

# Revealing Subtle Phenotypes in Small Microscopy Datasets Using Latent Diffusion Models

Anonymous CVPR submission

Paper ID

## Abstract

001 *Identifying subtle phenotypic variations in cellular images*  
002 *is critical for advancing biological research and acceler-*  
003 *ating drug discovery. These variations are often masked*  
004 *by the inherent cellular heterogeneity, making it challeng-*  
005 *ing to distinguish differences between experimental condi-*  
006 *tions. Recent advancements in deep generative models have*  
007 *demonstrated significant potential for revealing these nu-*  
008 *anced phenotypes through image translation, opening new*  
009 *frontiers in cellular and molecular biology as well as the*  
010 *identification of novel biomarkers. Among these generative*  
011 *models, diffusion models stand out for their ability to pro-*  
012 *duce high-quality, realistic images. However, training dif-*  
013 *fusion models typically requires large datasets and substan-*  
014 *tial computational resources, both of which can be limited*  
015 *in biological research. In this work, we propose a novel*  
016 *approach that leverages pre-trained latent diffusion mod-*  
017 *els to uncover subtle phenotypic changes. We validate our*  
018 *approach qualitatively and quantitatively on several small*  
019 *datasets of microscopy images. Our findings reveal that our*  
020 *approach enables effective detection of phenotypic varia-*  
021 *tions, capturing both visually apparent and imperceptible*  
022 *differences. Ultimately, our results highlight the promising*  
023 *potential of this approach for phenotype detection, espe-*  
024 *cially in contexts constrained by limited data and compu-*  
025 *tational capacity.*

## 026 1. Introduction

027 In recent years, generative models have undergone rapid  
028 and accelerating advancements [4, 14, 19, 33, 40], resulting  
029 in their widespread adoption across a variety of fields. No-  
030 tably, these models have made significant contributions to  
031 biological research. For example, they have been employed  
032 in protein design [41], predicting protein structures [22],  
033 integrating cancer data [37], synthesizing biomedical im-  
034 ages [13, 23], predicting molecular structures [5, 32], and  
035 identifying phenotypic cell variations [2, 3, 27].

Identifying phenotypic variations in biological images 036  
is crucial for advancing our understanding of biological 037  
processes. Detecting these differences can be particularly 038  
challenging due to the high degree of biological variabil- 039  
ity, yet it holds immense potential for enhancing disease 040  
understanding, discovering novel biomarkers, and develop- 041  
ing new therapeutics and diagnostics [7, 28, 31]. Traditi- 042  
onal methods for identifying these phenotypes often rely 043  
on cell segmentation and the quantification of features such 044  
as intensity, shape, and texture [31]. Recently, deep learn- 045  
ing techniques, particularly generative models [2, 3, 27], 046  
have been applied to automate and refine this process, en- 047  
abling the identification of more interpretable and biolog- 048  
ically meaningful features. Among these approaches, dif- 049  
fusion models have emerged as state-of-the-art generative 050  
models [8], achieving remarkable results in tasks such as 051  
image synthesis. However, training diffusion models, like 052  
other deep learning models requires large datasets, which is 053  
often difficult to obtain in biological applications. 054

In this work, we propose **Phen-LDiff** a method to detect 055  
cellular variations in small biological datasets by leveraging 056  
pre-trained Latent Diffusion Models (LDMs) [35]. 057

## 2. Related Work 058

**Diffusion Models.** Diffusion Models (DMs)[19, 40] are 059  
generative models that have recently achieved remarkable 060  
results in various tasks. DMs are latent variable models that 061  
operate through two key processes: a fixed forward process 062  
that gradually adds noise to the data and a learned back- 063  
ward process that denoises it, reconstructing the data distri- 064  
bution [19, 40]. Recently, these models have seen several 065  
advancements [8, 18, 26, 40], making them state-of-the-art 066  
in image synthesis, surpassing traditional generative mod- 067  
els like GANs [8]. One of the notable improvements is 068  
the introduction of Latent Diffusion Models (LDMs) [35], 069  
where images are first compressed into a latent space us- 070  
ing a variational autoencoder, and then the diffusion pro- 071  
cess occurs within this compressed latent space. This ap- 072  
proach enables more efficient scaling to higher-resolution 073

074 images and accelerates training times. Additionally, LDMs  
075 incorporate a conditioning mechanism, allowing for tasks  
076 such as text-conditioned image generation, inpainting, and  
077 super-resolution. These innovations in LDMs have facil-  
078 itated their training on massive datasets [36], resulting in  
079 powerful pre-trained models such as Stable Diffusion [35],  
080 which have demonstrated exceptional performances in var-  
081 ious generative tasks.

082 **Identifications of Phenotypes in Biological Images.**  
083 Identifying phenotypic variations in biological images is es-  
084 sential in biology and drug discovery [7, 31], yet it presents  
085 significant challenges. One of the key difficulties is the  
086 biological variability among cells within the same condi-  
087 tion, which can obscure the differences between distinct  
088 conditions. Recently, generative models have been em-  
089 ployed to cancel this natural variability in order to visu-  
090 alize and explain cellular phenotypes in microscopy im-  
091 ages [2, 11, 27]. In [2], cellular variations between con-  
092 ditions were identified through an image-to-image transla-  
093 tion task between two classes, following methodologies sim-  
094 ilar to those in [21, 44]. In Phenexplain [27], a conditional  
095 StyleGAN2 [25] was trained to detect cellular changes by  
096 performing translations between synthetic images within  
097 the latent space of StyleGAN2, allowing for training across  
098 multiple conditions, unlike the approach in [2]. A sim-  
099 ilar method was presented in [11], but instead of utilizing  
100 the latent space of GANs, the authors proposed learning  
101 a representation space using self-supervised learning tech-  
102 niques [15]. In [3], conditional diffusion models were ap-  
103 plied to identify phenotypes in real images. This approach  
104 consists of two stages: first, the source class image is in-  
105 verted into a latent code, which is then used to generate an  
106 image from the target class. This method provides a power-  
107 ful alternative for phenotype detection using real biological  
108 data. However, all of these models require a large number  
109 of images to be properly trained.

