

CLOOME: contrastive learning unlocks bioimaging databases for queries with chemical structures

Ana Sanchez-Fernandez* Elisabeth Rumetshofer*

Sepp Hochreiter * † Günter Klambauer*

* ELLIS Unit Linz and LIT AI Lab, Institute for Machine Learning,

Johannes Kepler University Linz, Austria

† Institute of Advanced Research in Artificial Intelligence (IARAI), Vienna, Austria

Abstract

Currently, bioimaging databases cannot be queried by chemical structures that induce the phenotypic effects captured by the image. We present a novel retrieval system based on contrastive learning that is able to identify the chemical structure inducing the phenotype out of ~2,000 candidates with a top-1 accuracy >70 times higher than a random baseline.

Brief Communication

Biological and chemical databases and their querying mechanisms are at the heart of research in molecular biology. Sequence databases, such as RefSeq [1] or UniProt [2], contain DNA or protein sequences, and are often queried with a given sequence using BLAST [3] or its variants. Genome databases [4] usually allow for multiple types of querying methods, such as genetic location, gene names, or accession numbers. Protein structure databases, for example, the Protein Data Bank (PDB) [5], offer a range of querying approaches from sequence similarities to structural queries based on 3D shape. The chemical databases ChEMBL [6] and PubChem are huge corpora of chemical structures that contain billions of small molecules. The International Chemical Identifier (InChI) [7] was designed to facilitate querying for chemical structures in such databases, which is difficult because of the graph matching problem. While BLAST, the structural search in PDB, and the InChI-based queries can be considered as associative or content-based querying, bioimaging databases still rely on manual annotation and text-based search. However, querying large bioimaging databases by a chemical structure that induces the phenotypic effect captured by the image could considerably empower biomedical research. Additionally, unlocking chemical databases for queries with a microscopy image capturing the phenotypic effects of a chemical structure could be equally important (see Figure 1A,B).

Recently, contrastive learning has emerged as a powerful paradigm to learn rich representations [8]. The contrastive learning methods CLIP and CLOOB embed natural language and images into the same representation space [9, 10]. Contrastive learning enforces that images and their matching captions are close to one another in this embedding space, while un-matched images and captions are separated. Therefore, text prompts can query an image database by extracting nearby images in the embedding space and vice versa [9]. These text-image embedding spaces enabled the generation of realistic images from short text prompts and led to the recent boom of "AI art" [11]. In this work, we use these powerful contrastive learning paradigms to enable retrieval or querying systems for microscopy images.

In order to characterize cell phenotypes, tissues, or cellular processes, microscopy imaging has been used as an informative and time- and cost-efficient biotechnology [12, 13]. Consequently, there have been efforts by the scientific community to use high-throughput microscopy imaging [14] as informative read-out and characterization of cellular systems and phenotypes under diverse perturbations [13, 15]. In addition to the wealth of information that is comprehensible and informative for human experts, these microscopy images also contain large amounts of biological information inaccessible to humans, but which can be successfully extracted by computational methods, such as Deep Learning [16]. The immense amount of microscopy imaging data are stored in large databases, many of which are publicly available. Their querying procedures, however, are still limited to queries by textual annotations. A common embedding space of (a) microscopy images capturing phenotypic effects of perturbations, and (b) chemical structures inducing those effects would allow for content-based or associative querying of both imaging and chemical databases. Such an embedding space would represent cellular processes both in terms of the chemical structures that induce them and in terms of images that capture the cell phenotypes caused by these processes. New applications such as the detection of novel cell phenotypes are possible through such embedding spaces (see Figure 4C,F).

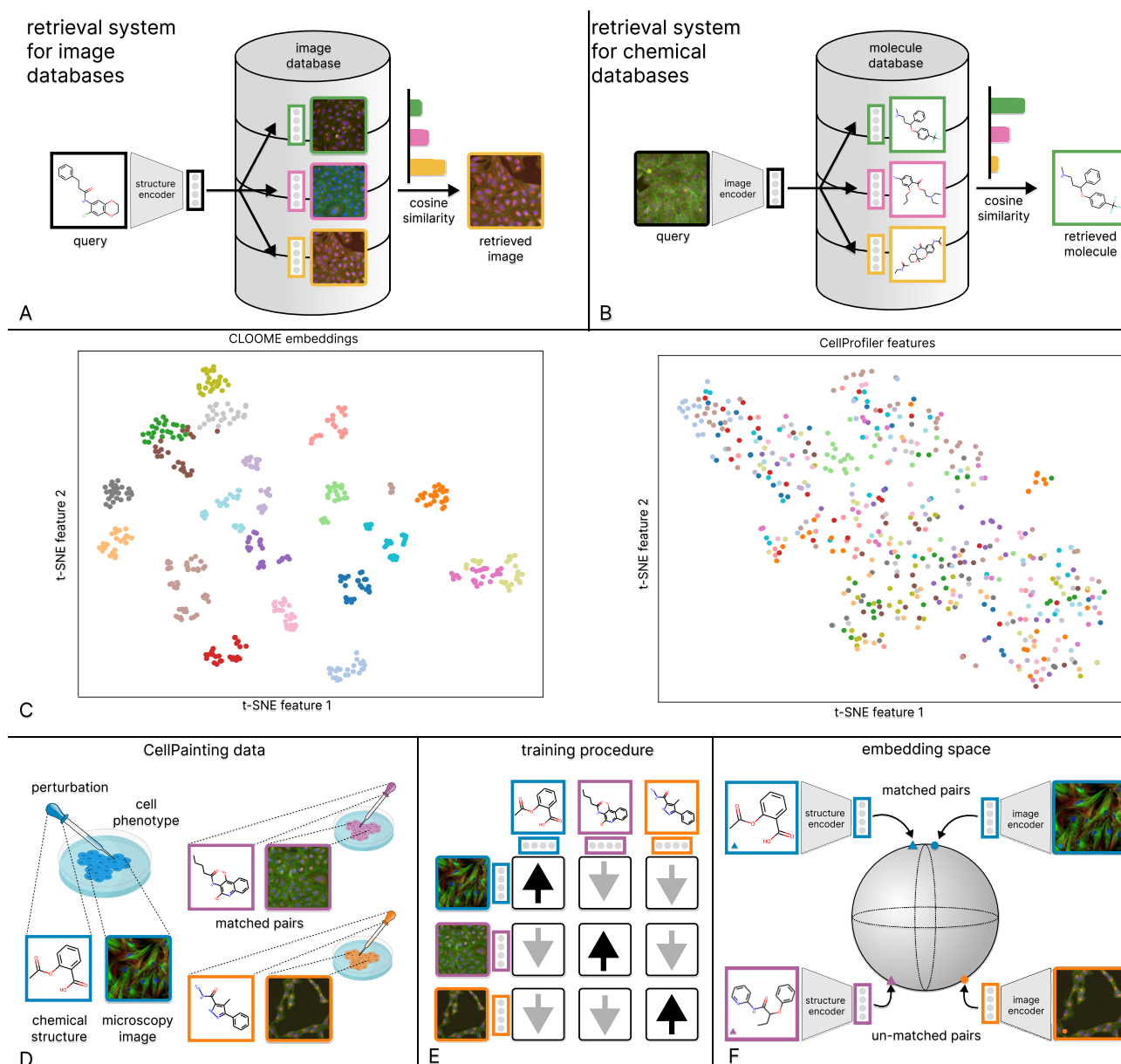


Figure 1: **A**, **B** The CLOOME encoders can be used to query a microscopy imaging database (**A**) by a chemical structure, and vice versa, query a chemical database by a microscopy image (**B**). **C** Visualization of the embedding space in terms of a t-SNE projection of image embeddings of new cell phenotypes. Each point represents a microscopy image from a hold-out set. The color indicates the cell phenotype, which was also withheld from training. The CLOOME embeddings (left) are indicative of the cell phenotype (clustered colors). CellProfiler features are less indicative of cell phenotypes (only few colors cluster together). **D** A multi-modal setting for imaging cell phenotypes. Small molecules are administered to cells which are then imaged to capture potential phenotypic changes. In this way, matched image-structure pairs are obtained. **E** Schematic depiction of the training procedure of CLOOME. During training, the similarity of matched image-structure pairs is increased (black arrows), while the similarity of un-matched image-structure pairs is decreased (grey arrows). **F** The encoders of CLOOME map chemical structures and microscopy images to the same embedding space using a structure and a microscopy image encoder. Both encoders are deep neural networks. Matched pairs of chemical structures and microscopy images are mapped to embeddings that are close together, whereas un-matched pairs are mapped to embeddings that are separated.

43 CLIP and CLOOB models have been constructed via contrastive learning on large image-text datasets [9].
44 Analogously to image-text pair datasets, the Cell Painting dataset [15] contains image-structure pairs (see
45 Figure 1D). Therefore, we were able to use contrastive learning to jointly train a microscopy image encoder
46 and a chemical structure encoder to construct a common embedding space of microscopy images capturing cell
47 phenotypes and chemical structures representing the perturbations (see Figure 1E, F). We propose a contrastive
48 learning framework for image-structure pairs that we call CLOOME (see Online Methods). The training process
49 of CLOOME would yield a) an image encoder that can map microscopy images to an informative embedding
50 space, b) a structure encoder that can map chemical structures to the same embedding space (see Figure 1F).
51 Both encoders are deep neural networks that build the basis of a search engine, which we also call CLOOME, for
52 microscopy images and chemical structures (Figure 1F). If CLOOME achieves similar results for microscopy
53 images to CLIP or CLOOB for natural images [9, 10], the image encoder should produce features, or equivalently
54 cell profiles [17], that are highly transferable and robust to distribution shifts.

