1     Full title

2     **Unsupervised phenotypic analysis of cellular images with multi-scale convolutional neural**

3     **networks**

4     Short title

5     **Unsupervised deep learning for cellular image analysis**

6     Authors and affiliations

7     William J. Godinez[1,2,*], Imtiaz Hossain[1], Xian Zhang[1,*]

8     [1]Novartis Institutes for BioMedical Research, Basel, Switzerland.

9     [2]Current address: Novartis Institutes for BioMedical Research, Emeryville, CA, USA

10    *To whom correspondence should be addressed.

11    **Abstract**

12    Large-scale cellular imaging and phenotyping is a widely adopted strategy for understanding

13    biological systems and chemical perturbations. Quantitative analysis of cellular images for

14    identifying phenotypic changes is a key challenge within this strategy, and has recently seen

15    promising progress with approaches based on deep neural networks. However, studies so far

16    require either pre-segmented images as input or manual phenotype annotations for training, or

17    both. To address these limitations, we have developed an unsupervised approach that exploits the

18    inherent groupings within cellular imaging datasets to define surrogate classes that are used to

19    train a multi-scale convolutional neural network. The trained network takes as input full-

20    resolution microscopy images, and, without the need for segmentation, yields as output feature

21    vectors that support phenotypic profiling. Benchmarked on two diverse benchmark datasets, the

22    proposed approach yields accurate phenotypic predictions as well as compound potency

23    estimates comparable to the state-of-the-art. More importantly, we show that the approach

24    identifies novel cellular phenotypes not included in the manual annotation nor detected by

25    previous studies.

26    **Author summary**

27    Cellular microscopy images provide detailed information about how cells respond to genetic or

28    chemical treatments, and have been widely and successfully used in basic research and drug

29    discovery. The recent breakthrough of deep learning methods for natural imaging recognition

30    tasks has triggered the development and application of deep learning methods to cellular images

31    to understand how cells change upon perturbation. Although successful, deep learning studies so

32    far either can only take images of individual cells as input or require human experts to label a

33    large amount of images. In this paper, we present an unsupervised deep learning approach that,

34    without any human annotation, analyzes directly full-resolution microscopy images displaying

35    typically hundreds of cells. We apply the approach to two benchmark datasets, and show that the

36    approach identifies novel visual phenotypes not detected by previous studies.

37    **Introduction**

38    Image-based high-throughput cellular assays allow meticulous monitoring of chemical or genetic

39    perturbations of cellular systems at large scale(1–4). Quantitative analysis of the collections of

40    image data generated by these assays is pivotal for an objective assessment of the phenotypic

41    diversity observed within the data. Conventional workflows developed for image analysis

42    involve a series of disjoint data-processing tasks, such as detection of cellular objects, numerical

43    characterization of these objects via feature engineering, as well as classification of cellular

2

44    objects based on their features into different phenotypes(5,6). Many of these steps have been

45    addressed with the *deep learning* methodology(7,8), which has previously yielded state-of-the-

46    art results for such computer vision tasks(9–13). Approaches(14–16) based on deep learning for

47    analyzing high-content cellular images follow primarily a *supervised learning* paradigm,

48    whereby images annotated with phenotypic labels are used to train a deep neural network model

49    that maps images to one of the labels. The predictions of supervised approaches are therefore

50    constrained to the set of phenotypes defined during training, and therefore do not naturally

51    support the identification of additional phenotypes. The acquisition of these phenotypic labels

52    through manual annotation of the image data is also time-consuming (e.g., requiring

53    crowdsourcing efforts(17)), and error-prone(18). The applicability of supervised approaches is

54    thus contingent upon the availability and quality of the manual annotation.

55        Strategies to escape the limitations imposed by the a priori definition and acquisition of

56    phenotypic labels include *transfer learning* as well as *unsupervised learning*. In the former, a

57    neural network classification model trained in a supervised manner on a non-cellular image

58    dataset is applied to a cellular image dataset(19). Since the categories defined in the *source* non-

59    cellular dataset do not match those of the *target* cellular dataset, the aim of this strategy is to map

60    cellular images to a continuous coordinate system, i.e., a *feature space*, by treating the activation

61    of the hidden layers of the pre-trained deep model as a feature vector. While this strategy has

62    been shown to work well for extracting biologically informative features(19), there are no

63    guarantees that models trained on non-cellular data generalize well to arbitrary cellular image

64    data. Technical issues such as different channel encodings (e.g., RGB channels in non-cellular

65    images compared with an arbitrary number of fluorescence channels in cellular images) and

66    noise models (e.g., additive Gaussian noise models in non-cellular images(20) compared with

67    Mixed-Poisson-Gaussian statistics(21) in fluorescence images) also hinder the applicability of

68    approaches based on transfer learning.

69    Approaches following an unsupervised learning paradigm are, in contrast, typically

70    optimized on the specific cellular dataset of interest. The aim of unsupervised learning is to map

71    images to a feature space where biologically relevant patterns within the dataset might emerge.

72    While in the supervised learning paradigm deep models are designed to predict an *extrinsic*

73    characteristic or attribute of the data, e.g., the phenotypic label manually assigned to the images,

74    in the unsupervised learning paradigm deep models are designed to predict an *intrinsic*

75    characteristic of the data. The most inherent property of each image is the pixel data itself. The

76    training process of both *autoencoder networks*(22) as well as *generative adversarial networks*

77    (23)(GANs) therefore typically involves the optimization of an image synthesis function aiming

78    to reconstruct an image's raw pixel data from a low dimensional representation of the input

79    image. This type of approaches has been able to map single-cell images with small dimensions

80    (e.g., $40 \times 40$ pixels) to a low-dimensional space (e.g., 64-D) where aberrant morphologies

81    during cell division as induced by siRNAs may be identified(24). Because of the high spatial and

82    phenotypical variability found in multi-cellular images with larger dimensions (e.g., $1280 \times 1024$

83    pixels over three fluorescent channels, i.e., a 3932160-D space), the compression and

84    reconstruction of high-content images via a neural network is currently a computationally

85    prohibitive task.

86    Here we present an unsupervised approach based on the *exemplar convolutional neural*

87    *network*(25) (Exemplar-CNN) training methodology that optimizes a network model to

88    discriminate among *surrogate classes*, which, in our case, are automatically defined through the

89    intrinsic groupings of images (e.g., images belonging to the same treatment) typically found in

4

90  high-content imaging studies. The proposed approach uses exclusively dataset-specific multi-

91  cellular images, and requires no phenotypic annotations or optimization of computationally-

92  expensive image reconstruction functions. Using this unsupervised strategy, we train our multi-

93  scale convolutional neural network architecture (M-CNN(16)) on the multi-cellular images of the

94  KiMorph(26) and BBBC021(27,28) datasets, which involve genetic and chemical perturbations

95  of cellular systems at scale, respectively. Our approach, without any user-provided phenotypic

96  labels and without any object segmentation, is able to map images to a feature space that enables

97  prediction of phenotypes that match well with held-out labels. In addition, we show that the

98  approach identifies novel phenotypes in the benchmark datasets not detected by previous studies.