110 **Fine-tuning Diffusion Models.** Fine-tuning [16, 20, 30,  
111 38, 42], a well-established strategy for training deep learn-  
112 ing models on limited data, involves adapting pre-trained  
113 models. It involves adapting a pre-trained model’s weights  
114 to fit a smaller dataset. Fine-tuning methods can be catego-  
115 rized into three main groups: adaptive methods [34, 38],  
116 where the entire model’s weights are adjusted; selective  
117 methods [1, 12, 43], where only a subset of the model’s  
118 parameters are modified; and additive methods [16, 20],  
119 where additional networks are incorporated to refine the  
120 weights. These techniques have proven effective for dis-  
121 criminative models and have recently been extended to gen-  
122 erative models, such as GANs, autoregressive generative  
123 models [20], and diffusion models [16]. Fine-tuning tech-  
124 niques for diffusion models have gained attention, particu-

larly due to the availability of models pre-trained on large  
125 datasets. Recently, several approaches have been proposed  
126 for fine-tuning diffusion models [16, 20, 30], driven by the  
127 popularity of pre-trained models like Stable Diffusion [35].  
128 In [30], it was demonstrated that modifying a subset of pa-  
129 rameters can lead to efficient fine-tuning. Low-Rank Adap-  
130 tation (LoRA)[20], a technique originally developed for  
131 fine-tuning large language models (LLMs) [29], can also be  
132 applied to diffusion models. LoRA freezes the pre-trained  
133 model’s weights and learns low-rank matrices that are in-  
134 jected into each layer of the network. In[16], the authors  
135 introduced SVDiff, a fine-tuning method for diffusion mod-  
136 els that focuses on learning shifts in the model’s singular  
137 values.  
138

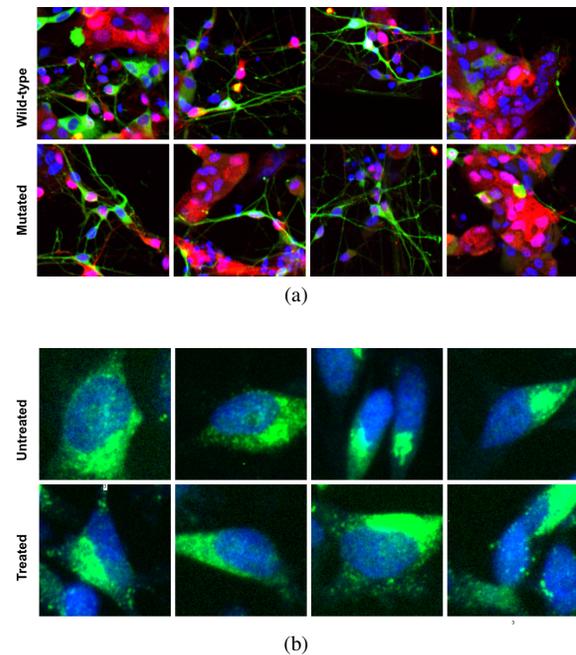


Figure 1. **Top:** Real images from the LRRK2 dataset, displaying wild-type images in the first row and images of mutated neurons in the second row. **Bottom:** Real images from the Golgi dataset, with untreated images in the first row and Nocodazole-treated images in the second row. In both (a) and (b), identifying and interpreting differences between the two classes by eye is highly challenging. However, it is essential for understanding the disease in (a) and assessing the treatment effects in (b)

### 3. Method

In this section we first provide an overview of DMs and the  
140 methods used for fine-tuning them, then we dive into the  
141 details of our approach.  
142

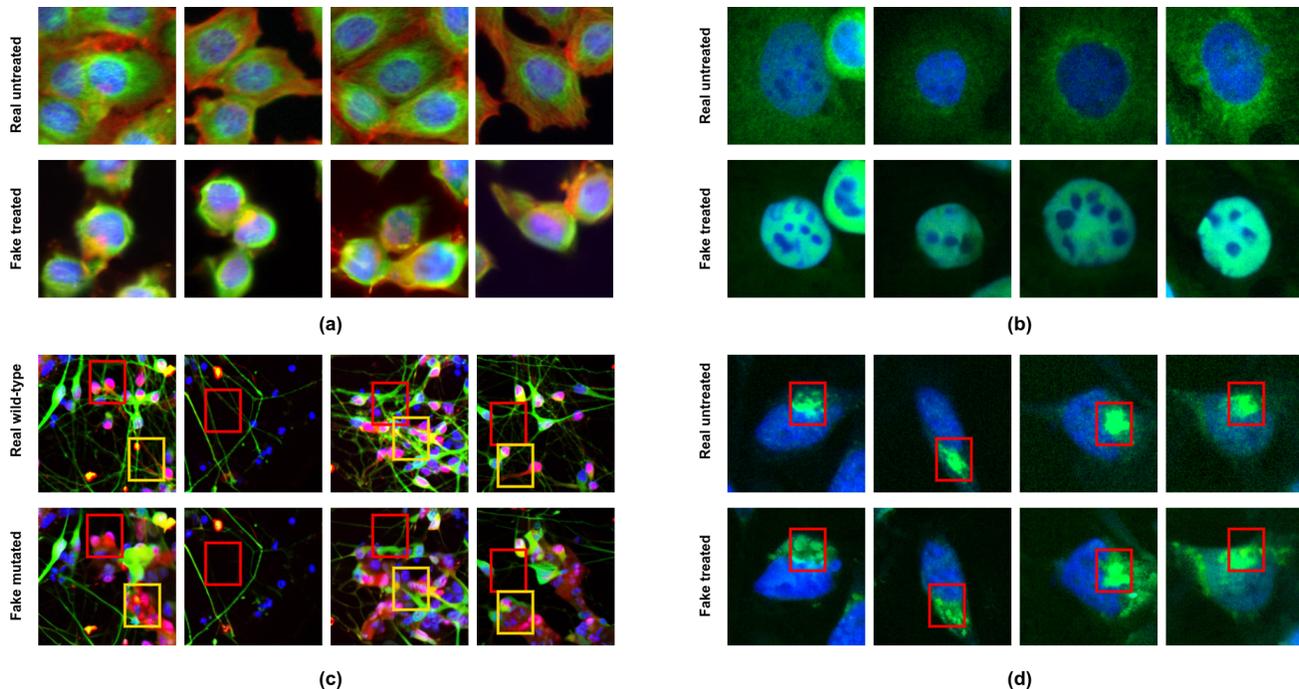


Figure 2. We fine-tuned diffusion models on four different microscopy image datasets and performed translations from the source class to the target class. We observed the following: In (a), the translated images of untreated BBBC021 samples successfully replicated the effects of Latrunculin B treatment, where we observed a decrease in cell count and the disappearance of the cytoplasmic skeleton, likely due to the toxicity of the treatment. In (b), TNF treatment on cells and its translocation effect was well recapitulated by image translation. In (c), we translated images of wild-type cells to images of LRRK2 mutated cells and noticed a reduction in neuron density and complexity (red squares) and an increase of  $\alpha$ -synuclein (yellow squares), recapitulating known effects of the mutation. Finally, in (d), we observed the correct replication of the effect of Nocodazole treatment causing the scattering of the Golgi apparatus (red squares). Note how pronounced ((a), (b)) as well as subtle ((c), (d)) phenotypic changes are well captured by our model. In any case seeing the same cell before and after treatment allowed us to assess the effect of the perturbation. Real images of both conditions of the four datasets can be seen in Appendix A.1.

