VITALLY CONSISTENT: SCALING BIOLOGICAL REP RESENTATION LEARNING FOR CELL MICROSCOPY

Anonymous authors

Paper under double-blind review

ABSTRACT

Large-scale cell microscopy screens are used in drug discovery and molecular biology research to study the effects of millions of chemical and genetic perturbations on cells. To use these images in downstream analysis, we need models that can map each image into a feature space that represents diverse biological phenotypes *consistently*, in the sense that perturbations with similar biological effects have similar representations. In this work, we present the largest foundation model for cell microscopy data to our knowledge, a new 1.9 billion-parameter ViT-G/8 MAE trained on over 8 billion microscopy image crops. Compared to a previous published ViT-L/8 MAE, our new model achieves a 60% improvement in linear separability of genetic perturbations and obtains the best overall performance on whole-genome biological relationship recall and replicate consistency benchmarks. We also show these performance trends hold on a public benchmark for measuring compound activity against target genes. Beyond scaling, we developed two key methods that improve performance: (1) training on a curated and diverse dataset; and, (2) using biologically motivated linear probing tasks to search across each transformer block for the best candidate representation of whole-genome screens. We find that many self-supervised vision transformers, pretrained on either natural or microscopy images, yield significantly more biologically meaningful representations of microscopy images in their intermediate blocks than in their typically used final blocks, therefore enabling significant cost and energy savings when deploying these large models in real-world applications. More broadly, our approach and results provide insights toward a general strategy for successfully building foundation models for large-scale biological image data.

032 033 034

004

010 011

012

013

014

015

016

017

018

019

020

021

024

025

026

027

028

029

031

1 INTRODUCTION

Large-scale cell microscopy assays are used to discover previously unknown biological processes 037 (Przybyla & Gilbert, 2022; Bock et al., 2022; Rood et al., 2024) and identify novel drug candi-038 dates and targets (Vincent et al., 2022). Labs are now able to achieve extremely high throughput by leveraging high content screening (HCS) systems that combine automated microscopy with robotic liquid handling (Boutros et al., 2015). Extracting meaningful features from microscopy images in 040 large-scale screens has become increasingly difficult as this scale has increased. Public datasets 041 like RxRx3 (Fay et al., 2023) and JUMP-CP (Chandrasekaran et al., 2023) now include millions of 042 cellular images across 100,000s of unique chemical and genetic perturbations. In addition to limita-043 tions in expressiveness of the features that can be derived from them, traditional methods relying on 044 customized pipelines for segmentation, feature extraction, and downstream analysis (Caicedo et al., 045 2017) struggle to handle this scale effectively (Chandrasekaran et al., 2021; Carpenter et al., 2006a).

The size and complexity of large-scale microscopy data demands image models that can extract rich biological features and do so consistently across experimental replicates, both of which are crucial for downstream biomedical applications. Rich, biologically meaningful representations reveal relationships between genes or compounds to drive the discovery of novel targets and drug candidates, while consistency in features extracted across batches and replicates ensures that findings are reproducible and reliable for therapeutic development.

Foundation models have been developed for representing high-dimensional unstructured biological data such as protein structures (Jumper et al., 2021) and transcriptomics (Hao et al., 2024), but the



Figure 1: (A) Overview of performance gain from different MAE pretraining and inference strate-068 gies. (B) Example whole-genome results for replicate consistency and biological relationship recall on StringDB for models trained with different combinations of strategies, by model name and dataset (left to right):

scale and dimensionality of large-scale microscopy data present unique challenges for generating 073 representations that are both biologically informative and consistent across replicates. HCS datasets 074 are often confounded by complex noise known as batch effects (Caicedo et al., 2017), stemming 075 from differences between experimental batches and biological variability. These batch effects -076 including natural variation in cell populations - obscure the biological effects of perturbations and 077 make it challenging to isolate the specific effects of the perturbations applied (Yang et al., 2019). 078 Overcoming these obstacles with a model capable of generating robust, biologically meaningful 079 representations can empower HCS to systematically interrogate gene function and identify novel 080 drug candidates (Rood et al., 2024).

081 State-of-the-art (SOTA) deep learning methods for microscopy leverage Vision Transformers 082 (ViT) (Dosovitskiy et al., 2020) trained with self-supervised learning (SSL) techniques (Balestriero 083 et al., 2023) to learn unbiased representations from large-scale screens (Doron et al., 2023; Kim 084 et al., 2023; Bourriez et al., 2024). Recent studies have demonstrated that ViTs trained as Masked 085 Autoencoders (MAEs) (He et al., 2022; Singh et al., 2023) can effectively scale beyond previous approaches and outperform various supervised and smaller SSL models in capturing biologically 087 informative representations of cell images (Kraus et al., 2024). However, the level of consistency found in these representations across a large number of experimental replicates was not previously reported. Furthermore, compared to recent multi-billion parameter transformers developed for natu-089 ral images (Dehghani et al., 2023) and natural language (Llama3, 2024), model scale in microscopy 090 lags behind (Kraus et al., 2024; Chen et al., 2023a) despite the existence of massive datasets. 091

092 This work offers the following contributions:

067

069

071

094

096

098

- We demonstrate that training on a **curated microscopy dataset** of statistically significant positive samples, named Phenoprints-16M, improves both recall of known gene-gene relationships and consistency of embeddings for gene knockout perturbations (Figure 1A). We describe components of this curation strategy that can be generalized to other scientific datasets (§ 3.1).
- We present a new foundation model, MAE-G/8, a 1.86 billion parameter ViT-G/8 MAE trained on Phenoprints-16M over 48,000 H100 GPU hours on more than 8 billion samples from the curated dataset (Figure 1A, § 3.2).
- 102 • We propose a new set of **biological linear probing tasks** to evaluate representations 103 learned by intermediate ViTs blocks for microscopy data (§ 4). Performance on these linear probing tasks are strongly correlated with performance on important whole-genome scale evaluation metrics while requiring significantly less resources to compute (Figure 4). 105 These trends hold when evaluating on a public benchmarking dataset called RxRx3-core, 106 where our new ViT-G/8 obtains the best average precision for zero-shot cosine similarity 107 prediction of target gene activity across thousands of compounds (§ 6).



struction are inspired by the image retrieval community (Weinzaepfel et al., 2022; Radenović et al., 2018; Berman et al., 2019). Existing methods often utilize pre-trained models for filtering and prun-150 ing, such as vision-language models to discard irrelevant pairs (Schuhmann et al., 2021), semantic 151 deduplication to remove redundancy (Abbas et al., 2023), and prototypicality-based approaches to 152 retain representative data (Sorscher et al., 2022). However, these techniques are less effective for 153 HCS, where redundancy, variability, and subtle morphological differences make conventional filter-154 ing challenging. Our work addresses these limitations by building on Celik et al. (2024)'s perturba-155 tion consistency framework to curate a balanced dataset of images across semantic classes, which is 156 vital for effective learning under the masked objectives (Zhang et al., 2022).

157

Layer-wise Analysis of Deep Neural networks. Recent work suggests that intermediate layers
(or, blocks) in large ViTs may achieve superior performance on certain linear probing tasks compared to the final encoder layer (Evci et al., 2022; Dehghani et al., 2023). Alkin et al. (2024) reported
that intermediate layers in large MAE-ViTs (ViT-L, ViT-H) have superior ImageNet-1K *k*-NN accuracy, likely because later encoder layers become more optimized for the reconstruction task.

162 3 VISION TRANSFORMERS FOR MICROSCOPY IMAGES

We train and evaluate various vision transformers (ViTs, Table 4) as encoders to extract feature embeddings from $256 \times 256 \times 6$ (HxWxC) microscopy image crops (Figure 2).

166 167

168

3.1 TRAINING DATASET CURATION

169 Many academic and industry labs have adopted the Cell Painting imaging protocol (Bray et al., 170 2016), which multiplexes fluorescent dyes to reveal eight broadly relevant cellular components. 171 The datasets used here contain a six-channel implementation of Cell Painting (Figure 2), as well as brightfield images, spanning 100,000s of chemical and genetic perturbations applied to dozens 172 of cell types (Kraus et al., 2024). In these datasets, cells that look like unperturbed cells tend to 173 be very over-represented because many perturbations do no induce a morphological change. Some 174 morphological changes are also far more common (e.g. many perturbations will kill cells, resulting 175 in a relatively high proportion of dead cell morphological phenotype). This results in significant 176 imbalance in the morphological phenotypes that the models learn to reconstruct. 177

To address this, we constructed an aggressively curated training dataset (§ A.1). To learn an initial 178 representation, we began by reproducing the MAE-L/8 model of Kraus et al. (2024) on a dataset of 179 similar size consisting of 93 million HCS images. Using this representation, we first filtered pertur-180 bations that did not induce consistent morphological changes to cells. To perform this filtering, we 181 utilized Celik et al. (2024)'s non-parametric perturbation consistency test (§ A.3) after correcting 182 for batch effects using Typical Variation Normalization (Ando et al., 2017; Kraus et al., 2024). This 183 test was applied within each experiment for computational efficiency, and we restricted the anal-184 ysis to wells containing single perturbations. This consistency was computed for CRISPR guides, 185 siRNAs, and particular concentrations of small molecules across replicates of the same perturbation. P-values were computed for each gene and each (perturbation, concentration) pair. When multiple 187 experiments existed for the same condition, we combined p-values using the Cauchy Combination 188 test (Liu & Xie, 2018).

