
Three million images and morphological profiles of cells treated with matched chemical and genetic perturbations

Srinivas Niranj Chandrasekaran¹, Beth A. Cimini¹, Amy Goodale¹, Lisa Miller¹

Maria Kost-Alimova¹, Nasim Jamali¹, John Doench¹, Briana Fritchman¹, Adam Skepner¹

Michelle Melanson¹, Daniel Kuhn², Desiree Hernandez¹, Jim Berstler¹, Hamdah Abbasi¹

David Root¹, Susanne E. Swalley³

Shantanu Singh¹, Anne E. Carpenter¹

{shsingh, anne}@broadinstitute.org

¹ Broad Institute of MIT and Harvard, 415 Main St, Cambridge, MA, 02142

² Merck Healthcare KGaA, Frankfurter Str. 250, 64293 Darmstadt, Germany

³ Biogen, Inc., 125 Broadway Street, Cambridge, MA 02139.

Abstract

1 We present a new, carefully designed and well-annotated dataset of images and
2 image-based profiles of cells that have been treated with chemical compounds and
3 genetic perturbations. Each gene that is perturbed is a known target of at least two
4 compounds in the dataset. The dataset can thus serve as a benchmark to evaluate
5 methods for predicting similarities between compounds and between genes and
6 compounds, measuring the effect size of a perturbation, developing style-transfer
7 methods to predict one experimental condition from another, and more generally,
8 learning effective representations for measuring cellular state from microscopy
9 images.

10 1 Introduction

11 Computer vision has benefitted dramatically from the revolution in deep learning. Biomedical
12 research is an exceptionally satisfying domain on which to apply advances in machine learning, and
13 yet deep learning applied to images in the biomedical domain has been relatively limited to medical
14 imaging from patients, including X rays and MRI, PET, and CT scans. By comparison, deep-learning
15 based image analysis for cell biology has generally focused on segmentation [1, 2]; whereas feature
16 extraction and applications have lagged behind [3].

17 One cell biology method – image-based profiling of cell samples – is proving increasingly useful
18 for the discovery of disease underpinnings and useful drugs [4]. In image-based profiling, human
19 cells are cultured in samples of a few hundred cells, each sample treated with a different chemical or
20 genetic perturbation. The resulting morphology (visual appearance) of each sample is compared by
21 microscopy to identify meaningful differences and similarities. Among many others, applications

22 include: (a) identifying the mechanisms of a disease by comparing cells from patients with a disease
23 to those without the disorder, (b) identifying the impact of a drug by comparing cells treated with it
24 to untreated cells, (c) identifying gene functions or the impact of chemicals on cells by unsupervised
25 clustering of large sets of samples to determine relationships among the perturbations tested in the
26 experiment. Thus, image-based profiling can reveal new targets for diseases, potential therapeutics,
27 and toxicities for particular compounds.

28 The vast majority of research using image-based profiling uses classical segmentation and feature
29 extraction; deep learning methods are beginning to be explored [3] and there is much room for
30 advancement. Historically, the lack of ground truth has been a major limiting factor in the field, as
31 the “correct” high-dimensional profile of a given sample is unknown, and the “correct” relationships
32 among most genes and compounds are unknown. Image-based profiling applications typically can be
33 described as representation learning tasks; if samples are represented optimally and ideal distance
34 metrics are applied, then biologically meaningful differences between samples will be detectable and
35 technical artifacts will be suppressed.

36 To push forward advancements in this field, we assembled a consortium of ten pharmaceutical
37 companies, two non-profit institutions, and several supporting companies, known as the JUMP-Cell
38 Painting Consortium (Joint Undertaking in Morphological Profiling). After extensive optimization
39 of the main assay used in image-based profiling, called the Cell Painting assay [5], this Consortium
40 created a ground truth dataset to move methods in the field forward. We selected and curated a set
41 of genes and compounds with (relatively) known relationships among each other, and designed an
42 experimental layout to enable testing and comparing methods to quantify their relationships.