55 We trained the encoders of CLOOME on 674,357 pairs of microscopy images and chemical structures of
56 the Cell Painting dataset, setting aside 28,632 for validation and 56,793 pairs for testing, ensuring that no
57 data leakage occurred between these sets (see Online Methods). We used a Residual Network [18] to encode
58 microscopy images [19], and a fully-connected neural network to encode chemical structures [20]. For each
59 training step, 256 image-structure pairs are randomly drawn from the training set and the encoders are updated
60 to increase the cosine similarities of the matched pairs and decrease the cosine similarity of un-matched pairs
61 (Figure 1E). We trained the CLOOME encoders for the retrieval system for 51 epochs, based on validation
62 performance (details in Online Methods). After the training process, we investigated CLOOME models for
63 a) the use as a retrieval systems for microscopy images and chemical structures, b) the quality of the image
64 embeddings to predict bioactivities c) the expressiveness of the image embeddings to distinguish between unseen
65 cell phenotypes.

66 a) The CLOOME encoders as retrieval system for microscopy images and chemical structures.

67 On a hold-out set of new 2,115 molecules and images, we tested whether CLOOME is able to correctly identify
68 the chemical structure with which the cells have been treated, and vice versa. This, to our knowledge, is a
69 task that is considered almost impossible for human experts. The trained CLOOME system is able to identify
70 the matched microscopy image given the chemical structure, and vice versa, with an top-10 accuracy of 8.4%
71 (95%-CI 7.3-9.7%) an 7.9% (95%-CI 6.8-9.1%), respectively. The task is extremely difficult because there are
72 many chemical structures that induce similar cell phenotypes or no phenotypic changes at all and because
73 one correct image or chemical structure has to be selected from a large set of $\sim 2,000$ candidates. Therefore,
74 the performance of human experts would likely be close to random, over which CLOOME exhibits a 70-fold
75 improvement for image retrieval and 64-fold improvement for structure retrieval over this random baseline (see
76 Table 1). This means that by just investigating the cell phenotype displayed on the microscopy image, CLOOME
77 is able to identify the matched chemical structure from a large database, and vice versa (see Table 1). Therefore,
78 the CLOOME encoders can be used as a content-based or associative retrieval system for microscopy images
79 and chemical structures ¹.

Method	Top-k accuracy (%)					
	Top-1 (rel.)	95%-CI	Top-5 (rel.)	95%-CI	Top-10 (rel.)	95%-CI
CLOOME (structure retr.)	3.03 (64.4x)	[2.34, 3.85]	6.62 (25.4x)	[5.60, 7.76]	8.42 (17.8x)	[7.27, 9.68]
CLOOME (image retr.)	3.31 (70.4x)	[2.59, 4.16]	6.24 (24.0x)	[5.24, 7.36]	7.90 (16.7x)	[6.78, 9.13]
Random	0.0473	[0.0012, 0.263]	0.236	[0.0768, 0.551]	0.473	[0.227, 0.868]

Table 1: Results for the database retrieval (retr.) task. Given a molecule-perturbed microscopy image, the matched molecule, i.e. chemical structure, must be selected from a set of $\sim 2,000$ candidate molecules (first row). Vice versa, given a chemical structure, the matched microscopy image, capturing the phenotype induced by the chemical perturbation, has to be selected from $\sim 2,000$ candidate images (second row). Top-1, top-5 and top-10 accuracy in percentage are shown for a hold-out test set, along with the upper and lower limits for a 95% confidence interval on the proportion.

80 **b) Bio-activity prediction.** Next, we investigated whether the embeddings are transferable to other tasks.
81 To this end, we used bioactivity prediction tasks as these have been approached before with cell profiling [16]
82 and convolutional neural networks (CNNs) [19]. We found that without the need for re-training or fine-tuning
83 any neural network, the CLOOME image embeddings could predict 209 activity prediction tasks with an AUC
84 of 0.714 ± 0.20 , which is on par with the best method CNNs trained in an end-to-end fashion. For CLOOME, we
85 just trained a logistic regression model, and, thus, we can conclude that the embeddings are highly transferable
86 and predictive for a diverse set of activity prediction tasks.

87 **c) Zero-shot image classification.** Lastly, we investigated how well the CLOOME image embeddings
88 characterize new cell phenotypes and whether new phenotypes could potentially be detected. We again used a

¹The search engine is available here <https://huggingface.co/spaces/anasanchezf/cloome>

89 hold-out set of microscopy images. We took the simplified assumption that each chemical structure induces a
90 separate new cell phenotype. If different structures in reality induce the same phenotype, this would even make
91 the prediction task simpler. Since multiple different samples, i.e. images, were treated with the same chemical,
92 they should produce similar embeddings or cell profiles. Therefore, a t-SNE plot, in which each point represents
93 a microscopy image embedding colored by phenotype, should show clusters of datapoints of the same color (see
94 Figure 1C left). For comparison, we also produced the same plot using the cell profiles computed with the
95 CellProfiler [17] software (see Figure 1C right). Besides the visual confirmation, the classification accuracy can
96 also be quantified (see Online Methods). Indeed, the CLOOME embeddings are similar for images capturing the
97 same, but previously unknown cell phenotype, and are thus highly expressive features of cells.

98 To conclude, we have demonstrated that self-supervised contrastive learning methods can be readily used for
99 multi-modal data arising from informative biotechnologies, such as microscopy images. Our contrastive learning
100 framework CLOOME was used to construct a common embedding space for microscopy images capturing cell
101 phenotypes and chemical structures inducing cellular processes. This enabled us to build a content-based retrieval
102 system for microscopy images and chemical structures. Furthermore, the learned microscopy image encoder has
103 been shown to produce highly transferable and expressive embeddings, or cell profiles, that can be efficiently
104 used to predict bioactivities or detect new phenotypes. We envision that our work paves the way for retrieval
105 systems for other pairs of modalities, for example querying microscopy imaging databases with transcriptomics
106 signatures or vice versa. We provide the CLOOME search engine as a free web application, the CLOOME
107 framework, the encoders and all code on github <https://github.com/ml-jku/cloome>.

108 Acknowledgements

109 ASF's research position was funded by the European Union's Horizon 2020 research and innovation program
110 under the Marie Skłodowska-Curie Innovative Training Network - European Industrial Doctorate grant agreement
111 No. 956832, "Advanced machine learning for Innovative Drug Discovery". The ELLIS Unit Linz, the LIT AI Lab,
112 the Institute for Machine Learning, are supported by the Federal State Upper Austria. IARAI is supported by
113 Here Technologies. We thank the projects AI-MOTION (LIT-2018-6-YOU-212), AI-SNN (LIT-2018-6-YOU-214),
114 DeepFlood (LIT-2019-8-YOU-213), Medical Cognitive Computing Center (MC3), INCONTROL-RL (FFG-
115 881064), PRIMAL (FFG-873979), S3AI (FFG-872172), DL for GranularFlow (FFG-871302), AIRI FG 9-N
116 (FWF-36284, FWF-36235), ELISE (H2020-ICT-2019-3 ID: 951847). We thank Audi.JKU Deep Learning Center,
117 TGW LOGISTICS GROUP GMBH, Silicon Austria Labs (SAL), FILL Gesellschaft mbH, Anyline GmbH, Google,
118 ZF Friedrichshafen AG, Robert Bosch GmbH, UCB Biopharma SRL, Merck Healthcare KGaA, Verbund AG,
119 Software Competence Center Hagenberg GmbH, TÜV Austria, Frauscher Sensonic and the NVIDIA Corporation.

120 References

- 121 1. O'Leary, N. A., Wright, M. W., Brister, J. R., Ciuffo, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse,
122 B., Smith-White, B., Ako-Adjei, D., *et al.* Reference sequence (RefSeq) database at NCBI: current status,
123 taxonomic expansion, and functional annotation. *Nucleic Acids Research* **44**, D733–D745 (2016).
- 124 2. UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Research* **47**, D506–
125 D515 (2019).
- 126 3. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool.
127 *Journal of Molecular Biology* **215**, 403–410 (1990).
- 128 4. Karolchik, D., Baertsch, R., Diekhans, M., Furey, T. S., Hinrichs, A., Lu, Y., Roskin, K. M., Schwartz, M.,
129 Sugnet, C. W., Thomas, D. J., *et al.* The UCSC genome browser database. *Nucleic Acids Research* **31**,
130 51–54 (2003).
- 131 5. Burley, S. K., Berman, H. M., Kleywegt, G. J., Markley, J. L., Nakamura, H. & Velankar, S. Protein Data
132 Bank (PDB): the single global macromolecular structure archive. *Protein Crystallography*, 627–641 (2017).
- 133 6. Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S.,
134 Michalovich, D., Al-Lazikani, B. & Overington, J. P. ChEMBL: a large-scale bioactivity database for drug
135 discovery. *Nucleic Acids Research* **40**, D1100–D1107 (2011).
- 136 7. Heller, S., McNaught, A., Stein, S., Tchekhovskoi, D. & Pletnev, I. InChI - the worldwide chemical structure
137 identifier standard. *Journal of Cheminformatics* **5**, 1–9 (2013).
- 138 8. van den Oord, A., Li, Y. & Vinyals, O. Representation Learning with Contrastive Predictive Coding. arXiv:
139 1807.03748 (2018).
- 140 9. Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin,
141 P., Clark, J., Krueger, G. & Sutskever, I. Learning Transferable Visual Models From Natural Language
142 Supervision. *International Conference on Machine Learning (ICML)* (2021).