189 We repeated this procedure with a weakly supervised learning (WSL) model trained on RxRx1 190 (Sypetkowski et al., 2023) and filtered to perturbations where any condition had a p-value < 0.01191 in either the MAE-L/8 or WSL model. This process reduced our original dataset of 93M samples 192 to 16M, which we refer to as Phenoprints-16M. While some redundancy remains when distinct 193 perturbations have the same effect, the proportion of samples with that differ from negative controls increased substantially with little decrease in overall diversity. We believe that iteratively repeating 194 this process with the best models from previous iterations to guide data selection for subsequent 195 models may be a viable strategy. 196

198 3.2 MODELS

199 **Baselines.** We compare to several non-finetuned baseline ViT image encoders: three different 200 Dino-v2 backbones (Oquab et al., 2024) (with 4 register tokens (Darcet et al., 2024)) trained on a 201 curated non-biological natural image dataset; a weakly supervised (WSL) classifier ViT-L/16 trained 202 on Imagenet-21k (Ridnik et al., 2021); a MAE ViT-L/16 trained on Imagenet-21k (He et al., 2022); 203 and an untrained ViT-S/16. We found that channel-wise self-standardization worked best as the 204 image normalization preprocessing for these baselines, and that the class token was slightly better 205 than the global pool of the patch tokens (except for MAE). Convolutional weights in the patch 206 embedding layer were repeated to embed 6 channel images when using models trained on RGB 207 datasets (Wightman, 2019).

208

197

Prior work. Our primary point of comparison is with respect to the best pretrained foundation model presented by Kraus et al. (2024), the MAE-ViT-L/8+ trained on RPI-93M. This MAE-L/8 was trained for approximately 40 epochs, learning from over 3.5 billion image crops, using the L2 mean squared error loss function plus an additional Fourier domain reconstruction loss term.

213

CA-MAE-S/16 trained on RxRx3. We trained a new channel-agnostic MAE (Kraus et al., 2024)
 ViT-S/16 on the RxRx3 dataset (Fay et al., 2023) for 100 epochs. Channel-agnostic ViTs tokenize each image channel separately with shared patch embedding weights and leverage the dynamic



(b) Anax functional gene group classification.

Figure 3: Block-wise validation set linear probe results comparing ViT models pretrained on cell microscopy images (left) versus natural images (right). (a) 1139-class RxRx1 SiRNA knockdown classification (Sypetkowski et al., 2023); (b) 40-class Anax functional gene group classification on HUVEC cell images from RxRx3 CRISPR knockouts (Fay et al., 2023).

sequence length of transformers with repeated positional encodings to train ViTs that can process images with varying numbers of channels (Bao et al., 2024; Bourriez et al., 2024; Kraus et al., 2024). Kraus et al. (2024) demonstrate that the large MAEs with 8x8 patch size perform either better or the same as the 16x16 channel-agnostic variants for consistently 6-channel data, so we opted to train standard MAEs for the following two new models since they require fewer tokens at inference time.

MAE-L/8 trained on Phenoprints-16M. Holding the model backbone constant compared to the MAE-ViT-L/8 by Kraus et al. (2024), we assess the impact of our curated dataset in contrast to the 93M dataset by training a new ViT-L/8 MAE for 500 epochs on Phenoprints-16M.

MAE-G/8 trained on Phenoprints-16M. Holding the dataset constant compared to MAE-L/8 above, we assess the impact of increased model scale in terms of parameters by training a new ViT-Gigantic MAE with nearly 1.9 billion parameters for 500 epochs on Phenoprints-16M. Training this model required 256 H100 GPUs running in parallel for over 1 week. See § A.2 for other hyperparameter settings we used for model training.

LINEAR PROBING REPRESENTATION LEARNING ACROSS VIT BLOCKS

We improve the quality of our learned image representations by leveraging previous findings that suggest intermediate blocks within an encoder can provide better representation compared to the final block (Alkin et al., 2024). Unfortunately, it is infeasible to search for the best block by simply performing whole-genome evaluation on each block of a large model because the evaluation is extremely time-consuming and resource intensive. For example, evaluating the final block of MAE-G/8 required 4,000 L4 GPU hours just for inference (§ 5). We demonstrate that using block-wise linear probes provides insights into the quality of biological features extracted by these models



Figure 4: Correlations between validation set linear probing (Figure 3) on Anax and RxRx1 for
best and last blocks (Eq. 1) compared to downstream whole-genome benchmarks (Table 1) for biological relationship recall on StringDB at 0.05-0.95 threshold and replicate consistency KS statistic.
Models with **bold** borders are *trimmed*, red are natural image baseline models and blue are trained
on microscopy.

291 292

295

296

297

298

299

300301302303304

305

306 307

in their intermediate blocks, allowing us to trim the model to an earlier block to both reduce inference costs and improve representation quality.

Our block-wise search consists of training a logistic regression model (linear probe) on the output features of each transformer block to predict either the gene that was perturbed or the functional group that the gene belongs to, and test performance on held-out experiments (§ A.4). We define the optimal block b^* for a probing task as the block whose output features achieve the highest test balanced accuracy when trained on the probing task, across all N blocks of the encoder,

$$b^* = \underset{b \in \{1,2,\dots,N\}}{\operatorname{arg\,max}} \operatorname{BalancedAccuracy}(\mathbf{z}^{(b)}), \tag{1}$$

where $z^{(b)}$ are output features from block b of a ViT. Performance on our linear probing tasks can be viewed as a measure of linear separability of a feature space across experimental batches.

RxRx1 1139-class siRNA genetic perturbation classification. We expect high quality representations of cell images to generate similar embeddings for cells with the same perturbation, hence a simple linear probe should be able to predict gene perturbation from these representation reasonably well. We train linear probes on the publicly-available RxRx1 dataset Sypetkowski et al. (2023) which consists of 125,510 high-resolution fluorescence microscopy images of human cells under 1,138 siRNA-induced gene knockdowns (plus unperturbed controls) across four cell types (HepG2, HUVEC, U2OS, RPE). These gene knockdowns produce strong phenotypes which makes the prediction task more feasible.

315 We found that, for MAE-G/8, the best features came from intermediate block $b^* = 38$ (out of 48) of 316 the encoder, achieving a balanced accuracy (0.51) that is 8.5% greater compared to its final block's 317 output features (Figure 3a, left). Additionally, these features achieved 60% greater accuracy than 318 the typically used final block of MAE-L/8+ (Kraus et al., 2024). We observed similar trends for 319 ViT models pretrained on natural images. For example, DINO-G/14 and ViT-L/16 MAE trained 320 on non-biological natural image data have their best features at blocks that are positioned within 321 the first half of the encoder. For ViT-L/16 MAE, the performance of the best block is 27% higher compared to its final block output features that are typically used for downstream tasks. The higher 322 performance observed for intermediate blocks does not appear to be an intrinsic feature of the ViT 323 architecture as an untrained ViT did not exhibit such a parabolic trend (Figure 3a, right).

324 Anax 40-class functional gene group classification. Biologically meaningful representation of 325 microscopy images of genetically perturbed cells should capture functional relationships between 326 genes, hence a simple linear probe should be able to predict functional gene groups when trained 327 on these representations. We curated a small subset of 80,000 wells from RxRx3 (Fay et al., 2023) 328 to evaluate linear probes on functional group prediction. We also evaluated similar whole genome knockout screens with ARPE-19 and an additional population of HUVEC cells with soluble TNF- α 329 added to all wells. We manually curated Anax, a set of 40 functionally-diverse gene groups con-330 taining 348 genes, with details provided in (§ A.8). Examples of groups include major protein 331 complexes (e.g. proteasome, ribosome-small/large), metabolic pathways (e.g. Krebs cycle) and sig-332 naling pathways (e.g. calcium signaling) (Figure 2). These groups span broad biological processes 333 that are conserved across cell types - linear separability of these groups would likely indicate that 334 representations are biologically meaningful regardless of cell type. 335

As shown in Figure 3b, MAE-G/8 significantly outperforms other models in Anax group linear probe
 classification. The best representations once again are obtained from an intermediate block, achiev ing a balanced accuracy (0.32) that is 5% greater compared to its final block. We observed similar
 trends for ViT models pretrained on natural images and representations computed from microscopy
 images of other cell types/conditions (§ A.5, Figure 7).

In Figure 4, we observe that performance on this novel linear probing task correlates strongly with downstream whole-genome benchmarks across all models (Table 1), whether they are trained on microscopy data or natural images, achieving an overall rank correlation $\rho = 0.97$ with whole-genome StringDB recall and $\rho = 0.91$ with whole-genome replicate consistency. This strong correlation is crucial as it allows us to trim our model to the block with the best linear probe performance as a way to improve the quality of our representations for the whole-genome (Table 1).