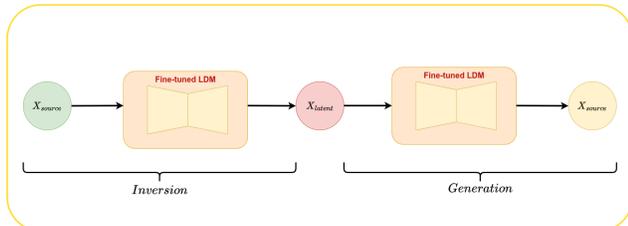


Figure 3. Phen-LDiff leverages fine-tuned LDMs to perform image-to-image translation, identifying phenotypic variations between the images of two conditions. First, a fine-tuned model is used to invert an image from the source class into a latent code, which is then used to generate an image in the target class.

### 143 3.1. Background

#### 144 3.1.1 Diffusion Models

145 Denoising Diffusion Probabilistic Models (DDPMs) are latent  
146 variable models that utilize two Markov processes: a  
147 fixed forward process that gradually adds noise to the data,

and a learned reverse process that removes noise to recover  
the data distribution. Formally, given data  $x_0 \sim q(x_0)$ , the  
forward process iteratively adds Gaussian noise over  $T$  time  
steps following a forward transition kernel given by:

$$q(x_t, |, x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}) \quad (1)$$

In the reverse process, noise is gradually removed using a  
learnable transition kernel:

$$p_\theta(x_{t-1}, |, x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad (2)$$

While DDPMs generate high-quality images, they require  
many iterations during inference, making the process computationally  
intensive. To accelerate inference, *Denoising Diffusion Implicit Models* (DDIMs) [40] can be employed. Notably, DDIMs offer deterministic sampling, allowing for *exact* inversion, a property that is crucial for our approach to observe phenotypic changes in real images.

Latent Diffusion Models (LDMs) [35] extend DDPMs by introducing a latent space to improve both efficiency and

165 flexibility in high-dimensional data generation tasks. In-  
166 stead of operating directly in the data space, LDMs learn to  
167 encode images into a lower-dimensional latent space  $\mathcal{E}(x)$ ,  
168 where the diffusion process occurs. This significantly re-  
169 duces computational overhead, as the diffusion steps are  
170 performed on a smaller latent representation rather than on  
171 the full-resolution image. This approach not only acceler-  
172 ates inference but also makes it feasible to train LDMs on  
173 very large datasets.

$$174 \quad L_{LDM} = \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_{\theta}(z_t, t, c)\|_2^2 \right] \quad (3)$$

175 where:  $\mathcal{E}$  is the encoder,  $c$  is the condition and  $\epsilon_{\theta}$  is the  
176 parameterized noise predictor.

### 177 3.1.2 Low Rank Adaptation (LoRA)

178 Low-Rank Adaptation [20] is a technique designed to effi-  
179 ciently fine-tune large pre-trained models by significantly  
180 reducing the number of trainable parameters. Instead of  
181 updating the entire weight matrix  $W$  during fine-tuning,  
182 LoRA introduces trainable low-rank matrices to approxi-  
183 mate the weight updates. Specifically, the weight update  
184  $\Delta W$  is decomposed into a product of two low-rank mat-  
185 rices  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{k \times r}$ , where  $r \ll \min(d, k)$ . The  
186 adapted weight matrix during training is computed as fol-  
187 lows:

$$188 \quad W' = W + BA^{\top} \quad (4)$$

189 This method can be either applied to all or a subset of the  
190 model layers.

### 191 3.1.3 SVDiff

192 SVDiff is a method developed to efficiently fine-tune large  
193 diffusion models by performing a **singular value decom-**  
194 **position (SVD)** on the weight matrices  $W$ .

$$195 \quad W = U\Sigma V^{\top}$$

196 During fine-tuning, instead of updating the entire weight  
197 matrix  $W$ , SVDiff updates only the singular values of this  
198 matrix. This significantly reduces the number of param-  
199 eters that need to be trained, leading to faster training times  
200 and reduced computational resources. By operating in this  
201 lower-dimensional space, SVDiff helps prevent overfitting  
202 and makes it more practical to adapt large diffusion models  
203 to specific tasks or datasets.

### 204 3.2. Datasets

205 In this work, we evaluated the proposed method on several  
206 biological datasets. In some of them, cell variations are pro-  
207 nounced to showcase our approach, while in others, the dif-  
208 ferences are more subtle illustrating the usefulness of the  
209 method to display them. The datasets used are as follows:

**BBBC021:** The BBBC021 dataset [10] is a publicly avail- 210  
able collection of fluorescent microscopy images of MCF-7, 211  
a breast cancer cell line treated with 113 small molecules at 212  
eight different concentrations. For our research, we focused 213  
on images of untreated cells and cells treated with the high- 214  
est concentration of the compound Latrunculin B. In Fig. 2, 215  
the green, blue and red channels label for B-tubulin, DNA 216  
and F-actin respectively. 217

**Golgi:** Fluorescent microscopy images of HeLa cells un- 218  
treated (DMSO) and treated with Nocodazole. In Fig. 8b, 219  
the green and blue channels label for B-tubulin and DNA 220  
respectively. 221

**LRKK2:** This dataset contains images of dopaminergic 222  
neurons derived from iPSCs reprogrammed from fibroblasts 223  
of a Parkinson’s disease patient affected by the LRRK2- 224  
G2019S mutation. It also includes images where the muta- 225  
tion was genetically corrected using CRISPR-cas9, provid- 226  
ing a rescued isogenic control [27]. In Fig. 8b the bleu, 227  
green and red label for DNA, dopaminergic neurons and 228  
alpha-synuclein (SNCA) respectively. 229

**Translocation:** Fluorescent microscopy images depicting 230  
the subcellular localization of the NF $\kappa$ B (nuclear factor 231  
kappa B) protein, either untreated or treated with TNF $\alpha$  (the 232  
pro-inflammatory cytokine tumor necrosis factor alpha). In 233  
Fig 2 (b), the blue and green channels labels for DNA and 234  
NF $\kappa$ B protein respectively. 235

### 236 3.3. Proposed Approach

237 In this work, we introduce **Phen-LDiff**, a method that 238  
leverages pre-trained Latent Diffusion Models (LDMs) for 239  
image-to-image translation on small biological datasets to 240  
identify phenotypic differences. Our approach begins by 241  
conditionally fine-tuning a general-purpose LDM on mi- 242  
croscopy images from different experimental conditions 243  
(e.g., treated vs. untreated, wild-type vs. mutant, as illus- 244  
trated in Fig.1). To perform the translation from one class 245  
to another, we first invert an image from the source class 246  
into its latent representation, which is then used to generate 247  
a corresponding image in the target class.