- 143 10. Fürst, A., Rumetshofer, E., Lehner, J., Tran, V., Tang, F., Ramsauer, H., Kreil, D. P., Kopp, M., Klambauer,
144 G., Bitto-Nemling, A. & Hochreiter, S. CLOOB: Modern Hopfield Networks with InfoLOOB Outperform
145 CLIP. *Advances in Neural Information Processing Systems (NeurIPS)* (2022).
- 146 11. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C. & Chen, M. Hierarchical text-conditional image generation
147 with CLIP latents. arXiv: 2204.06125 (2022).
- 148 12. Zanella, F., Lorens, J. B. & Link, W. High content screening: seeing is believing. *Trends in biotechnology*
149 **28**, 237–245 (2010).
- 150 13. Thul, P. J., Åkesson, L., Wiking, M., Mahdessian, D., Geladaki, A., Ait Blal, H., Alm, T., Asplund, A.,
151 Björk, L., Breckels, L. M., *et al.* A subcellular map of the human proteome. *Science* **356**, eaal3321 (2017).
- 152 14. Pepperkok, R. & Ellenberg, J. High-throughput fluorescence microscopy for systems biology. *Nature Reviews*
153 *Molecular Cell Biology* **7**, 690–696 (2006).
- 154 15. Bray, M.-A., Singh, S., Han, H., Davis, C. T., Borgeson, B., Hartland, C., Kost-Alimova, M., Gustafsdottir,
155 S. M., Gibson, C. C. & Carpenter, A. E. Cell Painting, a high-content image-based assay for morphological
156 profiling using multiplexed fluorescent dyes. *Nature Protocols* **11**, 1757–1774 (2016).
- 157 16. Simm, J., Klambauer, G., Arany, A., Steijaert, M., Wegner, J. K., Gustin, E., Chupakhin, V., Chong, Y. T.,
158 Vialard, J., Buijsters, P., *et al.* Repurposing high-throughput image assays enables biological activity
159 prediction for drug discovery. *Cell Chemical Biology* **25**, 611–618 (2018).
- 160 17. Carpenter, A. E., Jones, T. R., Lamprecht, M. R., Clarke, C., Kang, I. H., Friman, O., Guertin, D. A.,
161 Chang, J. H., Lindquist, R. A., Moffat, J., *et al.* CellProfiler: image analysis software for identifying and
162 quantifying cell phenotypes. *Genome Biology* **7**, 1–11 (2006).
- 163 18. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. *IEEE Conference on*
164 *Computer Vision and Pattern Recognition (CVPR)* (2016).
- 165 19. Hofmarcher, M., Rumetshofer, E., Clevert, D.-A., Hochreiter, S. & Klambauer, G. Accurate Prediction
166 of Biological Assays with High-Throughput Microscopy Images and Convolutional Networks. *Journal of*
167 *Chemical Information and Modeling* **59**, 1163–1171 (2019).
- 168 20. Mayr, A., Klambauer, G., Unterthiner, T., Steijaert, M., Wegner, J. K., Ceulemans, H., Clevert, D.-A. &
169 Hochreiter, S. Large-scale comparison of machine learning methods for drug target prediction on ChEMBL.
170 *Chemical Science* **9**, 5441–5451 (2018).

171 A Online Methods

172 "CLOOME" is a contrastive learning framework for multi-modal microscopy imaging data. Within the CLOOME
173 framework, a microscopy image encoder and a chemical structure encoder are learned by contrasting representa-
174 tions of matched image-structure pairs against un-matched examples from other pairs. Because our framework
175 extends the contrastive learning methods CLIP [1] and CLOOB [2] to image-structure pairs, we call it Contrastive
176 Learning and leave-One-Out-boost for Molecule Encoders (CLOOME). In the following we provide details on
177 the method, data, training, assessment and evaluation. Concretely,

- 178 • we introduce a new contrastive learning approach for image- and structure-based representations of
179 molecules,
- 180 • we show that the learned representations are highly transferable to relevant downstream tasks in drug
181 discovery by linear probing on activity prediction tasks;
- 182 • we demonstrate that our approach learns rich representations of molecules which allow to retrieve potential
183 bioisosteres from image or chemical databases.

184 A.1 CLOOME: Contrastive Learning and Leave-One-Out Boost for Molecule 185 Encoders

186 We propose contrastive learning of representations from pairs of microscopy images and chemical structures
187 to obtain a common embedding space of these two modalities, and to obtain highly transferable encoders (see
188 Figure 2). In contrast to previous approaches, in which chemical structure encoders learned representations using
189 activity data [3, 4] or microscopy image encoders used hand-crafted representations [5, 6], CLOOME optimizes
190 representations without activity data or human expertise.

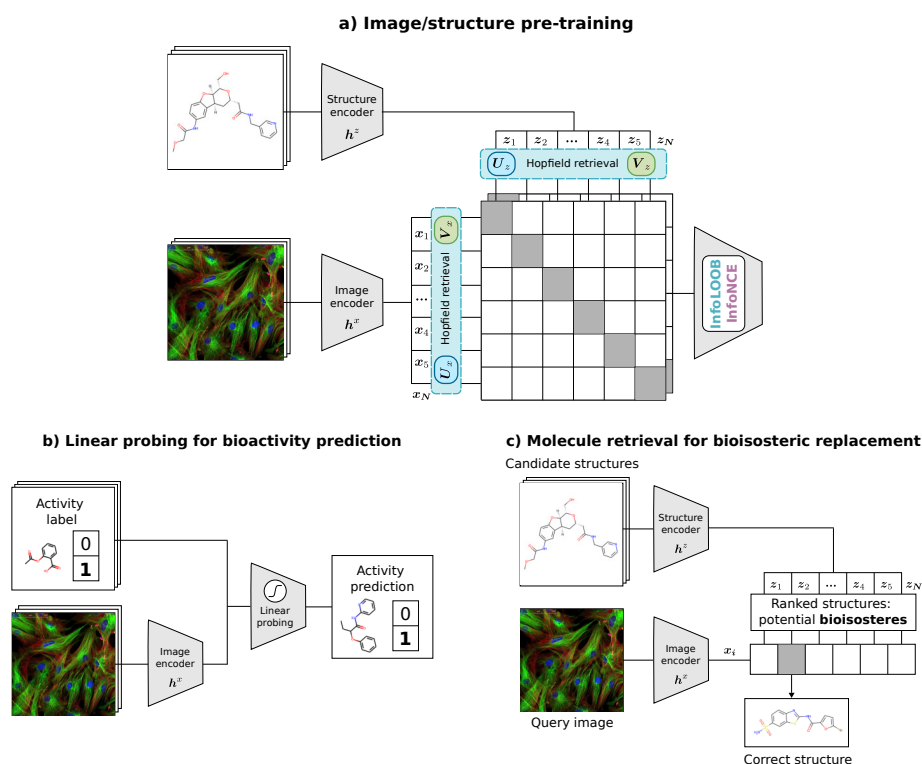


Figure 2: Schematic representation of CLOOME. Contrastive pre-training of embeddings of the two modalities, microscopy image and chemical structure, of a molecule using the CLOOB [2] approach. **b)** Using the CLOOME embeddings for activity prediction. A logistic regression model is trained for activity prediction tasks. **c)** The resulting embeddings can be used to rank chemical structures that induce similar phenotypic effects, which can be considered a bioisosteric replacement task.

191 The training dataset consists of N pairs of microscopy images of molecule-perturbed cells and chemical
192 structures of molecules $\{(x_1, z_1), \dots, (x_N, z_N)\}$. We assume that an adaptive image-encoder $h^x(\cdot)$ and an
193 adaptive structure-encoder $h^z(\cdot)$ are available that map the microscopy images and chemical structures to their

194 embeddings $\mathbf{x}_n = \mathbf{h}^x(x_n)$ and $\mathbf{z}_n = \mathbf{h}^z(z_n)$, respectively. Note that the original image is denoted as x_n , which is
195 mapped to an image embedding \mathbf{x}_n by a neural network $\mathbf{h}^x(\cdot)$, e.g. a ResNet. The stacked microscopy image
196 embeddings are denoted as $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and the stacked structure embeddings as $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_N)$. The
197 embeddings are normalized such that $\|\mathbf{x}_n\| = \|\mathbf{z}_n\| = 1 \forall n$. For notation, see also Table 6.

In a contrastive learning setting, methods aim at increasing the similarity of matched pairs and decrease the similarity of un-matched pairs. This task has often been approached by maximizing the mutual information of the embeddings using the InfoNCE loss [1, 7, 8], which is also used in the CLIP approach [1]. The InfoNCE objective function has the following form:

$$L_{\text{InfoNCE}} = -\frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{z}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{z}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{z}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{z}_i)}, \quad (1)$$

198 where τ^{-1} is the inverse temperature parameter, which is a hyperparameter of the method.

199 The contrastive learning method CLIP has the problem of "explaining away" [2, 9, 10]. Explaining away
200 describes the effect in which few features are over-represented while others are neglected. This effect can be
201 present a) when learning focuses only on few features and/or b) when the covariance structure in the data is
202 insufficiently extracted. Explaining away can be caused by saturation of the InfoNCE objective [2, 11, 12]. To
203 ameliorate these drawbacks, CLOOB [2] has introduced the InfoLOOB objective together with Hopfield networks
204 as a promising method for contrastive learning. Our contrastive learning framework CLOOME comprises both
205 methods CLIP [1] and CLOOB [2].

