1 Systematic morphological profiling of human gene and allele function via Cell Painting

2

MH Rohban¹, S Singh¹, X Wu¹, JB Berthet², M-A Bray^{1,3}, Y Shrestha¹, X Varelas², JS Boehm¹, AE
 Carpenter^{1*}

- 5
- 6 Affiliations:
- 7 1) Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA
- 8 2) Department of Biochemistry, Boston University School of Medicine, Boston MA USA
- 9 3) Present address: Novartis Institutes for BioMedical Research, Cambridge MA USA
- 10 * Corresponding author: AE Carpenter: anne@broadinstitute.org
- 11

12 Abstract

13 We hypothesized that human genes and disease-associated alleles might be systematically functionally annotated using morphological profiling of cDNA constructs, via a microscopy-based 14 15 Cell Painting assay. Indeed, 50% of the 220 tested genes yielded detectable morphological 16 profiles, which grouped into biologically meaningful gene clusters consistent with known functional annotation (e.g., the RAS-RAF-MEK-ERK cascade). We used novel subpopulation-based 17 visualization methods to interpret the morphological changes for specific clusters. This unbiased 18 morphologic map of gene function revealed TRAF2/c-REL negative regulation of YAP1/WWTR1-19 20 responsive pathways. We confirmed this discovery of functional connectivity between the NF-κB pathway and Hippo pathway effectors at the transcriptional level, thereby expanding knowledge of 21 these two signaling pathways that critically regulate tumor initiation and progression. We make the 22 23 images and raw data publicly available, providing an initial morphological map of major biological 24 pathways for future study.

- 25
- 26

27 Introduction

The dramatic increase in human genome sequence data has created a significant bottleneck. The number of genes and variants known to be associated with most human diseases has increased dramatically (Amberger et al. 2015). Unfortunately, the next step - understanding the function of each gene and the mechanism of each allele in the disease - typically remains non-systematic and labor-intensive. Most commonly, researchers painstakingly design, develop, and apply a diseasespecific or biological process-specific assay.

Over 30% of genes in the human genome are of unknown function (Leonetti et al. 2016) and even
 annotated genes have additional functions yet to be uncovered. Furthermore, even when a gene's
 normal functions are known, methods are lacking to predict the functional impact of the millions of

37 genetic variants found in patients. These gaps must be filled in order to convert the promise of38 human genome sequence data into clinical treatments.

39 Therefore, there is a widespread need for systematic approaches to functionally annotate genes 40 and variants therein, regardless of the biological process or disease of interest. One general 41 approach depends on guilt-by-association, linking unannotated genes to annotated ones based on 42 properties such as protein-protein interaction data, sequence similarity, or, most convincingly, 43 functional similarity (Shehu, Barbará, and Molloy 2016). In the latter category are profiling 44 techniques, where dozens to hundreds of measurements are made for each gene perturbation and 45 the resulting profile is compared against profiles for annotated genes. Various data sources can be 46 used for profiling; gene expression is one that can be performed in relatively high-throughput and it 47 has been proven useful in predicting gene function (Lamb et al. 2006). In fact, high-throughput 48 mRNA profiles were recently used to cluster alleles found in lung adenocarcinoma based on their 49 functional impact, a precursor to therapeutic strategy for variants of previously unknown 50 significance (Berger et al. 2016).

51 Images are a less mature data source for profiling but show tremendous promise. Morphological 52 profiling data is complementary to transcriptional profiling data (Wawer et al. 2014) and is less 53 expensive. Morphological profiling has succeeded across several applications, including grouping 54 small-molecule perturbations based on their mechanism of action (Caicedo, Singh, and Carpenter 55 2016; Bougen-Zhukov et al. 2016), and grouping genes based on morphological profiles derived 56 from cells perturbed by RNA interference (RNAi) (Mukherji et al. 2006; Boutros and Ahringer 2008; 57 Fuchs et al. 2010; Pau et al. 2013). One limitation of RNAi for morphological profiling is that the 58 number of measurements must be limited or else the resulting profiles are dominated by off-target 59 effects, especially seed effects (Singh et al. 2015). Some computational solutions have shown 60 some promise in overcoming this problem for gene expression profiling (Schmich et al. 2015), but their utility is unproven for image-based profiling, and regardless RNAi does not permit analysis of 61 62 gene variants, only knockdown. Modification of genes via CRISPR will require new libraries of 63 reagents and is as yet untested in morphological profiling.

64 In the proof-of-concept work presented here, we tested morphological profiling using 65 overexpression in human cells as a general approach to annotate gene and allele function. We 66 profiled a reference series of well-known genes, and a small number of variants thereof, by Cell 67 Painting. In particular, we wondered whether the information content of this strategy would 68 outweigh potential limitations (e.g., due to cellular context or expression level). We found that the 69 approach successfully clustered genes and alleles based on functional similarity, revealed specific 70 morphological changes even when present in only a subpopulation of heterogeneous cells, and 71 uncovered novel functional connections between important biological pathways.

72 Results

73 Morphological profiles from Cell Painting of expression constructs are sensitive and 74 reproducible

75 To profile each exogenously expressed gene (or allele therein), we used our previously developed

76 image-based profiling assay, called Cell Painting (Gustafsdottir et al. 2013; Bray et al. 2016). This

- 77 microscopy-based assay consists of six stains imaged in five channels and revealing eight cellular
- 78 components: DNA, mitochondria, endoplasmic reticulum, Golgi, cytoplasmic RNA, nucleoli, actin,

79 and plasma membrane (Fig. 1A). In five replicates in 384-well plate format, we infected U-2 OS 80 cells (human bone osteosarcoma cells), chosen for their flat morphology and previous validation in 81 the assay, with an arrayed "reference" expression library of 323 open reading frame (ORF) 82 constructs of partially characterized functions (Supplementary file 1A), a subset of which have 83 been previously described (E. Kim et al. 2016). Of these, we prioritized analysis of the 220 84 constructs that were most closely representative of the annotated full length transcripts (see Methods). Morphological profiles were extracted using CellProfiler for image processing, yielding 85 86 1,384 morphological features per cell, and Python/R scripts for data processing, including feature 87 selection and dimensionality reduction (Fig. 1B, and see Methods). This computational pipeline 88 yielded a 158-dimensional profile for each of 5 replicates for each gene or allele tested.





89

Figure 1: Morphological profiling by Cell Painting. (A) Example Cell Painting images from each of the five channels for a negative control sample (no gene introduced). (B) From left to right : Cell and nucleus outlines found by segmentation in CellProfiler; raw profiles (2,769 dimensional) containing median and median absolute deviation of each of 1,384 measurements over all the cells in a sample, plus cell count; processed profiles which are made less redundant by feature selection and Principal Component Analysis; dendrogram constructed based on the processed profiles (see Fig. 3). Replicates are merged to produce a profile for each gene which is then compared against others in the experiment to look for similarities and differences.

97 Not all genes are likely to impact cellular morphology given the limitation of our experiment--using a 98 single cell line at a single time point under a single set of conditions and stained with six fluorescent labels. We therefore first asked what fraction of these ORFs impacted morphology. 99 Surprisingly, we found that 50% (110/220) of these ORF constructs induced reproducible 100 101 morphological profiles distinct from negative control profiles (Fig. 2A, and see Methods). Next, we 102 ruled out the possibility that position artifacts may have artificially inflated this result by taking an 103 alternative pessimistic null distribution which takes well position into account (Fig. 2-figure 104 supplement 1). Therefore, we conclude that a single "generic" morphological profiling assay can 105 detect signal from a substantial proportion of genes in our reference set. We next turned to testing 106 whether those signals are biologically meaningful and can lead to novel, unbiased discoveries 107 about gene function.



110 Figure 2: (A) 50% of the gene overexpression constructs produced a detectable phenotype by image-111 based profiling. Constructs yielding a reproducible phenotype ought to have a median correlation among 112 replicates that is higher than the 95th percentile of correlations seen for pairs of different constructs; this is 113 true for 51% (112 out of 220) of the constructs (as shown). Additionally, we removed two constructs that 114 passed that filter but whose profiles were highly similar to negative control profiles (not shown), leaving 110 115 constructs (50%) for further analysis. (B) Of wild-type ORF pairs that both yielded a distinguishable 116 phenotype, 96% showed significant correlation to each other. Correlations between the 23 pairs of 117 constructs that are clones of the same gene (although with potential sequence variation or possibly different 118 isoforms) were almost always much higher than correlations between pairs of constructs related to different 119 genes. The threshold, shown as the dashed line, is set to 95th percentile of profile correlation for pairs of 120 different genes. Profile correlation of these 23 pairs lie above the threshold. (C) Genes in pathways 121 thought to regulate morphology were more likely to yield detectable phenotypes vs. the remainder of 122 genes in the experiment. The same cutoff as in (A) is used to identify percentage of genes with a detectable 123 phenotype. This percentage is 87% for the genes hypothesized to change morphology, while it is 48% for the 124 other genes.

125

126 Morphological profiling is robust, showing expected relationships

127 Given that technical replicates produce similar morphological profiles, we next evaluated whether 128 similarities between profiles induced by different constructs are meaningful. We began with the 129 simplest case: for a subset of genes in the experiment, a "wild-type" sequence (see Methods for 130 important definitions) was captured in more than one ORF construct (23 pairs). These pairs either 131 correspond to different physical cloning events and preparations but with highly similar full-length 132 sequence (as defined in Methods; category a: 9 pairs), or a substantive difference in their nucleotide sequence, for example, isoforms (category b: 14 pairs). We found that, as expected, the 133 phenotypes of over-expressed wild-type ORFs of the same gene were more similar to each other, 134 135 on average, than to randomly selected genes. Of the 23 pairs for which both wild-type ORFs 136 yielded a phenotype distinguishable from negative controls, 22 (~96%) of the pairs' profiles were 137 correlated more than expected by chance (Fig. 2B, the one pair not meeting that threshold was in 138 category b), confirming that different constructs with biological similarity indeed produce similar 139 morphological profiles.

This result also confirms that the sequence differences seen in separately cloned wild-type constructs do not generally have a major functional impact, but we caution that any individual construct of interest may have an impactful mutation; thus the raw sequence data should be examined and testing alternate constructs for a gene may be recommended. Note that if, for example, only 50% of wild-type pairs showed high profile correlation, it would remain ambiguous whether it was caused by the poor assay quality or the constructs' sequence mismatches. But in this particular case the mentioned near perfect consistency rules out either of the two possibilities. We also note that the 23 pairs analyzed here are located in different well locations on each plate; this result therefore also rules out widespread artifacts, such as plate position effects or metadata errors.

We suspected that the small number of engineered constitutively activating alleles for certain genes would, on average, yield a stronger phenotype than their wild-type counterparts. We indeed found that correlations between replicates of the constitutively activating allele were typically higher than correlations between replicates of the wild-type version of a gene (Supplementary file 1B; pvalue = 0.012, one-sided paired t-test).

We hypothesized that genes in pathways known to affect cellular morphology (RAC1, KRAS, 155 156 CDC42, RHOA, PAK1, and genes related to the Hippo pathway) would be more likely to yield a 157 morphological phenotype distinguishable from negative controls than other genes in the analysis. Indeed, we found this to be true (Fisher's test p-value = 3.7×10^{-3}) (Fig. 2C). Reassured by this 158 validation, we were curious which pathways would be most and least likely to yield detectable 159 morphological phenotypes, recognizing that "pathways" are neither separate nor well-defined 160 161 entities. We found genes manually annotated as being in the Hippo, Hedgehog, cytoskeletal 162 reorganization, and Mitogen-activated protein kinases (MAPK) pathways were more likely to result 163 in a phenotype, whereas genes annotated as belonging to the JAK/STAT, hypoxia, and BMP pathways were among the least likely to yield a phenotype under the conditions tested (Fig. 2-164 figure supplement 2 and Supplementary file 1C). Nevertheless, the majority of pathways could be 165 166 interrogated by morphological profiling.

167 Morphological signature similarity captures known gene-gene relationships

Given the caveats and limitations of overexpressing genes (see Discussion), we next tested whether image-based profiling of expression constructs could capture relationships among genes known to be functionally related. Because a reliable and complete map of all gene-gene connections is not available, we evaluated the accuracy of our results via two approaches.