### 248 4. Results

249 In this work, we utilized Stable Diffusion 2, which was pre- 250  
trained on the LAION-5B dataset [36]. LAION-5B is a 251  
large-scale collection of web-scraped image-text pairs, en- 252  
compassing a wide variety of general image sources across 253  
the internet. We fine-tuned this model on the BBBC021 254  
dataset using several strategies: (1) full fine-tuning, where 255  
all model parameters are updated; (2) attention fine-tuning, 256  
where only the attention layers of the model are modified; 257  
and (3) LoRA and SVDiff, two techniques designed to ef- 258  
ficiently reduce the number of trainable parameters while 259  
preserving model performance.

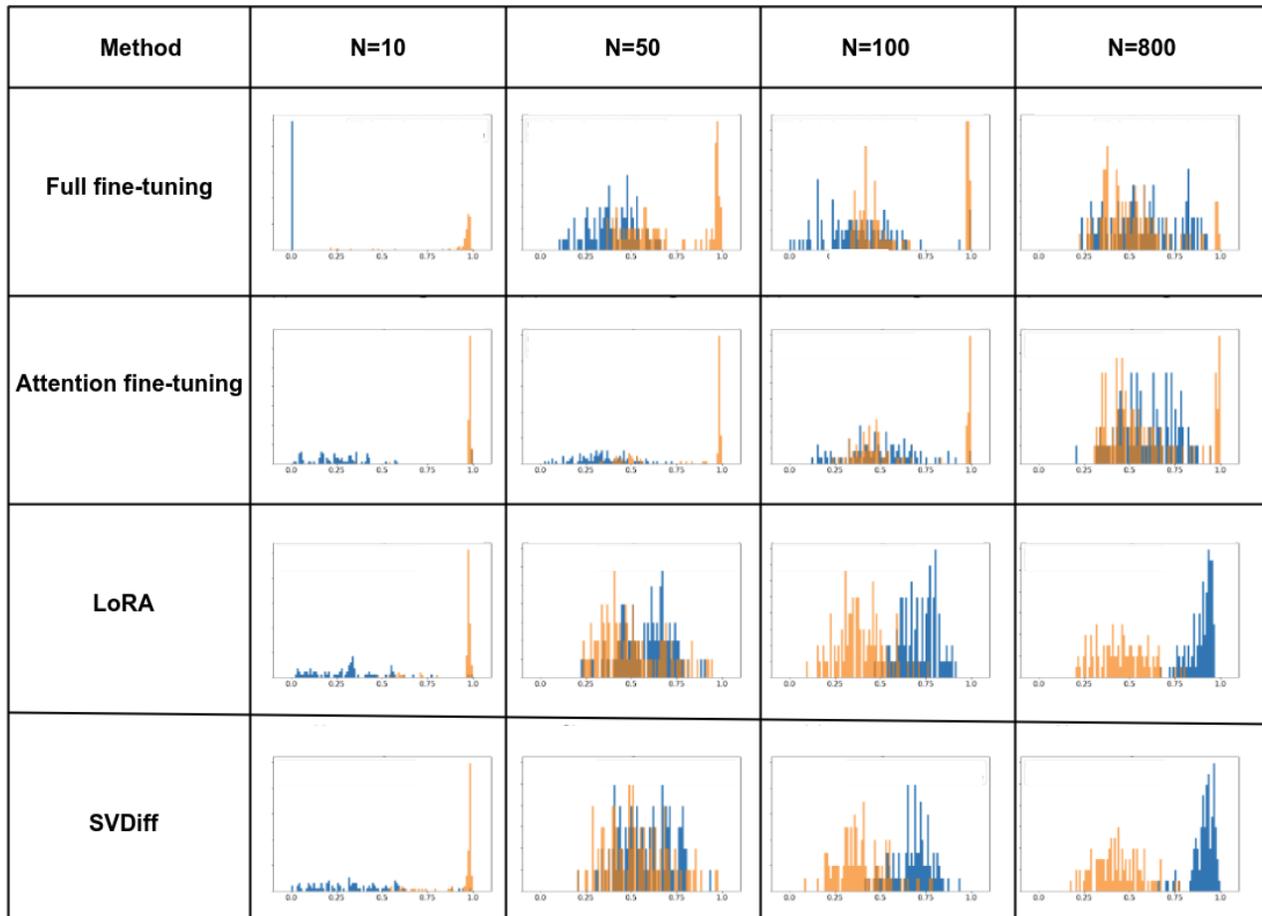


Figure 4. Visualizing the **generalization** and **memorization** of fine-tuned diffusion models on subsets of different sizes from the BBBC021 dataset. Each plot shows two histograms: the blue histogram represents the cosine similarity between images generated using the same seed by two fine-tuned models trained on distinct, **non-overlapping** subsets of the same size. If the model has achieved generalization, the **blue** histogram should be close to one, indicating that the two images generated by the models are very similar. The orange histogram represents the cosine similarity between a generated sample and its closest image from the training dataset. A well-generalized model would produce an **orange** histogram far from one, indicating that the generated images have low similarity to any specific training example.

260

#### 4.1. Domain adaptation of fine-tuned LDMs

261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271

As shown in Fig. 5, the fine-tuned Stable Diffusion 2 model demonstrates the ability to generate high-quality biological images. This highlights the model’s capability to shift its original distribution, from natural images to those closely aligned with the specific characteristics of biological data. Furthermore, the results indicate that the generated images maintain good quality across various biological datasets, even when trained on a limited number of images (100 images per dataset in our case). This suggests that pre-trained models can be effectively leveraged to learn new biological image distributions, even with a small training dataset.