206 For our extension of the CLOOB method, first image- and structure-embeddings are retrieved from stored
207 image embeddings \mathbf{U} and structure embeddings \mathbf{V} . $\mathbf{U}_{\mathbf{x}_i}$ denotes an image-retrieved image embedding, $\mathbf{U}_{\mathbf{z}_i}$ a
208 structure-retrieved image embedding, $\mathbf{V}_{\mathbf{x}_i}$ an image-retrieved structure embedding and $\mathbf{V}_{\mathbf{z}_i}$ a structure-retrieved
209 structure embedding. In analogy to CLOOB, these retrievals from continuous modern Hopfield networks are
210 computed as follows:

$$\mathbf{U}_{\mathbf{x}_i} = \mathbf{U} \text{softmax}(\beta \mathbf{U}^T \mathbf{x}_i), \quad (2) \quad \mathbf{V}_{\mathbf{x}_i} = \mathbf{V} \text{softmax}(\beta \mathbf{V}^T \mathbf{x}_i), \quad (4)$$

$$\mathbf{U}_{\mathbf{z}_i} = \mathbf{U} \text{softmax}(\beta \mathbf{U}^T \mathbf{z}_i), \quad (3) \quad \mathbf{V}_{\mathbf{z}_i} = \mathbf{V} \text{softmax}(\beta \mathbf{V}^T \mathbf{z}_i), \quad (5)$$

where β is a scaling parameter of the Hopfield network which is considered a hyperparameter. These retrieved embeddings $\mathbf{U}_{\mathbf{x}_i}, \mathbf{U}_{\mathbf{z}_i}, \mathbf{V}_{\mathbf{x}_i}, \mathbf{V}_{\mathbf{z}_i}$ are also normalized to unit norm. By default, we store the current minibatch in the modern Hopfield networks, that is, $\mathbf{U} = \mathbf{X}$ and $\mathbf{V} = \mathbf{Z}$. Note that \mathbf{X} contains the image embeddings (\mathbf{Z} the structure embeddings) and we use N ambiguously both as dataset size, but also as mini-batch size to keep the notation uncluttered. The choice that $\mathbf{U} = \mathbf{X}$ and $\mathbf{V} = \mathbf{Z}$ is mostly taken because of computational constraints, while \mathbf{U} and \mathbf{V} could hold the whole dataset or, alternatively, exemplars. For further details on notation, see Table 6. Then, the InfoLOOB objective [2, 13] for the retrieved embeddings is used as objective function:

$$L_{\text{InfoLOOB}} = -\frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{U}_{\mathbf{x}_i}^T \mathbf{U}_{\mathbf{z}_i})}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{U}_{\mathbf{x}_i}^T \mathbf{U}_{\mathbf{z}_j})} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{V}_{\mathbf{x}_i}^T \mathbf{V}_{\mathbf{z}_i})}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{V}_{\mathbf{x}_j}^T \mathbf{V}_{\mathbf{z}_i})}. \quad (6)$$

211 **Microscopy image encoder.** Microscopy images differ from natural images in several aspects, for example the
212 variable number of channels that depends on the staining procedure [3, 14]. Although standard image encoders,
213 such as Residual Networks [15] could be in principle used with minor adjustments, alternative architectures,
214 such as multiple instance learning approaches, could be required for very high resolution datasets [16]. In all
215 our experiments, we use a ResNet-50 encoder with five input channels and downsized the microscopy images to
216 320x320 pixels.

217 **Molecule structure encoder.** Since the advent of Deep Learning, a large number of architectures to encode
218 molecules have been suggested [17–21]. In contrast to computer vision and natural language processing, in which
219 only few prominent architectures have emerged, there is yet no standard choice for chemical structure encoders.
220 Because of their computational efficiency and good predictive performance, CLOOME uses a descriptor-based
221 fully-connected network [22, 23] with 4 hidden layers of 1024 units with ReLU activations and batch normalization
222 (for further details see Sec. A.2 and Sec. A.6). However, also any graph [20, 24–26], message-passing [27], or
223 sequence-based [28] neural network with an appropriate pooling operation can be used as structure encoder.

224 A.2 Experiments

225 **Dataset and preprocessing.** *Cell Painting.* We use pairs of microscopy images and molecules from the
226 Cell Painting [6, 29] dataset. This dataset is a collection of high-throughput fluorescence microscopy images
227 of U2OS cells treated with different small molecules [29]. The dataset consists of 919,265 five-channel images

228 corresponding to 30,616 different molecules. The experiment to obtain the microscopy images was conducted
229 using 406 multi-well plates, and each one of the before mentioned individual images are views from a sample
230 spanning the space in the corresponding well, so that six adjacent views belong to one single sample. After
231 disregarding erratic images (out of focus or containing high fluorescence material) as well as images of untreated
232 cells that were used as controls, our final dataset comprises 759,782 microscopy images treated with 30,404
233 different molecules.

234 *Pre-processing.* We followed the pre-processing protocol of Hofmarcher *et al.* [3], which consisted of converting
235 the original TIF images from 16-bit to 8-bit, simultaneously removing the 0.0028% of pixels with highest values.
236 Moreover, the images were normalized using the mean and standard deviation calculated for the training split.
237 Concerning molecules, their corresponding SMILES strings were transformed to 1024-bit Morgan fingerprints
238 with a radius of 3, taking chirality into account [30, 31].

239 *Data splits.* We split our dataset into training, validation, and test set, using the splits of Hofmarcher *et al.*
240 [3]. Samples which have not been used in the previous study due to missing activity data, are assigned to
241 the training split. Note that all images belonging to the same molecular structure are moved into the same
242 set. Finally, training, validation and test set consist of 674,357, 28,632 and 56,793 image and molecule pairs,
243 respectively.

244 **Pre-training, architecture and hyperparameters.** We use the suggested hyperparameters of OpenCLIP
245 [32] and CLOOB [2] wherever applicable, and tuned a few critical hyperparameters, such as learning rate and
246 the β parameter of the Hopfield layer on a validation set. The architecture of the structure encoder was inspired
247 by previous successful models [23] and was not subject to substantial hyperparameter optimization. Due to
248 computational constraints, the search was limited to the hyperparameters shown in Table 7. We used the Adam
249 optimizer [33] with decoupled weight decay regularization [34]. The value for weight decay was 0.1. For the
250 learning rate scheduler, we used cosine annealing with a warm-up of 20,000 steps and hard restarts every 7
251 epochs [35]. We set the dimension of the embedding space to $d = 512$, which determines the size of the output of
252 both encoders. We use a batch size of 256 as default due to computational constraints. For activity prediction
253 as downstream task, the inverse temperature parameter $\tau^{-1} = 30$ was used. For the Hopfield layers, the scaling
254 hyperparameter $\beta = 22$ was selected, and the model was trained for 63 epochs based on linear probing results
255 in the corresponding validation set. For data augmentation and to allow large batch sizes, the images were
256 cropped and re-scaled from the original 520x696 pixel resolution to 320x320 during training. For the retrieval and
257 zero-shot image classification tasks, a higher validation performance was achieved by a CLIP-like architecture
258 directly using the embeddings returned from the image and structure encoders and the InfoNCE loss. In this
259 case, the inverse temperature parameter τ^{-1} was set to 14.3, and the model was trained for 51 epochs based
260 on the top-1 accuracy in validation. In this case, images were cropped and re-scaled to a pixel resolution of
261 520x520, based on performance in the validation set. Hence, different pre-training settings have been found
262 to yield best results for bioactivity prediction and for both the retrieval and zero-shot image classification
263 task, respectively. However, the large majority of hyperparameters were shared in both strategies. Because
264 of the limited exploration of the vast hyperparameter space, we expect potential improvements from further
265 investigations. For further details on the hyperparameter selection, see Sec. A.6.

266 **a) A retrieval system for imaging and chemical databases to enable bioisosteric replacement and**
267 **scaffold hopping.** In this experiment, we assessed the ability of CLOOME to correctly retrieve the matched
268 chemical structure given a microscopy image of cells treated with this molecule. Notably, this is an extremely
269 challenging task for human experts: given a microscopy image of cells, the task is to select the chemical structure
270 with which they have been treated from a set of thousands of candidate structures. Since cells often do not
271 exhibit any or only subtle phenotypic changes, this task is highly ambitious.

272 This image-based retrieval task can also be understood as a bioisosteric replacement task [36]: Bioisosteres
273 are molecules with roughly the same biological properties or activities, which is highly relevant in drug discovery
274 when a chemical scaffold should be replaced with another, but at the same time its biological activity should
275 be kept. With this experiment, we evaluate the ability of CLOOME to correctly rank the matched molecular
276 structure given the image. Other high-ranked structures could be potential bioisosteres, which makes this
277 experiment a proxy for the bioisosteric replacement problems (see Figure 2 b)).

278 On hold-out data of 2,115 image and molecule pairs, CLOOME ranked the matched molecule in the first
279 place for 3% of the cases. A random method would achieve a value of $1/2,115 \approx 0.047\%$, which indicates a
280 ~ 70 -fold improvement of CLOOME. For this task, different hyperparameters and model were selected based
281 on the appropriate validation metric (see Sec. A.6). The top-1, top-5, top-10 accuracy are given in Table 2 for
282 retrieving from a database of 2,115 instances. Additionally, we report the same metrics for a sampling rate of
283 1%, or equivalently, 1 matched example together with 99 un-matched ones – a setting often used to evaluate
284 retrieval systems, see Table 3. Further, some examples are displayed in Figure 3. This is, to our knowledge, the
285 first system of cell-image-based retrieval of molecular structures.

