# ViTally Consistent: Scaling Biological Representation Learning for Cell Microscopy

Kian Kenyon-Dean[1]     Zitong Jerry Wang[1]     John Urbanik[1]     Konstantin Donhauser[2]

Jason Hartford[2,3]     Saber Saberian[1]     Nil Sahin[1]     Ihab Bendidi[2]     Safiye Celik[1]

Marta Fay[1]     Juan Sebastián Rodríguez Vera[1]     Imran S Haque[1]

Oren Kraus[1]

[1]Recursion  [2]Valence Labs  [3]University of Manchester

## 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 date, 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. 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. More broadly, our approach and results provide insights toward a general strategy for successfully building foundation models for large-scale biological data.[1]

## 1 Introduction

Large-scale cell microscopy assays are used to discover previously unknown biological processes (Przybyla & Gilbert, 2022; Bock et al., 2022; Rood et al., 2024) and identify novel drug candidates 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 large-scale screens has become increasingly difficult as this scale has increased. Public datasets like RxRx3 (Fay et al., 2023) and JUMP-CP (Chandrasekaran et al., 2023) now include millions of cellular images across 100,000s of unique chemical and genetic perturbations. In addition to limitations in expressiveness of the features that can be derived from them, traditional methods relying on customized pipelines for segmentation, feature extraction, and downstream analysis (Caicedo et al., 2017) struggle to handle this scale effectively (Chandrasekaran et al., 2021; Carpenter et al., 2006).

---

[1]Correspondence: `kian.kd@recursion.com`, `info@rxrx.ai`

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

State-of-the-art (SOTA) deep learning methods for microscopy leverage Vision Transformers (ViT) (Dosovitskiy et al., 2020) trained with self-supervised learning (SSL) techniques (Balestriero et al., 2023) to learn unbiased representations from large-scale screens (Doron et al., 2023; Kim et al., 2023; Bourriez et al., 2024). Recent studies have demonstrated that ViTs trained as Masked Autoencoders (MAEs) (He et al., 2022) can effectively scale beyond previous approaches and outperform various supervised and smaller SSL models in capturing biologically 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 natural images (Dehghani et al., 2023) and natural language (Llama3, 2024), model scale in microscopy lags behind (Kraus et al., 2024; Chen et al., 2023a) despite the existence of massive datasets.

In this work, we developed the largest foundation model to date for cell microscopy images, achieving SOTA results in both replicate consistency and biological recall of known gene-gene relationships. Specifically our work offers the following contributions:

- 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).

- We propose new set of **biological linear probing tasks** to evaluate representations 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).

- We find that **using intermediate layers leads to better performance** on these downstream whole-genome benchmarks at a lower computational inference cost, across SSL ViTs trained on microscopy or natural images. By taking advantage of our linear probing proxy task, we are able to cheaply find the best performing intermediate block (Eq. 1).

Our results indicate that the biological scaling properties first identified by Kraus et al. (2023) extend to the multi-billion parameter regime (§ A.9). We show that our MAE-G/8 model produces a nearly 60% more phenotypically linearly separable latent space compared to previous approaches using the final block of MAE-L/8 (Figure 4), correlating with significant improvements in both recall and replicate consistency when benchmarking across the whole genome (Figure 1B).
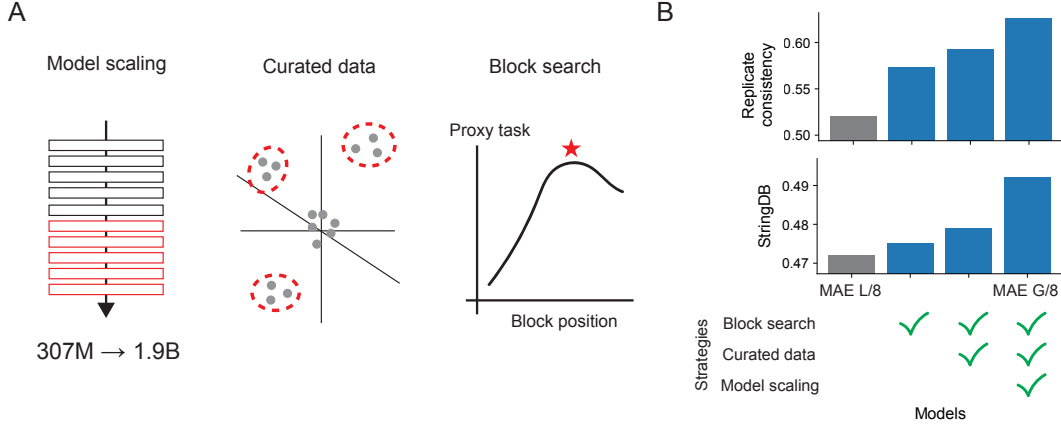
Figure 1: (A) Overview of performance gain from different foundation model pretraining and inference strategies. (B) Example whole-genome results for **replicate consistency** and **biological relationship recall** on StringDB for models trained with different combinations of strategies, by model name and dataset (left to right): MAE-L/8 (RPI-93M, block $b = 24$), MAE-L/8 trimmed to block $b^* = 15$, MAE-L/8 (Phenoprints-16M, block $b^* = 20$), MAE-G/8 (Phenoprints-16M, $b^* = 38$), where $b^*$ is the optimal block according to linear probes as defined in Equation 1.
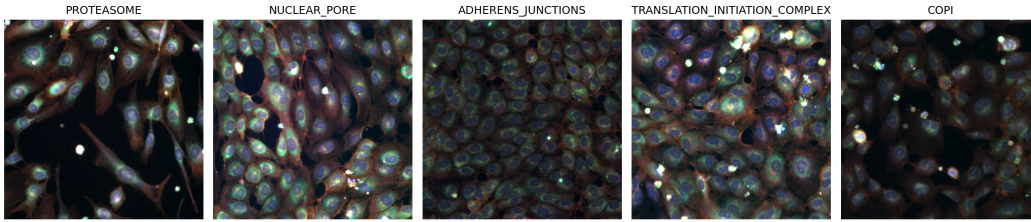


Figure 2: Samples for subset of groups in Anax 40-class functional gene group classification task.