347 348

5 WHOLE-GENOME BENCHMARKING

349 350

351 Table 1 presents our benchmarks computed across the whole-genome. These evaluate the genomic 352 representations obtained for each model by aggregating millions of embeddings of cell images span-353 ning >100,000 of genetic knockout perturbations (17,063 genes \times 6 single guide RNAs each) on HUVEC cells from RxRx3 (Fay et al., 2023). Computing these benchmarks for HCS screens typi-354 cally requires inferring 140 million crops from the genome-wide RxRx3 microscopy screen (Kraus 355 et al., 2023) (64 tiled crops per each of the 2.2 million wells), but, to reduce compute costs, we 356 discard the outer ring of crops, leaving the 36 center non-edge crops for each well. This requires 80 357 million forward passes to comprehensively evaluate a new encoder. After inference, we use typical 358 variation normalization (Ando et al., 2017) and chromosome arm bias correction (Lazar et al., 2024) 359 to post-process the embeddings and aggregate them to the gene-level. 360

We present the multivariate biological relationship recall benchmarks proposed by Celik et al. 361 (2024) and originally evaluated for MAEs by Kraus et al. (2023; 2024). These metrics evaluate 362 how many annotated pair-wise relationships are recalled from public databases (CORUM, hu.MAP, 363 Reactome-PPI, StringDB) in the extremities of a ranked list of cosine similarities of all pair-wise 364 post-processed embeddings (details in \S A.6). To ensure embeddings represent technical replicates of perturbations consistently, we also evaluate model performance on replicate consistency based 366 on the experimental design used in the RxRx3 dataset. Specifically, we compare the similarity of the 367 embedding for corresponding wells across different experiments via a non-parametric statistical test. 368 The test statistic measures the difference between the perturbation replicates' similarity distribution and an empirical null distribution, with larger values indicating greater consistency (details in § A.7). 369

370 In order to compare models, we summarize the resulting statistics over all technical replicates in 371 RxRx3 by taking their median, as reported in columns KS and CVM in Table 1. Even the small-372 est CA-MAE-S/16 trained on microscopy data outperforms all of the large baselines trained on 373 natural images. Furthermore, training on the Phenoprints-16M dataset improves the performance 374 of the MAEs, and MAE-G/8 achieves the best overall performance. Compared to the best pub-375 lished result for whole-genome benchmarks (MAE-L/8 trained on RPI-93M (Kraus et al., 2023)), MAE-G/8 obtains a 48% improvement in replicate consistency CVM (12.3 \rightarrow 18.2) and 4.3% im-376 provement in StringDB recall (.472 \rightarrow .492). Using linear probes to select block $b^* = 15$ (Equation 1) 377 for that MAE-L/8, our improvement changes to 20% in CVM and 3.5% in StringDB recall.

Model backbone	b	CORUM	hu.MAP	React.	StringDB	KS	CVM
Baseline ViTs							
ViT-S/16, Untrained	12	.452	.343	.205	.359	.30	4.3
ViT-L/16, ImageNet WSL	24	.518	.351	.210	.394	.34	5.5
ViT-L/16, ImageNet MAE	24	.526	.355	.215	.397	.34	5.1
trimmed	11	.532	.359	.218	.403	.35	5.8
Baseline Dino ViTs							
ViT-S/14, Dino-V2	12	.484	.345	.203	.380	.34	5.6
trimmed	5	.514	.359	.213	.396	.35	6.0
ViT-L/14, Dino-V2	24	.492	.339	.210	.383	.34	5.3
trimmed	12	.549	.367	.220	.413	.36	5.9
ViT-G/14, Dino-V2	40	.442	.312	.199	.351	.29	3.8
trimmed	16	.529	.354	.220	.398	.33	5.2
MAEs for microscopy							
CA-MAE-S/16, RxRx3	12	.549	.374	.229	.429	.47	10.4
MAE-L/8, RPI-93M	24	.609	.434	.251	.472	.52	12.3
trimmed	15	.602	.427	.255	.475	.57	15.2
MAE-L/8, PP-16M	24	.600	.432	.255	.479	.59	16.2
trimmed	20	.600	.435	.260	.482	.59	16.2
MAE-G/8, PP-16M	48	.621	.438	.263	.488	.60	16.4
trimmed	38	.615	.437	.263	.492	.63	18.2

Table 1: Multivariate **known biological relationship recall** and univariate **replicate consistency** benchmarks by model, encoding block *b*, benchmark database (CORUM, hu.MAP, Reactome-PPI, and StringDB), and consistency test statistics (KS and CVM). The *trimmed* models used linear probes to select an earlier block as the feature encoder (Fig. 3). Results are computed over all >17,000 whole-genome CRISPR knockout perturbation images in RxRx3, after applying TVN and chromosome arm bias correction. Over all benchmarks, higher is better, and best overall result is in **bold**, and best result among baselines is in *italics*. For relationship recall we report the mean @ 0.05-0.95 cosine threshold over 3 random seeds sampling the null distribution for each benchmark run (standard deviation for each result is $\leq \pm .0015$).

Model backbone	b	Pretraining data	CORUM	hu.MAP	Reactome	StringDB
CellProfiler	-	N/A	.219	.184	.131	.191
CA-MAE-S/16	12	RxRx3	.233	.199	.154	.214
MAE-L/8	24	RPI-93M	.248	.208	.160	.226
MAE-G/8	38	Phenoprints-16M	.264	.215	.165	.235

Table 2: Biological relationship recall benchmarks at 0.05-0.95 cosine threshold on public JUMP-CP image data (Chandrasekaran et al., 2023) generated by completely different labs and assay protocols compared to the data used for pretraining. Each result has a standard deviation $\leq \pm .0023$, and spans nearly 8,000 gene-knockouts and are computed after applying PCA with center-scaling for embedding post-processing alignment.

421 422

400

401

402

403

404

405

406

407

Similarly, linear probing to select optimal ViT blocks led to significant improvements even when applied to Dino-V2 based models pretrained on natural images. Dino-V2 ViT-G obtains a nearly 20% improvement in recall on CORUM (.44 \rightarrow .53) by using the embeddings extracted at $b^* = 16$ (chosen by linear probes) rather than the final embedding from b = 40 (which performs worse than a random untrained ViT-S). Dino-V2 ViT-S also observes improvements by using $b^* = 5$ rather than b = 12 and outperforms Dino-V2 ViT-G in replicate consistency. We also attempted to train Dino-V2 on microscopy data, but preliminary results were worse than CA-MAE-S/16 (§ A.9).

Additional results indicate that these MAEs effectively generalize to novel microscopy data generated in different labs with different assays (Table 2, described in § A.11), and that scaling trends for performance on biologically relevant tasks continues from MAE-L/8 to MAE-G/8 (§ A.10).

Model	Avg. prec.	Comp. z-score \uparrow	Energy dist.	Energy z-score ↑
Random baseline	$0.222\pm.007$	0.00	0.728 ± 0.05	0.00
CA-MAE-S/16, RxRx3	$0.273 \pm .016$	2.90	3.319 ± 0.83	3.10
MAE-L/8, RPI-93M	$0.290\pm.017$	3.77	4.856 ± 1.25	3.30
trimmed	$0.299 \pm .016$	4.49	4.514 ± 1.16	3.26
MAE-G/8, PP-16M	$0.302 \pm .015$	4.79	6.053 ± 1.47	3.63
trimmed	$\textbf{0.309}\pm.015$	5.38	$\textbf{6.586} \pm 1.52$	3.85

Table 3: Performance on the public **RxRx3-core compound-gene benchmark** and **RxRx3-core perturbation magnitude benchmark**, measuring mean average precision (\pm STD over 100 random seeds benchmarking with different negative samples) in predicting compound activity against target genes, and mean energy distance (\pm MAD over all perturbations) separating perturbation embeddings from controls with corresponding z-scores of improvement over a random baseline.

444 445 446

447 448

440

441

442

443

6 RXRX3-CORE BENCHMARKING

RxRx3-core¹ is a publicly available benchmarking dataset for assessing biological capabilities of 449 computer vision models. RxRx3-core includes labeled images (compressed to JPEG-2000) of 735 450 genetic knockouts and 1,674 small-molecule perturbations across eight concentrations drawn from 451 222,601 wells ($512 \times 512 \times 6$ pixel center-crops) drawn from the larger RxRx3 dataset. 452

We evaluate a random embedding baseline, the CA-MAE-S/16 model, the MAE-L/8 model from 453 previous work (Kraus et al., 2023), and the MAE-G/8. We evaluate both the trimmed and full-length 454 version of the latter two models to determine the impact of our model-trimming strategy in this 455 context. To evaluate each model, we first inference all $222,601 \times 4$ crops and then average the 4 456 embeddings to the well-level. Then, to perform standard batch correction alignment, we use the 457 "EMPTY", unperturbed wells as our control population. We fit PCA on those control embeddings, 458 use it transform the rest, and then fit a separate standard-scaler on each batch's controls to transform 459 the rest. This simplified alignment strategy empirically performed better TVN only on this dataset. 460

We present results for the benchmark measuring zero-shot prediction of compound-gene activity us-461 ing cosine similarities between embeddings (Figure 5). This measures, for each compound, whether 462 the cosine similarities from a model's embeddings correctly rank the compound's known target 463 genes higher than a randomly sampled set of other genes from a ground truth dataset. Table 3 pro-464 vides exact values along the max axis, which captures the strongest potential interaction regardless of 465 concentration. The relative ranking of model performance, holds as expected from the results in § 4 466 and § 5, and trimming benefits both MAE-L/8 and MAE-G/8, with MAE-G/8 offering a 42% (3.77 467 \rightarrow 5.38) improvement over the method from previous work in predicting compound-gene activity.