Method	Top-k accuracy (%)					
	Top-1 (rel.)	95%-CI	Top-5 (rel.)	95%-CI	Top-10 (rel.)	95%-CI
CLOOME (structure retr.)	3.03 (64.4x)	[2.34, 3.85]	6.62 (25.4x)	[5.60, 7.76]	8.42 (17.8x)	[7.27, 9.68]
CLOOME (image retr.)	3.31 (70.4x)	[2.59, 4.16]	6.24 (24.0x)	[5.24, 7.36]	7.90 (16.7x)	[6.78, 9.13]
Random	0.0473	[0.0012, 0.263]	0.236	[0.0768, 0.551]	0.473	[0.227, 0.868]

Table 2: Results for the retrieval task **among 2,115** candidates. Given a molecule-perturbed microscopy image, the matched molecule must be selected from a set of candidates, and vice versa. Top-1, top-5 and top-10 accuracy in percentage are shown for a hold-out test set, along with the upper and lower limits for a 95% confidence interval on the proportion.

Method	Top-k accuracy (%)					
	Top-1	95% CI	Top-5	95% CI	Top-10	95% CI
CLOOME (structure retrieval)	10.4	[9.10, 11.7]	21.3	[19.6, 23.1]	30.6	[28.7, 32.7]
CLOOME (image retrieval)	9.64	[8.42, 11.0]	20.7	[19.0, 22.4]	29.0	[27.1, 31.0]
Random	0.992	[0.616, 1.51]	5.01	[4.12, 6.03]	10.0	[8.78, 11.4]

Table 3: Results for the retrieval task **among 100** candidates. Given a molecule-perturbed microscopy image, the matched molecule must be selected from a set of candidates, and vice versa. Top-1, top-5 and top-10 accuracy in percentage are shown for a hold-out test set, along with the upper and lower limits for a 95% confidence interval on the proportion.

b) Bio-activity prediction as downstream tasks. In this experiment, we tested whether the representations learned by CLOOME are transferable by linear probing on 209 downstream activity prediction tasks. The *linear probing* test [8, 37] on downstream tasks is often performed for contrastive learning approaches to check the transferability of learned features. In such experiments, the representations of the pretrained encoders are used, and only a single-layer network, such as logistic regression, is fit to the given labels for the supervised task. If the linear probing test yields good predictive quality, usually below a fully supervised approach [8], the representations are considered transferable.

Linear probing evaluation. The prediction tasks that we employed for linear probing evaluation is the same as used in Hofmarcher *et al.* [3]. It is a subset of the Cell Painting dataset, consisting of 284,035 images for which the activity labels of the compound treatments were retrieved from ChEMBL. The retrieved labels correspond to 10,574 compounds across 209 activity prediction tasks, which are binary classification problems. However, activity data points are not available for all compounds in all of the tasks, which results in a sparse label matrix. The data was split into 70% training, 10% validation, and 20% test sets. This split had been carried out by grouping views from samples treated with the same molecule.

We use image features taken from the penultimate layer of the image encoder, omitting the classification layer. We train a logistic regression classifier, and report the corresponding metric for each task. The L2 regularization strength λ was tuned individually for each one of the tasks, considering the values $\{10^{-6}, 10^{-5}, \dots, 10^6\}$.

In order to evaluate model performance for this downstream task, we use the area under the ROC curve (AUC), which is one of the most prevalent metrics for drug discovery [3, 4], as it considers the order of the molecules regarding their activity. We also show the number of tasks for which this metric is higher than the thresholds 0.9, 0.8 and 0.7, respectively. These thresholds have been used in previous studies [3, 4] because models within those categories lead to certain levels of enrichment of hit rates in drug discovery projects.

Baselines. As baselines we consider methods reported in Hofmarcher *et al.* [3]. They are the best performing methods for bioactivity prediction using microscopy images to date and consist of different convolutional neural network architectures, used in a fully supervised setting, and a method ("FNN") that uses expert-designed cell features [4–6]. The compared methods were trained in a multi-task setting to predict activity labels for 209 tasks, extracted from ChEMBL.

Results. The predictive performance on the downstream activity prediction tasks is reported in Table 4. CLOOME reached an average AUC of 0.714 ± 0.20 across prediction tasks, which indicates that the learned representations are indeed transferable since no activity data had been used to train the CLOOME encoders. CLOOME even outperformed fully supervised methods, such as M-CNN [38] and SC-CNN [3], with respect to AUC.

c) Zero-shot microscopy image classification. The goal of this analysis is to evaluate the potential of the CLOOME image embeddings to distinguish between cell phenotypes. Classifying the phenotype captured by the microscopy image is a highly relevant biological question. Especially for drug discovery, where phenotypes are induced by chemical perturbations, embeddings that can identify novel phenotypes would provide some understanding about its possible mode of action and therefore its potential adverse effects.

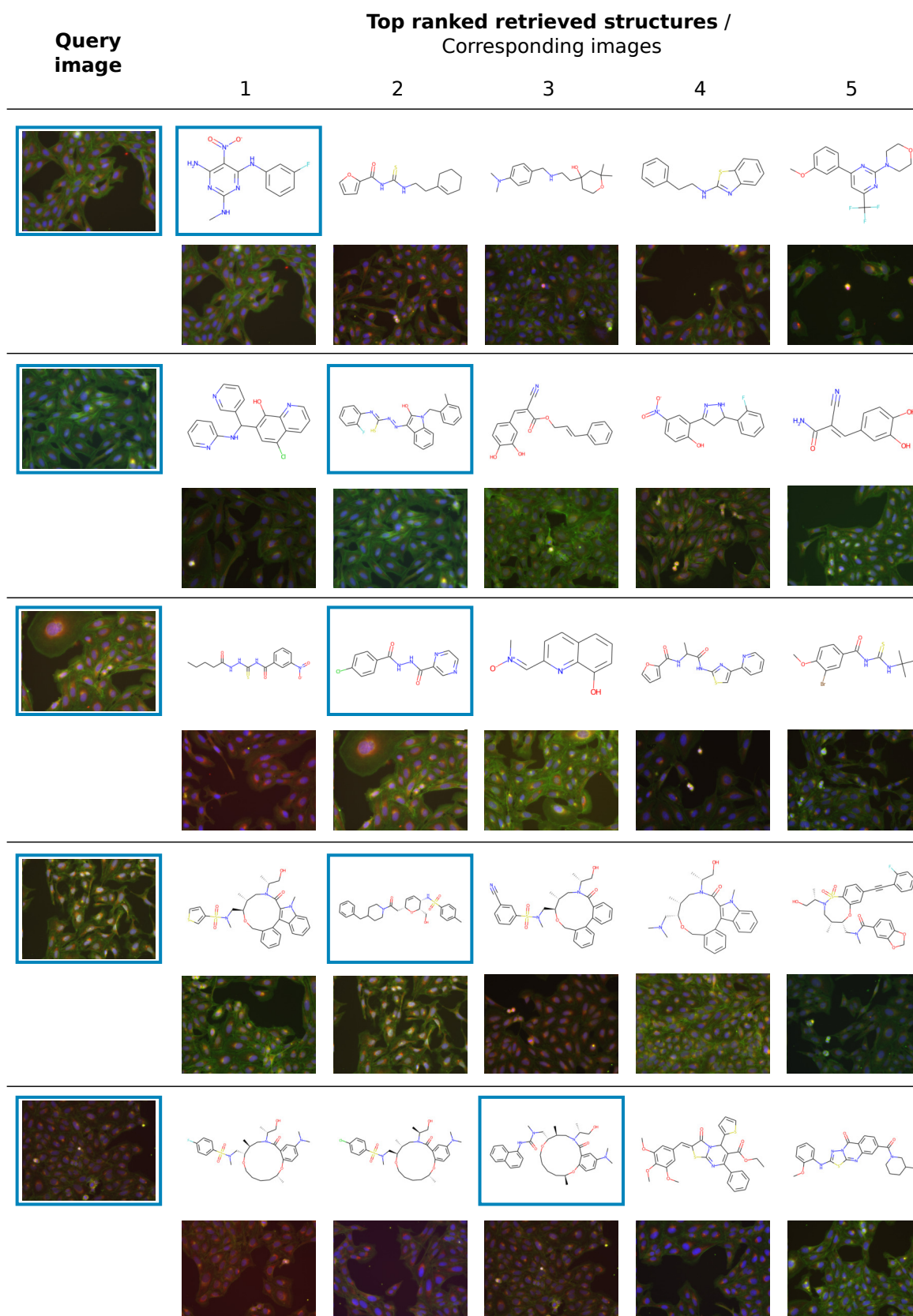


Figure 3: Example results for the retrieval task. On a hold-out test set, the five molecules for which representations are the most similar to the query image are shown along with their corresponding images. Blue boxes mark the query image and its matching molecular structure, i.e. the matching pair. CLOOME can be used to retrieve molecules that could produce similar biological effects on treated cells, i.e. bioisosteres.