99  **Results**

100  **Training and validating a deep neural network with the Exemplar-CNN methodology**

101  To train a neural network model without any user-provided phenotypic annotation, we used the

102  Exemplar-CNN(25) training methodology (see **Methods** for details). When applied to high-

103  content cellular images, the trained network model maps in one step an image to a continuous

104  feature space representing the phenotypic homogeneity and variability observed in the data (see

105  **Fig. 1** for a schematic overview of the approach). We validated this unsupervised strategy on the

106  Kimorph and BBBC021 datasets, which include images of cells subjected to siRNA and

107  compound treatments, respectively. On each dataset, we trained a multi-scale convolutional

108  neural network (M-CNN(16)) model with the Exemplar-CNN training methodology. To define

109  the surrogate classes required by this methodology, we hypothesized that images belonging to

110  the same well (KiMorph) or compound treatment (BBBC021) defined a single surrogate class.

111  No phenotypic categories and annotations were therefore needed for training. Note that different

5

112   surrogate classes might belong to the same phenotypic category, but this information is not

113   known to the network. We trained the neural network exclusively with the pixel data of

114   annotated images; the annotations were removed during training, and only used subsequently to

115   validate the performance of the approach.

116   **Fig. 1. Schematic overview of the Exemplar-CNN approach**. Images are grouped into

117   surrogate classes based on intrinsic information (such as treatment information) instead of

118   external annotation. Taking the full-resolution images of the surrogate classes as input, an M-

119   CNN model is trained with the objective of separating the different surrogate classes. Once

120   trained, as input images are fed to the network, the neural activation values are extracted as

121   feature vectors, thus mapping the input images to a low-dimensional feature space. Finally,

122   distance calculation and clustering analysis enable the identification of novel phenotypes.

123

124          Once the dataset-specific network models were trained, we validated the performance of

125   the approach in two steps. First, we built a nearest-neighbor classifier based on the feature

126   vectors computed by the trained M-CNN model to predict the phenotype of annotated images.

127   We used the classification accuracy of the nearest neighbor classifier to evaluate whether the

128   feature vectors computed by the network encoded relevant phenotypic information. Second, we

129   applied the trained M-CNN model to the dataset's entire image collection, including images with

130   no annotation and not used during training, thereby obtaining one feature vector for each image

131   in the entire collection. We performed hierarchical clustering analysis on the feature vectors, and,

132   through visual inspection of images in selected clusters, identified novel phenotypes. The next

133   two sections describe in detail the results for each dataset.

**KiMorph analysis and results**

The Kimorph dataset comes from an RNAi screen where HeLa cells were reverse transfected with siRNAs targeting ca. 800 kinases in duplicate. siRNA-mediated perturbations of the UBC, CLSPN and TRAPPC3 genes were used as positive controls while the Renilla luciferase (*Rluc*) siRNA treatment was used as a neutral control. After transfection, cells were fixed and labeled for DNA, F-actin, and B-tubulin, and imaged through an automated microscope (experimental details can be found in the original publication(26)). Each of the four control siRNA treatments (UBC, CLSPN, TRAPPC3, and Rluc) was spotted across 12 wells in duplicate. We declared all fields-of-view (FOVs) coming from each well to define a single surrogate class, which amounted to 48 surrogate classes (i.e., one class per replicate well). We then trained an M-CNN model to maximize separation among these classes (see **Methods** for training details). Note that some surrogate classes belong to the same control siRNA treatment but this information is not known to the network. Previous studies(26,29) have shown that each of the four control siRNA treatments induces a consistent phenotype across wells. If the network learned to identify invariant and discriminative features reflecting the control phenotypes, we hypothesized that these would remain relatively similar within surrogate classes belonging to the same control while varying more strongly across surrogate classes belonging to different controls.

After training, we used the M-CNN model and PCA to extract a 94-dimensional feature vector for each original FOV image. We aggregated the feature vectors at the well level, and calculated cosine distance values between each pair of wells (see **Methods** for details). The resulting distance matrix, with rows ordered by control groups, is shown in **Fig. 2a**. We observe square blocks (submatrices) along the diagonal of the matrix, which reflect the low distance values (i.e., high similarity) within wells from the same control siRNA treatment as entailed by

7

157   the feature vectors computed by the network. To quantitatively verify the invariant and

158   discriminative properties of the feature vectors within and across the four different phenotypes,

159   we tested whether we could predict the phenotype of each well based on the phenotype of the

160   well's nearest neighbor in feature space. Over 50 repetitions of a random hold-out cross-

161   validation strategy, this nearest-neighbor classification approach identifying the four control

162   phenotypes yielded 100% classification accuracy (**Supplementary Table 1**). The results show

163   that the feature vectors computed by the M-CNN model, which was trained without any

164   phenotypic annotation, remain relatively similar within the same phenotype yet vary across

165   different ones, thus encoding phenotypic information that enables the identification of distinct

166   phenotypes.

167   **Fig. 2. Exemplar-CNN training and clustering analysis results for the KiMorph dataset. (a)**

168   Cosine distance matrix between pairs of wells of control siRNA treatments. Rows are ordered by

169   control groups. Blue indicates a small distance value while yellow corresponds to a large

170   distance value. (**b**) Sample images from clusters including the control siRNA treatments (left

171   column) as well as from clusters including phenotypes distinct from the control treatments (right

172   column).

173

174        We next tested whether the M-CNN model trained with four control siRNA treatments

175   could generalize to images and phenotypes not used during training. We therefore fed each

176   image of the entire KiMorph dataset through the trained M-CNN model and PCA, and obtained a

177   94-D feature vector per image. Feature vectors of images belonging to the same siRNA treatment

178   were aggregated onto a single vector (see **Methods** for details). We calculated pairwise cosine

179   distances among all vectors corresponding to all 781 siRNA perturbations (see **Supplementary**

180   **Table 2**), performed hierarchical clustering, and grouped all siRNA perturbations onto 27

181   clusters (see **Supplementary Table 3**). On each cluster, we carried out a Gene Ontology (GO)

182   enrichment analysis for biological processes through the topGO R package(30), with all kinases

183   in the library taken as the background set. After Benjamini–Hochberg correction for multiple

184   testing, we found that 22 clusters out of the 24 clusters that included more than one gene were

185   enriched with two or more GO terms. This indicates that the feature vectors computed by the

186   network supported the identification of shared biological functions of groups of genes (see

187   **Supplementary Table 4**). To validate the results, we first inspected clusters 17, 16, and 18,

188   which included the three positive siRNA controls (viz. UBC, CLSPN, and TRAPPC3),

189   respectively (**Fig. 2b**). The UBC cluster, which the enrichment analysis associates with DNA

190   damage and integrity checkpoints, includes essential genes such as COPB2 and PLK1 that, when

191   knocked down, cause a lethal phenotype akin to that of the UBC treatment. Likewise, the

192   CLSPN cluster comprises genes such as CDC7 and CDK3, which are associated with cell cycle

193   control, and whose knockdowns induce an enlarged cell phenotype resembling that of the

194   CLSPN treatment. In the TRAPPC3 cluster, which is enriched with the 'integrin-mediated

195   signaling pathway' GO term, we typically observe an elongated cell phenotype that is likewise

196   triggered by the CARD10 and SYK knockdowns. Overall, the clustering results and the

197   recapitulation of known biological functions of groups of genes suggest the feature vectors

198   learned by the network capture phenotypic information.

