
Molphenix: A Multimodal Foundation Model for PhenoMolecular Retrieval

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Predicting molecular impact on cellular function is a core challenge in therapeutic
2 design. Phenomic experiments, designed to capture cellular morphology, utilize
3 microscopy based techniques and demonstrate a high throughput solution for un-
4 covering molecular impact on the cell. In this work, we learn a joint latent space
5 between molecular structures and microscopy phenomic experiments, aligning
6 paired samples with contrastive learning. Specifically, we study the problem of
7 *Contrastive PhenoMolecular Retrieval*, which consists of zero-shot molecular struc-
8 ture identification conditioned on phenomic experiments. We assess challenges
9 in multi-modal learning of phenomics and molecular modalities such as experi-
10 mental batch effect, inactive molecule perturbations, and encoding perturbation
11 concentration. We demonstrate improved multi-modal learner retrieval through
12 (1) a uni-modal pre-trained phenomics model, (2) a novel inter sample similarity
13 aware loss, and (3) models conditioned on a representation of molecular concentra-
14 tion. Following this recipe, we propose *MolPhenix*, a molecular phenomics model.
15 MolPhenix leverages a pre-trained phenomics model to demonstrate significant
16 performance gains across perturbation concentrations, molecular scaffolds, and
17 activity thresholds. In particular, we demonstrate an $8.1\times$ improvement in zero shot
18 molecular retrieval of active molecules over the previous state-of-the-art, reaching
19 77.33% in top-1% accuracy. These results open the door for machine learning to
20 be applied in virtual phenomics screening, which can significantly benefit drug
21 discovery applications.

22 1 Introduction

23 Quantifying cellular responses elicited by genetic and molecular perturbations represents a core
24 challenge in medicinal research [4, 48]. Out of an approximate 10^{60} druglike molecule designs,
25 a small number are able to alter cellular properties to reverse the course of diseases [5, 22]. In
26 recent years, microscopy-based cell morphology screening techniques, demonstrated potential for
27 quantitative understanding of a molecule’s biological effects. Experimental techniques such as
28 cell-painting are used to capture cellular morphology, which correspond to physical and structural
29 properties of the cell [6, 7]. Cells treated with molecular perturbations can change morphology,
30 which is captured by staining and high throughput microscopy techniques. Perturbations with similar
31 cellular impact induce analogous morphological changes, allowing to capture underlying biological
32 effects in phenomic experiments. Identifying such perturbations with similar morphological changes
33 can aid in discovery of novel therapeutic drug candidates [42, 24, 19].

34 Determining molecular impact on the cell can be formulated as a multi-modal learning problem,
35 allowing us to build on a rich family of methods [35, 52, 45]. Similar to text-image models,
36 paired data is collected from phenomic experiments along with molecules used to perturb the cells.

37 Contrastive objectives have been used as an effective approach in aligning paired samples from
 38 different modalities [35, 27]. A model that has learned a cross-modal joint latent space must be
 39 able to retrieve a molecular perturbant conditioned on the phenomic experiment. We identify this
 40 problem as *contrastive phenomolecular retrieval* (see Figure 2). Addressing this problem can allow
 41 for identification of molecular impact on cellular function, however, this comes with its own set of
 42 challenges. [15, 2, 54].

43 **(1)** Firstly, multi-modal paired phenomics molecular data suffers from lower overall dataset sizes and
 44 is subject to batch effects. Challenges with uniform processing and prohibitive costs associated with
 45 acquisition of paired data, leads to an order of magnitude fewer data points compared to text-image
 46 datasets [41, 9]. Furthermore, data is subject to random batch effects that capture non-biologically
 47 meaningful variation [28, 46]. **(2)** Paired phenomic-molecular data contains inactive perturbations
 48 that do not have a biological effect or do not perturb cellular morphology. It is difficult to infer a
 49 priori whether a molecule has a cellular effect, leading to the collection of paired molecular structures
 50 with unperturbed cells. These data-points are challenging to filter out without an effective phenomics
 51 embedding, as morphological effects are rarely discernible. These samples can be interpreted as
 52 misannotated, under the assumption of all collected pairs having biologically meaningful interactions.
 53 **(3)** Finally, a complete solution for capturing molecular effects on cells must capture molecular
 54 concentration. The same molecule can have drastically different effects along its dose response curve,
 55 thus making concentration an essential component for learning molecular impact.

56 In this work, we explore the problem of contrastive phenomolecular retrieval by addressing the above
 57 challenges circumvented in prior works. Our key contributions are as follows:

- 58 • We demonstrate significantly higher phenomolecular retrieval rates by utilizing a pretrained uni-
 59 modal phenomic encoder. Thus alleviating the data availability challenge, reducing the impact of
 60 batch effects, and identifying molecular activity levels.
- 61 • We propose a novel soft-weighted sigmoid locked loss (S2L) that addresses the effects of inactive
 62 molecules. This is done by leveraging distances computed in the phenomic embedding space to
 63 learn inter-sample similarities.
- 64 • We explore *explicit* and *implicit* methods to encode molecular concentration, assessing the model’s
 65 ability to perform retrieval in an inter-concentration setting and generalize to unseen concentrations.

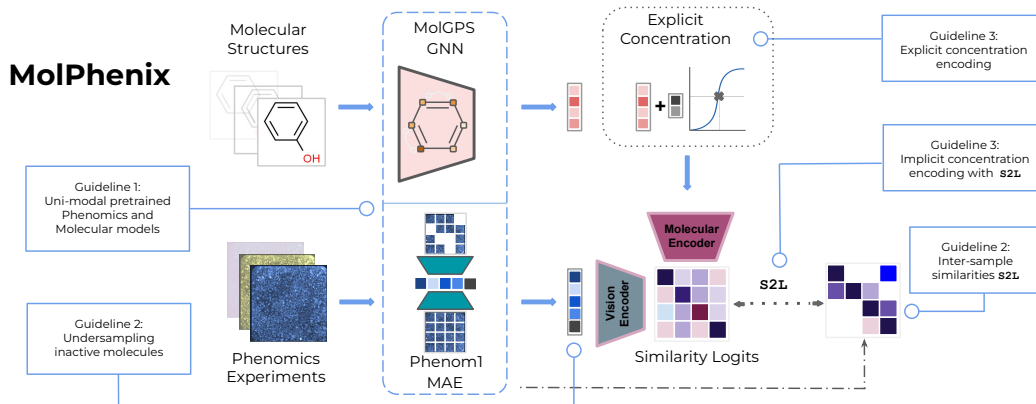


Figure 1: Illustration of proposed guidelines when incorporated in our *MolPhenix* contrastive phenomolecular retrieval framework. We address challenges by utilizing uni-modal pretrained MAE & MPNN models, inter-sample weighting with a dosage aware S2L loss, undersampling inactive molecules, and encoding molecular concentration.

66 Following these principles, we build *MolPhenix*, a multi-modal molecular phenomics model address-
 67 ing contrastive phenomolecular retrieval (Figure 1). *MolPhenix* demonstrates large and consistent
 68 improvements in the presence of batch effects, generalizing across different concentrations, molecules,
 69 and activity thresholds. Additionally, *MolPhenix* outperforms baseline methods in zero-shot setting,
 70 achieving 77.33% top-1% retrieval accuracies on active molecules, which corresponds to a **8.1**×
 71 improvement over the previous state-of-the-art (SOTA) [40].

72 2 Related Work

73 **Uni-modality Pretraining:** Self-supervised methods have demonstrated success across a variety
74 of domains such as computer vision, natural language processing and molecular representations
75 [3, 36, 51]. In vision, contrastive methods have been used to minimize distance in the model’s
76 latent space of two views of the same sample [10, 43, 16, 18]. Reconstruction objectives have
77 also permeated computer vision, such as masked autoencoders (MAE). MAEs typically utilize
78 vision transformers to partition the image into learnable tokens and reconstruct masked patches
79 [17, 14, 8, 12]. These methods have been extended to microscopy experimental data designed
80 to capture cell morphology [50, 23]. Phenom1 utilizes a masked autoencoder with a ViT-L/8+
81 architecture and a custom Fourier domain reconstruction loss, yielding informative representations of
82 phenomic experiments [23, 11]. From a representational perspective, Graph Neural Networks (GNN)
83 have been used to predict molecular properties by reasoning over graph structures. A combination
84 of reconstruction and supervised objectives have led to models generalizing to a diverse range of
85 prediction tasks [31, 55, 47, 39]. Our work leverages uni-modal foundation models, which are used
86 to generate embeddings of phenomic images and molecular graphs.

87 **Multi-Modal Objectives:** Multi-modal models combine samples from two or more domains, to
88 learn rich representations and demonstrate flexible ways to predict sample properties [35, 1, 20].
89 Contrastive methods minimize distances between paired samples, traditionally in text-image domains.
90 However, training these models is computationally expensive, requiring large datasets. Multiple
91 contributions have allowed for a reduction in compute and data budgets by an order of magnitude. In
92 *LiT*, the authors demonstrate that utilizing uni-modal pretrained models for one or both modalities
93 matches zero-shot performance with an order of magnitude fewer paired examples seen [53]. Zhai
94 et al. (2023) demonstrate that by replacing the softmax operation over cosine similarities with an
95 element wise sigmoid loss, allows contrastive learners to improve performance under label noise
96 regime [52]. By using a uni-modal pre-trained modal to calculate similarities between samples from
97 one of the modalities, Srinivasa et al. (2023) have demonstrated improved performance on zero-shot
98 evaluation [45]. In our work, we build along these directions in molecular phenomic multi-modal
99 training.

100 **Molecular-Phenomic Contrastive Learning:** Prior works in contrastive phenomic retrieval have
101 utilized the InfoNCE objective as a pre-training technique to construct uni-modal representations
102 [32]. Recent methods have attempted to improve retrieval by using the InfoLOOB objective [34].
103 Specifically, CLOOME utilizes the InfoLOOB loss with hopfield networks for zero-shot retrieval
104 on unseen data samples [37, 40]. Our work is parallel to the above directions, demonstrating a
105 significant increase in molecular-phenomic retrieval by building on algorithmic improvements from
106 the multi-modality literature.

107 3 Methodology

108 In this section, we explain key challenges facing phenomolecular retrieval and provide guidelines
109 that are key methodological improvements behind the success of MolPhenix 1.

110 **Preliminaries:** Our setting studies the problem of learning multi-modal representations of molecules
111 and phenomic experiments of treated cells [40]. The aim of this work is to learn a joint latent space
112 which maps phenomic experiments of treated cells and the corresponding molecular perturbations
113 into the same latent space. We consider a set of lab experiments \mathcal{E} defined as the tuple $(\mathbf{X}, \mathbf{M}, \mathbf{C}, \Psi)$.
114 Each experiment $\epsilon \in \mathcal{E}$ consists of data samples $\mathbf{x}_i \in \mathbf{X}$ (such as images) and perturbations $\mathbf{m}_i \in \mathbf{M}$
115 (such as molecules) which are obtained at a specific dosage concentration $\mathbf{c}_i \in \mathbf{C}$, while $\psi \in \Psi$
116 denotes molecular activity threshold.

117 Figure 2 describes the problem of contrastive phenomolecular retrieval, where for a single image \mathbf{x}_i ,
118 the challenge consists of identifying the matching perturbation, \mathbf{m}_i , and concentration, \mathbf{c}_i , used to
119 induce morphological effects. This can be accomplished in a zero-shot way by generating embeddings
120 for $(\mathbf{m}_1, \mathbf{c}_1), \dots, (\mathbf{m}_j, \mathbf{c}_j)$ and \mathbf{x}_i using functions $f_{\theta_m}(\mathbf{m}, \mathbf{c})$, $f_{\theta_x}(\mathbf{x})$ which map samples into \mathbb{R}^d .

121 Then, by defining a similarity metric between generated embeddings \mathbf{z}_{x_i} and \mathbf{z}_{m_i} , f_{sim} , we can
122 rank $(\mathbf{m}_1, \mathbf{c}_1), \dots, (\mathbf{m}_j, \mathbf{c}_j)$ based on computed similarities. An effective solution to the contrastive
123 phenomolecular retrieval problem would learn $f_{\theta_m}(\mathbf{m}, \mathbf{c})$ and $f_{\theta_x}(\mathbf{x})$ that results in consistently high
124 retrieval rates of $(\mathbf{m}_i, \mathbf{c}_i)$ used to perturb \mathbf{x}_i .

125 In practice, the image embeddings are gener-
 126 ated using a phenomics microscopy foundation
 127 MAE model [23, 17]. We use phenomic embed-
 128 dings to marginalize batch effects, infer inter-
 129 sample similarities, and undersample inactive
 130 molecules. Activity is determined using consis-
 131 tency of replicate measurements for a given
 132 perturbation. For each sample, a p value cutoff
 133 $\psi \in \Psi$ is used to quantify molecular activity.
 134 Only molecules below the p value cutoff ψ are
 135 considered active.

136 Prior methods in multi-modal contrastive learn-
 137 ing utilize the InfoNCE loss, and variants thereof
 138 [32] to maximize the joint likelihood of \mathbf{x}_i and
 139 \mathbf{m}_i . Given a set of $N \times N$ random samples
 140 $(\mathbf{x}_1, \mathbf{m}_1, \mathbf{c}_1), \dots, (\mathbf{x}_N, \mathbf{m}_N, \mathbf{c}_N)$ containing N
 141 positive samples at k^{th} index and $(N - 1) \times N$
 142 negative samples, optimizing Equation 1 maximizes the likelihood of positive pairs while minimizing
 143 the likelihood of negative pairs:

$$\mathcal{L}_{\text{InfoNCE}} = -\frac{1}{N} \sum_{i=1}^N \left[\log \frac{\exp(\langle \mathbf{z}_{x_i}, \mathbf{z}_{m_i} \rangle / \tau)}{\sum_{k=1}^N \exp(\langle \mathbf{z}_{x_i}, \mathbf{z}_{m_k} \rangle / \tau)} + \log \frac{\exp(\langle \mathbf{z}_{x_i}, \mathbf{z}_{m_i} \rangle / \tau)}{\sum_{k=1}^N \exp(\langle \mathbf{z}_{m_i}, \mathbf{z}_{x_k} \rangle / \tau)} \right]. \quad (1)$$

144 Where $\mathbf{z}_x, \mathbf{z}_m$ correspond to phenomics and molecular embeddings respectively, τ is softmax
 145 temperature, and $\langle \cdot \rangle$ corresponds to cosine similarity.

146 Challenge 1: Phenomic Pretraining and Generalization

147 We find that using a phenomics foundation model to embed microscopy images allows for mitigation
 148 of batch effects, reduces the required number of paired data points, and improves generalization in the
 149 process. While CLIP, a hallmark model in the field of text-image multi-modality, was trained on 400
 150 million curated paired data points, there is an order of magnitude fewer paired molecular-phenomic
 151 molecule samples [35]. Cost and systematic pre-processing of data make large scale data generation
 152 efforts challenging, and resulting data is affected by experimental batch effects. **Batch effects** induce
 153 noise in the latent space as a result of random perturbations in the experimental process, while
 154 biologically meaningful variation remains unchanged [33, 44]. Limited dataset sizes and batch effects
 155 make it challenging for contrastive learners to capture molecular features affecting cell morphology,
 156 yielding low retrieval rates [40].