Type	Method	AUC	F1	AUC >0.9	AUC >0.8	AUC >0.7
Linear probing on self-supervised	CLOOME	0.714±0.20	0.395±0.32	57	84	109
	ResNet	0.731±0.19	0.508±0.30	68	94	119
Supervised	DenseNet	0.730±0.19	0.530±0.30	61	98	121
	GapNet	0.725±0.19	0.510±0.29	63	94	117
	MIL-Net	0.711±0.18	0.445±0.32	61	81	105
	M-CNN	0.705±0.19	0.482±0.31	57	78	105
	SC-CNN	0.705±0.20	0.362±0.29	61	83	109
	FNN	0.675±0.20	0.361±0.31	55	71	90

Table 4: Comparison of the linear probing evaluation of the learned representations against fully supervised methods [3]. Note that the CLOOME encoders do not have access to any activity data. The features produced by the CLOOME encoder are still predictive for activity data as shown by fitting a logistic regression model, considered as linear probing. CLOOME reaches the performance of the several supervised methods, which indicates transferability of the learned representations [8]. The best method in each category is marked bold.

323 In the same hold-out test used for this zero-shot image classification task, each molecule is assumed to cause
 324 a different phenotype. While this assumption is not true, and distinct chemical structures can induce the same
 325 phenotype, this is the most difficult setting and the more realistic setting would make the classification task even
 326 easier. To provide details on this classification task, one image for each of the molecules was randomly selected,
 327 resulting in 2,115 classes. Then, samples corresponding to both the same molecule and plate as those from the
 328 class set were removed in order to ensure that the classification was not due to plate effects. This yielded a
 329 44,102 image test set.

330 We compared the CLOOME embeddings, to embeddings of a microscopy image encoder trained in supervised
 331 fashion and to profiles from CellProfiler. Regarding GapNet embeddings, the images were encoded using the
 332 model weights provided by Hofmarcher *et al.* [3], removing the last layer of its classifier, which resulted in a
 333 1024-dimension embedding space. CellProfiler embeddings consist in 148 features aggregated in one vector per
 334 image, as made available in Bray *et al.* [6].

Method	Accuracy[%]					
	Top-1	95% CI	Top-5	95% CI	Top-10	95% CI
CLOOME	17.8	[17.4, 18.2]	40.6	[40.2, 41.1]	55.3	[54.8, 55.8]
GapNet (CNN)	0.363	[0.309, 0.423]	1.07	[0.981, 1.18]	1.80	[1.67, 1.92]
CellProfiler	0.497	[0.433, 0.567]	1.75	[1.63, 1.88]	2.89	[2.74, 3.05]

Table 5: Results for the zero-shot microscopy image classification. Given a molecule-perturbed microscopy image, the image corresponding to the matched molecule must be selected from a set of candidates. Top-1, top-5 and top-10 accuracy in percentage are shown for a hold-out test set, along with the upper and lower limits for a 95% confidence interval on the resulting proportion.

335 A.3 Related work.

336 **Contrastive learning has had a strong impact on computer vision and natural language processing.**
 337 Over the last decade, supervised deep learning methods have achieved outstanding results in the field of computer
 338 vision [15, 39]. These supervised methods require large amounts of labeled data, which may be very costly or
 339 unfeasible to obtain, and they have limited generalization abilities [40, 41]. This has led to the exploration of new
 340 methods that are able to learn robust representations of the data which can be transferred to different downstream
 341 tasks [8, 42]. With contrastive learning methods [43] and self-supervision these meaningful representations can be
 342 obtained without the need for large amounts of expensive manually-provided labels [8, 44–46]. While uni-modal
 343 methods typically use pre-text tasks [8], for multi-modal methods the self-supervision arises from the availability
 344 of two modalities of an instance, such as image and text [1, 47]. Both uni-modal and multi-modal contrastive
 345 learning methods have recently had a substantial impact in computer vision and natural language processing
 346 [48].

347 **CLIP for multi-modal data yields spectacular performance at zero-shot transfer learning and has**
 348 **recently been improved by CLOOB.** An outstanding multi-modal approach is Contrastive Language-

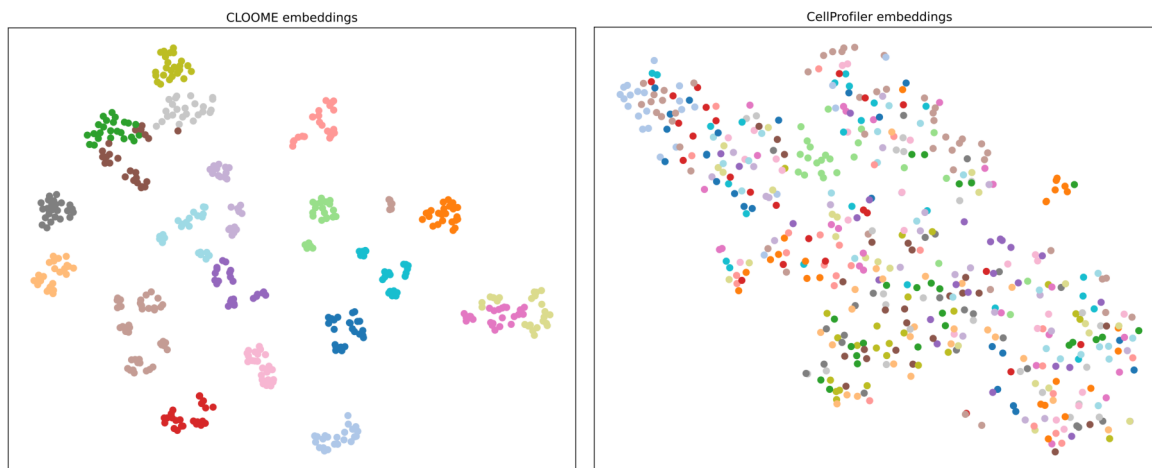


Figure 4: t-SNE downprojection of the CLOOME learned embeddings (left) and the CellProfiler extracted features (right) of all microscopy images corresponding to 20 molecules randomly selected from the test set. The colors represent different molecules.

349 Image Pre-training (CLIP) [1], which learns both image- and text-representations simultaneously. CLIP shows
350 comparable performance to methods that are solely image-based and yields highly transferable representations,
351 which is shown by its high performance at zero-shot transfer learning. However, CLIP has recently been shown
352 to suffer from the "explaining away" effect [2, 9, 10] (details in Section A.1). Considering this caveat, the
353 "Contrastive Leave One Out Boost" (CLOOB) method has been proposed [2]. CLOOB uses a different objective,
354 the "InfoLOOB" (LOOB for "Leave One Out Bound") objective [13], which does not include the positive pair
355 in the denominator to avoid saturation effects [2]. Moreover, continuous modern Hopfield networks [49] are
356 used to reinforce the covariance structure of the data. As a result, CLOOB has further improved zero-shot
357 transfer learning. The ability to learn transferable representation from multi-modal data makes CLOOB the
358 prime candidate for learning representations of molecules in drug discovery.

359 **Contrastive learning for molecule representations in drug discovery.** In drug discovery, the effect of
360 the limited availability of data on molecules is even more severe, since the acquisition of a single bioactivity
361 data point can cost several thousand dollars and take several weeks or months [50, 51]. Therefore, methods
362 that can learn transferable representations from unlabelled data are highly demanded. Thus, several contrastive
363 learning approaches have been recently developed for different tasks in drug discovery. MolCLR [52] uses
364 contrastive molecule-to-molecule training by augmenting molecular graphs. Stärk *et al.* [53] contrastively learn
365 3D and 2D molecule representations to inform the learned molecule encoder with 3D information. Lee *et al.* [54]
366 and Seidl *et al.* [55] use contrastive learning for molecules and chemical reactions, and Vall *et al.* [56] utilizes
367 text representations of wet-lab procedures to enable zero-shot predictions. However, none of these methods
368 have exploited the wealth of information contained in microscopy images of molecule-perturbed cells [29] and
369 demonstrated strong transferability of the learned molecule encoders.

370 **Image-based profiling of small molecules has strongly improved the drug discovery process.**
371 Characterizing a small molecule by the phenotypic changes it induces to a cell, is considered promising for
372 accelerating drug discovery [4, 29, 57, 58]. The advantages of this biotechnology are that it is time- and
373 cost-effective as compared to standard activity measurements. Measuring the effects of a molecule on a biological
374 system early in the drug discovery process might be useful to improve clinical success rates [59]. Particularly,
375 microscopy image-based profiles of small molecules have been suggested to be effective together with deep learning
376 methods [58]. However, the current efforts are still in standard supervised learning settings based on extracted
377 features [4] or deep architectures [3]. The amount of labeled images is in the range of few tens of thousands,
378 although international efforts are currently building datasets which are magnitudes larger [60]. Instead of the
379 currently used activity measurements as labels [3, 4], we propose a self-supervised contrastive learning strategy
380 of image- and structure-based molecule encoders: Contrastive Leave One Out boost for Molecule Encoders
381 (CLOOME). CLOOME extends recent successful contrastive learning methods to the fields of biological imaging
382 and drug discovery. Our approach intends to overcome the limited transferability of current molecule encoders
383 [61, 62].

384 A.4 Discussion and conclusion

385 We have introduced a contrastive learning method for learning representations of microscopy images and chemical
 386 structures. On the largest available dataset of this type, we demonstrate that CLOOME is able to learn
 387 transferable representations. This opens the possibility to re-use the learned representations for activity or
 388 property prediction and for other tasks, such as retrieval tasks from microscopy image or chemical databases.