## 2   Related work

**Evaluating Representations for Drug Discovery.**   Evaluating the quality of biological representation learning methods for drug discovery remains challenging, as ground truth data is sparse, noisy, biased to well-studied diseases and pathways, and poorly annotated. Metrics have been proposed that use mean average precision (Kalinin et al., 2024) or AUC ROC (Sivanandan et al., 2023) to assesses how similar related samples are represented, including replicates of the same perturbation or different perturbations with similar annotated biological activities. Recently, Celik et al. (2024) introduced terminology for describing perturbative "maps of biology", in which representations of perturbations in HCS data can be placed in unified, relatable embedding spaces allowing for the generation of genome-scale sets of pairwise comparisons. Here we leverage the *biological relationship recall* benchmark proposed by Celik et al. (2024), which assess how well known relationships between pairs of perturbations are recalled among the most similar or dissimilar embeddings. Computing reliable versions of these relationship benchmarks with HCS data is particularly expensive as they require genome-wide embeddings to be inferred for hundreds of millions of image crops from the genome-wide RxRx3 microscopy screen (Fay et al., 2023).

**Dataset Curation for Foundation Models.**   Dataset curation is crucial for enhancing the efficiency of foundation models, especially in large-scale contexts. Usual approaches to dataset construction 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 pruning, such as vision-language models to discard irrelevant pairs (Schuhmann et al., 2021), semantic deduplication to remove redundancy (Abbas et al., 2023), and prototypicality-based approaches to retain representative data (Sorscher et al., 2022). However, these techniques are less effective for HCS, where redundancy, variability, and subtle morphological differences make conventional filter-

ing challenging. Our work addresses these limitations by building on Celik et al. (2024)'s *perturbation consistency* framework to curate a balanced dataset of images across semantic classes, which is vital for effective learning under the masked objectives (Zhang et al., 2022).

**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). For example, Alkin et al. (2024) reported that intermediate layers in large MAE-ViTs (ViT-L, ViT-H) have superior ImageNet-1K $k$-NN accuracy. They attributed this property to the later encoder layers becoming more optimized for the reconstruction task.

# 3  Vision Transformers for Microscopy Images

We train and evaluate various vision transformers (ViTs, Table 1) as encoders to extract feature embeddings from $256 \times 256 \times 6$ (HxWxC) microscopy image crops (Figure 2).

## 3.1  Training Dataset Curation

Many academic and industry labs have adopted the Cell Painting imaging protocol (Bray et al., 2016), which multiplexes fluorescent dyes to reveal eight broadly relevant cellular components. 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 of cell types (Kraus et al., 2024). In these datasets, cells that look like unperturbed cells tend to be very over-represented because many perturbations do no induce a morphological change. Some morphological changes are also far more common (e.g. many perturbations will kill cells, resulting in a relatively high proportion of dead cell morphological phenotype). This results in significant imbalance in the morphological phenotypes that the models learn to reconstruct.

To address this, we constructed an aggressively curated training dataset (§ A.1). To learn an initial representation, we began by reproducing the MAE-L/8 model of Kraus et al. (2024) on a dataset of similar size consisting of 93 million HCS images. Using this representation, we first filtered perturbations that did not induce consistent morphological changes to cells. To perform this filtering, we utilized Celik et al. (2024)'s non-parametric perturbation consistency test (§ A.3) Typical Variation Normalization (Ando et al., 2017a; Kraus et al., 2024). This test was applied within each experiment for computational efficiency, and we restricted the analysis to wells containing single perturbations. This consistency was computed for CRISPR guides, 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 experiments existed for the same condition, we combined p-values using the Cauchy Combination test (Liu & Xie, 2018).

We repeated this procedure with a weakly supervised learning (WSL) model trained on RxRx1 (Sypetkowski et al., 2023) and filtered to perturbations where any condition had a p-value $< 0.01$ in either the MAE-L/8 or WSL model. This process reduced our original dataset of 93M samples to 16M, which we refer to as Phenoprints-16M. While some redundancy remains when distinct 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 this process with the best models from previous iterations to guide data selection for subsequent models may be a viable strategy.

## 3.2  Models

**Baselines.** We compare to several non-finetuned baseline ViT image encoders: three different Dino-v2 backbones (Oquab et al., 2024) (with 4 register tokens (Darcet et al., 2024)) trained on a curated non-biological natural image dataset; a weakly supervised (WSL) classifier ViT-L/16 trained on Imagenet-21k (Ridnik et al., 2021); a MAE ViT-L/16 trained on Imagenet-21k (He et al., 2022); and an untrained ViT-S/16. Preliminary investigations found that channel-wise self-standardization worked best as the image normalization preprocessing for these baselines, and that the class token was slightly better than the global pool of the patch tokens (except for MAE). Convolutional weights

Table 1: Overview of vision transformer (ViT) encoders used and evaluated in this work.

