

1 **Cell morphology-guided de novo hit design by**
2 **conditioning generative adversarial networks on**
3 **phenotypic image features**

4
5 Oscar Méndez-Lucio,^{1,2} Paula Marin-Zapata,³ Joerg Wichard,⁴ David Rouquié,
6 ¹ Djork-Arné Clevert,³

7 ¹ Bayer SAS, Bayer CropScience, 355 rue Dostoïevski, CS 90153, 06906
8 Sophia Antipolis, France

9 ² Bloomoon, 13 Avenue Albert Einstein, 69100 Villeurbanne, France

10 ³ Department of Machine Learning Research, Bayer AG, 13353 Berlin, Germany

11 ⁴ Department of Genetic Toxicology, Bayer AG, 13353 Berlin, Germany

12
13 E-mail: oscar.mendezlucio.ext@bayer.com, paula.marinzapata@bayer.com,
14 joerg.wichard@bayer.com, david.rouquie@bayer.com, djork-
15 arne.clevert@bayer.com,

16 **Abstract**

17 **Developing new small molecules that are bioactive is time-consuming,**
18 **costly and rarely successful. As a mitigation strategy, we apply, for the**
19 **first time, generative adversarial networks to de novo design of small**
20 **molecules using a phenotype-based drug discovery approach. We trained**
21 **our model on a set of 30,000 compounds and their respective**
22 **morphological profiles extracted from high content images; no target**
23 **information was used to train the model. Using this approach, we were**
24 **able to automatically design agonist-like compounds of different**
25 **molecular targets.**

26

27 **Introduction**

28 Recent studies have categorized pharmaceutical research and development
29 (R&D) as a very inefficient process in terms of the high cost and small number
30 of approved molecules^{1,2}. Given that the cost of bringing a new drug to market
31 is doubling approximately every 9 years, the current R&D model might not be
32 sustainable in a few more decades¹. This concern raises the question of how to
33 improve the current pharmaceutical R&D efficiency. Some authors propose to
34 focus more into the final desired biological response, such as in a systems-
35 based approach, rather than the current target-based approach³⁻⁵. The latter is
36 characterized for focusing only on a well characterized target or mode of action.
37 On the other hand, systems-based drug discovery aims to identify or optimize a
38 molecule with little knowledge of the biological target or mode of action and
39 relying more on phenotypic changes. Recent studies have compared the
40 efficiency of these approaches based on the number of first-in-class drugs
41 approved by the FDA but it is too soon to draw a final conclusion due to the
42 long time frame of drug discovery projects^{6,7}.

43 Despite the many advantages of the systems-based approach³, only a few
44 experimental settings have been used for large-scale screenings. One of these
45 is transcriptomic analysis where the change in gene expression levels caused
46 by perturbation (either chemical or biological) is used as readout to select active
47 compounds. This approach has been used in the Connectivity Map project^{8,9} to
48 connect disease, genes and drugs and has been successfully applied to identify
49 new active molecules¹⁰⁻¹², find new uses for known drugs^{13,14}, and even for
50 mode of action identification^{15,16}. An alternative is to observe the effect of the
51 perturbation at morphological level rather than at transcriptomic level. In this
52 regard, the Cell Painting approach¹⁷, a new technique based on High-Content
53 Imaging (HCI), has been extremely useful. Cell Painting uses different dyes to
54 simultaneously stain several organelles and cellular components, thus capturing
55 information of the complete cellular state¹⁸. In the context of profiling, images
56 are usually processed by computational pipelines to extract feature
57 representations (morphological profiles)¹⁹⁻²³, which serve as phenotype
58 descriptors in further tasks. This technique has been used to cluster small
59 molecules by similar phenotypic effect¹⁸, for drug repurposing²⁴, to map cellular

60 morphology to gene function²⁵, and to predict results from biochemical
61 assays^{26,27} among other applications²⁸. Interestingly, most cell painting
62 applications to small molecules have focused on clustering or classification
63 tasks and, to our best knowledge, no use-case has combined these
64 information-rich assays with a molecular generative model^{29,30} to directly guide
65 molecular design.