157 We address data availability and generalization challenges by utilizing representations from a large **uni-**
 158 **modal pre-trained phenomic model**, θ_{ph} , trained to capture representations of cellular morphology.
 159 θ_{ph} is pretrained on microscopy images using a Fourier modified MAE objective, utilizing the
 160 ViT-L/8 architecture with methodology similar to Kraus et al. (2024) [17, 12, 23]. For simplicity
 161 in future sections, we refer to this model as *Phenom1*. This pretrained model allows a drastic
 162 reduction in the required number of paired multi-modal samples [53]. In addition, using phenomic
 163 representations alleviates the challenge of batch effects by averaging samples, \mathbf{z}_x , generated with the
 164 same perturbation \mathbf{m}_i over multiple lab experiments ϵ_i . Averaging model representations $\frac{1}{N} \sum_{i \in N} \mathbf{z}_{x_i}$
 165 allows marginalizing batch effect induced by individual experiments.

Guideline 1 Utilizing pre-trained uni-modal encoder, θ_{ph} , can be used to reduce the num-
 ber of paired data-points compared to training θ without prior optimization. In addition,
 averaging phenomic embeddings \mathbf{z}_x from matched perturbations can alleviate batch effects.

166

167 To reason over molecular structures, we make use of features learned from GNNs trained on molecular
 168 property prediction [29]. We utilize a pretrained MPNN foundational model up to the order of 1B
 169 parameters for extracting molecular representations following a similar procedure to Sypetkowski et
 170 al. (2024) [47]. We refer to this model as *MolGPS*.

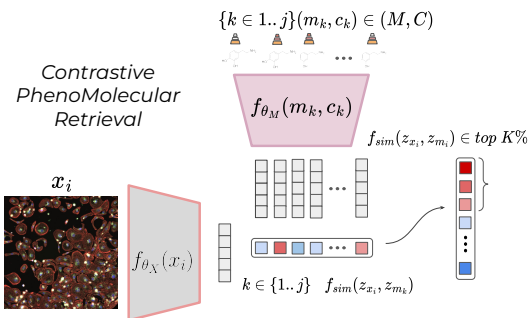


Figure 2: Illustration of the contrastive phenomolecular retrieval challenge. Image \mathbf{x}_i and a set of molecules and corresponding concentrations $(\mathbf{m}_k, \mathbf{c}_k)$ get mapped into a \mathbb{R}^d latent space. Their similarities get computed with f_{sim} and ranked to evaluate whether the paired perturbation appears in the top $K\%$.

171 **Challenge 2: Inactive Molecular Perturbations**

172 The phenomics-molecular data collection process can
 173 result in pairing of molecular structures with unper-
 174 turbed cells in cases where the molecule has no effect
 175 on cell morphology (Figure 3)

176 Since the morphological effects observed in cell \mathbf{x}_i
 177 is conditioned on the perturbation, in the absence of
 178 a molecular effect $P(\mathbf{x}_i|\mathbf{x}_i^0, \mathbf{c}_i, \mathbf{m}_i) \sim P(\mathbf{x}_i|\mathbf{x}_i^0)$. In
 179 these samples, phenomic data will be independent,
 180 from paired molecular data, which results in misanno-
 181 tation under the assumption of data-pairs having an
 182 underlying biological relationship. We demonstrate how utilizing Phenom1 to undersample inactive
 183 molecules and learn continuous similarities between samples can alleviate this challenge.

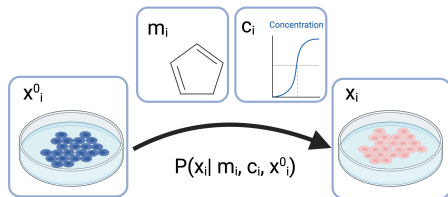


Figure 3: Data generation process of a phenomic experiment on cells \mathbf{x}_i with molecular perturbations \mathbf{m}_i and concentrations \mathbf{c}_i .

184 To **undersample inactive molecules**, we extract the embeddings from Phenom1 and calculate the
 185 relative activity of each perturbation $(\mathbf{m}_i, \mathbf{c}_i) \in (\mathbf{M}, \mathbf{C})$. This is done using the rank of cosine simi-
 186 larities between technical replicates produced for a molecular perturbation against a null distribution.
 187 The null distribution is established by calculating cosine similarities from random pairs of Phenom1
 188 embeddings generated with perturbation $(\mathbf{m}_j, \mathbf{c}_j)$, $(\mathbf{m}_k, \mathbf{c}_k)$. Hence, we can compute a p-value and
 189 filter out samples likely to belong to the null distribution with an arbitrary threshold ψ .

190 In addition, by utilizing an inter-sample aware S2L **training objective**, the model can learn similarities
 191 between inactive molecules. S2L is grounded in previous work which demonstrates improved
 192 robustness to label noise (SigLip) and learnable inter-sample associations (CWCL) [52, 45]. Continu-
 193 ous Weighted Contrastive Loss (CWCL) provides better multi-modal alignment using a uni-modal
 194 pretrained model to suggest sample distances, relaxing the negative equidistant assumption present in
 195 InfoNCE [45]:

$$\mathcal{L}_{\text{CWCL}, \mathcal{M} \rightarrow \mathcal{X}} = -\frac{1}{N} \sum_{i=1}^N \left[\frac{1}{\sum_{j=1}^N \mathbf{w}_{i,j}^{\mathcal{X}}} \sum_{j=1}^N \mathbf{w}_{i,j}^{\mathcal{X}} \log \frac{\exp(\langle \mathbf{z}_{\mathbf{x}_i}, \mathbf{z}_{\mathbf{m}_j} \rangle / \tau)}{\sum_{k=1}^N \exp(\langle \mathbf{z}_{\mathbf{x}_j}, \mathbf{z}_{\mathbf{m}_k} \rangle / \tau)} \right]. \quad (2)$$

196 CWCL weights logits with a continuous measure of similarity $\mathbf{w}^{\mathcal{X}}$, resulting in better alignment of
 197 embeddings $\mathbf{z}_{\mathbf{x}_i}$ and $\mathbf{z}_{\mathbf{m}_j}$ across modalities. In equation 2, $\mathbf{w}^{\mathcal{X}}$ is computed using a within modality
 198 similarity function such as $\mathbf{w}_{i,j}^{\mathcal{X}} = \langle z_{\mathbf{x}_i}, z_{\mathbf{x}_j} \rangle / 2 + 0.5$. Note, the above formula is used only for
 199 mapping samples from modality \mathcal{M} to \mathcal{X} for which a pre-trained model θ_{ph} is available.

200 Another work, SigLIP, demonstrates robustness to label noise and reduces computational requirements
 201 during contrastive training [52]. It does so by avoiding computation of a softmax over the entire set
 202 of in-batch samples, instead relying on element-wise sigmoid operation:

$$\mathcal{L}_{\text{SigLIP}} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \left[\log \frac{1}{1 + \exp(\mathbf{y}_{i,j}(-\alpha \langle \mathbf{z}_{\mathbf{x}_i}, \mathbf{z}_{\mathbf{m}_j} \rangle + b))} \right]. \quad (3)$$

203 In equation 3, α and b are learned, calibrating the model confidence conditioned on the ratio of
 204 positive to negative pairs. $\mathbf{y}_{i,j}$ is set to 1 if $i = j$ and -1 otherwise.

205 Inspired by prior works, we introduce S2L for molecular representation learning, which leverages
 206 inter-sample similarities and robustness to label noise to mitigate weak or inactive perturbations.

$$\mathcal{L}_{\text{S2L}} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N \log \left[\frac{\mathbf{w}_{i,j}^{\mathcal{X}}}{1 + \exp(-\alpha \mathbf{z}_{\mathbf{x}_i} \cdot \mathbf{z}_{\mathbf{m}_j} + b)} + \frac{(1 - \mathbf{w}_{i,j}^{\mathcal{X}})}{1 + \exp(\alpha \mathbf{z}_{\mathbf{x}_i} \cdot \mathbf{z}_{\mathbf{m}_j} + b)} \right]. \quad (4)$$

207 In the equation above, $\mathbf{z}_{\mathbf{x}_i}$ and $\mathbf{z}_{\mathbf{m}_j}$ correspond to latent representations of images and molecules,
 208 respectively. α and b correspond to learnable temperature and bias parameters for the calibrated
 209 sigmoid function. $\mathbf{w}_{i,j}^{\mathcal{X}}$ is an inter-sample similarity function computed from images using the
 210 pretrained model θ_{ph} . To compute $\mathbf{w}_{i,j}^{\mathcal{X}}$, we use the arctangent of L2 distance instead of cosine

211 similarity, as was the case for Equation 2 (more details in Appendix D.3). Intuitively, S2L can be
212 thought of as shifting from a multi-class classification to a soft multi-label problem. In our problem
213 setting, the labels are continuous and determined by sample similarity in the phenomics space.

Guideline 2 *When training a molecular-phenomic model, mitigating the effect of inactive molecules in training data distribution can be carried out by undersampling inactive molecules and using an inter-sample similarity aware, S2L loss (equation 4).*

214

215 **Challenge 3: Variable Concentrations**

216 Perturbation effect on a cell is determined by both molecular structure and corresponding concentra-
217 tion [49]. A model capturing molecular impact on cell morphology must be able to generalize across
218 different doses, since variable concentrations can correspond to different data distributions.

219 We note that providing concentrations c_i as input to the model would benefit performance, as this
220 would indicate the magnitude of molecular impact. However, we find that simply concatenating
221 concentrations does not result in effective training due to its compressed dynamic range. To that end,
222 we add concentration information in two separate ways: *implicit* and *explicit* formulations.

223 We add **implicit concentration** as molecular perturbation classes by using the S2L loss (Equation 4)
224 to treat perturbation m_i with concentrations c_i and c_j as distinct classes. This pushes samples apart
225 in the latent space proportionally to similarities between phenomic experiments.

226 We add **explicit concentration** c_i by passing it to the molecular encoder. We explore different
227 formulation for dosage concentrations, $f'(c_i)$, where f' maps $c_i \rightarrow \mathbb{R}$. Encoded representations $f'(c_i)$
228 are concatenated at the initial layer of the model. We find simple functional encodings f' (such as
229 one-hot and `logarithm`) to work well in practice.

Guideline 3 *When training a molecular-phenomic model, conditioning on an (implicit and explicit) representation of concentration $f'(c_i)$ aids in capturing molecular impacts on cell morphology and improves generalization to previously unseen molecules and concentrations.*

230

231 **4 Experimental Setup**

232 In this section, we describe evaluation datasets used, and descriptions of the underlying data modalities.
233 To assess phenomolecular retrieval, we use 1% recall metric unless stated otherwise, as it allows
234 direct comparison between datasets with different number of samples. Additional implementation
235 and evaluation details can be found in Appendix D.

236 **Datasets:** Our training dataset consists of fluorescent microscopy images paired with molecular
237 structures and concentrations, which are used as perturbants. We assess models' phenomolecular
238 retrieval capabilities on three datasets of escalating generalization complexity. First dataset, consisting
239 of unseen microscopy images and molecules present in the training dataset. Second, a dataset consist-
240 ing of previously unseen phenomics experiments and molecules split by the corresponding molecular
241 scaffold. Finally, we evaluate on an open source dataset with a different data generating distribution
242 [13]. In the case of the latter two datasets, the model is required to perform zero-shot classification,
243 as it has no access to those molecules in the training data. This requires the model to reason over
244 molecular graphs to identify structures inducing corresponding cellular morphology changes. Using
245 methodology described in guideline 2 we report retrieval results for all molecules as well as on an
246 active subset. Finally, all datasets are comprised of molecular structures at multiple concentrations
247 (.01, .1, 1.0, 10, etc.) Additional details regarding the datasets can be found in Appendix C.

248 **Modality Representations:** In our evaluations, we consider different representations for molecular
249 perturbations and phenomic experiments and quantitatively evaluate their impact.

250 • **Images:** Image encoders utilize 6-channel fluorescent microscopy images of cells representing
251 phenomic experiments. Images are 2048×2048 pixels, capturing cellular morphology changes post
252 molecular perturbation. We downscale each image to 256×256 using block mean downsampling.

- 253 • Phenom1: We characterize phenomic experiments by embedding high resolution microscopy
254 images in the latent space of a phenomics model θ_{Ph} as described in guideline 1.
- 255 • Fingerprints: Molecular fingerprints utilize RDKit [26], MACCS [25] and MORGAN3 [38] bit
256 coding, which represent binary presence of molecular substructures. Additional information such
257 as atomic identity, atomic radius and torsional angles are included in the fingerprint representations.
- 258 • MolGPS: We generate molecular representations from a large pretrained GNN. Specifically, we
259 obtain molecular embeddings from a 1B parameter MPNN [29].

260 5 Results and Discussion

261 To evaluate the effectiveness of Guidelines 1, 2, and
262 3 we carry out evaluation in two different settings:
263 (1) cumulative concentrations, and (2) held-out concen-
264 trations, testing the models’ ability to generalize
265 to new molecular doses. Finally, we perform compre-
266 hensive ablations testing model performance with
267 varying data, model, and optimization parameters.
268 The comprehensive set of results can be found in
269 Tables 10, 11, 12, and 13.

270 5.1 Evaluation on cumulative concentrations:

271 We demonstrate improvements in phenomolecular
272 recall due to usage of a phenomics pre-trained founda-
273 tion model, identify that MolPhenix set of design
274 choices results in higher final performance, and more
275 data efficient learning. Figure 4 demonstrates recall
276 accuracy on all molecules and an active subset for
277 CLOOME and MolPhenix models, as a function of training samples seen.

278 We observe a large performance gap between models trained on Phenom1 embeddings as opposed to
279 images, emphasizing the utility of using a pre-trained encoder for microscopy images (Table 1). We
280 note that provision of Phenom1 (CLOOME-Phenom1 Vs CLOOME-Images) significantly improves
281 both active and all molecule retrieval by $5.69\times$ and, $4.75\times$ respectively (Table 1).

282 Furthermore, we identify that while all molecules retrieval stagnates throughout training, the per-
283 formance on an active subset keeps improving, underscoring the importance of identification of the
284 active subset. Finally, we compare CLOOME and MolPhenix trained using Phenom1 embeddings
285 and find there is a consistent retrieval performance gap, throughout training, with a $1.26\times$ final
286 improvement (Figure 4, Table 1). Compared to CLOOME [40] trained directly on images, MolPhenix
287 achieves an average improvement of $8.78\times$ on active molecules on the unseen dataset. These results
288 verify the effectiveness of Guideline 1 in accelerating training, and the importance of Guidelines 2
289 and 3 in recall improvements over CLOOME.