| Model Name | Parameters | Blocks | Model Dim | Pretraining Data |
|---|---|---|---|---|
| **Baselines** | | | | |
| Untrained ViT-S/16 | 25M | 12 | 384 | N/A |
| Dino-V2 ViT-S/14 | 25M | 12 | 384 | Natural images |
| Dino-V2 ViT-L/14 | 307M | 24 | 1024 | Natural images |
| Dino-V2 ViT-G/14 | 1,100M | 40 | 1536 | Natural images |
| ViT-L/16 WSL | 307M | 24 | 1024 | Imagenet-21k |
| ViT-L/16 MAE | 307M | 24 | 1024 | Imagenet-21k |
| **MAEs for microscopy** | | | | |
| CA-MAE-S/16 | 25M | 12 | 384 | RxRx3 |
| MAE-L/8 | 307M | 24 | 1024 | RPI-93M |
| MAE-L/8 | 307M | 24 | 1024 | Phenoprints-16M |
| MAE-G/8 | 1,860M | 48 | 1664 | Phenoprints-16M |

in the patch embedding layer were repeated to embed 6 channel images when using models trained on RGB datasets (Wightman, 2019).

**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-93Mrained 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.

**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 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.

## 4 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 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

(a) RxRx1 siRNA knockdown classification.



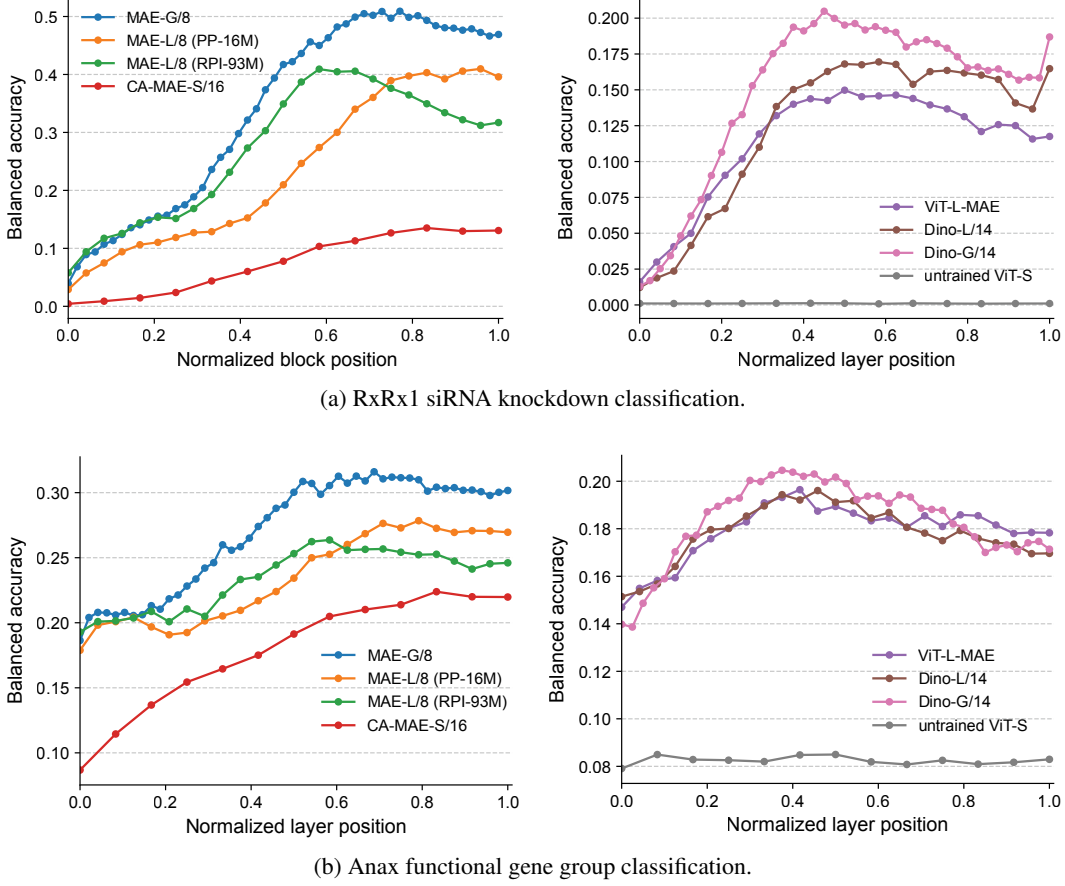(b) Anax functional gene group classification.

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).

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,

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

where $\mathbf{z}^{(b)}$ 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 in 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.