389 *Limitations.* Our method currently has several limitations. Our trained networks are restricted to a particular
 390 type of microscopy images, which are acquired with the Cell Painting protocol [29]. This protocol has been
 391 published and currently there are community efforts [60] to increase the amount of available data. Large and
 392 more diverse datasets of molecule-perturbed cells or internal pharmaceutical company datasets will likely improve
 393 the learned representations, both image and structure encoder [63]. Due to the computational complexity, the
 394 hyperparameter and architecture space is currently under-explored such that we expect our method to further
 395 improve with better hyperparameters or encoder architectures. Furthermore, it has not escaped our notice
 396 that the learned structure encoder can also be used for transfer learning on molecular activities and properties.
 397 Also, it is worth noting that, although linear probing has been extensively used for the purpose of evaluating
 398 the quality of representations [1, 2], if the latter are very high dimensional, this method presents the risk of
 399 overfitting [37]. Having addressed these limitations, we nevertheless believe that the representations obtained
 400 with CLOOME could be highly useful for both the community using bioimaging as well as for drug discovery.

401 A.5 Notation overview

Definition	Symbol/Notation	Dimension
molecule-perturbed microscopy image	x	image dimension, e.g. $320 \times 320 \times 5$
chemical structure of molecule	z	symbolic, e.g. graph
image embedding	\mathbf{x}	d
structure embedding	\mathbf{z}	d
stacked image embeddings	\mathbf{X}	$d \times N$
stacked structure embeddings	\mathbf{Z}	$d \times N$
stored image embeddings	\mathbf{U}	$d \times N$
stored structure embeddings	\mathbf{V}	$d \times N$
image-retrieved image embedding	\mathbf{U}_{x_i}	d
structure-retrieved image embedding	\mathbf{U}_{z_i}	d
image-retrieved structure embedding	\mathbf{V}_{x_i}	d
structure-retrieved structure embedding	\mathbf{V}_{z_i}	d
microscopy image encoder	$\mathbf{h}^x(\cdot)$	$\mathbb{R}^{320 \times 320 \times 5} \rightarrow d$
molecule structure encoder	$\mathbf{h}^z(\cdot)$	$\mathcal{M} \rightarrow d$
temperature parameter of the loss functions	τ	
scaling parameter of Hopfield net	β	
embedding dimension	d	
batch or dataset size	N	
chemical space	\mathcal{M}	
indices	i, j, n	

Table 6: Symbols and notations used in this paper.

402 A.6 Hyperparameter search space

	Hyperparameter	Explored space
Learning	Optimizer	{AdamW}
	Learning rate	{0.0005, 0.001 , 0.005}
	Scheduler	{Cosine annealing with restarts}
	Weight decay	{0.1}
	Batch size	{ 256 , 512}
	Warm-up iterations	{10000, 20000 }
	Inverse temperature	{ 30 }
Image encoder	Image resolution	{ 320 , 520}
	Model	{ResNet50}
Structure encoder	Number of layers	{4}
	Layer dimension	{1024}
	Activation	{ReLU}
	Batch normalization	{False, True }
Hopfield layers	β	{8, 14.3, 22 }
Embedding space	Number of dimensions	{ 512 }

Table 7: Considered hyperparameter space of CLOOME models. The selected configurations for downstream activity prediction based on manual search on validation set shown in bold.

403 References

- 404 1. Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin,
405 P., Clark, J., Krueger, G. & Sutskever, I. Learning Transferable Visual Models From Natural Language
406 Supervision. *International Conference on Machine Learning (ICML)* (2021).
- 407 2. Fürst, A., Rumetshofer, E., Lehner, J., Tran, V., Tang, F., Ramsauer, H., Kreil, D. P., Kopp, M., Klambauer,
408 G., Bitto-Nemling, A. & Hochreiter, S. CLOOB: Modern Hopfield Networks with InfoLOOB Outperform
409 CLIP. *Advances in Neural Information Processing Systems (NeurIPS)* (2022).
- 410 3. Hofmarcher, M., Rumetshofer, E., Clevert, D.-A., Hochreiter, S. & Klambauer, G. Accurate Prediction
411 of Biological Assays with High-Throughput Microscopy Images and Convolutional Networks. *Journal of*
412 *Chemical Information and Modeling* **59**, 1163–1171 (2019).
- 413 4. Simm, J., Klambauer, G., Arany, A., Steijaert, M., Wegner, J. K., Gustin, E., Chupakhin, V., Chong, Y. T.,
414 Vialard, J., Buijsters, P., *et al.* Repurposing high-throughput image assays enables biological activity
415 prediction for drug discovery. *Cell Chemical Biology* **25**, 611–618 (2018).
- 416 5. Carpenter, A. E., Jones, T. R., Lamprecht, M. R., Clarke, C., Kang, I. H., Friman, O., Guertin, D. A.,
417 Chang, J. H., Lindquist, R. A., Moffat, J., *et al.* CellProfiler: image analysis software for identifying and
418 quantifying cell phenotypes. *Genome Biology* **7**, 1–11 (2006).
- 419 6. Bray, M.-A., Gustafsdottir, S. M., Rohban, M. H., Singh, S., Ljosa, V., Sokolnicki, K. L., Bittker, J. A.,
420 Bodycombe, N. E., Dancik, V., Hasaka, T. P., *et al.* A dataset of images and morphological profiles of 30
421 000 small-molecule treatments using the Cell Painting assay. *Gigascience* **6**, 1–5 (2017).
- 422 7. van den Oord, A., Li, Y. & Vinyals, O. Representation Learning with Contrastive Predictive Coding. arXiv:
423 1807.03748 (2018).
- 424 8. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A Simple Framework for Contrastive Learning of Visual
425 Representations. *International Conference on Machine Learning (ICML). Proceedings of Machine Learning*
426 *Research (PMLR)* 1597–1607 (2020).
- 427 9. Pearl, J. Embracing causality in default reasoning. *Artificial Intelligence* **35**, 259–271 (1988).
- 428 10. Wellman, M. P. & Henrion, M. Explaining 'Explaining Away'. *IEEE Transactions on Pattern Analysis and*
429 *Machine Intelligence* **15**, 287–292 (1993).
- 430 11. Yeh, C.-H., Hong, C.-Y., Hsu, Y.-C., Liu, T.-L., Chen, Y. & LeCun, Y. Decoupled Contrastive Learning.
431 arXiv: 2110.06848 (2021).
- 432 12. Zhang, C., Zhang, K., Pham, T. X., Niu, A., Qiao, Z., Yoo, C. D. & Kweon, I. S. Dual Temperature Helps
433 Contrastive Learning Without Many Negative Samples: Towards Understanding and Simplifying MoCo.
434 arXiv: 2203.17248 (2022).
- 435 13. Poole, B., Ozair, S., van den Oord, A., Alemi, A. A. & Tucker, G. On Variational Bounds of Mutual
436 Information. *Proceedings of Machine Learning Research (PMLR)* **97**, 5171–5180 (2019).
- 437 14. Pepperkok, R. & Ellenberg, J. High-throughput fluorescence microscopy for systems biology. *Nature Reviews*
438 *Molecular Cell Biology* **7**, 690–696 (2006).
- 439 15. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. *IEEE Conference on*
440 *Computer Vision and Pattern Recognition (CVPR)* (2016).
- 441 16. Ilse, M., Tomczak, J. M. & Welling, M. in *Handbook of Medical Image Computing and Computer Assisted*
442 *Intervention* 521–546 (Elsevier, 2020).
- 443 17. Lusci, A., Pollastri, G. & Baldi, P. Deep architectures and deep learning in chemoinformatics: the prediction
444 of aqueous solubility for drug-like molecules. *Journal of Chemical Information and Modeling* **53**, 1563–1575
445 (2013).
- 446 18. Dahl, G. E., Jaitly, N. & Salakhutdinov, R. Multi-task Neural Networks for QSAR Predictions. *ArXiv.*
447 *eprint: 1406.1231* (2014).
- 448 19. Unterthiner, T., Mayr, A., Klambauer, G., Steijaert, M., Ceulemans, H., Wegner, J. K. & Hochreiter, S.
449 Deep Learning as an Opportunity in Virtual Screening. *Advances in Neural Information Processing Systems*
450 *(NeurIPS), Workshop on Deep Learning and Representation Learning* (2014).
- 451 20. Kearnes, S., McCloskey, K., Berndl, M., Pande, V. & Riley, P. Molecular graph convolutions: moving
452 beyond fingerprints. *Journal of Computer-Aided Molecular Design* **30**, 595–608 (2016).
- 453 21. Jiang, D., Wu, Z., Hsieh, C.-Y., Chen, G., Liao, B., Wang, Z., Shen, C., Cao, D., Wu, J. & Hou, T. Could
454 graph neural networks learn better molecular representation for drug discovery? A comparison study of
455 descriptor-based and graph-based models. *Journal of Cheminformatics* **13**, 1–23 (2021).