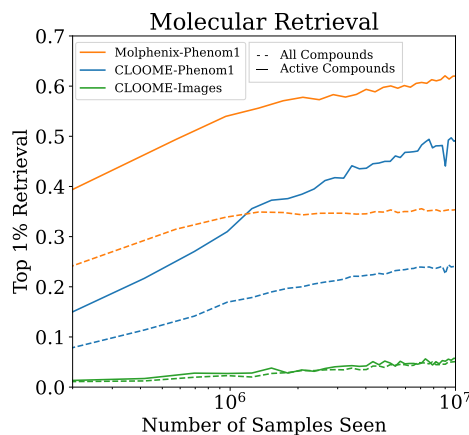


Figure 4: Comparison of training phenomic encoder from scratch and utilizing pre-trained Phenom1 unseen dataset. X-axis plotted on logarithmic scale.

Table 1: Impact of pre-trained Phenom1 and MolGPS on CLOOME and MolPhenix for a matched number of seen samples (Top), where we observe an $8.1\times$ improvement of MolPhenix over the CLOOME baseline for active unseen molecules. SOTA results trained with a higher number of steps by utilizing the best hyperparameters (Bottom *). We note that MolPhenix’s main components such as S2L and embedding averaging relies on having a pre-trained uni-modal phenomics model.

Method	Modality	Active Molecules			All Molecules		
		Unseen Im.	Unseen Im. + Mol.	Unseen Dataset	Unseen Im.	Unseen Im. + Mol.	Unseen Dataset
CLOOME	Images & Multi-FPS	.0756 ± .0042	.0787 ± .0065	.0528 ± .0057	.0547 ± .0028	.0661 ± .0020	.0223 ± .0014
CLOOME	Phenom1 & Multi-FPS	.4659 ± .0042	.5057 ± .0014	.2065 ± .0146	.3009 ± .0053	.2474 ± .0013	.1337 ± .0045
MolPhenix	Phenom1 & Multi-FPS	.7807 ± .0025	.6365 ± .0014	.3545 ± .0097	.5253 ± .0029	.3655 ± .0017	.2163 ± .0021
MolPhenix	Phenom1 & MolGPS	.7646 ± .0014	.6387 ± .0056	.4160 ± .0016	.5012 ± .0002	.3511 ± .0004	.2508 ± .0026
MolPhenix*	Phenom1 & MolGPS	.9689 ± .0017	.7733 ± .0036	.5860 ± .0082	.5583 ± .0007	.3824 ± .0016	.2809 ± .0060

Table 2: Top-1% recall accuracy with use of the proposed MolPhenix guidelines, such as Phenom1 and embedding averaging. We omit explicit concentration from this experiment.

Loss	Active Molecules			All Molecules		
	Unseen Images	Unseen Im. + Mol.	Unseen Dataset	Unseen Images	Unseen Im. + Mol.	Unseen Dataset
CLIP	.3373 ± .0043	.4228 ± .0008	.1514 ± .0038	.1761 ± .0043	.1867 ± .0022	.0734 ± .0022
Hopfield-CLIP	.2578 ± .0042	.3559 ± .0042	.1256 ± .0092	.1531 ± .0046	.1709 ± .0029	.0673 ± .0020
InfoLOOB	.3351 ± .0011	.4206 ± .0031	.1563 ± .0028	.1746 ± .0003	.1860 ± .0029	.0745 ± .0019
CLOOME	.3572 ± .0026	.4348 ± .0039	.1658 ± .0063	.1968 ± .0029	.2005 ± .0026	.0911 ± .0022
DCL	.6363 ± .0025	.6177 ± .0047	.3184 ± .0087	.3277 ± .0047	.2562 ± .0008	.1364 ± .0067
CWCL	.7091 ± .0045	.6529 ± .0020	.3556 ± .0094	.3635 ± .0064	.2696 ± .0019	.1526 ± .0058
SigLip	.7763 ± .0045	.6401 ± .0065	.3396 ± .0042	.3729 ± .0039	.2544 ± .0014	.1470 ± .0038
S2L (ours)	.9097 ± .0020	.6759 ± .0012	.4181 ± .0012	.4688 ± .0009	.2852 ± .0001	.1838 ± .0007

Table 3: Top-1% recall accuracy across different concentration encoding choices with use of the proposed MolPhenix guidelines, such as Phenom1 and embedding averaging.

Implicit Concentration	Explicit Concentration	Active Molecules			All Molecules		
		Unseen Im.	Unseen Im. + Mol.	Unseen Dataset	Unseen Im.	Unseen Im. + Mol.	Unseen Dataset
✗	✗	.7350 ± .0071	.6509 ± .0104	.3333 ± .0004	.3610 ± .0025	.2668 ± .0034	.1532 ± .0007
✓	✗	.9097 ± .0020	.6759 ± .0012	.4181 ± .0012	.4688 ± .0009	.2852 ± .0001	.1838 ± .0007
✓	sigmoid	.9423 ± .0011	.7155 ± .0016	.4573 ± .0022	.5071 ± .0024	.3441 ± .0026	.2144 ± .0026
✓	logarithm	.9426 ± .0066	.7451 ± .0050	.4727 ± .0056	.5183 ± .0027	.3700 ± .0036	.2275 ± .0032
✓	one-hot	.9430 ± .0029	.7490 ± .0052	.4850 ± .0020	.5433 ± .0030	.3819 ± .0032	.2384 ± .0049

290 We evaluate the impact of different loss objectives on the proposed MolPhenix training frame-
 291 work. Table 2 presents top-1% retrieval accuracy across different contrastive losses utilized to
 292 train molecular-phenomics encoders on cumulative concentrations. Compared to prior methods, the
 293 proposed S2L loss demonstrates improved retrieval rates in cumulative concentration setting. Label
 294 noise and inter-sample similarity aware losses such as CWCL and SigLip also demonstrate improved
 295 performance. The effectiveness of S2L can be attributed to smoothed inter-sample similarities and
 296 implicit concentration information.

297 Finally, in Table 3, we observe recall improvements when considering both molecular structures and
 298 concentration. We note the importance of the addition of implicit concentration, further confirming the
 299 importance of considering molecular effects at different concentrations as different classes. Explicitly
 300 encoding molecular concentration with one-hot, logarithm and sigmoid yields improved recall
 301 performance, where one-hot performs the best in a cumulative concentration setting. These findings
 302 verify the efficacy of implicit and explicit concentration encoding outlined in Guideline 3.

Table 4: Top-1% recall accuracy of different loss objectives while using the proposed MolPhenix guidelines, such as Phenom1 and embedding averaging.

Loss	Unseen Im.	Unseen Im. + Mol.	Unseen Dataset
CLIP	.2109	.2425	.1519
Hopfield-CLIP	.1581	.2034	.1198
InfoLOOB	.2122	.2496	.1501
CLOOME	.2164	.2461	.1479
DCL	.4717	.4027	.2841
CWCL	.5731	.4403	.3232
SigLip	.5718	.4217	.3021
S2L (ours)	.8334	.4615	.3792

Table 5: Top-1% recall accuracy across different concentration encoding choices while using the proposed MolPhenix guidelines, such as Phenom1 and embedding averaging.

Implicit Concentration	Explicit Concentration	Unseen Im.	Unseen Im. + Mol.	Unseen Dataset
✗	✗	.5942	.4315	.3129
✓	✗	.8334	.4615	.3792
✓	sigmoid	.8256	.4692	.3765
✓	logarithm	.7953	.4466	.3664
✓	one-hot	.7489	.4088	.3379

Results are averaged across experiments for each dropped concentration, and across three seeds. Recall is reported for active molecules, while the results for all molecules can be found in Table 13.

303 5.2 Evaluation on held-out concentrations:

304 Next, we evaluate recall on held-out concentrations to obtain a measure of generalization performance.
 305 This evaluation allows us to capture the utility of our models for prediction of unseen concentrations,

306 hence resembling *in-silico* testing. We omit concentrations from the training set and evaluate recall
 307 at the excluded data, where we observe a drop in retrieval performance for unseen concentrations.
 308 Similar to cumulative concentration results, we find that using S2L improves recall over other losses
 309 and outperforms CLOOME by up to 126% (Table 4). While one-hot encoding exhibits significant
 310 improvements in cumulative concentrations, its expressivity on unseen concentrations is limited
 311 (Table 5) and sigmoid encoding provides a sufficient representation of concentration.

312 5.3 Ablation Studies

313 We assess the importance of our design decisions by conducting an ablation study over our proposed
 314 guidelines. Figure 5 presents the variation of top-1% recall accuracy across key components such as
 315 cutoff p value, fingerprint type, and embedding averaging. We observe that employing a lower cutoff
 316 p value yields improved generalization for unseen dataset, while employing a higher cutoff appears to
 317 be optimal for unseen images + unseen molecules. For molecular structure representations, we find
 318 that using embeddings from the large pretrained MPNN graph based model (e.g., MolGPS) surpasses
 319 traditional fingerprints. Finally, utilization of embedding averaging demonstrates improved recall.
 320 More ablations over model size, projector dimension, and batch size can be found in Appendix E.5.

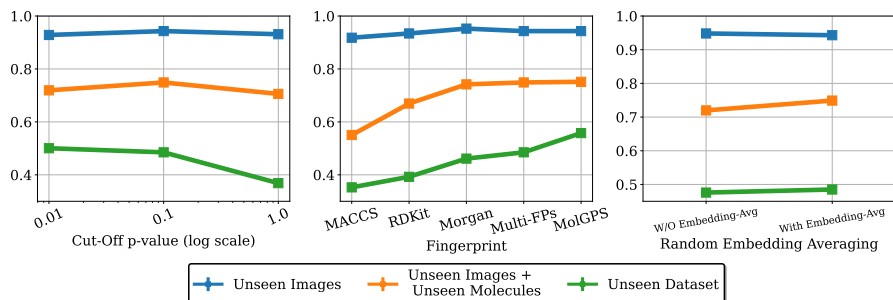


Figure 5: Ablations of top-1 % recall accuracy with **(bottom-left)** cutoff p value, **(bottom-center)** fingerprint type, and **(bottom-right)** embedding averaging.

321 6 Conclusion

322 In this work, we investigate the problem of *contrastive phenomolecular retrieval* by constructing a
 323 joint multi-modal embedding of phenomic experiments and molecular structures. We identify a set of
 324 challenges afflicting molecular-phenomic training and proposed a set of guidelines for improving
 325 retrieval and generalization. Empirically, we observed that contrastive learners demonstrate higher
 326 retrieval rates when using representations from a high-capacity uni-modal pretrained model. Use
 327 of inter-sample similarities with a label noise resistant loss such as S2L allows us to tackle the
 328 challenge of inactive molecules. Finally, adding implicit and explicit concentrations allows models to
 329 generalize to previously unseen concentrations. MolPhenix demonstrates an $8.1\times$ improvement in
 330 zero shot retrieval of active molecules over the previous state-of-the-art, reaching 77.33% in top-1%
 331 accuracy. In addition, we conduct a preliminary investigation on MolPhenix’s ability to uncover
 332 biologically meaningful properties (activity prediction, zero-shot biological perturbation matching,
 333 and molecular property prediction in Appendix E.1, E.2, and E.3, respectively.). We expect a wide
 334 range of applications for MolPhenix, particularly in drug discovery. While there’s a remote chance of
 335 misuse for developing chemical weapons, such harm is unlikely, with our primary focus remaining
 336 on healthcare improvement.

337 **Limitations and Future Work:** While our study covers challenges in phenomolecular recall, we
 338 leave three research directions for future work. (1) Future investigations could consider studying
 339 additional modalities such as text, genetic perturbations and chemical multi-compound interventions.
 340 (2) While we propose and evaluate our guidelines on previously conducted phenomic experiments,
 341 we note that a rigorous evaluation would evaluate model predictions in a wet-lab setting. (3) In
 342 addition, our work makes the assumption that the initial unperturbed cell state x_i^0 can be marginalized
 343 by utilizing a single cell line with an unperturbed genetic background. Future works can relax this
 344 assumption, aiming to capture innate intercellular variation.