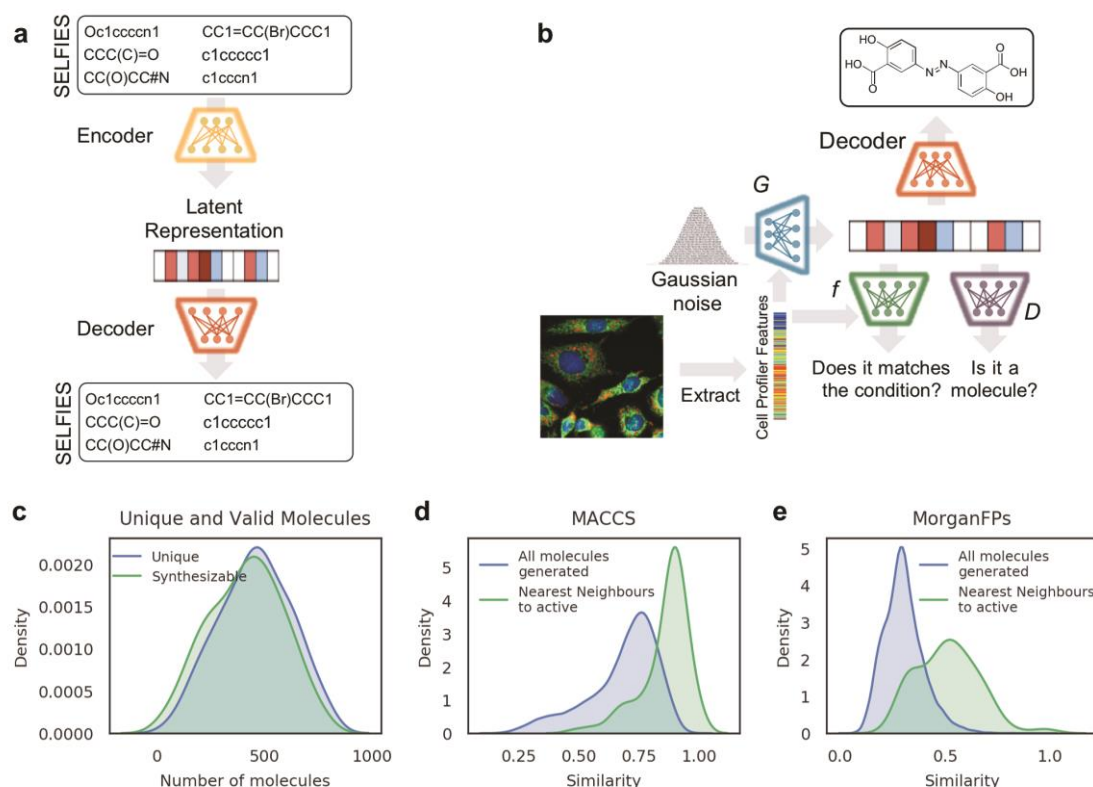
66 In the current work we combine cell painting morphological profiles with
67 generative adversarial networks (GANs)³¹ in order to automatically design
68 compounds that induce a specific morphological effect. GANs are generative
69 models that have been extensively used for molecular design, specifically for
70 the design of new active molecules and photovoltaic materials^{32,33}. In this paper
71 we show that by conditioning a GAN using morphological profiles we can
72 automatically design compounds with high molecular similarity to known
73 agonists of particular targets without having target or biological activity
74 annotations in the training set.

75

76 **Results**

77 The framework used in this study is divided in two stages. The first stage
78 consists in training a model capable of encoding molecular structures into a
79 continuous latent space, but also able to transform any point of the latent space
80 back into a valid molecular structure (Figure 1a). This encoder-decoder model
81 was designed following a similar implementation of the one proposed by Winter
82 et al³⁴, which is inspired on machine translation. The second stage is a
83 conditional generative adversarial network (GAN) capable of designing
84 compounds that can induce specific cell morphology changes (Figure 1b). More
85 specifically, the GAN is composed of three neural networks, namely generator,
86 discriminator and conditional. The task of the generator is to suggest points of
87 the latent space that correspond to molecules able to induce a desired
88 morphological profile. On the other hand, the discriminator evaluates if the
89 generated point of the latent space indeed corresponds to a proper molecule
90 whereas the conditional network calculates the probability of this molecule to
91 induce the desired morphological profile. The GAN is trained in an adversarial

92 manner meaning that at every training step the generator comes with better
 93 suggestions to fulfil the criteria of the discriminator, but at the same time the
 94 evaluation by the discriminator becomes stricter each time. This conditional
 95 GAN was trained with 126,779 morphological profiles induced by 30,616
 96 compounds in the U2OS cell line reported by Bray *et al.*³⁵ (see Methods for
 97 more details).