172 First, we compared our data to protein-protein interaction data from BioGRID (Stark et al. 2006). 173 This is imperfect ground truth for judging our predictions because two proteins might physically 174 interact without producing the same morphological phenotype when overexpressed, and genes in 175 the same pathway might regulate the same phenotype without any physical interaction. Nevertheless, we expect that the corresponding proteins of gene pairs with highest profile similarity 176 177 are more likely than average to physically interact. Indeed, looking at wild-type versions of genes 178 showing a detectable phenotype (the 73 genes represented in the 110 constructs), the ratio of 179 verified gene connections among the top 5% correlated gene pairs (9%, 13 verified out of 143 possible combinations) is significantly higher than that of other gene pairs (5%, 128 verified out of 180 181 2,485 possible; Fisher's test p-value = 0.04; Supplementary file 1D).

Second, we manually annotated each gene for the pathway with which it is associated. This approach is based on expert opinion and thus imperfect knowledge of all genes' function; furthermore many pathways interrelate, and genes in the same pathway are not expected to have 185 identical phenotypes given that their functions are rarely identical (most notably, overexpression of 186 some may activate while others suppress a biological pathway or process). Nonetheless, we 187 expect pairs of genes whose morphological profiles correlate highly to be more likely than average 188 to be annotated in the same pathway vs. different pathways. Using the same 73 genes as in the 189 previous analysis, the ratio of gene connections with the same-pathway annotation in the top 5% 190 most-correlated gene pairs was 20% (29 pairs out of 143), significantly higher than the ratio for the remaining pairs (6%, 139 pairs out of 2,485; Fisher's test p-value = 7.53×10^{-9} ; Supplementary file 191 192 1E).

193 An initial morphological map of gene function

194 Having quantitatively established that morphological profiling is sensitive, robust, and captures 195 known gene-gene relationships, we explored these relationships in a correlation matrix (Fig. 3 196 bottom left and Fig. 3-figure supplement 1). The overall structure, with multiple groupings along the 197 diagonal, is consistent with the fact that the 110 constructs (73 unique genes) that showed a 198 phenotype had been annotated as representing 19 different pathways. That is, we did not see 199 large, homogeneous clusters, as would be expected if morphological profiling was sensitive to 200 perturbation but not highly specific. This rules out uniform toxicity induced by a large number of 201 genes, for example. Neither did we see only signal along the diagonal, which would have indicated 202 no strong similarity between any gene pairs.

We next created a dendrogram (Fig. 3) and defined 25 clusters (see Methods and Fig. 3-figure supplement 2) to explore the similarities among genes. Pairs of wild-type ORFs almost always clustered adjacently, consistent with our quantitative analysis described above (Fig. 2B). After retaining only one copy of replicate ORFs, we found that the majority of clusters (19 out of the 22 clusters containing more than one gene) were enriched for one or more Gene Ontology terms (Supplementary file 1F), indicating shared biological functions within each cluster.

Using this dendrogram, we began by interrogating three clusters that conformed well to prior biological knowledge. First, we analyzed Cluster 20, containing the two canonical Hippo pathway members YAP1 and WWTR1 (more detail in Supplementary file 2-20A and Supplementary file 2-20B PDFs, and in a later section of the text). Both are known to encode core transcriptional effectors of the Hippo pathway (Johnson and Halder 2014), and a negative regulator of these proteins, STK3 (also known as MST2), is the strongest anti-correlating gene for the cluster (Supplementary file 2-20A PDF, panel c1).

Second, we noted Cluster 21 is comprised of the two phosphatidylinositol 3-kinase signaling/Akt (PI3K) regulating genes, PIK3R1 and PTEN, both frequently mutated across 12 cancer types in The Cancer Genome Atlas (TCGA) (Kandoth et al. 2013). These results are consistent with previous observations that certain isoforms of PIK3R1 reduce levels of activated Akt, a dominant negative effect (Abell et al. 2005) AKT3 is in a cluster anti-correlated to the Cluster 21 (Supplementary file 2-21A PDF, panel b1).



224 Figure 3: Morphological relationships among overexpressed genes/alleles, determined by Cell 225 Painting. Correlations between pairs of genes/alleles were calculated and displayed in a correlation matrix 226 (bottom left inset, full resolution is available as Fig. 3-figure supplement 1). Only the 110 genes/alleles with a 227 detectable morphological phenotype were included. The rows and columns are ordered based on a 228 hierarchical clustering algorithm such that each blue submatrix on the diagonal shows a cluster of genes 229 resulting in similar phenotypes. The correlations were then used to create a dendrogram (main panel) where 230 the radius of the subtree containing a cluster shows the strength of correlation. The 25 clusters containing at 231 least two constructs are printed on the dendrogram in arbitrary colored fonts, while gene names colored gray 232 and marked by asterisks are those that do not correlate as strongly with their nearest neighbors (i.e., they 233 are singletons or fall below the threshold used to cut the dendrogram for clustering). Each colored arc 234 corresponds to a cell subpopulation as noted in the legend. Line thickness indicates the strength of 235 enrichment of the subpopulation in the cluster samples compared to the negative control. Solid vs. dashed 236 lines indicate the over- vs. under-representation of the corresponding subpopulation in a cluster, 237 respectively. Note that the number next to each cluster in the dendrogram is referenced in the main text and 238 corresponds to the numbered supplemental data file for each cluster.

240 Third, we examined three Clusters (19, 6 and 3) that included many MAPK-related genes. Cluster 19 is the largest example of a tight cluster of genes already known to be associated; it includes 241 242 four activators in the RAS-RAF-MEK-ERK cascade: KRAS, RAF1 (CRAF), BRAF, and MOS. Notably, two constitutively active alleles of these genes, BRAF^{V600E} (H. Davies et al. 2002) and 243 RAF1^{L613V} (Wu et al. 2011), form a separate cluster (Cluster 6) adjacent to their wild-type 244 counterparts. Furthermore, the constitutively active RAS alleles HRAS^{G12V} and KRAS^{G12V} (McCOY, 245 Bargmann, and Weinberg 1984) are in the next-closest cluster (Cluster 3), which also contains 246 247 MAP2K4 and MAP2K3 (known to be activated by Ras (Shin et al. 2005)), as well as CDKN1A (Jalili 248 et al. 2012). By contrast, MAPKs that are known to be unrelated to the RAS-RAF-MEK-ERK cascade, such as MAPK14 in Cluster 5, are far away in the dendrogram. 249

250 Overall, these results support the notion that connections between genes can be efficiently 251 discovered using our approach.

252 Visualization approaches to assist interpretation of morphological signatures

We hypothesized that the specific morphologic features that segregated each of the clusters would provide insight into gene function. Examining images (Supplementary file 2-19A PDF, panel 3) or rank-ordered lists of features that distinguish individual profiles or clusters (Supplementary file 1G) is tedious and lacks sensitivity for all but the most obvious of phenotypes, confirming that quantitative morphological profiling is more sensitive than the human visual system.

258 We therefore devised several strategies to enhance biological interpretability from these 259 experiments and applied these in combination. First, we grouped features into meta-features 260 based on their type of measurement, i.e., shape, texture, intensity, etc., and the cell constituents to 261 which they are related, to create a Feature Grid (Fig. 4A). Second, we performed unsupervised 262 grouping of features by mapping the top 20 most-distinguishing features for each cluster onto a 263 plane, creating a Feature Map (Fig. 4B), in which highly correlated features are mapped nearby 264 each other (see "Feature Interpretation" in Methods for an explanation of individual feature names). 265 In certain cases, these visualizations revealed the nature of the morphological phenotype (e.g., 266 nuclear shape abnormalities distinguishing Supplementary file 2-7A PDF), but for others these 267 approaches did not suffice to yield an obvious phenotypic conclusion (e.g., for Cluster 19, Fig. 4A 268 and 4B).

269 Third, we hypothesized that leveraging the single-cell resolution of image-based profiling might be 270 highly sensitive in enhancing interpretation, particularly for cases where only a subset of cells is 271 distinctive from negative controls. To test this, for each given cluster of genes together with 272 negative controls we identified 20 subpopulations using k-means clustering on single cell data. We 273 calculated the abundance of cells in each of the 20 subpopulations to determine which are 274 over/under-represented relative to controls for the given cluster (corresponding images are shown; Supplementary file 2-1B, 2B, ... 25B PDFs. For example, the MAPK pathway activators in Cluster 275 276 19 show increased prevalence of a subpopulation of cells with strongly asymmetric ER, 277 mitochondria, and Golgi staining, indicating a cell polarization phenotype (Fig. 4C, and 278 Supplementary file 2-19B PDF, Categories 1 and 2), for which there is evidence in the literature 279 (Šamaj, Baluška, and Hirt 2004; Elsum, Martin, and Humbert 2013; Godde et al. 2014). This 280 phenotype was not captured by manual inspection nor the first two approaches (e.g., 281 Supplementary file 2-19A PDF, panels a2 and b2).

Encouraged by this, we supplemented the morphological map by compiling these and other visualizations into PDF files for each cluster, summarized in Fig. 5 and provided in full as Supplementary file 2. We also noticed that certain subpopulations were similar across several clusters (Fig. 3-figure supplement 3 shows sample cell images of each such subpopulation); we annotated their enrichment/de-enrichment on the dendrogram (Fig. 3).



287 288 Figure 4: Visualizations used to interpret morphology of Cluster 19 (for other clusters, see 289 Supplementary file 2-1A, ..., 25A PDF files). (A) Feature Grid. RNA and AGP (actin, Golgi, plasma 290 membrane) intensity contribute most to the genes in Cluster 19 (KRAS, RAF1, BRAF, and MOS). Dark blue 291 colors indicate higher median z-score of the relevant measurements for genes in the cluster relative to 292 negative controls. As "RadialDistribution" features do not exist for the DNA channel, it is colored in black. (B) 293 Feature Map. The feature names showing the greatest difference between the cluster and negative controls 294 are shown, based on largest absolute value of z-scores (full resolution version is available in Cluster 19A 295 PDF). They are mapped in 2D space such that features that are highly correlated with each other across all 296 genes' profiles are placed close together and thus can be interpreted together. Blue/red colored names 297 indicate positive/negative sign of the z-score (i.e., blue indicates that the cluster shows higher values than 298 controls). According to this map, the average intensity of AGP, RNA and Mito shows high variation for cells 299 within samples in Cluster 19 (e.g., large mad Cytoplasm Intensity MeanIntensity AGP, where the prefix 300 "mad" refers to median absolute deviation, a robust form of standard deviation). (C) Sample images of a 301 subpopulation of cells enriched and de-enriched for all genes in Cluster 19. Cells with asymmetric 302 organelle distribution are highly over-represented for genes in the cluster, and cells with more even 303 distribution of organelles are less abundant. Note that the exemplar cells are shown at the center of the 304 patches. This explains the duplications observed in some patches. Scale bars are 39.36 µm long. Pixel 305 intensities are multiplied by 5 for display.



307 Figure 5: Data and visualizations supporting the morphological map for each cluster. For all 25 308 clusters, there are two corresponding Supplemental PDF files. Left: Supplementary file 2-type A PDFs 309 (e.g., "1A.pdf") provide an overview of data about the cluster. Panel a1 lists the genes/alleles in the cluster 310 as well as expert annotations regarding related pathways and the cell count (as a z-score) for each 311 gene/allele. Panel b1 contains the average correlation of the cluster to other clusters, indicating uniqueness 312 of the cluster's morphological phenotype. Panel c1 lists the top five negatively correlated gene/alleles to the 313 cluster. Panel a2 shows the Feature Grid summarizing categories of morphological features distinguishing 314 the cluster from the negative control. Panel b2 shows the Feature Map displaying the names of the top 20 315 morphological features distinguishing the cluster from the negative control, positioned based on similarity. 316 Explanations for feature names can be found in the Methods section. Panel c2 shows a correlation matrix for 317 just those genes/alleles in the cluster. Panel 3 contains sample images of fields of view of cells expressing 318 each gene/allele in the cluster, along with images of the control for comparison. Right: Supplementary file 319 2-type B PDFs contain multiple plots aiming to illustrate the phenotype based on single-cell data, including 320 cell subpopulation enrichment/suppression in the cluster. First, a histogram of single-cell DNA content is 321 shown for all cells from all genes/allele treatments in the cluster, indicating the overall cell cycle distribution. 322 Next, bar plots show (for the cluster overall and for each gene in the cluster) which of 20 subpopulations of 323 cells are enriched and suppressed relative to negative controls. Finally, each subsequent page of the PDF is 324 devoted to the subpopulations whose representation differs from negative controls in a statistically significant 325 way, whether enriched or suppressed (subpopulations which are very small in both the cluster and negative 326 control samples are omitted). For each subpopulation, a bar plot shows the top 10 most-distinguishing 327 feature names (versus negative control cells). Then, sample images are shown of individual representative 328 cells from each subpopulation.