199       The main advantage of an unsupervised approach is its ability to discover novel phenotypes.

200   To verify this premise, we identified clusters that were relatively distant (in terms of the cosine

201   distance values) from the four siRNA controls (see **Methods** for details). One of these distant

9

202    clusters (viz. cluster 3) only includes images from the LAK perturbation. Visual inspection of the

203    images reveals an experimental artifact in images from replicate 1 (**Fig. 2b** top right; compare to

204    images from replicate 2). The artifact is not visible in images of any other treatment, and so the

205    approach correctly grouped images with this artifact onto a separate cluster. Images from cluster

206    27, which includes genes such as ACVR1 and GALK2, display a phenotype of enlarged nuclei

207    and cells with a strong actin signal (**Fig. 2b** middle right) that do not resemble any of the control

208    phenotypes. Likewise, in cluster 20, which includes genes such as RAC1 and PDGFRB, and

209    shows enrichment for the 'positive regulation of Rho protein signal transduction' GO term, we

210    observe a reduced cell count and cell size in the images. These results suggest that the proposed

211    unsupervised strategy supports the identification of novel imaging phenotypes.

212    **BBBC021 analysis and results**

213    In the BBBC021 dataset, MCF-7 breast cancer cells were treated with 113 compounds at eight

214    concentrations in triplicate, before being fixed and labeled for DNA, F-actin, and B-tubulin.

215    Images were captured from each channel with four fields per well(28). A subset of 103

216    compound-concentration pairs (hereafter defined as treatments) covering 38 compounds was

217    previously inspected and annotated for one of twelve mechanisms-of-action (MoAs)(31). We

218    declared all images coming from each treatment to belong to the same surrogate class, which

219    resulted in 103 surrogate classes. The M-CNN model was trained to discriminate among all 103

220    surrogate classes (see **Methods** for training details). Note that the network is not aware that

221    certain surrogate classes (i.e., treatments) belong to the same compound or the same MoA. If the

222    network learned MoA-relevant features, we posited that these would remain relatively invariant

223    within each MoA yet vary across different MoAs.

10

224    Once trained, we used the M-CNN model and PCA to extract an 8-D vector for each

225    input image used during training. We aggregated the feature vectors of images belonging to the

226    same treatment onto a single feature vector and computed cosine distances between all pairs of

227    treatments (**Fig. 3a)**. Each row of the matrix corresponds to one treatment. Treatments are

228    ordered by MoA, and then by compound and concentration. Overall we observe sub-matrices

229    (squares) along the diagonal indicating that the network learned features that remain relatively

230    invariant within each MoA. To quantitatively verify the homogeneity and variation of the learned

231    features within and across phenotypes, respectively, we tested whether the MoA of a treatment

232    could be identified based on the MoA label of the treatment's nearest neighbor in feature space.

233    Here we adopted the same leave-one-compound-out cross-validation strategy used in previous

234    benchmarking studies(31) that prevents matching treatments from the same compound (see

235    **Methods** for details). With this nearest-neighbor classification strategy, we achieve a median

236    accuracy over all classes of 88%. The confusion matrix is shown in **Supplementary Table 5**.

237    For certain MoAs (e.g., aurora kinase inhibitors, cholesterol-lowering, protein degradation, and

238    protein synthesis), the approach achieved 100% classification accuracy. For other MoAs, the

239    accuracy ranged from 75% to 89%. The performance of the approach is comparable to our

240    previous supervised approach(16), which is explicitly optimized to distinguish these twelve

241    MoAs, as well as to other non-supervised approaches(19,31). Overall the results show that, while

242    the MoA categories are unknown to the network, it manages to learn features which remained

243    relatively invariant within each MoA yet varied across MoAs.

244    **Fig. 3. Exemplar-CNN training results for the BBBC021 annotated subset. (a)** Cosine

245    distance matrix between pairs of treatments. Labels on the left are the MoA annotations. Rows

246    are ordered by MoA annotation, compound, and concentration. Blue indicates a small distance

11

247 value while yellow corresponds to a large distance value. (**b**) Zoomed-in view of the red box in

248 (a), KI for kinase inhibitors and MD for microtubule destabilizers. Labels on the right are

249 compounds and concentrations in μM. (**c**) Sample images corresponding to the treatments in (b).

250 Colchicine, which is annotated as a microtubule destabilizer (MD), visually looks more similar

251 to treatments annotated as kinase inhibitors (KI).

252

253       While the approach is able to recapitulate known information about the annotated data,

254 we tested further its ability to reveal phenotypic information beyond the annotation. To this aim,

255 we conducted a closer examination of the distance matrix. While treatments annotated with the

256 same MoA are mapped by the network to nearby positions in feature space, and are therefore

257 distinguishable via a nearest-neighbor classifier, the sub-matrices along the diagonal of the

258 distance matrix in **Fig. 3a** reveal a certain heterogeneity within individual MoAs. For example,

259 for the aurora kinase inhibitors (Aur) MoA, the corresponding sub-matrix reveals three groups

260 corresponding to the three compounds annotated with this MoA (viz. AZ-A, AZ258 and AZ841),

261 which suggests that the compounds caused slightly different sub-phenotypes. A similar

262 observation can be made for the actin disruptors (Act), protein degradation (PD), and protein

263 synthesis (PS) MoAs. The microtubule destabilizers (MD) MoA is comprised by four

264 compounds (14 treatments, sub-matrix highlighted in red in **Fig. 3a**. and zoomed in view in **Fig.**

265 **3b**). Three of the four compounds, Demecolcine, Nocodazole, and Vincristine, are relatively

266 similar to each other as well as distant to the kinase inhibitor (KI) group, although Nocodazole

267 shows a sub-phenotype different from Demecolcine and Vincristine. The fourth compound,

268 Colchicine at 0.03μM, which is annotated as MD, instead seems to be closer to the kinase

269 inhibitors (KI) treatments than to the other MD treatments, and is accordingly predicted as KI by

12

270    the nearest-neighbor classification scheme. Visual examination of the corresponding images also

271    confirms Colcichine's similarity to the KI treatments (**Fig. 3c**). Although only one concentration

272    (0.03μM) of Colcichine is included in the annotation subset, there are seven concentrations

273    (0.001 – 3.0 μM) in the entire BBBC021 dataset. Only at 3.0μM, Colchicine causes phenotypes

274    similar to other microtubule destabilizers (**Supplementary Fig. 2**). The proposed unsupervised

275    approach is thus able to detect phenotype information in the data beyond the manual annotation,

276    which is not feasible with a supervised method.

277        Next, we set out to verify the applicability of the approach to the entire BBBC021

278    dataset. Here we first trained an M-CNN model using all images from the 103 annotated

279    treatments (amounting to 103 surrogate classes) plus images from the neutral control (DMSO),

280    which were grouped into an additional surrogate class. Once trained, we applied the model to all

281    13200 images included in the dataset. Using PCA, we obtained a 77-D feature vector per image.