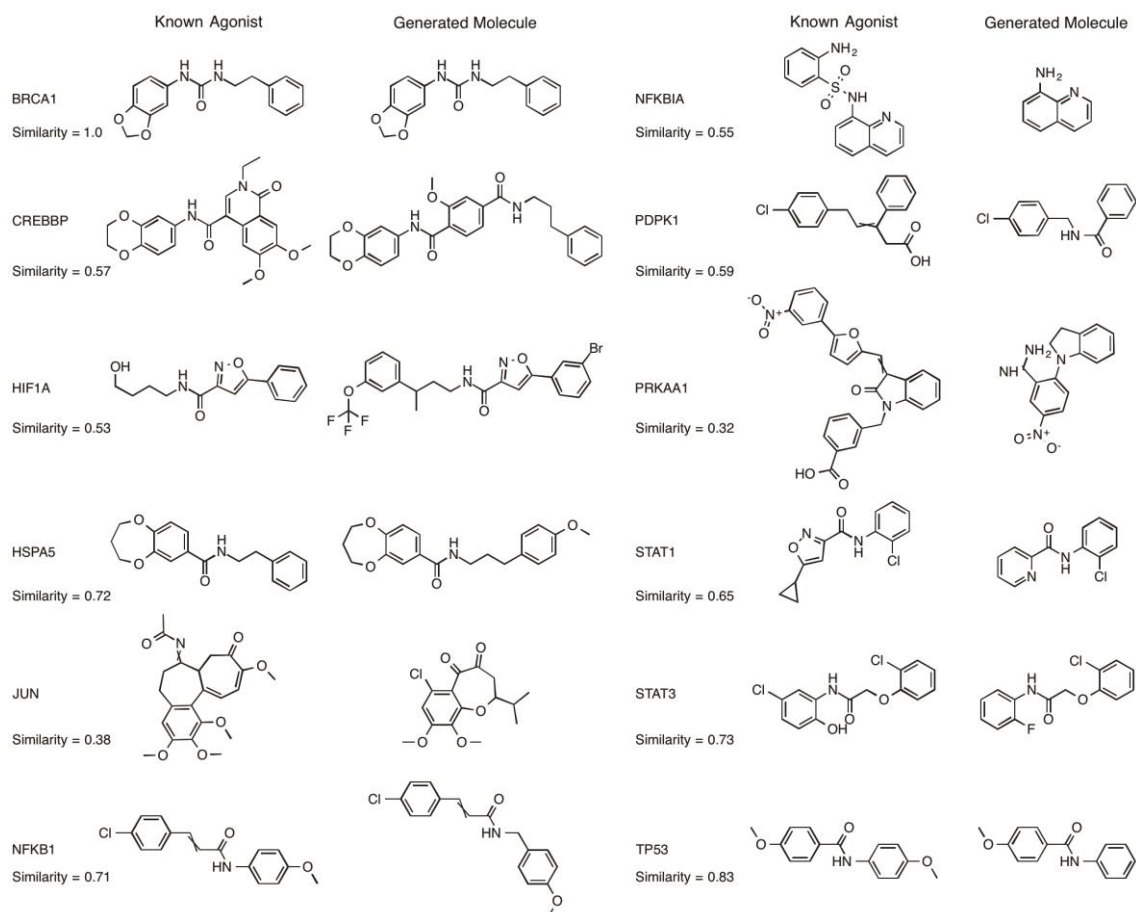


98

99 Figure 1. Graphical representation of the models and distribution of similarity.
 100 Molecules were encoded using a model that transforms SELFIES³⁶ of a
 101 molecule into a latent representation that can be later decoded back into
 102 SELFIES (a). The generative adversarial network in (b) has a generator (*G* in
 103 blue) that takes the desired morphological profile together with a vector of
 104 random noise to produce a molecular representation that can be decoded into
 105 SELFIES using the decoder (in red). The discriminator (*D* in purple) calculates
 106 the probability of the molecular representation to be a real molecule and the
 107 conditional network (*f* in green) calculates the probability of the molecular
 108 representation to match the morphological profile. Distributions shown in (c)
 109 represent the number of unique (blue) and synthesizable (green) molecules
 110 generated across the 114 morphological features from overexpressed genes.

111 Distributions of similarity measurements between all generated molecules and
112 their most similar known agonist were calculated using MACCS (d) and Morgan
113 fingerprints (e).

114 We used this framework to generate compounds that mimic the morphological
115 changes observed when overexpressing a specific target gene in the U2OS cell
116 line using cDNA. We focused on 12 genes that were individually
117 overexpressed by Rohban *et al.*²⁵ in U2OS using different cDNA open reading
118 frames (ORFs), namely *BRCA1*, *CREBBP*, *HIF1A*, *HSPA5*, *JUN*, *NFKB1*,
119 *NFKBIA*, *PDPK1*, *PRKAA1*, *STAT1*, *STAT3*, and *TP53*. Considering all ORFs
120 per gene, we computed a total of 114 morphological profiles and generated
121 1,000 molecules for each of them. On average, each morphological profile
122 produced 46.7% of valid molecules from which most of them (45.6%)
123 correspond to unique compounds and 41.4% presented a synthetic accessibility
124 score³⁷ < 4.5, meaning they have high chances to be synthesized (Figure 1c).
125 To test the ability of our model to generate phenotype-tailored molecules, we
126 evaluated the similarity between the generated compounds and known agonists
127 of the 12 selected targets. Agonists were retrieved from ExCAPE database³⁸
128 making sure to exclude all of them that were present in the training set. Overall,
129 generated molecules share similar chemical groups with the known agonist
130 compounds (median MACCS³⁹ similarity = 0.72, IQR = 0.18) but also share a
131 good proportion of molecular substructures (median Morgan Fingerprints
132 similarity = 0.30, IQR = 0.12) as shown in Figures 1d and 1e. Figure 2 shows
133 examples of the generated molecules presenting the highest similarity to a
134 known active molecule for each of the overexpressed targets.



135

136 Figure 2. Chemical structures of molecules generated by the model. These
 137 molecules were obtained by conditioning the model using the microscopy
 138 images of 12 different overexpressed genes and are depicted next to their
 139 closest active nearest neighbour.

140 A scaffold analysis of the generated molecules revealed that the model can
 141 generate molecules containing the same scaffolds as the known active
 142 molecules. Tables 1 and 2 show the number of shared scaffolds per target.
 143 Interestingly, more than 20% of molecules generated for *BRCA1* and *NFKB1*
 144 contain a Murcko scaffold that is known to be active. This rate is considerably
 145 increased when using Murcko generic scaffolds, which is an abstraction form of
 146 scaffold where all atoms are converted into carbon atoms and all bonds
 147 transform into single bonds. In this case, more than 50% of molecules
 148 generated for *BRCA1* and *NFKB1* and more than 25% for *TP53* and *HSPA5*
 149 contain a known active generic Murcko scaffold. Interestingly, none of the
 150 molecules generated for *JUN* and *PRKAA1* contain a known active Murcko
 151 scaffold or generic scaffold. Nonetheless, it is worth noting that only less than

152 five Murcko scaffolds have been reported in the ExCAPE database to have an
 153 agonist effect on these targets.

