

[Click here to view linked References](#)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## **Title: A dataset of images and morphological profiles of 30,000 small-molecule treatments using the Cell Painting assay**

Mark-Anthony Bray<sup>1</sup> (equal contributor) [mbray@broadinstitute.org](mailto:mbray@broadinstitute.org)

Sigrun M. Gustafsdottir<sup>2</sup> (equal contributor) [sigrun.gustafsdottir@gmail.com](mailto:sigrun.gustafsdottir@gmail.com)

Vebjorn Ljosa<sup>1</sup> (equal contributor) [vebjorn@ljosa.com](mailto:vebjorn@ljosa.com)

Shantanu Singh<sup>1</sup> [shsingh@broadinstitute.org](mailto:shsingh@broadinstitute.org)

Katherine L. Sokolnicki<sup>1</sup> [kate.sokolnicki@gmail.com](mailto:kate.sokolnicki@gmail.com)

Joshua A. Bittker<sup>2</sup> [jbittker@broadinstitute.org](mailto:jbittker@broadinstitute.org)

Nicole E. Bodycombe<sup>2</sup> [nemmith@gmail.com](mailto:nemmith@gmail.com)

Vlado Dančik<sup>2</sup> [vdancik@broadinstitute.org](mailto:vdancik@broadinstitute.org)

Thomas P. Hasaka<sup>2</sup> [thasaka@gmail.com](mailto:thasaka@gmail.com)

C. Suk- Yee Hon<sup>2</sup> [cindyhon@broadinstitute.org](mailto:cindyhon@broadinstitute.org)

Melissa M. Kemp<sup>2</sup> [melissak.broad@gmail.com](mailto:melissak.broad@gmail.com)

Kejie Li<sup>2</sup> [kejie.li@biogen.com](mailto:kejie.li@biogen.com)

Deepika Walpita<sup>2</sup> [walpitad@janelia.hhmi.org](mailto:walpitad@janelia.hhmi.org)

Mathias J. Wawer<sup>2</sup> [mwawer@broadinstitute.org](mailto:mwawer@broadinstitute.org)

Todd R. Golub<sup>3</sup> [golub@broadinstitute.org](mailto:golub@broadinstitute.org)

Stuart L. Schreiber<sup>2</sup> [schreiber@broadinstitute.org](mailto:schreiber@broadinstitute.org)

Paul A. Clemons<sup>2</sup> [pclemons@broadinstitute.org](mailto:pclemons@broadinstitute.org)

Alykhan F. Shamji<sup>2</sup> [ashamji@broadinstitute.org](mailto:ashamji@broadinstitute.org)

Anne E. Carpenter<sup>1\*</sup> [anne@broadinstitute.org](mailto:anne@broadinstitute.org), <http://www.broadinstitute.org/~anne/>

<sup>1</sup> Imaging Platform, Broad Institute of Harvard and MIT, Cambridge, MA, USA

<sup>2</sup> Center for the Science of Therapeutics, Broad Institute of Harvard and MIT, Cambridge, MA, USA

<sup>3</sup> Cancer Program, Broad Institute of Harvard and MIT, Cambridge, MA USA

\*To whom correspondence should be addressed

## ABSTRACT

**Background:** Large-scale image sets acquired by automated microscopy of perturbed samples enable a detailed comparison of cell states induced by each perturbation, such as a small molecule from a diverse library. Highly multiplexed measurements of cellular morphology can be extracted from each image and subsequently mined for a number of applications.

**Findings:** This microscopy data set includes 919,874 five-channel fields of view representing 30,616 tested compounds, available at 'The Cell Image Library' repository. It also includes data files containing morphological features derived from each cell in each image, both at the single-cell level and population-averaged (i.e., the per-image level); the image analysis workflows that generated the morphological features are also provided. Quality-control metrics are provided as metadata, indicating fields of view that are out-of-focus (blurry) or containing highly fluorescent material or debris. Lastly, chemical annotations are supplied for the compound treatments applied.

