plates were returned to the stage of
1080 the spinning-disk confocal microscope and all 3D FOVs, at pre-selected positions, were acquired with
1081 a 100X/1.25 NA objective at a final magnification of 120X. Four channels were acquired at each z-
1082 step (interwoven channels) in the following order: bright field, mEGFP or mTagRFP-T, CMDR and
1083 NucBlue Live with laser powers and exposure times as stated above. For the single camera system
1084 acquisition only, empty channels were acquired between each channel with 0.3 ms exposure time to
1085 reduce noise introduced during the filter position change. This was necessary due to the long travel
1086 range of the filter wheel moving between four different positions at each z-step. Pipeline 4.4 3D FOV
1087 acquisition was performed with two cameras using two interwoven sets of simultaneous acquisitions.
1088 In this case, bright field and CMDR channels were acquired on the back camera and all other
1089 channels acquired on the left camera. Resulting images from either single or dual camera systems
1090 were of 16 bits and 924 x 624 pixel² in x-y dimension after 2x2 binning. Images from the dual-camera
1091 system required channel alignment (see below). Following channel alignment, the final images were
1092 cropped into a final size of 900 x 600 pixels² in the x-y dimension. FOVs from both single and dual
1093 camera systems had a final x-y pixel size of 0.108 μm and z-stacks composed of 50–75 z-slices (to
1094 encompass the full height of the cells within an FOV) acquired at a z interval of 0.29 μm .

1095

1096

1097 **Post-acquisition FOV image processing**

1098

1099 ***Channel alignment for dual camera acquired images (Pipeline 4.4 only)***

1100 Optical control images were acquired at the start of each data acquisition day to monitor
1101 microscope performance. Optical control images of TetraSpeck microspheres or the “field of
1102 ring” pattern on the Argolight HM slide (Argolight) were used to register and align the appropriate
1103 channel images of an FOV acquired with two cameras. A z-stack of 10 to 30 z-slices of these patterns

1104 was acquired at 120X with all four fluorescent channels. Channel images from the 638 nm laser line
1105 and 706/95 nm BP filter (back camera) and 488 nm laser line and 525/50 nm BP filter (left camera)
1106 were used to generate an affine transformation matrix identifying the shift in x-y, rotation and scaling
1107 factor between the 638 nm (from the back camera) and 488 nm (from the left camera) wavelength
1108 channels. We used the z-slice with maximum focus along the z-axis. The two channel images were
1109 pre-processed separately by normalizing the intensities and applying gaussian smoothing prior to
1110 segmenting the objects such as individual beads or rings with intensity thresholding. Due to the nature
1111 of the sample preparation of TetraSpeck beads, which randomly adhere to the glass, we excluded
1112 some beads based on the following criteria: 1) overlapping beads, 2) beads that are outside of the
1113 range of an expected bead size and intensity, and 3) beads that have inconsistent centroid location
1114 (mass versus peak intensity). Centroid locations of segmented objects (beads or rings) from both 638
1115 nm and 488 nm channels were compared and only objects in close proximity (within 5 pixels) between
1116 the two channels were kept. The exclusion steps were not necessary with the stable and consistent
1117 field of ring pattern of the Argolight HM slide. Using the two sets of centroid locations of objects, the
1118 method estimated a similarity transform matrix with the “estimate_transform” function in scikit-image
1119 that transforms the image with translation, rotation and scaling. The values from this matrix were also
1120 used to identify any deviations from the normal trend indicating potential system performance issues
1121 over time. The affine transformation matrix was applied on every z-slice of the channel acquired on
1122 the back camera (bright field and 638 nm) and as such aligned to the reference channel images
1123 acquired on the left camera (405 nm, 488 nm and 561 nm) with a Warp function (scikit-image). FOVs
1124 were then cropped in x-y for a final dimension of 900 by 600 pixels² to remove empty pixels
1125 introduced in the bright field and 638nm channel images by the alignment.

1126

1127 ***3D FOV image quality control***

1128 All 3D FOV images were visually inspected by experts for obvious issues related to the
1129 experimental settings. Typical exclusion criteria were related to microscope acquisition system failures
1130 (laser, exposure time, z-slice positioning in relation to cell height, empty or out of order channels), or
1131 any other issues that would cause downstream processing to fail or analysis steps to identify outliers.
1132 Some of these QC steps were also automated with a series of Python scripts to ensure a more
1133 systematic and standardized way to catch problematic FOVs and exclude any outliers. To do so,
1134 intensity metrics were extracted from each channel of each FOV and trends and averages were used
1135 to determine exclusion thresholds or cutoff values. Overall, three main automated exclusion FOV QC
1136 steps were applied to the hiPSC Single-Cell Image Dataset; channel intensity out of range, z-stacks
1137 with incomplete cell height, and z-slice empty or out of order.

1138

1139

1140 *FOV channel intensity quality control*

1141 The FOV channel intensity QC script calculated the minimum, maximum, median, 99.5th and
1142 0.5th percentile pixel intensity value in each channel for each FOV. FOVs were flagged if the median
1143 intensity of one channel was outside a predetermined range (low and high cutoffs, see values below).
1144 These cutoff values were based on offset, noise and maximum intensity values of the microscopes,
1145 fluorescent tags and dyes imaged.

1146 **Cutoffs for median intensity**

Channels	Low cutoff (AU)	High cutoff (AU)
bright field	0	50000
405 nm	400	430
488 nm	400	1600
561 nm	400	700
638 nm	400	8000

1147

1148 Low cutoffs for the maximum intensity in the 405 nm (DNA) and 638 nm (cell membrane)
1149 channels were also applied to ensure the minimum required contrast in the images for successful
1150 single cell segmentation.

1151 **Cutoffs for max intensity**

Channels	Low cutoff (AU)
405 nm	500
638 nm	635

1152

1153 Given a normal distribution of FOV intensities, we also excluded from the dataset all individual
1154 FOVs with a channel median intensity within the bottom 0.5th percentile of the whole dataset. We
1155 calculated a z-score for each channel of each FOV and excluded all FOVs that had a channel
1156 intensity with z-score of 2.58 below the mean.

1157

1158 *Automated detection of z-stack with incomplete cell height quality control*

1159 We automated the detection of z-stacks with incomplete cell height in a FOV due to mis-
1160 positioning or mis-sampling of a z-stack acquisition. The cell membrane channel (638 nm) was used
1161 to determine whether the top and/or bottom of the cell were included in the z-stack image. The

1162 intensity of the cell membrane channel image was first normalized to the maximum intensity of the cell
1163 membrane channel image. Next, the median intensity and contrast (maximum intensity-background
1164 intensity)/maximum intensity) for each z-slice were calculated to generate an intensity and contrast
1165 profile along the z-axis. Local maxima of the intensity profile were detected with a “peak detection”
1166 method described in SciPy-image ((Virtanen, 2020); scipy.org), where the lower peak corresponds to
1167 the bottom of a z-stack and the higher peak corresponds to the top of a z-stack. In the scenario where
1168 more than 2 peaks were detected, the method used the top-most peak and the bottom-most peak and
1169 the contrast profile to refine the measured range of the z-stack.

1170 Thresholds of contrast values for bottom (0.2) and top (0.19) of a z-stack were estimated from
1171 data trends of the entire dataset. Using these threshold values, the method iterated from the
1172 top/bottom peaks detected to the full range of the z-stack and reported the closest z-slice to reach the
1173 thresholds as the detected top/bottom of the cells for this z-stack. We also measured the rate of
1174 change in contrast in the detected top and bottom 5 slices of each z-stack and flagged z-stacks as
1175 incomplete cell height if the rate of change in the top 5 slices were smaller than -0.015 and the bottom
1176 5 slices were smaller than -0.01 (in contrast units, see contrast definition above). To ensure cell height
1177 completeness any FOV with detected top/bottom z-slices within 5 slices of the first and last slice of the
1178 z-stack were flagged as either an “incomplete -missing top” or a “incomplete -missing bottom” and
1179 excluded these FOVs from the datasets.

1180

1181 *Out of order z-stacks in FOV quality control*

1182 Out of order z-stacks were also observed. We generated an algorithm capable of identifying if
1183 the z-stack first z-slice had the highest median intensity, indicating that the z-stacks were placed in
1184 improper order by the acquisition software. We excluded any FOV with a first z-slice registering the
1185 maximum intensity and flagged the FOV as “z-stacks out of order”.

1186

1187

1188 **3D segmentation**

1189

1190 ***Cell and nuclear segmentation***

1191 To segment each individual cell and its corresponding DNA from the membrane dye and DNA
1192 dye channels of each 3D z-stack, we used the deep learning-based cell and nuclear instance
1193 segmentation algorithm developed as part of Allen Cell and Structure Segmenter (Chen et al., 2018).
1194 We combined the Segmenter’s Iterative Deep Learning workflow and the Training Assay approach to
1195 ensure accurate and robust segmentation for downstream quantitative analysis. Complete step-by-
1196 step details of this algorithm are described in (Chen et al., 2018).The code, trained models, and demo
1197 Jupyter notebooks have been released at https://github.com/AllenCell/segmenter_model_zoo.

1198 In the Training Assay approach, a secondary experimental assay that is more amenable to accurate
1199 segmentation is linked to the primary assay for the purpose of training segmentation models. The
1200 secondary assay is used to generate accurate segmentations, which are then imposed as the target
1201 for training the model to segment the images of the primary assay. As a result, the final segmentation
1202 model can achieve better accuracy and robustness even when running on the poorer-quality primary
1203 assay images. We applied two training assays to develop the cell and nuclear segmentation
1204 algorithm.

1205 The first training assay (**Figure S1E**) addressed the challenge that the membrane dye images
1206 suffered from very weak signal near the top of cells due to both dye labeling of a very thin membrane
1207 and photobleaching even during a single z-stack acquisition via 3D spinning-disk confocal
1208 microscopy. The secondary assay in this training assay used the CAAX cell line containing the
1209 membrane-targeting domain of K-Ras tagged with mTagRFP-T, which made it possible to accurately
1210 delineate cell boundaries, even near the top of cells. This training assay is described in detail in (Chen
1211 et al., 2018). In brief, the first step of this training assay is to obtain the initial semantic (whole FOV)
1212 segmentation of tagged CAAX signal on ten sample images using a semi-automatic algorithm based
1213 on a seeded watershed. Seven images were sorted as having good segmentation and used to train a
1214 CAAX segmentation model. We then applied this CAAX segmentation model on 312 CAAX images to
1215 create a CAAX-based cell segmentation ground truth set, which we then used together with the
1216 membrane dye images to train a membrane dye-based segmentation model. This model robustly
1217 segmented cells including their dimly visible top boundaries, from the membrane dye images in all
1218 18,186 FOV's in the dataset.