282    Vectors of images belonging to the same treatment replicate (well) were aggregated onto a single

283    vector. We then determined the *similarity* of each replicate vector to each of the 12 MoAs and

284    DMSO based on the cosine distances of each vector to all replicate vectors belonging to the 103

285    MoA-annotated treatments as well as to all DMSO wells (see **Methods** for details). The 13

286    similarity values of each treatment replicate are shown in **Supplementary Table 6**.

287        In our previous supervised analysis of the BBBC021 dataset, four compounds were

288    selected as representative concentration-response curves (16). In the current study, we selected

289    the same four compounds and plotted the similarity values to each MoA and DMSO as a

290    function of the concentration (**Supplementary Fig. 1**). In the previous supervised approach, the

291    y-axis was the classification probability and for each compound there was only one or two

292    dominant MoAs across concentrations. In the current unsupervised approach, the y-axis is the

13

293    similarity to each MoA which ranges from 0 to 2, and the gaps between curves are much less

294    pronounced. To compare the overall trend, we simplified the plot by only showing MoAs shown

295    as dominant in the previous supervised approach (**Fig. 4a**). Data points highlighted by dashed

296    circles correspond to concentrations annotated with the curve's MoA (and therefore achieving

297    maximum similarity). For Floxuridine, consistent with the supervised approach, the DNA

298    replication (DR) MoA is the top MoA prediction for all concentrations. For Nocodazole, the

299    curve shows a similar trend to the supervised approach, with DMSO as the top MoA in low

300    concentrations (0.001-0.01µM) and microtubule destabilizers (MD) as the top MoA in high

301    concentrations (0.1-4.0µM). For Alsterpaullone, in the current unsupervised analysis, kinase

302    inhibitor (KI) and DMSO are on the same level until the higher concentrations. DNA-damage

303    (DD) increases at the last concentration but does not pass the level of KI. In the previous

304    supervised analysis the differences among the MoAs were much more obvious although with

305    larger error bars. Finally, for Hydroxyurea, for which none of the concentrations was included in

306    the training data, the trend of the curve is consistent with the supervised approach, where DMSO

307    decreases over concentration while DNA-damage (DD) increases and takes over at the two

308    highest concentrations.

309    **Fig. 4. Example concentration-response curves and clustering analysis for the BBBC021**

310    **dataset.** (**a**) Similarity-vs-concentration plots for four compounds. The similarity (y-axis) to

311    selected MoAs and DMSO over concentration (x-axis) computed using the features vectors

312    yielded by the proposed approach is shown. The dots and error bars represent the median and

313    MAD over the experimental replicates ($n = 2$ for Alsterpaullone and $n = 3$ for the other three

314    compounds). Data points marked by dashed circles are annotated with the curve's MoA and

14

315    therefore achieve maximum similarity. (**b**) Sample images of clusters including distinct

316    phenotypes not related to the annotated phenotypes.

317    Finally, we tested the ability of the approach to detect novel phenotypes. To this end, we

318    further aggregated the replicate vectors belonging to the same treatment onto a single vector.

319    Likewise, DMSO vectors stemming from the same plate were aggregated onto a single vector.

320    We calculated pairwise cosine distances among all treatments, including DMSO, and applied a

321    hierarchical clustering procedure that yielded 79 clusters (see **Methods** as well as

322    **Supplementary Table 7** for the complete distance matrix). We inspected visually clusters that

323    included exclusively compound-concentration treatments without any MoA annotation (see **Fig.**

324    **4b**). For example, in cluster 37, we found images from Mitoxantrone at 10μM and Staurosporine

325    at 0.1μM and 0.3μM that induced a strong toxic phenotype. In cluster 20, which included images

326    from AZ-841 at 30μM only, we found images that displayed an unusual purple phenotype that

327    could hint at a tubulin toxin/disruptor MoA for this treatment.  Finally, in cluster 41, we found

328    images from Staurosporine at 0.0003μM, Bryostatin at 3.0μM, as well as Valproic Acid at

329    150μM where groups of elongated cells with thin protrusions forming a networked pattern were

330    visible. The results underscore the ability of the proposed unsupervised approach to identify

331    novel phenotypes not previously known and not included during training.

332    **Discussion**

333    Deep learning has been successfully pioneered in the field of image-based high-throughput

334    screening(14–16,19,24). The majority of approaches based on deep neural networks adopt a

335    supervised learning paradigm that requires manual definition and acquisition of phenotypic

336    labels. As such, supervised approaches do not support naturally the discovery of new

15

337   phenotypes. In this work, instead of relying on predefined phenotypic labels, we developed an

338   unsupervised learning approach that exploits the inherent variation across treatments typically

339   found in imaging-based studies to learn phenotypically relevant features that enable the

340   discovery of novel phenotypes.

341     The proposed approach obviates the need for manually specified phenotypic categories by

342   defining automatically surrogate categories through the inherent grouping of images (e.g.,

343   images belonging to the same well) found in the experimental design of high-content studies.

344   The fact that multiple surrogate categories may belong to a (known) phenotypic class remains

345   explicitly held-out to the neural network model throughout. Our results on two benchmark

346   datasets demonstrate that the feature vectors extracted from the images through the trained

347   models support the recognition of known phenotypes included within the surrogate categories.

348   By testing the models on images outside of the surrogate categories, we also showed that the

349   models generalize to phenotypes beyond those used during training. With a straightforward

350   clustering analysis of the feature vectors, we managed to pinpoint novel phenotypes, which is

351   one of the main goals of image-based high-content screening studies, where genetic or chemical

352   perturbations may potentially induce a range of unexpected phenotypes.

353     Certainly, one could identify novel phenotypes with conventional image analysis

354   approaches, which typically require segmentation and manual feature engineering(26,28,32,33).

355   It is however encouraging to see that the proposed unsupervised approach, which requires no

356   segmentation, no manual feature engineering, and no phenotypic categories and annotations, also

357   supports the identification of novel phenotypes in a more automated fashion. The proposed

358   approach does not provide single-cell readouts, and therefore does not replace single-cell

359   analyses(34,35).

16

360    With the proposed unsupervised learning strategy, the inferred network models depend on

361    the phenotypic data included within the surrogate classes. In our study, we restricted the

362    surrogate classes to images that had a phenotypic annotation. This strategy facilitated the

363    validation of the approach, as it allowed testing whether the approach supported the recovery of

364    known phenotypic classes. Additional work is however needed to decide which images and

365    phenotypes should be included within the surrogate classes. One possibility would be to adapt an

366    active learning approach, where surrogate classes would be iteratively added based on a certain

367    performance criterion.

368    **Methods**

369    **Exemplar Convolutional Neural Networks**

370    We use the Exemplar-CNN optimization strategy(25) to train a convolutional neural network

371    without relying on any phenotypic label annotation. In contrast to a typical supervised learning

372    approach, where the neural network is trained to discriminate among a set of predefined

373    phenotypic classes, the proposed approach is trained to discriminate among a set of *surrogate*

374    *classes*. The main idea underlying the Exemplar-CNN methodology is to learn image features

375    that are both *invariant* within each surrogate class as well as *discriminative* across surrogate

376    classes. In the original strategy, each *exemplar* (i.e., a region-of-interest within an image) and

377    transformed versions thereof (obtained through extreme data augmentation schemes) defined a

378    single surrogate class. This strategy was shown to work well with a large number of surrogate

379    classes (e.g., up to 4000). However, when the number of (exemplar) images is very large and the

380    images look very similar, discrimination among the surrogate classes becomes more challenging.