**Conclusions:** Because computational algorithms and methods for handling single-cell morphological measurements are not yet routine, the dataset serves as a useful resource for the wider scientific community applying morphological (image-based) profiling. The data set can be mined for many purposes, including small-molecule library enrichment and chemical mechanism-of-action studies, including target identification. Integration with genetically-perturbed datasets could enable identification of small-molecule mimetics of particular disease- or gene-related phenotypes that could be useful as probes or potential starting points for development of future therapeutics.

## KEYWORDS

phenotypic profiling, high-content screening, image-based screening, cellular morphology, small-molecule library, U2OS

## DATA DESCRIPTION

### Purpose of data acquisition

High-throughput quantitative analysis of cellular image data has led to critical insights across many fields in biology[1, 2]. While microscopy has enriched our understanding of biology for centuries, only recently has robotic sample preparation and microscopy equipment become widely available, together with large libraries of chemical and genetic perturbagens. Concurrently, the advent of high-throughput imaging has also become an engine for pharmacological screening and basic research, by allowing multiparametric image-based interrogation of physiological processes at a large scale[3, 4].

A typical imaging assay uses several fluorescent probes (or fluorescently-tagged proteins) simultaneously to stain cells, each labeling distinct cellular components in each sample. In this way, the morphological characteristics (or "phenotype") of cells, tissues, or even whole organisms can be examined, along with the concomitant changes induced by the perturbants of choice[5–7].

Phenotypic profiling has emerged as a powerful tool to discern subtle differences among treated samples in a relatively unbiased manner. In contrast to a screening strategy, where a usually limited number of features are quantified to select for a known cellular phenotype, profiling relies on collecting a large suite of per-cell

1 morphological features and then using statistical analysis to uncover latent morphological patterns  
2 (“signatures”) by which the perturbations can be characterized. The “Cell Painting” assay used for the dataset  
3 presented here uses fluorescent markers to broadly stain a number of cellular structures in high-throughput  
4 format; automated software extracts the single-cell image-based morphological features. Further analysis then  
5 aggregates the data as multivariate profiles of these features to compare signatures among sample  
6 treatments.  
7  
8  
9

10 The applications of image-based profiling are many and diverse. A dataset comprising small-molecule  
11 perturbations, as presented here, can be used for small-molecule library enrichment (to create smaller libraries  
12 while retaining high diversity of phenotypic impact) and small-molecule mechanism-of-action studies, including  
13 target identification. Integration of this dataset with datasets resulting from other types of perturbations (e.g.,  
14 patient cell samples or genetically-perturbed samples) enables identification of small-molecule mimetics of  
15 particular disease- or gene-related phenotypes that could be useful as probes or potential starting points for  
16 development of future potential therapeutics.  
17  
18  
19  
20

## 21 **Data acquisition protocol and quality control**

22

23 To maximize the morphological information extracted from a single assay, we sought to “paint the cell” with as  
24 many distinct fluorescent morphological markers as possible simultaneously. Balancing technical and cost  
25 considerations, we developed the Cell Painting assay protocol in which cells are stained for eight major  
26 organelles and sub-compartments, using a mixture of six well-characterized fluorescent dyes suited for use in  
27 high-throughput (Fig. 1).  
28  
29  
30