154

155 Table 1. Comparison of Murcko scaffolds in generated molecules and known
 156 agonists from ExCAPE database.

	Valid and synthetic accessible molecules	Unique scaffolds in known agonists	Generated molecules containing an active scaffold	Unique active scaffolds present in generated molecules
<i>PDPK1</i>	4624	1	0	0
<i>BRCA1</i>	3444	4632	861	97
<i>NFKB1</i>	3343	256	703	39
<i>STAT1</i>	4304	78	61	10
<i>NFKBIA</i>	2435	4	93	1
<i>STAT3</i>	6824	56	359	9
<i>PRKAA1</i>	2410	3	0	0
<i>CREBBP</i>	2591	80	0	0
<i>HIF1A</i>	4665	38	1	1
<i>JUN</i>	5787	4	0	0
<i>TP53</i>	4250	7522	718	126
<i>HSPA5</i>	2592	551	133	24

157

158 Table 2. Comparison of Murcko generic scaffolds in generated molecules and
 159 known agonists from ExCAPE database.

	Valid and synthetic accessible molecules	Unique generic scaffolds in known agonists	Generated molecules containing an active generic scaffold	Unique active generic scaffolds present in generated molecules
<i>PDPK1</i>	4624	1	19	1
<i>BRCA1</i>	3444	2146	2248	237
<i>NFKB1</i>	3343	201	1783	66
<i>STAT1</i>	4304	62	562	36
<i>NFKBIA</i>	2435	4	109	1
<i>STAT3</i>	6824	51	835	21
<i>PRKAA1</i>	2410	2	0	0
<i>CREBBP</i>	2591	76	100	14
<i>HIF1A</i>	4665	32	141	7
<i>JUN</i>	5787	4	0	0
<i>TP53</i>	4250	4627	1851	374
<i>HSPA5</i>	2592	406	677	90

160

161 In order to choose the most promising compounds for experimental evaluation
 162 we filtered all generated compounds based on physicochemical properties,
 163 namely LogP between 0 and 7, molecular weight less than 750 daltons, number
 164 of hydrogen bond donors and acceptors below 10 and less than 10 rotatable
 165 bonds. This filtering step reduced the total number of molecules to 39,801.
 166 Additionally, we removed those compounds that contained non medicinal
 167 chemistry friendly SMART and toxalert^{40,41} leaving a total number of 33,341
 168 compounds. Since we are interested in novel molecules, we filtered out those
 169 compounds with Morgan Fingerprints similarity higher than 0.8 compared to
 170 compounds in the training set, which reduced the number of compounds to a
 171 total of 33,309. In a further step, we used in-house models to rank generated
 172 molecules based on their probability to cross the cellular membrane in a Caco-2
 173 assay. These helped us to choose at least 100 molecules for each target that
 174 were examined by an experienced chemist. Some of these molecules are
 175 currently being synthesized and tested on a cell painting assay.

176

177 Table 3. Number of generated molecules per target after each of the filtering
 178 steps.

	Valid and synthetic accessible molecules	After physicochemical filters	After removing compounds with ToxAlerts	Compounds with low similarity to the training set*
<i>BRCA1</i>	3444	2843	2347	2342
<i>CREBBP</i>	2591	2185	1905	1902
<i>HIF1A</i>	4665	4107	3564	3562
<i>HSPA5</i>	2592	2373	2104	2102
<i>JUN</i>	5787	4071	3405	3403
<i>NFKB1</i>	3343	3184	2634	2626
<i>NFKBIA</i>	2435	2136	1699	1699
<i>PDPK1</i>	4624	4059	3447	3443
<i>PRKAA1</i>	2410	1883	1489	1486
<i>STAT1</i>	4304	3849	3337	3335
<i>STAT3</i>	6824	5589	4664	4663
<i>TP53</i>	4250	3522	2746	2746

179 *Similarity threshold set to 0.8 using Morgan Fingerprints and Dice distance

180

181

182 **Conclusions**

183 In this work, we showed that it is possible to combine artificial intelligence with
184 morphological cellular images in order to design new compounds. In this
185 specific proof of concept, the model was able to use the information contained
186 in a genetically induced phenotypic profile to produce molecules that present
187 high structural similarity to known agonists of a particular target. We expect that
188 this work will help to attract more attention to systems-based approach,
189 particularly to high throughput imaging screening. In addition, we expect that
190 more computational tools, like the one presented in this paper, will help
191 increase the efficiency of the systems-based approaches in order to accelerate
192 drug discovery R&D.