We found that, for MAE-G/8 , the best features came from intermediate block $b^* = 38$ (out of 48) of the encoder, achieving a balanced accuracy (0.51) that is $8.5\%$ greater compared to its final block's output features (Figure 3a, left). Additionally, these features achieved $60\%$ greater accuracy than
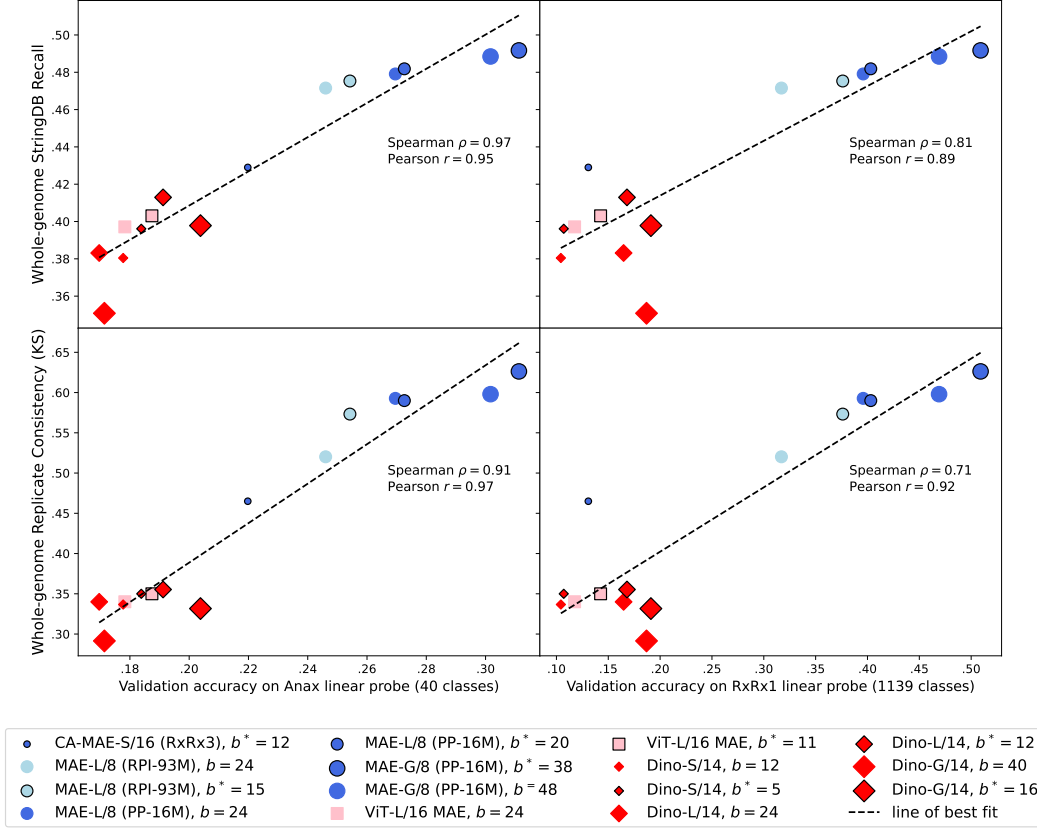
6

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 2) for biological relationship recall on StringDB at 0.05-0.95 threshold and replicate consistency KS statistic.

the typically used final block of MAE-L/8+ (Kraus et al., 2024). We observed similar trends for ViT models pretrained on natural images. For example, DINO-G/14 and ViT-L/16 MAE trained on non-biological natural image data have their best features at blocks that are positioned within 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 performance observed for intermediate blocks does not appear to be an intrinsic feature of the ViT architecture as an untrained ViT did not exhibit such a parabolic trend (Figure 3a, right).

**Anax 40-class functional gene group classification.** Biologically meaningful representation of microscopy images of genetically perturbed cells should capture functional relationships between genes, hence a simple linear probe should be able to predict functional gene groups when trained on these representations. We curated a small subset of 80,000 wells from RxRx3 (Fay et al., 2023) 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-$\alpha$ added to all wells. We manually curated Anax, a set of 40 functionally-diverse gene groups containing 348 genes, with details provided in (§ A.8). Examples of groups include major protein complexes (e.g. proteasome, ribosome-small/large), metabolic pathways (e.g. Krebs cycle) and signaling pathways (e.g. calcium signaling) (Figure 2). These groups span broad biological processes that are conserved across cell types – linear separability of these groups would likely indicate that representations are biologically meaningful regardless of cell type.

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, achieving a balanced accuracy (0.32) that is $5\%$ greater compared to its final block's output features. We

Table 2: Multivariate **known biological relationship recall** and univariate **replicate consistency** benchmarks by model, encoding block $b$, benchmark database, and test statistic. The *trimmed* models used linear probes to select an earlier block as the feature encoder (Fig. 3). Results are computed over all whole-genome CRISPR knockout perturbation images in RxRx3, after applying TVN and chromosome arm bias correction. For relationship recall, we report results over four databases (React stands for Reactome-PPI (Gillespie et al., 2021)). Best overall result is in **bold**.

| Model backbone | $b$ | CORUM | hu.MAP | React | StringDB | KS | CVM |
|---|---|---|---|---|---|---|---|
| **Baseline ViTs** | | | | | | | |
| ViT-S/16, Untrained | 12 | .45 | .34 | .205 | .36 | .30 | 4.3 |
| ViT-S/14, Dino-V2 | 12 | .48 | .345 | .20 | .38 | .34 | 5.6 |
| *trimmed* | 5 | .51 | .36 | .21 | .40 | .35 | 6.0 |
| ViT-L/16, ImageNet WSL | 24 | .52 | .35 | .21 | .39 | .34 | 5.5 |
| ViT-L/14, Dino-V2 | 24 | .49 | .34 | .21 | .38 | .34 | 5.3 |
| *trimmed* | 12 | .55 | .37 | .22 | .41 | .36 | 5.9 |
| ViT-L/16, ImageNet MAE | 24 | .53 | .355 | .215 | .40 | .34 | 5.1 |
| *trimmed* | 11 | .53 | .36 | .22 | .40 | .35 | 5.8 |
| ViT-G/14, Dino-V2 | 40 | .44 | .31 | .20 | .35 | .29 | 3.8 |
| *trimmed* | 16 | .53 | .35 | .22 | .40 | .33 | 5.2 |
| **MAEs for microscopy** | | | | | | | |
| CA-MAE-S/16 , RxRx3 | 12 | .55 | .37 | .23 | .43 | .47 | 10.4 |
| MAE-L/8 , RPI-93M | 24 | .61 | .43 | .25 | .47 | .52 | 12.3 |
| *trimmed* | 15 | .60 | .43 | .255 | .475 | .57 | 15.2 |
| MAE-L/8 , PP-16M | 24 | .60 | .43 | .255 | .48 | .59 | 16.2 |
| *trimmed* | 20 | .60 | .435 | .26 | .48 | .59 | 16.2 |
| MAE-G/8 , PP-16M | 48 | **.62** | **.44** | **.26** | .49 | .60 | 16.4 |
| *trimmed* | 38 | .615 | **.44** | **.26** | **.49** | **.63** | **18.2** |

observed similar trends for ViT models pretrained on natural images and representations computed from microscopy images of other cell types/conditions (§ A.5, Figure 6).