345 References

- 346 [1] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican,
347 M. Reynolds, R. Ring, E. Rutherford, S. Cabi, T. Han, Z. Gong, S. Samangooei, M. Monteiro,
348 J. Menick, S. Borgeaud, A. Brock, A. Nematzadeh, S. Sharifzadeh, M. Binkowski, R. Barreira,
349 O. Vinyals, A. Zisserman, and K. Simonyan. Flamingo: a visual language model for few-shot
350 learning, 2022.
- 351 [2] S. Albelwi. Survey on self-supervised learning: auxiliary pretext tasks and contrastive learning
352 methods in imaging. *Entropy*, 24(4):551, 2022.
- 353 [3] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes, A. Bardes,
354 G. Mialon, Y. Tian, A. Schwarzschild, A. G. Wilson, J. Geiping, Q. Garrido, P. Fernandez,
355 A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum. A cookbook of self-supervised learning,
356 2023.
- 357 [4] C. Bock, P. Datlinger, F. Chardon, M. A. Coelho, M. B. Dong, K. A. Lawson, T. Lu, L. Maroc,
358 T. M. Norman, B. Song, G. Stanley, S. Chen, M. Garnett, W. Li, J. Moffat, L. S. Qi, R. S.
359 Shapiro, J. Shendure, J. S. Weissman, and X. Zhuang. High-content crispr screening. *Nature*
360 *Reviews Methods Primers*, 2(1), Feb. 2022.
- 361 [5] R. S. Bohacek, C. McMartin, and W. C. Guida. The art and practice of structure-based drug
362 design: A molecular modeling perspective. *Medicinal Research Reviews*, 16(1):3–50, Jan. 1996.
- 363 [6] M. Boutros, F. Heigwer, and C. Laufer. Microscopy-based high-content screening. *Cell*,
364 163(6):1314–1325, 2015.
- 365 [7] M.-A. Bray, S. Singh, H. Han, C. T. Davis, B. Borgeson, C. Hartland, M. Kost-Alimova,
366 S. M. Gustafsdottir, C. C. Gibson, and A. E. Carpenter. Cell painting, a high-content image-
367 based assay for morphological profiling using multiplexed fluorescent dyes. *Nature protocols*,
368 11(9):1757–1774, 2016.
- 369 [8] S. Cao, P. Xu, and D. A. Clifton. How to understand masked autoencoders. *arXiv preprint*
370 *arXiv:2202.03670*, 2022.
- 371 [9] S. N. Chandrasekaran, J. Ackerman, E. Alix, D. M. Ando, J. Arevalo, M. Bennion, N. Boisseau,
372 A. Borowa, J. D. Boyd, L. Brino, P. J. Byrne, H. Ceulemans, C. Ch’ng, B. A. Cimini, D.-A.
373 Clevert, N. Deflaux, J. G. Doench, T. Dorval, R. Doyonnas, V. Dragone, O. Engkvist, P. W.
374 Faloon, B. Fritchman, F. Fuchs, S. Garg, T. J. Gilbert, D. Glazer, D. Gnuttt, A. Goodale,
375 J. Grignard, J. Guenther, Y. Han, Z. Hanifehlu, S. Hariharan, D. Hernandez, S. R. Horman,
376 G. Hormel, M. Huntley, I. Icke, M. Iida, C. B. Jacob, S. Jaensch, J. Khetan, M. Kost-Alimova,
377 T. Krawiec, D. Kuhn, C.-H. Lardeau, A. Lembke, F. Lin, K. D. Little, K. R. Lofstrom, S. Lotfi,
378 D. J. Logan, Y. Luo, F. Madoux, P. A. Marin Zapata, B. A. Marion, G. Martin, N. J. McCarthy,
379 L. Mervin, L. Miller, H. Mohamed, T. Monteverde, E. Mouchet, B. Nicke, A. Ogier, A.-L.
380 Ong, M. Osterland, M. Otrocka, P. J. Peeters, J. Pilling, S. Precht, C. Qian, K. Rataj, D. E.
381 Root, S. K. Sakata, S. Scrace, H. Shimizu, D. Simon, P. Sommer, C. Spruiell, I. Sumia, S. E.
382 Swalley, H. Terauchi, A. Thibaudeau, A. Unruh, J. Van de Waeter, M. Van Dyck, C. van Staden,
383 M. Warchoř, E. Weisbart, A. Weiss, N. Wiest-Daessle, G. Williams, S. Yu, B. Zapiec, M. Żyła,
384 S. Singh, and A. E. Carpenter. Jump cell painting dataset: morphological impact of 136,000
385 chemical and genetic perturbations. *bioRxiv*, 2023.
- 386 [10] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. E. Hinton. Big self-supervised models
387 are strong semi-supervised learners. *Advances in neural information processing systems*,
388 33:22243–22255, 2020.
- 389 [11] M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. P. Steiner, M. Caron,
390 R. Geirhos, I. Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In
391 *International Conference on Machine Learning*, pages 7480–7512. PMLR, 2023.
- 392 [12] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
393 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
394 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- 395 [13] M. M. Fay, O. Kraus, M. Victors, L. Arumugam, K. Vuggumudi, J. Urbanik, K. Hansen, S. Celik,
396 N. Cernek, G. Jagannathan, et al. Rrx3: Phenomics map of biology. *Biorxiv*, pages 2023–02,
397 2023.
- 398 [14] C. Feichtenhofer, Y. Li, K. He, et al. Masked autoencoders as spatiotemporal learners. *Advances*
399 *in neural information processing systems*, 35:35946–35958, 2022.
- 400 [15] A. Fürst, E. Rumetshofer, J. Lehner, V. T. Tran, F. Tang, H. Ramsauer, D. Kreil, M. Kopp,
401 G. Klambauer, A. Bitto, et al. Cloob: Modern hopfield networks with infoloob outperform clip.
402 *Advances in neural information processing systems*, 35:20450–20468, 2022.
- 403 [16] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A.
404 Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko. Bootstrap
405 your own latent: A new approach to self-supervised learning, 2020.
- 406 [17] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable
407 vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
408 *recognition*, pages 16000–16009, 2022.
- 409 [18] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual
410 representation learning, 2020.
- 411 [19] M. Hofmarcher, E. Rumetshofer, D.-A. Clevert, S. Hochreiter, and G. Klambauer. Accurate
412 prediction of biological assays with high-throughput microscopy images and convolutional
413 networks. *Journal of chemical information and modeling*, 59(3):1163–1171, 2019.
- 414 [20] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed,
415 B. Patra, Q. Liu, K. Aggarwal, Z. Chi, J. Bjorck, V. Chaudhary, S. Som, X. Song, and F. Wei.
416 Language is not all you need: Aligning perception with language models, 2023.
- 417 [21] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Rad-
418 ford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint*
419 *arXiv:2001.08361*, 2020.
- 420 [22] C. Knox, M. Wilson, C. M. Klinger, M. Franklin, E. Oler, A. Wilson, A. Pon, J. Cox, N. E. L.
421 Chin, S. A. Strawbridge, M. Garcia-Patino, R. Kruger, A. Sivakumaran, S. Sanford, R. Doshi,
422 N. Khetarpal, O. Fatokun, D. Doucet, A. Zubkowski, D. Y. Rayat, H. Jackson, K. Harford,
423 A. Anjum, M. Zakir, F. Wang, S. Tian, B. Lee, J. Liigand, H. Peters, R. Q. R. Wang, T. Nguyen,
424 D. So, M. Sharp, R. da Silva, C. Gabriel, J. Scantlebury, M. Jasinski, D. Ackerman, T. Jewison,
425 T. Sajed, V. Gautam, and D. S. Wishart. Drugbank 6.0: the drugbank knowledgebase for 2024.
426 *Nucleic Acids Research*, 52(D1):D1265–D1275, Nov. 2023.
- 427 [23] O. Kraus, K. Kenyon-Dean, S. Saberian, M. Fallah, P. McLean, J. Leung, V. Sharma, A. Khan,
428 J. Balakrishnan, S. Celik, D. Beaini, M. Sypetkowski, C. V. Cheng, K. Morse, M. Makes,
429 B. Mabey, and B. Earnshaw. Masked autoencoders for microscopy are scalable learners of
430 cellular biology, 2024.
- 431 [24] O. Z. Kraus, B. T. Grys, J. Ba, Y. Chong, B. J. Frey, C. Boone, and B. J. Andrews. Automated
432 analysis of high-content microscopy data with deep learning. *Molecular systems biology*,
433 13(4):924, 2017.
- 434 [25] H. Kuwahara and X. Gao. Analysis of the effects of related fingerprints on molecular similarity
435 using an eigenvalue entropy approach. *Journal of Cheminformatics*, 13:1–12, 2021.
- 436 [26] G. Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry, and
437 predictive modeling. *Greg Landrum*, 8(31.10):5281, 2013.
- 438 [27] F. Lanusse, L. Parker, S. Golkar, M. Cranmer, A. Bietti, M. Eickenberg, G. Krawezik, M. Mc-
439 Cabe, R. Ohana, M. Pettee, B. R.-S. Blancard, T. Tesileanu, K. Cho, and S. Ho. Astroclip:
440 Cross-modal pre-training for astronomical foundation models, 2023.
- 441 [28] J. T. Leek, R. B. Scharpf, H. C. Bravo, D. Simcha, B. Langmead, W. E. Johnson, D. Geman,
442 K. Baggerly, and R. A. Irizarry. Tackling the widespread and critical impact of batch effects in
443 high-throughput data. *Nature Reviews Genetics*, 11(10):733–739, 2010.

- 444 [29] D. Masters, J. Dean, K. Klaser, Z. Li, S. Maddrell-Mander, A. Sanders, H. Helal, D. Beker,
445 A. Fitzgibbon, S. Huang, et al. Gps++: Reviving the art of message passing for molecular
446 property prediction. *arXiv preprint arXiv:2302.02947*, 2023.
- 447 [30] D. Mendez, A. Gaulton, A. P. Bento, J. Chambers, M. De Veij, E. Félix, M. P. Magariños, J. F.
448 Mosquera, P. Mutowo, M. Nowotka, et al. ChEMBL: towards direct deposition of bioassay data.
449 *Nucleic acids research*, 47(D1):D930–D940, 2019.
- 450 [31] O. Méndez-Lucio, C. Nicolaou, and B. Earnshaw. Mole: a molecular foundation model for
451 drug discovery, 2022.
- 452 [32] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding.
453 *arXiv preprint arXiv:1807.03748*, 2018.
- 454 [33] H. S. Parker and J. T. Leek. The practical effect of batch on genomic prediction. *Statistical
455 applications in genetics and molecular biology*, 11(3), 2012.
- 456 [34] B. Poole, S. Ozair, A. Van Den Oord, A. Alemi, and G. Tucker. On variational bounds of mutual
457 information. In *International Conference on Machine Learning*, pages 5171–5180. PMLR,
458 2019.
- 459 [35] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
460 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.
461 In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 462 [36] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are
463 unsupervised multitask learners. 2019.
- 464 [37] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M. Holzleit-
465 ner, M. Pavlović, G. K. Sandve, et al. Hopfield networks is all you need. *arXiv preprint
466 arXiv:2008.02217*, 2020.
- 467 [38] D. Rogers and M. Hahn. Extended-connectivity fingerprints. *Journal of chemical information
468 and modeling*, 50(5):742–754, 2010.
- 469 [39] Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang, and J. Huang. Self-supervised graph
470 transformer on large-scale molecular data, 2020.
- 471 [40] A. Sanchez-Fernandez, E. Rumetshofer, S. Hochreiter, and G. Klambauer. Cloome: contrastive
472 learning unlocks bioimaging databases for queries with chemical structures. *Nature*, 2023.
- 473 [41] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes,
474 A. Katta, C. Mullis, M. Wortsman, et al. Laion-5b: An open large-scale dataset for training
475 next generation image-text models. *Advances in Neural Information Processing Systems*,
476 35:25278–25294, 2022.
- 477 [42] J. Simm, G. Klambauer, A. Arany, M. Steijaert, J. K. Wegner, E. Gustin, V. Chupakhin, Y. T.
478 Chong, J. Vialard, P. Buijnsters, et al. Repurposing high-throughput image assays enables
479 biological activity prediction for drug discovery. *Cell chemical biology*, 25(5):611–618, 2018.
- 480 [43] K. Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Proceedings
481 of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, page
482 1857–1865, Red Hook, NY, USA, 2016. Curran Associates Inc.
- 483 [44] C. Sonesson, S. Gerster, and M. Delorenzi. Batch effect confounding leads to strong bias in
484 performance estimates obtained by cross-validation. *PloS one*, 9(6):e100335, 2014.
- 485 [45] R. S. Srinivasa, J. Cho, C. Yang, Y. M. Saidutta, C.-H. Lee, Y. Shen, and H. Jin. Cwcl: Cross-
486 modal transfer with continuously weighted contrastive loss. *Advances in Neural Information
487 Processing Systems*, 36, 2023.
- 488 [46] M. Sypetkowski, M. Rezanejad, S. Saberian, O. Kraus, J. Urbanik, J. Taylor, B. Mabey,
489 M. Victors, J. Yosinski, A. R. Sereshkeh, et al. Rxrx1: A dataset for evaluating experimental
490 batch correction methods. In *Proceedings of the IEEE/CVF Conference on Computer Vision
491 and Pattern Recognition*, pages 4284–4293, 2023.

- 492 [47] M. Sypetkowski, F. Wenkel, , F. Poursafaei, N. Dickson, K. Suri, P. Fradkin, and D. Beaini. On
493 the scalability of foundational models for molecular graphs. *arxiv*, 2024.
- 494 [48] F. Vincent, A. Nueda, J. Lee, M. Schenone, M. Prunotto, and M. Mercola. Phenotypic drug
495 discovery: recent successes, lessons learned and new directions. *Nature Reviews Drug Discovery*,
496 21(12):899–914, 2022.
- 497 [49] R. M. Walmsley and N. Billinton. How accurate is in vitro prediction of carcinogenicity?
498 *British Journal of Pharmacology*, 162(6):1250–1258, Feb. 2011.
- 499 [50] R. Xie, K. Pang, G. D. Bader, and B. Wang. Maester: Masked autoencoder guided segmentation
500 at pixel resolution for accurate, self-supervised subcellular structure recognition. In *2023*
501 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2023.
- 502 [51] S. Zaidi, M. Schaarschmidt, J. Martens, H. Kim, Y. W. Teh, A. Sanchez-Gonzalez, P. Battaglia,
503 R. Pascanu, and J. Godwin. Pre-training via denoising for molecular property prediction, 2022.
- 504 [52] X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer. Sigmoid loss for language image pre-
505 training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
506 11975–11986, 2023.
- 507 [53] X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyer. Lit:
508 Zero-shot transfer with locked-image text tuning, 2022.
- 509 [54] Y. Zhong, H. Tang, J. Chen, J. Peng, and Y.-X. Wang. Is self-supervised learning more robust
510 than supervised learning? *arXiv preprint arXiv:2206.05259*, 2022.
- 511 [55] G. Zhou, Z. Gao, Q. Ding, H. Zheng, H. Xu, Z. Wei, L. Zhang, and G. Ke. Uni-mol: A universal
512 3d molecular representation learning framework. In *The Eleventh International Conference on*
513 *Learning Representations*, 2023.

514 A NeurIPS Paper Checklist

515 1. Claims

516 Question: Do the main claims made in the abstract and introduction accurately reflect the
517 paper’s contributions and scope?

518 Answer: [Yes]

519 Justification: In the abstract, we claim that we build a multi-modal molecular-phenomics
520 model and demonstrate improvements over prior works. This is done by taking using a
521 uni-modal pre-trained phenomics model, tackling inactive molecules by undersampling and
522 learning inter-sample similarities. In addition, we take into account concentration in our
523 model training. We demonstrate comprehensive results supporting these claims.

524 Guidelines:

- 525 • The answer NA means that the abstract and introduction do not include the claims
526 made in the paper.
- 527 • The abstract and/or introduction should clearly state the claims made, including the
528 contributions made in the paper and important assumptions and limitations. A No or
529 NA answer to this question will not be perceived well by the reviewers.
- 530 • The claims made should match theoretical and experimental results, and reflect how
531 much the results can be expected to generalize to other settings.
- 532 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
533 are not attained by the paper.

534 2. Limitations

535 Question: Does the paper discuss the limitations of the work performed by the authors?

536 Answer: [Yes]

537 Justification: In the conclusion, we have a limitations subsection discussing future research
538 directions and assumptions in our work.

539 Guidelines:

- 540 • The answer NA means that the paper has no limitation while the answer No means that
541 the paper has limitations, but those are not discussed in the paper.
- 542 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 543 • The paper should point out any strong assumptions and how robust the results are to
544 violations of these assumptions (e.g., independence assumptions, noiseless settings,
545 model well-specification, asymptotic approximations only holding locally). The authors
546 should reflect on how these assumptions might be violated in practice and what the
547 implications would be.
- 548 • The authors should reflect on the scope of the claims made, e.g., if the approach was
549 only tested on a few datasets or with a few runs. In general, empirical results often
550 depend on implicit assumptions, which should be articulated.
- 551 • The authors should reflect on the factors that influence the performance of the approach.
552 For example, a facial recognition algorithm may perform poorly when image resolution
553 is low or images are taken in low lighting. Or a speech-to-text system might not be
554 used reliably to provide closed captions for online lectures because it fails to handle
555 technical jargon.
- 556 • The authors should discuss the computational efficiency of the proposed algorithms
557 and how they scale with dataset size.
- 558 • If applicable, the authors should discuss possible limitations of their approach to
559 address problems of privacy and fairness.
- 560 • While the authors might fear that complete honesty about limitations might be used by
561 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
562 limitations that aren’t acknowledged in the paper. The authors should use their best
563 judgment and recognize that individual actions in favor of transparency play an impor-
564 tant role in developing norms that preserve the integrity of the community. Reviewers
565 will be specifically instructed to not penalize honesty concerning limitations.

566 3. Theory Assumptions and Proofs

567 Question: For each theoretical result, does the paper provide the full set of assumptions and
568 a complete (and correct) proof?