43 Here, we describe our design and creation of this dataset from a single large experiment comprising
44 nearly three million images and over seventy five million single cells, called CPJUMP1, which
45 contains chemical and genetic perturbation pairs that target the same genes in cells. It allows
46 exploring a number of technical and biological parameters that might affect matching ability and
47 testing computational strategies to match samples to each other and thus uncover valuable biological
48 relationships.

49 **2 Related datasets**

50 We are not aware of any other Cell Painting image-based datasets that include pairs of genetic
51 and chemical perturbations with their relationships to each other annotated, and executed in par-
52 allel so as to minimize technical variations that may confound the signal. Nevertheless, other
53 Cell Painting datasets are public and may be useful to the community, for example as training
54 data for self-supervised feature extraction methods. These single-perturbation-type experi-
55 ments include several datasets from the Carpenter-Singh laboratory (available through the Im-
56 age Data Resource [6] at [https://idr.openmicroscopy.org/search/?query=Publication%
57 20Authors:Carpenter](https://idr.openmicroscopy.org/search/?query=Publication%20Authors:Carpenter) and the 2018 CytoData challenge [https://github.com/cytodata/
58 cytodata-hackathon-2018](https://github.com/cytodata/cytodata-hackathon-2018)), one from the New York Stem Cell Foundation [7] and several from
59 Recursion, a clinical-stage biotechnology company (available at <http://rxrx.ai>).

60 **3 Data acquisition**

61 **3.1 Compound and gene selection**

62 Our dataset consists of images and profiles of cells that were perturbed separately by chemical
63 and genetic perturbations, where both sets were chosen based on known relationships among them.
64 Chemical perturbations are small molecules (i.e. chemical compounds) that modulate the function
65 of cells while the genetic perturbations are either open reading frames (ORFs) that can overexpress
66 genes (i.e. yield more of the gene’s product in the cell) or guide RNAs that mediate CRISPR-Cas9
67 (clustered regularly interspaced short palindromic repeats) that can knockdown gene function (i.e.
68 yield less of the gene’s product in the cell). Most compounds are thought to inhibit the function of
69 their target gene’s product, so we expect CRISPRs to generally correlate to (mimic) the corresponding
70 compound’s profile, whereas ORFs are generally expected to anti-correlate (oppose) the corresponding
71 small molecule’s profile, and ORFs and CRISPRs targeting the same gene should generally yield
72 opposite (anti-correlated) effects on the cells’ profiles. However, we strongly note that there will
73 be numerous exceptions given the non-linear behavior of many biological systems and a number of

74 distinct mechanisms by which these general principles may not hold. In fact, one aim of generating
75 this dataset is to quantify how often the expected relationships and directionalities occur.

76 We derived the list of compounds from Broad’s Drug Repurposing Hub dataset [8], a curated and
77 annotated collection of FDA-approved drugs, clinical trial drugs, and pre-clinical tool compounds.
78 The genes perturbed by genetic perturbations were chosen because they are the annotated targets of
79 the compounds. We filtered the Repurposing Hub compounds using several criteria, of which three
80 are important:

- 81 1. The compounds should target genes that belong to diverse gene families (Table 1). This is
82 because the ideal methods would work well for many different biological pathways, not just
83 a few that are well-characterized and/or easy to predict.
- 84 2. Each gene should be targeted by at least two compounds, so that gene-compound matching
85 and compound-compound matching can both be performed using the dataset.
- 86 3. We additionally considered applying the constraint that each compound should target only a
87 single gene. However, this criterion is difficult to achieve due to polypharmacology (Table
88 2), which is the property for compounds to bind and impact many different gene products in
89 the cell; this is especially common for protein kinase inhibitors in the dataset. Instead, we
90 only filtered out the so-called “historical compounds” listed in the Chemical Probes Portal
91 [9], comprising compounds that are known to be quite non-selective (or not sufficiently
92 potent) compared with other available chemical probes.

