

Scoring diverse cellular morphologies in image-based screens with iterative feedback and machine learning

Thouis R. Jones^{a,b,c,1}, Anne E. Carpenter^{a,b,1,2}, Michael R. Lamprecht^b, Jason Moffat^{b,3}, Serena J. Silver^a, Jennifer K. Grenier^a, Adam B. Castoreno^d, Ulrike S. Eggert^d, David E. Root^a, Polina Golland^c, and David M. Sabatini^{a,b,e}

^aThe Broad Institute of MIT and Harvard, 7 Cambridge Center, Cambridge, MA 02142; ^bWhitehead Institute for Biomedical Research, Massachusetts Institute of Technology, 9 Cambridge Center, Cambridge, MA 02142; ^cComputer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, 32 Vassar Street, Cambridge, MA 02139; ^dDana-Farber Cancer Institute and Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, 250 Longwood Avenue, Boston, MA 02115; and ^eDepartment of Biology, Massachusetts Institute of Technology, 31 Ames Street, Cambridge, MA 02139

Edited by Edward M. Scolnick, The Broad Institute, Cambridge, MA, and approved December 12, 2008 (received for review September 8, 2008)

Many biological pathways were first uncovered by identifying mutants with visible phenotypes and by scoring every sample in a screen via tedious and subjective visual inspection. Now, automated image analysis can effectively score many phenotypes. In practical application, customizing an image-analysis algorithm or finding a sufficient number of example cells to train a machine learning algorithm can be infeasible, particularly when positive control samples are not available and the phenotype of interest is rare. Here we present a supervised machine learning approach that uses iterative feedback to readily score multiple subtle and complex morphological phenotypes in high-throughput, image-based screens. First, automated cytological profiling extracts hundreds of numerical descriptors for every cell in every image. Next, the researcher generates a rule (i.e., classifier) to recognize cells with a phenotype of interest during a short, interactive training session using iterative feedback. Finally, all of the cells in the experiment are automatically classified and each sample is scored based on the presence of cells displaying the phenotype. By using this approach, we successfully scored images in RNA interference screens in 2 organisms for the prevalence of 15 diverse cellular morphologies, some of which were previously intractable.

high-content screening | high-throughput image analysis | phenotype

The history of biology has been dramatically shaped by classic visual screens in model organisms, including *Drosophila melanogaster* (1–3), *Saccharomyces cerevisiae* (4), *Caenorhabditis elegans* (5), and the zebrafish *Danio rerio* (6, 7). In each case, biological pathways were discovered because researchers were intrigued by groups of peculiar-looking mutants and identified the genes underlying their phenotypes. Because researchers have favored the extensive study of relatively few genes (8), classic, wide-net approaches like screening are as relevant as ever to probe known biological pathways and discover new ones. Modern technology now enables large-scale experiments in cultured cells to identify human genes that underlie biological processes via RNAi. Automation also allows the screening of chemical libraries to identify perturbants useful as research tools or drugs.

Despite these advances, scoring cells in images for rare and unusual morphologies has, in general, remained a significant bottleneck (9–12). Cell image analysis allows accurate identification and measurement of cells' features, enabling automated analysis of certain phenotypes that were previously intractable (13–26). However, many interesting phenotypes require the assessment of several measured features of cells. Machine learning methods that select and combine multiple features for automated cell classification have been used to score many phenotypes (15–26). These methods require the provision of example cells that do and do not display the morphology of interest (i.e., positive and negative cells). Finding positive cells is straightforward when positive control samples are available and most of the cells therein show the phenotype. However, when this is not the case, as in classic exploratory screens, finding a sufficient number of positive cells can be prohibitively difficult.

Even when positive control samples are available, using positive example cells from only those samples can lead to inaccurate scoring because of overfitting of the machine learning algorithm.

Here we describe our approach to scoring multiple complex and subtle phenotypes in large-scale, image-based screens. It is particularly effective when positive control samples are not available or not highly penetrant, as is often the case in RNAi and chemical screens. Our approach uses: (a) a biologist's ability to identify an "interesting" phenotype, (b) automatic measurement of multiple features for each cell, (c) a computer's ability to rapidly test multiple combinations of features using machine learning algorithms, and (d) a computer's ability to quickly and objectively classify millions of individual cells based on their measured features. We used our approach to score 15 diverse cellular phenotypes in large-scale RNAi screens in human and *D. melanogaster* cells, demonstrating that automated scoring for image-based chemical and genetic screens for multiple complex, low-penetrance phenotypes is now feasible.