329 330

331 Using these visualizations, we began by interrogating three adjacent and correlating clusters 332 (Clusters 4, 7, and 11) contain wild-type and mutant alleles of CDC42, a gene encoding a Rho 333 family GTPase with diverse roles in cell polarity, morphology, and migration (Melendez, Grogg, and 334 Zheng 2011; Martin 2015). Cluster 4 contains the constitutively active mutant CDC42 Q61L (Nobes 335 and Hall 1999) as well as MAP3K2 and MAP3K9. The highly similar Cluster 7 contains the dominant negative alleles CDC42 T17N (Nobes and Hall 1999) and RAC1 T17N (S. Zhang et al. 336 337 1995), a related RAS superfamily member. That activating and inhibiting alleles would yield similar phenotypes when overexpressed is not surprising for CDC42 (Melendez, Grogg, and Zheng 2011). 338 339 Cluster 7 also contains isoforms and alleles of AKT: specifically, AKT3 and the constitutively active E17K alleles of both AKT1 and AKT3 (M. S. Kim et al. 2008; M. A. Davies et al. 2008). Akt is 340 341 known to be essential for certain Cdc42-regulated functions (Higuchi et al. 2001) and vice versa

342 (Stengel and Zheng 2012). Finally, the nearby Cluster 11 (which is discussed in more detail later) 343 contains the wild-type form of CDC42 as well as TRAF2. a canonical NF-κB activator: these two 344 are known to interact and share functions in actin remodeling (Marivin et al. 2014). We also note 345 that anti-correlating genes to these clusters (generally in Clusters 13 and 21) are consistent with 346 existing knowledge, including (a) AKT family member AKT1S1 (a Proline rich AKT substrate, 347 PRAS40 (Kovacina et al. 2003; Wiza et al. 2014), Supplementary file 2-7A PDF panels b1 and c1) 348 (b) CDK2 (a known target of Akt (Maddika et al. 2008)), (c) PIK3R1 and PTEN in Cluster 21, 349 described previously, which have known interactions with AKT (Cheung and Mills 2016; Hemmings 350 and Restuccia 2015). Thus, all of these connections have previously been identified.

Subpopulation visualization revealed that Clusters 4, 7, and 11 are enriched in cells that are huge and binucleate (Fig. 3, example images shown in Supplementary file 2-4B PDF). Genes in all three clusters also show irregularities in DNA content, namely, an enrichment in cells with sub-2N DNA content, a decrease in cells with 2N DNA content, and, for most genes, a decrease in cells with S phase and 4N DNA content, indicating a significant amount of DNA fragmentation and thus apoptosis (DNA histograms in Supplementary file 2-4B, 7B, and 11B PDFs). These phenotypes are consistent with these genes' known role in the cell cycle and cell polarity (Chircop 2014).

As a second test case, we examined Cluster 8, which contains PRKACA (the catalytic subunit α of protein kinase A, PKA) and two of its known substrates: GLI1 (a transcription factor mediating Hedgehog signaling)(Asaoka 2012), and RHOA^{Q63L} (a Ras homolog gene family member)(Lang et al. 1996; Rolli-Derkinderen et al. 2005). The highly similar Cluster 10 contains the wild-type RHOA, as well as ELK1 which is also linked to the Rho GTPase family and PKA (Bachmann et al. 2013; Murai and Treisman 2002).

364 We investigated the morphological changes causing these genes to cluster. RhoA is a known 365 regulator of cell morphology and cell rounding is a known related phenotype (Oishi et al. 2012). We 366 found that indeed all members of Clusters 8 and 10 significantly induce cell rounding 367 (Supplementary file 1H). Although cell count is lower for genes Clusters 8 and 10, the degree 368 varies greatly (from z-score -0.67 to -3.02, Supplementary file 2-8A and 10A PDFs, panel a1), 369 ruling out that simple sparseness of cells explains their high similarity in the assay. As well, the 370 overall DNA content distribution of the cell populations appears relatively normal (Supplementary 371 file 2-8B and 10B PDFs). Subpopulation extraction provides a satisfying biological explanation for 372 these clusters' distinctive phenotype: the increased roundness and strong variation in intensity 373 levels (per the Feature Grid) across the population stems from an increased proportion of 374 telophase, anaphase, and apoptotic cells (Fig. 3 and Supplementary file 2-8B and 10B PDFs).

We therefore conclude that the morphological map can link related genes to each other and that the morphological data can provide insight into their functions, particularly with the help of subpopulation visualization.

An unexpected relationship between the Hippo pathway and regulators of NF-κB signaling (Clusters 11, 20, and 22)

We wondered whether novel relationships might emerge from our unbiased classification of gene and allele function based on morphologic profiling. We noticed that the known regulator of NF-κB signaling, TRAF2 (in Cluster 11, together with CDC42) (Grech et al. 2004; Tada et al. 2001), yields a signature strongly anti-correlated to YAP1/WWTR1 (Cluster 20), which encode the transcriptional

. _ /_

12

384 effectors of the Hippo pathway, YAP (Yes-associated protein) and TAZ (Transcriptional coactivator with a PDZ-domain). The Hippo pathway and NF-kB signaling are critical regulators of cell 385 386 survival and differentiation, and dysregulation of these pathways is implicated in a number of 387 cancers (Varelas 2014; Hoesel and Schmid 2013; Tornatore et al. 2012), but we found no 388 evidence in the literature (in particular through BioGRID) of physical interaction between the 389 proteins encoded by Cluster 11 genes and Cluster 20 genes. Confirming our approach, a 390 functional connection between CDC42 (Cluster 11) and YAP1 (Cluster 20) has been identified: 391 deletion of CDC42 phenocopies the loss of YAP1 in kidney-specific conditional knockouts in mice 392 (Reginensi et al. 2013). Still, the NF- κ B pathway (and in particular the Cluster 11 member TRAF2), 393 has not been closely tied to YAP and TAZ in human cells (see Discussion).

394 We first wanted to characterize Clusters 11 and 20 to confirm that relationships within each cluster 395 are supported in the literature. Indeed we found evidence for most of the within-cluster 396 connections. CDC42 and TRAF2 (Cluster 11) physically interact and share functions in actin 397 remodeling (Marivin et al. 2014). As described in a prior section YAP/TAZ (Cluster 20) are known 398 to share functional similarities in the Hippo pathway, being regulated by, and also regulating, 399 cytoskeletal dynamics. Consistent with these known functions, we found that core effector of the 400 Hippo pathway which functions to restrict YAP/TAZ nuclear activity, STK3 (which encodes the Mst2) 401 kinase) (Meng, Moroishi, and Guan 2016), has a morphological signature strongly anti-correlated 402 to YAP1/WWTR1 (Supplementary file 2-20A PDF, panel c1). We note that although STK3 and 403 TRAF2 are both moderately anti-correlated with YAP/TAZ (Cluster 20), STK3 and TRAF2 are not 404 themselves highly correlated, indicating each has a different subset of phenotypes that anti-405 correlate to YAP/TAZ. We also note that two clones of another protein known to influence YAP 406 activity, STK11, form Cluster 22 which falls nearby YAP1/WWTR1; a connection between STK11 407 and YAP has been identified (albeit with opposite directionality, identified via knockdown of STK11 408 (Mohseni et al. 2014)). Further, YAP1 is among the highest anti-correlating genes to REL (data not 409 shown; REL is a singleton in the dendrogram and thus not in a cluster), whose protein product, c-410 Rel, has a known connection to TRAF2 (Jin et al. 2015). These results reaffirm that the Cell 411 Painting-based morphological signatures are a useful reporter of biologically meaningful 412 connections among genes in these pathways.

Given the striking inverse correlation between YAP1/WWTR1 and TRAF2, we sought to confirm a
negative regulatory relationship between the Hippo and NF-κB pathways by multiple orthogonal
methods.

First, we explored the observed inverse morphological impact using the Cell Painting data. The morphological impact of genes in Cluster 11 and 20 is quite strong (median replicate correlation is at the 74th and 81st percentile, and average within-group correlations are 0.66 and 0.73). Subpopulation analysis showed that Cluster 20 (YAP1, WWTR1) is enriched for cells that are slightly large, slightly elongated, and have disjoint, bright mitochondria patterns, whereas Cluster 11 (TRAF2, CDC42) is de-enriched for those subpopulations and instead enriched for binucleate cells, very large cells, and small cells with asymmetric organelles (Fig. 3 and 6A and 6B).



424 Figure 6: Morphological and transcriptional cross-talk between the Hippo pathway and regulators of 425 NF-KB signaling. (A) The TRAF2/CDC42 cluster (Cluster 11) is enriched for bi-nucleate cells, small cells 426 with asymmetric organelles, and huge cells. Note that exemplar images shown are not labeled with the 427 actual gene they are associated with. Rather they are only supposed to provide a visual insight of the cell 428 morphologies which are enriched in the gene cluster. (B) The YAP1/WWTR1 cluster (Cluster 20) is enriched 429 for cells with bright disjoint mitochondria patterns, slightly large cells, and slightly elongated cells. Scale bars 430 are 39.36 μm long. Pixel intensities are multiplied by 5 for display. (C) Gene Set Enrichment Analysis 431 (GSEA) reveals that gene overexpression leading to down-regulation of YAP1 targets (CTGF, CYR61, and 432 BIRC5) are enriched for regulators of the NF- κ B pathway (Enrichment Score p-value = 8.19×10^{-5}). The 433 horizontal axis gives the index of ORFs sorted based on the average amount of down-regulation of the YAP1 434 targets. Each blue hash mark on this axis indicates an NF-kB pathway member. The running enrichment 435 score, which can range from -1 to 1, is plotted on the vertical axis and quantifies the accumulation of NF-κB 436 pathways members on the sorted list of ORFs. (D) TRAF2 and REL suppress YAP and TAZ transcriptional 437 activity. REL and TRAF2 suppress the ability of wild-type (D1) YAP and (D2) TAZ to drive the expression of 438 a TEAD-regulated luciferase reporter. Activity of nuclear active mutants of (D3) YAP (5SA) and (D4) TAZ 439 (4SA) are similarly suppressed. Luciferase reporter activity was measured in HEK293T cells co-transfected 440 with expression constructs as indicated and a TEAD luciferase reporter was used to measure YAP-directed 441 transcription. (* p-value < 0.05, ** p-value = 0.001, *** p-value < 0.0001)

442 Second, given that YAP/TAZ are transcriptional regulators, we analyzed gene expression data. 443 Using the same constructs as in our Cell Painting experiment, we found an anti-correlated 444 relationship at the mRNA level, consistent with the anti-correlation we had seen in morphological 445 space. To do this, we used Gene Set Enrichment Analysis (Subramanian et al. 2005) and publicly 446 available data, which includes data from four to nine different cell lines at one to four time points 447 (http://lincscloud.org). Time point refers to the duration of treating the cells with over-expression 448 constructs until the time gene expression readouts are made. This analysis revealed that the NF-449 κB pathway is the pathway most enriched among genes whose overexpression results in down-450 regulation of known YAP1 targets, CTGF, CYR61, and BIRC5 (Zhao et al. 2008) (Benjamini and Hochberg (BH) adjusted p-value = 2×10^{-8} in Supplementary file 1I, and Fig. 6C), with TRAF2 451 452 being among the genes contributing to this enrichment (Supplementary file 11). We also saw 453 enrichment of NF-κB pathway members when testing a data-driven set of targets of YAP1/TAZ 454 (Fig. 6-figure supplement 1, see Methods). In the inverse analysis, genes that alter the levels of 455 TRAF2/REL common targets are weakly enriched in Hippo pathway members (Fig. 6-figure 456 supplement 2, see Methods). This is consistent with the hypothesis that NF-kB members can 457 downregulate YAP/TAZ targets but not strongly vice versa.