31 The protocols for staining and imaging have been described in detail elsewhere[8, 9]. Briefly, U2OS cells were  
32 plated in 384-well plates, then treated with each of 30,616 compounds in quadruplicate. Of these compounds,  
33 10,162 compounds came from the Molecular Libraries Small Molecule Repository (MLSMR), 2,222 were  
34 drugs, natural products, and small- molecule probes that are part of the Broad Institute known bioactive  
35 compound collection, 274 were confirmed screening hits from the Molecular Libraries Program (MLP), and  
36 19,137 were novel compounds derived from diversity-oriented synthesis (DOS). Live cell staining was first  
37 performed to stain the mitochondria using MitoTracker. After incubation, the cells were fixed with  
38 formaldehyde, permeabilized with Triton X-100, and stained with the remaining dyes to identify the nucleus  
39 (Hoechst), nucleoli and cytoplasmic RNA (SYTO 14), endoplasmic reticulum (concanavalin A), Golgi and  
40 plasma membrane (wheat germ agglutinin), and the actin cytoskeleton (phalloidin). Each of the 413 multi-well  
41 plates was imaged using an ImageXpress Micro XLS automated microscope (Molecular Devices), with five  
42 fluorescent channels at 20× magnification, and 6 fields of view (sites) imaged per well (Table 1). Each image  
43 channel was then stored as a separate, grayscale image file in 16-bit TIF format. All raw image data is publicly  
44 available at ‘The Cell Image Library’ repository[10].  
45  
46  
47  
48  
49

50 The dataset available at GigaDB consists of the processed data derived from the acquired raw image data  
51 (Table 2; see also Additional File and “Availability of supporting data” Section). The quantitative analysis of the  
52 images used a three-step workflow using the modular open-source software CellProfiler[11]. First, an  
53 illumination pipeline estimated the heterogeneities in the spatial fluorescence distribution introduced by the  
54 microscope optics. This approximation was calculated on a per-plate basis for each channel and yielded a  
55 collection of illumination correction functions (ICFs) for later use in intensity correction; we have found that this  
56 approach not only aids in cell identification but also improves accuracy in signature classification[12]. Second,  
57 a quality control pipeline identified and labeled images with aberrations such as saturation artifacts and focal  
58 blur as described previously[13, 14] (see also Additional File 1). Finally, a feature-extraction pipeline applied  
59 the ICFs to correct each channel, identified the nuclei, cell body and cytoplasm, and extracted the  
60  
61  
62  
63  
64  
65

1 morphological features for each cell, depositing the results into a database for downstream analysis. The  
2 extracted features include a broad array of cellular shape and adjacency statistics, as well as intensity and  
3 texture statistics that are measured in each channel. The pipelines, ICFs, and extracted morphological data are  
4 provided the GigaDB repository.  
5

6  
7 Many approaches exist to creating per-sample profiles based on the per-cell data from each replicate; we have  
8 found that producing profiles simply by averaging the cellular features across all cells for each well yielded  
9 good results in characterizing compounds[15]. These profiles are provided in GigaDB along with a listing of  
10 chemical annotations for the compounds applied. The downstream analysis of morphological profiling data is a  
11 field very much in flux at present; our own laboratory has a GitHub repository of R scripts for this purpose at  
12 <https://github.com/CellProfiler/cytominr>.  
13  
14

## 15 16 **Potential uses**

17  
18 Phenotypic profiling provides a powerful means for assessing the biological impact of molecular or genetic  
19 perturbagens in the context of a biological system. The images and annotations provided in this Data Note have  
20 already been used in two published analyses from our own group; unsupervised clustering of a subset of 1,601  
21 bioactive compounds in a proof-of-principle study of compound mechanism of action  
22 (<https://www.broadinstitute.org/bbbc/BBBC022/>)[16] and small-molecule library enrichment based on the full  
23 set of 31,795 small molecules, a study in which morphological profiles successfully selected compound  
24 subsets with higher performance diversity than randomly-selected compounds[8].  
25  
26  
27  
28

## 29 30 **AVAILABILITY AND REQUIREMENTS**

- 31
- 32 ● **Project name:** CellProfiler (RRID:nif-0000-00280)
- 33 ● **Project home page:** <http://www.cellprofiler.org>
- 34 ● **Operating system(s):** Platform independent
- 35 ● **Programming language:** Python
- 36 ● **Other requirements:** None
- 37 ● **License:** BSD
- 38 ● **Any restrictions to use by non-academics:** None
- 39
- 40
- 41
- 42