381    A prior grouping (e.g., through clustering) of similar images was suggested as an approach to

382    reduce the number of classes as well as to group very similar images into a single surrogate class.

17

383    In our case, instead of taking each image and its variations as a single surrogate class, we take

384    advantage of the intrinsic grouping of images provided by the experimental design of each study

385    to define the surrogate classes. For example, for each well, multiple fields-of-view (FOVs) are

386    typically acquired. We may therefore define all FOVs from a single well to define a single

387    surrogate class. Similarly, each treatment combination (e.g., a compound at a specific

388    concentration) is typically replicated. Images from these replicates may be therefore declared as

389    a single surrogate class. The definition of surrogate classes depends on the experimental details

390    in each study. After defining $N_{\mathrm{s}}$ surrogate classes in such a way, we associate a numerical label

391    $y_{\mathrm{surrogate}}$ with each surrogate class and its images.

392        We use a multi-scale convolutional neural network (M-CNN) architecture to solve the

393    task of surrogate class discrimination. The last two layers of our M-CNN architecture include a

394    fully connected layer with 128 hidden units, as well as a soft-max output layer, which yields a

395    vector $\boldsymbol{\rho}$ with elements $\rho_k$ that encode a probability score for each of the $N_{\mathrm{s}}$ surrogate classes to

396    be identified (all architectural details are provided in **Supplementary Table 8**). Using $N_{\mathrm{t}}$ images

397    associated with surrogate classes and their numerical labels, we optimize the parameters of the

398    M-CNN by minimizing the following error function:

399    $$\frac{1}{N_{\mathrm{t}}}\sum_{i=1}^{N_{\mathrm{t}}} f\left(\boldsymbol{\rho}^{(i)}, y_{\mathrm{surrogate}}^{(i)}\right) + \lambda\|\mathbf{w}\|_2$$

400    where $f(\cdot,\cdot)$ is the cross-entropy error function evaluating the agreement between the network's

401    soft-max output $\boldsymbol{\rho}^{(i)}$ and the surrogate (true) label $y_{\mathrm{surrogate}}^{(i)}$ for the *i-th* training example, $\|\cdot\|_2$ is

402    the L2 norm, $\mathbf{w}$ is a vector including all weights of the network, and $\lambda$ is a coefficient that

403    regulates the influence of the magnitude of the weight vector on the error function. We use the

18

404    stochastic gradient descent (SGD) algorithm via backpropagation and drop-out to approximate a

405    solution.

**Learning details**

407    Generally, we used the same strategy and parameter values that we used previously to train the

408    M-CNN architecture in a supervised way(16). In this study, we however increased the number of

409    training epochs to 27. The step size over which the learning rate is held constant was also

410    increased to 9 epochs. One epoch is equal to the number of iterations needed to evaluate all

411    images in the training dataset. We additionally used the dropout technique(36) on the

412    penultimate layer of the M-CNN architecture to encourage a better exploration of the available

413    activation space.

**Feature extraction, projection, and aggregation**

415    Once trained, the application of the M-CNN model to any input image yields a 128-dimensional

416    activation vector $\mathbf{z}$ with elements $z_i$ corresponding to the activation values of each hidden unit

417    within the fully connected layer (second-to-last layer) that are recorded as the input image is

418    passed through the network. We subsequently project all activation vectors onto an orthogonal

419    basis computed via principal component analysis (PCA) that takes exclusively into consideration

420    the activation vectors of (non-augmented) images used during training. Principal components

421    explaining 99% of the variance define the new feature sub-space onto which all activation

422    vectors are typically projected.

423        Feature vectors belonging to the fields-of-view (FOVs) of a well are aggregated by taking

424    the element-wise median of the vectors. The resulting vector is taken as the feature vector

425    representing the corresponding well. Likewise, to construct the feature vector for a given

19

426    treatment, feature vectors of the treatment's replicate wells are summarized by taking the

427    element-wise median of the vectors.

428

429    **Distance and similarity calculations**

430    To compare treatments, we use the cosine distance between two feature vectors. The cosine

431    distance is defined as one minus the cosine of the angle between the vectors. The values thus

432    range from 0 (denoting an identical direction for both vectors) to 2 (denoting opposite

433    directions). To obtain a measure of *similarity* between treatments within the same numerical

434    range, we subtract each cosine distance value from two.

435

436    **Clustering**

437    We compute cosine distances among all pairs of treatments in a dataset. We use a hierarchical

438    clustering algorithm to group treatments based on these pairwise cosine distance values. The

439    resulting hierarchical tree is partitioned with a threshold value equivalent to the cosine distance

440    entailed by an angle of $\pi / 3$.

441

442    **Nearest neighbor classifier**

443    Using pairwise cosine distances, we build a nearest neighbor classifier to investigate whether the

444    feature vectors obtained via the unsupervised model encoded information that supported the

445    retrieval of known phenotypic categories that had been manually assigned to a subset of

446    treatments. Evaluation of the classifier's performance requires splitting the feature vectors onto a

447    training set and a test set. Given a feature vector from the test set, we determine its closest

448    feature vector (i.e., its nearest-neighbor) within the training set, and assign the nearest neighbor's

449    phenotypic or MoA category to the test feature vector.

450         In the KiMorph dataset, we use a random hold-out cross-validation strategy where we

451    randomly group all feature vectors into a training set and test set. The proportion of treatments

452    assigned to the training set is 90%. Using the nearest-neighbor classifier, we predict the

453    phenotype of the feature vectors in the test set, and evaluate the classification performance. We

454    repeat the partitioning and evaluation process 50 times. The confusion matrix aggregating the

455    results over the 50 repeats is shown in **Supplementary Table 1**.

456

457    In the BBBC021 dataset, we use a leave-one-compound-out validation strategy, where the

458    training dataset excludes feature vectors of treatments (i.e., compound-concentration pairs)

459    sharing the same compound as the test feature vector. We use all 103 treatments as test feature

460    vectors once, obtain a nearest-neighbor prediction for the MoA, and compare the prediction with

461    treatments' known MoA. The resulting confusion matrix is shown in **Supplementary Table 4**.

462

463    **Image pre-processing**

464    All image intensities are subjected to an Anscombe transform. Histogram normalization of each

465    image is carried out on per-plate basis as described previously(16). All image intensities are

466    mapped to an 8-bit range.

467

468    **Image datasets**

469    The KiMorph dataset is available from the Wolfgang Huber Group EBI website at

470    https://www.ebi.ac.uk/huber-srv/cellmorph/kimorph/.

471    The BBBC021 version 1 image dataset is available from the Broad Bioimage Benchmark

472    Collection at http://www.broadinstitute.org/bbbc/BBBC021/.

473

474    Detailed description of the datasets can be found on their corresponding webpages.

475

**Acknowledgements**

477    We would like to thank Florian Fuchs for fruitful discussions as well as for providing the

478    KiMorph image data via the Wolfgang Huber Group EBI website. We also acknowledge the

479    BBBC021 dataset provided by Peter D. Caie via the Broad Bioimage Benchmark Collection.

480

**Author contributions**

482    All authors conceived jointly the study. WJG designed the neural network architecture as well as

483    training scheme, and performed the analysis. IH set up the deep learning computational

484    framework. WJG and XZ wrote the manuscript.