Results

Overview of the Approach. We have developed and validated a method for researchers to rapidly train a computer to score unusual cell morphologies automatically (Fig. 1). First, we automatically identify and measure every cell in every image in the experiment by using the cell-image analysis software CellProfiler (13), which generates a cytological profile (27), or cytoprofile, for each cell. This cytoprofile consists of a set of numbers that describe the cell's characteristics, including size, shape, and the intensity and texture of various stains in various compartments (Fig. 1A). Next, the researcher initiates the training phase by identifying a few positive example cells that display a phenotype of interest and negative example cells without the phenotype (Fig. 1B). These cells can be from control samples if the screen has been designed to address a particular phenotype, or selected at random if the screen's goal is to uncover previously uncharacterized phenotypes in an exploratory screen. Most commonly, these example cells are taken from the full population without reference to the particular sample from

Author contributions: T.R.J., A.E.C., D.E.R., P.G., and D.M.S. designed research; T.R.J., A.E.C., M.R.L., J.M., S.J.S., J.K.G., A.B.C., and U.S.E. performed research; T.R.J., A.E.C., and P.G. analyzed data; and T.R.J. and A.E.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

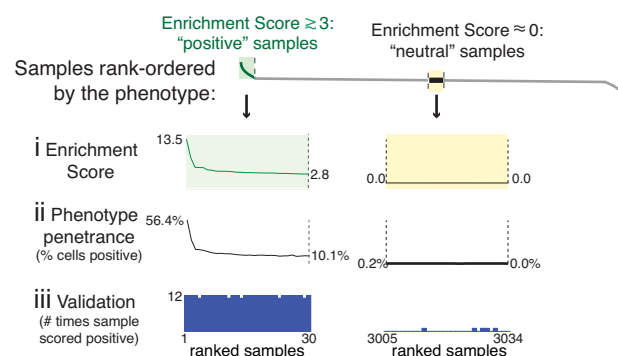
¹T.R.J. and A.E.C. contributed equally to this work.

²To whom correspondence should be addressed. E-mail: anne@broad.mit.edu.

³Present address: Banting and Best Department of Medical Research, Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, 160 College Street, Room 802, Toronto, Ontario, Canada M5S 3E1.

This article contains supporting information online at www.pnas.org/cgi/content/full/0808843106/DCSupplemental.

© 2009 by The National Academy of Sciences of the USA



correction and rule refinement, the researcher has classified a few hundred cells, and these are used to produce a rule specific to the phenotype of interest. In the final step (Fig. 1C), the rule is applied to the cytoprofiles of every cell in the experiment, classifying each cell as positive or negative. Ultimately, the goal of the screen is to score each sample, which is a population of cells subjected to a particular RNAi or chemical treatment. Because simply ranking samples by the percentage of cells that are positive can be misleading for samples with few cells, we developed an “enrichment score” to rank each sample (see Fig. 2 and *Methods*). The researcher may continue to conduct further rounds of error correction and rule refinement based on images from samples with many positive cells, ultimately producing a rule with satisfactory accuracy. Although highly dependent on the complexity of the phenotype and the scarcity of positive example cells, the entire process of training for a phenotype typically takes a few hours.

Nearly every phenotype we attempted to score could be scored accurately without customization of the image processing. That is, the standard cytoprofiles were sufficient for accurate classification in all but the Peas in a Pod phenotype. We added one feature (angle between a nucleus' 2 nearest neighbors) to the image-analysis step to better identify this phenotype (Fig. 4C). Also, we abandoned attempts to train a rule to identify a “sparkly actin” phenotype (Fig. S1); few positive example cells could be found, and it is possible that our cytoprofiles did not contain appropriate texture measurements.

Features from the cytoprofiles that were used to classify cells for each phenotype usually included a mixture of measurements of

which they are derived. This action is taken to avoid overfitting the machine learning algorithm to a few particular samples.