## 43 44 **AVAILABILITY OF SUPPORTING DATA**

45  
46 The raw image data described in this article is available at ‘The Cell Image Library’ repository as Plates 24278-  
47 26794[10]. The remainder of the dataset supporting the results of this article is available in the GigaScience  
48 repository, GigaDB, [INSERT DOI, HYPERLINK][REF]. All data relating to a plate are within a folder named  
49 after the unique 5-digit identifier for each plate. This includes illumination correction functions, metadata related  
50 to sample treatment and image quality control, extracted morphological features, and profiles, with each data  
51 type existing in a separate sub-folder within the parent plate-ID folder (Table 2). Each of the plate folders has  
52 been packed as tape archive (TAR, .tar), before being compressed using GNU Gzip (.gz) and can be  
53 downloaded individually. Updates to the pipelines (e.g., to accommodate updated software versions or updated  
54 versions of the protocol) can be found at our Cell Painting wiki  
55 ([https://github.com/carpenterlab/2016\\_bray\\_natprot](https://github.com/carpenterlab/2016_bray_natprot)). Any publication arising from the use of the deposited data  
56 must acknowledge the source of the dataset.  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 **ETHICS APPROVAL AND CONSENT TO PARTICIPATE**

2  
3 Not applicable.

4  
5  
6 **CONSENT FOR PUBLICATION**

7  
8  
9 Not applicable.

10  
11  
12 **COMPETING INTERESTS**

13  
14  
15 The authors declare that they have no competing interests.

16  
17  
18 **FUNDING**

19  
20  
21 Research reported in this publication was supported in part by NSF CAREER DBI 1148823 (AEC).

22  
23  
24  
25 **AUTHOR CONTRIBUTIONS**

26  
27  
28 MAB and AEC drafted the manuscript. MJW, SMG, CSYH, JAB, TRG, AEC, AFS, SLS, and PAC designed  
29 research. SMG, VL, MAM, KLS, MMK, TPH, and JAB performed research. MJW, KL, VL, NEB, MAB, VD,  
30 AEC, AFS, SLS, PAC, SS and MAB analyzed data. CSYH served as a Project Manager.

31  
32  
33  
34 **ACKNOWLEDGMENTS**

35  
36  
37 The authors thank David Orloff and Willy Wong from 'The Cell Image Library' for their efforts in assisting  
38 upload and annotation of the image portion of the dataset. We also thank Mohammad Hossein Rohban for  
39 expanding and contributing the compound annotations.

40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## FIGURES

**Figure 1:** Sample images from the small-molecule Cell Painting experiment using U2OS cells. Images are shown from a DMSO well (negative control, top row) and a parabendazole well (bottom row). The columns display the five channels imaged in the Cell Painting assay protocol; see Table 1 for details about the stains and channels imaged.

1 **TABLES AND CAPTIONS**  
 2  
 3  
 4

5 **Table 1:** Details of dyes, stained cellular sub-compartments and channels imaged in the Cell Painting assay  
 6

Dye	Organelle or cellular component	Channel name	
		CellProfiler	ImageXpress
Hoechst 33342	Nucleus	DNA	w1
Concanavalin A/Alexa Fluor 488 conjugate	Endoplasmic reticulum	ER	w2
SYTO 14 green fluorescent nucleic acid stain	Nucleoli, cytoplasmic RNA	RNA	w3
Phalloidin/Alexa Fluor 568 conjugate, wheat germ agglutinin (WGA)/Alexa Fluor 555 conjugate	F-actin cytoskeleton, Golgi, plasma membrane	AGP	w4
MitoTracker Deep Red	Mitochondria	Mito	w5

36 The CellProfiler channel name refers to the name given by the software to each channel; this nomenclature  
 37 also applies to the naming of the extracted morphological features. The ImageXpress channel name refers to  
 38 the text in the raw image file name identifying the acquired wavelength.  
 39  
 40  
 41  
 42  
 43  
 44  
 45  
 46  
 47  
 48  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65

**Table 2:** Summary of the raw and intermediately processed data included in this Data Descriptor, and nomenclature in GigaDB.