93 Our list of compounds and genes also includes both negative and positive controls. The negative
94 controls for each perturbation modality are:

- 95 • Compounds: DMSO (Dimethyl sulfoxide), which is the solvent for all the compounds
96 studied. In other words, all samples will have DMSO added at the same concentration but
97 the negative controls have no additional compound added.
- 98 • ORFs: 15 ORFs with the weakest signature in previous image-based profiling experiments
99 (Rohban et al., 2017).
- 100 • CRISPRs: 30 CRISPR guides that target an intergenic site (cutting controls, $n = 3$) or don’t
101 have a target sequence that exists in human cells (non-cutting controls, $n = 27$).

102 There are three types of compound positive controls in our list. First, we included chemical probes
103 that are very well-studied and (unlike most compounds) are known to very selectively modulate the
104 genes that they target [9]. Second, we included compounds that strongly correlate with the correct
105 genetic perturbation in previous image-based profiling experiments with ORFs [10] and compounds
106 [11]. Finally, we included a set of maximally diverse pairs of compounds with strong intra-pair and
107 weak inter-pair correlations.

108 A complete description of the filtering criteria and the procedure for selecting positive and negative
109 controls is available at <https://github.com/jump-cellpainting/JUMP-Target/>.

110 In the future, a commercial vendor may offer the compound set so others can test the same perturba-
111 tions in other contexts for comparison.

112 3.2 Plate layout design

113 After applying the filters and including positive controls, we selected a total of 306 compounds and
114 160 genes such that they could fit into three 384-well plates, one plate per perturbation modality
115 (compounds, ORFs and CRISPRs). Apart from a dozen or so compounds, most compounds are in
116 singlicate. All plates included negative controls as discussed above: $n=4$ replicates of the 15 ORF
117 negative controls in the ORF plate, $n=2$ replicates of the 30 CRISPR negative controls in the CRISPR
118 plate, and $n=64$ replicates of DMSO in the compound plate. On the CRISPR plate, there are two
119 guides per gene, each arrayed in its own well and kept separate, with no within-plate replicates. In
120 the case of the ORF plate, for which there was only one perturbation reagent per gene, there are two
121 replicates per plate.

122 We also considered the impact of edge effects, or plate-layout effects, in our design. Edge effects are
123 the technical artifact whereby different samples will yield different behavior depending on where

Table 1: **Number of gene families with a given number of gene targets.** To maximize the diversity of genes, the genes were chosen such that most gene families (n=92) have only a single gene targeted in the final list.

Number of gene targets (N)	Number of gene families with N gene targets in the final list
1	92
2	16
3	2

Table 2: **Number of compounds with a given number of gene targets.** The compounds were chosen such that most compounds (n=218) in the final list are annotated as having only a single target.

Number of gene targets (N)	Number of compounds in the final list targeting N gene targets
1	218
2	49
3	23
4	7
5	4
6	3
7	1
8	1

124 they are located on a plate; generally this is most observed in the outer two rows and columns of the
 125 plate, and the problem persists despite efforts to mitigate it experimentally (Lundholt, Scudder and
 126 Pagliaro, 2003). While designing the plate layout, we divided the plate into outer and inner wells
 127 where the outer wells are the two rows and columns closest to the edge of the plate and the inner
 128 wells are the rest of the wells on the plate. Then we applied the following constraints in order to
 129 minimize the impact of edge effects:

- 130 1. Both of the compounds that target the same gene will either be in the inner wells or in the
 131 outer wells. They will not be split such that one of the compounds is in the inner well while
 132 the other is in the outer well.
- 133 2. The gene target of outer well compounds will be in the outer wells of the genetic perturbation
 134 plate.
- 135 3. All the positive control compounds are in the inner wells.

136 If preferable, an analysis can be constrained to the inner wells only, to ensure that edge effects have
 137 minimal influence on the results.