1219 The second training assay (**Figure S1F**) was to use images of mEGFP-tagged lamin B1 cells
1220 for segmenting interphase nuclei and mEGFP-tagged H2B cells for segmenting mitotic DNA during
1221 mitosis (representing the “nucleus” during nuclear envelope break-down). Lamin B1 and H2B both
1222 provided more biologically accurate detection of the nuclear boundary. The shell of intensity around
1223 the nucleus in tagged lamin B1 cells was more directly detectable in 3D than the DNA dye images and
1224 both endogenously tagged structure cell lines had better signal to noise compared to both dyes. This
1225 training assay is described in detail in (Chen et al., 2018). Briefly, we began with classic image
1226 segmentation results for lamin B1 where the “shell” of lamin B1 is filled to represent the nucleus. We
1227 sorted eight out of 80 images and used these to train a deep learning model to segment “nuclei” (i.e.
1228 filled shells) from lamin B1 images. We then applied this model on 1,017 lamin B1 images to create a
1229 lamin B1-based nuclear segmentation ground truth set, which we then used together with the DNA
1230 dye images to train a DNA dye-based segmentation model for interphase nuclei. Regions containing
1231 mitotic cells in these images were automatically identified and excluded from training, see (Chen et
1232 al., 2018) for details. In parallel we used H2B images and a classic segmentation workflow to
1233 generate a cleaner segmentation target for training a mitotic DNA segmentation model. We generated

1234 a set of 28 merged segmentation targets (mitotic DNA segmentation from H2B images and interphase
1235 nuclei segmentation by applying the first interphase model on DNA dye images) to train an overall
1236 DNA dye-based nuclear segmentation model. This is the model we applied to all 18,186 FOVs in the
1237 dataset.

1238 To convert the cell and nuclear segmentation model outputs into individual cells (i.e., instance
1239 segmentation), we had to train two additional models: a “cell seeding” model and a “mitotic pair
1240 detection” model. We took advantage of the DNA dye-based nuclear segmentation model to create a
1241 deep learning based “cell seeding” model. This used a subset of 600 images from the same training
1242 images as for the interphase DNA dye segmentation model, but with a modified segmentation target
1243 obtained by shrinking the mask for interphase nuclei (and very early and very late mitosis DNA), and
1244 generating a convex hull for the mask for other mitotic DNA. The binarized membrane segmentation
1245 model output was used to cut the potentially falsely merged seeds from tightly touching nuclei and the
1246 resultant seeds were applied back on the cell membrane segmentation output for use in a seeded
1247 watershed to identify individual cells. We also trained a FasterRCNN-based mitotic pair detection
1248 model, which permitted us to identify mitotic cells that were in anaphase and telophase/cytokinesis
1249 and make sure they were segmented as one cell. Several other steps were performed to enhance the
1250 robustness of the cell and nuclear segmentation for application at scale to the 18,186 FOVs in the
1251 hiPSC Single-Cell Image Dataset. These are described in detail in (Chen et al., 2018) and included
1252 training and applying a label-free segmentation model of nuclei and cell membrane to boost the
1253 robustness when the signals in the DNA dye or membrane dye channel were extremely dim, as well
1254 as several minor steps such as morphological refinement on the segmented nuclei and refinement of
1255 the bottom of the cell. The very bottom surface of the cells protrudes out into the tightly packed
1256 neighboring cells and the z-resolution does not permit proper disentangling of the overlapping parts.
1257 We therefore automatically identified a z-slice with a reasonable cell segmentation near the bottom
1258 and propagated it downward through all other z-slices to the bottom of the cell.

1259 To validate the performance of the cell and nuclear segmentation results, we selected and
1260 inspected a representative set of images (576 images from 22 different cell lines) at the single-cell
1261 instance level. From this validation, we estimated the percentage of well-segmented cells and the
1262 percentage of FOVs for which the segmentation of all cells and nuclei were successful without
1263 obvious errors along the segmented boundaries. We developed an in-house scoring interface in
1264 Python using Napari that allows for overlaying the segmentations on the original images and
1265 inspecting them slice by slice in 3D. Each image was manually scored by at least two human experts.
1266 We found that over 98% of individual cells were well-segmented and over 80% of images generated
1267 successful cell and nuclear segmentations for all cells in the entire FOV. Based on these validation
1268 results, we decided the cell and nuclear instance segmentation algorithm was sufficiently reliable to
1269 be applied to all of the FOVs in the dataset. For quality control purposes, single-cell visualizations of

1270 all segmented cells were generated using <https://github.com/AllenCellModeling/actk> as a set of
1271 contact sheets and all cells in the final dataset were manually reviewed for basic quality criteria such
1272 as only one nucleus per cell except later in mitosis, no obviously chopped nuclei, and no especially
1273 aberrant cell shapes due to segmentation errors.

1274

1275 ***Cellular structure segmentation***

1276 We applied a collection of modular segmentation workflows from the Classic Segmentation
1277 component of the Segmenter, each optimized for the particular morphological features of the target
1278 cellular structures (Chen et al., 2018). Representative examples for each of the 25 FP-tagged cellular
1279 structures are shown in **Figure S1**. For each structure, results of the segmentation workflow were
1280 evaluated on sets of images representing the variation observed across imaged cells (e.g. different
1281 regions of colonies) to ensure consistent segmentation quality across all images for each structure.
1282 The algorithms for all but two of the 25 cellular structures were classic image segmentation workflows.
1283 The exceptions were the plasma membrane (via CAAX) and the nuclear envelope (via lamin B1). All
1284 Classic Segmentation workflows contain three parts: pre-processing, core segmentation algorithms,
1285 and post-processing. For each part, there exists a set of algorithms to choose from, each with
1286 restricted numbers of parameters to tune. All workflows are accessible at
1287 <https://github.com/AllenCell/aics-segmentation>. For the plasma membrane, a deep learning-based
1288 segmentation model was developed as part of the training assay for cell segmentation described
1289 above. For the nuclear envelope, we developed an algorithm that combines multiple deep learning
1290 models including (1) the lamin B1 filled segmentation model we developed for nuclear segmentation
1291 training assay, (2) an overall lamin B1 segmentation model, (3) a lamin B1 seeding model, and (4) the
1292 plasma membrane segmentation model developed for cell and nuclear segmentation model. Briefly,
1293 we used the plasma membrane segmentation model output to cut the lamin B1 seeding model
1294 outputs to generate one seed per interphase nucleus. Then, the seeds were applied on the overall
1295 lamin B1 segmentation model via seeded watershed to obtain a one-voxel thick “shell” for each
1296 interphase nucleus. The “shells” were merged with the overall lamin B1 segmentation as the final
1297 lamin B1 segmentation result, which contained both complete nuclear envelope and properly segment
1298 invaginations and lamin B1 during mitosis. More details can be found in (Chen et al., 2018). Both
1299 models and code for CAAX and lamin B1 can be accessed via
1300 https://github.com/AllenCell/segmenter_model_zoo.

1301 We performed an additional validation step to determine whether a given target structure
1302 segmentation was sufficient for interpretation in the cellular structure volume analysis (**Figure 6**). We
1303 identified ten structures for which there were obvious caveats to the ability to use their target structure
1304 segmentation for biological interpretations of how much of the target structure was present in each
1305 cell and thus these ten structures were excluded from the structure volume analysis (**Figure S1B-D**).

1306 These three types of caveats included: (1) The cell boundary segmentation may have potential
1307 segmentation errors in the very top slices of the cell. This type of error has a minor effect on the
1308 overall segmentation of the cell but for structures localizing to the cell periphery at the very top of
1309 cells, this caveat can cause structures to be miss-assigned to neighboring cells (including tight
1310 junctions (ZO-1), gap junctions (connexin-43), desmosomes (desmoplakin; **Figure S1B**), adherens
1311 junctions (beta-catenin), actin filaments (beta-actin), actin bundles (alpha-actinin-1), and actomyosin
1312 bundles (non-muscle myosin IIB)). Therefore, these seven structures were not validated for cellular
1313 structure volume analyses. (2) Structures localizing or partially localizing to a thin 3D surface (such as
1314 the cell or nuclear periphery), especially when that surface is slanted, may suffer from non-uniform
1315 accuracy between the middle and the top/bottom of that structure due to the anisotropic resolution of
1316 the images. The accuracy of the nuclear pores target segmentation was sufficient to identify the
1317 general location of nuclear pores in the cell for the location-based analyses but not sufficient to be
1318 validated for use in the cellular structure volume analysis and thus this structure was excluded
1319 (**Figure S1C**). This nuclear periphery caveat was also observed for perinuclear ER (both Sec61 beta
1320 and SERCA2) and the nuclear lamina enriched localization of histones (H2B). However, these
1321 structures were still well segmented for the cytoplasmic ER localized throughout the cell and for
1322 histones localized throughout the nucleoplasm, each of which contributed more to overall structure
1323 volume. Therefore, those structures were not excluded from the cellular structure volume analysis.
1324 This caveat was also observed for structures that localize to the cell periphery (listed in the first
1325 caveat), which were excluded from the structure volume analysis. (3) The segmentation result for
1326 cohesins (via SMC1A) can depend on how far along a cell is in interphase and works well for most,
1327 but not all, of interphase (**Figure S1D**). Therefore, this structure was excluded from the structure
1328 volume analysis. Matrix adhesions (paxillin) localized to the very bottom of the cells where the
1329 membrane dye signal does not permit accurate identification of cell boundaries (see “Cell and nuclear
1330 segmentation” in Methods). Therefore, due to high likelihood of misassignment of matrix adhesions to
1331 neighboring cells, they were excluded from the structure volume analysis.

1332

1333

1334 **Single cell dataset generation**

1335

1336 ***Single cell image generation***

1337 To build the single-cell version of the image dataset for downstream analysis, we extracted all
1338 complete individual cells in each FOV automatically from the cell segmentation results of the image,
1339 ignoring any cells that were not at least 4 pixels away from the image border in the xy-plane (~12
1340 complete cells per FOV, on average). All images were rescaled to isotropic voxel sizes by
1341 interpolating along the z dimension to upscale the voxel size from 0.108333 μm x 0.108333 μm x 0.29

1342 μm to $0.108333 \mu\text{m} \times 0.108333 \mu\text{m} \times 0.10833 \mu\text{m}$. For each cell, a cropping region of interest (ROI)
1343 was calculated by extending the 3D bounding box of the cell by 40 voxels in each direction in both x
1344 and y and by 10 voxels in each direction in z.

1345 This same cropping ROI was applied to the original intensity z-stacks to extract the DNA,
1346 membrane and tagged structure for each cell. Similarly, the cropping ROI was used to extract the cell,
1347 nuclear and structure segmentations for each cell within this ROI. These extracted segmentations
1348 were then each masked by the cell segmentation result such that all voxels outside of the segmented
1349 boundary of the cell was set to zero. A roof-augmented version of the cell segmentation was also
1350 calculated for each cell to ensure proper inclusion of structures within the cell due to limited resolution
1351 and accuracy near the top of the cells (see “Single cell basic feature extraction” section). The roof-
1352 augmented cell segmentation is created by applying a morphological dilation (voxels only along the z-
1353 axis) at the top 25% of the cell segmentation mask. Each individual cell is thus associated with five
1354 segmentations: DNA segmentation, cell segmentation, roof-augmented cell segmentation, structure
1355 segmentation, and roof-augmented structure segmentation, which is masked by the roof-augmented
1356 cell segmentation after ROI cropping.