468 In Table 3 we also present the results from the perturbation magnitude benchmark (Celik et al., 469 2024), which measures the energy distance between perturbation embeddings and control embed-470 dings (visualized in Figure 6). Unlike on the rest of the benchmarks, the MAE-L/8 trained on 471 RPI-93M does not benefit from trimming here. But, MAE-G/8 obtains the best performance overall 472 and improves when trimmed to the best layer as detected by our linear probes.

473 474 475

DISCUSSION AND CONCLUSIONS 7

476 This work demonstrates that: (1) within the context of biological imaging, trimming many ViTs to an 477 earlier block leads to stronger biological linearity and improved performance on downstream tasks in 478 addition to cheaper inference costs (Figure 3); (2) linear probing performance on a subset of genetic 479 perturbations correlates strongly with downstream performance on whole-genome benchmarks and 480 can be used to optimize which block is selected for representing the whole-genome (Figure 4); (3) 481 the most scaled model, MAE-G/8, obtains the overall best performance across all benchmarks and 482 linear probes, providing further evidence for the scaling hypothesis in biological image data (Table 1, 483 Table 3, § A.10). This demonstrates that intentionally scaling training compute and parameters of 484 SSL models for microscopy can benefit a wide variety of biologically relevant tasks.

¹huggingface.co/datasets/recursionpharma/rxrx3-core



Figure 5: Mean average precision performance on RxRx3-core public benchmark in predicting compound activity against annotated gene targets, across all compound concentrations with error bars for 100 runs of the benchmark with different random seeds (Table 3).



Figure 6: Distribution of perturbation magnitudes for different models as measured by energy distance between model embeddings of perturbations versus controls on RxRx3-core (Table 3).

More broadly, this work proposes a reusable recipe for training and extracting optimal representations from fully self-supervised models trained on experimental data. The pattern we use can be applied to other domains that contain data from repeated experiments but without accurate ground truth labels. Specifically, we recommend: (1) curating the training set by identifying diverse sets of samples that are represented consistently, e.g., by using a pre-existing model to select such samples; (2) training a scaled transformer-based model using a self-supervised learning technique, such as masked autoencoding; and, (3) evaluating the performance of the trained transformer at every block to identify the optimal layer for representing the data.

526 527 528

529

516

517 518 519

520

521

522

523

524

525

486

500

501

502

503

LIMITATIONS AND REPRODUCIBILITY

530 In this work, we evaluated baselines and new linearly probed MAE models trained on a specially 531 curated microscopy dataset. Our preliminary attempts to train ViTs with DINO on this microscopy 532 data encountered suboptimal performance (§ A.9). Consequently, we allocated our time and com-533 pute budget to investigate scaling MAE to ViT-G/8. We recognize the potential for other SSL train-534 ing regimes and fine-tuning strategies (Singh et al., 2023; Lehner et al., 2024; Hondru et al., 2024; 535 Khan & Fang, 2024; Alkin et al., 2024) oriented for microscopy data to lead to future improve-536 ments on these tasks. With this work, we can publicly release the training, inference, reconstruction 537 visualization, and benchmarking code², with the full weights³ for CA-MAE ViT-S/16.

⁵³⁸

²github.com/[redacted]

³huggingface.co/[redacted]

540 REFERENCES

546

547

548

551

561

- Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. Semdedup: Dataefficient learning at web-scale through semantic deduplication. *arXiv preprint arXiv:2303.09540*, 2023.
 - Benedikt Alkin, Lukas Miklautz, Sepp Hochreiter, and Johannes Brandstetter. Mim-refiner: A contrastive learning boost from intermediate pre-trained representations. *arXiv preprint arXiv:2402.10093*, 2024.
- D. Michael Ando, Cory Y. McLean, and Marc Berndl. Improving Phenotypic Measurements in High-Content Imaging Screens. *bioRxiv*, pp. 161422, 2017. doi: 10.1101/161422.
- Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and Micah Goldblum. A Cookbook of Self-Supervised Learning. *arXiv*, 2023. doi: 10.48550/arxiv.2304.12210.
- Yujia Bao, Srinivasan Sivanandan, and Theofanis Karaletsos. Channel vision transformers: An image is worth 1 x 16 x 16 words. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=CK5Hfb5hBG.
- Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain:
 a unified image embedding for classes and instances, 2019. URL https://arxiv.org/ abs/1902.05509.
- Christoph Bock, Paul Datlinger, Florence Chardon, Matthew A. Coelho, Matthew B. Dong, Keith A. Lawson, Tian Lu, Laetitia Maroc, Thomas M. Norman, Bicna Song, Geoff Stanley, Sidi Chen, Mathew Garnett, Wei Li, Jason Moffat, Lei S. Qi, Rebecca S. Shapiro, Jay Shendure, Jonathan S. Weissman, and Xiaowei Zhuang. High-content CRISPR screening. *Nature Reviews Methods Primers*, 2(1):8, 2022. doi: 10.1038/s43586-021-00093-4.
- Nicolas Bourriez, Ihab Bendidi, Ethan Cohen, Gabriel Watkinson, Maxime Sanchez, Guillaume
 Bollot, and Auguste Genovesio. Chada-vit : Channel adaptive attention for joint representation
 learning of heterogeneous microscopy images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- 575
 576
 576
 577
 Michael Boutros, Florian Heigwer, and Christina Laufer. Microscopy-Based High-Content Screening. *Cell*, 163(6):1314–1325, 2015. ISSN 0092-8674. doi: 10.1016/j.cell.2015.11.007.
- Mark-Anthony Bray, Shantanu Singh, Han Han, Chadwick T Davis, Blake Borgeson, Cathy Hartland, Maria Kost-Alimova, Sigrun M Gustafsdottir, Christopher C Gibson, and Anne E Carpenter. Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nature Protocols*, 11(9):1757–1774, 2016. ISSN 1754-2189. doi: 10.1038/nprot.2016.105.
- Juan C Caicedo, Sam Cooper, Florian Heigwer, Scott Warchal, Peng Qiu, Csaba Molnar, AliakStatise S Vasilevich, Joseph D Barry, Harmanjit Singh Bansal, Oren Kraus, Mathias Wawer, Lassi
 Paavolainen, Markus D Herrmann, Mohammad Rohban, Jane Hung, Holger Hennig, John Concannon, Ian Smith, Paul A Clemons, Shantanu Singh, Paul Rees, Peter Horvath, Roger G Linington, and Anne E Carpenter. Data-analysis strategies for image-based cell profiling. *Nature Methods*, 14(9):849–863, 2017. ISSN 1548-7091. doi: 10.1038/nmeth.4397.
- Anne E Carpenter, Thouis R Jones, Michael R Lamprecht, Colin Clarke, In Han Kang, Ola
 Friman, David A Guertin, Joo Han Chang, Robert A Lindquist, Jason Moffat, Polina Golland,
 and David M Sabatini. CellProfiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biology*, 7(10):R100, 2006a. ISSN 1465-6906. doi: 10.1186/
 gb-2006-7-10-r100.