569 Answer: [NA]

570 Justification: Our work does not contain proofs or theorems.

571 Guidelines:

- 572 • The answer NA means that the paper does not include theoretical results.
- 573 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
574 referenced.
- 575 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 576 • The proofs can either appear in the main paper or the supplemental material, but if
577 they appear in the supplemental material, the authors are encouraged to provide a short
578 proof sketch to provide intuition.
- 579 • Inversely, any informal proof provided in the core of the paper should be complemented
580 by formal proofs provided in appendix or supplemental material.
- 581 • Theorems and Lemmas that the proof relies upon should be properly referenced.

582 4. Experimental Result Reproducibility

583 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
584 perimental results of the paper to the extent that it affects the main claims and/or conclusions
585 of the paper (regardless of whether the code and data are provided or not)?

586 Answer: [Yes]

587 Justification: Our work documents our design decisions in detail and has comprehensive
588 details about the underlying dataset. We document all our hyperparameter choices and model
589 architectural decisions. Our evaluation is performed on a publicly accessible dataset RXX3,
590 allowing for benchmarking of other methods. To reproduce the pre-trained phenomics model,
591 we base our architecture on the work from [23], for which they have also provided access
592 to a snacker model, namely Phenom-Beta via a web platform hosted on the BioNeMo
593 platform <https://www.rxx.ai/phenom>. To reproduce the pre-trained molecular model,
594 we based our architecture on [47], for which the authors provide all the code and data
595 needed to reproduce it. We further note that the molecular model can be replaced by simple
596 molecular fingerprints with only a slight drop in performance.

597 Guidelines:

- 598 • The answer NA means that the paper does not include experiments.
- 599 • If the paper includes experiments, a No answer to this question will not be perceived
600 well by the reviewers: Making the paper reproducible is important, regardless of
601 whether the code and data are provided or not.
- 602 • If the contribution is a dataset and/or model, the authors should describe the steps taken
603 to make their results reproducible or verifiable.
- 604 • Depending on the contribution, reproducibility can be accomplished in various ways.
605 For example, if the contribution is a novel architecture, describing the architecture fully
606 might suffice, or if the contribution is a specific model and empirical evaluation, it may
607 be necessary to either make it possible for others to replicate the model with the same
608 dataset, or provide access to the model. In general, releasing code and data is often
609 one good way to accomplish this, but reproducibility can also be provided via detailed
610 instructions for how to replicate the results, access to a hosted model (e.g., in the case
611 of a large language model), releasing of a model checkpoint, or other means that are
612 appropriate to the research performed.
- 613 • While NeurIPS does not require releasing code, the conference does require all submis-
614 sions to provide some reasonable avenue for reproducibility, which may depend on the
615 nature of the contribution. For example
 - 616 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
617 to reproduce that algorithm.
 - 618 (b) If the contribution is primarily a new model architecture, the paper should describe
619 the architecture clearly and fully.

- 620 (c) If the contribution is a new model (e.g., a large language model), then there should
621 either be a way to access this model for reproducing the results or a way to reproduce
622 the model (e.g., with an open-source dataset or instructions for how to construct
623 the dataset).
- 624 (d) We recognize that reproducibility may be tricky in some cases, in which case
625 authors are welcome to describe the particular way they provide for reproducibility.
626 In the case of closed-source models, it may be that access to the model is limited in
627 some way (e.g., to registered users), but it should be possible for other researchers
628 to have some path to reproducing or verifying the results.

629 5. Open access to data and code

630 Question: Does the paper provide open access to the data and code, with sufficient instruc-
631 tions to faithfully reproduce the main experimental results, as described in supplemental
632 material?

633 Answer: [No]

634 Justification: As part of the submission, we are unable to provide code to reproduce model
635 training due to use of its proprietary nature. The training dataset is also an asset of a private
636 institution, meaning that we are unable to be made publicly accessible. The unseen dataset
637 RXX3 is, however, open source and can be used by the community to evaluate public
638 phenomics models.

639 Guidelines:

- 640 • The answer NA means that paper does not include experiments requiring code.
- 641 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
642 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 643 • While we encourage the release of code and data, we understand that this might not be
644 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
645 including code, unless this is central to the contribution (e.g., for a new open-source
646 benchmark).
- 647 • The instructions should contain the exact command and environment needed to run to
648 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
649 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 650 • The authors should provide instructions on data access and preparation, including how
651 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 652 • The authors should provide scripts to reproduce all experimental results for the new
653 proposed method and baselines. If only a subset of experiments are reproducible, they
654 should state which ones are omitted from the script and why.
- 655 • At submission time, to preserve anonymity, the authors should release anonymized
656 versions (if applicable).
- 657 • Providing as much information as possible in supplemental material (appended to the
658 paper) is recommended, but including URLs to data and code is permitted.

659 6. Experimental Setting/Details

660 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
661 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
662 results?

663 Answer: [Yes]

664 Justification: We provide details regarding our hyperparameter choices in the Appendix
665 C. In addition we document the use of scaffold splitting for Unseen Molecules & Images
666 dataset. Unseen Dataset RXX3 is publicly accessible.

667 Guidelines:

- 668 • The answer NA means that the paper does not include experiments.
- 669 • The experimental setting should be presented in the core of the paper to a level of detail
670 that is necessary to appreciate the results and make sense of them.
- 671 • The full details can be provided either with the code, in appendix, or as supplemental
672 material.

673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Our reported results are averaged over 3 random seeds used to initialize the model and dictating stochasticity during model training. We report most standard deviations in the main text, and the remaining ones are all present in the Appendix E.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide details on compute time for each experiment in Appendix D.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research described does not violate the NeurIPS Code of Ethics. Our experiments do not include human subjects, we follow fair use of data, privacy, and do not release model weights for mitigating impact measures.

Guidelines:

- 724
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - 725 • If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - 726
 - 727 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
 - 728

729 10. Broader Impacts

730 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

731 Answer: [Yes]

732 Justification: Our work discusses the potential in which MolPhenix can have positive societal impact and we touch on the extenralities in our concluding statements.

733 Guidelines:

- 734 • The answer NA means that there is no societal impact of the work performed.
- 735 • If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- 736 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 737 • The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- 738 • The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- 739 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757

758 11. Safeguards

759 Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

760 Answer: [NA]

761 Justification: In our work we do not release model weights or the underlying code.

762 Guidelines:

- 763 • The answer NA means that the paper poses no such risks.
- 764 • Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- 765 • Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 766 • We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774

775 12. Licenses for existing assets

776 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
777 the paper, properly credited and are the license and terms of use explicitly mentioned and
778 properly respected?

779 Answer: [Yes]

780 Justification: Assets used are referenced and licenses checked or otherwise not released
781 publicly.

782 Guidelines:

- 783 • The answer NA means that the paper does not use existing assets.
- 784 • The authors should cite the original paper that produced the code package or dataset.
- 785 • The authors should state which version of the asset is used and, if possible, include a
786 URL.
- 787 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 788 • For scraped data from a particular source (e.g., website), the copyright and terms of
789 service of that source should be provided.
- 790 • If assets are released, the license, copyright information, and terms of use in the
791 package should be provided. For popular datasets, `paperswithcode.com/datasets`
792 has curated licenses for some datasets. Their licensing guide can help determine the
793 license of a dataset.
- 794 • For existing datasets that are re-packaged, both the original license and the license of
795 the derived asset (if it has changed) should be provided.
- 796 • If this information is not available online, the authors are encouraged to reach out to
797 the asset's creators.

798 13. New Assets

799 Question: Are new assets introduced in the paper well documented and is the documentation
800 provided alongside the assets?

801 Answer: [NA]

802 Justification: The paper does not release new assets.

803 Guidelines:

- 804 • The answer NA means that the paper does not release new assets.
- 805 • Researchers should communicate the details of the dataset/code/model as part of their
806 submissions via structured templates. This includes details about training, license,
807 limitations, etc.
- 808 • The paper should discuss whether and how consent was obtained from people whose
809 asset is used.
- 810 • At submission time, remember to anonymize your assets (if applicable). You can either
811 create an anonymized URL or include an anonymized zip file.

812 14. Crowdsourcing and Research with Human Subjects

813 Question: For crowdsourcing experiments and research with human subjects, does the paper
814 include the full text of instructions given to participants and screenshots, if applicable, as
815 well as details about compensation (if any)?

816 Answer: [NA]

817 Justification: The paper does not involve crowdsourcing not human subject research.

818 Guidelines:

- 819 • The answer NA means that the paper does not involve crowdsourcing nor research with
820 human subjects.
- 821 • Including this information in the supplemental material is fine, but if the main contribu-
822 tion of the paper involves human subjects, then as much detail as possible should be
823 included in the main paper.
- 824 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
825 or other labor should be paid at least the minimum wage in the country of the data
826 collector.

827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing not human subject research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

846 B Assumption of the Initial Cell State

847 There is an important distinction between phenomics - molecule and text - image contrastive training
848 although there are initial similarities. In the text - image domain the two modalities are directly
849 generated by the same latent variable which is the underlying semantic class. Whereas in phenomics -
850 molecule, the observed phenomics variable is actually conditioned on molecular structure and the
851 initial state. There are two important conclusions from this: (1) This indicates that if molecular
852 structure has no effect on the initial cell state, there will not be a positive pairing between the
853 molecular structure and morphological patterns captured by phenomics, making it indistinguishable
854 from a control image. (2) There is an underlying assumption that the initial cell state x_i^0 is constant.
855 In accordance with this assumption we utilize experiments with a fixed cell line, *HUVEC-19*, and
856 a constant genetic background. Future works can relax this assumption by taking into account
857 phenomics experiments of the cells prior to the perturbation. This can allow the models to generalize
858 beyond a single cell line and to diverse genetic backgrounds.

859 C Dataset

860 Models have been trained using our in house training set and we have conducted our evaluation on
861 two novel datasets and an open-source molecule dataset [13]:

- 862 • **Training Set:** Our training dataset comprises 1,316,283 pairs of molecules and concentration
863 concentration combinations, complemented by fluorescent microscopy images generated through
864 over 2,150,000 phenomic experiments.
- 865 • **Evaluation set 1 - Unseen Images + Seen Molecules:** The first set consists of unseen
866 images and seen molecules. Unseen microscopy images are associated with 15,058 pairs of
867 molecules and concentrations from the training set and selected randomly.
- 868 • **Evaluation set 2 - Unseen Images + Unseen Molecules:** The second set includes previ-
869 ously unseen molecules, and images (consisting of 45,771 molecule and concentration pairs).
870 Predicting molecular identities of previously unseen molecular perturbations corresponds to zero-
871 shot prediction. Scaffold splitting was used to split this validation dataset from training ensuring
872 minimal information leakage.
- 873 • **Evaluation set 3 - Unseen Dataset:** Finally, we utilize the RXXR3 dataset [13], an open-
874 source out of distribution (OOD) dataset consisting of 6,549 novel molecule and concentration
875 pairs associated with phenomic experiments. The distribution of molecular structures differs from
876 previous datasets, making this a challenging zero-shot prediction task.

877 C.1 Concentration Details

878 Additional details regarding the number of molecules at significant concentrations of each evaluation
879 set are available in Table 6.

Table 6: Separated number of molecules for different concentrations at various pvalue cut-offs

Concentration	pvalue=1.0			pvalue=.1			pvalue=.01		
	Unseen Im.	Unseen Im. + Mol.	Unseen Data	Unseen Im.	Unseen Im. + Mol.	Unseen Data	Unseen Im.	Unseen Im. + Mol.	Unseen Data
.1	1497	1109	0	387	170	0	161	68	0
.25	1775	1111	1638	600	203	237	334	121	165
1.0	2721	11392	1639	1259	734	390	672	390	268
2.5	1787	4018	1636	1329	644	516	929	413	375
3.0	74	10454	0	12	1540	0	4	729	0
5.0	3	50	0	0	27	0	0	20	0
1.0	2712	11392	1636	2544	8117	792	2116	4815	625
25.0	0	2916	0	0	1734	0	0	950	0
Unique molecules	3026	14256	1639	2729	9857	823	2309	5778	642

880 D Implementation Details

881 In our experiments we report the top 1% recall metric as it is agnostic to the size of the dataset used.
882 Across different datasets, top 1 metric can correspond to varying levels of difficulty due to the number
883 of negatives evaluated. Top 1% can be used to compare models with different batch sizes, datasets,
884 and evaluations with different number of concentrations.

885 D.1 Hyperparameters

886 Our design choices and utilized hyperparameters for is presented in Table 7. We set batch size to 512
887 through experiments presented in top section of Table 1 and Figure 4 since training CLOOME model
888 on images is not efficient compared to using pretrained models. In addition, results presented at
889 bottom section of Table 1 are based on the best parameters found through described ablation studies
890 (section E.5).

Table 7: Hyperparameter values utilized in our proposed MolPhenix training framework.

Parameter	Value
number of seeds	3
learning rate	1e-3
weight decay	3e-3
optimizer	AdamW
training batch size	8192
validation batch size	12000
embedding dim	512
model size	medium (38.7 M)
model structure	6 ResNet Blocks + 1 Linear layer + 1 ResNet Block + 1 Linear layer
epochs	100
self similarity clip val	.75
learnable temperature initialization	2.302
learnable bias initialization	-1.0
Distance function	arctangent of l2 distance

891 D.2 Resource Computation

892 We utilized an NVIDIA A100 GPU to train Molphenix using Phenom1 and MolGPS embeddings,
893 which takes approximately ~ 4.75 hours each. For loss comparison experiments, we run each model
894 using 3 different seeds and 8 different losses, resulting in a total of 114 hours of GPU processing
895 time. For concentration experiments we train 7 runs, one for each concentration, with 3 seeds each
896 totaling 21 runs per set of parameters. With 25 sets of parameters evaluated (13), that amounts to
897 2,500 A100 compute hours. Moreover, we employed 8 NVIDIA A100 GPUs to train CLOOME
898 model on phenomics images, with an average of 40 hour usage per run. Across three seeds, that
899 amounts to ~ 1000 hours of A100 GPU usage (8 GPUs for 40 hours 3 times).

900 Note that, without accounting for the time to train Phenom1, MolPhenix is $8.4 \times$ faster than the
901 CLOOME baseline.

902 D.3 S2L Distance function

903 To calculate inter sample distances, we utilize arctangent of l2 distances between Phenom1 embed-
904 dings. More specifically, we calculate distances with

$$\arctan(\|z_{\mathbf{x}_i} - z_{\mathbf{x}_j}\|_2^2 / c) * \frac{4}{\pi} - 1, \quad (5)$$

905 where c is a constant indicating the median l2 distance between a null set of embeddings. Empirically,
906 we’ve found that setting similarities below a threshold k to 0 improves model performance: $\lceil w \rceil^k$.