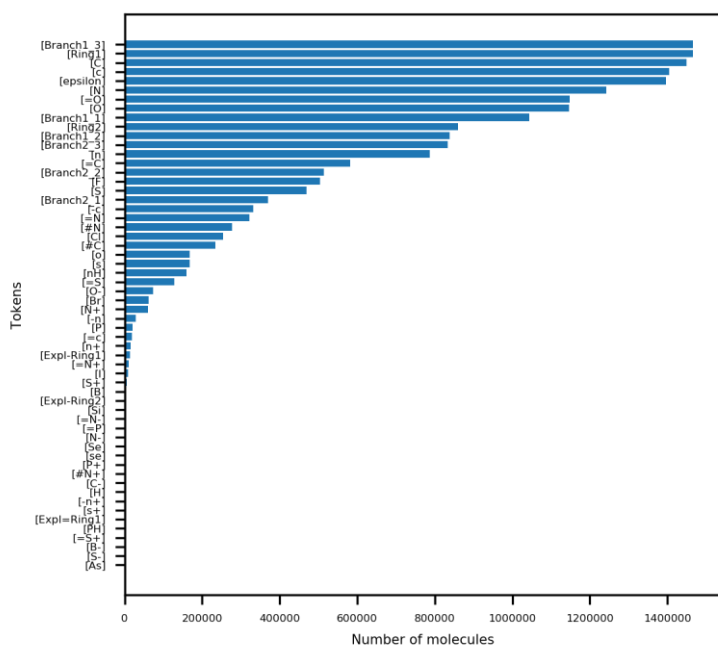
193

194 **Methods**

195 **Data set:** For this study we used a data set coming from a high-content
196 screening of 30,616 compounds that followed the cell painting protocol.³⁵ More
197 specifically, Human Bone Osteosarcoma Epithelial cells (U2OS cells) were
198 plated and treated with each of these compounds (details can be found
199 elsewhere¹⁷) in quadruplicate. Each well was imaged using 5 fluorescent
200 channels and 6 fields of view. The different dyes used in cell painting were
201 chosen to stain different organelles or cellular components including: nucleus,
202 endoplasmic reticulum, nucleoli, cytoplasmic RNA, cytoskeleton, Golgi, plasma
203 membrane and mitochondria.^{17,18} Raw images from different fields and
204 channels were processed using CellProfiler²³ to extract morphological profiles
205 composed by 1,783 quantitative features (e.g. cellular shape, intensity, texture,
206 etc.) able to capture slight morphological changes induced by the treatment.
207 More details on the experimental and processing protocols can be found in the
208 original publication.³⁵ In the end, this data set is composed by 126,779
209 morphological profiles corresponding to 30,616 compounds. Notice that each
210 compound is associated with at least 4 population-averaged morphological
211 profiles coming from replicates, which were treated independently.

212 As a test set we used morphological profiles calculated from images taken
213 during a high-content profiling campaign of gene gain-of-function recently
214 published.²⁵ In that work, Rohban *et al.* profiled the overexpression of 323 open
215 reading frames (ORF) constructs also following the cell painting protocol.
216 Illumination corrected images for the test set were obtained from The Image
217 Data Resource (IDR) web API. To assure comparability of profiles between test
218 and training sets, we closely followed the feature extraction protocol from Bray
219 *et al.*³⁵. Single-cell profiles were calculated with CellProfiler 2.2.0 using the
220 analysis.cppipe pipeline, which was slightly modified to take input .csv files.
221 Population-averaged profiles were computed as the median among all cells
222 from all fields of view per well. As sanity check, we confirmed that our feature
223 extraction methodology was able to reproduce the single-cell and population-
224 averaged profiles reported by Bray *et al.*³⁵ for a randomly chosen plate (plate
225 number 25690).

226 All molecular structures were pre-processed using the MolStandardize module
227 from RDKit⁴². First, molecules as SMILES codes were standardized by
228 removing hydrogens, sanitizing the molecules, disconnecting metals and
229 normalizing the structure. Then, only the largest fragment of each molecule was
230 kept, and their charges were neutralized. Finally, the uncharged molecules were
231 transformed back into a non-isomeric form of canonical SMILES. Despite
232 SMILES can capture information regarding the molecular structure, they do not
233 capture the molecular graph connectivity. In order to overcome this issue, we
234 further transform SMILES into self-referencing embedded strings (SELFIES)³⁶
235 which are able to encode molecular graphs in a simple and deterministic way. In
236 order to use SELFIES as input for our model, they were transformed into a one-
237 hot encoding format just keeping those tokens that appear in at least 100
238 molecules from the 1.5 million contained in ChEMBL22⁴³ (Figure 3). Molecules
239 containing different SELFIES tokens to the ones selected previously or SMILES
240 strings larger than 120 character were removed from the dataset.



247 case it learns to encode SELFIES of a molecule (as one-hot encoding) into a
248 continuous latent vector, which then can be decoded back into SELFIES. The
249 architecture of this model is based on the previous work of Winter et al.^{34,45},
250 where the encoder is comprised by three stacked Gate Recurrent Unit (GRU)
251 cells where their resulting cell states are concatenated and fed into two
252 independent fully connected layers which predict the mean and standard
253 deviation of a distribution from which the latent vector will be sampled for the
254 variational approach. After sampling, the latent vectors are compressed by a
255 tanh activation function to produce the final latent vectors of 256 dimensions
256 with values between -1 and 1. The decoder takes as input the latent vectors
257 which are fed into a fully connected layer two expand them from 256 to 768
258 dimensions. This vector is split into three vectors of 256 dimensions and used
259 as initial states of another three stacked GRU cells. The output of the stacked
260 GRU was followed by a dropout layer (rate of 0.2) and a dense layer (with a
261 softmax activation function), which generates a probability distribution over all
262 possible SELFIES tokens for each time step. This variational molecular
263 autoencoder was trained following a teacher-forcing⁴⁷ scheme during 10 epochs
264 on 1.25 million molecules extracted from ChEMBL 22⁴³ complemented by the
265 molecules in the training set.

266 **Generator architecture (G(z,c)):** The generator receives as input the condition,
267 in this case a morphological profile (of 1,449 features), and a 1,000-
268 dimensional noise vector sampled from a normal distribution. The input noise
269 is processed by a 2-layer multilayer perceptron (MLP) with 512 and 256 nodes,
270 respectively, where each layer uses LeakyRelu as activation function. At the
271 same time, the morphological features are processed by an MLP of size [1024,
272 512, 256] also using LeakyRelu after each layer. The two resulting tensors are
273 concatenated and used as input for another 2-layer MLP, where the first layer
274 has 256 nodes and LeakyRelu activation function. The second layer acts as an
275 output layer (i.e. the number of nodes is equal to the dimensionality of the latent
276 space) and is followed by a tanh activation function.