1357

1358 ***FOV-based feature extraction***

1359 FOV-based features were calculated for each cell. Specifically, we calculated (1) the
1360 Euclidean distance from the nucleus of each cell to the nucleus of each complete neighboring cell
1361 within the FOV, (2) the lowest and highest z position of all cells in this FOV, and (3) whether a cell is
1362 located on edge of a colony, for those cells within colony edge FOVs (mode C edge; see image
1363 acquisition methods). All details are released via <https://github.com/AllenCell/cvapipe>.

1364

1365 ***Colony-based feature extraction***

1366 In addition to FOV-based and single-cell-based features, we extracted colony-based features.
1367 For each 12X overview image and 120X FOV image taken, we extracted the name of the well in the
1368 96-well plate and the stage coordinates at which the image was taken from the file metadata. To
1369 obtain colony segmentations from the 12X overview images, we applied the same segmentation
1370 method used for automated position selection (see “Automated well position selection” section
1371 above). We then associated each segmented colony with a set of colony or well features including the
1372 confluency of the well, the size of the colony, the centroid location of the colony in the overview
1373 image, and whether the colony was touching the boundary of the overview image. We mapped the
1374 position where the 120X FOV image was taken relative to the 12X overview image by using the
1375 microscope stage coordinates, identified the colony in which that 120X FOV image was taken and
1376 added colony features to this 120X FOV image. We also calculated the Euclidean distance between
1377 the center of the FOV image and the nearest edge location of the colony. We added QC methods to

1378 ensure data accuracy and usability by flagging 120X FOVs with: (1) poor colony segmentations,
1379 detected as well confluency less than 10%, (2) 120X FOV images that were taken outside of the 12X
1380 overview image FOV and (3) 120X FOV images that were in a colony touching the edge of the 12X
1381 overview image. These colony-based features were not only linked to each 120X FOV but also to all
1382 of the individual cells associated with that FOV.

1383

1384 ***Deep learning based single cell annotation***

1385 Each cell in the hiPSC Single-Cell Image dataset was automatically annotated by a deep
1386 learning-based classifier into one of the following 7 annotation categories: interphase, prophase, early
1387 prometaphase, prometaphase/metaphase, anaphase/telophase (unpaired cell), anaphase/telophase
1388 (paired cell) or other (e.g. failed segmentations, dead cell segments, or dye blobs). Note: unpaired
1389 cells in anaphase/telophase refer to cells where it was impossible to find the other member of the pair
1390 (e.g. the other pair member is outside of the FOV). The automated classifier is a combination of a
1391 rule-based classifier and an ensemble of three 9-class 3D ResNet50 models. First, a cell is annotated
1392 as category anaphase/telophase (pair) if the nuclear segmentation satisfies the following three
1393 criteria: (1) contains at least two connected components, (2) the ratio of the sizes of the largest two
1394 connected components is greater than 0.64, an empirically determined value, (3) the distance
1395 between the centroid of the largest two connected components is greater than 85 voxels. Otherwise,
1396 the 9-class ResNet50 models are used. To train the ResNet50 models, we created a training set
1397 consisting of 5,664 cells from the main dataset and through expert-annotation assigned these into 9
1398 classes: 1-interphase, 2-prophase, 3-early prometaphase, 4-prometaphase/metaphase, 5-
1399 anaphase/telophase unpaired, 6-anaphase/telophase paired (but not necessarily satisfying all three
1400 criteria), 7-failed segmentation, 8-dead cell segments and 9-dye blobs. Class 1 (interphase)
1401 accounted for 43.5% of the data to ensure a balanced training set, while the total of classes 7, 8, and
1402 9 accounted for 2.9%. Three ResNet50 models were trained with different training/validation splits. An
1403 ensemble of these three models was used to make the final class predictions. These ResNet50
1404 models were validated by testing on 100 cells that were held out from the training set. The model
1405 generated eight incorrect predictions, but all were either incorrectly predicting mitotic stages (3/100) or
1406 incorrectly predicting a cell in interphase to be in mitosis (5/100). The recall rate for interphase cells
1407 was 100%. Cells that were predicted to be of classes 7, 8 or 9, or that generated prediction of low
1408 confidence, were annotated as belonging to the “other” category and removed from the hiPSC Single-
1409 Cell Image dataset. The confidence score of a prediction was approximated as the highest probability
1410 among all 9 classes. Confidence scores lower than 0.677 were considered low confidence and these
1411 cells removed. The final automated 3D image classifier code (for both training and testing) and all
1412 trained models are available at https://github.com/AllenCell/image_classifier_3d.

1413

1414 **Single cell basic feature extraction**

1415 Methods described here are implemented in several repositories for different parts of the
1416 analysis and begin with the hiPSC Single-Cell Image Dataset (see Data and code availability above).
1417 The data table downloaded from Quilt contains 216,062 rows and 47 columns, where each row
1418 corresponds to a single cell uniquely identified by its ID that is specified by the column CellId. The
1419 columns contain both the necessary information about each cell (e.g., path to segmentations, path to
1420 images, and important meta data) to calculate single cell features, as well as the calculated features.
1421 These features are included ready to download for ease of use. Demos in
1422 https://github.com/AllenCell/cvapipe_figure_notebooks can be used to produce the main figures from
1423 these released features, while all these features can be reproduced with the released code (see Data
1424 and Code availability above, (McKinney, 2011; Nicholas Sofroniew et al., 2019; Paszke et al., 2019;
1425 Pedregosa et al.; Walt et al., 2014)). The cell segmentation, the DNA segmentation and the roof-
1426 augmented structure segmentation were used to extract basic cell, nuclear and cellular structure
1427 features, respectively. The analysis dataset (see just below) contains only interphase cells, so the
1428 DNA segmentation represents the nucleus. The cell or nuclear segmentation for each cell is used as
1429 the input for calculating the following basic cell and nuclear features, respectively: (i) cell or nuclear
1430 volume as the number of non-zero voxels in the input image. The single-voxel volume $(0.108 \mu\text{m})^3$
1431 was used to rescale this feature for further analysis. (ii) cell or surface area as the number of voxel
1432 sides facing the background in the input image. According to this metric, an isolated voxel has all its 6
1433 sides facing the background and therefore a surface area equal to 6. The single-voxel-side area of
1434 $(0.108 \mu\text{m})^2$ was used to rescale this feature for further analysis (iii) cell or nuclear height as the
1435 distance in voxels along the z-axis between the bottom-most and top-most voxels in the input image.
1436 The single-voxel height of $0.108 \mu\text{m}$ was used to rescale this feature for further analysis. To calculate
1437 the volume of each cellular structure within a cell, the roof-augmented structure segmentation of that
1438 cell was used as the input. This ensures proper inclusion of structures within the cell due to limited
1439 resolution and accuracy near the top of the cells (see “Single cell image generation” section). The
1440 volume of the cellular structure in the cell is calculated as the number of non-zero voxels in the input
1441 image. We used the single-voxel volume $(0.108 \mu\text{m})^3$ to scale this feature for further analysis. These
1442 single cell basic features were merged into the hiPSC Single-Cell Image Dataset as additional
1443 columns and used in subsequent quantification and analysis.

1444

1445 **QUANTIFICATION AND STATISTICAL ANALYSIS**

1446

1447 **Analysis dataset generation**

1448

1449

1450 ***Mitotic cells removal***

1451 The first operation performed on the full dataset to create the analysis dataset was the
1452 removal of all of the 11,238 mitotic cells. This is done by removing all rows of the data table for which
1453 the column cell_stage is different from M0, the value used to flag interphase cells, resulting in a table
1454 with 204,824 rows (cells).

1455

1456 ***Outlier detection***

1457 In total, 1,087 (~0.5%) cells were identified and removed from the dataset, resulting in a table
1458 with 203,737 rows that we refer to as the analysis dataset throughout the paper. These outliers fall
1459 into two classes. First, there were 670 cells for which the structure volume was zero. Cells with an
1460 empty structure segmentation could be real outliers (e.g., no FP signals within that specific cell) or
1461 could indicate errors in either structure segmentation or cell and nuclear segmentation (see caveats in
1462 the “Structure Segmentation” section). Since cells with zero structure segmentation only account for
1463 ~0.3% of the whole population, we considered all such cells with potential segmentation errors, even
1464 minor, in cell and/or nuclear shapes as outliers. Second, we identified 417 cells, which were identified
1465 as outliers by an automated bi-variate outlier detection algorithm. Here, “bi-variate” refers to the notion
1466 that we looked at pairs of two variables to detect outliers and not at a single variable. As an example,
1467 the outlier detection procedure identified cells as outliers that have a very large cell volume (first
1468 variable) but very small nuclei (second variable), and clearly fall outside of the typical distribution of
1469 cell and nuclear volume. The outlier detection algorithm uses Gaussian kernel density estimation on
1470 the 2D space spanned by two variables, thereby assigning a probability to each of the cells. We use
1471 density estimation in the same way for visualization of bi-variate associations in scatter plots (see
1472 “Visualization of bi-variate association” section below and **Figure 6**). Cells with an extremely low
1473 probability were identified as outliers. We applied this outlier detection to the 21 pairs of variables that
1474 can be made of the seven main cellular and nuclear metrics: cell volume (μm^3), cell surface area
1475 (μm^2), cell height (μm), nuclear volume (μm^3), nuclear surface area (μm^2), nuclear height (μm),
1476 cytoplasmic volume (μm^3). Cells with resultant probabilities smaller than $1\text{e-}20$ were identified as
1477 outliers ($n=177$). This outlier analysis was also applied to pairs of variables for the four following cell
1478 and nuclear metrics (cell volume and surface area, nuclear volume and surface area) each with
1479 cellular structure volume for the 15 structures validated for structural volume analysis, totaling
1480 $15 \times 4 = 60$ scatter plots. Cells with a probability smaller than $1\text{e-}10$ in any of these 60 scenarios were
1481 identified as outliers ($n=240$). The thresholds mentioned above were identified manually after
1482 inspection of the scatter plots and visual inspection of many cells identified as outliers. The majority of
1483 the inspected cells clearly showed imaging or segmentation artifacts.