Data item	Location	Description
Raw fluorescence images	The Cell Image Library[10]	Five fluorescence channels, recorded at 6x fields of view per well at 20X magnification. The experiment comprises 413 plates (Plates 24278-26794)
CellProfiler pipelines	GigaDB: pipelines.zip	CellProfiler software was used to correct for uneven illumination, perform quality control and segment cells into nuclei, cell body and cytoplasmic sub-compartments and measure morphological features for each cell.
Illumination correction functions (ICFs)	GigaDB: <plate_ID>/illumination_correction_functions	An ICF is an estimation of the spatial illumination distribution introduced by the microscopy optics. There is one image per channel, per plate.
Quality control metadata	GigaDB: <plate_ID>/quality_control	Each field of view is assessed for the presence of two artifacts (focal blur and saturated objects), and assigned a label of 1 if present, and 0 if not.
Extracted morphological features	GigaDB: <plate_ID>/extracted_features	Two data tables consisting of (a) per-image cellular statistics (e.g. ,cell count) and experimental metadata, and (b) per-cell size, shape, intensity, textural and adjacency statistics measured for the nuclei, cytoplasm, and cell body. Includes a MySQL dump file for importing the data tables into a MySQL database.
Morphological profiles	GigaDB: <plate_ID>/profiles	Per-well averages of each extracted morphological feature computed across the cells.
Image curation statistics	GigaDB: image_curation_statistics.csv	A summary of image statistics, such as the number of images, wells, and sites in the plates archived at The Cell Image Library, the number of sites with quality measures and the number of wells with morphological profiles.
Chemical annotations	GigaDB: chemical_annotations.csv	Chemical annotations including the compound names, SMILES, and PubChem identifiers (CID/SID)

<plate\_ID> refers to the 5-digit plate ID assigned by the ImageXpress microscope system.