458 Finally, we more directly confirmed negative crosstalk between NF-KB effectors and YAP/TAZ 459 using a synthetic TEAD luciferase reporter that is YAP/TAZ responsive (Dupont et al. 2011). 460 Importantly, these confirmatory experiments used different cellular contexts and perturbation 461 constructs versus the original Cell Painting data. Co-expression of the NF-κB pathway effectors 462 TRAF2 or C-REL with YAP or TAZ led to significantly lower reporter activity than expression of 463 YAP or TAZ alone (Fig. 6D1 and 6D2). Intriguingly, mutants of YAP or TAZ that are insensitive to 464 negative regulation by the Hippo pathway (YAP-5SA and TAZ-4SA; (Zhao et al. 2008)) remained 465 sensitive to suppression of transcriptional activity by TRAF2 and C-REL, indicating that the 466 negative relationship we identified may be independent of canonical upstream Hippo pathway 467 signals (Fig. 6D3 and 6D4).

468

469 Discussion

We conclude that connections among genes can be profitably analyzed using morphological profiling of overexpressed genes via the Cell Painting assay. In a single inexpensive experiment, we were able to rediscover a remarkable number of known biological connections among the genes tested. Further, we found that morphological data from the Cell Painting assay, together with novel subpopulation visualization methods, can be used to flesh out the functionality of particular genes and/or clusters of interest.

By adopting a two-pronged approach, merging this Cell Painting morphological analysis with transcriptional data, we were able to identify an unexpected relationship in human cells between two major signaling pathways, Hippo and NF-κB, both under intense study recently for their involvement in cancer. Through validation of these clustered genes, we have identified that YAP/TAZ-directed transcription is negatively regulated by NF-κB pathway effectors and our data suggests a novel regulatory mechanism that is independent of upstream Hippo kinases.

To date, there has been little evidence of the intersection between these important signaling pathways. Recent work examining osteoclast-osteoblast differentiation has suggested that Hippo 484 pathway kinases, such as Mst2, may affect the NF-κB pathway through phosphorylation of IkB 485 proteins, thereby promoting nuclear translocation of NF-κB transcription factors (Lee et al. 2015). 486 TAZ was found to be a direct target of NF-kB transcription factors and its expression is regulated 487 via NF-KB signaling (Cho et al. 2010). Our work, however, supports a possible additional mode of 488 interaction, whereby regulators of NF-κB signaling directly regulate the function of Yap and Taz as 489 transcriptional co-factors. Recent work has demonstrated, in Drosophila, that NF-kB activation via 490 Toll receptor signaling negatively regulates the transcriptional activity of Yorkie, the homolog of 491 YAP/TAZ, through activation of canonical hippo pathway kinases (B. Liu et al. 2016). The work 492 described here identifies, for the first time in a mammalian system, that a negative regulatory 493 relationship exists between NF-κB activation and YAP/TAZ transcriptional function. Furthermore, 494 we have identified that this regulation of YAP/TAZ occurs in a manner that is independent of Hippo 495 pathway-mediated phosphorylation events on YAP/TAZ, suggesting a more direct relationship 496 between NF-kB and YAP/TAZ signaling.

497 In this work, we tested quantitatively and explored qualitatively the connections among genes 498 revealed by morphological profiling. Our underlying hypothesis was that functionally similar genes 499 would generally yield morphologically similar cells when overexpressed, and indeed we found this 500 to be the case. Still, some discussion of this point is warranted. Most commonly, gene 501 overexpression will result in activation of the corresponding pathway via amplification of the 502 endogenous gene's function. However, it is important to note that the profiling strategy to discover 503 functional relationships does not assume or require this. For example, overexpression could also 504 disrupt a protein complex, producing a trans-dominant negative effect that results in precisely the 505 opposite phenotypic effect (Veitia 2007). In still other cases, overexpression of a particular gene 506 may not affect any of the normal functions of the gene (producing a false negative signal), or 507 trigger a stress response (vielding a confounded profile), or produce a complicated response, due 508 to feedback loops. Further, artifactual phenotypes could be seen, e.g., if overexpression yields a 509 non-physiological interaction among proteins or toxic aggregates. Nevertheless, despite these 510 caveats and complications, our results indicate that valuable information could be gleaned from the 511 similarity and dissimilarity of the morphological perturbations induced by gene overexpression. 512 Using overexpression avoids the complications of RNAi off-target effects (often due to seed 513 effects), which were far more prevalent (impacting 90% of constructs in our recent study (Singh et 514 al. 2015)).

In addition to functionally annotating genes, as demonstrated here, one particularly appealing application enables personalized medicine: it should be feasible to use morphological profiling to predict the functional impact of various disease alleles, particularly rare variants of unknown significance. This has recently been successful using mRNA profiles (Berger et al. 2016). Thus, an even more exciting prospect would be to combine mRNA profiles with morphological profiles to better predict groups of alleles of similar mechanism, and ultimately to predict effective therapeutics for each group of corresponding patients.

522 We make all raw images, extracted cellular features, calculated profiles, and interpretive 523 visualizations publicly available, providing an initial morphological map for several major signaling 524 pathways, including several unexplored connections among genes for further study (see 525 Supplementary file 2). Expanding this map to full genome scale could prove an enormously fruitful 526 resource.

528 Materials and Methods

529 cDNA constructs used for expression

530 The Reference Set of human cDNA clones utilized here has been previously described (E. Kim et 531 al. 2016); ~90% of these constructs induce expression of the intended gene greater than 2 532 standard deviations above the control mean. Briefly, wild-type ORF constructs were obtained as 533 Entry clones from the human ORFeome library version 8.1 (http://horfdb.dfci.harvard.edu) with 534 additional templates generously provided by collaborating laboratories, and cloned into the 535 pDONR223 Gateway Entry vector. In addition, here, to maximize coverage of cellular pathways, 536 we included additional clones with minimal sequence deviations from the intended templates. 537 Sanger sequencing of Entry clones verified the intended transcripts and, if applicable, the intended 538 mutation. Entry constructs and associated sequencing data will be publicly available via 539 www.addgene.org and may also be available via members of the ORFeome Collaboration 540 (http://www.orfeomecollaboration.org/), including the Dana-Farber/Harvard Cancer Center 541 (DF/HCC) DNA Resource Core DNA Repository (http://www.dfhcc.harvard.edu/core-facilities/dna-542 resource/) and the DNASU Plasmid Repository at ASU Biodesign Institute 543 (http://dnasu.asu.edu/DNASU/Home.jsp). Clone requests must include the unique clone identifier 544 numbers provided in the last column of Supplementary file 1A (e.g. ccsbBroadEn 12345 as an 545 example for a specific entry clone and ccsbBroad304 12345 as an example for a specific 546 expression clone). ORFs were transferred to the pLX304 lentiviral expression vector (Yang et al. 547 2011) by LR (attL x attR) recombination.

548 For simplicity, throughout this paper "wild-type" refers to ORFs found in the original collection 549 without a particular known mutation intentionally engineered. Due to natural human variation, and 550 occasional cloning artifacts, there are often non-identical matches of such constructs to reference 551 sequence; these differences are fully documented for each construct and sequence data will be 552 publicly available through AddGene, in addition to the sequencing data for the original Entry clones 553 for the genome-scale library (Yang et al. 2011).

554 Cell lines

555 U-2 OS cells (human bone osteosarcoma cells), RRID:CVCL_0042, were obtained from ATCC and 556 propagated in the William Hahn lab; they were not additionally authenticated prior to this 557 experiment. The cell line tested negative for mycoplasma prior to this experiment. HEK293T cells, 558 RRID:CVCL_0063, were obtained from ATCC. The cell line was validated by STR profiling 559 (Genetica DNA Laboratories) and was negative for mycoplasma as measured by MycoAlert 560 Mycoplasma Detection Kit (Lonza).

561 Lentiviral transduction for morphological profiling

562 We followed our previously described protocol (E. Kim et al. 2016; Berger et al. 2016) except for 563 durations of some steps. Briefly, cells were plated in 384-well plates and transduced with lentiviral 564 particles carrying ORF constructs the next day. Viral particles were removed 18–24 hr post-565 infection and cells cultured for 48 hr until staining and imaging (72 hr total post-transduction). The 566 experiment was conducted in five replicates, each in a different plate. The number of replicates 567 being five was decided based on prior experiments (Bray et al. 2016).

568 Cell staining and imaging

569 The Cell Painting assay followed our previously published protocol (Bray et al. 2016). Briefly, eight 570 different cell components and organelles were stained with fluorescent dyes: nucleus (Hoechst 571 33342), endoplasmic reticulum (concanavalin A/AlexaFluor488 conjugate), nucleoli and 572 cytoplasmic RNA (SYTO14 green fluorescent nucleic acid stain), Golgi apparatus and plasma 573 membrane agglutinin/AlexaFluor594 (wheat germ conjugate, WGA). F-actin 574 (phalloidin/AlexaFluor594 conjugate) and mitochondria (MitoTracker Deep Red). WGA and 575 MitoTracker were added to living cells, with the remaining stains carried out after cell fixation with 576 3.2% formaldehyde. Images from five fluorescent channels were captured at 20x magnification on 577 an ImageXpress Micro epifluorescent microscope (Molecular Devices): DAPI (387/447 nm), GFP 578 (472/520 nm), Cy3 (531/593 nm), Texas Red (562/624 nm), Cy5 (628/692 nm). Nine sites per well were acquired, with laser based autofocus using the DAPI channel at the first site of each well. 579

580 Image processing and feature extraction

581 The workflow for image processing and cellular feature extraction has been described elsewhere (Bray et al. 2016), but we describe it briefly here. CellProfiler (Carpenter et al. 2006) software 582 583 version 2.1.0 was used to correct the image channels for uneven illumination, and identify, 584 segment, and measure the cells. An image quality workflow (Bray et al. 2012) was applied to 585 exclude saturated and/or out-of focus wells; six wells containing blurry images were excluded, 586 retaining 1,914 plate/well combinations in the experiment. Cellular morphological, intensity, textural 587 and adjacency statistics were then measured for the cell, nuclei and cytoplasmic sub-588 compartments. The 1,402 cellular features thus extracted were normalized as follows: For each 589 feature, the median and median absolute deviation were calculated across all untreated cells within 590 a plate; feature values for all the cells in the plate were then normalized by subtracting the median 591 and dividing by the median absolute deviation (MAD) times 1.4826 (Chung et al. 2008). Features 592 having MAD = 0 in any plate were excluded, retaining 1,384 features in all. The image data along 593 with the extracted morphological features at the per-cell level were made publicly available in the 594 Image Data Repository under DOI http://dx.doi.org/10.17867/10000105.

595 Profiling and data preprocessing

596 The code repository for the profiling and all the subsequent analysis will be made publicly available 597 at https://github.com/carpenterlab/2016 rohban submitted. We will next explain details of each analysis step implemented in the code. Single cell measurements in each well and plate position 598 are summarized into the profiles by taking their median and median absolute deviation 599 600 (abbreviated as "MAD" or "mad" in some tables) over all the cells. Although this method does not 601 explicitly capture population heterogeneity, no alternate method has yet been proven more 602 effective (Ljosa et al. 2013). We also include the cell count in a sample as an additional feature. 603 This results in a vector of 2,769 elements describing the summarized morphology of cells in a 604 sample. We then use the median polishing algorithm immediately after obtaining the summarized 605 profiles, to remove and correct for any plate position artifacts. For each feature, the algorithm de-606 trends the rows, i.e. by subtracting the row median from the corresponding feature of each profile 607 in that particular row. Next, it de-trends the columns in a similar way using column medians. The 608 row and column de-trending is repeated until convergence is reached in all the features. For the 609 rest of the analysis we considered only the constructs which have more than 99% sequence 610 identity to both the intended protein and gene transcript, to avoid testing uncharacterized 611 mutations/truncations.