1484

1485

1486 **Statistical analysis for quality control of the hiPSC Single-Cell Image Dataset**

1487 To be able to map cells from the 25 cell lines into the same shape space and cluster similar
1488 cells to integrate the location of their separately imaged structures we must first ensure that the cell
1489 lines themselves are not an experimental source of cell and nuclear shape variation. Further, this
1490 extensive dataset was acquired over a period of three years, including changes in the extent of
1491 pipeline automation, necessary adjustments to the microscopes, the lots of Matrigel, and other such
1492 experimental factors over the course of the imaging pipeline timeline (see “Imaging workflows”
1493 section). Therefore, we performed an extensive analysis to identify and account for any potential
1494 experimental contributions to cell shape variation (**Figure S7**). An analysis of how each of the Shape
1495 Modes varied with respect to the timeline of the imaging pipeline revealed that only Shape Modes 1
1496 and 2, representative of cell height and cell volume, showed any signs of possible systematic
1497 experimental variation (**Figure S7A**). For cell height, we observed variation between cell lines
1498 throughout the pipeline timeline, while for cell volume we only observed a possible systematic
1499 difference between Pipeline 4.4 and the rest of the pipeline workflows (**Figure S7B&C**). The greatest
1500 systematic effect on cell height over the pipeline timeline was visible in the sequential imaging of the
1501 last two structures (nuclear speckles via SON and cohesins via SMC1A), which both contained flatter
1502 cells. These differences were attributable to a change in both the lot of Matrigel and an adjustment to
1503 the glass bottom well-plate Matrigel coating protocol as described above. This can be seen in a
1504 control experiment comparing the tagged actomyosin bundles (via non-muscle myosin IIB) cell line
1505 before and after this protocol change (**Figure S7B**). We separated the pipeline timeline into three
1506 periods, the period before Pipeline 4.4, and then within Pipeline 4.4, the period before and after the
1507 change in Matrigel coating protocol and compared both cell height and cell volume between these
1508 periods. We found that while the adjusted Matrigel coating protocol decreased cell height significantly,
1509 it did not affect cell volume. However, both cell height and cell volume were slightly and consistently
1510 decreased during the entire Pipeline 4.4. Further investigation into possible causes revealed a
1511 systematic inaccuracy in z spacing due to the use of a piezo z stage, which leads to an approximate
1512 10% reduction in the z-step size and thus also in the overall height of the cell. When we corrected the
1513 Pipeline 4.4 z-step size by this approximate amount, we found this could account for the cell height
1514 difference. Cell volumes cannot be directly corrected by one single factor adjustment due to the varied
1515 cell shapes. However, the slight, yet significant and consistent decrease in average volumes of all cell
1516 lines imaged during Pipeline 4.4 can be accounted for by the same piezo-dependent problem.
1517 Unfortunately, we could not retroactively determine the exact adjustment to the z-step size for each
1518 independent image acquisition that was performed during Pipeline 4.4 and thus did not correct the
1519 data for this issue. However, the magnitude of the effect was much smaller than the variation of cell
1520 volumes and heights within the cell line datasets.

1521 In addition to these two systematic experimental sources of variation during Pipeline 4.4, we
1522 observed variation in average cell height throughout the entire pipeline timeline. This suggested
1523 additional possible sources of variation. We had experimentally observed that cell height seemed to
1524 vary both with colony area and the location of cells within a colony (**Figure 1A**), suggesting that cell
1525 height variation might be part of normal changes to cell packing behavior within a growing colony. To
1526 test this observation quantitatively, we measured the cell area of a subset of colonies with accurate
1527 colony segmentations as well as both the distance from the center of the FOV to the edge of the
1528 colony and the average height of all the cells within that FOV. We transformed colonies and the
1529 locations of FOVs within them into circular representations and compared the location patterns, cell
1530 heights, and colony areas (**Figure S7D**). We found that smaller colonies tended to contain taller cells
1531 while in larger colonies, cells closer to the colony periphery were taller than those towards the center
1532 of colonies. Other than Shape Mode 1, representing cell height, none of the other shape modes
1533 showed any colony-specific patterns within the dataset (**Figure S7E**).

1534 We next investigated how much of the variation in cell height (median height of the cells in an
1535 FOV) was explained by a set of eleven experimental variables including the distance of an FOV to the
1536 colony edge representing the position of cells in a colony, the colony area, the cell line identity, and
1537 several imaging pipeline settings (**Figure S7F**). We performed a Random Forest regression analysis
1538 (Liaw and Wiener, 2002) and found we could predict cell height with moderate accuracy ($R^2 = 0.52$)
1539 based on this combination of eleven variables. When we removed cell line identity as a variable within
1540 this regression analysis, the accuracy of cell height prediction barely change ($R^2 = 0.51$). The feature
1541 “FOV to colony edge distance” had the largest feature importance. Importantly, we found that cell line
1542 identity was statistically correlated with several imaging pipeline settings that varied throughout the
1543 imaging pipeline timeline.. All of the results above together confirm that cell line identity can contribute
1544 to cell height variation due to the fact that each cell line was imaged under a particular set of imaging
1545 conditions which varied throughout the imaging pipeline timeline, but that cell line identity itself does
1546 not greatly contribute to the variation in cell height observed in the hiPSC Single-Cell Image Dataset.

1547

1548 ***Circular colony mapping***

1549 We took advantage of the fact that many cells ($n=104,269$) of our hiPSC Single-Cell Image
1550 Dataset could be associated with information relative to the colony where they came from (see
1551 “Colony-based feature extraction” section), to visualize radially dependent spatial patterns of our cells.
1552 This is achieved by mapping the location of cells in a colony into a unit circle, as illustrated in **Figure**
1553 **S7D**. First, the distance from the center of the FOV to the closest edge point (d) is normalized by the
1554 effective radius of the colony (R_{eff}) to determine the relative distance $\ell=d/R_{\text{eff}}$. Then, all cells in the
1555 FOV are mapped into a unit circle at radial distance ℓ from the edge of the circle. Each cell is assigned
1556 to an angular location drawn from a uniform distribution of angles in the range $[0,2\pi]$.

1557 **Random forest regression model to predict cell height from experimental features**

1558 A multivariate Random Forest regression model was trained to predict the median cell height
1559 of all cells in an FOV from experimental, assay-dependent variables, including (1) cell growth
1560 information from the confluency of cells in the well, Matrigel-coating protocol, the FP-tagged protein
1561 name, and two cell passaging numbers, (2) colony features from the size of the colony the FOV was
1562 imaged at and the distance between the FOV and the nearest colony edge, and (3) instrument
1563 hardware configurations including the pipeline workflow information, the ID of the microscope which
1564 the FOV was taken with and the piezo configuration of the microscope. We first calculated the median
1565 cell height of an FOV from the single-cell segmentation that provides the height of each cell in the
1566 FOV. We then pre-processed the continuous variables (FOV to colony edge distance, confluency,
1567 colony area, total passages and passages post-thaw) with z-normalization, and labeled categorical
1568 variables (cell line via its FP-tagged protein name, imaging mode, workflow ID, Matrigel protocol,
1569 piezo setting, microscope ID) in R Studio. We added a control variable by randomly generating a
1570 number that ranges from -3 to 3 for each FOV.

1571 This analysis was based on the 20 cell lines that contain >100 FOVs each, with a total of 7,914
1572 FOVs. The cell lines with the following tagged proteins (**Table 1**) were included: alpha-actinin-1,
1573 alpha-tubulin, beta-actin, CAAX, centrin-2, connexin-43, desmoplakin, fibrillarin, H2B, lamin B1,
1574 LAMP-1, non-muscle myosin IIB, Nup153, paxillin, Sec61 beta, sialyltransferase 1, SMC-1A, SON,
1575 Tom20, and ZO-1. For each cell line, we randomly selected 90 FOVs for training, resulting in a
1576 training dataset of $90 \times 20 = 1,800$ FOVs and used the remainder of the FOVs ($n = 6,114$) to evaluate
1577 the model. We trained a Random Forest model with using all variables in R Studio with the
1578 RandomForest package (Liaw and Wiener, 2002) with 500 trees. We also trained another Random
1579 Forest model with all variables except cell line identify, again using 500 trees. We evaluated the
1580 model by calculating the Coefficient of Determination (R^2) on the test set ($n = 6,114$ FOVs). Feature
1581 importance scores were calculated as the difference in mean squared error (MSE) between a model
1582 including the feature in question and a model where the values of that feature were randomly
1583 permuted across the samples. We repeated the sampling and model training 100 times to obtain
1584 confidence intervals of model performance and feature importance as shown in **FigureS7F** (left).

1585

1586

1587 **Spherical harmonics expansion (SHE) of cell and nuclear shapes**

1588 In addition to the basic features described above, we also used SHE coefficients as shape
1589 descriptors for cell and nuclear shape (Ruan and Murphy, 2019; Shen et al., 2009). We created a
1590 publicly available open-source Python package, aics-shparam (see “Data and Code Availability”
1591 section) to extract SHE coefficients from segmented images of cells and nuclei.

1592

1593 ***Cell and/or nuclear alignment***

1594 SHE coefficients are sensitive to the orientation of the shape they are extracted from.
1595 Therefore, a given set of cells and nuclei can be used to create different versions of a shape space,
1596 depending on how they are pre-aligned. To create the cell and nuclear joint shape space (**Figure 2**)
1597 we wanted to preserve the apical basal axis of the cell, which is the z-axis in the lab frame of
1598 reference. Therefore, we only aligned cells by rotation in the xy-plane. Cells and nuclei were rotated
1599 such that the longest cell axis falls along the x-axis. The cell segmentation was used to estimate the
1600 longest axis of the cell through a principal component analysis of the x and y coordinates of
1601 foreground voxels. The longest axis was defined as the direction of the first principal component and
1602 the alignment angle defined as the smallest angle between the longest axis and the x-axis. That cell
1603 was then rotated by the alignment angle such that the longest axis was aligned with the x-axis. Cells
1604 were rotated by using the function rotate from Python package scikit-image (Walt et al., 2014) with
1605 zero order interpolation. The input image was also resized as necessary to fit the whole rotated cell.
1606 The alignment procedure was implemented by the function align_image_2d in aics-shparam using
1607 default parameters. This function returns the final alignment angle, which is then used to align other
1608 images related to that cell, in this case the segmented images of the nucleus and the particular
1609 cellular structure in the cell as well as the three channels of the z-stack containing the original images
1610 of the membrane dye, DNA dye and FP-tagged structure. This was done using the function
1611 apply_image_alignment_2d available in the same Python package.

1612 1613 ***From segmented, aligned images to SHE coefficients and 3D meshes***

1614 Once a segmented image of a cell and nucleus is aligned, it is used as input for the function
1615 get_shcoeffs from aics-shparam. This function first converts the input binary image into a 3D
1616 triangular mesh using a traditional marching cubes algorithm from VTK Python library (Schroeder et
1617 al., 2018). To improve the quality of the output mesh, the binary input image is convolved with a
1618 Gaussian kernel with size $\sigma_x=\sigma_y=\sigma_z=2$, which is enough to smooth the image while retaining the
1619 overall cell and nuclear shape. Next, the mesh is translated to the origin and the coordinates of the
1620 mesh points are converted from cartesian to geographic coordinates (latitude, longitude and altitude).
1621 Altitude coordinates are then interpolated, using nearest neighbor, over a (lat,lon) spherical grid where
1622 each cell has a resolution of $\pi/128$. At this point, aics-shparam uses the Python package pyshtools
1623 (Wieczorek and Meschede, 2018) to expand, up to degree Lmax, the equally spaced grid into
1624 spherical harmonics coefficients using Driscoll and Healy's sampling theorem (Driscoll et al., 1994).
1625 We used Lmax=16 as the SHE degree expansion to parameterize both cell and nuclear segmentation
1626 images. This was enough to guarantee a high fidelity mesh reconstruction, which can be quantified by
1627 the average distance between closest points in the original and reconstructed 3D meshes. We
1628 observed average distances of $0.33 \mu\text{m} \pm 0.1 \mu\text{m}$ for cells (n=300 randomly selected samples) and

1629 0.12 μm +/- 0.02 μm for nucleus (n=300 randomly selected samples). Compared to the voxel size of
1630 our images (0.108 μm), we can say that $L_{\text{max}}=16$ yields single pixel level precision for the nucleus
1631 and about three voxels precision for the cell, in average. This degree of expansion results in 289
1632 coefficients for each input. Therefore, the shape of each cell in our dataset can be represented by a
1633 total of 578 coefficients (**Figure 2A**).