617

628

637

- Anne E Carpenter, Thouis R Jones, Michael R Lamprecht, Colin Clarke, In Han Kang, Ola Friman, David A Guertin, Joo Han Chang, Robert A Lindquist, Jason Moffat, et al. Cellprofiler: image analysis software for identifying and quantifying cell phenotypes. *Genome biology*, 7:1–11, 2006b.
- Safiye Celik, Jan-Christian Hütter, Sandra Melo Carlos, Nathan H. Lazar, Rahul Mohan, Conor Tillinghast, Tommaso Biancalani, Marta M. Fay, Berton A. Earnshaw, and Imran S. Haque. Building, benchmarking, and exploring perturbative maps of transcriptional and morphological data. *PLOS Computational Biology*, 20(10):1–24, 10 2024. doi: 10.1371/journal.pcbi.1012463. URL https://doi.org/10.1371/journal.pcbi.1012463.
- Srinivas Niranj Chandrasekaran, Hugo Ceulemans, Justin D. Boyd, and Anne E. Carpenter. Image based profiling for drug discovery: due for a machine-learning upgrade? *Nature Reviews Drug Discovery*, 20(2):145–159, 2021. ISSN 1474-1776. doi: 10.1038/s41573-020-00117-w.
- Srinivas Niranj Chandrasekaran, Jeanelle Ackerman, Eric Alix, D Michael Ando, John Arevalo, Melissa Bennion, Nicolas Boisseau, Adriana Borowa, Justin D Boyd, Laurent Brino, et al. Jump cell painting dataset: morphological impact of 136,000 chemical and genetic perturbations. *bioRxiv*, pp. 2023–03, 2023.
- Richard J Chen, Tong Ding, Ming Y Lu, Drew F K Williamson, Guillaume Jaume, Bowen Chen, Andrew Zhang, Daniel Shao, Andrew H Song, Muhammad Shaban, Mane Williams, Anurag Vaidya, Sharifa Sahai, Lukas Oldenburg, Luca L Weishaupt, Judy J Wang, Walt Williams, Long Phi Le, Georg Gerber, and Faisal Mahmood. A General-Purpose Self-Supervised Model for Computational Pathology. *arXiv*, 2023a. doi: 10.48550/arxiv.2308.15474.
- Kiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, et al. Symbolic discovery of optimization algorithms.
 arXiv preprint arXiv:2302.06675, 2023b.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2024. URL https://arxiv.org/abs/2309.16588.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer,
 Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling
 vision transformers to 22 billion parameters. In *International Conference on Machine Learning*,
 pp. 7480–7512. PMLR, 2023.
- Michael Doron, Théo Moutakanni, Zitong S. Chen, Nikita Moshkov, Mathilde Caron, Hugo Touvron, Piotr Bojanowski, Wolfgang M. Pernice, and Juan C. Caicedo. Unbiased single-cell morphology with self-supervised vision transformers. *bioRxiv*, 2023. doi: 10.1101/2023.06.16.
 545359. URL https://api.semanticscholar.org/CorpusID:259213557.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im age is worth 16x16 words: Transformers for image recognition at scale. In *International Confer- ence on Learning Representations (ICLR)*, 2020.
- Kevin Drew, Chanjae Lee, Ryan L Huizar, Fan Tu, Blake Borgeson, Claire D McWhite, Yun Ma, John B Wallingford, and Edward M Marcotte. Integration of over 9,000 mass spectrometry experiments builds a global map of human protein complexes. *Molecular Systems Biology*, 13(6): 932, 2017. ISSN 1744-4292. doi: 10.15252/msb.20167490.
- Utku Evci, Vincent Dumoulin, Hugo Larochelle, and Michael C Mozer. Head2toe: Utilizing intermediate representations for better transfer learning. In *International Conference on Machine Learning*, pp. 6009–6033. PMLR, 2022.
- Marta M Fay, Oren Kraus, Mason Victors, Lakshmanan Arumugam, Kamal Vuggumudi, John Urbanik, Kyle Hansen, Safiye Celik, Nico Cernek, Ganesh Jagannathan, et al. Rxrx3: Phenomics map of biology. *bioRxiv*, pp. 2023–02, 2023.

648 649	Marc Gillespie, Bijay Jassal, Ralf Stephan, Marija Milacic, Karen Rothfels, Andrea Senff-Ribeiro, Johannes Griss, Cristoffer Sevilla, Lisa Matthews, Chugiao Gong, Chuan Deng, Thawfeek Varu-
650	sai, Eliot Ragueneau, Yusra Haider, Bruce May, Veronica Shamoysky, Joel Weiser, Timothy Brun-
651	son, Nasim Sanati, Liam Beckman, Xiang Shao, Antonio Fabregat, Konstantinos Sidiropoulos.
652	Julieth Murillo, Guilherme Viteri, Justin Cook, Solomon Shorser, Gary Bader, Emek Demir, Chris
653	Sander, Robin Haw, Guanming Wu, Lincoln Stein, Henning Hermjakob, and Peter D'Eustachio.
654	The reactome pathway knowledgebase 2022. Nucleic Acids Research, 50(D1):D687–D692, 2021.
655	ISSN 0305-1048. doi: 10.1093/nar/gkab1028.
656	Madalina Giurgiu, Julian Reinhard, Barbara Brauner, Irmtraud Dunger, Kaltenhach, Gisela Foho
657	Goar Frishman Corinna Montrone and Andreas Ruepp CORUM the comprehensive resource
658	of mammalian protein complexes—2019. Nucleic Acids Research, 47(Database issue):D559–
659	D563, 2019. ISSN 0305-1048. doi: 10.1093/nar/gky973.
660	Minshana Hao, Jina Cana, Vin Zana, Chimina Liu, Yashana Cua, Vinasi Chana, Taifana Wana
661	Jianzhu Ma, Yuagong Zhang, and La Song. Large scale foundation model on single cell tran
662	scriptomics. Nature Methods, pp. 1, 11, 2024
663	scriptonnes. Nature methods, pp. 1–11, 2024.
664	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
665	toencoders are scalable vision learners. In <i>Proceedings of the IEEE/CVF conference on computer</i>
666	vision and pattern recognition, pp. 16000–16009, 2022.
667	Vlad Hondru, Florinel Alin Croitoru, Shervin Minaee, Radu Tudor Ionescu, and Nicu Sebe. Masked
668	image modeling: A survey. arXiv preprint arXiv:2408.06687, 2024.
669	Jahn Lumman Dishard Essens, Alamandar Drival, Tim Casar, Mishard Eisenman, Olaf Danasharaan
670	John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olar Ronneberger,
671	Katnryn Tunyasuvunakool, Russ Bates, Augustin Zidek, Anna Potapenko, et al. Hignly accurate
672	protein structure prediction with appharoid. <i>nature</i> , 590(7875).585–589, 2021.
673	Alexandr A. Kalinin, John Arevalo, Loan Vulliard, Erik Serrano, Hillary Tsang, Michael Bornholdt,
674	Bartek Rajwa, Anne E. Carpenter, Gregory P. Way, and Shantanu Singh. A versatile information
675	retrieval framework for evaluating profile strength and similarity. <i>bioRxiv</i> , pp. 2024.04.01.587631,
676	4 2024. doi: 10.1101/2024.04.01.587631.
677	Muhammad Osama Khan and Yi Fang. What is the best way to fine-tune self-supervised medical
678	imaging models? In Annual Conference on Medical Image Understanding and Analysis, pp.
679	267–281. Springer, 2024.
680	Vladislav Kim Nikolaos Adaloglou Marc Osterland Flavio M Morelli and Paula A Marin Zapata
681	Self-supervision advances morphological profiling by unlocking powerful image representations.
682	<i>bioRxiv</i> , 2023. doi: 10.1101/2023.04.28.538691.
683	Oren Kreue Kien Kennen Deen Sehen Sehenien Memory Felleh Deter Mellen, Jese Lever Me
684	suday Sharma, Ayla Khan, Jia Balakrishnan, Safiya Calik, et al. Maskad autoencoders are scalable
685	learners of cellular morphology. In Neural Information Processing Systems Workshop on Gener-
686	ative AI and Biology (NeurIPS GenBio). 2023.
687	
688	Oren Kraus, Kian Kenyon-Dean, Saber Saberian, Maryam Fallah, Peter McLean, Jess Leung, Va-
689	sudev Sharma, Ayla Khan, Jia Balakrishnan, Safiye Celik, et al. Masked autoencoders for mi-
690	croscopy are scalable learners of cellular biology. In <i>Proceedings of the IEEE/CVF Conference</i> on Computer Vision and Pattern Personition, pp. 11757, 11768, 2024
691	on Computer vision and Fattern Recognition, pp. 11757–11708, 2024.
692	Nathan H Lazar, Safiye Celik, Lu Chen, Marta M Fay, Jonathan C Irish, James Jensen, Conor A
693	Tillinghast, John Urbanik, William P Bone, Christopher C Gibson, et al. High-resolution genome-
094	wide mapping of chromosome-arm-scale truncations induced by crispr–cas9 editing. <i>Nature Ge</i> -
095	<i>neucs</i> , pp. 1–12, 2024.
096	Johannes Lehner, Benedikt Alkin, Andreas Fürst, Elisabeth Rumetshofer, Lukas Miklautz. and Sepp
600	Hochreiter. Contrastive tuning: A little help to make masked autoencoders forget. In Proceedings
600	of the AAAI Conference on Artificial Intelligence, volume 38, pp. 2965–2973, 2024.
700	Vaowu Liu and Jun Xie. Cauchy combination test: A nowerful test with analytic p value calculation
700	under arbitrary dependency structures <i>Journal of the American Statistical Association</i> 115:303
101	- 402, 2018. URL https://api.semanticscholar.org/CorpusID:56320647.