138 3.3 Experimental conditions

139 We acquired our data under the following experimental conditions:

- 140 1. Four replicate plates of compounds and CRISPRs and two replicate plates of ORFs (which,
 141 as mentioned, contain two replicates within each plate) at two time points and two cell lines
 142 each. The short and long time points were different for each perturbation type: compounds
 143 (24-hour, 48-hour), ORFs (48-hour, 96-hour) and CRISPRs (96-hour, 144-hour). The two
 144 cell lines were U2OS and A549.
- 145 2. One plate of the A549 96-hour ORF plate where the cells have been additionally treated
 146 with Blasticidin (a drug that kills cells that have not been properly infected with the genetic
 147 reagent).
- 148 3. Two replicate plates of the A549 144-hour CRISPR plate where the cells have been addi-
 149 tionally treated with Puromycin (a drug that kills cells that have not been properly infected
 150 with the genetic reagent).
- 151 4. Two replicate plates of the A549 48-hour compound plate with 20% higher cell seeding
 152 density than the baseline.

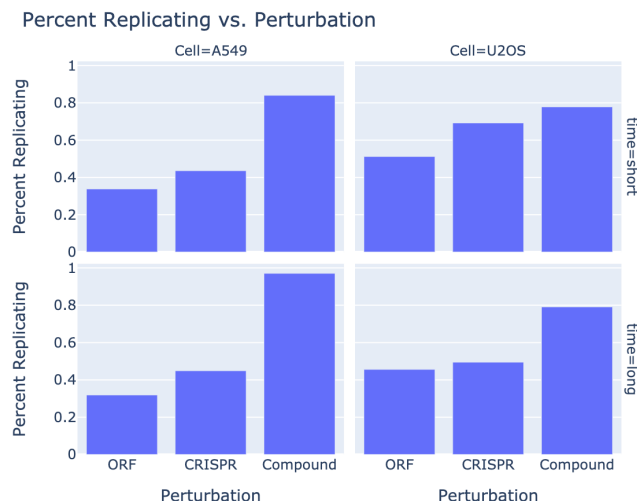


Figure 1: *Percent Replicating vs. perturbation modality*. Compounds have a stronger within-replicate correlation compared to ORFs and CRISPRs.

- 153 5. Two replicate plates of the A549 48-hour compound plate with 20% lower cell seeding
 154 density than the baseline.
- 155 6. Four replicate plates of the A549 24-hour compound plate were imaged six additional times
 156 to test photobleaching from repeated imaging.
- 157 7. Two replicates of the ORF plates in U2OS and A549 at 96-hour and 144-hour were imaged
 158 four additional times, once on each of days 1, 4, 14, 28 after the first imaging, to test the
 159 stability of samples over time.

160 4 Potential uses

161 The CPJUMP1 dataset was designed to test several experimental conditions to determine which yield
 162 the highest signals and best matching ability. We will establish best practices for the laboratory work
 163 based on our analysis of these results, not further detailed here (Cimini et al., in preparation). Here
 164 we focus on the applications that are most of interest to a machine learning audience.

165 4.1 Benchmarking perturbation-detection methods

166 Detecting which samples are measurably different from negative controls is one task that often
 167 precedes other useful applications, and is equivalent to measuring the effect size. For example, a set
 168 might be filtered by this criterion before embarking on subsequent laboratory experiments, or prior to
 169 training a model, or other analysis that could be confounded by noisy signals. It can also be useful
 170 for determining what experimental protocol or computational analysis pipeline to use among several
 171 alternatives. It should be noted that even given perfect computational methods for feature extraction,
 172 batch correction, and profile comparison, not all samples will be detectably different from negative
 173 controls for several biological reasons. For example, a drug or genetic perturbation may only impact
 174 cell morphology in a particular cell type, under particular environmental conditions, at a particular
 175 time, or if particular stains were used, conditions which may not have been met in the experiment.