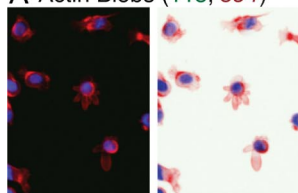
Once a few dozen individual cells have been classified by the researcher, a machine learning algorithm is used to determine a tentative rule (i.e., a classifier) that distinguishes the cytoprofiles of the positive and negative example cells, using the GentleBoosting algorithm applied to regression stumps (28). Other machine learning methods are likely to be equally effective, based on their performance in previous work (15–24). The system then presents the researcher with a new batch of cells, which it has classified based on the tentative rule, and the researcher corrects errors. The corrections are used to refine the rule. After several rounds of error

Phenotype (# pos, # neg
example cells in training set)

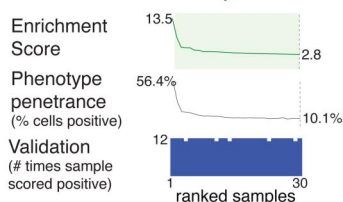
Validation

Penetrance histogram

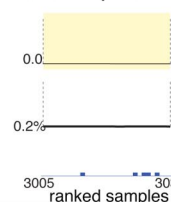
A Actin Blebs (115, 394)



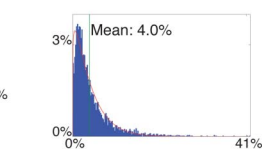
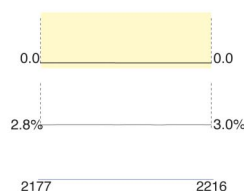
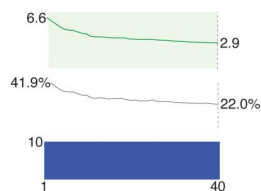
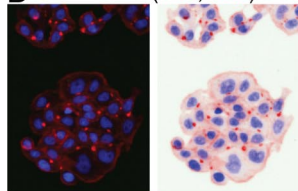
Positive samples



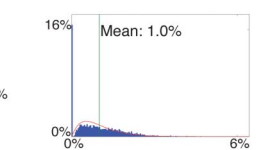
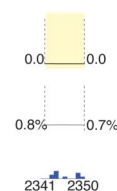
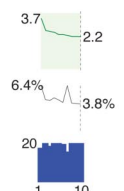
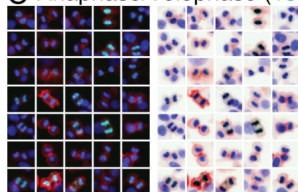
Neutral samples



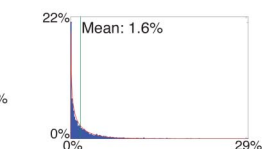
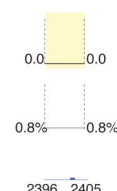
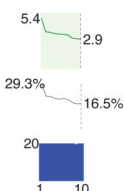
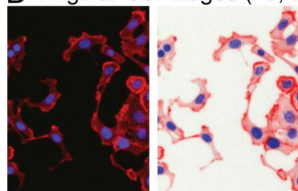
B Actin Dots (120, 182)



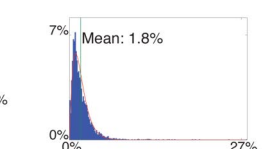
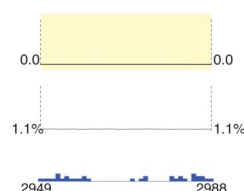
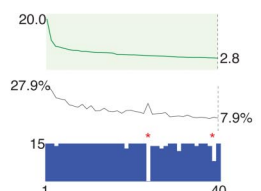
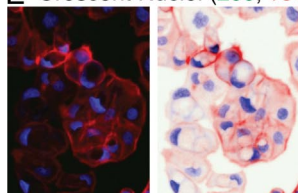
C Anaphase/Telophase (187, 1653)



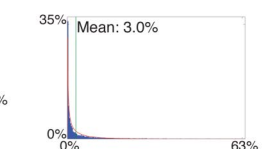
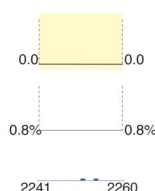
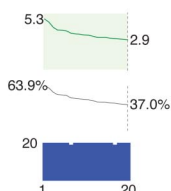
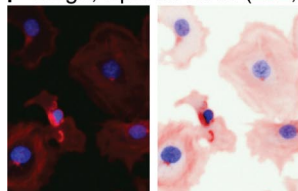
D Angular Cell Edges (73, 321)