702	Team Llama3. The Llama 3 Herd of Models. arXiv, 2024. doi: 10.48550/arxiv.2407.21783.
703	Maxime Oquab. Timothée Darcet. Théo Moutakanni, Huy Vo. Marc Szafraniec. Vasil Khalidov.
705	Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nico-
706	las Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael
707	Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Ar-
708	mand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision,
709	2024. UKL https://arxiv.org/abs/2304.0/193.
710	Laralynne Przybyla and Luke A. Gilbert. A new era in functional genomics screens. <i>Nature Reviews</i>
711	Genetics, 23(2):89-103, 2022. ISSN 1471-0056. doi: 10.1038/s41576-021-00409-w.
712	Filin Radenović Giorgos Tolias and Ondřej Chum Fine-tuning cnn image retrieval with no human
713	annotation, 2018. URL https://arxiv.org/abs/1711.02512.
714	
715	Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the message. In Thirty fifth Conference on Neural Information Processing Systems Datasets and
716	Benchmarks Track (Round 1) 2021
717	Denchmarks Track (Round 1), 2021.
/18	Jennifer E Rood, Anna Hupalowska, and Aviv Regev. Toward a foundation model of causal cell and
719	tissue biology with a perturbation cell and tissue atlas. Cell, $18/(17):4520-4545$, 2024.
720	Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
721	Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
723	clip-filtered 400 million image-text pairs. ArXiv, abs/2111.02114, 2021. URL https://api.
724	semanticscholar.org/CorpusID:241033103.
725	Mannat Singh, Quentin Duval, Kalyan Vasudev Alwala, Haoqi Fan, Vaibhav Aggarwal, Aaron Ad-
726	cock, Armand Joulin, Piotr Dollár, Christoph Feichtenhofer, Ross Girshick, et al. The effec-
727	tiveness of mae pre-pretraining for billion-scale pretraining. In Proceedings of the IEEE/CVF
728	International Conference on Computer Vision, pp. 5484–5494, 2023.
729	Srinivasan Sivanandan, Bobby Leitmann, Eric Lubeck, Mohammad Muneeb Sultan, Panagiotis
730	Stanitsas, Navpreet Ranu, Alexis Ewer, Jordan E. Mancuso, Zachary F Phillips, Albert Kim,
731	John W. Bisognano, John Cesarek, Fiorella Ruggiu, David Feldman, Daphne Koller, Eilon
732	Sharon, Ajamete Kaykas, Max R. Salick, and Ci Chu. A Pooled Cell Painting CRISPR Screen-
733	hig Prationin Enables de novo interence of Gene Function by Sen-supervised Deep Learning. higRriv np. 2023 08 13 553051, 2023, doi: 10.1101/2023 08 13 553051
734	<i>biolowy</i> , pp. 2025.00.15.555051, 2025. doi: 10.1101/2025.00.15.555051.
736	Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural
737	scaling laws: beating power law scaling via data pruning. ArXiv, abs/2206.14486, 2022. URL
738	neep3.//ap1.3emane1e3en01a1.019/c01pu510:2301132/3.
739	Maciej Sypetkowski, Morteza Rezanejad, Saber Saberian, Oren Kraus, John Urbanik, James Tay-
740	lor, Ben Mabey, Mason Victors, Jason Yosinski, Alborz Rezazadeh Sereshkeh, et al. Rxrx1: A
741	Conference on Computer Vision and Pattern Recognition pp. 4284_4293, 2023
742	Conjerence on Computer vision and Futtern Recognition, pp. 4204-4275, 2025.
743	Damian Szklarczyk, Annika L Gable, Katerina C Nastou, David Lyon, Rebecca Kirsch, Sampo
744	Pyysalo, Nadezhda T Doncheva, Marc Legeay, Tao Fang, Peer Bork, Lars J Jensen, and Chris-
745	uan von wiering. The STKING database in 2021: customizable protein–protein networks, and functional characterization of user-unloaded gena/massurement sets. Nucleic Acids Passarch 40
746	(D1):D605-D612, 2020, ISSN 0305-1048, doi: 10.1093/nar/gkaa1074.
747	
740	Fabien vincent, Arsenio Nueda, Jonathan Lee, Monica Schenone, Marco Prunotto, and Mark
750	Nature Reviews Drug Discovery 21(12):809_914 2022 ISSN 1474_1776 doi: 10.1038/
751	s41573-022-00472-w.
752	
753	Philippe weinzaeptel, Thomas Lucas, Diane Larlus, and Yannis Kalantidis. Learning super-features
754	tor image reuteval, 2022. UKL https://arxiv.org/abs/2201.13182.
755	Ross Wightman. Pytorch image models. https://github.com/rwightman/ pytorch-image-models, 2019.

758	Model Name	Parameters	Blocks	Model Dim	Pretraining Data
759					
760	Baselines				
761	Untrained ViT-S/16	25M	12	384	N/A
/01	Dino-V2 ViT-S/14	25M	12	384	Natural images
762	Dino-V2 ViT-L/14	307M	24	1024	Natural images
763	Dino-V2 ViT-G/14	1,100M	40	1536	Natural images
764	ViT-L/16 WSL	307M	24	1024	Imagenet-21k
765	ViT-L/16 MAE	307M	24	1024	Imagenet-21k
766	MAEs for microsco	nv			
767	CA-MAE-S/16	ру 25М	12	384	RxRx3
768	MAF-L/8	307M	$\frac{12}{24}$	1024	RPI-93M
769	MAE-L/8	307M	24	1024	Phenoprints-16M
770	MAE-G/8	1,860M	48	1664	Phenoprints-16M
771					

Table 4: Overview of vision transformer (ViT) encoders used and evaluated in this work.

Samuel J. Yang, Scott L. Lipnick, Nina R. Makhortova, Subhashini Venugopalan, Minjie Fan, Zan Armstrong, Thorsten M. Schlaeger, Liyong Deng, Wendy K. Chung, Liadan O'Callaghan, Anton Geraschenko, Dosh Whye, Marc Berndl, Jon Hazard, Brian Williams, Arunachalam Narayanaswamy, D. Michael Ando, Philip Nelson, and Lee L. Rubin. Applying Deep Neural Network Analysis to High-Content Image-Based Assays. Slas Discovery, 24(8):829-841, 2019. ISSN 2472-5552. doi: 10.1177/2472555219857715.

Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. 779 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12104-12113, 2022.

Qi Zhang, Yifei Wang, and Yisen Wang. How mask matters: Towards theoretical understandings of masked autoencoders. In NeurIPS, 2022.

APPENDIX А

TRAINING DATASET CURATION DETAILS A.1

In order to produce Phenoprint-16M, we curated 93M using the following steps:

- 1. Filtering out data that did not pass data quality filters related to the focus of the image, quantity of dead cells, assay conditions, and presence of strong anomalous imaging artifacts.
- 2. Filtering out data with missing information about the perturbations applied, data with more than 3 perturbations applied, and data of unusual size (in the image dimension or number of channels).
- 3. Filtering out perturbation conditions that had been in less than 3 distinct experiments or 20 distinct wells so as to capture a variety of batch effects and have a broad sample of positives per class.
- 4. Under-sampling perturbation conditions that were clearly over-represented in the dataset. Our experiment designs contain positive controls, negative controls, and wells without perturbation within each experiment. At this step, we keep 10% of positive controls and wells without any perturbation, 30% of negative controls, and all other perturbation conditions.
- 5. Filtering out wells where none of the perturbation conditions had a phenoprint (§A.3) (across different map types) in any experiment it had been run in.
- 805 806

808

756

774

775

776

777

778

781 782

783

784 785

786 787

788 789

790 791

792

793 794

796 797

798

799

800

801

802

- A.2 TRAINING HYPERPARAMETERS
- Table 5 provides the hyperparameters used for training the new vision transformers presented in this 809 work. Each model was trained using a 75% mask ratio and the standard decoder architecture for

Table 5: Training hyperparameters for the new models presented in this work. Each used a onecycle cosine learning rate decay schedule with 10% warm-up using the Lion optimizer from Chen
et al. (2023b) with betas (0.9, 0.95) and weight decay of 0.05, with additional ViT settings such as
LayerScale as proposed by Dehghani et al. (2023). *Note that MAE-G/8 had multiple restarts during
training due to challenges associated with massive model training on large-scale shared distributed
compute clusters.