176 To detect the number of samples with a measurably distinct phenotype, we estimated *Percent*
 177 *Replicating* (Figure 1), which is the proportion of samples that are distinct from the null distribution
 178 built from samples that are non-replicates. A sample is considered to have a detectable signature if
 179 the median of the correlation between the replicates of the sample is greater than the 95th percentile
 180 of the null distribution. In other words, *Percent Replicating* is the True Positive Rate if the False
 181 Positive Rate is set to 5%, for a binary classification problem where replicates make up the positive
 182 class and non-replicates make up the negative class.

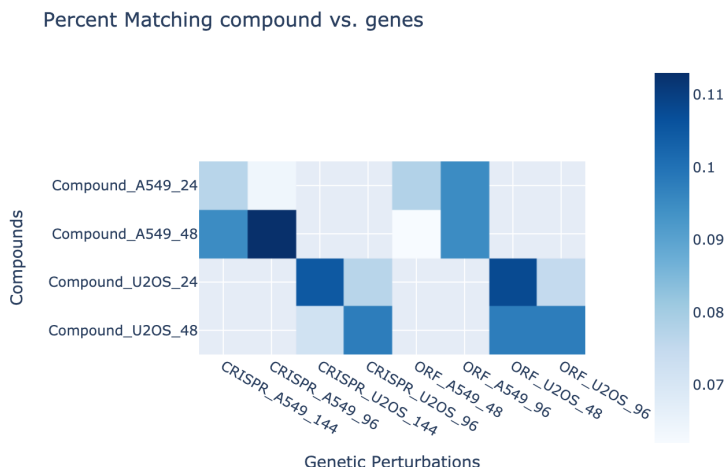


Figure 2: **Percent Matching between compounds and genetic perturbations.** Axis labels include the cell lines (A549 or U2OS) and the timepoints (48 hour, 96 hour, and 144 hour).

183 4.2 Benchmarking gene-compound matching methods

184 This dataset presents a unique opportunity to match profiles of perturbations across modalities
 185 (chemical versus genetic), because genes in this dataset that are targeted by two types of genetic
 186 perturbations (ORF and CRISPR) are also targeted by two compounds. To establish a baseline
 187 approach to match profiles across modalities, we computed the Pearson correlation between all
 188 chemical and genetic perturbation pairs. We then evaluated the performance of our approach by
 189 estimating *Percent Matching* (Figure 2), which is the proportion of “true” connections (chemical-
 190 genetic perturbation pairs that target the same gene) that are distinct from a null distribution built
 191 from “false” connections (chemical-genetic perturbation pairs that are not known to target the same
 192 gene). A true connection is considered to be correctly detected if its correlation is greater than the
 193 95th percentile of the null distribution. In other words, *Percent Matching* is the True Positive Rate if
 194 the False Positive Rate is set to 5%, for a binary classification problem where the true connections
 195 make up the positive class and the false connections make up the negative class.

196 The baseline results show that there is a signal in this dataset for matching chemical and genetic
 197 perturbations that target the same gene (7-11%, against a false positive rate of 5%), but there is much
 198 room for improvement. It should be strongly noted, though, that significant time and resources can
 199 be required to identify the target of a compound, and similarly to identify compounds that target a
 200 particular gene. Therefore, these low rates may already be highly meaningful, and improvements in
 201 image representations and measuring similarities could have a major impact on the pharmaceutical
 202 industry.

203 Given this dataset also has pairs of compounds targeting the same gene, it can also be used to test
 204 compound-compound matching.

205 4.3 Benchmarking style transfer methods

206 The design of CPJUMP1 included multiple cell types, timepoints, modalities (compound, ORF, and
 207 CRISPR), imaging conditions, and selection conditions. This allows the unusual opportunity to
 208 attempt prediction of one experimental condition from another. There are many potential combinations
 209 here, so we do not provide a baseline but simply point out this possibility to the interested researcher.