277 **Discriminator architecture:** The discriminator is composed of a 4-layer MLP of
278 [256, 256, 256, 1] hidden units with LeakyRelu activation function in the first
279 three layers. In order to reduce overfitting, dropout with rate of 0.4 was used

280 between the second and third hidden layers and between the third and the last
281 layer of the MLP.

282 **Conditional network architecture:** This network estimates the probability of a
283 molecular representation to match a morphological profile. Here the
284 morphological profile is processed by a MLP of 3 layers with 1024, 512 and 256
285 units, respectively, and then regularized by a dropout layer with rate of 0.4. In a
286 similar way, the latent space coordinates of the compound are also fed into a 2-
287 layer MLP with dimension [256, 256] and finalized with a dropout layer. The
288 outputs of these two MLP, corresponding to the processed morphological profile
289 and compound information, were concatenated and used as input of an MLP of
290 size [256, 1] using LeakyRelu and sigmoid activation functions, which estimates
291 the probability of a molecule to produce a certain morphological profile.

292 **Training:** The conditional generative adversarial network was trained during
293 500 epochs using a batch size of 256. Here, an epoch was composed by 496
294 steps where the weights of the discriminator were updated after each step,
295 whereas those of the generator every 10 steps. The network was trained using
296 the RMSprop optimizer with learning rate of 1×10^{-4} for both the generator and
297 discriminator. The weights of the conditional network were pretrained and
298 frozen during the GAN training process. All neural networks were built and
299 trained using Tensorflow 1.14⁴⁸.

300 **Model evaluation:** During training, we generated a molecular representation for
301 each morphological profile in the training set at the end of each epoch. These
302 were used to evaluate the similarity between the generated and real molecular
303 representations using Fréchet distance. The Fréchet distance measures the
304 similarity between two distributions (in this case the real and generated) and
305 was recently proposed as an efficient way to evaluate the efficacy of generative
306 models.⁴⁹ This metric takes the mean and covariance of the real distribution (μ_r
307 and C_r , respectively), together with the mean and covariance of the generated
308 distribution (μ_g and C_g) in the following formulation:

$$309 \quad d^2 \left((\mu_r, C_r), (\mu_g, C_g) \right) = \|\mu_g - \mu_r\|_2^2 + Tr(C_g + C_r - 2(C_g C_r)^{1/2})$$

310 **Generating molecules from features of morphological images:** We used
311 114 morphological profiles from the ones reported by Rohban et al. after
312 overexpressing open reading frames (ORFs) corresponding to 12 different
313 genes: *BRCA1*, *CREBBP*, *HIF1A*, *HSPA5*, *JUN*, *NFKB1*, *NFKBIA*, *PDPK1*,
314 *PRKAA1*, *STAT1*, *STAT3*, and *TP53*. Known agonists of the corresponding
315 protein expressed by these genes were extracted from the ExCAPE database³⁸.
316 For these 12 targets we generated 1000 molecular representation for each
317 morphological profile (114 in total) and evaluated the model by comparing the
318 generated molecules to the known agonists of these targets. Dice similarity
319 between the generated compounds and the molecules with known agonistic
320 effect was computed using MACCS keys and Morgan Fingerprints (radius = 3,
321 1024 bits) with RDKit.⁴² Despite that the similarity between each generated
322 molecule and each known agonist was calculated, only the one corresponding
323 to the closest known agonist was recorded for each generated molecule and
324 used for further analysis. It is worth mentioning that during these validation
325 tasks both the number of valid molecules (validity measure) and the number of
326 unique molecules (uniqueness measure) were recorded as sanity check for the
327 generative model.

328 **Author contributions**

329 O.M.L. performed the computational experiments and designed the deep
330 learning model. P.M.Z pre-processed all images in the data set, calculated
331 features from morphological images and helped with the manuscript
332 preparation. D.R. and J.W. provided guidance and helped with the manuscript
333 preparation. D.A.C. conceived the study and supervised the work. O.M.L.,
334 P.M.Z, J.W., D.R. and D.A.C. read and approved the manuscript.

335 **Acknowledgments**

336 Authors thank Joerg Tiebes for providing chemistry expert feedback and for his
337 useful comments. We are also grateful to Arwa Al-Dilaimi and Angela Becker for
338 supporting the project and for insightful discussions.

339 **Competing interests**

340 D.A.C., P.M.Z and J.W. are employees of Bayer AG. O.M.L and D.R. work

341 directly or indirectly for Bayer SAS.

342 **Data availability**

343 The data that support the findings of this study will be available on request from
344 the corresponding author after the final publication of this manuscript.

345 **Code availability**

346 The code used to generate results shown in this study will be available from the
347 corresponding author upon request after the final publication of this manuscript.