E Crescent Nuclei (200, 1510)



F Large, Spread Cells (202, 316)



G Long Projections (59, 345)

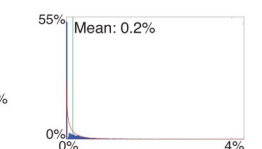
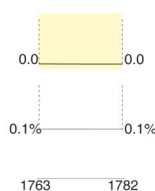
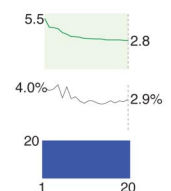
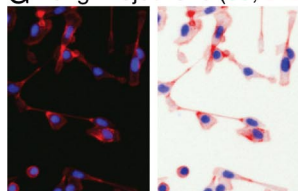


Fig. 3. Results of the phenotype-scoring system, for diverse cellular morphologies in human cells. Each row shows images and data for a different cellular morphology that the system was trained to recognize and score. The phenotype column shows the name of each phenotype along with the number of positive and negative example cells in the training set after all rounds of iteration were completed by the researcher. Images for each phenotype follow a color scheme: blue, DNA (contrast-stretched); red, actin (contrast-stretched); green, phospho-histone H3 (absolute scale). (*Left*) Traditional pseudocoloring of the fluorescence microscopy images. (*Right*) Color-adjustment using the “Invert For Printing” module of CellProfiler. The width of each image (or montage, for multiframe images) is 102 μm . For details on the validation column, see Fig. 2. The penetrance histogram column shows the distribution of per-sample penetrance for each phenotype, along with the mean (shown as text and with a green line) and the model fit to the data (red line).

frequency in WT cells. Factors like cell cycle, local environment, stochastic noise, and epigenetics all play a role in generating nonuniform populations of cells (29, 30). We therefore wondered whether any samples would have an unusually high proportion of cells showing these naturally occurring rare morphologies. Interestingly, every phenotype we pursued yielded at least some RNAi samples in which the phenotype was significantly enriched. This is consistent with the possibility that the number of phenotypic states that are possible for a cell is fairly limited, and natural variation in mRNA expression levels can push cells into one of these states, even without the influence of RNAi. In any event, the system enabled us to indulge our curiosity by pursuing unusual and uncharacterized cellular morphologies, as in classic genetic screens.

Validation, Comparison to Previous Methods, and Flexibility. We tested our method's accuracy at ranking samples by having researchers score samples (that is, images showing a population of cells) by eye. The biologically relevant score for a sample is enrichment of cells that display the phenotype, rather than a hard "positive" or "negative" label, because samples in screens typically do not fall into clear positive and negative classes (particularly when judged by different researchers), but instead fall along a continuum (31). Our goal is to bring highly enriched samples to the attention of the researcher; therefore, our validation design (forced choice, described in *Methods*) (32) aimed to test whether top-ranking samples were indeed enriched relative to samples scored as neutral.

The results for actin blebs are shown in detail in Fig. 2, and data for all human cell phenotypes are shown in the validation column in Figs. 3 and 4. For each phenotype, we rank-ordered the 5,000 puromycin-treated samples by enrichment score (Fig. 2A), as would be done in a typical screen. For validation, researchers were forced to choose between pairs of samples. One sample in each pair had been scored by the computer as highly enriched for the phenotype and the other as neutral. We recorded the number of times each sample was chosen as positive by the researchers (bar chart, Fig. 2C).

Among all 360 samples identified as "hits" across the different phenotypes (Figs. 3 and 4, positive samples column), there were 0 false negatives among the 360 samples identified as neutral and 2 potential false positives (red stars in Fig. 3E). Note that false positives can be readily weeded out by eye after analysis and that we cannot estimate the actual false-negative rate without knowing a priori the number of true positive samples, which is not possible in this screen. Agreement between humans was comparable with that between humans and automated scoring (Table S1), indicating sufficient accuracy to bring samples enriched for each phenotype to the attention of the researcher.