485

489

490

491 **Competing interests**

492 None declared

493

## References

494

495    1.    Götte M, Hofmann G, Michou-Gallani AI, Glickman JF, Wishart W, Gabriel D. An

496        imaging assay to analyze primary neurons for cellular neurotoxicity. J Neurosci Methods.

497        2010;192(1):7–16.

498    2.    Liberali P, Snijder B, Pelkmans L. Single-cell and multivariate approaches in genetic

499        perturbation screens. Nat Rev Genet. 2014;16(1):18–32.

500    3.    Finkbeiner S, Frumkin M, Kassner PD. Cell-Based Screening: Extracting Meaning from

501        Complex Data. Neuron. 2015;86(1):160–74.

502    4.    Boutros M, Heigwer F, Laufer C. Microscopy-Based High-Content Screening. Cell.

503        2015;163(6):1314–25.

504    5.    Sommer C, Gerlich DW. Machine learning in cell biology - teaching computers to

505        recognize phenotypes. J Cell Sci. 2013;126(24):5529–39.

506    6.    Caicedo JC, Cooper S, Heigwer F, Warchal S, Qiu P, Molnar C, et al. Data-analysis

507        strategies for image-based cell profiling. Nat Methods. 2017;14(9):849–63.

508    7.    LeCun Y, Bengio Y, Hinton G. Deep learning. Nature. 2015;521(7553):436–44.

509    8.    Angermueller C, Pärnamaa T, Parts L, Stegle O. Deep learning for computational biology.

510        Mol Syst Biol. 2016;12(7):878.

511    9.    Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image

512        Segmentation. In: Proc Medical Image Computing and Computer-Assisted Intervention.

513        2015. p. 238–45.

514    10.    Ciresan D, Meier U, Schmidhuber J. Multi-column Deep Neural Networks for Image

515        Classification. In: Proc IEEE Conference on Computer Vision and Pattern Recognition.

516        2012. p. 3642–9.

517  11.  Krizhevsky A, Sutskever I, Hinton G. ImageNet Classification with Deep Convolutional

518      Neural Networks. In: Proc Advances in Neural Information Processing Systems 25. 2012.

519      p. 1097–105.

520  12.  Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, et al. Going Deeper with

521      Convolutions. In: Proc IEEE Conference on Computer Vision and Pattern Recognition.

522      2015. p. 1–9.

523  13.  Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM, et al. Dermatologist-level

524      classification of skin cancer with deep neural networks. Nature. 2017;542(7639):115–8.

525  14.  Dürr O, Sick B. Single-Cell Phenotype Classification Using Deep Convolutional Neural

526      Networks. J Biomol Screen. 2016;21(9):3–8.

527  15.  Kraus OZ, Grys BT, Ba J, Chong Y, Frey BJ, Boone C, et al. Automated analysis of high-

528      content microscopy data with deep learning. Mol Syst Biol. 2017;13(4):924.

529  16.  Godinez WJ, Hossain I, Lazic SE, Davies JW, Zhang X. A Multi-Scale Convolutional

530      Neural Network for Phenotyping High-Content Cellular Images. Bioinformatics.

531      2017;33(13):2010–9.

532  17.  Helmstaedter M. Cellular-resolution connectomics: challenges of dense neural circuit

533      reconstruction. Nat Methods. 2013;10(6):501–7.

534  18.  Li J, Newberg JY, Uhlén M, Lundberg E, Murphy RF. Automated Analysis and

535      Reannotation of Subcellular Locations in Confocal Images from the Human Protein Atlas.

536      PLoS One. 2012;7(11).

537  19.  Pawlowski N, Caicedo JC, Singh S, Carpenter AE, Storkey A. Automating Morphological

538      Profiling with Generic Deep Convolutional Networks. In: Proc 2016 NIPS Workshop on

539      Machine Learning in Computational Biology. 2016.

25

bioRxiv preprint doi: https://doi.org/10.1101/361410; this version posted July 3, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

540   20.   Sheikh HR, Bovik AC. Image information and visual quality. IEEE Trans Image Process.

541         2006;15(2):430–44.

542   21.   Zhang B, Fadili JM, Starck J, Variance- JOM. Multiscale Variance-stablizing Transform

543         for Mixed-Poisson-Gaussian Processes and its Applications in Bioimaging. In: Proc

544         International Conference on Image Processing. 2007. p. 233–6.

545   22.   Hinton GE, Salakhutdinov R. Reducing the Dimensionality of Data with Neural

546         Networks. Science (80- ). 2006 Jul 28;313(5786):504–7.

547   23.   Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al.

548         Generative Adversarial Nets. In: Proc Advances in Neural Information Processing

549         Systems 27. 2014. p. 2672–80.

550   24.   Sommer C, Hoefler R, Samwer M, Gerlich DW, Biocenter V. A deep learning and novelty

551         detection framework for rapid phenotyping in high-content screening. Mol Biol Cell.

552         2017;28(23):3133–470.

553   25.   Dosovitskiy A, Fischer P, Springenberg JT, Riedmiller M, Brox T. Discriminative

554         Unsupervised Feature Learning with Exemplar Convolutional Neural Networks. IEEE

555         Trans Pattern Anal Mach Intell. 2016;38(9):1734–47.

556   26.   Fuchs F, Pau G, Kranz D, Sklyar O, Budjan C, Steinbrink S, et al. Clustering phenotype

557         populations by genome-wide RNAi and multiparametric imaging. Mol Syst Biol.

558         2010;6:370.

559   27.   Ljosa V, Sokolnicki KL, Carpenter AE. Annotated high-throughput microscopy image

560         sets for validation. Nat Methods. 2012;9(7):637–637.

561   28.   Caie PD, Walls RE, Ingleston-Orme A, Daya S, Houslay T, Eagle R, et al. High-content

562         phenotypic profiling of drug response signatures across distinct cancer cells. Mol Cancer

563         Ther. 2010;9(6):1913–26.

564    29.  Zhang X, Boutros M. A novel phenotypic dissimilarity method for image-based high-

565         throughput screens. BMC Bioinformatics. 2013;14(1):336.

566    30.  Alexa A, Rahnenführer J, Lengauer T. Improved scoring of functional groups from gene

567         expression data by decorrelating GO graph structure. Bioinformatics. 2006;22(13):1600–

568         7.

569    31.  Ljosa V, Caie PD, ter Horst R, Sokolnicki KL, Jenkins EL, Daya S, et al. Comparison of

570         Methods for Image-Based Profiling of Cellular Morphological Responses to Small-

571         Molecule Treatment. J Biomol Screen. 2013;18(10):1321–9.

572    32.  Bakal C, Church G, Perrimon N. Quantitative Morphological Signatures Define Local

573         Signaling Networks Regulating Cell Morphology. Science (80- ). 2007;316(5832):1753–

574         6.

575    33.  Loo L-H, Wu LF, Altschuler SJ. Image-based multivariate profiling of drug responses

576         from single cells. Nat Methods. 2007;4(5):445–53.

577    34.  Carpenter AE, Jones TR, Lamprecht MR, Clarke C, Kang IH, Friman O, et al.