907 Usage of arctan-l2 distances is motivated by an observation that cosine similarities do not effectively
908 separate inactive molecules from other molecular pairs (Figure 6). To alleviate inactive molecule
909 challenge, we require significant separation of CDF curves of inactive perturbations (p value $> .9$) and
910 active molecules (p $< .01$). We observe that in both the plots using arctangent and cosine similarities
911 achieves this purpose. However, if we compare high p-value curves with high-low, we find that
912 in the case of cosine similarities they are almost identical. This indicates that the distribution of

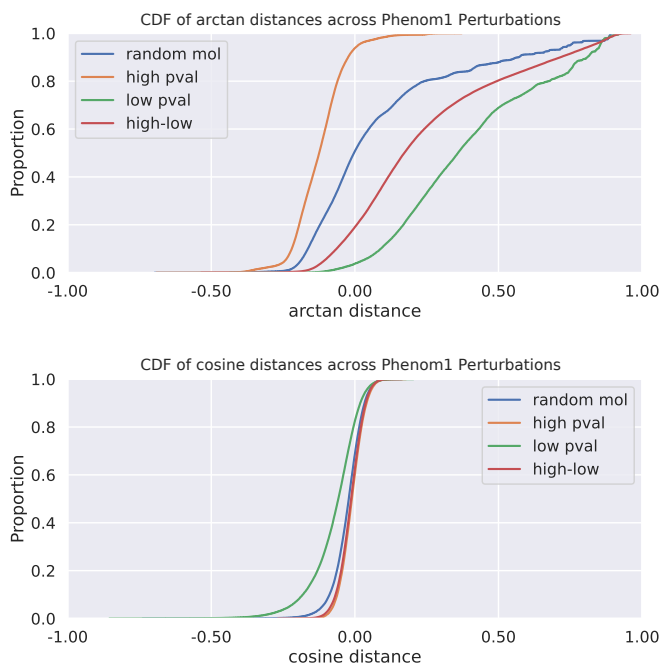


Figure 6: Plotted are cumulative densities of distance metrics for cosine similarity and arctangent of l_2 distances between embeddings. Random mol corresponds to Phenom1 distances between random molecules, high pval corresponds to distances between molecules with high p-values, low pval corresponds to distances between active molecules with low p-values, finally high-low corresponds to distances between active and inactive molecules.

913 cosine similarities between active - inactive molecules is almost identical to that of inactive - inactive
 914 molecules. In contrast, when using arctangent similarities, we observe that the two CDF curves are
 915 well separated.

916 This property of l_2 distances can inform our model training to identify inactive-inactive molecules.
 917 These results informed our decision to utilize arctangent of l_2 distances to specify sample similarities
 918 for the S2L loss.

919 E Additional Results

920 E.1 Predicting molecular activity

921 Given the significance of identifying active molecules, we evaluate the ability of the chemical encoder
 922 to predict molecular activity. To do so, we assessed whether embeddings generated from the chemical
 923 encoder can be used to predict calculated p-values for unseen molecules. We fit a logistic regression
 924 on molecular embeddings from the training set, classifying whether a molecular perturbation and
 925 concentration have a p-value below .01. We find that the trained logistic regression is capable of
 926 predicting molecular activity on two downstream datasets with a non-overlapping set of molecules,
 927 Figure 8. In addition, we provide a u-map of molecular embedding for the unseen dataset RXX3,
 928 colored by p-value. We qualitatively observe a clustering of active molecules using a U-map (Figure
 929 7). It demonstrates that predicting compounds activity is possible using MolPhenix chemical encoder
 930 as molecules representations are distinct, independent of the experimented dosage concentration.

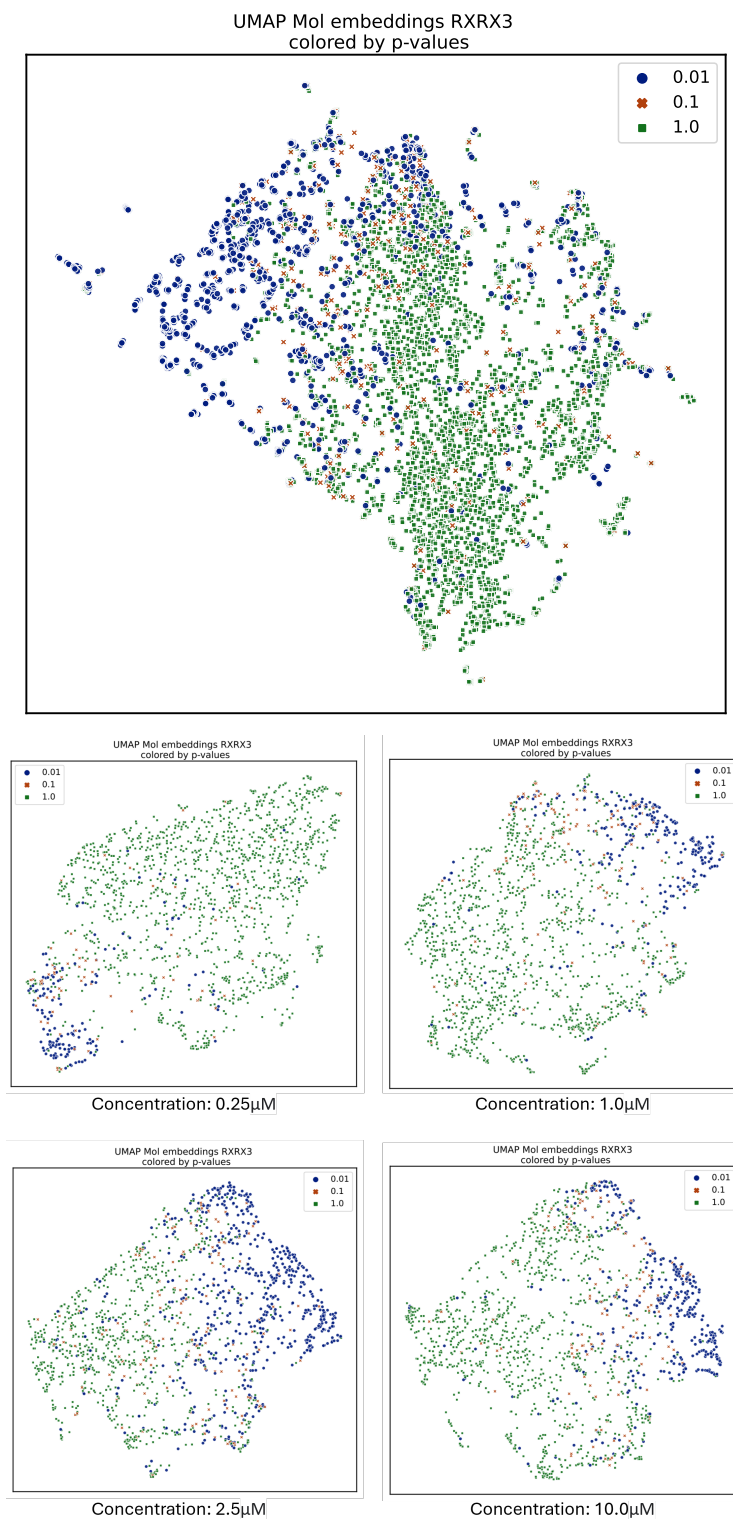


Figure 7: U-map demonstrating dimensionality reduction of the chemical embeddings of unseen dataset RXX3. First two dimensions are visualized and points are colored corresponding to their activity measured in phenomics experiments. Activity is evaluated using p-values calculated using technical replicability of Phenom1 embeddings. Top plot shows the u-map figure of all chemical embeddings, and bottom figure contains u-map figure of representations at specific concentrations.

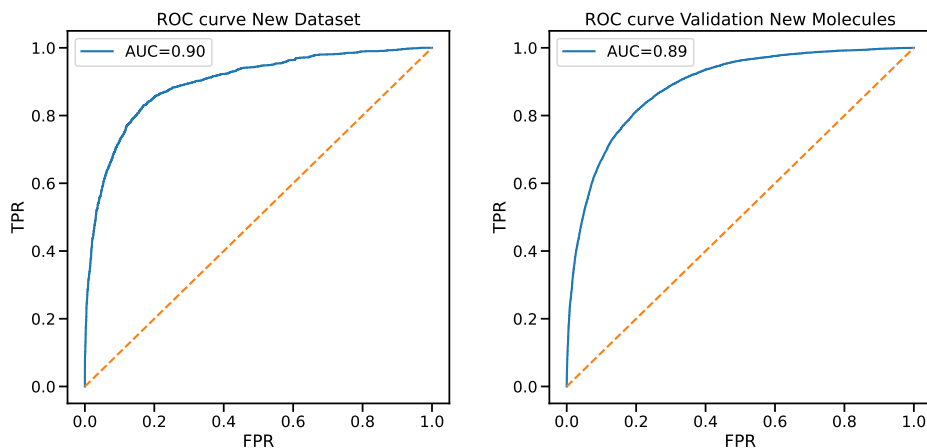


Figure 8: **Top left:** ROC AUC of logistic regression predicting molecular activity on new dataset. **Top right:** ROC AUC of logistic regression predicting molecular activity on validation dataset with new molecules and new images.

931 E.2 Zero Shot Biological Validation

932 We conduct a preliminary investigation into whether MolPhenix can be used to identify biological
 933 relationships without the need for conducting the underlying experiments. To this end, we evaluate
 934 on a subset of ChEMBL with curated pairs of gene knockouts and molecular perturbants [30]. These
 935 pairs of perturbations were curated due to the similarity of their effects on cells, although these might
 936 not always be captured through phenomic experiments. Thus, there is maximum performance that
 937 can be reached through just phenomic data, which we assume to be achieved by experimental data
 938 embedded using Phenom1.

939 To evaluate MolPhenix’s ability to identify previously known biological associations directly from
 940 data, we embed phenomics experiments from gene knockouts using the vision encoder. To perform
 941 in-silico screening, we then embed the molecular structures associated with positive pairs using the
 942 chemical encoder. Generating molecular embeddings and the corresponding concentrations does not
 943 utilize any experimental data. We then calculate cosine similarities between embeddings of phenomics
 944 experiments evaluating gene knockouts, and representations of the chemical representations along
 945 with encoded concentrations. Using the computed cosine similarities we are then able to assess
 946 whether MolPhenix is capable of identifying known associations between gene knockouts and
 947 molecular structures. Since there is no information on molecular concentration at which the cells
 948 must be treated with, we repeat the experiment across 4 concentrations. To get a null distribution of
 949 cosine similarities we take pairs of genes knockouts and molecules for which there are no annotated
 950 relationships. We calculate a cut-off for a low and high percentiles, and then evaluate what percentage
 951 of pairs of genes and molecules with known relationships exceed the set thresholds.

952 Figure 9 demonstrates that in-silico screening using MolPhenix Molecular encoder is capable of
 953 recovering a significant portion of known interactions. This is performed without the use of exper-
 954 imental data on the molecular encoder. It is difficult to estimate an upper bound on the expected
 955 performance due to uncertainty in the quality of curation of known pairs, presence of unknown
 956 associations between genes and molecules, and uncertainty regarding molecular concentration. There
 957 is a clear trend however that MolPhenix molecular encoder is capable of recovering a meaningful
 958 fraction of these interactions.

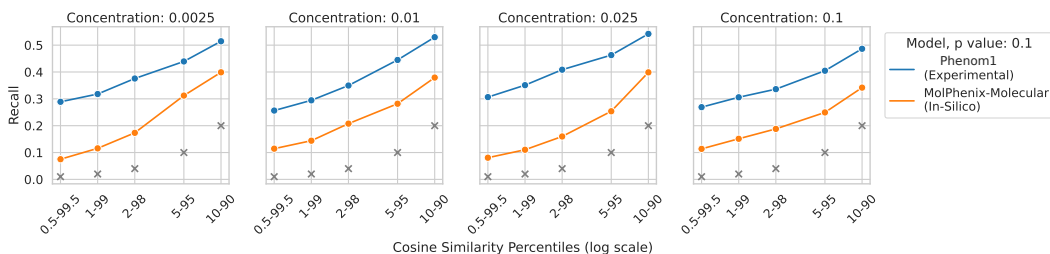


Figure 9: Evaluation of 0-shot ChEMBL identification of gene knockout and molecular phenomic similarities. On the X axis are percentile ranges, at which points the threshold is computed for cosine similarities. On the y axis is plotted total recall of recovered known interactions. Grey x plotted for each range indicate baseline recall. Orange line indicates MolPhenix-Molecular encoding of chemical compounds and MolPhenix-Vision for encoding gene knockout phenomics experiment. Blue line indicates Phenom1 encoding of phenomics experiments for both the molecular perturbation and gene knockouts. In-silico encoding of molecular perturbation, as well as the corresponding concentration, recovers a significant fraction of observed interactions.

959 E.3 Molecular Property Prediction

960 We expand our evaluation with additional experiments supporting the utility of MolPhenix beyond
 961 retrieval. We conduct a KNN evaluation of the MolPhenix latent space, assessing the learned
 962 embedding on 35 molecular property prediction tasks across the Polaris and TDC datasets (Table
 963 8 and 9). We find that MolPhenix trained with fingerprint embeddings consistently outperforms
 964 standalone input fingerprints, demonstrating that the MolPhenix latent space effectively clusters
 965 molecules according to their biological properties. We observed an interesting effect where prediction
 966 quality is positively correlated with implied dosage, indicating that MolPhenix learns dosage-specific
 967 effects. In addition, utilizing

Table 8: Comparison of a KNN applied on MolPhenix molecular embedding with **traditional fingerprints** on different tasks of TDC and Polaris datasets. Mean results for TDC, Polaris and together are available in the last three columns. Binary fingerprints use tanimoto similarity, while floating-point fingerprints use cosine similarity.