The phenotypes we chose were particularly challenging because their average penetrance was low (0.2–6.1%), and even the strongest hits for some phenotypes contained <5% positive cells. All phenotypes were, nonetheless, readily scored by our method. Previous approaches (15, 19, 20) have succeeded on highly penetrant phenotypes where positive control samples are known, but none of the phenotypes in our study had positive control samples available, and most were low-penetrance. We chose 4 of the phenotypes in this study and retrospectively tested a positive control-based method on them (Fig. S4). The method worked well on the most highly penetrant, straightforward phenotype, large spread cells (Fig. S4A), but was inferior on the other 3 phenotypes of greater morphological complexity and lower penetrance, in some cases even failing to highly rank the training samples (Fig. S4 B–D).

Overfitting is a concern when using machine learning algorithms, but boosting variants are fairly resistant to it (28). Cross-validation results (Fig. S5) show that this is also the case for our approach. The classification accuracy is typically not significantly reduced as the number of individual regression stumps forming a rule for a phenotype increases. To increase the coverage of the training set

and guard against training to a too-narrow definition of a phenotype, it is useful to inspect images of the top-ranked samples (or positive control samples, if available), in which positively classified cells are marked. From these images, it is easy to identify false-negative cells and add them to the training set during the iterative training phase.

We considered whether a rule will generalize to new experiments. A rule trained on one experiment is unlikely to be applicable to experiments involving different assay protocols, cellular stains, or image acquisition instrumentation, although with our approach, the time required to generate a new rule for the new experiment is minimal. For replicate experiments, creating a training set from one replicate and applying the rule it generates to another replicate risks negatively impacting its accuracy because of undetected experimental variation (Fig. S6B). The more robust approach is to create a training set spanning all replicates (Fig. S6A).

Lastly, we tested our method's flexibility by applying it to another large-scale image set. Previously, 288 genes were screened for a metaphase phenotype by RNAi in *Drosophila* by using living-cell microarrays (33). In our previous work, we identified cells in metaphase by empirically applying sequential gates based on 4 measured features of the DNA stain of each cell. This process took more than a week. With our new approach, we identified metaphase nuclei and accurately scored the entire screen within 4 h, of which only 1 h was hands-on time (Fig. S7 and Fig. S8). The top of the rank-ordered list of genes from the screen (*SI Text*) contained widerborst (CG5643, the one hit in our original study), as well as other cell-cycle-related genes, e.g., polo (CG12306) and microtubule star (CG7109). The gene most deenriched for metaphase nuclei was Nima-related kinase 2 (Nek2, CG17256; "Nima" derives from "never in mitosis"). As was the case for complex human phenotypes (Fig. S4), providing the positive control sample images directly to the machine learning algorithm was unsuccessful (Fig. S9).

Discussion

Together, this work indicates that automated scoring of a wide variety of morphologies can be accomplished quickly and easily, even when a phenotype is rare in the WT population and positive control samples are not available. Specifically, the approach is scalable to large-scale, image-based screens (chemical or genetic) in which multiple complex phenotypes are examined. Whereas screening for perturbations of general cellular functions like cell division has yielded large networks of genes (14, 34), the ability to identify more subtle and rare cellular morphologies should yield more tightly focused families of genes worthy of study (35). In particular, morphologies of unknown biological significance are likely to lead to the study of entirely new pathways in the spirit of classic genetic screens.

The approach described here is compatible with automated image analysis systems and, importantly, is robust to the occasional segmentation errors produced by such systems. Previous work has demonstrated that machine learning algorithms can be successfully trained by using all cells from positive and negative control samples to create a training set, even for some phenotypes that cannot be visually distinguished by humans (25). Here we showed that, whereas this approach can be successful for highly penetrant phenotypes (Fig. S4A), it is not suitable when the phenotype is less penetrant (Fig. S4 B–D and Fig. S9). We have addressed these challenging situations, thus enabling screens for low-penetrance phenotypes that lack positive control samples. Even when positive control samples are available, leveraging the user's visual perception to select individual example cells helps prevent the machine learning algorithm from focusing on aspects of morphology that are irrelevant to the biological question at hand or from becoming tuned to cells that display some complex combination of phenotypes as the positive control samples (i.e., pleiotropic effects) rather than the specific phenotype of interest.