#### 4.2. Assessing generalization and memorization in fine-tuned LDMs

272  
273

Recently, some studies have observed that diffusion models can memorize samples from the training set, leading to their replication during inference [6, 39]. This behavior was particularly noted in [24], where diffusion models trained on small datasets exhibited memorization. In contrast, it was demonstrated that the same models do not exhibit this memorization when trained on sufficiently large datasets. To ensure that our fine-tuned models do not merely memorize the training datasets but instead learn the underlying distribution of the images, we adopted the approach proposed in [24]. Specifically, we fine-tuned two models using two non-overlapping subsets from the same datasets

274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285

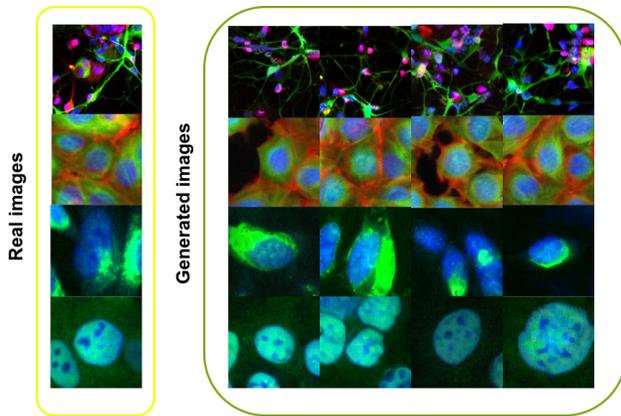
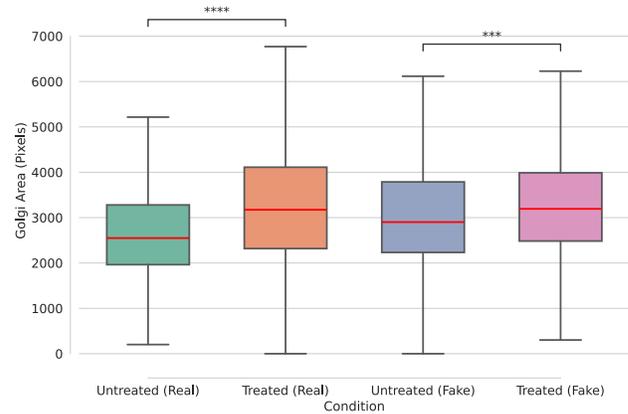


Figure 5. The images generated by a diffusion model fine-tuned on 100 images using LoRA on different biological datasets, we can see that the generated samples resemble the real ones.

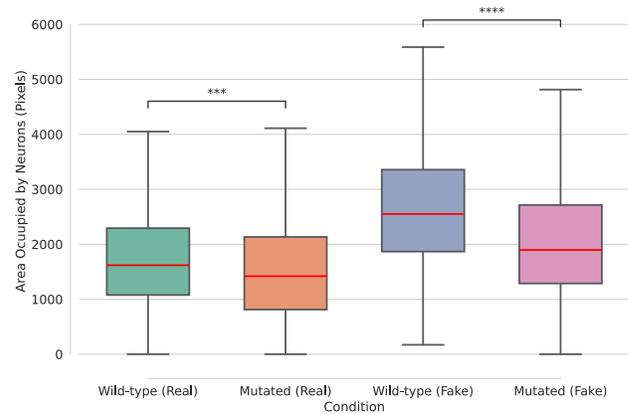
286 (thus two different samples from the same distribution) and  
 287 measured the cosine similarity between images generated  
 288 from the same seed, as well as the correlation between each  
 289 generated image and its closest match from the training  
 290 dataset. This evaluation was conducted across four different  
 291 fine-tuning methods: full fine-tuning, attention fine-tuning,  
 292 SVDiff, and LoRA, as illustrated in Fig. 4. From the results,  
 293 we observe that with only 10 training images, all fine-tuning  
 294 methods tend to memorize the training dataset, resulting in  
 295 high correlation values between the generated images and  
 296 the closest ones from the training set. Furthermore, we notice  
 297 that full and attention fine-tuning struggle to generalize  
 298 effectively, even as the number of training images increases.  
 299 In contrast, for LoRA and SVDiff, we see that with just 50  
 300 training images, the blue and orange histograms begin to  
 301 shift toward 1 and 0, respectively, indicating greater general-  
 302 ization and reduced memorization. Although no significant  
 303 differences were observed in the quality of the generated  
 304 images across the methods, we chose to use LoRA for  
 305 the remaining experiments due to the more optimized and  
 306 faster implementation available to us.

### 307 4.3. Identifying subtle cellular variations with 308 image-to-image translation

309 So far, we have demonstrated that fine-tuning Latent Diffu-  
 310 sion Models (LDMs) is feasible even on limited biological  
 311 datasets. However, our primary goal is to detect subtle cel-  
 312 lular variations in biological samples. In Fig. 2, we illustrate  
 313 the image-to-image translation performed on small datasets:  
 314 100 images per class for BBBC021, Golgi, and LRRK2,  
 315 and for translocation. In Fig. 2 (a) and (b), the effects of  
 316 treatment are visible. Specifically, for the BBBC021 dataset  
 317 Fig. 2 (a), the phenotypic changes induced by Latrunculin B  
 318 are evident. The actin cytoskeleton (red channel) has largely  
 319 disappeared and a significant decrease in cell count is ob-



(a) The measurement of the Golgi apparatus area performed on real and synthetic images for both conditions indicates a difference in the area occupied by the Golgi apparatus, confirming the observation made by Phen-LDiff. Specifically, it appears more scattered in the treated case, which explains its larger size.



(b) The measurement of the area occupied by neurons (green channel) on real and synthetic images for both conditions indicates a reduced neuron count in the mutated case, confirming the observation made by Phen-LDiff. Indeed, the mutation that causes Parkinson's disease leads to a decrease in both the number and complexity of neurons

Figure 6. An image analysis measurement using CellProfiler [9] on the Golgi and LRRK2 datasets, performed on real and synthetic images for both conditions, led to the same quantitative conclusions, indicating that Phen-LDiff can detect subtle cellular variations in models fine-tuned on datasets with as few as 100 images per class.

served, indicating the toxicity of the treatment. In Fig. 2  
 (b), upon treatment with  $TNF\alpha$ , the transcription factor  
 translocates to the nucleus, causing the fluorescence signal  
 to shift from the cytoplasm to the nuclear region, resulting  
 in cells displaying brightly fluorescent green nuclei. These  
 phenotypic changes are prominent and easily recognizable.  
 Conversely, the second row showcases more subtle pheno-  
 types, which may be challenging to detect, even for special-

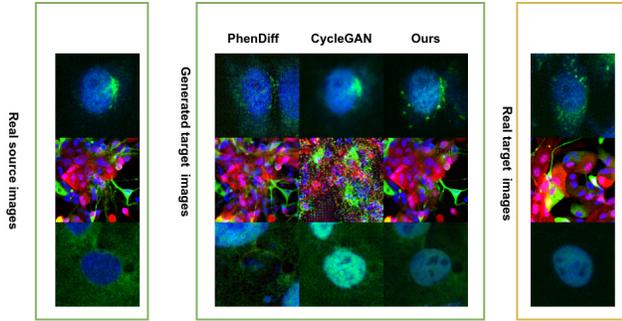


Figure 7. We translated real untreated (Wild-type) images to the treated (mutated) condition using PhenDiff, CycleGAN, and PhenLDiff, all the models were trained on datasets of 100 images. For PhenDiff, we can see that the translated images do not resemble the cell images in the source class but are rather new samples from the target distribution than translated cells. For CycleGAN, the translated images are very similar to the source class, but the quality is somewhat lower and the image does not recapitulate well the target class phenotype. In contrast, for the images translated with our method, we can see that they produce the desired phenotypes for the cells that were present in the provided image from the source class, indicating a successful translation.