612 Not all of the morphological features contain useful reproducible information. We first filter out 613 features for which their replicate correlation across all samples (except the negative controls) is 614 less than 0.30, retaining 2,200 features. Subsequently, a feature selection method is used (Fischer 615 et al. 2015). Briefly, starting with features (measurements) that we identify as essential, a new 616 feature that contributes the most information with respect to those that have been chosen, is added 617 to the set. The contribution of each feature to the already-selected features is measured by the 618 replicate correlation of the residue when the feature is regressed on the already selected features. 619 This is repeated until the incremental information added drops below a threshold. The original 620 method proposed in (Fischer et al. 2015) overfits in its regression step when the original data is 621 very high dimensional. As a remedy, in the regression step we only use features that have a 622 Pearson correlation of more than 0.50 with the selected features thus far. This prevents overfitting 623 of regression when the dimensionality of selected features grows. We stop feature selection when 624 the maximum replicate correlation of residue is less than 0.30.

The feature selection method greatly removes redundancy, but because of the non-optimal "greedy" strategy, some redundancy remains. Principal component analysis is then applied to keep 99% of variance in data, resulting in 158 principal components being selected.

628 Feature interpretation

The features measured using CellProfiler follow a standard naming convention. Each feature name is made up of several tokens separated by underscores, in the following order:

- Prefix which could be either empty or "mad". This means that the feature is calculated
 either by taking median (no prefix) or median absolute deviation ("mad" prefix) of the
 relevant measurement over all the cells in a sample.
- Cellular compartment in which the measurement related to the feature is made, i.e.,
 "Cells", "Cytoplasm", or "Nuclei". Note that features labeled "Nuclei" are based on
 segmentation of nuclei using Hoechst staining, "Cells" are based on segmentation of the
 cell edges using the RNA channel, and "Cytoplasm" is the subtraction of the
 aforementioned compartments.
- Measurement type, which can be either "Intensity", "Texture", "RadialDistribution",
 "AreaShape", "Correlation_Correlation", "Granularity", and "Neighbors". Note that
 "Correlation_Correlation" measures, within a cellular compartment, the correlation between
 gray level intensities of corresponding pixel pairs across two channels (specified in the next
 tokens in the feature name). Note also that the relative positioning of a cell is measured in
 the "Neighbors" category.
- Name(s) of channels in which the measurement is made, if appropriate (omitted for AreaShape and Neighbors).
- Feature name. The precise measurement name appears at the end. A description of each metric can be found in the CellProfiler manual (http://d1zymp9ayga15t.cloudfront.net/CPmanual/index.html)

650 Identifying ORF constructs that are distinguishable from negative controls

651 Our method to identify which genes produce a discernable profile involves first normalizing each 652 profile to the negative controls, such that a treatment's median replicate correlation becomes a 653 surrogate for phenotype strength. In the case that a treatment does not show a phenotype different 654 from the negative control, its replicates would center around the origin in the feature space. This 655 would consequently decrease the median replicate correlation. On the other hand, a phenotype 656 which is consistently observed in the replicates and is significantly different from the controls 657 results in the replicates to concentrate in a region far from the origin in the feature space, and 658 hence a high median replicate correlation value.

The cutoff for "discernible" is set based on the top 5th percentile of a null distribution. The null distribution is defined based on the correlations between non-replicates (that is, different constructs) in the experiment. Treatments whose replicate correlations are greater than the 95th percentile of the null distribution are considered as "hits" that have a morphological phenotype that is highly reproducible (Fig. 2A).

At this point, for strong treatments, all profiles of the replicates are collapsed by taking the average of individual features. 110 out of the 112 selected ORFs were significantly different from the untreated profiles in the feature space. That is, their average Euclidean distances to the untreated profiles were higher than 95th percentile of untreated profile distances to themselves. This shows these two alternative notions of phenotype strength–replicate reproducibility and distance to negative control–are consistent. We restrict all the remaining analyses to the 110 ORFs.

- 670 *Comparison of morphological connections between genes to protein-protein* 671 *interaction data and pathway annotations*
- 672 In this analysis, mutant alleles were removed and we considered only one wild-type allele for each 673 gene with a detectable phenotype, retaining 73 genes. We calculated a threshold to identify 674 significantly correlated gene pairs. We picked the threshold to minimize the probability of error in 675 classifying wild-type clone pairs versus different-gene pairs. To do so, we found the value at which 676 the probability density functions of the two groups intersect; this value (here, 0.43) can be proved 677 to have the desired property (Duda, Hart, and Stork 2012). This approach results in about 5% of 678 the gene pairs being categorized as highly correlated. We next formed a 2 by 2 contingency table, 679 where the rows correspond to two groups of gene pairs, determined by whether they have high 680 profile correlation or not. Similarly, the columns also correspond to two groups of gene pairs, 681 determined by whether the corresponding proteins have been reported to interact in BioGRID (or 682 alternatively have been annotated to be in the same pathway; Supplementary files 1C and 1D). 683 This table was then used to perform a one-tailed Fisher's exact test.
- 684 *Creation of a dendrogram relating genes to each other, and agglomerative* 685 *clustering by cutting the dendrogram*
- A dendrogram was created based on the Pearson correlation distance and average linkage, usingthe hclust function in R (Fig. 3).
- 688 Gene clusters were formed by cutting the dendrogram at a fixed correlation level, 0.522, which was 689 chosen using a stability-based measure. The measure is defined as follows: the local clustering

690 stability is measured for a range of candidate cutoffs, from 0.43 (used earlier to test consistency to 691 protein interaction data) to 0.70. The point with highest stability was chosen (Fig. 3-figure 692 supplement 2), and the stability measure was defined as the proportion of treatments whose 693 clusters do not change if the cutoff is slightly changed by a small amount, $\epsilon = .002$.

694 Subpopulation extraction

695 In order to extract cell categories (subpopulations) and subpopulation enrichment laid over the 696 dendrogram in Fig. 3, we applied k-means clustering on the normalized single cell data for each 697 gene cluster and the control. Data normalization was carried out on a plate-wise basis by z-scoring 698 each feature using the control samples as reference. In order to avoid curse of dimensionality, we 699 restricted the dataset to the features obtained from the feature selection step mentioned earlier. 700 We set k=20 to be the number of subpopulations. The algorithm was run for at most 5000 701 iterations. Each cell was assigned to the subpopulation for which it has the shortest Euclidean 702 distance to its center. Then, the number of cells belonging to each cell subpopulation was counted 703 and the proportion in each subpopulation for genes in the cluster was compared against that of the 704 control. If the change in proportion of a cell category was consistent across the genes in the 705 cluster, the cell category is shown in the Supplementary file 2-type B PDFs. To quantify this 706 consistency, we used the inverse coefficient of variation of the change in a category proportion. If 707 this quantity exceeded 1, we called the change consistent and included the corresponding cell 708 category in the PDFs. Images of cells which have highest similarity to the category center in the 709 feature space are then used to interpret and give name to each cell category (Fig. 3-figure 710 supplement 3)

711 Identifying targets of a gene using a data-driven approach

712 For this purpose, we used a replicate of the original experiment but with L1000 gene-expression 713 readouts, which is provided in the supplemental data; i.e. cell line, time point, and ORF constructs 714 are the same. This data is different from the data used in creating GSEA plots, which entails 715 multiple cell lines and time points. The mRNA levels are all normalized with respect to the negative 716 control. For each replicate of the overexpression construct, we sort the expression levels of landmark genes and take the list of top and bottom 50 landmark genes. Then, to find targets of the 717 718 gene related to the construct, we find the landmark genes among this list which has shown up at 719 least in p% of replicates/clones of the gene. In particular, we set p to 33% for YAP1, 50% for 720 WWTR1, TRAF2, and REL. Then, we simply take the intersection of predicted targets of YAP1 and 721 WWTR1 (and similarly TRAF2 and REL, separately) to get their common targets. These targets 722 are then used to produce Fig. 6-figure supplements 1-2.

723 Gene Set Enrichment Analysis

In order to produce Fig. 6C, we specified the three known targets of YAP/WWTR1 (CYR61, CTGF, and BIRC5) and queried for ORFs resulting in down-regulation of these genes. This scores each ORF (out of the 430 in the dataset) based on the observed change in mRNA level of the specified YAP/WWTR1 targets, across between four to nine different cell lines and between one to four time points. For each ORF, we then sought the summarized score which takes the mean of 4 largest scores across time point/cell line combinations. Finally, the ORFs were sorted based on the summarized score, and top 30 ORFs were tested for enrichment in different pathways (Supplementary file 1I). We used the "clusterProfiler" package in R and the KEGG pathwayenrichment analysis implemented in it for creating the GSEA plot (Yu et al. 2012).

733 Luciferase Reporter assay

734 Wild-type and mutant sequences of WWTR1 (TAZ) (4SA: S66A, S89A, S117A, and S311A) and 735 YAP1 (5SA: S61A, S109A, S127A, S164A, and S397A) were previously generated and cloned into 736 the pCMV5 backbone; these constructs are distinct from those used in the original Cell Painting 737 data set. TRAF2 and REL were cloned from the original constructs (using Broad ID# 738 ccsbBroadEn 01710 and ID# ccsbBroadEn 11094, respectively), into pCMV5 expression vectors. 739 These were sequenced and confirmed to BLAST against the appropriate Broad clone ID. The 740 empty pCMV5 backbone was used as the control condition. The Tead luciferase reporter construct, 741 8xGTIIC-luciferase was a gift from Stefano Piccolo (Addgene plasmid # 34615).

742 HEK293T cells, RRID:CVCL 0063, were transfected using Turbofect (ThermoFisher Scientific) 743 according to manufacturer's protocol. All cells were co-transfected with a β-galactosidase reporter 744 plasmid (pCMV-LacZ from Clontech) as a transfection control. Cells were lysed 48 hours following 745 transfection. Lysates were mixed with firefly luciferase (Promega) according to the manufacturer's 746 protocol and luminescence was measured using a luminometer (BioTek). Lysates were mixed with 747 o-nitrophenyl- β -D-galactoside (ONPG) and β -galactosidase expression was determined 748 spectrophotometrically by measurement of absorbance at 405nm following ONPG cleavage. All 749 luciferase readings were normalized to β -galactosidase expression for the sample. Statistical 750 analysis was conducted using a two tailed unpaired Student's t test. The data shown in Fig. 6D are 751 from triplicate samples within a single experiment and is representative of replicate experiments.

752 Acknowledgements

The authors gratefully acknowledge contributions from members of the Carpenter laboratory, especially Steven A. Moore. Funding for this work was provided by the National Science Foundation (NSF CAREER DBI 1148823 to AEC), a BroadNext10 grant from the Broad Institute, and the Slim Initiative for Genomic Medicine, a project funded by the Carlos Slim Foundation in Mexico.

758

759 Competing interests

760 The authors declare no competing financial or non-financial interests.

761 References

Abell, Kathrine, Antonio Bilancio, Richard W. E. Clarkson, Paul G. Tiffen, Anton I. Altaparmakov, Thomas G.
Burdon, Tomoichiro Asano, Bart Vanhaesebroeck, and Christine J. Watson. 2005. "Stat3-Induced
Apoptosis Requires a Molecular Switch in PI (3) K Subunit Composition." *Nature Cell Biology* 7 (4).
Nature Publishing Group: 392–98. http://www.nature.com/ncb/journal/v7/n4/abs/ncb1242.html.

Alexa, Adrian, and Jörg Rahnenführer. 2009. "Gene Set Enrichment Analysis with topGO." Available.
 https://bioconductor.riken.jp/packages/3.2/bioc/vignettes/topGO/inst/doc/topGO.pdf.

Amberger, Joanna S., Carol A. Bocchini, François Schiettecatte, Alan F. Scott, and Ada Hamosh. 2015.
"OMIM. Org: Online Mendelian Inheritance in Man (OMIM®), an Online Catalog of Human Genes and Genetic Disorders." *Nucleic Acids Research* 43 (D1). Oxford Univ Press: D789–98.
http://nar.oxfordjournals.org/content/43/D1/D789.short.