In Figure 4, we observe that performance on this novel linear probing task correlates strongly with downstream whole-genome benchmarks across all models (Table 2), 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 2).

## 5 Whole-genome benchmarking

Table 2 presents our benchmarks computed across the whole-genome. These evaluate the genomic representations obtained for each model by aggregating millions of embeddings of cell images spanning >100,000 of genetic knockout perturbations (17,063 genes × 6 single guide RNAs each) on HUVEC cells from RxRx3 (Fay et al., 2023). Computing these benchmarks for HCS screens typically requires inferring 140 million crops from the genome-wide RxRx3 microscopy screen (Kraus et al., 2023) (64 tiled crops per each of the 2.2 million wells), but, to reduce compute costs, we discard the outer ring of crops, leaving the 36 center non-edge crops for each well. This requires 80 million forward passes to comprehensively evaluate a new encoder. After inference, we use typical variation normalization (Ando et al., 2017b) and chromosome arm bias correction (Lazar et al., 2023) to post-process the embeddings and aggregate them to the gene-level.

We present the multivariate **biological relationship recall** benchmarks proposed by Celik et al. (2024) and originally evaluated for MAEs by Kraus et al. (2023, 2024). These metrics evaluate how many annotated pair-wise relationships are recalled from public databases (CORUM, hu.MAP, Reactome-PPI, StringDB) in the extremities of a ranked list of cosine similarities of all pair-wise post-processed embeddings (details in § A.6). To ensure embeddings represent technical replicates of perturbations consistently, we also evaluate model performance on **replicate consistency** based
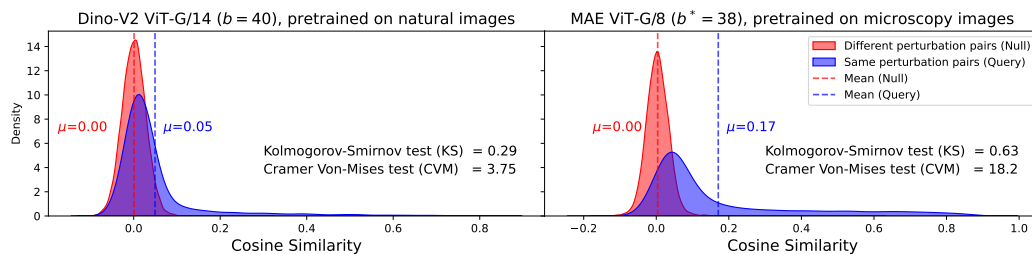
Figure 5: **Replicate consistency** whole-genome results between perturbations from cosine similarity distributions on RxRx3 post-TVN. Comparing baseline Dino-V2 ViT-G/14 at its typically used final block (left) versus MAE-G/8 at the best block found via linear probing (right).

on the experimental design used in the RxRx3 dataset. Specifically, we compare the similarity of the embedding for corresponding wells across different experiments via a non-parametric statistical test. 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).

In order to compare models, we summarize the resulting statistics over all technical replicates in RxRx3 by taking their median, as reported in columns KS and CVM in Table 2, and visualized in Figure 5. MAEs pretrained on microscopy data show improved performance compared the baseline models. Furthermore, training on the Phenoprints-16M dataset improves the performance of the MAEs significantly and *trimmed* MAE-G/8 achieves the best overall performance.

Our linear probing analysis (Figure 3) allowed us to *trim* our models to better encoding blocks. Comparing models on their best respective blocks, MAE-G/8 improves on MAE-L/8 with a 16% improvement in Anax functional gene group classification (.27→.31) and a 24% improvement in RxRx1 perturbation classification (.41→.51). Compared to the best published result for whole-genome benchmarks (MAE-L/8 trained on RPI-93M (Kraus et al., 2023)), MAE-G/8 obtains a 20% improvement in replicate consistency KS (.52→.63) and 4.3% improvement in StringDB recall (.472→.492). When using our linear probes to select outputs from block $b^* = 15$ (Equation 1) from that MAE-L/8, the gain for MAE-G/8 changes to 9.2% in KS and 3.5% in StringDB recall.

Similarly, linear probing to select optimal ViT blocks led to significant improvements even when applied to frozen Dino-V2 based models pretrained on natural images. Dino-V2 ViT-G obtains a nearly 20% improvement CORUM recall (.44→.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.

## 6   Discussion and Conclusions

This work demonstrates that: (1) within the context of biological imaging, trimming many ViTs to an earlier block leads to stronger biological linearity and improved performance on downstream tasks in addition to cheaper inference costs (Figure 3); (2) linear probing performance on a subset of genetic perturbations correlates strongly with downstream performance on whole-genome benchmarks and can be used to optimize which block is selected for representing the whole-genome (Figure 4); (3) the most scaled model, MAE-G/8 , obtains the overall best performance across all benchmarks and linear probes, providing further evidence for the scaling hypothesis in biological image data (Table 2). This demonstrates that intentionally scaling training compute and parameters of SSL models for microscopy can benefit downstream biological relationship recall, whole-genome replicate consistency, and biological linear separability on smaller datasets (see § A.9 for scaling plots).

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

9

masked autoencoding; and, (3) evaluating the performance of the trained transformer at every block to identify the optimal layer for representing the data.