- 456 22. Mayr, A., Klambauer, G., Unterthiner, T. & Hochreiter, S. DeepTox: toxicity prediction using deep learning.
457 *Frontiers in Environmental Science* **3**, 80 (2016).
- 458 23. Mayr, A., Klambauer, G., Unterthiner, T., Steijaert, M., Wegner, J. K., Ceulemans, H., Clevert, D.-A. &
459 Hochreiter, S. Large-scale comparison of machine learning methods for drug target prediction on ChEMBL.
460 *Chemical Science* **9**, 5441–5451 (2018).
- 461 24. Merkwirth, C. & Lengauer, T. Automatic generation of complementary descriptors with molecular graph
462 networks. *Journal of Chemical Information and Modeling* **45**, 1159–1168 (2005).
- 463 25. Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M. & Monfardini, G. The graph neural network model.
464 *IEEE Transactions on Neural Networks* **20**, 61–80 (2008).
- 465 26. Xu, K., Hu, W., Leskovec, J. & Jegelka, S. How Powerful are Graph Neural Networks? *International
466 Conference on Learning Representations (ICLR)* (2018).
- 467 27. Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O. & Dahl, G. E. Neural message passing for quantum
468 chemistry. *International Conference on Machine Learning (ICML)*, 1263–1272 (2017).
- 469 28. Alperstein, Z., Cherkasov, A. & Rolfe, J. T. All SMILES variational autoencoder. arXiv: 1905.13343
470 (2019).
- 471 29. Bray, M.-A., Singh, S., Han, H., Davis, C. T., Borgeson, B., Hartland, C., Kost-Alimova, M., Gustafsdottir,
472 S. M., Gibson, C. C. & Carpenter, A. E. Cell Painting, a high-content image-based assay for morphological
473 profiling using multiplexed fluorescent dyes. *Nature Protocols* **11**, 1757–1774 (2016).
- 474 30. Morgan, H. L. The Generation of a Unique Machine Description for Chemical Structures - A Technique
475 Developed at Chemical Abstracts Service. *Journal of Chemical Documentation* **5**, 107–113 (1964).
- 476 31. Rogers, D. & Hahn, M. Extended-Connectivity Fingerprints. *Journal of Chemical Information and Modeling*
477 **50**, 742–754 (2010).
- 478 32. Wortsman, M., Ilharco, G., Li, M., Kim, J. W., Hajishirzi, H., Farhadi, A., Namkoong, H. & Schmidt, L.
479 Robust fine-tuning of zero-shot models. arXiv: 2109.01903 (2021).
- 480 33. Kingma, D. P., Mohamed, S., Rezende, D. . & Welling, M. Semi-supervised Learning with Deep Generative
481 Models. *Advances in Neural Information Processing Systems (NeurIPS)*, 3581–3589 (2014).
- 482 34. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. *International Conference on Learning
483 Representations (ICLR)* (2019).
- 484 35. Loshchilov, I. & Hutter, F. SGDR: Stochastic Gradient Descent with Warm Restarts. *International
485 Conference on Learning Representations (ICLR)* (2017).
- 486 36. Lipinski, C. A. Bioisosterism in drug design. *Annual Reports in Medicinal Chemistry* **21**, 283–291 (1986).
- 487 37. Alain, G. & Bengio, Y. Understanding intermediate layers using linear classifier probes. arXiv: 1610.01644
488 (2016).
- 489 38. Godinez, W. J., Hossain, I., Lazic, S. E., Davies, J. W. & Zhang, X. A multi-scale convolutional neural
490 network for phenotyping high-content cellular images. *Bioinformatics* **33**, 2010–2019 (2017).
- 491 39. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet Classification with Deep Convolutional Neural
492 Networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 1097–1105 (2012).
- 493 40. Sun, C., Shrivastava, A., Singh, S. & Gupta, A. Revisiting Unreasonable Effectiveness of Data in Deep
494 Learning Era. arXiv: 1707.02968v2 (2017).
- 495 41. Marcus, G. Deep learning: A critical appraisal. arXiv: 1801.00631 (2018).
- 496 42. Luo, Z., Zou, Y., Hoffman, J. & Fei-Fei, L. F. Label efficient learning of transferable representations across
497 domains and tasks. *Advances in Neural Information Processing Systems (NeurIPS)* (2017).
- 498 43. Gutmann, M. & Hyvärinen, A. Noise-contrastive estimation: A new estimation principle for unnormalized
499 statistical models. *International Conference on Artificial Intelligence and Statistics. Proceedings of Machine
500 Learning Research* **9**, 297–304 (2010).
- 501 44. He, K., Fan, H., Wu, Y., Xie, S. & Girshick, R. B. Momentum Contrast for Unsupervised Visual Rep-
502 resentation Learning. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*
503 (2020).
- 504 45. Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P. & Joulin, A. Unsupervised Learning of Vi-
505 sual Features by Contrasting Cluster Assignments. *Advances in Neural Information Processing Systems
506 (NeurIPS)*, 9912–9924 (2020).

- 507 46. Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P. H., Buchatskaya, E., Doersch, C., Pires, B. Á.,
508 Guo, Z. D., Azar, M. G., Piot, B., Kavukcuoglu, K., Munos, R. & Valko, M. Bootstrap Your Own Latent - A
509 New Approach to Self-Supervised Learning. *Advances in Neural Information Processing Systems (NeurIPS)*
510 **33**, 21271–21284 (2020).
- 511 47. Devillers, B., Bielawski, R., Choski, B. & VanRullen, R. Does language help generalization in vision models?
512 arXiv: 2104.08313 (2021).
- 513 48. Jaiswal, A., Babu, A. R., Zadeh, M. Z., Banerjee, D. & Makedon, F. A survey on contrastive self-supervised
514 learning. *Technologies* **9**, 2 (2020).
- 515 49. Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlović, M.,
516 Sandve, G. K., Greiff, V., Kreil, D., Kopp, M., Klambauer, G., Brandstetter, J. & Hochreiter, S. *Hopfield*
517 *networks is all you need* in *International Conference on Learning Representations (ICLR)* (2021).
- 518 50. MacArron, R., Banks, M. N., Bojanic, D., Burns, D. J., Cirovic, D. A., Garyantes, T., Green, D. V.,
519 Hertzberg, R. P., Janzen, W. P., Paslay, J. W., Schopfer, U. & Sittampalam, G. S. Impact of high-throughput
520 screening in biomedical research. *Nature Reviews Drug Discovery* **10**, 188–195. ISSN: 14741776 (2011).
- 521 51. Knight, A., Bailey, J. & Balcombe, J. Animal carcinogenicity studies: 3. Alternatives to the bioassay.
522 *Alternatives to Laboratory Animals* **34**, 39–48 (2006).
- 523 52. Wang, Y., Wang, J., Cao, Z. & Barati Farimani, A. Molecular contrastive learning of representations via
524 graph neural networks. *Nature Machine Intelligence*, 1–9 (2022).
- 525 53. Stärk, H., Beaini, D., Corso, G., Tossou, P., Dallago, C., Günnemann, S. & Liò, P. 3D Infomax improves
526 GNNs for Molecular Property Prediction. arXiv: 2110.04126 (2021).
- 527 54. Lee, H., Ahn, S., Seo, S.-W., Song, Y. Y., Yang, E., Hwang, S.-J. & Shin, J. RetCL: A Selection-based
528 Approach for Retrosynthesis via Contrastive Learning. arXiv: 2105.00795 (2021).
- 529 55. Seidl, P., Renz, P., Dyubankova, N., Neves, P., Verhoeven, J., Wegner, J. K., Segler, M., Hochreiter, S.
530 & Klambauer, G. Improving Few-and Zero-Shot Reaction Template Prediction Using Modern Hopfield
531 Networks. *Journal of Chemical Information and Modeling* (2022).
- 532 56. Vall, A., Hochreiter, S. & Klambauer, G. BioassayCLR: Prediction of biological activity for novel bioassays
533 based on rich textual descriptions. *ELLIS ML4Molecules workshop* (2021).
- 534 57. Caicedo, J. C., Cooper, S., Heigwer, F., Warchal, S., Qiu, P., Molnar, C., Vasilevich, A. S., Barry, J. D.,
535 Bansal, H. S., Kraus, O., *et al.* Data-analysis strategies for image-based cell profiling. *Nature Methods* **14**,
536 849–863 (2017).
- 537 58. Chandrasekaran, S. N., Ceulemans, H., Boyd, J. D. & Carpenter, A. E. Image-based profiling for drug
538 discovery: due for a machine-learning upgrade? *Nature Reviews Drug Discovery* **20**, 145–159 (2021).
- 539 59. Bender, A. & Cortés-Ciriano, I. Artificial intelligence in drug discovery: what is realistic, what are illusions?
540 Part 1: ways to make an impact, and why we are not there yet. *Drug Discovery Today* **26**, 511–524 (2021).
- 541 60. Chandrasekaran, S. N., Cimini, B. A., Goodale, A., Miller, L., Kost-Alimova, M., Jamali, N., Doench,
542 J., Fritchman, B., Skepner, A., Melanson, M., Arevalo, J., Caicedo, J. C., Kuhn, D., Hernandez, D.,
543 Berstler, J., Shafqat-Abbasi, H., Root, D., Swalley, S., Singh, S. & Carpenter, A. E. Three million images
544 and morphological profiles of cells treated with matched chemical and genetic perturbations. bioRxiv:
545 10.1101/2022.01.05.475090 (2022).
- 546 61. Cai, C., Wang, S., Xu, Y., Zhang, W., Tang, K., Ouyang, Q., Lai, L. & Pei, J. Transfer learning for drug
547 discovery. *Journal of Medicinal Chemistry* **63**, 8683–8694 (2020).
- 548 62. Stanley, M., Bronskill, J. F., Maziarz, K., Misztela, H., Lanini, J., Segler, M., Schneider, N. & Brockschmidt,
549 M. FS-Mol: A few-shot learning dataset of molecules. *Conference on Neural Information Processing Systems*
550 *(NeurIPS), Datasets and Benchmarks Track (Round 2)* (2021).
- 551 63. Sturm, N., Mayr, A., Le Van, T., Chupakhin, V., Ceulemans, H., Wegner, J., Golib-Dzib, J.-F., Jeliaskova,
552 N., Vandriessche, Y., Böhm, S., *et al.* Industry-scale application and evaluation of deep learning for drug
553 target prediction. *Journal of Cheminformatics* **12**, 1–13 (2020).