772 Asaoka, Yoshinari. 2012. "Phosphorylation of Gli by cAMP-Dependent Protein Kinase." Vitamins and

- 773 *Hormones* 88: 293–307. doi:10.1016/B978-0-12-394622-5.00013-4.
- 774 Bachmann, Verena A., Anna Riml, Roland G. Huber, George S. Baillie, Klaus R. Liedl, Taras Valovka, and 775 Eduard Stefan. 2013. "Reciprocal Regulation of PKA and Rac Signaling." Proceedings of the National 776 Academv of Sciences of the United States of America 110 (21): 8531-36. 777 doi:10.1073/pnas.1215902110.
- Berger, Alice H., Angela N. Brooks, Xiaoyun Wu, Yashaswi Shrestha, Candace Chouinard, Federica
 Piccioni, Mukta Bagul, et al. 2016. "High-Throughput Phenotyping of Lung Cancer Somatic Mutations."
 Cancer Cell 30 (2): 214–28. doi:10.1016/j.ccell.2016.06.022.
- Bougen-Zhukov, Nicola, Sheng Yang Loh, Hwee Kuan Lee, and Lit-Hsin Loo. 2016. "Large-Scale Image Based Screening and Profiling of Cellular Phenotypes." *Cytometry. Part A: The Journal of the International Society for Analytical Cytology*, July. doi:10.1002/cyto.a.22909.
- Boutros, Michael, and Julie Ahringer. 2008. "The Art and Design of Genetic Screens: RNA Interference."
 Nature Reviews. Genetics 9 (7): 554–66. doi:10.1038/nrg2364.
- Bray, Mark-Anthony, Adam N. Fraser, Thomas P. Hasaka, and Anne E. Carpenter. 2012. "Workflow and
 Metrics for Image Quality Control in Large-Scale High-Content Screens." *Journal of Biomolecular Screening* 17 (2): 266–74. doi:10.1177/1087057111420292.
- Bray, Mark-Anthony, Shantanu Singh, Han Han, Chadwick T. Davis, Blake Borgeson, Cathy Hartland, Maria
 Kost-Alimova, Sigrun M. Gustafsdottir, Christopher C. Gibson, and Anne E. Carpenter. 2016. "Cell
 Painting, a High-Content Image-Based Assay for Morphological Profiling Using Multiplexed Fluorescent
 Dyes." *Nature Protocols* 11 (9): 1757–74. doi:10.1038/nprot.2016.105.
- 793 Caicedo, Juan C., Shantanu Singh, and Anne E. Carpenter. 2016. "Applications in Image-Based Profiling of 794 Perturbations." *Current Opinion in Biotechnology* 39 (June): 134–42. doi:10.1016/j.copbio.2016.04.003.
- Cantwell-Dorris, Emma R., John J. O'Leary, and Orla M. Sheils. 2011. "BRAFV600E: Implications for
 Carcinogenesis and Molecular Therapy." *Molecular Cancer Therapeutics* 10 (3): 385–94.
 doi:10.1158/1535-7163.MCT-10-0799.
- Carpenter, Anne E., Thouis R. Jones, Michael R. Lamprecht, Colin Clarke, In Han Kang, Ola Friman, David
 A. Guertin, et al. 2006. "CellProfiler: Image Analysis Software for Identifying and Quantifying Cell
 Phenotypes." *Genome Biology* 7 (10): R100. doi:10.1186/gb-2006-7-10-r100.
- Cheung, Lydia Wt, and Gordon B. Mills. 2016. "Targeting Therapeutic Liabilities Engendered by PIK3R1
 Mutations for Cancer Treatment." *Pharmacogenomics* 17 (3): 297–307. doi:10.2217/pgs.15.174.
- Chircop, Megan. 2014. "Rho GTPases as Regulators of Mitosis and Cytokinesis in Mammalian Cells." *Small GTPases* 5 (July). doi:10.4161/sgtp.29770.
- 805 Cho, Hyun Hwa, Keun Koo Shin, Yeon Jeong Kim, Ji Sun Song, Jong Myung Kim, Yong Chan Bae, Chi Dae 806 Kim, and Jin Sup Jung. 2010. "NF-KB Activation Stimulates Osteogenic Differentiation of Mesenchymal 807 Stem Cells Derived from Human Adipose Tissue by Increasing TAZ Expression." Journal of Cellular 808 Physiology Online Library: 168-77. 223 (1). Wilev 809 http://onlinelibrary.wiley.com/doi/10.1002/jcp.22024/full.
- Chung, Namjin, Xiaohua Douglas Zhang, Anthony Kreamer, Louis Locco, Pei-Fen Kuan, Steven Bartz, Peter
 S. Linsley, Marc Ferrer, and Berta Strulovici. 2008. "Median Absolute Deviation to Improve Hit Selection
 for Genome-Scale RNAi Screens." *Journal of Biomolecular Screening* 13 (2): 149–58.
 doi:10.1177/1087057107312035.
- Bavies, Helen, Graham R. Bignell, Charles Cox, Philip Stephens, Sarah Edkins, Sheila Clegg, Jon Teague,
 et al. 2002. "Mutations of the BRAF Gene in Human Cancer." *Nature* 417 (6892): 949–54.
 doi:10.1038/nature00766.
- Bavies, M. A., K. Stemke-Hale, C. Tellez, T. L. Calderone, W. Deng, V. G. Prieto, A. J. F. Lazar, J. E.
 Gershenwald, and G. B. Mills. 2008. "A Novel AKT3 Mutation in Melanoma Tumours and Cell Lines."
 British Journal of Cancer 99 (8): 1265–68. doi:10.1038/sj.bjc.6604637.
- 820 O.. Ρ. E. Hart, and G. Stork. 2012. Pattern Duda. R. D. Classification. Wiley. 821 https://books.google.com/books?id=Br33IRC3PkQC.
- Bupont, Sirio, Leonardo Morsut, Mariaceleste Aragona, Elena Enzo, Stefano Giulitti, Michelangelo
 Cordenonsi, Francesca Zanconato, et al. 2011. "Role of YAP/TAZ in Mechanotransduction." *Nature* 474 (7350): 179–83. doi:10.1038/nature10137.
- Elsum, Imogen A., Claire Martin, and Patrick O. Humbert. 2013. "Scribble Regulates an EMT Polarity
 Pathway through Modulation of MAPK-ERK Signaling to Mediate Junction Formation." *Journal of Cell Science* 126 (Pt 17): 3990–99. doi:10.1242/jcs.129387.
- Eser, S., A. Schnieke, G. Schneider, and D. Saur. 2014. "Oncogenic KRAS Signalling in Pancreatic Cancer."
 British Journal of Cancer 111 (5): 817–22. doi:10.1038/bjc.2014.215.
- 830 Fischer, Bernd, Thomas Sandmann, Thomas Horn, Maximilian Billmann, Varun Chaudhary, Wolfgang

- Huber, and Michael Boutros. 2015. "A Map of Directional Genetic Interactions in a Metazoan Cell." *eLife*4 (March). doi:10.7554/eLife.05464.
- Fuchs, Florian, Gregoire Pau, Dominique Kranz, Oleg Sklyar, Christoph Budjan, Sandra Steinbrink, Thomas
 Horn, Angelika Pedal, Wolfgang Huber, and Michael Boutros. 2010. "Clustering Phenotype Populations
 by Genome-Wide RNAi and Multiparametric Imaging." *Molecular Systems Biology* 6 (June): 370.
 doi:10.1038/msb.2010.25.
- 837 Godde, Nathan J., Julie M. Sheridan, Lorey K. Smith, Helen B. Pearson, Kara L. Britt, Ryan C. Galea, Laura
 838 L. Yates, Jane E. Visvader, and Patrick O. Humbert. 2014. "Scribble Modulates the MAPK/Fra1 Pathway
 839 to Disrupt Luminal and Ductal Integrity and Suppress Tumour Formation in the Mammary Gland." *PLoS*840 *Genetics* 10 (5): e1004323. doi:10.1371/journal.pgen.1004323.
- Grech, Adrian P., Michelle Amesbury, Tyani Chan, Sandra Gardam, Antony Basten, and Robert Brink. 2004.
 "TRAF2 Differentially Regulates the Canonical and Noncanonical Pathways of NF-kappaB Activation in Mature B Cells." *Immunity* 21 (5): 629–42. doi:10.1016/j.immuni.2004.09.011.
- Gustafsdottir, S. M., V. Ljosa, K. L. Sokolnicki, J. Anthony Wilson, D. Walpita, M. M. Kemp, K. Petri Seiler, et
 al. 2013. "Multiplex Cytological Profiling Assay to Measure Diverse Cellular States." *PloS One* 8: e80999. doi:10.1371/journal.pone.0080999.
- Hemmings, Brian A., and David F. Restuccia. 2015. "The PI3K-PKB/Akt Pathway." Cold Spring Harbor
 Perspectives in Biology 7 (4). doi:10.1101/cshperspect.a026609.
- Higuchi, M., N. Masuyama, Y. Fukui, A. Suzuki, and Y. Gotoh. 2001. "Akt Mediates Rac/Cdc42-Regulated
 Cell Motility in Growth Factor-Stimulated Cells and in Invasive PTEN Knockout Cells." *Current Biology:* CB 11 (24): 1958–62. https://www.ncbi.nlm.nih.gov/pubmed/11747822.
- Hoesel, Bastian, and Johannes A. Schmid. 2013. "The Complexity of NF-κB Signaling in Inflammation and Cancer." *Molecular Cancer* 12 (August): 86. doi:10.1186/1476-4598-12-86.
- Jalili, Ahmad, Christine Wagner, Mikhail Pashenkov, Gaurav Pathria, Kirsten D. Mertz, Hans R. Widlund,
 Mathieu Lupien, et al. 2012. "Dual Suppression of the Cyclin-Dependent Kinase Inhibitors CDKN2C and
 CDKN1A in Human Melanoma." *Journal of the National Cancer Institute* 104 (21): 1673–79.
 doi:10.1093/jnci/djs373.
- Jin, Jin, Yichuan Xiao, Hongbo Hu, Qiang Zou, Yanchuan Li, Yanpan Gao, Wei Ge, Xuhong Cheng, and
 Shao-Cong Sun. 2015. "Proinflammatory TLR Signalling Is Regulated by a TRAF2-Dependent
 Proteolysis Mechanism in Macrophages." *Nature Communications* 6 (January): 5930.
 doi:10.1038/ncomms6930.
- Johnson, Randy, and Georg Halder. 2014. "The Two Faces of Hippo: Targeting the Hippo Pathway for
 Regenerative Medicine and Cancer Treatment." *Nature Reviews. Drug Discovery* 13 (1): 63–79.
 doi:10.1038/nrd4161.
- Kandoth, Cyriac, Michael D. McLellan, Fabio Vandin, Kai Ye, Beifang Niu, Charles Lu, Mingchao Xie, et al.
 2013. "Mutational Landscape and Significance across 12 Major Cancer Types." *Nature* 502 (7471): 333–
 39. doi:10.1038/nature12634.
- Kim, Eejung, Nina Ilic, Yashaswi Shrestha, Lihua Zou, Atanas Kamburov, Cong Zhu, Xiaoping Yang, et al.
 2016. "Systematic Functional Interrogation of Rare Cancer Variants Identifies Oncogenic Alleles."
 Cancer Discovery 6 (7): 714–26. doi:10.1158/2159-8290.CD-16-0160.
- Kim, M. S., E. G. Jeong, N. J. Yoo, and S. H. Lee. 2008. "Mutational Analysis of Oncogenic AKT E17K
 Mutation in Common Solid Cancers and Acute Leukaemias." *British Journal of Cancer* 98 (9): 1533–35.
 doi:10.1038/sj.bjc.6604212.
- Kovacina, Kristina S., Grace Y. Park, Sun Sik Bae, Andrew W. Guzzetta, Erik Schaefer, Morris J. Birnbaum,
 and Richard A. Roth. 2003. "Identification of a Proline-Rich Akt Substrate as a 14-3-3 Binding Partner." *The Journal of Biological Chemistry* 278 (12): 10189–94. doi:10.1074/jbc.M210837200.
- Lamb, Justin, Emily D. Crawford, David Peck, Joshua W. Modell, Irene C. Blat, Matthew J. Wrobel, Jim
 Lerner, et al. 2006. "The Connectivity Map: Using Gene-Expression Signatures to Connect Small
 Molecules, Genes, and Disease." *Science* 313 (5795): 1929–35. doi:10.1126/science.1132939.
- Lang, P., F. Gesbert, M. Delespine-Carmagnat, R. Stancou, M. Pouchelet, and J. Bertoglio. 1996. "Protein Kinase A Phosphorylation of RhoA Mediates the Morphological and Functional Effects of Cyclic AMP in Cytotoxic Lymphocytes." *The EMBO Journal* 15 (3): 510–19. http://www.ncbi.nlm.nih.gov/pubmed/8599934.
- Lee, Jongwon, Bang Ung Youn, Kabsun Kim, Jung Ha Kim, Da-Hye Lee, Semun Seong, Inyoung Kim, et al.
 2015. "Mst2 Controls Bone Homeostasis by Regulating Osteoclast and Osteoblast Differentiation."
 Journal of Bone and Mineral Research: The Official Journal of the American Society for Bone and
 Mineral Research 30 (9): 1597–1607. doi:10.1002/jbmr.2503.
- 888 Leonetti, Manuel D., Sayaka Sekine, Daichi Kamiyama, Jonathan S. Weissman, and Bo Huang. 2016. "A