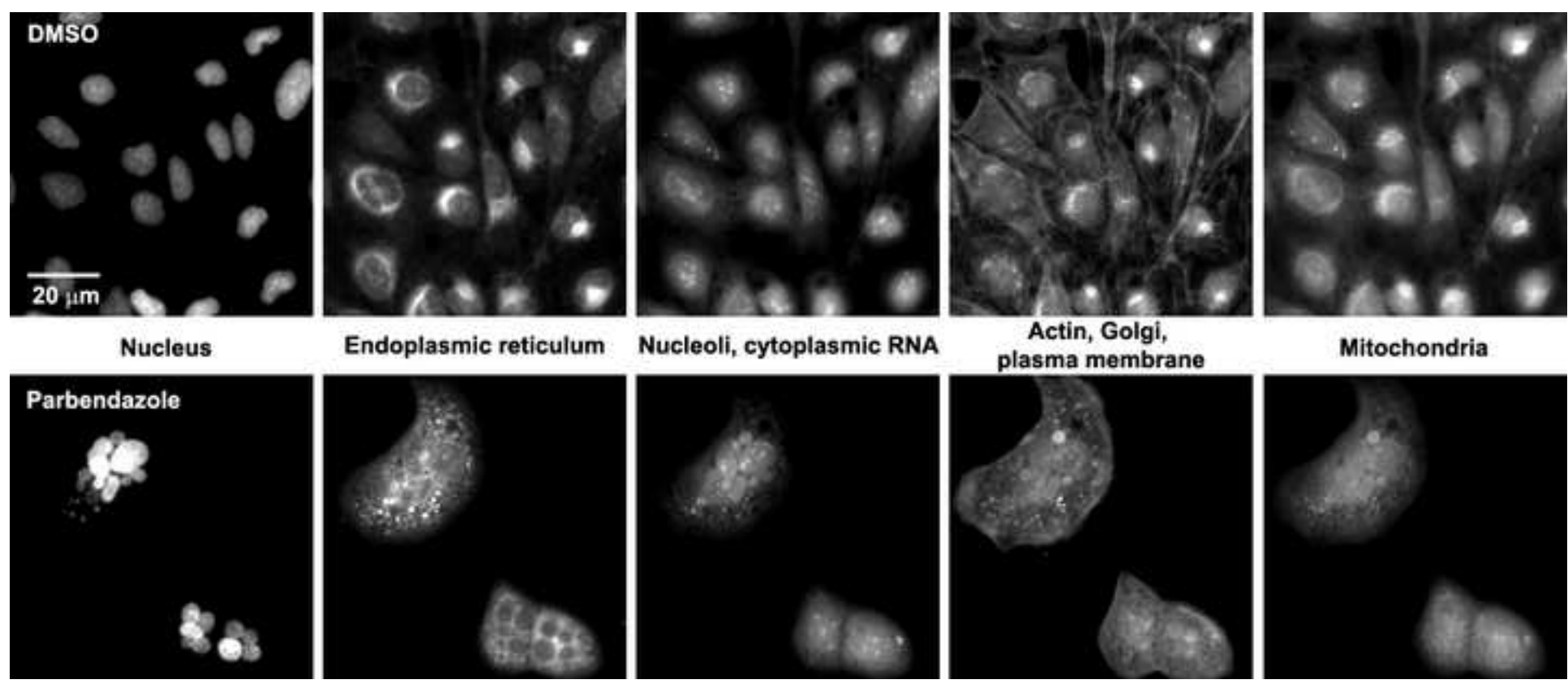


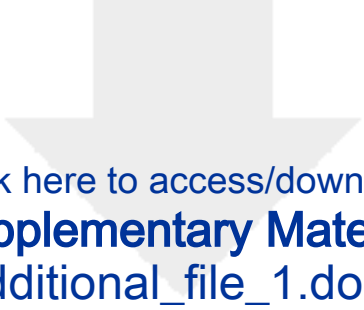
## REFERENCES

1. Conrad C, Gerlich DW: **Automated microscopy for high-content RNAi screening.** *J Cell Biol* 2010, **188**:453–461.
2. Thomas N: **High-content screening: a decade of evolution.** *J Biomol Screen* 2010, **15**:1–9.
3. Bickle M: **The beautiful cell: high-content screening in drug discovery.** *Anal Bioanal Chem* 2010, **398**:219–226.
4. Boutros M, Heigwer F, Laufer C: **Microscopy-Based High-Content Screening.** *Cell* 2015, **163**:1314–1325.
5. Levsky JM, Singer RH: **Gene expression and the myth of the average cell.** *Trends Cell Biol* 2003, **13**:4–6.
6. Snijder B, Pelkmans L: **Origins of regulated cell-to-cell variability.** *Nat Rev Mol Cell Biol* 2011, **12**:119–125.
7. Altschuler SJ, Wu LF: **Cellular heterogeneity: do differences make a difference?** *Cell* 2010, **141**:559–563.
8. Wawer MJ, Li K, Gustafsdottir SM, Ljosa V, Bodycombe NE, Marton MA, Sokolnicki KL, Bray M-A, Kemp MM, Winchester E, Taylor B, Grant GB, Hon CS-Y, Duvall JR, Wilson JA, Bittker JA, Dančik V, Narayan R, Subramanian A, Winckler W, Golub TR, Carpenter AE, Shamji AF, Schreiber SL, Clemons PA: **Toward performance-diverse small-molecule libraries for cell-based phenotypic screening using multiplexed high-dimensional profiling.** *Proc Natl Acad Sci U S A* 2014, **111**:10911–10916.
9. Bray M-A, Singh S, Han H, Davis CT, Borgeson B, Gustafsdottir SM, Gibson CC, Carpenter AE: **Cell Painting, an image-based assay for morphological profiling.** *Nat. Protoc.* .
10. **Human U2OS cells - compound cell-painting experiment**  
[[http://www.cellimagelibrary.org/pages/project\\_20269](http://www.cellimagelibrary.org/pages/project_20269)]
11. Kametsky L, Jones TR, Fraser A, Bray M-A, Logan DJ, Madden KL, Ljosa V, Rueden C, Eliceiri KW, Carpenter AE: **Improved structure, function and compatibility for CellProfiler: modular high-throughput image analysis software.** *Bioinformatics* 2011, **27**:1179–1180.
12. Singh S, Bray M-A, Jones TR, Carpenter AE: **Pipeline for illumination correction of images for high-throughput microscopy.** *J Microsc* 2014, **256**:231–236.
13. Bray M-A, Fraser AN, Hasaka TP, Carpenter AE: **Workflow and metrics for image quality control in large-scale high-content screens.** *J Biomol Screen* 2012, **17**:266–274.
14. Bray M-A, Carpenter AE: **Quality control for high-throughput imaging experiments using machine learning in CellProfiler.** In *Methods in Molecular Biology Series: High Content Imaging, Analysis and Screening: Applications in Basic Science and Drug Discovery.* Edited by Paul A. Johnston PA, Trask OJ Jr. Humana Press,; .
15. Ljosa V, Caie PD, Ter Horst R, Sokolnicki KL, Jenkins EL, Daya S, Roberts ME, Jones TR, Singh S, Genovesio A, Clemons PA, Carragher NO, Carpenter AE: **Comparison of methods for image-based profiling of cellular morphological responses to small-molecule treatment.** *J Biomol Screen* 2013, **18**:1321–1329.
16. Gustafsdottir SM, Ljosa V, Sokolnicki KL, Anthony Wilson J, Walpita D, Kemp MM, Petri Seiler K, Carrel HA, Golub TR, Schreiber SL, Clemons PA, Carpenter AE, Shamji AF: **Multiplex cytological profiling assay**