1634 We can also recreate the 3D mesh representation of a particular set of SHE coefficients with
1635 aics-shparam. The Driscoll and Healy's sampling theorem allows one to obtain a spherical grid from
1636 pre-computed SHE coefficients. These points on the spherical grid can be radially translated to their
1637 actual values in the grid to give rise to a 3D non-spherical shape.

1638

1639

1640 **Building the cell and nuclear shape space**

1641

1642 ***Principal component analysis for dimensionality reduction***

1643 We used principal component analysis (PCA) to reduce the dimensionality of our joint vectors
1644 for all cells (578 SHE coefficients) down to eight principal components. We used the PCA
1645 implementation from the Python library scikit-learn (Pedregosa et al.) with default parameters (**Figure**
1646 **2B**). Since the sign of a given principal component (PC) is arbitrary, we flipped the sign to ensure that
1647 the average volume of cells with negative PC values was less than that of cells with positive PC
1648 values. This was done independently for each PC.

1649

1650 ***Identifying the primary modes of shape variation***

1651 To prevent cells with extreme shapes from affecting the interpretation of the PCs, we excluded
1652 all cells that fell into the range 0th to 1st or 99th to 100th percentiles of each PC from subsequent
1653 analysis. These percentile ranges are shown by the vertical red lines in **Figure S3C**. The total number
1654 of cells left in the dataset was 175,935. We z-scored all PCs independently by dividing the PC values
1655 by the standard deviation (σ) of that PC. The probability distribution of each z-scored PC is shown in
1656 **Figure S3C**. The z-scored principal components are referred to as "shape modes" and the
1657 combination of the first 8 shape modes creates the 8-dimensional generative shape space used
1658 throughout this paper. We used the inverse of the PCA transform generated above to map shapes
1659 from the shape space back into SHE coefficients, which in turn, can be used to reconstruct the
1660 corresponding 3D shape. For example, the 8-components vector (0,0,0,0,0,0,0,0) represents the
1661 origin of the shape space and its corresponding 3D shape is called the mean cell shape throughout
1662 the paper (**Figure 2C**).

1663 To systematically explore the shape space along each of the eight orthogonal axes, we let the
1664 elements of the 8-component array vary, one at the time, over discrete map points with values -2σ , -

1665 1.5σ , -1.0σ , -0.5σ , 0 , 0.5σ , 1.0σ , 1.5σ and 2.0σ . The combination of all eight shape modes and nine
1666 map points generates a grid of 8×9 3D shapes. Three different 2D views are used to visualize the 3D
1667 shapes. Top views represent the intersection of the 3D reconstructed mesh with the xy -plane, the
1668 equivalent of a single xy -slice through the center of the cell. In the same way, side views 1 and 2
1669 represent the intersection of the 3D reconstructed shape with the xz - and yz -plane. To easily assign
1670 real cells to map points in the shape space, each shape mode is binned into nine bins of width 0.5σ ,
1671 each centered around one map point, as represented by the black vertical lines in **Figure S3C**.

1672 Cell and nuclear 3D mesh reconstructions using the inverse PCA transform are centered at
1673 the origin. Therefore, a few extra steps are required to translate the nuclear mesh back to its correct
1674 location relative to the center of the cell. We average all of the nuclear locations relative to their cell
1675 center for all the real cells within particular shape mode bin (**Figure S3C**). For example, to correct the
1676 nuclear location of the 3D mesh corresponding to the 8-components vector $(0,0,0,0,0,-1.5\sigma,0,0)$ of
1677 Shape Mode 6 (**Figure 2C**), one would use the average location of all real cells that fall into the bin
1678 highlighted in blue in **Figure S3C**. Both the cell meshes and nuclear meshes with corrected locations
1679 for all shape modes are saved in VTK polydata format (Schroeder et al., 2018) for further analysis.

1680

1681 ***Alternative versions of the shape space***

1682 In addition to the joint cell and nuclear shape space, we also generated independent cell-only
1683 and nucleus-only shape spaces. For the cell-only shape space, the PCA was applied only on the cell
1684 SHE coefficients to reduce the data dimensionality from 289 to 8. For the nucleus-only shape space,
1685 images of DNA segmentation were aligned independently from any cell information. Nuclei were
1686 rotated such that the longest nuclear axis fell along the x -axis. The DNA segmentation was used to
1687 estimate the longest axis of a nucleus through a principal component analysis of the x and y
1688 coordinates of foreground voxels. The longest axis was defined as the direction of the first principal
1689 component and the alignment angle defined as the smallest angle between the longest axis and the x -
1690 axis. That nucleus was then rotated by the alignment angle such that the longest axis was aligned
1691 with the x -axis. Aligned images of nuclei were used as input for SHE coefficients calculation. PCA was
1692 applied only on the nuclear SHE coefficients to reduce the data dimensionality from 289 down to 8.
1693 After dimensionality reduction through PCA, these two alternative shape spaces were analyzed
1694 identically to the joint cell and nuclear shape space to identify the main modes of shape variation
1695 shown in **Figure SB&C**.

1696

1697

1698 **SHE coefficient-based parameterization and 3D morphing to build integrated average cells**

1699

1700

1701 ***Cytoplasmic and nuclear mapping***

1702 Pre-computed SHE coefficients were interpolated to morph the nuclear centroid mesh into the
1703 nuclear surface mesh and the nuclear surface mesh into the cell surface mesh. First, the nuclear
1704 centroid of each cell is described by the SHE coefficients representing a one-pixel radius (0.108 μm)
1705 3D spherical mesh. These SHE coefficients representing the nuclear centroid and the pre-computed
1706 cell and nuclear SHE coefficients are concatenated and computationally described by a 3x289 matrix.
1707 This matrix is linearly interpolated to generate a 64x289 matrix. The interpolation is done by the
1708 function `interp1d` from `scikit-learn` in such a way that it guarantees that 1st, 3-th and 64th rows of the
1709 output matrix correspond exactly to SHE coefficients of centroid, nuclear and cell. SHE coefficients of
1710 each row of the interpolated matrix can be used to reconstruct corresponding 3D meshes. Meshes
1711 corresponding to rows 32 to 64 in the interpolated matrix are translated to a location that corresponds
1712 to a linear interpolation between nucleus and cell centroid. The visualization of subsequent 3D
1713 meshes (subsequent rows) causes the effect of mesh interpolation, as shown in **Figure 3A**, where we
1714 show only eight out of the 64 possible meshes (differently colored regions), including centroid (black
1715 dot) nuclear and cell meshes (represented by dashed lines).

1716

1717 ***Parameterized Intensity representation***

1718 Each of the 3D meshes is composed of points with xyz-coordinates. As the meshes are being
1719 generated from the interpolated matrix point by point, we can visit the corresponding xyz location in
1720 the aligned images that were used to generate the cell and nuclear SHE coefficients in the first place
1721 and associate the intensity value of that location with the mesh xyz coordinate. We can record either
1722 the original intensity values or the segmented intensity values since both types of images were
1723 aligned. The results can be organized as a matrix as shown in **Figure 3A** for the original FP signal.
1724 This matrix encodes a parameterized intensity representation of the cell.

1725 This parameterized intensity representation can be used to reconstruct the aligned image that
1726 was used as the original input. We start with an empty image. We assign the value of each element of
1727 the parameterized intensity matrix to its closest xyz location in the empty image. We call this
1728 procedure voxelization and it produces a sparse representation of the original aligned image as
1729 shown in **Figure 3A**. The gaps in this image are due to the fact that our parameterized intensity
1730 representation samples only as many voxels of the original image as we have points in the 3D mesh.
1731 The gaps can be filled in by a nearest neighbor interpolation to produce an image that looks very
1732 similar to the original aligned image, as shown at the top of **Figure 3A**. We used the function
1733 `NearestNDInterpolator` from `scikit-learn` to perform the multidimensional nearest neighbor
1734 interpolation. We used the voxel-wise Pearson correlation coefficient in 3D to evaluate the similarity
1735 between reconstructed and original aligned images. We also performed an analysis between
1736 reconstructed and original aligned images on 32 randomly selected cells of all 25 cellular structures

1737 when the parameterized intensity representation is used to encode either original FP or segmented
1738 intensities (**Figure S3A&B**).

1739

1740 ***Generating morphed cells***

1741 The cellular mapping procedure described above only requires cell and nuclear SHE
1742 coefficients. Therefore, it can be applied to cell and nuclear shapes obtained for all map points of all
1743 shape modes. This is illustrated in **Figure 3B** for map point $(0,0,1.5\sigma,0,0,0,0)$ of Shape Mode 3. The
1744 parameterized intensity representation of any given real cell can now be voxelized into any map point
1745 shape that underwent cellular mapping, to generate a morphed version of the real cell into that shape.
1746 This is illustrated in **Figure 3** by morphing the FP signal from the real cell shown in panel (A) into the
1747 shape of map point $(0,0,1.5\sigma,0,0,0,0)$ of Shape Mode 3 shown in panel (B). To prevent morphed
1748 cells from containing overly distorted signal intensity locations compared to the real cells, for instance
1749 by morphing a very flat real cell into a very tall shape (e.g. map point $(2\sigma,0,0,0,0,0,0)$ of Shape
1750 Mode 1), we restrict our ourselves to apply the morphing only when the real cell shape and the map
1751 point shape are similar. This is achieved throughout the paper by allowing only cells of a given map
1752 point bin (**Figure S3C**) to be morphed into the corresponding map point shape. For example, only
1753 cells that fall into the bin highlighted in blue in **Figure S3C** are allowed to be morphed into the shape
1754 corresponding to map point $(0,0,0,0,0,-1.5\sigma,0,0)$ of Shape Mode 6. The number of cells per structure
1755 available in each bin of each shape mode is shown in **Table S1**. We selected 300 randomly chosen
1756 cells (or the maximum number of cells available) per cellular structure and per shape mode and
1757 morphed these cells into their corresponding map point shapes. These morphed cells are stored as
1758 multichannel TIFF files and were used further for stereotypy analysis as described below.

1759

1760 ***Aggregating morphed cells***

1761 We compute the average and standard deviation of parameterized original FP and segmented
1762 intensity representations for all cells of each structure across map points of all shape modes of the
1763 shape space. This computation produces average and standard deviation parameterized intensity
1764 representations that could also be morphed into map point shapes of shape modes as described
1765 above. Results of these average and standard deviation images for all 25 cellular structures are
1766 shown as the first three columns of **Figure S3C** for map point $(0,0,1.5\sigma,0,0,0,0)$ of Shape Mode 3.
1767 To quantify the location variation of each cellular structure, we normalized the standard deviation
1768 images by the average images to create coefficient of variation images. To prevent areas with very
1769 low average values (effectively very low original FP or segmented intensities) from greatly impacting
1770 the coefficient of variation, we defined the structure-localized coefficient of variation. The structure-
1771 localized coefficient of variation is computed as the coefficient of variation limited to a set of voxels
1772 containing intensities above a set threshold. The threshold was chosen to be the median of all non-

1773 zero voxels in the average image. Structure-localized coefficients of variation for all 25 cellular
1774 structures are shown as the 4th column of **Figure S3C** for map point (0,0,1.5 σ ,0,0,0,0) of Shape
1775 Mode 3. Aggregated cells are saved as 5D hyperstacks and used for visualization and concordance
1776 analysis, as described below.