328 ists. For instance, in Fig.2(d), untreated cell images from  
 329 the Golgi dataset were translated to resemble treated cells.  
 330 Changes in Golgi apparatus morphology due to Nocodazole  
 331 treatment are noticeable, with the apparatus fragmenting  
 332 into smaller stacks. In Fig.2(c), when translating rescued  
 333 WT images to diseased ones, we observed a decrease  
 334 in dopaminergic neurons and dendritic complexity, as well  
 335 as an increase in alpha-synuclein (red channel), more ex-  
 336 amples of translations can be found in Appendix A.2. To  
 337 confirm these subtle observations, we used CellProfiler [9]  
 338 to quantify the changes detected by **Phen-LDiff**. For exam-  
 339 ple, to confirm that the Golgi apparatus is more scattered in  
 340 the treated case, we measured the area it occupies in both  
 341 conditions. Similarly, for the LRRK2 dataset, we measured  
 342 the area occupied by neurons (green channel) in both syn-  
 343 thetic and real image. In Fig. 6, the measurements align  
 344 with the observed changes spotted by **Phen-LDiff**. Indeed,  
 345 there is a significant difference between the measurements  
 346 in the treated (WT) versus treated (mutated) cases, suggest-  
 347 ing that we are identifying meaningful changes. All these  
 348 now-visible differences can assist biologists in better un-  
 349 derstanding these diseases and the effects of treatments.

#### 350 4.4. Comparing our method to the existing ones

351 Using generative models to identify cellular variations is a  
 352 growing area of research due to their potential in advancing  
 353 biological studies [2, 3, 27]. Although methods like Phen-  
 354 Explain [27] can identify these variations in synthetic im-  
 355 ages, they struggle with real images due to the difficulty

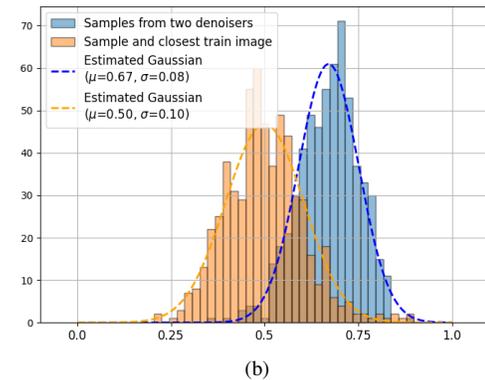
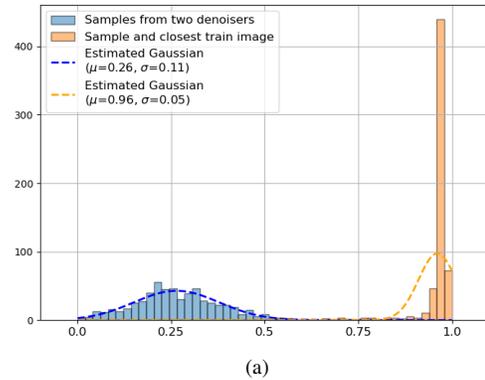


Figure 8. In this figure, we trained both PhenDiff and PhenLDiff on a subset of 50 images from the BBBC021 dataset. **Top:** The memorization histogram is close to 1, indicating very strong memorization for PhenDiff. **Bottom:** Phen-LDiff shows less memorization and achieves better generalization compared to PhenDiff.

of inverting images using GANs. This challenge was over-  
 come in PhenDiff [3] by leveraging the inversion properties  
 of DDIM. However it still necessitated large datasets that  
 are hard to get in biology. Our approach proposes the use of  
 a pretrained latent diffusion model to enable effective per-  
 formance even with limited data availability.

We compare our method to two representative models:  
 PhenDiff, which uses diffusion models (DMs) trained from  
 scratch, and CycleGAN [44], which is based on GANs.  
 As shown in Fig. 7, our method effectively highlights phe-  
 notypic cellular changes induced by the target conditions.  
 Specifically, the Golgi apparatus appears more scattered,  
 there is an increase in  $\alpha$ -synuclein, and the transcription fac-  
 tor translocates to the nucleus in the translocation datasets.  
 These observations are less apparent with PhenDiff and Cy-  
 cleGAN. For instance, in CycleGAN, the translation quality  
 is lower, likely due to limited data, which makes learning  
 the target distribution challenging. In the case of PhenDiff,  
 although some phenotypic variations are reconstructed, the  
 translated images differ substantially from the original ones,

Table 1. Performance Metrics Across Different Datasets to evaluate

Method	BBBC021		Translocation		LRKK2		Golgi	
	FID	Cycle loss	FID	Cycle loss	FID	Cycle loss	FID	Cycle loss
CycleGAN	75.98	<b>528.83</b>	40.56	<b>643.12</b>	71.23	<b>428.48</b>	32.28	<b>341.36</b>
Phendiff	33.31	2555.38	60.65	1704.54	74.23	2633.73	<b>23.66</b>	958
Ours	<b>24.30</b>	1707.38	<b>32.79</b>	1021	<b>18.57</b>	923.98	30.31	773.26

376 making direct comparison with real images difficult. Addi-  
 377 tional translation examples are provided in the Appendix  
 378 B.1.

379 To quantitatively compare the performance of each trans-  
 380 lation method, we evaluated the quality of the translated  
 381 images using FID [17] and assessed similarity to the origi-  
 382 nal images using cycle loss. For the cycle loss, an image  
 383 is translated from the original domain to the target domain  
 384 and back, and we compute the  $L_2$  norm between the origi-  
 385 nal and reconstructed images. As shown in Table 1, our  
 386 method achieves a better FID score on almost all datasets.  
 387 However, CycleGAN shows a lower cycle consistency loss  
 388 while producing lower-quality translations compared to the  
 389 other models. This is primarily due to the cycle consis-  
 390 tency loss used in CycleGAN training, which helps in re-  
 391 constructing images but fails to produce accurate translation  
 392 and thus identify phenotypic changes. Our method offers  
 393 the best trade-off between capturing phenotypic variations  
 394 and maintaining proximity to the initial target distribution.