1 **to measure diverse cellular states.** *PLoS One* 2013, **8**:e80999.

2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65





Click here to access/download  
**Supplementary Material**  
additional\_file\_1.docx



415 Main Street  
Cambridge, MA 02142  
T 617-714-7000 F 617-714-8972  
[www.broadinstitute.org](http://www.broadinstitute.org)

Dear editors,

We are pleased to submit the manuscript, "A dataset of images and morphological profiles of 30,000 small-molecule treatments using the Cell Painting assay" for your consideration in GigaScience as a Data Note.

It describes a valuable image set, associated metadata, and extracted numerical morphological features, based on ~1 million fields of view from ~30,000 tested small molecules. We are unaware of any comparably sized publicly available image set involving small-molecule perturbations in a single experiment; we therefore think this Data Note will be quite valuable to the field.

Indeed, we have received numerous requests for the raw image data and derived cellular measurements from multiple academic institutions and biotech/pharma companies. This interest is in part derived from the wide variety of applications for which morphological profiling of images from a rich assay like ours can be used. We therefore expect this Data Note to be of broad interest, and the manuscript we are submitting will now make these data publicly available and suitably annotated.

We suggest the following researchers to review the Data Note for their understanding of the challenges and value of large-scale image sets: Jason Swedlow ([jason@lifesci.dundee.ac.uk](mailto:jason@lifesci.dundee.ac.uk)), Peter Horvath ([peter.horvath@bc.biol.ethz.ch](mailto:peter.horvath@bc.biol.ethz.ch)), Michael Boutros ([m.boutros@dkfz.de](mailto:m.boutros@dkfz.de)), and Jan Ellenberg ([jan.ellenberg@embl.de](mailto:jan.ellenberg@embl.de)).

The authors declare that they have no competing interests; all authors have approved the manuscript for submission. The content of the manuscript has not been published, or submitted for publication elsewhere.

Sincerely,

Anne E. Carpenter, Ph.D.  
Director, Imaging Platform  
Broad Institute of Harvard and MIT