# References

Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. Semdedup: Data-efficient 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. M. J. Ando, Cory Y. McLean, and Marc Berndl. Improving phenotypic measurements in high-content imaging screens. *bioRxiv*, 2017a. URL https://api.semanticscholar.org/CorpusID:26552204.

D. Michael Ando, Cory Y. McLean, and Marc Berndl. Improving Phenotypic Measurements in High-Content Imaging Screens. *bioRxiv*, pp. 161422, 2017b. 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.

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, Aliaksei 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, 2006. ISSN 1465-6906. doi: 10.1186/gb-2006-7-10-r100.

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.

Xiangning 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 image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference 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.

Marc Gillespie, Bijay Jassal, Ralf Stephan, Marija Milacic, Karen Rothfels, Andrea Senff-Ribeiro, Johannes Griss, Cristoffer Sevilla, Lisa Matthews, Chuqiao Gong, Chuan Deng, Thawfeek Varusai, Eliot Ragueneau, Yusra Haider, Bruce May, Veronica Shamovsky, Joel Weiser, Timothy Brunson, Nasim Sanati, Liam Beckman, Xiang Shao, Antonio Fabregat, Konstantinos Sidiropoulos, Julieth Murillo, Guilherme Viteri, Justin Cook, Solomon Shorser, Gary Bader, Emek Demir, Chris Sander, Robin Haw, Guanming Wu, Lincoln Stein, Henning Hermjakob, and Peter D'Eustachio. The reactome pathway knowledgebase 2022. *Nucleic Acids Research*, 50(D1):D687–D692, 2021. ISSN 0305-1048. doi: 10.1093/nar/gkab1028.

Madalina Giurgiu, Julian Reinhard, Barbara Brauner, Irmtraud Dunger-Kaltenbach, Gisela Fobo, Goar Frishman, Corinna Montrone, and Andreas Ruepp. CORUM: the comprehensive resource of mammalian protein complexes—2019. *Nucleic Acids Research*, 47(Database issue):D559–D563, 2019. ISSN 0305-1048. doi: 10.1093/nar/gky973.

Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. Large-scale foundation model on single-cell transcriptomics. *Nature Methods*, pp. 1–11, 2024.

Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.

John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.

Alexandr A. Kalinin, John Arevalo, Loan Vulliard, Erik Serrano, Hillary Tsang, Michael Bornholdt, Bartek Rajwa, Anne E. Carpenter, Gregory P. Way, and Shantanu Singh. A versatile information retrieval framework for evaluating profile strength and similarity. *bioRxiv*, pp. 2024.04.01.587631, 4 2024. doi: 10.1101/2024.04.01.587631.

Vladislav Kim, Nikolaos Adaloglou, Marc Osterland, Flavio M Morelli, and Paula A Marin Zapata. Self-supervision advances morphological profiling by unlocking powerful image representations. *bioRxiv*, 2023. doi: 10.1101/2023.04.28.538691.

Oren Kraus, Kian Kenyon-Dean, Saber Saberian, Maryam Fallah, Peter McLean, Jess Leung, Vasudev Sharma, Ayla Khan, Jia Balakrishnan, Safiye Celik, et al. Masked autoencoders are scalable learners of cellular morphology. In *Neural Information Processing Systems Workshop on Generative AI and Biology (NeurIPS GenBio)*, 2023.

Oren Kraus, Kian Kenyon-Dean, Saber Saberian, Maryam Fallah, Peter McLean, Jess Leung, Vasudev Sharma, Ayla Khan, Jia Balakrishnan, Safiye Celik, et al. Masked autoencoders for microscopy are scalable learners of cellular biology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11757–11768, 2024.

Nathan H Lazar, Safiye Celik, Lu Chen, Marta Fay, Jonathan C Irish, James Jensen, Conor A Tillinghast, John Urbanik, William P Bone, Genevieve HL Roberts, et al. High-resolution genome-wide mapping of chromosome-arm-scale truncations induced by crispr-cas9 editing. *bioRxiv*, pp. 2023–04, 2023.

Yaowu Liu and Jun Xie. Cauchy combination test: A powerful test with analytic p-value calculation under arbitrary dependency structures. *Journal of the American Statistical Association*, 115:393 – 402, 2018. URL https://api.semanticscholar.org/CorpusID:56320647.

Team Llama3. The Llama 3 Herd of Models. *arXiv*, 2024. doi: 10.48550/arxiv.2407.21783.

Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. URL https://arxiv.org/abs/2304.07193.

Laralynne Przybyla and Luke A. Gilbert. A new era in functional genomics screens. *Nature Reviews Genetics*, 23(2):89–103, 2022. ISSN 1471-0056. doi: 10.1038/s41576-021-00409-w.

Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human annotation, 2018. URL `https://arxiv.org/abs/1711.02512`.

Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.

Jennifer E Rood, Anna Hupalowska, and Aviv Regev. Toward a foundation model of causal cell and tissue biology with a perturbation cell and tissue atlas. *Cell*, 187(17):4520–4545, 2024.

Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *ArXiv*, abs/2111.02114, 2021. URL `https://api.semanticscholar.org/CorpusID:241033103`.

Srinivasan Sivanandan, Bobby Leitmann, Eric Lubeck, Mohammad Muneeb Sultan, Panagiotis Stanitsas, Navpreet Ranu, Alexis Ewer, Jordan E. Mancuso, Zachary F Phillips, Albert Kim, John W. Bisognano, John Cesarek, Fiorella Ruggiu, David Feldman, Daphne Koller, Eilon Sharon, Ajamete Kaykas, Max R. Salick, and Ci Chu. A Pooled Cell Painting CRISPR Screening Platform Enables de novo Inference of Gene Function by Self-supervised Deep Learning. *bioRxiv*, pp. 2023.08.13.553051, 2023. doi: 10.1101/2023.08.13.553051.

Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *ArXiv*, abs/2206.14486, 2022. URL `https://api.semanticscholar.org/CorpusID:250113273`.

Maciej Sypetkowski, Morteza Rezanejad, Saber Saberian, Oren Kraus, John Urbanik, James Taylor, Ben Mabey, Mason Victors, Jason Yosinski, Alborz Rezazadeh Sereshkeh, et al. Rxrx1: A dataset for evaluating experimental batch correction methods. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4284–4293, 2023.

Damian Szklarczyk, Annika L Gable, Katerina C Nastou, David Lyon, Rebecca Kirsch, Sampo Pyysalo, Nadezhda T Doncheva, Marc Legeay, Tao Fang, Peer Bork, Lars J Jensen, and Christian von Mering. The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Research*, 49 (D1):D605–D612, 2020. ISSN 0305-1048. doi: 10.1093/nar/gkaa1074.

Fabien Vincent, Arsenio Nueda, Jonathan Lee, Monica Schenone, Marco Prunotto, and Mark Mercola. Phenotypic drug discovery: recent successes, lessons learned and new directions. *Nature Reviews Drug Discovery*, 21(12):899–914, 2022. ISSN 1474-1776. doi: 10.1038/s41573-022-00472-w.

Philippe Weinzaepfel, Thomas Lucas, Diane Larlus, and Yannis Kalantidis. Learning super-features for image retrieval, 2022. URL `https://arxiv.org/abs/2201.13182`.

Ross Wightman. Pytorch image models. `https://github.com/rwightman/pytorch-image-models`, 2019.

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. 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.

# A   Appendix

## A.1   Training dataset curation details

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.

## A.2   Training hyperparameters

Table 3: Training hyperparameters for the new models presented in this work. Each used a one-cycle 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.

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

Table 3 provides the hyperparameters used for training the new vision transformers presented in this work. Each model was trained using a 75% mask ratio and the standard decoder architecture for MAEs (He et al., 2022). Each model was trained with the standard L2 MAE loss and the Fourier-space loss function implemented by Kraus et al. (2024) with a weight of $\alpha = 0.01$. We note, 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 with the 2D Fourier transform, the permute operation resulted in adjacent pixels being from different channels of the input, resulting in the high frequency components of the loss being a function of the relationships between input channels. An initial investigation with a ViT-L/8 showed that changing the implementation to the one described in the paper did not dramatically change probing results. As such, we used the implementation as-is and leave additional analysis of loss function design for MAEs to future work.

## 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 $x_{g,1}, x_{g,2}, \cdots, x_{g,n}$ 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 $x_g$.

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

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 $e$, $\{\bar{s}'_1, \cdots, \bar{s}'_K\}$. The p-value for $\bar{s}_g^e$ is computed as follows

$$p_g = \frac{\max\left\{ \#\{\bar{s}'_k \geq \bar{s}_g^e\}, 1 \right\}}{K}. \tag{3}$$

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.
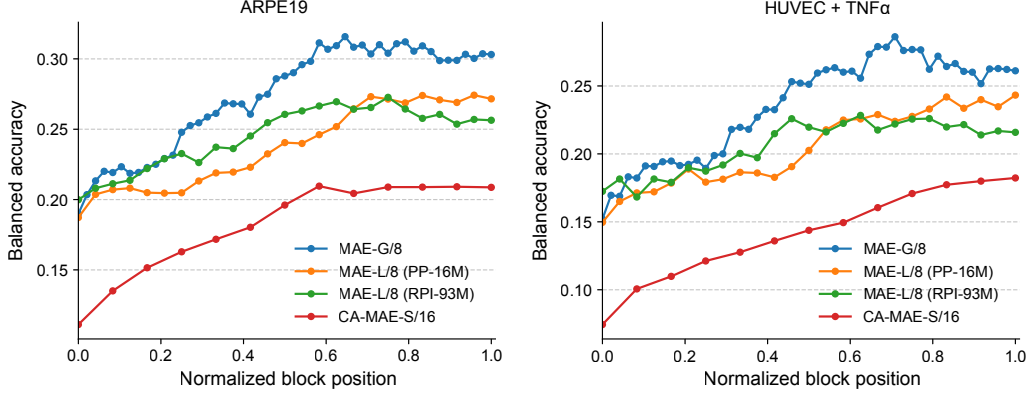
Figure 6: Layerwise validation set linear probe performance on Anax functional gene group classification beyond RxRx3: CRISPR knockouts in the ARPE-19 immortalized epithelial cell-line (left), and in HUVEC cells with a TNF-$\alpha$ background (right).

## 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., 2017b), and also correct for possible chromosome arm biases known to exist in CRISPR-Cas9 HCS data (Lazar et al., 2023). To infer biological relationships, we compute the aggregate embedding of each perturbation by taking the spherical mean over its replicate embeddings across experiments. We use the cosine similarity of a pair of perturbation representations as the relationship metric, setting the origin of the space to the mean of negative controls. We compare these similarities with the relationships 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% combined score). Table 2 reports the recall of known relationships amongst the top and bottom 5% of all cosine similarities between CRISPR knockout representations in RxRx3 (Fay et al., 2023).