395 To better understand the good translation performance of  
 396 our method, we compared the memorization and generaliza-  
 397 tion abilities of PhenDiff and our model on 50 images per  
 398 class from the BBBC021 dataset. Following the same strat-  
 399 egy as previously described, generalization was assessed by  
 400 calculating the cosine similarity between images generated  
 401 from the same seed by two models trained on two independ-  
 402 ent datasets of 50 images each. Memorization was evalu-  
 403 ated by calculating the cosine similarity between a gener-  
 404 ated image and its closest match from the training dataset.  
 405 In Fig. 8, we can clearly see that PhenDiff falls into a mem-  
 406 orization regime, whereas Phen-LDiff shows less memo-  
 407 rization and greater generalization. Further comparisons us-  
 408 ing other datasets and sizes are presented in Appendix B.2.  
 409 These results suggest that fine-tuned models achieve better  
 410 generalization in low-data regimes, which explains the good  
 411 translation performance of our method.

412 Additionally, we compared the training time of Phen-  
 413 Diff and Phen-LDiff on two NVIDIA L40S GPUs using  
 414 the BBBC021 dataset. Training took approximately 6 hours  
 415 for PhenDiff and around 2 hours for Phen-LDiff. This dif-  
 416 ference would be even more significant with larger train-  
 417 ing images, demonstrating the computational efficiency of  
 418 Phen-LDiff.

## 5. Conclusion 419

420 In this work, we propose **Phen-LDiff**, a method for  
 421 image-to-image translation using fine-tuned Latent Diffu-  
 422 sion Models (LDMs) to identify phenotypic variations from  
 423 limited microscopy data. Our approach demonstrates that  
 424 LDMs can be effectively fine-tuned on biological datasets,  
 425 capturing their underlying distributions even when data is  
 426 limited. We found that certain fine-tuning approaches, such  
 427 as full model fine-tuning and attention fine-tuning, can lead  
 428 to memorization. In contrast, methods like LoRA and SVD-  
 429 iff promote better generalization, even with small datasets  
 430 containing as few as 100 images per class. Our method en-  
 431 ables image-to-image translation by first inverting an im-  
 432 age into a latent space, followed by conditional genera-  
 433 tion to highlight phenotypic variations between conditions.  
 434 We tested this approach across multiple biological datasets,  
 435 showing its capability to reveal both apparent and subtle  
 436 differences between experimental conditions. When com-  
 437 pared to other representative methods, Phen-LDiff outper-  
 438 formed them in translation quality, even with limited image  
 439 datasets. Furthermore, our method avoids memorization  
 440 and is computationally more efficient than diffusion mod-  
 441 els trained from scratch, reducing training time significantly  
 442 without compromising quality.

443 We anticipate that Phen-LDiff can contribute to biolog-  
 444 ical research and drug discovery by enabling experts to  
 445 gain deeper insights into disease mechanisms and treatment  
 446 effects, especially in low-data regimes where traditional  
 447 methods struggle. This efficiency and ability to generalize  
 448 make Phen-LDiff a promising tool for advancing precision  
 449 in phenotypic analysis.

## References 450

- 451 [1] Samyadeep Basu, Daniela Massiceti, Shell Xu Hu, and So-  
 452 heil Feizi. Strong baselines for parameter efficient few-shot  
 453 fine-tuning, 2023. 2
- 454 [2] Anis Bourou and Auguste Genovesio. Unpaired image-to-  
 455 image translation with limited data to reveal subtle pheno-  
 456 types. In *IEEE 20th International Symposium on Biomedical  
 457 Imaging (ISBI)*, 2023. 1, 2, 7
- 458 [3] Anis Bourou, Thomas Boyer, Marzieh Gheisari, Kévin  
 459 Daupin, Véronique Dubreuil, Aurélie De Thonel, Valérie  
 460 Mezger, and Auguste Genovesio. PhenDiff: Revealing Sub-