		concentration																																			TDC Standardized Mean			Polaris Standardized Mean		
		admc-fang-fcClint-1	admc-fang-RPPB-1	admc-fang-FERM-1	admc-fang-fcClint-1	admc-fang-RPPB-1	admc-fang-SOUL-1	ams	bbb_sarantis	bioavailability_ma	encod_wang	clearance_lepatocyte_az	clearance_mitrosome_az	cyt2d9_substrate_carbonmangels	cyt2d9_veth	cyt2d6_substrate_carbonmangels	cyt2d6_veth	cyt3a4_substrate_carbonmangels	cyt3a4_veth	dhl	half_life_obach	berg	bia_hou	ld50_zhu	lipophilicity_nitrazocin	ppp_brocateali	pk1a2-egf-wt-e-1	pk1a2-egf-wt-e-1	pk1a2-kt-wt-e-1	pk1a2-kt-wt-e-1	pk1a2-ret-wt-e-1	pk1a2-ret-wt-e-1	ppbr_az	solubility_asefalb	vdhs_lumbarido							
metric		pearson	pearson	pearson	pearson	pearson	pearson	mae	mae	mae	pearman	pearman	supr	supr	supr	supr	supr	supr	supr	pearman	pearman	mae	mae	supr	supr	pearson	pearson	pearson	pearson	pearson	pearson	mae	mae	pearman								
rdkit		0.32	0.34	0.48	0.23	0.38	0.29	0.69	0.72	0.58	0.54	0.25	0.45	0.31	0.45	0.45	0.29	0.54	0.59	0.71	0.26	0.61	0.71	-0.70	0.84	0.76	0.20	0.42	0.33	0.53	0.36	0.45	13.03	-1.63	0.22	2.62	-0.51	-1.88				
ecfp		0.46	0.60	0.49	0.43	0.60	0.39	0.69	0.75	0.48	-0.43	0.37	0.50	0.32	0.52	0.44	0.33	0.60	0.64	0.67	0.47	0.73	0.65	-0.73	-0.78	0.79	0.41	0.57	0.33	0.51	0.40	0.55	-9.91	-1.27	0.47	-1.96	0.03	-1.28				
maccs		0.37	0.56	0.52	0.22	0.43	0.44	0.71	0.77	0.53	-0.47	0.35	0.42	0.32	0.49	0.45	0.32	0.62	0.61	0.75	0.43	0.66	0.70	-0.66	-0.80	0.79	0.21	0.35	0.25	0.32	0.44	0.49	-10.13	-1.47	0.46	-1.91	0.45	-1.40				
Concatenated fps		0.41	0.66	0.58	0.33	0.40	0.37	0.70	0.77	0.58	-0.43	0.38	0.52	0.33	0.54	0.42	0.33	0.57	0.62	0.74	0.45	0.70	0.72	-0.67	-0.80	0.84	0.36	0.56	0.34	0.57	0.44	0.57	-10.94	-1.46	0.48	-1.78	0.00	-1.15				
Molphenix fingerprint	1	0.57	0.75	0.57	0.55	0.72	0.57	0.70	0.74	0.54	-0.48	0.29	0.46	0.32	0.57	0.47	0.38	0.59	0.64	0.77	0.55	0.67	0.69	-0.71	-0.70	0.80	0.20	0.41	0.30	0.43	0.31	0.39	-8.93	-1.10	0.55	-1.64	0.14	-0.21				
Molphenix fingerprint	25	0.64	0.71	0.65	0.62	0.67	0.58	0.69	0.78	0.54	-0.42	0.30	0.45	0.32	0.56	0.40	0.42	0.60	0.67	0.77	0.38	0.69	0.74	-0.67	-0.68	0.84	0.17	0.42	0.32	0.39	0.37	0.45	-8.43	-1.02	0.50	-1.40	0.28	-0.82				

Table 10: **Evaluation on cumulative concentrations for active molecules:** Average Top-1% and Top-5% recall accuracies of methods utilizing different contrastive learning loss functions and concentration encoding information. We evaluate all methods on unseen images, unseen images and unseen molecules and an unseen dataset for zero-shot retrieval. Entries in **bold** denote best performance when the loss function is fixed while entries in **highlight** denote best performance across all guidelines.

Method	Explicit Concentration (ours)	Modality	top-1%				top-5%			
			Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
CLIP	\times	Phenom1	.3373	.4228	.1514	3038	.6162	.7182	.3660	.5668
Hopfield-CLIP	\times	Phenom1	.2578	.3559	.1256	2464	.5457	.6751	.3270	.5159
InfoLOOB	\times	Phenom1	.3351	.4206	.1563	3040	.6128	.7204	.3730	.5687
CLOOME	\times	Phenom1	.3572	.4348	.1658	3193	.6330	.7259	.3918	.5836
CLOOME	sigmoid	Phenom1	.5813	.4968	.2360	4380	.8748	.7658	.4859	.7088
CLOOME	logarithm	Phenom1	.6057	.5255	.2445	.4586	.8858	.8117	.4957	.7310
CLOOME	one-hot	Phenom1	.5967	.5255	.2380	4534	.8800	.8120	.4829	.7250
DCL	\times	Phenom1	.6363	.6177	.3184	5241	.8638	.8180	.5632	.7483
DCL	sigmoid	Phenom1	.8858	.6694	.4527	6693	.9600	.8472	.6845	.8305
DCL	logarithm	Phenom1	.8934	.6952	.4511	6799	.9581	.8788	.6889	.8419
DCL	one-hot	Phenom1	.8901	.7002	.4601	.6834	.9591	.8770	.6873	.8411
CWCL	\times	Phenom1	.7091	.6529	.3556	5725	.9018	.8368	.6027	.7804
CWCL	sigmoid	Phenom1	.9138	.6985	.4810	6977	.9681	.8643	.7070	.8464
CWCL	logarithm	Phenom1	.9141	.7248	.4815	7068	.9651	.8920	.7131	.8567
CWCL	one-hot	Phenom1	.9128	.7261	.4850	.7079	.9665	.8927	.6998	.8530
SigLip	\times	Phenom1	.7763	.6401	.3396	5853	.9361	.8308	.5714	.7792
SigLip	sigmoid	Phenom1	.9463	.6911	.4576	6990	.9816	.8606	.6759	.8393
SigLip	logarithm	Phenom1	.9493	.7256	.4868	.7205	.9814	.8926	.7019	.8586
SigLip	one-hot	Phenom1	.9489	.7210	.4859	.7186	.9823	.8868	.7045	.8578
MolPhenix (ours)	\times	Phenom1	.9097	.6759	.4181	.6679	.9768	.8539	.6436	.8247
MolPhenix (ours)	sigmoid	Phenom1	.9423	.7155	.4573	.7050	.9808	.8775	.6778	.8453
MolPhenix (ours)	logarithm	Phenom1	.9426	.7451	.4727	.7201	.9808	.8964	.6952	.8574
MolPhenix (ours)	one-hot	Phenom1	.9430	.7490	.4850	.7256	.9816	.8984	.7040	.8613
MolPhenix (ours)	\times	Phenom1 + MolGPS	.9105	.6710	.4501	.6772	.9755	.8527	.7098	.8460
MolPhenix (ours)	sigmoid	Phenom1 + MolGPS	.9395	.7034	.5252	.7227	.9811	.8729	.7630	.8723
MolPhenix (ours)	logarithm	Phenom1 + MolGPS	.9413	.7505	.5473	.7463	.9811	.9085	.7878	.8924
MolPhenix (ours)	one-hot	Phenom1 + MolGPS	.9430	.7514	.5577	.7507	.9830	.9043	.7821	.8898

Table 9: Comparison of a KNN applied on MolPhenix molecular embedding with **MolGPS** on different tasks of TDC and Polaris datasets. Mean results for TDC, Polaris and together are available in the last three columns.

metric	concentration																		TDC Standardized Mean	Polaris Standardized Mean	Standardized Mean																	
	adme-fang-HCLint-1	adme-fang-HPPE-1	adme-fang-PERM-1	adme-fang-RCLint-1	adme-fang-RPPB-1	adme-fang-SOLL-1	anes	bbb_martins	bioavailability_ma	en02_wmg	clearance_bpatoctoye_ax	clearance_microsome_ax	cydatl_substrate_carbonmangels	cydatl_veith	cydatl_substrate_carbonmangels	cydatl_veith	cydatl_substrate_carbonmangels	cydatl_veith																				
MolGPS	0.54	0.66	0.70	0.56	0.64	0.55	0.69	0.76	0.49	-0.50	0.40	0.57	0.30	0.62	0.50	0.41	0.66	0.68	0.81	0.52	0.70	0.74	-0.69	-0.71	0.84	0.34	0.51	0.44	0.55	0.30	0.48	-0.71	0.98	0.63	-1.19	0.38	-0.65	
MolPhenix with Molgps	1.00	0.78	0.69	0.61	0.68	0.65	0.70	0.79	0.59	-0.49	0.36	0.51	0.29	0.62	0.55	0.42	0.58	0.67	0.72	0.45	0.74	0.79	-0.71	-0.65	0.83	0.14	0.33	0.34	0.44	0.32	0.42	-8.28	-1.00	0.63	-1.23	0.29	-0.69	
MolPhenix with Molgps	25	0.68	0.74	0.70	0.67	0.77	0.63	0.71	0.78	0.60	-0.47	0.38	0.53	0.33	0.62	0.50	0.43	0.66	0.67	0.70	0.40	0.73	0.85	-0.70	-0.62	0.84	0.12	0.29	0.41	0.45	0.29	0.43	-8.46	-0.97	0.62	-0.62	0.38	-0.46

968 E.4 Addressing Challenges in Contrastive Phenomic Retrieval

969 Table 10 and 12 show the complete Top 1% and 5% results of evaluation on cumulative concentrations
 970 on only active and all molecules, respectively. Similarly, Table 11 and 13 presents the full retrieval
 971 results of held-out concentrations experiments. In comparison to prior loss functions, S2L loss
 972 objective demonstrates consistent high retrieval rate in all tasks and molecular groups (i.e. all or active
 973 molecules), while using the same modality (Phenom1) and with or without explicit concentration
 974 information.

975 E.5 Ablation Studies

976 Figure 10 and Table 15, 16, 17, 18 and 19 present top-1% recall accuracy across for the full ablation
 977 study on the variation of MolPhenix key components. We note that compact embedding sizes from
 978 pretrained models stabilize training. This indicates that embeddings are expressive and accurately
 979 capture intricate aspects of molecules. Larger batch sizes result in a greater number of negative
 980 samples, hence improving performance. This is in line with prior contrastive learning methods
 981 continuing to improve by increasing the batch size [10]. Increasing the number of parameters leads
 982 to more expressive models thereby enhancing retrieval performance. This result is in accordance with

Table 11: **Evaluation on held-out concentration for active molecules:** Average Top-1% and Top-5% recall accuracies of methods utilizing different contrastive learning loss functions and concentration encoding information. We evaluate all methods on unseen images, unseen images and unseen molecules and an unseen dataset for zero-shot retrieval. Entries in **bold** denote highest performance when the loss function is fixed while entries in **highlight** denote highest performance across all guidelines.

Method	Explicit Concentration (ours)	Modality	top-1%				top-5%			
			Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
CLIP	\times	Phenom1	.2109	.2425	.1519	.2018	.4458	.4968	.3591	.4339
Hopfield-CLIP	\times	Phenom1	.1581	.2034	.1198	.1604	.3783	.4413	.3045	.3747
InfoLOOB	\times	Phenom1	.2122	.2496	.1501	.2040	.4443	.5003	.3515	.4320
CLOOME	\times	Phenom1	.2164	.2461	.1479	.2035	.4590	.4956	.3528	.4358
CLOOME	sigmoid	Phenom1	.3338	.2681	.1801	.2606	.6037	.5202	.3879	.5039
CLOOME	logarithm	Phenom1	.3094	.2345	.1665	.2368	.5960	.4874	.3534	.4790
CLOOME	one-hot	Phenom1	.3073	.2040	.1670	.2261	.5997	.4246	.3657	.4633
DCL	\times	Phenom1	.4717	.4027	.2841	.3861	.7352	.6579	.5322	.6417
DCL	sigmoid	Phenom1	.7282	.4100	.3560	.4980	.9226	.6561	.6015	.7267
DCL	logarithm	Phenom1	.6903	.3558	.3211	.4557	.8869	.6146	.5667	.6894
DCL	one-hot	Phenom1	.6562	.3607	.3272	.4480	.8831	.5983	.5659	.6824
CWCL	\times	Phenom1	.5731	.4403	.3232	.4455	.8218	.6833	.5706	.6919
CWCL	sigmoid	Phenom1	.7780	.4425	.3777	.5327	.9386	.6844	.6244	.7491
CWCL	logarithm	Phenom1	.7452	.3989	.3523	.4988	.9117	.6482	.5962	.7187
CWCL	one-hot	Phenom1	.7048	.4009	.3593	.4883	.9037	.6284	.6061	.7127
SigLip	\times	Phenom1	.5718	.4217	.3021	.4318	.8104	.6602	.5176	.6627
SigLip	sigmoid	Phenom1	.8366	.4640	.3830	.5612	.9623	.7023	.6080	.7575
SigLip	logarithm	Phenom1	.8097	.4391	.3747	.5411	.9437	.6746	.6046	.7409
SigLip	one-hot	Phenom1	.7561	.4020	.3345	.4975	.9279	.6248	.5557	.7028
MolPhenix (ours)	\times	Phenom1	.8334	.4615	.3792	.5580	.9638	.6937	.6128	.7567
MolPhenix (ours)	sigmoid	Phenom1	.8256	.4692	.3765	.5571	.9638	.7068	.6115	.7607
MolPhenix (ours)	logarithm	Phenom1	.7953	.4466	.3664	.5361	.9466	.6889	.5924	.7426
MolPhenix (ours)	one-hot	Phenom1	.7489	.4088	.3379	.4985	.9325	.6465	.5644	.7144
MolPhenix (ours)	\times	Phenom1 & MolGPS	.8277	.4739	.4071	.5695	.9602	.7041	.6798	.7813
MolPhenix (ours)	sigmoid	Phenom1 & MolGPS	.8218	.4771	.4287	.5758	.9588	.7117	.7045	.7916
MolPhenix (ours)	logarithm	Phenom1 & MolGPS	.7836	.4757	.4297	.563	.9402	.7138	.7011	.7850
MolPhenix (ours)	one-hot	Phenom1 & MolGPS	.7391	.4307	.3940	.5212	.9198	.6724	.6698	.7540

Table 12: **Evaluation on cumulative concentrations for active and inactive perturbations** Average Top-1% and Top-5% Recall accuracy of methods utilizing different contrastive learning methods. Best performing methods are highlighted in **bold**.

Loss	Explicit Concentration	Modality	top-1%				top-5%			
			Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
CLIP	\times	Phenom1	.1761	.1867	.0734	.1454	.3710	.3769	.2065	.3181
Hopfield-CLIP	\times	Phenom1	.1531	.1709	.0673	.1304	.3464	.3637	.1942	.3014
InfoLOOB	\times	Phenom1	.1746	.1860	.0745	.1450	.3697	.3756	.2058	.3170
CLOOME	\times	Phenom1	.1968	.2005	.0911	.1628	.3938	.3888	.2321	.3383
CLOOME	sigmoid	Phenom1	.3875	.2592	.1415	.2627	.5662	.4601	.2940	.4401
CLOOME	logarithm	Phenom1	.4088	.3046	.1503	.2879	.5730	.5166	.3053	.4650
CLOOME	one-hot	Phenom1	.4080	.3123	.1496	.2900	.5801	.5306	.3054	.4720
DCL	\times	Phenom1	.3277	.2562	.1364	.2401	.4856	.4170	.2768	.3931
DCL	sigmoid	Phenom1	.4881	.3380	.2009	.3423	.6222	.5186	.3581	.4930
DCL	logarithm	Phenom1	.4983	.3615	.2122	.3573	.6311	.5581	.3587	.5160
DCL	one-hot	Phenom1	.5226	.3790	.2288	.3768	.6791	.5870	.3968	.5543
CWCL	\times	Phenom1	.3635	.2696	.1526	.2619	.5122	.4267	.2933	.4107
CWCL	sigmoid	Phenom1	.5070	.3457	.2101	.3542	.6378	.5272	.3462	.5037
CWCL	logarithm	Phenom1	.5146	.3725	.2246	.3706	.6437	.5733	.3660	.5277
CWCL	one-hot	Phenom1	.5401	.3849	.2336	.3862	.6882	.5991	.4001	.5625
SigLip	\times	Phenom1	.3729	.2544	.1470	.2581	.5200	.4179	.2838	.4072
SigLip	sigmoid	Phenom1	.5021	.3275	.2072	.3456	.6360	.5231	.3430	.5007
SigLip	logarithm	Phenom1	.5156	.3636	.2233	.3675	.6452	.5689	.3653	.5265
SigLip	one-hot	Phenom1	.5354	.3745	.2317	.3805	.6858	.5928	.3945	.5577
S2L (ours)	\times	Phenom1	.4688	.2852	.1838	.3126	.5970	.4519	.3171	.4554
S2L (ours)	sigmoid	Phenom1	.5071	.3441	.2144	.3552	.6428	.5315	.3554	.5099
S2L (ours)	logarithm	Phenom1	.5183	.3700	.2275	.3720	.6492	.5650	.3756	.5300
S2L (ours)	one-hot	Phenom1	.5433	.3819	.2384	.3879	.6954	.5895	.4030	.5626
S2L (ours)	\times	Phenom1 & MolGPS	.4688	.2729	.2001	.3139	.5956	.4374	.3430	.4587
S2L (ours)	sigmoid	Phenom1 & MolGPS	.4983	.3230	.2397	.3537	.6343	.5035	.3790	.5056
S2L (ours)	logarithm	Phenom1 & MolGPS	.5101	.3589	.2535	.3742	.6398	.5660	.3992	.5350
S2L (ours)	one-hot	Phenom1 & MolGPS	.5370	.3720	.2676	.3922	.6870	.5888	.4326	.5695

Table 13: **Evaluation on held-out concentrations for active and inactive perturbations** Average Top-1% and Top-5% Recall accuracy of methods utilizing different contrastive learning methods. Best performing methods are highlighted in **bold**.