210 5 Code and Data availability

211 Cell images, morphological profiles, image analysis pipelines, profile generation pipelines, plate maps
 212 and plate and compound metadata are available online at <https://broad.io/neurips-cpjump1>.

213 The data used to generate the figures are available online. Figure 1: [https://github.com/
214 jump-cellpainting/neurips-cpjump1/tree/main/analysis#percent-replicating](https://github.com/jump-cellpainting/neurips-cpjump1/tree/main/analysis#percent-replicating)
215 and Figure 2: [https://github.com/jump-cellpainting/neurips-cpjump1/tree/main/
216 analysis#percent-matching-across-modalities](https://github.com/jump-cellpainting/neurips-cpjump1/tree/main/analysis#percent-matching-across-modalities).

217 **6 Methods**

218 **6.1 Sample preparation and image acquisition**

219 The Cell Painting assay involves staining eight components of cells with six fluorescent dyes: nucleus
220 (Hoechst), nucleoli and cytoplasmic RNA (SYTO 14), endoplasmic reticulum (concanavalin A), Golgi
221 and plasma membrane (wheat germ agglutinin; WGA), mitochondria (MitoTracker), and the actin
222 cytoskeleton (phalloidin). We optimized the Cell Painting assay described in (Bray et al., 2016) by
223 changing the concentrations of Hoechst, phalloidin, concanavalin A and SYTO14 and combining dye
224 addition and dye permeabilization steps. These changes will be described in more detail in (Cimini et
225 al., in preparation) and are currently publicly available at [https://github.com/carpenterlab/
226 2016_bray_natprot/wiki#updates-to-the-cell-painting-protocol](https://github.com/carpenterlab/2016_bray_natprot/wiki#updates-to-the-cell-painting-protocol). The images were
227 acquired across five fluorescent channels using a Perkin Elmer Opera Phenix HCI microscope at 20x
228 magnification.

229 **6.2 Image processing**

230 We used the CellProfiler [12] bioimage analysis software to process the images. We corrected
231 for variations in background intensity, and then segmented cells, distinguishing between nuclei
232 and cytoplasm. Then, across the various channels captured, we measure various features of cells
233 across several categories including fluorescence intensity, texture, granularity, density, location (see
234 <http://cellprofiler-manual.s3.amazonaws.com/CellProfiler-3.0.0/index.html> for
235 more details). Following the image analysis pipeline, we obtain more than 75 million cells and 5792
236 feature measurements.

237 **6.3 Image-based profiling**

238 We used *cytominer* (<https://cytominer.github.io/profiling-handbook/>) and *pycy-*
239 *tominer* workflows (<https://github.com/jump-cellpainting/profiling-recipe>) to pro-
240 cess the single cell features. We aggregated the single cell profiles by computing the mean. We then
241 normalized the averaged profiles by subtracting the median and dividing by the median absolute
242 deviation (m.a.d.) of each feature. This was done in two ways: using the median and m.a.d. of (i)
243 the negative control wells on the plate (used in the analysis shown here), and (ii) all the wells on the
244 plate. Finally, we filtered out redundant features as well as features with low variance. All the steps
245 in the profiling workflow were performed for each individual plate separately.

246 **Acknowledgments and Disclosure of Funding**

247 The authors appreciate the more than 100 scientists who have contributed to the organization and
248 scientific direction of the JUMP Cell Painting Consortium. We thank Max Macaluso (operations) and
249 Tanaz Abid (technical) at the Broad Institute for their assistance as well.

250 The authors gratefully acknowledge a grant from the Massachusetts Life Sciences Center Bits to
251 Bytes Capital Call program for funding the data production. We appreciate funding to support data
252 analysis and interpretation from members of the JUMP Cell Painting Consortium and from the
253 National Institutes of Health (NIH MIRA R35 GM122547 to AEC). The authors also gratefully
254 acknowledge the use of the PerkinElmer Opera Phenix High-Content/High-Throughput imaging
255 system at the Broad Institute, funded by the S10 Grant NIH OD-026839-01.