348 **References**

- 349 1. Scannell, J. W., Blanckley, A., Boldon, H. & Warrington, B. Diagnosing
350 the decline in pharmaceutical R&D efficiency. *Nature Reviews Drug*
351 *Discovery* vol. 11 191–200 (2012).
- 352 2. Pammolli, F., Magazzini, L. & Riccaboni, M. The productivity crisis in
353 pharmaceutical R&D. *Nat. Rev. Drug Discov.* **10**, 428–438 (2011).
- 354 3. Butcher, E. C. *Can cell systems biology rescue drug discovery? Nature*
355 *Reviews Drug Discovery* vol. 4 www.nature.com/reviews/drugdisc (2005).
- 356 4. Zheng, W., Thorne, N. & McKew, J. C. Phenotypic screens as a renewed
357 approach for drug discovery. *Drug Discov. Today* **18**, 1067–1073 (2013).
- 358 5. Moffat, J. G., Vincent, F., Lee, J. A., Eder, J. & Prunotto, M. Opportunities
359 and challenges in phenotypic drug discovery: An industry perspective.
360 *Nat. Rev. Drug Discov.* **16**, 531–543 (2017).
- 361 6. Swinney, D. C. & Anthony, J. How were new medicines discovered? *Nat.*
362 *Rev. Drug Discov.* **10**, 507–519 (2011).
- 363 7. Eder, J., Sedrani, R. & Wiesmann, C. The discovery of first-in-class
364 drugs: Origins and evolution. *Nat. Rev. Drug Discov.* **13**, 577–587 (2014).
- 365 8. Lamb, J. *et al.* The connectivity map: Using gene-expression signatures
366 to connect small molecules, genes, and disease. *Science (80-.).* **313**,
367 1929–1935 (2006).
- 368 9. Subramanian, A. *et al.* A Next Generation Connectivity Map: L1000
369 Platform and the First 1,000,000 Profiles. *Cell* **171**, 1437-1452.e17
370 (2017).
- 371 10. Hieronymus, H. *et al.* Gene expression signature-based chemical
372 genomic prediction identifies a novel class of HSP90 pathway modulators.
373 *Cancer Cell* **10**, 321–330 (2006).
- 374 11. Wei, G. *et al.* Gene expression-based chemical genomics identifies
375 rapamycin as a modulator of MCL1 and glucocorticoid resistance. *Cancer*

- 376 *Cell* **10**, 331–342 (2006).
- 377 12. De Wolf, H. *et al.* High-Throughput Gene Expression Profiles to Define
378 Drug Similarity and Predict Compound Activity. *Assay Drug Dev. Technol.*
379 **16**, 162–176 (2018).
- 380 13. Aliper, A. *et al.* Deep learning applications for predicting pharmacological
381 properties of drugs and drug repurposing using transcriptomic data. *Mol.*
382 *Pharm.* **13**, 2524–2530 (2016).
- 383 14. Iorio, F., Rittman, T., Ge, H., Menden, M. & Saez-Rodriguez, J.
384 Transcriptional data: A new gateway to drug repositioning? *Drug Discov.*
385 *Today* **18**, 350–357 (2013).
- 386 15. Iwata, M., Sawada, R., Iwata, H., Kotera, M. & Yamanishi, Y. Elucidating
387 the modes of action for bioactive compounds in a cell-specific manner by
388 large-scale chemically-induced transcriptomics. *Sci. Rep.* **7**, 40164
389 (2017).
- 390 16. Wacker, S. A., Houghtaling, B. R., Elemento, O. & Kapoor, T. M. Using
391 transcriptome sequencing to identify mechanisms of drug action and
392 resistance. *Nat. Chem. Biol.* **8**, 235–237 (2012).
- 393 17. Bray, M. A. *et al.* Cell Painting, a high-content image-based assay for
394 morphological profiling using multiplexed fluorescent dyes. *Nat. Protoc.*
395 **11**, 1757–1774 (2016).
- 396 18. Gustafsdottir, S. M. *et al.* Multiplex cytological profiling assay to measure
397 diverse cellular states. *PLoS One* **8**, 1–7 (2013).
- 398 19. Scheeder, C., Heigwer, F. & Boutros, M. Machine learning and image-
399 based profiling in drug discovery. *Current Opinion in Systems Biology* vol.
400 10 43–52 (2018).
- 401 20. Moen, E. *et al.* Deep learning for cellular image analysis. *Nat. Methods*
402 doi:10.1038/s41592-019-0403-1 (2019) doi:10.1038/s41592-019-0403-1.
- 403 21. Jackson, P. T. *et al.* Phenotypic profiling of high throughput imaging
404 screens with generic deep convolutional features. *Proc. 16th Int. Conf.*