578         CellProfiler: image analysis software for identifying and quantifying cell phenotypes.

579         Genome Biol. 2006;7(10):R100.

580    35.  Matula P, Kumar A, Wörz I, Erfle H, Bartenschlager R, Eils R, et al. Single-cell-based

581         image analysis of High-throughput cell array screens for quantification of viral infection.

582         Cytom Part A. 2009;75(4):309–18.

583    36.  Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: A Simple

584         Way to Prevent Neural Networks from Overfitting. J Mach Learn Res. 2014;15:1929–58.

585

586

587

588 **Supplementary Material**

589 **Supplementary Table 1:** Confusion matrix on control siRNA treatments of the KiMorph

590 dataset.

591 **Supplementary Table 2:** Distance matrix of entire siRNA collection of the KiMorph dataset

592 with entries sorted according to an optimal leaf ordering for the hierarchical cluster tree

593 computed based on the distance values

594 **Supplementary Table 3:** Clustering results of entire kinase siRNA collection of the KiMorph

595 dataset

596 **Supplementary Table 4:** GO term enrichment analysis on the entire kinase siRNA collection of

597 the KiMorph dataset

598 **Supplementary Table 5:** Confusion matrix on annotated treatments of the BBBC021 dataset

599 **Supplementary Table 6:** Similarity values of all treatments to reference treatments of the

600 BBBC021 dataset

601 **Supplementary Table 7:** Distance matrix of entire compound collection of the BBBC021

602 dataset with entries sorted according to an optimal leaf ordering for the hierarchical cluster tree

603 computed based on the distance values

604 **Supplementary Table 8:** M-CNN architecture used to analyze both the KiMorph and BBBC021
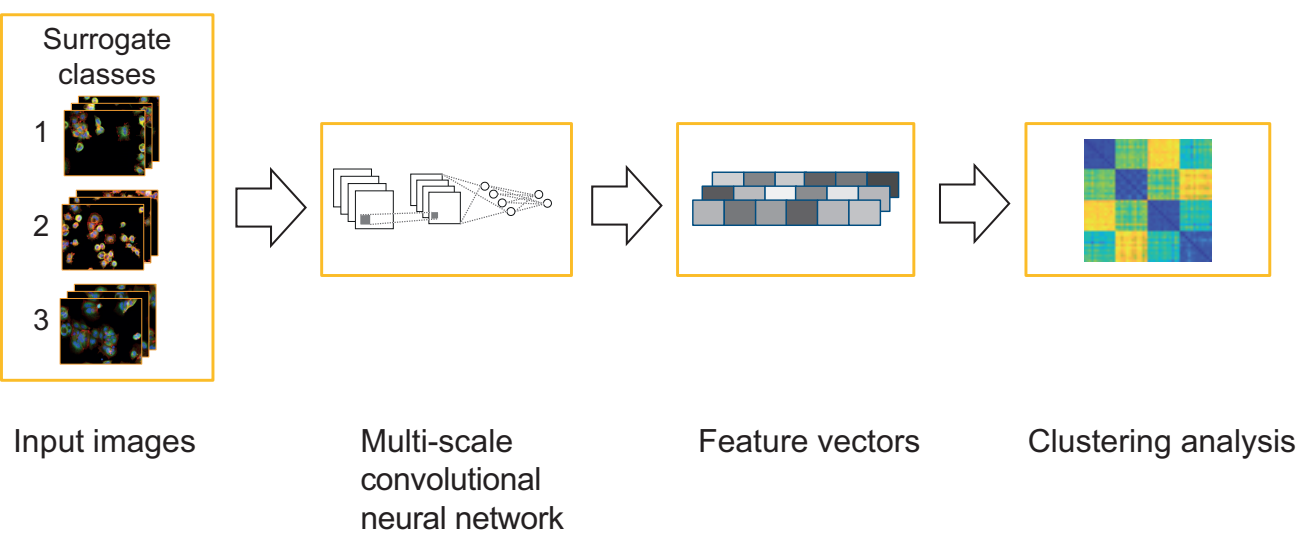
605 dataset

606

607 **Supplementary Figure 1:** Similarity-vs-concentration curves for selected compounds of the
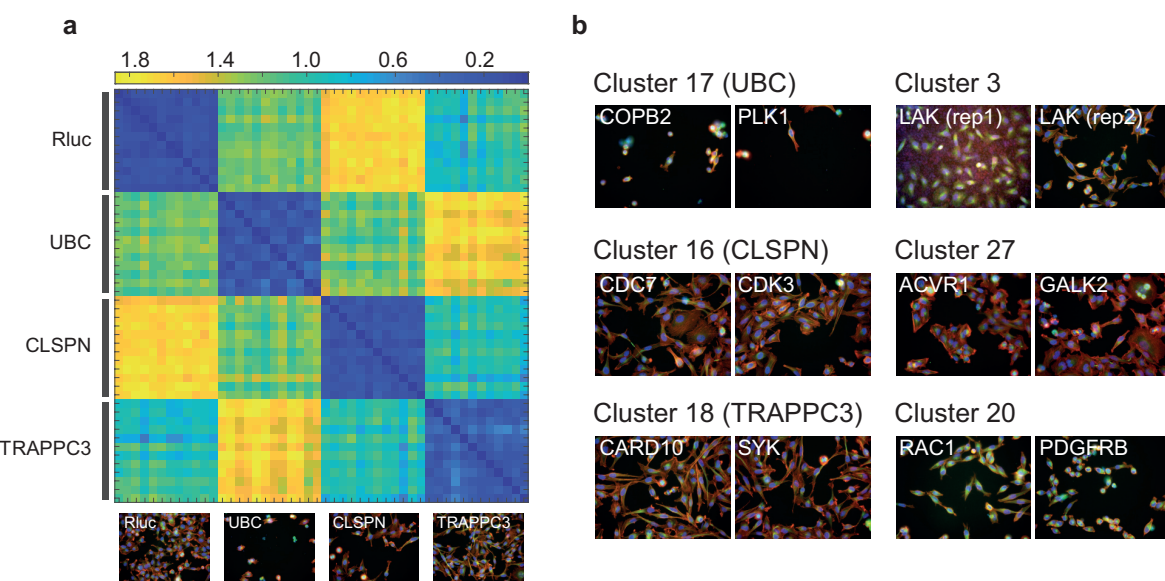
608 BBBC021 dataset.

609 **Supplementary Figure 2:** Example images of Cochicine at varying concentrations.