1777

1778 ***Visualizing integrated average morphed cells***

1779 Average images of cellular structures morphed into the same map point shape are rendered
1780 simultaneously to illustrate the spatial relationships of different structures based on their average
1781 location in cells of a particular shape. Each volumetric channel of the 5D hyperstacks generated in the
1782 previous section for Shape Mode 3 was segmented using the default Surface option found in the
1783 Volume Viewer window of ChimeraX (Pettersen et al., 2020). Thresholds for each channel were
1784 selected manually to clarify dominant localization patterns observed in the voxel intensities.

1785

1786

1787 **Stereotypy calculation from morphed cells**

1788 We used morphed cell images to quantify the location stereotypy of a given cellular structure
1789 across different cells with similar shape. All 300 morphed cells available for each shape mode map
1790 point for each cellular structure were used to generate unique pairs of images. We calculated the
1791 voxel-wise Pearson correlation between all pairs of images, as illustrated in **Figure 4A** for lamin B1
1792 (top) and mitochondria (bottom). The values of the resulting correlation coefficients represent a
1793 distribution of stereotypy values for each set of 300 cells. The distributions of stereotypy values for all
1794 25 cellular structures for the mean cell shape are represented by the box plots in **Figure 4B**. The
1795 mean of the distribution of stereotypy values is called the mean stereotypy. The mean stereotypy
1796 values calculated for all 25 cellular structures across map points of all shape modes are shown as
1797 heatmaps in **Figure 4C**. To highlight the difference between mean stereotypy values relative to the
1798 mean cell shape, we created difference heatmaps as shown in **Figure S4C**, where the mean
1799 stereotypy of the mean cell bin is subtracted from the mean stereotypy of other map points.

1800

1801

1802 **Concordance calculation from average morphed cells**

1803 We used the 5D hyperstacks of average morphed images generated as described above to
1804 quantify the location concordance of all 25 cellular structures. The average morphed image of each
1805 structure for a given map point of a particular shape mode was used to build a voxel-wise correlation
1806 matrix as shown in **Figure 5A**. The element (i,j) of this matrix gives the concordance between
1807 structures i and j. The 25x25 correlation matrix is used as input for a hierarchical clustering algorithm
1808 to cluster all 25 cellular structures according to their relative concordance. We used the function

1809 cluster.hierarchy.linkage of type “average” from the Python package scipy (Virtanen, 2020) to produce
1810 the clustering represented by the dendrogram in **Figure 5B** calculated for the mean cell.

1811 Concordance matrices were also calculated across map points for all shape modes. These
1812 matrices are represented by heatmaps across shape modes in **Figure 5C**, where the lower and upper
1813 triangle of each heatmap represent extreme opposite map points (see figure legend). To highlight the
1814 difference between concordance values relative to the mean cell, we create difference heatmaps as
1815 shown in **Figure S5B**.

1816

1817

1818 **Multiscale stereotypy and concordance analysis**

1819 Both stereotypy and concordance analysis were also performed across different spatial scales
1820 to investigate whether cellular structures display non-trivial behavior compared to what was observed
1821 in our initial analysis. Images for this analysis at different spatial scales were created by effectively
1822 downsampling the original images in all three dimensions by factors of 2 (**see Figure S4A**). The initial
1823 voxel-size of the morphed cell images was 0.108 μm . The downsampling process was repeated
1824 seven times to reach a voxel-size of $\sim 13.82 \mu\text{m}$.

1825

1826

1827 **Cellular structure size scaling analysis**

1828 Statistical associations between volumes and areas of cells, nuclei and 15 cellular structures
1829 show how strongly these metrics are coupled to each and how they scale with respect to each other.

1830

1831 ***Description of data used for cellular structure size scaling analysis***

1832 This statistical analysis uses six metrics: The cell volume (μm^3) and surface area (μm^2), the
1833 nuclear volume (μm^3) and surface area (μm^2), the cytoplasmic volume (μm^3), calculated by
1834 subtracting nuclear volume from cell volume, and the cellular structure volume (μm^3). The cell and
1835 nuclear metrics are available for all cells ($n=203,737$) and calculated based on the segmentation of
1836 the cell and nucleus, respectively. The cellular structure volume is based on the segmentation of the
1837 FP-tagged structure in the cell and is applied to the 15 cellular structures validated for structure
1838 volume analysis. (see “Structure segmentation” section; **Table S1**). If multiple pieces (connected
1839 components) of the structure are present in this cell, structure volume gives the total volume of all
1840 connected components.

1841

1842 ***Linear regression model to compute statistical coupling between metrics***

1843 We employed a simple linear regression model ($y = ax + b$) to compute the amount of
1844 explained variance in the dependent variable y by the independent variable x . Linear regression

1845 models were calculated with x as one of the five cell and nuclear metrics (cell volume and area,
1846 nuclear volume and area, cytoplasmic volume) and y as one of all six metrics (including cellular
1847 structure volume). In the case of structure volume, the model was computed for each structure
1848 separately, using only those cells that correspond to the structure in question. The explained variance
1849 in y due to x, or the R² statistic (coefficient of determination), was computed for all models and is
1850 shown in **Figure 6B**. We used a bootstrap analysis (n=100 bootstraps) to calculate the 5-95%
1851 confidence interval, visualized as horizontal error bars in **Figure 6H**.

1852

1853 ***Linear regression model to compute cellular structure scaling rates***

1854 Using the same simple linear regression model ($y = ax + b$), we calculated the “scaling rate” of
1855 each cellular structure relative to cell volume. The scaling rate gives the increase in volume (or area)
1856 of a cellular structure as cell size is doubled. In this case x is cell volume and y is one of the other five
1857 metrics. Using a histogram density estimation of cell volume, we determined the interval with the most
1858 cells where the cell volume doubles. This interval is from $x_0 = 1160 \mu\text{m}^3$ to $x_1 = 2320 \mu\text{m}^3$. These x
1859 values are then evaluated with the learned regression model to get the corresponding y values,
1860 termed y_0 and y_1 . The scaling rate is computed as $(y_1 - y_0) / y_0 * 100\%$. **Figure 6B** depicts this process
1861 to compute the scaling rate for nuclear volume. In this case y_0 is $346 \mu\text{m}^3$ and y_1 is $669 \mu\text{m}^3$, giving a
1862 scaling rate of 93%. The scaling rates across all metrics is given in **Figure 6A**. We used a bootstrap
1863 analysis (n=100 bootstraps) to calculate the 5-95% confidence interval, visualized as vertical error
1864 bars in **Figure 6H**.

1865

1866 ***Multivariate regression model to isolate the effect of cell and nuclear metrics in explaining*** 1867 ***structure volumes***

1868 Cell and nuclear metrics show a large degree of collinearity, which makes it non-trivial to
1869 isolate the effect of one particular cell or nuclear metric on structure volume. We used multivariate
1870 regression models to isolate the effect of cell and nuclear metrics. In contrast to univariate regression
1871 models ($y = ax + b$, where x is a vector and a is scalar), multivariate models have multiple dependent
1872 variables ($y = aX + b$, where X is a matrix with p columns and a is a vector with p entries). We first
1873 computed the explained variance in cellular structure volume using cell volume, cell surface area,
1874 nuclear volume and nuclear surface area as independent variables. Note that cytoplasmic volume is a
1875 linear combination of cell and nuclear volumes and does not need to be added to the model. Then, we
1876 remove a single metric or a pair of metrics from the independent variables and recalculate the model.
1877 The “unique explained variance” ascribed to the metric or pair of metrics is calculated as the
1878 difference in explained variance between the full model, i.e. containing the four metrics and the model
1879 where the metric or pair of metrics was left out. Specifically, the metrics (pairs) for which this unique
1880 explained variance was computed were cell volume and cell surface area (cell v+a), cell volume, cell

1881 surface area, nuclear volume and nuclear surface area (nuc v+a), nuclear volume, and nuclear
1882 surface area. The total explained variance (using all four metrics) as well as the unique explained
1883 variance portions are depicted in **Figure 6A**.

1884

1885 ***Non-linear regression models to compute statistical coupling between metrics***

1886 For each of the linear regression models described above, we also computed a more complex,
1887 non-linear model. Specifically, given a linear regression model $y = aXl + b$, where the design matrix Xl
1888 contains either a single vector or multiple columns, we expanded the design matrix Xl using two steps:
1889 1) for all pairs of columns in Xl , we computed the pointwise product and added these new columns to
1890 the design matrix; and 2) for each column in the design matrix we added four copies and raised the
1891 values of these new columns to the following four powers: $1/3$ (cube root), $1/2$ (square root), 2
1892 (square), 3 (cube). The resulting design matrix, Xc , was then used in the linear regression model $y =$
1893 $aXc + b$ to compute the explained variances. A visualization of the explained variances using simple
1894 regression models compared with the non-linear models with interaction effects is shown in **Figure**
1895 **S6B**.

1896

1897 ***Visualization of bi-variate association using scatter plots***

1898 Associations between pairs of metrics were visualized in scatter plots, where each cell is
1899 plotted as a point in the two-dimensional space spanned by the two metrics, x (on the x -axis) and y
1900 (on the y -axis). The number of cells is stated in the upper left corner. The regression model is
1901 depicted as a gray straight line ($y = ax + b$) and the explained variance in y due to x (the R^2 statistic)
1902 is also stated in the upper left corner. There are two additional graphical aspects to improve the
1903 interpretation of these bi-variate associations: 1) A green line is shown that depicts the running
1904 average. Briefly, the values of metric x are binned in 100 equally spaced bins. For each of these bins,
1905 the mean value for metric y is computed from all cells in that bin, i.e. unless the number of cells in the
1906 bin is below 50 in which case no value is recorded. The green line is the running average of metric y
1907 as a function of the bin centers. 2) Cells are colored according to a density estimate. Briefly, a kernel
1908 density estimate is performed in the two-dimensional space. Based on this estimation, each cell is
1909 assigned a probability. The probabilities are transformed to cumulative probabilities and normalized,
1910 such that the cell with the highest probability, i.e. the one within the highest density region, gets a
1911 value of 1. By aligning the probabilities with a colormap, cells are colored to convey the density. The
1912 use of cumulative probabilities ensures that the colors have the same interpretation across different
1913 plots, i.e. different metrics. See **Figure 6**.

1914

1915

1916

1917 **Generalizability analysis**

1918 To test the generalizability of our multi-part analysis approach of the locations, amounts, and
1919 their variation of the 25 cellular structures we ran the main analyses shown in this paper on subsets of
1920 the analysis dataset where we selected (by downsampling) a much smaller number of cells per
1921 structure.