## A.7 Replicate Consistency

In order to assess the reproducibility of the perturbations across their technical replicates, we compare the distributions of the similarities for same perturbations across replicates against an empirical null distribution. Specifically, for technical replicate experiments $e_a^i$ and $e_b^i$, we calculate the cosine similarity between the embeddings of perturbation $x_j$ in them, denoted as $s^{x_j}$. The query distribution $q^{e_i}$ is constructed by computing the cosine similarities for all perturbations that have a matching well on experiments $e_a^i$ and $e_b^i$. An empirical null distribution of identical cardinality is created by computing cosine similarity, $r^{x_k,x_l}$, between random pairs from $e_a^i$ and $e_b^i$ such that no pair corresponds 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 from the same distribution. Formally, let $Q^{e_i}(x)$ and $P_0^{e_i}(x)$ be the cumulative distribution functions for $q^{e_i}$ and $p_0^{e_i}$ respectively, then the KS statistic for the two-sample case of technical replicate experiments $e_a^i$ and $e_b^i$ is defined as:

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

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

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

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 $\text{CVM}^{e_i}$ and $\text{KS}^{e_i}$ over all technical replicate experiment pairs $e_i$.

Table 4: Anax groups and their associated genes. This table presents a comprehensive list of gene groups and their corresponding genes.

| Anax Group | Genes |
| --- | --- |
| Acyl Coa Biosynthesis | ELOVL2, ELOVL5, 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, MPHOSPH6 |
| 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, ATP5MF |
| 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 |
| V-ATPase | ATP6V1A, ATP6V, ATP6V1D, ATP6V1E1, ATP6V1F, ATP6V1H |

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, similar to what Figure 5 illustrates. 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.

## 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.

The gene groups we use for the 40-class Anax group classification task (§ A.4) are listed in Table 4.
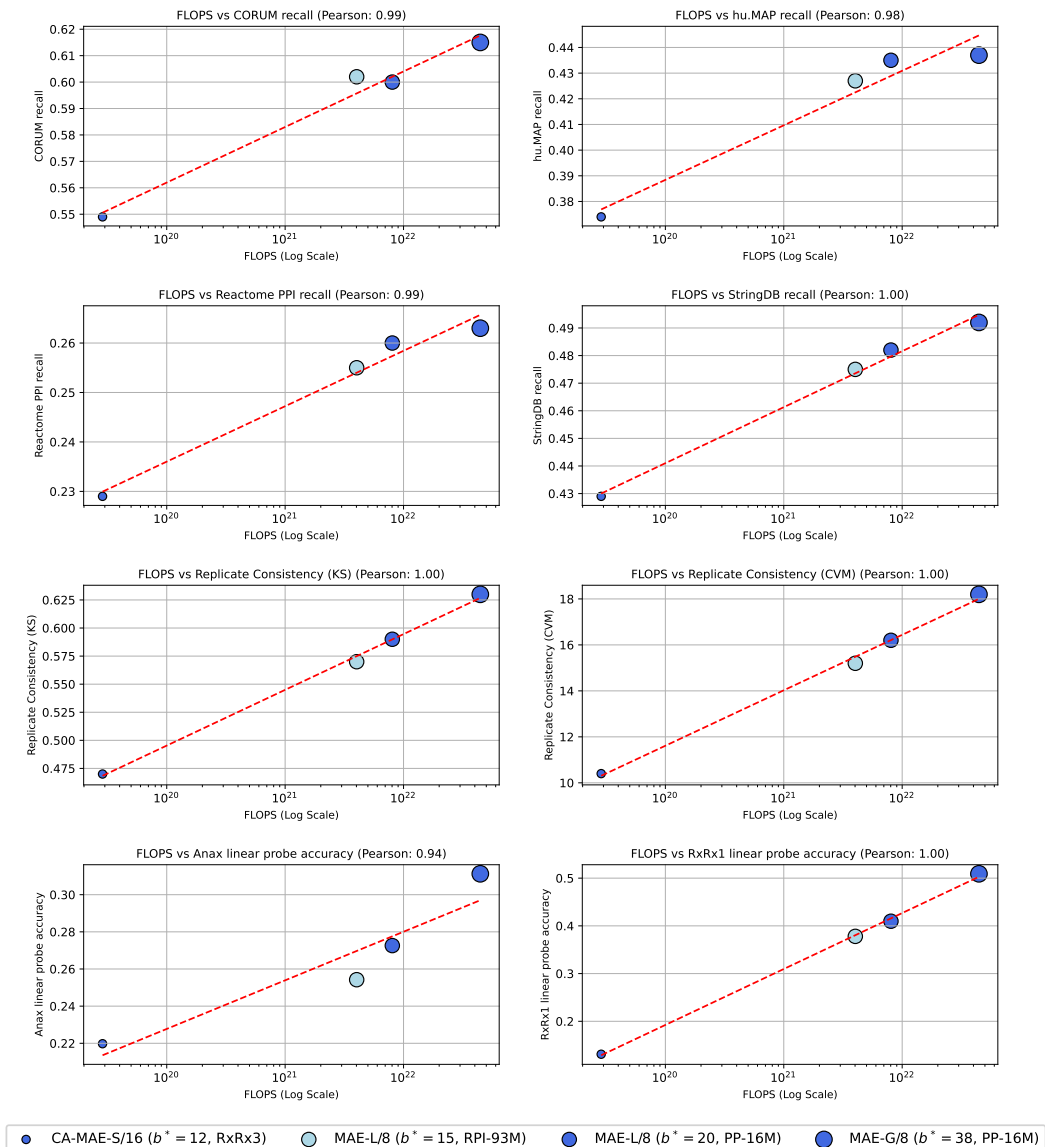
17

Figure 7: Relationship between FLOPs and benchmark evaluation results for the six whole-genome tasks (Table 2) and the two linear probing tasks (Figure 3).

## A.9 Correlation between model scale and benchmark results

In Figure 7 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.