256 AEC has optional ownership interest in Recursion, a public biotechnology company using image-
257 based profiling for drug discovery. SES is an employee of Dewpoint Therapeutics. Daniel Kuhn is an
258 employee of Merck Healthcare KGaA, Darmstadt, Germany.

259 References

- 260 [1] E. Moen, D. Bannon, T. Kudo, W. Graf, M. Covert, and D. Van Valen, “Deep learning for
261 cellular image analysis,” *Nat. Methods*, vol. 16, no. 12, pp. 1233–1246, Dec. 2019.
- 262 [2] J. C. Caicedo, A. Goodman, K. W. Karhohs, B. A. Cimini, J. Ackerman, M. Haghghi, C. Heng,
263 T. Becker, M. Doan, C. McQuin, M. Rohban, S. Singh, and A. E. Carpenter, “Nucleus segmen-
264 tation across imaging experiments: the 2018 data science bowl,” *Nat. Methods*, vol. 16, no. 12,
265 pp. 1247–1253, Dec. 2019.
- 266 [3] A. Pratapa, M. Doron, and J. C. Caicedo, “Image-based cell phenotyping with deep learning,”
267 *Curr. Opin. Chem. Biol.*, vol. 65, pp. 9–17, May 2021.
- 268 [4] S. N. Chandrasekaran, H. Ceulemans, J. D. Boyd, and A. E. Carpenter, “Image-based profiling
269 for drug discovery: due for a machine-learning upgrade?” *Nat. Rev. Drug Discov.*, pp. 1–15,
270 2020.
- 271 [5] M.-A. Bray, S. Singh, H. Han, C. T. Davis, B. Borgeson, C. Hartland, M. Kost-Alimova, S. M.
272 Gustafsdottir, C. C. Gibson, and A. E. Carpenter, “Cell painting, a high-content image-based
273 assay for morphological profiling using multiplexed fluorescent dyes,” *Nat. Protoc.*, vol. 11,
274 no. 9, pp. 1757–1774, Sep. 2016.
- 275 [6] E. Williams, J. Moore, S. W. Li, G. Rustici, A. Tarkowska, A. Chessel, S. Leo, B. Antal,
276 R. K. Ferguson, U. Sarkans, A. Brazma, R. E. C. Salas, and J. R. Swedlow, “The image data
277 resource: A bioimage data integration and publication platform,” *Nat. Methods*, vol. 14, no. 8,
278 pp. 775–781, Aug. 2017.
- 279 [7] L. Schiff, B. Migliori, Y. Chen, D. Carter, C. Bonilla, J. Hall, M. Fan, E. Tam, S. Ahadi, B. Fis-
280 chbacher, A. Geraschenko, C. J. Hunter, S. Venugopalan, S. DesMarteau, A. Narayanaswamy,
281 S. Jacob, Z. Armstrong, P. Ferrarotto, B. Williams, G. Buckley-Herd, J. Hazard, J. Goldberg,
282 M. Coram, R. Otto, E. A. Baltz, L. Andres-Martin, O. Pritchard, A. Duren-Lubanski, K. Reggio,
283 NYSCF Global Stem Cell Array Team, L. Bauer, R. S. Aiyar, E. Schwarzbach, D. Paull, S. A.
284 Noggle, F. J. Monsma, M. Berndl, S. J. Yang, and B. Johannesson, “Deep learning and auto-
285 mated cell painting reveal parkinson’s disease-specific signatures in primary patient fibroblasts,”
286 Nov. 2020.
- 287 [8] S. M. Corsello, J. A. Bittker, Z. Liu, J. Gould, P. McCarren, J. E. Hirschman, S. E. Johnston,
288 A. Vrcic, B. Wong, M. Khan, J. Asiedu, R. Narayan, C. C. Mader, A. Subramanian, and T. R.
289 Golub, “The drug repurposing hub: a next-generation drug library and information resource,”
290 *Nat. Med.*, vol. 23, no. 4, pp. 405–408, Apr. 2017.
- 291 [9] C. H. Arrowsmith, J. E. Audia, C. Austin, J. Baell, J. Bennett, J. Blagg, C. Bountra, P. E.
292 Brennan, P. J. Brown, M. E. Bunnage, C. Buser-Doepner, R. M. Campbell, A. J. Carter,
293 P. Cohen, R. A. Copeland, B. Cravatt, J. L. Dahlin, D. Dhanak, A. M. Edwards, M. Frederiksen,
294 S. V. Frye, N. Gray, C. E. Grimshaw, D. Hepworth, T. Howe, K. V. M. Huber, J. Jin, S. Knapp,
295 J. D. Kotz, R. G. Kruger, D. Lowe, M. M. Mader, B. Marsden, A. Mueller-Fahrnow, S. Müller,
296 R. C. O’Hagan, J. P. Overington, D. R. Owen, S. H. Rosenberg, B. Roth, R. Ross, M. Schapira,
297 S. L. Schreiber, B. Shoichet, M. Sundström, G. Superti-Furga, J. Taunton, L. Toledo-Sherman,
298 C. Walpole, M. A. Walters, T. M. Willson, P. Workman, R. N. Young, and W. J. Zuercher, “The
299 promise and peril of chemical probes,” *Nat. Chem. Biol.*, vol. 11, no. 8, pp. 536–541, Aug. 2015.
- 300 [10] M. H. Rohban, S. Singh, X. Wu, J. B. Berthet, M.-A. Bray, Y. Shrestha, X. Varelas, J. S. Boehm,
301 and A. E. Carpenter, “Systematic morphological profiling of human gene and allele function via
302 cell painting,” *Elife*, vol. 6, Mar. 2017.
- 303 [11] M.-A. Bray, S. M. Gustafsdottir, M. H. Rohban, S. Singh, V. Ljosa, K. L. Sokolnicki, J. A.
304 Bittker, N. E. Bodycombe, V. Dancík, T. P. Hasaka, C. S. Hon, M. M. Kemp, K. Li, D. Walpita,
305 M. J. Wawer, T. R. Golub, S. L. Schreiber, P. A. Clemons, A. F. Shamji, and A. E. Carpenter,
306 “A dataset of images and morphological profiles of 30 000 small-molecule treatments using the
307 cell painting assay,” *Gigascience*, vol. 6, no. 12, pp. 1–5, Dec. 2017.