- 461 tle Phenotypes with Diffusion Models in Real Images . In  
462 *proceedings of Medical Image Computing and Computer As-*  
463 *sisted Intervention – MICCAI*, 2024. 1, 2, 7
- [4] Anis Bourou, Valérie Mezger, and Auguste Genovesio. Gans  
464 conditioning methods: A survey, 2024. 1  
465
- [5] Nicola De Cao and Thomas Kipf. Molgan: An implicit generative  
466 model for small molecular graphs, 2022. 1  
467
- [6] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski,  
468 Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito,  
469 and Eric Wallace. Extracting training data from diffusion models,  
470 2023. 5  
471
- [7] Srinivas Niranj Chandrasekaran Chandrasekaran, Hugo Ceulemans,  
472 Jeffrey D Boyd, and Anne E Carpenter. Image-based profiling for  
473 drug discovery: due for a machine-learning upgrade? *Nature Reviews Drug Discovery*,  
474 20:145–159, 2021. 1, 2  
475
- [8] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on  
476 image synthesis, 2021. 1  
477
- [9] Anne Carpenter et al. Cellprofiler: Image analysis software for  
478 identifying and quantifying cell phenotypes. *Genome biology*,  
479 7:R100, 2006. 6, 7  
480
- [10] Peter D. Caie et al. High-Content Phenotypic Profiling of Drug  
481 Response Signatures across Distinct Cancer Cells. *Molecular Cancer Therapeutics*,  
482 9(6):1913–1926, 2010. 4  
483
- [11] Ruben Fonnegra, Mohammad Sanian, Zitong Chen, Lassi Paavola,  
484 and Juan Caicedo. Analysis of cellular phenotypes with image-based  
485 generative models. In *NeurIPS 2023 Generative AI and Biology (GenBio) Workshop*,  
486 2023. 2  
487
- [12] Jonathan Frankle, David J. Schwab, and Ari S. Morcos. Training  
488 batchnorm and only batchnorm: On the expressive power of random  
489 features in cnns, 2021. 2  
490
- [13] C. Gao et al. Synthetic data accelerates the development of  
491 generalizable learning-based algorithms for x-ray image analysis.  
492 *Nature Machine Intelligence*, 5:294–308, 2023. 1  
493
- [14] Harshvardhan GM, Mahendra Kumar Gourisaria, Manjusha  
494 Pandey, and Siddharth Swarup Rautaray. A comprehensive survey  
495 and analysis of generative models in machine learning. *Computer Science Review*,  
496 38:100285, 2020. 1  
497
- [15] Jie Gui, Tuo Chen, Jing Zhang, Qiong Cao, Zhenan Sun, Hao  
498 Luo, and Dacheng Tao. A survey on self-supervised learning: Algorithms,  
499 applications, and future trends, 2023. 2  
500
- [16] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris  
501 Metaxas, and Feng Yang. Svdiff: Compact parameter space for  
502 diffusion fine-tuning. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*,  
503 pages 7289–7300, 2023. 2  
504
- [17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard  
505 Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update  
506 rule converge to a local nash equilibrium, 2018. 8  
507
- [18] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance,  
508 2022. 1  
509
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion  
510 probabilistic models, 2020. 1  
511
- [20] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu,  
512 Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.  
513 Lora: Low-rank adaptation of large language models, 2021. 2, 4  
514
- [21] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros.  
515 Image-to-image translation with conditional adversarial networks, 2018. 2  
516
- [22] John Jumper, Richard Evans, Alexander Pritzel, Tim Green,  
517 Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool,  
518 Russ Bates, Augustin Zidek, Aidan Potapenko, Anna Bridgland,  
519 Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie,  
520 Basil Romera-Paredes, Stanislaw Nikolov, Rishub Jain, Jonas Adler,  
521 Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski,  
522 Mikolaj Steinegger, Martin Pacholska, Tomás Berghammer, Sebastian  
523 Bodenstein, David Silver, Oriol Vinyals, Alex J. W. Senior, Koray  
524 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate  
525 protein structure prediction with AlphaFold. *Nature*, 596:583–589,  
526 2021. 1  
527
- [23] E. Jung, M. Luna, and S. H. Park. Conditional gan with 3d  
528 discriminator for mri generation of alzheimer’s disease progression.  
529 *Pattern Recognition*, 133:109061, 2023. 1  
530
- [24] Zahra Kadhodaie, Florentin Guth, Eero P. Simoncelli, and Stéphane  
531 Mallat. Generalization in diffusion models arises from geometry-  
532 adaptive harmonic representations, 2024. 5  
533
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko  
534 Lehtinen, and Timo Aila. Analyzing and improving the image quality  
535 of StyleGAN. In *Proc. CVPR*, 2020. 2  
536
- [26] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating  
537 the design space of diffusion-based generative models, 2022. 1  
538
- [27] Alexis Lamiable, Tiphaine Champetier, Francesco Leonardi, Ethan  
539 Cohen, Peter Sommer, David Hardy, Nicolas Argy, Achille Massoug-  
540 bodji, Elaine Del Nery, Gilles Cottrell, Yong-Jun Kwon, and Auguste  
541 Genovesio. Revealing invisible cell phenotypes with conditional generative  
542 modeling. *Nature Communications*, 14(1):6386, 2023. 1, 2, 4, 7  
543
- [28] Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno,  
544 Leon Hetzel, Yuge Ji, Ignacio L Ibarra, Sanjay R Srivatsan, Mohsen  
545 Naghypourfar, Riza M Daza, Beth Martin, Jay Shendure, Jose L  
546 McFaline-Figueroa, Pierre Boyeau, F Alexander Wolf, Nafissa Yakubova,  
547 Stephan Günemann, Cole Trapnell, David Lopez-Paz, and Fabian J  
548 Theis. Predicting cellular responses to complex perturbations in high-  
549 throughput screens. *Molecular Systems Biology*, 19(6), 2023. 1  
550
- [29] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu,  
551 Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language  
552 models: A survey, 2024. 2  
553
- [30] Taehong Moon, Moonseok Choi, Gayoung Lee, Jung-Woo Ha, and  
554 Juho Lee. Fine-tuning diffusion models with limited data. In *NeurIPS  
555 2022 Workshop on Score-Based Methods*, 2022. 2  
556
- [31] Nikita Moshkov, Michael Bornholdt, Santiago Benoit, Matthew  
557 Smith, Claire McQuin, Allen Goodman, Rebecca Senft, Yu Han,  
558 Mehrtash Babadi, Peter Horvath, Beth A. Cimini, Anne E. Carpenter,  
559 Shantanu Singh, and Juan C. Caicedo. Learning representations for  
560 image-based profiling of perturbations. *bioRxiv*, 2022. 1, 2  
561

- 575 [32] Maria Popova, Olexandr Isayev, and Alexander Tropsha.  
576 Deep reinforcement learning for de novo drug design. *Sci-*  
577 *ence Advances*, 4(7):eaap7885, 2018. 1
- 578 [33] Danilo Rezende and Shakir Mohamed. Variational inference  
579 with normalizing flows. In *Proceedings of the 32nd Interna-*  
580 *tional Conference on Machine Learning*, pages 1530–1538,  
581 Lille, France, 2015. PMLR. 1
- 582 [34] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lih  
583 Zelnik-Manor. Imagenet-21k pretraining for the masses,  
584 2021. 2
- 585 [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
586 Patrick Esser, and Björn Ommer. High-resolution image syn-  
587 thesis with latent diffusion models, 2022. 1, 2, 3
- 588 [36] Christoph Schuhmann, Romain Beaumont, Richard Vencu,  
589 Cade Gordon, Ross Wightman, Mehdi Cherti, Theo  
590 Coombes, Aarush Katta, Clayton Mullis, Mitchell Worts-  
591 man, Patrick Schramowski, Srivatsa Kundurthy, Katherine  
592 Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia  
593 Jitsev. Laion-5b: An open large-scale dataset for training  
594 next generation image-text models, 2022. 2, 4
- 595 [37] Nikolay Simidjievski, Cristian Bodnar, Imran Tariq, Philipp  
596 Scherer, Henrik Andres Terre, Zahra Shams, Mateja Jam-  
597 nik, and Pietro Liò. Variational autoencoders for cancer data  
598 integration: Design principles and computational practice.  
599 *Frontiers in Genetics*, 10:1205, 2019. 1
- 600 [38] Karen Simonyan and Andrew Zisserman. Very deep convo-  
601 lutional networks for large-scale image recognition, 2015. 2
- 602 [39] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas  
603 Geiping, and Tom Goldstein. Diffusion art or digital forgery?  
604 investigating data replication in diffusion models, 2022. 5
- 605 [40] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denois-  
606 ing diffusion implicit models, 2022. 1, 3
- 607 [41] Alexey Strokach and Philip M. Kim. Deep generative mod-  
608 eling for protein design. *Current Opinion in Structural Biol-*  
609 *ogy*, 72:226–236, 2022. 1
- 610 [42] Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang,  
611 Chao Yang, and Chunfang Liu. A survey on deep transfer  
612 learning. In *Artificial Neural Networks and Machine Learn-*  
613 *ing – ICANN 2018*, pages 270–279, Cham, 2018. Springer  
614 International Publishing. 2
- 615 [43] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit:  
616 Simple parameter-efficient fine-tuning for transformer-based  
617 masked language-models, 2022. 2
- 618 [44] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A.  
619 Efros. Unpaired image-to-image translation using cycle-  
620 consistent adversarial networks, 2020. 2, 7