1922

1923 ***Shape space generalizability***

1924 Using the main analysis dataset (n = 203,737 cells) including the 578 SHE coefficients, we
1925 randomly selected 300 cells and we applied PCA on the 300x578 table to reduce its dimensionality
1926 down to eight principal components. The resulting shape space is analyzed identically to the main
1927 shape space.

1928

1929 ***Stereotypy and concordance generalizability***

1930 We calculate the mean stereotypy of all 25 cellular structures morphed into the mean cell for
1931 different numbers of pairs of morphed cells. We varied the number of pairs of morphed cells used to
1932 average the Pearson correlation scores from 2 to 300 with a step size of one (**Figure 7B**). By visual
1933 inspection we determined the minimum number of pairs of cells required to recover the ranking of
1934 cellular structure mean stereotypy from 300 pairs of cells to be 35 pairs. The morphed cells used here
1935 were randomly sampled from the 300 morphed cells available per cellular structure per map point of
1936 each shape mode generated as described in the “Generating morphed cells” section above.

1937 For location concordance, we selected a set of 300 cells chosen at random for each cellular
1938 structure within the mean cell shape bin (except n= 252 for nuclear speckles, see **DataFile S1**). We
1939 used the 5D hyperstacks of average morphed cell images from this downsampled dataset as
1940 described above to calculate the location concordance.

1941

1942 ***Downsampling cellular structure size scaling analysis***

1943 We created downsampled versions of the dataset with n cells per structure randomly selected
1944 (n=10, 20, 30, 50, 100, 200, 300, 500, 1000, 1500), each with three repeats. The regression models
1945 to compute explained variances and scaling rates were recalculated on these downsampled versions
1946 of dataset. **Figure 7D** shows these statistics for a single repeat of n=300. This figure also shows how
1947 the recalculated numbers differ from the original numbers as a function of the number of cells per
1948 structure.

1949 **ADDITIONAL RESOURCES**

1950 The Allen Cell Collection, the hiPSC Single-Cell Image Dataset, protocols, the Allen Cell Discussion
 1951 Forum and additional information can be found here: (<https://www.allencell.org/>)

1952

1953 **RESOURCES TABLE**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Chemicals, Peptides, and Recombinant Proteins		
mTeSR™ 1 Medium	STEMCELL Technologies	Cat#85850
mTeSR™ 1 Medium Without Phenol Red	STEMCELL Technologies	Cat#05876
L-Ascorbic acid	Sigma-Aldrich	SKU A4403
OxyFluor™	Oxyrase	Cat#OF-0005
DL-Sodium lactate 60% (w/w) in aqueous solution	VWR	Cat# AA41529-AK
Penicillin-Streptomycin	Thermo Fisher Scientific	Cat#15140122
Y-27632 (Dihydrochloride)	STEMCELL Technologies	Cat#72308
Matrigel Growth Factor Reduced (GFR) Basement Membrane Matrix, Phenol Red-free, LDEV-free	Corning	Cat#356231; Lot#5292003; Lot#9021357
DMEM/F12 (1:1) 1X	Thermo Fischer Scientific	Cat#11039021
DPBS (1X)	Thermo Fischer Scientific	Cat#14190144
StemPro™ Accutase™ Cell Dissociation Reagent	Thermo Fischer Scientific	Cat#A1110501

CellMask™ Deep Red Plasma Membrane Stain	Thermo Fisher Scientific	Cat#C10046; Lot#1813792 (5X final concentration); Lot#1853335 and #1900978 (3X final concentration)
NucBlue™ Live ReadyProbes™ Reagent	Thermo Fisher Scientific	Cat#R37605
Tetraspeck microsphere		
Deposited Data		
hiPSC Single-Cell Image Dataset; “hipsc_single_cell_image_dataset” - contents: 216,062 cells (includes 18,186 FOVs, 25 structures)	https://open.quiltdata.com/b/allencel/packages/aics/hipsc_single_cell_image_dataset/tree/1606093417/	
Supplementary MYH10 repeat dataset	https://open.quiltdata.com/b/allencel/packages/aics/hipsc_single_cell_image_dataset_supp_myh10	
12X colony dataset:	https://open.quiltdata.com/b/allencel/packages/aics/hipsc_12x_overview_image_dataset	
Experimental Models: Cell Lines		

AICS-0014 cl. 6, nucleoli (DFC)	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0014&Product=iPSC	CVCL_JM17
AICS-0057 cl. 50, nucleoli (GC)	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0057-050&PgId=166	CVCL_VK85
AICS-0094 cl. 24, nuclear speckles	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0094-024&PgId=166	CVCL_YU30
AICS-0068 cl. 9, cohesins	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0068-009&PgId=166	CVCL_UK04
AICS-0061 cl. 36, histones	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0061-036&PgId=166	CVCL_UD17

AICS-0013 cl. 210, nuclear envelope	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0013&Product=iPSC	CVCL_IR32
AICS-0069 cl. 88, nuclear pores	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0069-088&PgId=166	CVCL_UD18
AICS-0010 cl. 55, ER (Sec61 beta)	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0010&Product=iPSC	CVCL_JM14
AICS-0046 cl. 51, ER (SERCA2)	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0046-051&PgId=166	CVCL_UD14

AICS-0011 cl. 27, mitochondria	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0011&Product=iPSC	CVCL_IR33
AICS-0033 cl. 115, peroxisomes	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0033-115&PgId=166	CVCL_VK79
AICS-0040 cl. 35, endosomes	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0040-035&PgId=166	CVCL_VK82
AICS-0022 cl. 37, lysosomes	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0022-037&Product=iPSC	CVCL_LK42

AICS-0025 cl. 44, Golgi	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0025-044&Product=iPSC	CVCL_LK43
AICS-0032 cl. 19, centrioles	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0032-019&Product=iPSC	CVCL_LK45
AICS-0012 cl. 105, microtubules	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0012&Product=iPSC	CVCL_IR34
AICS-0054 cl. 91, plasma membrane	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0054-091&PgId=166	CVCL_VK84

AICS-0016 cl. 184, actin filaments	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0016&Product=iPSC	CVCL_JM16
AICS-0007 cl. 79, actin bundles	Not released yet	N/A
AICS-0024 cl. 80, actomyosin bundles	https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0024&Product=iPSC	CVCL_JM15
AICS-0058 cl. 67, adherens junctions	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0058-067&PgId=166	CVCL_VK86
AICS-0053 cl. 16, gap junctions	https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0053-016&PgId=166	CVCL_VK83

AICS-0023 cl. 20, tight junctions	https://catalog.corieell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0023&Product=iPSC	CVCL_JM18
AICS-0017 cl. 65, desmosomes	https://catalog.corieell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0017&Product=iPSC	CVCL_IR31
AICS-0005 cl. 50, matrix adhesions	https://catalog.corieell.org/0/Sections/Search/Sample_Detail.aspx?Ref=AICS-0005&Product=iPSC	CVCL_IR30
Software and Algorithms		
FlowJo version 10.2	Treestar	
Zen 2.3 (blue edition); version 23.69.1003; service pack 2.3.69.01000; hotfix 2.3.69.01003	Zeiss	
Code used for feature calculation: aicsfeature	https://github.com/AllenCell/aicsfeature	
Code used for feature calculation: spherical harmonics parameterization	https://github.com/AllenCell/aics-shparam	

cytoplasmic parameterization	https://github.com/AllenCell/aics-cytoparam	
Code used to perform organelle size-scaling analysis	https://github.com/AllenCell/stemcel_lorganellesizescaling	
Code used to perform morphing, compute Shape Modes, and calculate multi-resolution Pearson correlation analysis on 3D single cell images	https://github.com/AllenCell/cvapipe_analysis	
Mitotic classifier annotation: The final automated 3D image classifier code (for both training and testing) and all trained models	https://github.com/AllenCell/image_classifier_3d	
Original/source data for figures in the paper are available in Github	https://github.com/aics-int/cvapipe_figure_notebooks	
Tutorials and demo for how to access the data for different purposes	https://github.com/AllenCell/quilt-data-access-tutorials	
Segmentation code used to reproduce the deep learning cell and nuclear segmentations, trained models and demo Jupyter notebooks	https://github.com/AllenCell/segmenter_model_zoo	
Segmentation code used to reproduce structure segmentation from a set of algorithms to choose from, each with restricted numbers of parameters to tune	https://github.com/AllenCell/aics-segmentation	
Code used to generate the contact sheet quality control single-cell visualizations of all segmented cells	https://github.com/AllenCellModeling/actk	

Cell Feature Explorer – 216,016 cells (from 18,186 FOVs); 25 structures; 10 features +/- apical and radial proximity	https://cfe.allencell.org	
Allencell.org - Allen Institute for Cell Science website	https://allencell.org	
ChimeraX, developed by the Resource for Biocomputing, Visualization, and Informatics at the University of California, San Francisco, with support from National Institutes of Health R01-GM129325 and the Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases	https://www.cgl.ucsf.edu/chimerax/docs/credits.html	
Spinning-disk Confocal Microscope		
Observer.Z1 microscope stand	Zeiss	
10X/4.5 NA Plan-Apochromat objective	Zeiss	Cat#420640-9900-000
100X/1.25 W C-Apochromat Korr UV Vis IR objective	Zeiss	Cat#421797-9970-000
Spinning-disk scan head CSU-X with Primary dichroic RQFT 405/488/568/647 BP filter	Yokogawa	M1N-E/FBO/C101; 24V DC; P8X006; 95P900140
NucBlue Live dye: BP filter 450/50	Chroma	Cat#ET450/50m
mEGFP tag structure or bright field: BP filter 526/50	SEMROCK	Cat#FF03-525/50
CMDR dye: BP filter 690/50	Chroma	Cat#690/50m
mTagRFP-T tag structure: BP filter 600/50	Chroma	Cat#ET600/50m
Bright field: BP filter 706/95	Chroma	Cat#ET706/95m
Laser: LASOS 405 50mw	Zeiss	Part#400600-9011-000
Laser: LASOS 488 100mw	Zeiss	Part#400600-9061-000

Laser: LASOS 561 75mw	Zeiss	Part#400600-9111-000
Laser: LASOS 638 75mw	Zeiss	Part#400600-9121-000
Transmitted light red; LAMBDA TLED+ with 740 nm wavelength LED	Sutter Instruments	
Transmitted light white; Attachment lamp VIS-LED (400-700 nm) with collector for laser system	Zeiss	Part#423053-9060
Orca Flash 4.0 V2+ cameras	Hamamatsu	Part#C1144-22CU
Piezo drive: Prior NanoScan Z 100 μ m piezo z stage: NZ100ZM/a	Zeiss	Part#2802000 224
PECON Incubator XLmulti S1; (37°C with 5% CO2)	Zeiss	Part#2802000 224
Stage insert: H201 k frame slim profile model	Okolab	Part#H201 k
Other		
96-well glass bottom plate with high performance #1.5 cover glass	Cellvis	Cat#P96-1.5H-N
CELLSTAR™ Cell Culture Multi-well Plates for Suspension Cultures, (6-well plate)	Greiner Bio-One	Cat#657185
25cm ² Rectangular Canted Neck Cell Culture Flask with Vented Cap	Corning	Cat#9381M10
96-Well, Non-Treated, U-Shaped-Bottom Microplate	Thermo Fischer Scientific	Cat#08-772-54
CytoFLEX S V4-B2-Y4-R3 Flow Cytometer (13 Detectors, 4 Lasers)	Beckman Coulter	Cat#CO9766
Argolight HM Slide	Argolight	

1954

1955 **References**

1956

1957 Aldridge, S., and Teichmann, S.A. (2020). Single cell transcriptomics comes of age. *Nature*
1958 *Communications* 11, 4307.