308 [12] C. McQuin, A. Goodman, V. Chernyshev, L. Kametsky, B. A. Cimini, K. W. Karhohs, M. Doan,
309 L. Ding, S. M. Rafelski, D. Thirstrup, W. Wiegraebe, S. Singh, T. Becker, J. C. Caicedo, and
310 A. E. Carpenter, “CellProfiler 3.0: Next-generation image processing for biology,” *PLoS Biol.*,
311 vol. 16, no. 7, p. e2005970, Jul. 2018.

312 Checklist

- 313 1. For all authors...
- 314 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
315 contributions and scope? [Yes]
- 316 (b) Did you describe the limitations of your work? [Yes]
- 317 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 318 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
319 them? [Yes]
- 320 2. If you ran experiments...
- 321 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
322 mental results (either in the supplemental material or as a URL)? [Yes]
- 323 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
324 were chosen)? [N/A]
- 325 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
326 ments multiple times)? [No]
- 327 (d) Did you include the total amount of compute and the type of resources used (e.g., type
328 of GPUs, internal cluster, or cloud provider)? [Yes]

329 **A Appendix**

330 The landing page of the GitHub repository for this dataset has all the relevant additional information:
331 <https://broad.io/neurips-cpjump1>.

332 We have released the data with a CC0 licence and the code with a BSD 3-Clause license.

333 We have chosen GitHub as the hosting platform, and use GitLFS to store large files.