816				
817	Hyperparameter	CA-MAE-S/16	MAE-L/8	MAE-G/8
818	Vision transformer backbone	ViT-S	ViT-L	ViT-G (Zhai et al., 2022)
819	Pretraining Data	RxRx3	Phenoprints-16M	Phenoprints-16M
820	Training epochs	100	500	500*
821	Learning rate	1e-4	3e-5	3e-5
822	Global batch size	2048	16384	8192
823	Stochastic depth	0.1	0.3	0.6
824	# GPUs	16 A100s	128 H100s	256 H100s
825	# GPU-hours	400	15,360	48,000

826

839

840

853

854

855

856

858 859

827 828 MAEs (He et al., 2022). Each model was trained with the standard L2 MAE loss and the Fourier-829 space loss function implemented by Kraus et al. (2024) with a weight of $\alpha = 0.01$. We note, 830 however, that the details presented by Kraus et al. (2024) do not precisely correspond with the implementation provided in their Github repository; when reshaping the tokens to a shape compatible 831 with the 2D Fourier transform, the permute operation resulted in adjacent pixels being from different 832 channels of the input, resulting in the high frequency components of the loss being a function of the 833 relationships between input channels. An initial investigation with a ViT-L/8 showed that changing 834 the implementation to the one described in the paper did not dramatically change probing results. 835 As such, we used the implementation as-is and leave additional analysis of loss function design for 836 MAEs to future work. 837 838

A.3 PERTURBATION CONSISTENCY

In order to assess the consistency of the induced morphology on the cells by the perturbations, we used a non-parametric perturbation consistency test similar to the one introduced in Celik et al. (2024). Let $x_{g,1}, x_{g,2}, \dots, x_{g,n}$ be the embeddings for replicates of perturbation x_g on experiment (batch) *e*. As the test statistic for perturbation consistency, \bar{s}_g^e is defined as the mean of the cosine similarities across all pairs of replicates of x_g .

$$\bar{s}_{g}^{e} = \frac{1}{n^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{\langle x_{g,i}, x_{g,j} \rangle}{||x_{g,i}|| ||x_{g,j}||}.$$
(2)

where $\langle . \rangle$ and ||.|| denote dot product and L_2 norm.

Statistical significance of \bar{s}_g^e is assessed using a permutation test comparing it against an empirical null distribution generated using the same statistic for a set of randomly selected perturbations in experiment $e, \{\bar{s}'_1, \dots, \bar{s}'_K\}$. The p-value for \bar{s}_g^e is computed as follows

$$p_g = \frac{\max\left\{\#\{\bar{s}'_k \ge \bar{s}^e_g\}, 1\right\}}{K}.$$
(3)

861 862

863 When multiple experiments existed for the same perturbation, we combined p-values using the Cauchy Combination test (Liu & Xie, 2018).

A.4 TRAINING LINEAR PROBES

In this section, we provide details about the training process and preprocessing steps used in our
logistic regression models. These models were trained on output features derived from various
Vision Transformer (ViT) blocks.

The data was split by experiments, ensuring that the test data originated from experiments distinct
 from those used for training. This approach helps to validate the generalization performance of our
 models across different experimental conditions.

For both RxRx1 gene prediction and Anax group prediction, we apply StandardScaler from the scikit-learn library as the only preprocessing step to standardize the features prior to training linear probes. StandardScaler transformation was fitted on data from the train split. We trained the logistic regression models using scikit-learn's LogisticRegression class. The following parameters and settings were used during model optimization:

• Solver: lbfgs

- Maximum Iterations: 2000
- Class Weight: balanced

For RxRx1 gene prediction, we trained logistic regression models to predict one of 1139 possible perturbation labels (1138 genetic perturbation and non-perturbed control). For Anax group prediction, we trained logistic regression models to predict one of 40 possible function group labels (§ A.8). We report the balanced test accuracy as the main evaluation metric for all linear probing experiments.

A.5 ANAX CLASSIFICATION FOR OTHER CELL LINES/TREATMENT CONDITIONS: ARPE19 AND HUVEC WITH TNF-ALPHA BACKGROUND

We performed linear probing on imaging data obtained for a retinal pigment epithelia (RPE) cell line, ARPE19, and HUVEC cells treated with an inflammatory cytokine, $TNF\alpha$. We similarly observed that intermediate blocks often have the most linearly separate features compared to the final block.



Figure 7: Layerwise validation set linear probe performance on Anax functional gene group classification beyond RxRx3: CRISPR knockouts in the ARPE-19 immortalized epithelial cell-line (left), and in HUVEC cells with a TNF- α background (right).

A.6 BIOLOGICAL RELATIONSHIP RECALL

A valuable use of large-scale HCS experiments is to perform large-scale inference of biological relationships between genetic perturbations. We evaluate each model's ability to recall known relationships by using the *biological relationship recall* benchmark described in Celik et al. (2024).
First, we correct for batch effects using *Typical Variation Normalization* (TVN) (Ando et al., 2017), and also correct for possible chromosome arm biases known to exist in CRISPR-Cas9 HCS data

918 (Lazar et al., 2024). To infer biological relationships, we compute the aggregate embedding of each 919 perturbation by taking the spherical mean over its replicate embeddings across experiments. We use 920 the cosine similarity of a pair of perturbation representations as the relationship metric, setting the 921 origin of the space to the mean of negative controls. We compare these similarities with the rela-922 tionships found in the following public databases: CORUM (Giurgiu et al., 2019), hu.MAP (Drew et al., 2017), Reactome (Gillespie et al., 2021), and StringDB Szklarczyk et al. (2020) (with >95% 923 combined score). Table 1 reports the recall of known relationships amongst the top and bottom 5% 924 of all cosine similarities between CRISPR knockout representations in RxRx3 (Fay et al., 2023). 925

926 927

A.7 REPLICATE CONSISTENCY

928 In order to assess the reproducibility of the perturbations across their technical replicates, we com-929 pare the distributions of the similarities for same perturbations across replicates against an empirical 930 null distribution. Specifically, for technical replicate experiments e_a^i and e_b^i , we calculate the cosine 931 similarity between the embeddings of perturbation x_i in them, denoted as s^{x_j} . The query distribution 932 q^{e_i} is constructed by computing the cosine similarities for all perturbations that have a matching well 933 on experiments e_a^i and e_b^i . An empirical null distribution of identical cardinality is created by com-934 puting cosine similarity, r^{x_k,x_l} , between random pairs from e_a^i and e_b^i such that no pair corresponds 935 to the same perturbation, $p_0^{e_i}$. Using non-parametric statistical tests, namely Kolmogorov-Smirnov (KS) and Cramer Von-Mises (CVM), we can evaluate the hypothesis that q^{e_i} and $p_0^{e_i}$ are drawn 936 from the same distribution. Formally, let $Q^{e_i}(x)$ and $P_0^{e_i}(x)$ be the cumulative distribution func-937 tions for q^{e_i} and $p_{0_i}^{e_i}$ respectively, then the KS statistic for the two-sample case of technical replicate 938 experiments e_a^i and e_b^i is defined as: 939

$$KS^{e_i} = \sup_{x} |Q^{e_i}(x) - P_0^{e_i}(x)|.$$
(4)

The Cramér–von Mises test statistic (CVM) for experiments e_a^i and e_b^i is computed as:

946

940 941

$$CVM^{e_i} = \frac{1}{2N^2} \sum_{m=1}^{N} \left[(r_m - m)^2 + (s_m - m)^2 \right] - \frac{4N^2 - 1}{12N}.$$
(5)

where N is the cardinality of q^{e_i} and $p_0^{e_i}$ and s_m and r_m are ranks of similarities s^{x_j} and r^{x_k,x_l} in the combined distribution of q^{e_i} and $p_0^{e_i}$ when ordered. In order compare models, we use the median of CVM^{e_i} and KS^{e_i} over all technical replicate experiment pairs e_i .

Since the pairs are randomly selected for $p_0^{e_i}$, the embeddings would be mostly orthogonal thus the distribution would be centered around 0.Given that not all CRISPR knockouts would induce a morphological change in the cells, it's plausible for distribution q^{e_i} to exhibit a peak around 0. As the model approaches the precision of an oracle, we would anticipate the mass situated around this peak to shift towards higher cosine similarity values.

955 956 957

A.8 ANAX GROUP PREDICTION DETAILS

The Anax probing task introduced in this paper is intended to balance capturing a diverse range of
biology that is broadly conserved between cell types with a reduced cost of execution. The name
"Anax" is a reference to Anaximander, the 6th century B.C. philosopher credited with making the
first world map.

In curating these genes, we analyzed the sources listed in § A.6 as well as internal gene expression
data to produce "functional" groups corresponding to biological processes, cellular components, and
molecular functions. Not all genes within each group are expected to have the same knockout phenotype, but are classified by humans as having related function – linear separability of these genes
would indicate that a model has learned similar concepts to those deemed significant by biologists.

- 967 The gene groups we use for the 40-class Anax group classification task (§ A.4) are listed in Table 7.
- 968

- 969 A.9 DINO-V2 PRETRAINING ON MICROSCOPY DATA
- 971 We attempted to train two Dino-v2 models on microscopy data. One ViT-L/16 from scratch on RxRx3, and another attempt of fine-tuning the MAE-L/8 on RPI-93M with the Dino-v2 losses.



Figure 8: Loss curve when training Dino-v2 ViT-L/16 on RxRx3.

DinoV2 model	Avg. Prec.	Comp. Z-score	Energy dist.	Energy Z-score
ViT-L/16, RxRx3	$0.258 \pm .015$	2.13	4.818 ± 1.49	2.740
ViT-L/8 (ft) RPI-93M	$0.255 \pm .018$	1.76	2 705 ± 0 75	2.626

Table 6: RxRx3-core benchmarks for our initial attempts to train Dino-V2 models on microscopy data. The latter was finetuned from the MAE-L/8 trained on RPI-93M. Results compare to Table 3.