1959 Baghbaderani, A.A., Tian, X., Neo, B.H., Burkall, A., Dimezzo, T., Sierra, G., Zeng, X., Warren, K.,
1960 Kovarcik, D.P., Fellner, T., et al. (2015). cGMP-Manufactured Human Induced Pluripotent Stem Cells
1961 Are Available for Pre-clinical and Clinical Applications. *Stem Cell Reports* 5, 647–659.

1962 Caicedo, J.C., Cooper, S., Heigwer, F., Warchal, S., Qiu, P., Molnar, C., Vasilevich, A.S., Barry, J.D.,
1963 Bansal, H.S., Kraus, O., et al. (2017). Data-analysis strategies for image-based cell profiling. *Nat*
1964 *Methods* 14, 849–863.

1965 Chen, J., Ding, L., Viana, M.P., Hendershott, M.C., Yang, R., Mueller, I.A., and Rafelski, S.M. (2018).
1966 The Allen Cell Structure Segmenter: a new open source toolkit for segmenting 3D intracellular
1967 structures in fluorescence microscopy images. *BioRxiv* 491035.

1968 Coston, M.E., Gregor, B.W., Arakaki, J., Borensztejn, A., Do, T.P., Fuqua, M.A., Haupt, A.,
1969 Hendershott, M.C., Leung, W., Mueller, I.A., et al. (2020). Automated hiPSC culture and sample
1970 preparation for 3D live cell microscopy. *BioRxiv* 2020.12.18.423371.

1971 Drubin, D.G., and Hyman, A.A. (2017). Stem cells: the new “model organism.” *Mol Biol Cell* 28, 1409–
1972 1411.

1973 Falcon, W., and Cho, K. (2020). A Framework For Contrastive Self-Supervised Learning And
1974 Designing A New Approach. *ArXiv:2009.00104 [Cs]*.

1975 Fransen, M., Lismont, C., and Walton, P. (2017). The Peroxisome-Mitochondria Connection: How and
1976 Why? *Int J Mol Sci* 18.

1977 Gerbin, K.A., Grancharova, T., Donovan-Maiye, R., Hendershott, M.C., Brown, J., Dinh, S.Q.,
1978 Gehring, J.L., Hirano, M., Johnson, G.R., Nath, A., et al. (2020). Cell states beyond transcriptomics:
1979 integrating structural organization and gene expression in hiPSC-derived cardiomyocytes. *BioRxiv*
1980 2020.05.26.081083.

1981 Gut, G., Herrmann, M.D., and Pelkmans, L. (2018). Multiplexed protein maps link subcellular
1982 organization to cellular states. *Science* 361.

1983 Hao, F., Kondo, K., Itoh, T., Ikari, S., Nada, S., Okada, M., and Noda, T. (2018). Rheb localized on the
1984 Golgi membrane activates lysosome-localized mTORC1 at the Golgi–lysosome contact site. *J Cell Sci*
1985 131.

1986 Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser,
1987 E., Taylor, J., Berg, S., Smith, N.J., et al. (2020). Array programming with NumPy. *Nature* 585, 357–
1988 362.

1989 Hockemeyer, D., Soldner, F., Beard, C., Gao, Q., Mitalipova, M., DeKolver, R.C., Katibah, G.E.,
1990 Amora, R., Boydston, E.A., Zeitler, B., et al. (2009). Efficient targeting of expressed and silent genes
1991 in human ESCs and iPSCs using zinc-finger nucleases. *Nature Biotechnology* 27, 851–857.

1992 Hunter, J.D. (2007). Matplotlib: A 2D Graphics Environment. *Computing in Science Engineering* 9,
1993 90–95.

- 1994 Johnson, G.T., Autin, L., Al-Alusi, M., Goodsell, D.S., Sanner, M.F., and Olson, A.J. (2015). cellPACK:
1995 a virtual mesoscope to model and visualize structural systems biology. *Nature Methods* 12, 85–91.
- 1996 Kreitzer, F.R., Salomonis, N., Sheehan, A., Huang, M., Park, J.S., Spindler, M.J., Lizarraga, P.,
1997 Weiss, W.A., So, P.-L., and Conklin, B.R. (2013). A robust method to derive functional neural crest
1998 cells from human pluripotent stem cells. *Am J Stem Cells* 2, 119–131.
- 1999 Lauffenburger, D.A., and Horwitz, A.F. (1996). Cell migration: a physically integrated molecular
2000 process. *Cell* 84, 359–369.
- 2001 Liaw, A., and Wiener, M. (2002). Classification and Regression by randomForest. 2, 5.
- 2002 Macklin, D.N., Ahn-Horst, T.A., Choi, H., Ruggero, N.A., Carrera, J., Mason, J.C., Sun, G., Agmon, E.,
2003 DeFelice, M.M., Maayan, I., et al. (2020). Simultaneous cross-evaluation of heterogeneous E. coli
2004 datasets via mechanistic simulation. *Science* 369.
- 2005 Marshall, W.F. (2020). Scaling of Subcellular Structures. *Annu. Rev. Cell Dev. Biol.* 36, 219–236.
- 2006 Marshall, W.F., Dernburg, A.F., Harmon, B., Agard, D.A., and Sedat, J.W. (1996). Specific
2007 interactions of chromatin with the nuclear envelope: positional determination within the nucleus in
2008 *Drosophila melanogaster*. *Mol Biol Cell* 7, 825–842.
- 2009 McCormick, M.M., Liu, X., Ibanez, L., Jomier, J., and Marion, C. (2014). ITK: enabling reproducible
2010 research and open science. *Front. Neuroinform.* 8.
- 2011 McKinney, W. (2011). “pandas: a foundational Pythonlibrary for data analysis and statistics”. In
2012 In:Pythonfor High Performance and Scientific Computing, p. 14.9.
- 2013 Nicholas Sofroniew, Kira Evans, Juan Nunez-Iglesias, Ahmet Can Solak, Talley Lambert,
2014 kevin yamauchi, Jeremy Freeman, Loic Royer, Shannon Axelrod, Peter Boone, et al. (2019).
2015 napari/napari: 0.2.8 (Zenodo).
- 2016 Ocegüera-Yanez, F., Kim, S.-I., Matsumoto, T., Tan, G.W., Xiang, L., Hatani, T., Kondo, T., Ikeya, M.,
2017 Yoshida, Y., Inoue, H., et al. (2016). Engineering the AAVS1 locus for consistent and scalable
2018 transgene expression in human iPSCs and their differentiated derivatives. *Methods* 101, 43–55.
- 2019 Ounkomol, C., Seshamani, S., Maleckar, M.M., Collman, F., and Johnson, G.R. (2018). Label-free
2020 prediction of three-dimensional fluorescence images from transmitted-light microscopy. *Nature*
2021 *Methods* 15, 917–920.
- 2022 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,
2023 N., Antiga, L., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library.
2024 *Advances in Neural Information Processing Systems* 32, 8026–8037.
- 2025 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
2026 Prettenhofer, P., Weiss, R., Dubourg, V., et al. Scikit-learn: Machine Learning in Python. *MACHINE*
2027 *LEARNING IN PYTHON* 6.
- 2028 Pettersen, E.F., Goddard, T.D., Huang, C.C., Meng, E.C., Couch, G.S., Croll, T.I., Morris, J.H., and
2029 Ferrin, T.E. (2020). UCSF ChimeraX: Structure visualization for researchers, educators, and
2030 developers. *Protein Sci.*
- 2031 Pincus, Z., and Theriot, J.A. (2007). Comparison of quantitative methods for cell-shape analysis. *J*
2032 *Microsc* 227, 140–156.

- 2033 Roberts, B., Haupt, A., Tucker, A., Grancharova, T., Arakaki, J., Fuqua, M.A., Nelson, A., Hookway,
2034 C., Ludmann, S.A., Mueller, I.A., et al. (2017a). Systematic gene tagging using CRISPR/Cas9 in
2035 human stem cells to illuminate cell organization. *Mol Biol Cell* 28, 2854–2874.
- 2036 Roberts, B., Haupt, A., Tucker, A., Grancharova, T., Arakaki, J., Fuqua, M.A., Nelson, A., Hookway,
2037 C., Ludmann, S.A., Mueller, I.A., et al. (2017b). Systematic gene tagging using CRISPR/Cas9 in
2038 human stem cells to illuminate cell organization. *Mol Biol Cell* 28, 2854–2874.
- 2039 Rocklin, M. (2015). *Dask: Parallel Computation with Blocked algorithms and Task Scheduling*. (Austin,
2040 Texas), pp. 126–132.
- 2041 Roggiani, M., and Goulian, M. (2015). Oxygen-Dependent Cell-to-Cell Variability in the Output of the
2042 *Escherichia coli* Tor Phosphorelay. *Journal of Bacteriology* 197, 1976–1987.
- 2043 Ruan, X., and Murphy, R.F. (2019). Evaluation of methods for generative modeling of cell and nuclear
2044 shape. *Bioinformatics* 35, 2475–2485.
- 2045 Schroeder, W., Martin, K., and Lorensen, B. (2018). *The Visualization Toolkit: An Object-Oriented
2046 Approach To 3D Graphics*.
- 2047 Shen, L., Farid, H., and McPeck, M. (2009). Modeling three-dimensional morphological structures
2048 using spherical harmonics. *Evolution* 63, 1003–1016.
- 2049 Thul, P.J., Åkesson, L., Wiking, M., Mahdessian, D., Geladaki, A., Ait Blal, H., Alm, T., Asplund, A.,
2050 Björk, L., Breckels, L.M., et al. (2017). A subcellular map of the human proteome. *Science* 356.
- 2051 Valencia, P., Dias, A.P., and Reed, R. (2008). Splicing promotes rapid and efficient mRNA export in
2052 mammalian cells. *Proc Natl Acad Sci U S A* 105, 3386–3391.
- 2053 Virtanen, P. (2020). *SciPy 1.0: fundamental algorithms for scientific computing in Python* | *Nature
2054 Methods*.
- 2055 Walt, S. van der, Schönberger, J.L., Nunez-Iglesias, J., Boulogne, F., Warner, J.D., Yager, N.,
2056 Gouillart, E., and Yu, T. (2014). *scikit-image: image processing in Python*. *PeerJ* 2, e453.
- 2057 Wang, T., and Hong, W. (2002). Interorganellar Regulation of Lysosome Positioning by the Golgi
2058 Apparatus through Rab34 Interaction with Rab-interacting Lysosomal Protein. *Mol Biol Cell* 13, 4317–
2059 4332.
- 2060 Wieczorek, M.A., and Meschede, M. (2018). *SHTools: Tools for Working with Spherical Harmonics*.
2061 *Geochemistry, Geophysics, Geosystems* 19, 2574–2592.
- 2062