Loss	Explicit Concentration	Modality	top-1%				top-5%			
			Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
CLIP	\times	Phenom1	.1684	.1111	.0964	.1253	.3916	.2545	.2356	.2476
Hopfield-CLIP	\times	Phenom1	.1290	.0921	.0756	.0989	.3485	.2287	.2095	.2217
InfoLOOB	\times	Phenom1	.1715	.1114	.0948	.1259	.3944	.2578	.2349	.2495
CLOOME	\times	Phenom1	.1745	.1088	.0910	.1248	.4093	.2487	.2355	.2439
CLOOME	sigmoid	Phenom1	.2573	.1208	.1062	.1614	.5169	.2638	.2513	.3440
CLOOME	logarithm	Phenom1	.2379	.1081	.0992	.1484	.4958	.2444	.2324	.3242
CLOOME	one-hot	Phenom1	.2346	.0970	.0974	.1430	.5014	.2224	.2348	.3195
DCL	\times	Phenom1	.3516	.1655	.1533	.2235	.5693	.3125	.3006	.3082
DCL	sigmoid	Phenom1	.4741	.1725	.1726	.2731	.6637	.3261	.3105	.3204
DCL	logarithm	Phenom1	.4286	.1596	.1581	.2488	.6244	.3071	.3032	.3056
DCL	one-hot	Phenom1	.4308	.1495	.1600	.2468	.6244	.2938	.3015	.2966
CWCL	\times	Phenom1	.4126	.1801	.1667	.2531	.6128	.3266	.3066	.3194
CWCL	sigmoid	Phenom1	.5112	.1856	.1811	.2926	.6901	.3384	.3190	.3314
CWCL	logarithm	Phenom1	.4664	.1696	.1709	.2690	.6502	.3195	.3066	.3148
CWCL	one-hot	Phenom1	.4681	.1612	.1734	.2676	.6465	.3019	.3104	.3050
SigLip	\times	Phenom1	.3942	.1578	.1390	.2303	.5931	.3015	.2737	.2914
SigLip	sigmoid	Phenom1	.5392	.1828	.1710	.2977	.7102	.3399	.3121	.3298
SigLip	logarithm	Phenom1	.5022	.1698	.1669	.2796	.6841	.3240	.3068	.3177
SigLip	one-hot	Phenom1	.4657	.1443	.1451	.2517	.6544	.2879	.2790	.2847
S2L (ours)	\times	Phenom1	.5336	.1842	.1713	.2963	.6961	.3322	.3045	.3221
S2L (ours)	sigmoid	Phenom1	.5409	.1899	.1753	.3020	.7178	.3469	.3201	.3372
S2L (ours)	logarithm	Phenom1	.5036	.1791	.1727	.2851	.6925	.3342	.3157	.3275
S2L (ours)	one-hot	Phenom1	.4726	.1537	.1521	.2595	.6696	.2998	.2887	.2958
S2L (ours)	\times	Phenom1 & MolGPS	.5248	.1829	.1910	.2996	.6904	.3268	.3305	.3281
S2L (ours)	sigmoid	Phenom1 & MolGPS	.5338	.1897	.2029	.3088	.7098	.3427	.3495	.3452
S2L (ours)	logarithm	Phenom1 & MolGPS	.4900	.1839	.2031	.2923	.6776	.3354	.3511	.3411
S2L (ours)	one-hot	Phenom1 & MolGPS	.4622	.1569	.1762	.2651	.6578	.3030	.3187	.3087

983 recent advances in language modelling and scaling laws across different data and compute budgets
 984 [21].

Model size	Depth	Width	Unseen images	Unseen images + Unseen molecules	Unseen dataset (0-shot)
Tiny - 2.7m	4 ResBlocks	256	.8337	.7186	.4030
Small - 9.4m	6 ResBlocks	512	.9174	.7352	.4562
Medium - 38.7m	8 ResBlocks	1024	.9430	.7490	.485

Table 14: Ablations across different model sizes. Larger capacity models are found to be more expressive.

Batch size	Unseen images	Unseen images + Unseen molecules	Unseen dataset (0-shot)
128	.8600	.7163	.4044
512	.9252	.7511	.4657
2048	.9450	.7616	.4940
8192	.9489	.7563	.4966

Table 15: Ablation across different batch sizes. Larger batch sizes benefit contrastive learning.

985 E.6 Investigating Other Pre-trained Phenomic Encoders

986 To investigate the impact of pre-trained encoders, we perform additional experiments evaluating a
 987 supervised phenomic image encoder (Table 20). Instead of Phenom1, we trained Molphenix frame-
 988 work using AdaBN, a CNN-based supervised phenomic encoder, with an analogous implementation
 989 discussed in [46]. We find that the general trends between Phenom1 and AdaBN are consistent with a
 990 slight decrease in overall performance. These findings provide additional support to the generality of
 991 the proposed guidelines.

992 E.7 Integrating MolGPS Embeddings With Other Fingerprints

993 Molphenix architecture is flexible, allowing that the proposed components be replaced by other
 994 phenomic or molecular pretrained models. We leveraged from MolGPS, which is a MPNN based

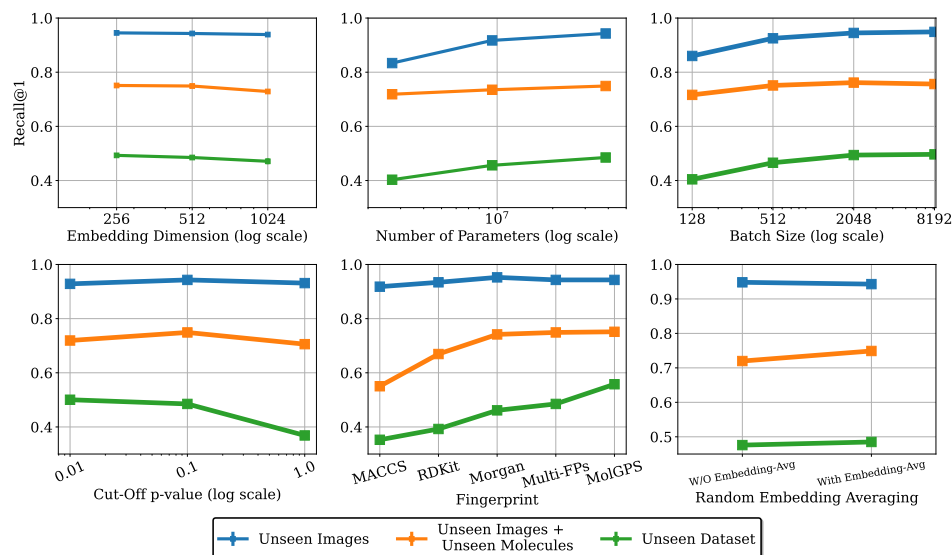


Figure 10: Ablations of top-1 % recall accuracy with **(top-left)** the size of embedding dimension, **(top-center)** number of parameters, **(top-right)** batch size, **(bottom-left)** cutoff p value, **(bottom-center)** fingerprint type, and **(bottom-right)** random batch averaging. Compact embedding sizes from pretrained models, larger number of parameters, larger batch sizes, lower cutoff p -values, pretrained MolGPS fingerprints and presence of random batch averaging improving retrieval of our MolPhenix framework.

Dim size	Unseen images	Unseen images + Unseen molecules	Unseen dataset (0-shot)
256	.9452	.7510	.4929
512	.9430	.7490	.4850
1024	.9392	.7288	.4710

Table 16: Ablation across different embedding dimensions. Compact embedding sizes capture more molecular information.

995 GNN model with 1B parameters which allows us to maximize architecture expressivity while
 996 minimizing the risk of overfitting [29, 47]. For additional investigation, we combine MolGPS
 997 molecular embeddings with RDKit, MACCS, and Morgan fingerprints and show that they can
 998 provide Molphenix with richer molecular information and yields overall higher performance of
 999 MolPhenix in both cumulative and held-out concentration scenarios. Results for active and all
 1000 molecules retrieval of Molphenix trained on the discussed combinational molecular embeddings are
 1001 available in table 21 and 22.

cut-off	Unseen images	Unseen images + Unseen molecules	Unseen dataset (0-shot)
p < 1.0	.9312	.7057	.3686
p < .1	.9430	.7490	.4850
p < .01	.9284	.7192	.5005

Table 17: Ablation across different p-value cutoff thresholds. p values < .1 benefit retrieval of active molecules.

fingerprint	unseen images	unseen images + unseen molecule	unseen dataset
MACCS	.9180	.5503	.3526
RDKit	.9341	.6693	.3925
Morgan	.9524	.7417	.4613
Multi-FPs	.9430	.7490	.485
Phenom1 + MolGPS	.9430	.7514	.5577

Table 18: Ablation across different fingerprint types. A combination of embeddings bootstrapped from Phenom1 and MolGPS significantly benefit retrieval.

	Unseen images	Unseen images + Unseen molecules	Unseen dataset (0-shot)
W/O Random Embedding Avg.	.9482	.7198	.4759
With Random Embedding Avg.	.9430	.7490	.485

Table 19: Ablation across random embedding averaging. Utilizing random batch averaging stabilizes training and benefits retrieval.

Table 20: Evaluation on **cumulative concentrations** while using **AdaBN**. Molphenix is trained on combination of RDKIT, MACCS, and Morgan fingerprints in this experiment

Method	Explicit Concentration	Modality	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
			top-1% active molecules				top-5% active molecules			
MolPhenix	-	AdaBN	.8568	.5336	.3525	.581	.9562	.7603	.5772	.7646
MolPhenix	sigmoid	AdaBN	.911	.5858	.4	.6323	.971	.7997	.6203	.797
MolPhenix	logarithm	AdaBN	.9155	.6106	.4242	.6501	.9729	.8332	.6503	.8188
MolPhenix	one-hot	AdaBN	.9187	.6125	.4225	.6512	.9744	.8302	.6419	.8155
			top-1% all molecules				top-5% all molecules			
MolPhenix	-	AdaBN	.4593	.2409	.1599	.2867	.5983	.4081	.285	.4305
MolPhenix	sigmoid	AdaBN	.5104	.3142	.1957	.3401	.6496	.5165	.331	.499
MolPhenix	logarithm	AdaBN	.5379	.3393	.2071	.3614	.6867	.5561	.3606	.5345
MolPhenix	one-hot	AdaBN	.5476	.3425	.2082	.3661	.7007	.5641	.3603	.5417

Table 21: Evaluation on **cumulative concentrations** while **combining MolGPS, RDKIT, MACCS, and Morgan fingerprints**.

Method	Explicit Concentration	Modality	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
			top-1% active molecules				top-5% active molecules			
MolPhenix	-	Phenom1 & MolGPS & 3 fps	.9185	.7212	.4717	.7038	.9784	.8805	.718	.859
MolPhenix	sigmoid	Phenom1 & MolGPS & 3 fps	.9395	.7408	.5119	.7307	.9817	.8932	.7458	.8736
MolPhenix	logarithm	Phenom1 & MolGPS & 3 fps	.9454	.7798	.5658	.7637	.9815	.9163	.7849	.8942
MolPhenix	one-hot	Phenom1 & MolGPS & 3 fps	.9419	.7687	.5526	.7544	.9807	.9113	.7681	.8867
			top-1% all molecules				top-5% all molecules			
MolPhenix	-	Phenom1 & MolGPS & 3 fps	.4764	.3011	.2068	.3281	.604	.4647	.3415	.4701
MolPhenix	sigmoid	Phenom1 & MolGPS & 3 fps	.5076	.342	.2382	.3626	.6383	.521	.3769	.512
MolPhenix	logarithm	Phenom1 & MolGPS & 3 fps	.525	.379	.2648	.3896	.658	.5743	.411	.5478
MolPhenix	one-hot	Phenom1 & MolGPS & 3 fps	.5355	.3845	.265	.395	.6862	.5916	.4233	.567

Table 22: Evaluation on **heldout concentrations** while **combining MolGPS, RDKIT, MACCS, and Morgan fingerprints**.

Method	Explicit Concentration	Modality	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.	Unseen Images	Unseen Images + Unseen Molecules	Unseen Dataset (0-shot)	Avg.
			top-1% active molecules				top-5% active molecules			
MolPhenix	-	Phenom1 & MolGPS & 3 fps	.8364	.5115	.4263	.5914	.9640	.7363	.6850	.7951
MolPhenix	sigmoid	Phenom1 & MolGPS & 3 fps	.8300	.5021	.4363	.5895	.9640	.7409	.6931	.7993
MolPhenix	logarithm	Phenom1 & MolGPS & 3 fps	.8112	.5107	.4376	.5865	.9544	.7406	.6866	.7939
MolPhenix	one-hot	Phenom1 & MolGPS & 3 fps	.7467	.4409	.3830	.5235	.9320	.6827	.6520	.7556
			top-1% all molecules				top-5% all molecules			
MolPhenix	-	Phenom1 & MolGPS & 3 fps	.5339	.1980	.1966	.3095	.6968	.2909	.4274	.4717
MolPhenix	sigmoid	Phenom1 & MolGPS & 3 fps	.5463	.2026	.2066	.3185	.7179	.3116	.4359	.4885
MolPhenix	logarithm	Phenom1 & MolGPS & 3 fps	.5247	.2009	.2078	.3111	.7067	.3133	.4319	.4840
MolPhenix	one-hot	Phenom1 & MolGPS & 3 fps	.4690	.1653	.1756	.2700	.6635	.2592	.4118	.4448