They had the following hyperparameter settings which were tuned on another dataset: output-dim
65536, 2 global crops, 4 local crops, dino loss weight 1.0, koleo loss weight 0.1, ibot loss weight 1.0,
Lion optimizer with max learning rate 1e-5, weight decay 0.05, betas 0.9 0.95, and cosine annealing.
In both cases, we observed significant over-fitting of the loss from the start (Figure 8).

997 In Table 6 we show that both models fail to improve on the RxRx3-core benchmark metrics (i.e., 998 z-scores over the random baseline) versus the CA-MAE-S/16 RxRx3 model which had Z-score of 999 2.90 on average precision for predicting compound activity and 3.10 on energy distance between 1000 perturbations and controls. We have not found an effective recipe for training Dino on microscopy 1001 data. However, we note that theoretical evidence exists arguing that MAE learning is in some ways 1002 equivalent to contrastive learning Hondru et al. (2024), so even if an appropriate Dino recipe is found it would remain to be seen if it differs substantially from MAEs for microscopy given the 1003 same training compute. As described in the Limitations section, we expect that future work would 1004 have to dedicate significant training ablations and creativity to determine the best possible training 1005 recipe for training Dino on microscopy data.

1007

981 982

989

990

991 992

> 1008 1009

A.10 CORRELATION BETWEEN MODEL SCALE AND BENCHMARK RESULTS

In Figure 9 we show the correlations between training FLOps (floating point operations) and downstream results. Over all benchmarks we observe a very strong consistent linear trend where scaling training FLOps improves overall pwerformance. This work provides the next log step in scale as we enter into the billion-parameter model regime with MAE-G/8. These results therefore provide additional evidence that the trend initially discovered by Kraus et al. (2023) between FLOps and relationship recall actually extends both to billion-parameter models and even moreso for other biologically meaningful benchmarks pertaining to linear probes on small experiments and to replicate consistency on the whole-genome.

1017

1018 A.11 PERFORMANCE ON NEW DATA (JUMP-CP)

In order to validate that the MAEs generalize to entirely novel data, we evaluated a subset of models
on completely external public data generated by different assays and from a variety of different labs
as produced by the JUMP-CP consortium (Chandrasekaran et al., 2023). Table 2 presents these
results using the relationship recall benchmarks of (Celik et al., 2024), noting that only a subset
of 7,976 gene-knockouts are covered by this dataset. For post-processing embedding alignment, we
use PCA with center-scaling. We observe that the MAEs perform better than the Cellprofiler manual
feature extraction baseline (Carpenter et al., 2006b), and that the general trend is maintained with the

1026	Table 7: Anax groups and their associated genes. This table presents a comprehensive list of gene
1027	groups and their corresponding genes.

Anax Group	Genes
Acyl Coa Biosynthesis	ELOVL2, ELOVL6, ELOVL6, HACD1, HACD2, HSD17B12, SCD, SCD5, TECR
Adherens Junctions	ACTB, ACTG1, AFDN, CDH1, CTNNA1, CTNNB1, CTNND1, NECTIN1, NECTIN3, NECTIN4
Amino Acid Metabolism	ALDH4A1, ARG2, CKB, CKMT2, CPS1, DAO, OTC, PYCR2, PYCR3, SAT1
Apoptosis	CFLAR, DFFB, CASP6, CASP3, FASLG, BCL2, DFFA, XIAP, TNFSF10, AKT3
Autophagy	ATG12, ATG3, ATG4B, ATG4C, ATG7, GABARAP, PIK3C3, PIK3R4, PRKAA1, ULK1
Beta Oxidation Of Fatty Acids	ACAA2, ACADL, ACADM, ACADS, ACADVL, ECHS1, ECI1, HADH, HADHA, HADHB
Calcium Signaling	ADCY1, ADCY2, ADCY3, CALM1, CAMK2B, CAMK2D, PDE1B, PDE1C, PRKACG, PRKX
Clathrin Coated Vesicles	AP2A1, AP2A2, AP2B1, AP2M1, AP2S1
COPI	ARCN1, COPA, COPB1, COPB2, COPE, COPG1, COPZ1
COPII Vesicles	SEC13, SEC23A, SEC24B, SEC24D, SEC31A
DNA Damage Repair	BLM, BRCA2, EME1, NBN, POLD2, RAD51B, RAD51C, RAD51D, RPA1, XRCC2
Dynein	DYNC1H1, DYNC1I2, DYNC1LI1, DYNC1LI2, DYNLT1
ER Protein Translocation	SPCS3, SEC61A1, SRP14, SRP72, SPCS1, SRPRA, SEC11A, SRP68, SRPRB, SRP54
Exosome	DIS3, EXOSC10, EXOSC3, EXOSC4, EXOSC5, EXOSC6, EXOSC7, EXOSC8, EXOSC9, MPHOSPH
Gap Junctions	ADCY8, DRD2, HTR2C, ITPR2, LPAR1, PDGFD, PDGFRB, PLCB3, TUBA1C, TUBB1
Golgi	ACTR10, ACTR1A, CAPZA3, COG4, CTSZ, PPP6C, RAB1B, SEC22C, SEC24C, TMED9
MAPK	DUSP4, EGF, FGF18, FGF20, HSPB1, MAP2K2, MAPKAPK5, RAC1, RAP1A, RASGRP3
Mitochondria Structure	APOOL, APOO, TMEM11, CHCHD6, ATP5ME, MICOS13, ATP5F1C, DNAJC11, DMAC2L, ATP5M
Mitochondrial Transport	ATP5F1A, COA4, COA6, COX17, HSPA9, IDH3G, PITRM1, PMPCA, PMPCB, SLC25A4
mTOR Pathway	CAB39, CAB39L, EIF4EBP1, MLST8, PRKAA2, RPS6KB1, RPTOR, STK11, STRADA, TSC1
Nonsense Mediated Decay	CASC3, EIF4A3, MAGOH, MAGOHB, RBM8A
Nuclear Pore	NUP107, NUP133, NUP153, NUP188, NUP205, NUP37, NUP85, NUP93
Nucleolus Structure	FBL, NAT10, NOLC1, NOP58, UTP20
Nucleotide Metabolism	ADSL, ADSS1, ADSS2, ATIC, GMPS, IMPDH1, IMPDH2, PAICS, PFAS, PPAT
P53 Stress Signaling	ATM, ATR, CCNG1, CDK1, CHEK1, CHEK2, MDM2, MDM4, TP53, TP73
Pentose Phosphate Pathway	G6PD, TALDO1, DERA, RPE, PGM2, RBKS, PGD, PGLS, RPEL1, PRPS2
Peroxisome Biology	ACOT8, AGPS, BAAT, HMGCL, HSD17B4, MLYCD, PAOX, PEX12, PEX6, PIPOX
Prespliceosome Complex	ALYREF, AQR, CRNKL1, DDX5, HNRNPK, LSM2, PLRG1, PRPF4, SMNDC1, SRSF4
Proteasome	PSMA1, PSMA4, PSMB1, PSMB2, PSMB7, PSMA6, PSMA3, PSMB4, PSMA5, PSMB3
Ribosome Large	RPL13A, RPL11, RPL10, RPL23A, RPL30, RPL7A, RPLP2, RPL28, RPL5, RPL27A
Ribosome Small	RPS2, RPS6, RPS8, RPS16, RPS11, RPS3A, RPS19, RPS15, RPS4X, RPS9
RNA Polymerase II	POLR2A, POLR2B, POLR2C, POLR2G, POLR2I, POLR2L
TCA Cycle	ACO2, DLST, FH, IDH2, IDH3B, MDH2, OGDH, SDHB, SUCLA2, SUCLG2
Tight Junctions	CLDN14, CLDN17, CLDN18, CLDN19, CLDN4, CLDN8, CLDN9, MPP5, PARD6B, PRKCI
Translation Initiation Complex	EIF3G, EIF3A, EIF3D, EIF3I, EIF3K, EIF3M, EIF3B, EIF3H, EIF3E, EIF3L
Transport Of Fatty Acids	APOD, LCN12, LCN15, LCN9, SLC27A1, SLC27A4, SLC27A6
Tubulin	TUBA3C, TBCC, TBCD, TUBA4B, TUBA8, TUBAL3, TUBA1A, TUBB4B. ARL2. TUBA1B
Unfolded Protein Response	CXXC1, DNAJB11, EIF2S3, KHSRP, MBTPS1, SHC1, TATDN2, TLN1, TSPYL2, YIF1A
VATPasa	$\Delta T D E V 1 \Delta \Delta T D E V 1 D \Delta T D E V 1 E 1 \Delta T D E V 1 E 1 \Delta T D E V 1 E 1 D E V 1 $

trimmed MAE-G/8 obtaining the best recall overall. Notably, recall on JUMP-CP is considerably lower than on RxRx3 (Table 1) likely due to different assay protocols and more variance in the data.



Figure 9: Relationship between FLOPs and benchmark evaluation results for the six whole-genome tasks (Table 1) and the two linear probing tasks (Figure 3).