- 405 *Mach. Vis. Appl. MVA 2019* 1–4 (2019)
406 doi:10.23919/MVA.2019.8757871.
- 407 22. Lu, A. X., Kraus, O. Z., Cooper, S. & Moses, A. M. Learning unsupervised
408 feature representations for single cell microscopy images with paired cell
409 inpainting. *PLoS Comput. Biol.* **15**, 1–27 (2019).
- 410 23. Carpenter, A. E. *et al.* CellProfiler: image analysis software for identifying
411 and quantifying cell phenotypes. *Genome Biol.* **7**, R100 (2006).
- 412 24. Gibson, C. C. *et al.* Strategy for identifying repurposed drugs for the
413 treatment of cerebral cavernous malformation. *Circulation* **131**, 289–299
414 (2015).
- 415 25. Rohban, M. H. *et al.* Systematic morphological profiling of human gene
416 and allele function via Cell Painting. *Elife* **6**, 1–23 (2017).
- 417 26. Simm, J. *et al.* Repurposing High-Throughput Image Assays Enables
418 Biological Activity Prediction for Drug Discovery. *Cell Chem. Biol.* **25**, 611-
419 618.e3 (2018).
- 420 27. Hofmarcher, M., Rumetshofer, E., Clevert, D. A., Hochreiter, S. &
421 Klambauer, G. Accurate Prediction of Biological Assays with High-
422 Throughput Microscopy Images and Convolutional Networks. *J. Chem.*
423 *Inf. Model.* **59**, 1163–1171 (2019).
- 424 28. Caicedo, J. C., Singh, S. & Carpenter, A. E. Applications in image-based
425 profiling of perturbations. *Current Opinion in Biotechnology* vol. 39 134–
426 142 (2016).
- 427 29. Butler, K. T., Davies, D. W., Cartwright, H., Isayev, O. & Walsh, A.
428 Machine learning for molecular and materials science. *Nature* vol. 559
429 547–555 (2018).
- 430 30. Sanchez-Lengeling, B. & Aspuru-Guzik, A. Inverse molecular design
431 using machine learning: Generative models for matter engineering.
432 *Science (80-.)*. **361**, 360–365 (2018).
- 433 31. Goodfellow, I. J. *et al.* Generative adversarial nets. in *Advances in Neural*

- 434 *Information Processing Systems* vol. 3 2672–2680 (Curran Associates,
435 Inc., 2014).
- 436 32. Guimaraes, G. L., Sanchez-Lengeling, B., Outeiral, C., Farias, P. L. C. &
437 Aspuru-Guzik, A. Objective-Reinforced Generative Adversarial Networks
438 (ORGAN) for Sequence Generation Models. Preprint at
439 <http://arxiv.org/abs/1705.10843> (2017).
- 440 33. Kadurin, A., Nikolenko, S., Khrabrov, K., Aliper, A. & Zhavoronkov, A.
441 DruGAN: An Advanced Generative Adversarial Autoencoder Model for de
442 Novo Generation of New Molecules with Desired Molecular Properties in
443 Silico. *Mol. Pharm.* **14**, 3098–3104 (2017).
- 444 34. Winter, R., Montanari, F., Noé, F. & Clevert, D. A. Learning continuous
445 and data-driven molecular descriptors by translating equivalent chemical
446 representations. *Chem. Sci.* **10**, 1692–1701 (2019).
- 447 35. Bray, M. A. *et al.* A dataset of images and morphological profiles of 30
448 000 small-molecule treatments using the Cell Painting assay.
449 *Gigascience* **6**, 1–5 (2017).
- 450 36. Krenn, M., Häse, F., Nigam, A., Friederich, P. & Aspuru-Guzik, A.
451 SELFIES: a robust representation of semantically constrained graphs with
452 an example application in chemistry. (2019).
- 453 37. Ertl, P. & Schuffenhauer, A. Estimation of synthetic accessibility score of
454 drug-like molecules based on molecular complexity and fragment
455 contributions. *J. Cheminform.* **1**, 1–11 (2009).
- 456 38. Sun, J. *et al.* ExCAPE-DB: An integrated large scale dataset facilitating
457 Big Data analysis in chemogenomics. *J. Cheminform.* **9**, 1–9 (2017).
- 458 39. Durant, J. L., Leland, B. A., Henry, D. R. & Nourse, J. G. Reoptimization
459 of MDL keys for use in drug discovery. *J. Chem. Inf. Comput. Sci.* **42**,
460 1273–1280 (2002).
- 461 40. Non MedChem-Friendly SMARTS.
462 <https://www.surechembl.org/knowledgebase/169485-non-medchem->

- 463 friendly-smarts.
- 464 41. Sushko, I., Salmina, E., Potemkin, V. A., Poda, G. & Tetko, I. V.
465 ToxAlerts: A web server of structural alerts for toxic chemicals and
466 compounds with potential adverse reactions. *J. Chem. Inf. Model.* **52**,
467 2310–2316 (2012).
- 468 42. Landrum, G. A. RDKit: Open-source cheminformatics.
469 <http://www.rdkit.org>.
- 470 43. Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Res.* **45**,
471 D945–D954 (2017).
- 472 44. Gómez-Bombarelli, R. *et al.* Automatic Chemical Design Using a Data-
473 Driven Continuous Representation of Molecules. *ACS Cent. Sci.* **4**, 268–
474 276 (2018).
- 475 45. Winter, R. *et al.* Efficient multi-objective molecular optimization in a
476 continuous latent space. *Chem. Sci.* **10**, 8016–8024 (2019).
- 477 46. Méndez-Lucio, O., Baillif, B., Clevert, D.-A., Rouquié, D. & Wichard, J. De
478 Novo Generation of Hit-like Molecules from Gene Expression Signatures
479 Using Artificial Intelligence. doi:10.26434/chemrxiv.7294388.v1.
- 480 47. Williams, R. J. & Zipser, D. A Learning Algorithm for Continually Running
481 Fully Recurrent Neural Networks. *Neural Comput.* **1**, 270–280 (1989).
- 482 48. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on
483 Heterogeneous Distributed Systems. Preprint at
484 <http://arxiv.org/abs/1603.04467> (2016).
- 485 49. Preuer, K., Renz, P., Unterthiner, T., Hochreiter, S. & Klambauer, G.
486 Fréchet ChemNet Distance: A Metric for Generative Models for Molecules
487 in Drug Discovery. *J. Chem. Inf. Model.* **58**, 1736–1741 (2018).

488