29

610

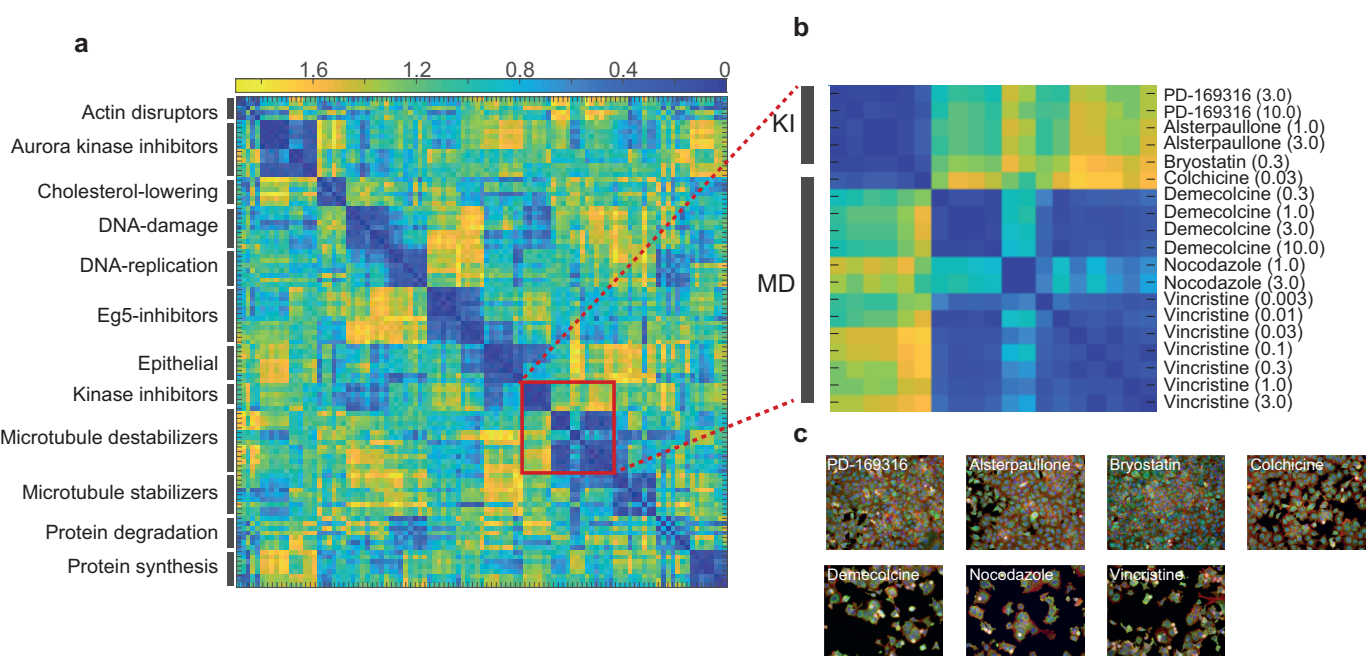611    **Supplementary Software**: Solver definition and network specification files.

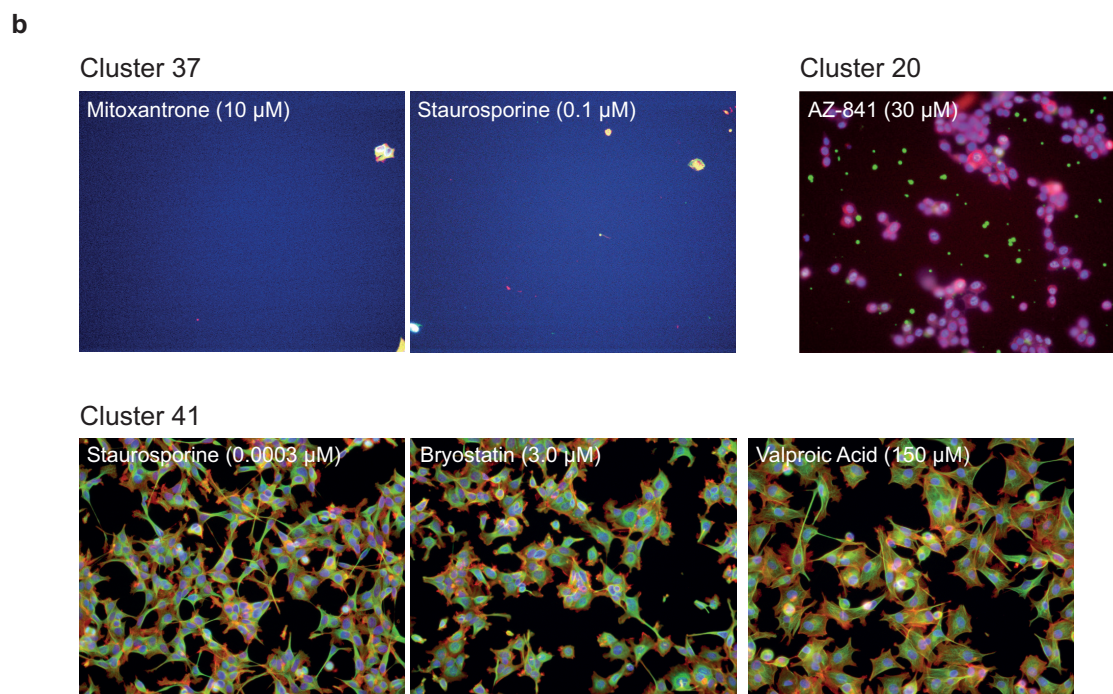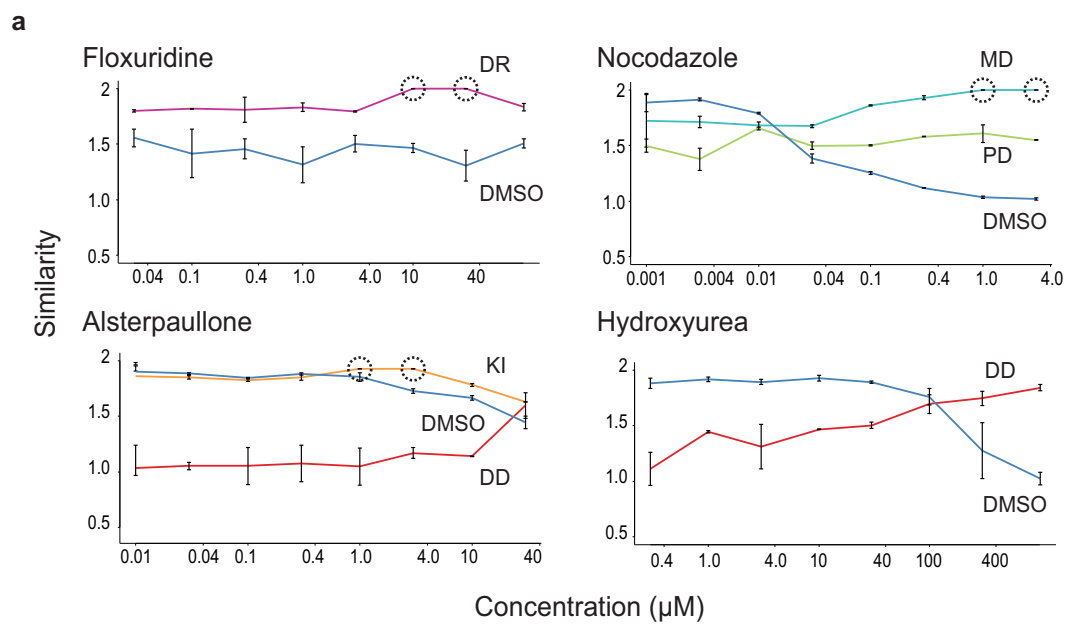Surrogate classes

1

2

3

Input images

Multi-scale convolutional neural network

Feature vectors

Clustering analysis

**Figure 1**

**a**



**b**



**Figure 2**

**Figure 3**

**a**



**b**



**Figure 4**