- 24
- Scalable Strategy for High-Throughput GFP Tagging of Endogenous Human Proteins." *Proceedings of the National Academy of Sciences of the United States of America* 113 (25): E3501–8.
 doi:10.1073/pnas.1606731113.
- Liu, Bo, Yonggang Zheng, Feng Yin, Jianzhong Yu, Neal Silverman, and Duojia Pan. 2016. "Toll Receptor Mediated Hippo Signaling Controls Innate Immunity in Drosophila." *Cell* 164 (3): 406–19.
 doi:10.1016/j.cell.2015.12.029.
- Liu, Chen-Ying, Xianbo Lv, Tingting Li, Yanping Xu, Xin Zhou, Shimin Zhao, Yue Xiong, Qun-Ying Lei, and
 Kun-Liang Guan. 2011. "PP1 Cooperates with ASPP2 to Dephosphorylate and Activate TAZ." *The Journal of Biological Chemistry* 286 (7): 5558–66. doi:10.1074/jbc.M110.194019.
- Ljosa, Vebjorn, Peter D. Caie, Rob Ter Horst, Katherine L. Sokolnicki, Emma L. Jenkins, Sandeep Daya,
 Mark E. Roberts, et al. 2013. "Comparison of Methods for Image-Based Profiling of Cellular
 Morphological Responses to Small-Molecule Treatment." *Journal of Biomolecular Screening* 18 (10):
 1321–29. doi:10.1177/1087057113503553.
- Maddika, Subbareddy, Sudharsana Rao Ande, Emilia Wiechec, Lise Lotte Hansen, Sebastian Wesselborg,
 and Marek Los. 2008. "Akt-Mediated Phosphorylation of CDK2 Regulates Its Dual Role in Cell Cycle
 Progression and Apoptosis." *Journal of Cell Science* 121 (Pt 7): 979–88. doi:10.1242/jcs.009530.
- Marivin, A., J. Berthelet, J. Cartier, C. Paul, S. Gemble, A. Morizot, W. Boireau, et al. 2014. "cIAP1 Regulates
 TNF-Mediated cdc42 Activation and Filopodia Formation." *Oncogene* 33 (48): 5534–45.
 doi:10.1038/onc.2013.499.
- Martin, Sophie G. 2015. "Spontaneous Cell Polarization: Feedback Control of Cdc42 GTPase Breaks
 Cellular Symmetry." *BioEssays: News and Reviews in Molecular, Cellular and Developmental Biology* 37 (11): 1193–1201. doi:10.1002/bies.201500077.
- McCOY, Melissa S., Cornelia I. Bargmann, and Robert A. Weinberg. 1984. "Human Colon Carcinoma Kiras2 Oncogene and Its Corresponding Proto-Oncogene." *Molecular and Cellular Biology* 4 (8). Am Soc Microbiol: 1577–82. http://mcb.asm.org/content/4/8/1577.short.
- Melendez, Jaime, Matthew Grogg, and Yi Zheng. 2011. "Signaling Role of Cdc42 in Regulating Mammalian
 Physiology." *The Journal of Biological Chemistry* 286 (4): 2375–81. doi:10.1074/jbc.R110.200329.
- Meng, Žhipeng, Toshiro Moroishi, and Kun-Liang Guan. 2016. "Mechanisms of Hippo Pathway Regulation."
 Genes & Development 30 (1): 1–17. doi:10.1101/gad.274027.115.
- Mohseni, Morvarid, Jianlong Sun, Allison Lau, Stephen Curtis, Jeffrey Goldsmith, Victor L. Fox, Chongjuan
 Wei, et al. 2014. "A Genetic Screen Identifies an LKB1-MARK Signalling Axis Controlling the Hippo-YAP
 Pathway." *Nature Cell Biology* 16 (1): 108–17. doi:10.1038/ncb2884.
- Mukherji, Mridul, Russell Bell, Lubica Supekova, Yan Wang, Anthony P. Orth, Serge Batalov, Loren Miraglia, et al. 2006. "Genome-Wide Functional Analysis of Human Cell-Cycle Regulators." *Proceedings of the National Academy of Sciences of the United States of America* 103 (40): 14819–24. doi:10.1073/pnas.0604320103.
- Murai, Kasumi, and Richard Treisman. 2002. "Interaction of Serum Response Factor (SRF) with the Elk-1 B
 Box Inhibits RhoA-Actin Signaling to SRF and Potentiates Transcriptional Activation by Elk-1." *Molecular and Cellular Biology* 22 (20): 7083–92. https://www.ncbi.nlm.nih.gov/pubmed/12242287.
- Nobes, C. D., and A. Hall. 1999. "Rho GTPases Control Polarity, Protrusion, and Adhesion during Cell
 Movement." *The Journal of Cell Biology* 144 (6): 1235–44.
 https://www.ncbi.nlm.nih.gov/pubmed/10087266.
- Oishi, Atsuro, Noriko Makita, Junichiro Sato, and Taroh Iiri. 2012. "Regulation of RhoA Signaling by the
 cAMP-Dependent Phosphorylation of RhoGDlα." *The Journal of Biological Chemistry* 287 (46): 38705–
 doi:10.1074/jbc.M112.401547.
- Pau, Gregoire, Thomas Walter, Beate Neumann, Jean-Karim Hériché, Jan Ellenberg, and Wolfgang Huber.
 2013. "Dynamical Modelling of Phenotypes in a Genome-Wide RNAi Live-Cell Imaging Assay." *BMC Bioinformatics* 14 (October): 308. doi:10.1186/1471-2105-14-308.
- Reginensi, Antoine, Rizaldy P. Scott, Alex Gregorieff, Mazdak Bagherie-Lachidan, Chaeuk Chung, Dae-Sik
 Lim, Tony Pawson, Jeff Wrana, and Helen McNeill. 2013. "Yap- and Cdc42-Dependent Nephrogenesis
 and Morphogenesis during Mouse Kidney Development." *PLoS Genetics* 9 (3): e1003380.
 doi:10.1371/journal.pgen.1003380.
- Rolli-Derkinderen, Malvyne, Vincent Sauzeau, Laurent Boyer, Emmanuel Lemichez, Céline Baron, Daniel Henrion, Gervaise Loirand, and Pierre Pacaud. 2005. "Phosphorylation of Serine 188 Protects RhoA from Ubiquitin/proteasome-Mediated Degradation in Vascular Smooth Muscle Cells." *Circulation Research* 96 (11): 1152–60. doi:10.1161/01.RES.0000170084.88780.ea.
- Šamaj, Jozef, František Baluška, and Heribert Hirt. 2004. "From Signal to Cell Polarity: Mitogen□activated
 Protein Kinases as Sensors and Effectors of Cytoskeleton Dynamicity." *Journal of Experimental Botany*

- 947 55 (395): 189–98. doi:10.1093/jxb/erh012.
- Schmich, Fabian, Ewa Szczurek, Saskia Kreibich, Sabrina Dilling, Daniel Andritschke, Alain Casanova,
 Shyan Huey Low, et al. 2015. "gespeR: A Statistical Model for Deconvoluting off-Target-Confounded
 RNA Interference Screens." *Genome Biology* 16 (October): 220. doi:10.1186/s13059-015-0783-1.
- Shehu, Amarda, Daniel Barbará, and Kevin Molloy. 2016. "A Survey of Computational Methods for Protein
 Function Prediction." In *Big Data Analytics in Genomics*, edited by Ka-Chun Wong, 225–98. Springer
 International Publishing. doi:10.1007/978-3-319-41279-5_7.
- Shin, Ilchung, Seonhoe Kim, Hyun Song, Hyeong-Reh Choi Kim, and Aree Moon. 2005. "H-Ras-Specific
 Activation of Rac-MKK3/6-p38 Pathway: Its Critical Role in Invasion and Migration of Breast Epithelial
 Cells." *The Journal of Biological Chemistry* 280 (15): 14675–83. doi:10.1074/jbc.M411625200.
- Singh, Shantanu, Xiaoyun Wu, Vebjorn Ljosa, Mark-Anthony Bray, Federica Piccioni, David E. Root, John G.
 Doench, Jesse S. Boehm, and Anne E. Carpenter. 2015. "Morphological Profiles of RNAi-Induced Gene Knockdown Are Highly Reproducible but Dominated by Seed Effects." *PloS One* 10 (7): e0131370.
 doi:10.1371/journal.pone.0131370.
- Stark, Chris, Bobby-Joe Breitkreutz, Teresa Reguly, Lorrie Boucher, Ashton Breitkreutz, and Mike Tyers.
 2006. "BioGRID: A General Repository for Interaction Datasets." *Nucleic Acids Research* 34 (Database issue): D535–39. doi:10.1093/nar/gkj109.
- Stengel, Kristy R., and Yi Zheng. 2012. "Essential Role of Cdc42 in Ras-Induced Transformation Revealed by Gene Targeting." *PloS One* 7 (6): e37317. doi:10.1371/journal.pone.0037317.
- Subramanian, Aravind, Pablo Tamayo, Vamsi K. Mootha, Sayan Mukherjee, Benjamin L. Ebert, Michael A.
 Gillette, Amanda Paulovich, et al. 2005. "Gene Set Enrichment Analysis: A Knowledge-Based Approach for Interpreting Genome-Wide Expression Profiles." *Proceedings of the National Academy of Sciences* of the United States of America 102 (43): 15545–50. doi:10.1073/pnas.0506580102.
- 970 Tada, K., T. Okazaki, S. Sakon, T. Kobarai, K. Kurosawa, S. Yamaoka, H. Hashimoto, et al. 2001. "Critical 971 Roles of TRAF2 and TRAF5 in Tumor Necrosis Factor-Induced NF-Kappa B Activation and Protection 972 from Cell Death." The Journal of Biological Chemistry 276 (39): 36530-34. 973 doi:10.1074/jbc.M104837200.
- 974 Tornatore, Laura, Anil K. Thotakura, Jason Bennett, Marta Moretti, and Guido Franzoso. 2012. "The Nuclear
 975 Factor Kappa B Signaling Pathway: Integrating Metabolism with Inflammation." *Trends in Cell Biology* 22
 976 (11): 557–66. doi:10.1016/j.tcb.2012.08.001.
- Varelas, Xaralabos. 2014. "The Hippo Pathway Effectors TAZ and YAP in Development, Homeostasis and
 Disease." *Development* 141 (8): 1614–26. doi:10.1242/dev.102376.
- Veitia, Reiner A. 2007. "Exploring the Molecular Etiology of Dominant-Negative Mutations." *The Plant Cell* 19 (12): 3843–51. doi:10.1105/tpc.107.055053.
- Wawer, Mathias J., Kejie Li, Sigrun M. Gustafsdottir, Vebjorn Ljosa, Nicole E. Bodycombe, Melissa A.
 Marton, Katherine L. Sokolnicki, et al. 2014. "Toward Performance-Diverse Small-Molecule Libraries for Cell-Based Phenotypic Screening Using Multiplexed High-Dimensional Profiling." *Proceedings of the National Academy of Sciences of the United States of America* 111 (30): 10911–16. doi:10.1073/pnas.1410933111.
- Wiza, Claudia, Alexandra Chadt, Marcel Blumensatt, Timo Kanzleiter, Daniella Herzfeld De Wiza, Angelika 986 987 Horrighs, Heidi Mueller, et al. 2014. "Over-Expression of PRAS40 Enhances Insulin Sensitivity in 988 Skeletal Muscle." Archives Physiology and Biochemistrv 120 64-72. of (2): 989 doi:10.3109/13813455.2014.894076.
- Wu, Xue, Jeremy Simpson, Jenny H. Hong, Kyoung-Han Kim, Nirusha K. Thavarajah, Peter H. Backx,
 Benjamin G. Neel, and Toshiyuki Araki. 2011. "MEK-ERK Pathway Modulation Ameliorates Disease
 Phenotypes in a Mouse Model of Noonan Syndrome Associated with the Raf1L613V Mutation."
 doi:10.1172/JCI44929.
- Yang, Xiaoping, Jesse S. Boehm, Xinping Yang, Kourosh Salehi-Ashtiani, Tong Hao, Yun Shen, Rakela
 Lubonja, et al. 2011. "A Public Genome-Scale Lentiviral Expression Library of Human ORFs." *Nature Methods* 8 (8): 659–61. doi:10.1038/nmeth.1638.
- Yu, Guangchuang, Li-Gen Wang, Yanyan Han, and Qing-Yu He. 2012. "clusterProfiler: An R Package for Comparing Biological Themes among Gene Clusters." *Omics: A Journal of Integrative Biology* 16 (5): 284–87. doi:10.1089/omi.2011.0118.
- 1000 Zhang, S., J. Han, M. A. Sells, J. Chernoff, U. G. Knaus, R. J. Ulevitch, and G. M. Bokoch. 1995. "Rho 1001 Family GTPases Regulate p38 Mitogen-Activated Protein Kinase through the Downstream Mediator 1002 Pak1." The Journal of Biological Chemistry 270 23934-36. (41): 1003 https://www.ncbi.nlm.nih.gov/pubmed/7592586.
- 1004 Zhang, Z-G, C. A. Lambert, S. Servotte, G. Chometon, B. Eckes, T. Krieg, C. M. Lapière, B. V. Nusgens, and

- 1005M. Aumailley. 2006. "Effects of Constitutively Active GTPases on Fibroblast Behavior." Cellular and1006Molecular Life Sciences: CMLS 63 (1): 82–91. doi:10.1007/s00018-005-5416-5.
- 1007 Zhao, Bin, Xin Ye, Jindan Yu, Li Li, Weiquan Li, Siming Li, Jianjun Yu, et al. 2008. "TEAD Mediates YAP1008 Dependent Gene Induction and Growth Control." *Genes & Development* 22 (14): 1962–71.
 1009 doi:10.1101/gad.1664408.



1013 Figure 2-figure supplement 1: Position artifacts do not contribute to the hit rate seen in the 1014 experiment. We were concerned that position artifacts may result in overestimating the replicate correlations 1015 because replicates of the same treatment are assigned to the same well location across different plates in 1016 the experiment (that is, it was infeasible to scramble well locations). We ruled out this possibility by taking an 1017 alternative pessimistic null distribution which takes well position into account. In contrast to Figure 2A, which 1018 shows a 51% hit rate, a more pessimistic alternative null distribution is shown here (left), calculated based on 1019 the replicate correlation of pairs of negative controls in the same position only. We consider this less reliable 1020 because the number of such pairs is small (26) and we excluded edge wells; nevertheless the hit rate 1021 increases slightly, to 60%.



Figure 2-figure supplement 2: Strength of morphological phenotypes, according to annotated pathway. Morphological phenotype strength is calculated as the average replicate correlation for genes that experts manually annotated genes as belonging to each pathway. The number of genes tested in each category is shown in parentheses after the pathway name (two wild-type clones of the same gene are only counted once, but a mutant allele is counted separately). The red line shows the threshold beyond which an individual gene's profile would be considered to yield a distinguishable phenotype. Error bars indicate the deviation in the replicate correlation among the genes associated with the pathway.



Figure 3-figure supplement 1: Correlation among the 110 genes/alleles with a detectable morphological phenotype. The rows and columns are ordered based on a hierarchical clustering algorithm such that each blue submatrix on the diagonal shows a cluster of genes resulting in similar phenotypes. The scale bar depicts Pearson correlation.



1041

1042 Figure 3-figure supplement 2: Smoothed stability score across different cutoffs, in order to choose a 1043 threshold for cutting the dendrogram to form clusters. The maximum occurs at threshold = 0.522. 1044 Smoothing is done by taking the moving average of order 0.02. The stability score is defined as the 1045 proportion of treatments whose clusters are not affected if the cutoff is increased or decreased by a small 1046 amount ($\epsilon = .002$).

- 1047
- 1048



Figure 3-figure supplement 3: Common cell subpopulations seen across more than one cluster. 1052 These names are used to annotate clusters of genes in Fig. 3. Example images shown are taken from 1053 individual clusters. Scale bar is 63 μm and image intensities are log normalized. References to size and 1054 shape in the subpopulation legends refer to both the nucleus and cell borders, unless otherwise noted.





Figure 6-figure supplement 1: Gene Set Enrichment Analysis (GSEA) reveals that overexpression constructs sorted based on their similarity to YAP1/WWTR1 overexpression (in terms of impact on particular mRNA targets), are enriched for regulators of the NF-KB pathway (Enrichment Score pvalue = 0.0019). mRNA targets common to both YAP1 and WWTR overexpression include INPP4B, MAP7, LAMA3, STMN1, and TRAM2, which are positively regulated, and SPP1, IER3, RAB31, and GPR56, which are negatively regulated.





1072 Figure 6-figure supplement 2: Gene Set Enrichment Analysis (GSEA) reveals that overexpression 1073 constructs sorted based on their similarity to TRAF2/REL overexpression (in terms of impact on 1074 particular mRNA targets), are weakly enriched for regulators of the Hippo pathway (Enrichment Score 1075 p-value = 0.024). mRNA targets common to both TRAF2 and REL overexpression include NFKBIA, IKBKE, 1076 AKAP8, and BIRC2, which are positively regulated and RPA3 which is negatively regulated. As compared to 1077 Fig. 6C and Fig. 6-figure supplement 1, this is a weaker/lower-confidence enrichment - note the lower 1078 maximum height (~ 0.44 compared to > 0.6) and higher p-value (0.024 compared to < 0.002). Still, we note 1079 that WWTR1 and PPP1CA are the top two matches among those annotated as related to the Hippo pathway 1080 in KEGG; PPP1CA (also known as PP-1A) activates TAZ (C.-Y. Liu et al. 2011). 1081

1082 Supplementary Files

1083

1085

1084 Supplementary File 1

1086 We have released the data tables related to graphs and analyses presented in the paper in the ZIP format.1087 This file includes several PDFs listed below:

- 1088
 1A: List of all the 323 constructs used in the experiment along with the target transcript and their public clone ID.
- 1090
 1B: Replicate correlation is higher in the constitutively active mutant allele compared to the wild-type allele, except for AKT3_E17K. Constitutively active mutant annotations were obtained by literature search for all the mutants in the experiment showing a detectable phenotype. Genes shown here are only those where either the wild-type gene or its constitutively activating allele yielded a phenotype distinct from controls.
- 1C: Pathways sorted based on proportion of their associated gene showing a detectable phenotype.
- 1096
 1D: Highly correlated proteins (according to morphology in the Cell Painting assay) that have also been reported to interact physically.
- 1E: Highly correlated genes (according to morphology in the Cell Painting assay) that have also been annotated to be related to the same pathway.
- 1F: Gene Ontology terms associated with each gene cluster.
- 1G: Rank ordered list of distinctive features based on their z-scores for Cluster 19.
- 1H: All genes/alleles in Cluster 8 and 10 induce cell rounding.
- 11: The NF-□B signaling pathway is the most enriched when searching for gene overexpressions that downregulate known YAP/TAZ targets (CYR61, CTGF, and BIRC5).
- 1105 Supplementary File 2
- Type A and B PDFs are collected in a ZIP file in supplementary file 2. The details of the contents have beendescribed in Fig. 5.
- 1108 Supplementary File 3
- 1109 The CellProfiler pipeline used to process the images is released as the supplementary file 3.







Nucleus (DNA)

Cytoplasmic RNA/ Nucleoli (RNA)

В



Samples

Cells and Nuclei Segmentation



Endoplasmic reticulum (ER)

Measurements

Raw profiles (2,769 dimensional)



Principal Components

Mitochondria (Mito) Actin/Golgi/ Plasma membrane (AGP)







		5%	60%
örrelation	1.0-		
	0.8-		
	0.6-		
	0.4-		
	02-		







1%		5%
	1.0-	
	0.8-	
	0.6-	
	0.4-	
	- 0.2-	
	0.0	
	-0.2-	
	-0.4 -	
ates of construct		pairs of different genes

n=220

n=5956



pairs of constructs containing the same gene n=23





							RAF1_WT.1 BAF1_WT.2
							BRAF_WT.1
							BRAF_WT.2 BCL2L11 WT
							CTNNB1_S33A.S37A.T41A.T45A
							BMP2 WT
							ATF6_WT.2
							CEBPA_W1.1 CEBPA_WT.2
							JUN_WT.2
							NOTCH1_ICN1.1
							JUN_WT.1
							STK3_WT.1 STK3_WT.2
							ELK1_WT
							RHOA_WT RHOA Q63L
							PRKACA_WT.2
							GLI1 WT
							REL_WT.1
							EIF4EBP1 WT
							DIABLO_WT
							PPARGC1A_W1.2 HSPA5_WT
							ERN1_WT.1
							SDHA WT
							ARAF_WT.1
							YAP1_WT.2 YAP1_WT.1
							YAP1_WT.3
							WWTR1 WT
							AKT1_WT.1
							RHOA_T19N RELB WT
							STK11_WT.1
							PKIA_WT
							DKK1_WT
							RPS6KB1_WT.2
							SMO_WT.2
							DLL1_WT
							SMAD3_WT.1
							PRKCZ_WT.1 PRKCZ_WT.2
							PRKCE_WT.1
							GRB10_WT.2
							CRY1_WT.1
							PRKACG_WT.1
							GSK3B_WT.1
							DVL3_WT
							TBK1_WT.1
							STAT3_C-C
							CDK4_R24C
							ATF4_WT.1
							ATF4_WT.2
							XBP1_WT.2
							MAPK14_WT.1
							RBPJ_WT.1
							RBPJ_WT.2
							PIK3R1_WT.1
							PTEN_WT
							FOXO1_WT.2
							CDK2_WT.2
							AKT1S1_WT.2
							CCND1_WT.2

binculeate huge DNA **RNA** ER Mito AGP

apoptotic

triangular

extremely elongated

slightly large

telophase

small cells (condensed)

slightly elongated

late telophase

disjoint bright mitochondria

DNA RNA

mad_Cells_Intensity_MedianIntensity_AGP mad_Cytoplasm_Intensity_UpperQuartileIntensity_AGP mad_Cytoplasm_Intensity_MeanIntensity_AGP mad_Cytoplasm_Intensity_MeanIntensity_Mito mad_Cytoplasm_Intensity_MedianIntensity_AGP mad_Cells_Intensity_MeanIntensity_Mito

Cells_Intensity_LowerQuartileIntensity_RNA

Cytoplasm_Intensity_LowerQuartileIntensity_RNA

mad_Cells_Intensity_UpperQuartileIntensity_RNA mad_Cells_Intensity_MeanIntensity_RNA mad_Cells_Intensity_LowerQuartileIntensity_RNA mad_Cells_Intensity_NetionIntensity_RNA LowerQuartileIntensity_RNA mad_Cytoplasm_Intensity_MeanIntensity_RNA mad_Cytoplasm_Intensity_MeanIntensityEdge_RNA mad_Cytoplasm_Intensity_LowerQuartileIntensity_RNA Cells_Intensity_MeanIntensityEdge_RNA mad_Cytoplasm_Intensity_MedianIntensity_RNA mad_Cytoplasm_Intensity_UpperQuartileIntensity_RNA

Mito ER AGP

Cells_Intensity_StdIntensityEdge_RNA

De-enriched

Α

Treatment Pathway Expert Allocation GRB10,WT.1 Cassesical Insulin Receptor Signaling Activator GRB10,WT.2 Cassesical Insulin Receptor Signaling Activator PBSCE.WT.1 Cassesical PSC Activator

inter a

subpopulation enrichment/ suppression

10 top most

Α

Β

