TRANSCRIPTOMICS-MORPHOLOGY GENERATION VIA TREATMENT CONDITIONING WITH RECTIFIED FLOW

Anonymous authors

000

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026027028

029

040 041

042 043

044

045

046 047

048

051

052

Paper under double-blind review

ABSTRACT

Predicting cellular responses to drug perturbations requires capturing complex dependencies between transcriptomic and morphological changes that singlemodality approaches cannot adequately model. We introduce **PertFlow**, the first unified framework that jointly predicts gene expression profiles and generates cellular morphology images in response to drug treatments, conditioned on control cellular states. Our method integrates control transcriptomic and imaging data through multi-head cross-modal attention mechanisms, learning a shared latent representation that incorporates drug compound features, background cellular profiles, and treatment specifications. From this unified representation, PertFlow employs a regression head for RNA-seq prediction and rectified flow dynamics for stable morphological image generation, with cross-modal consistency losses ensuring coherent molecular and phenotypic predictions. PertFlow enables accurate predictions from either complete multi-modal inputs or single-modality data alone, demonstrating robust cross-modal learning. Our evaluation on paired RNA-seq and Cell Painting fluorescent imaging datasets demonstrates that Pert-Flow achieves stronger cross-modal consistency and accurate prediction of druginduced changes compared to diffusion baselines.

1 Introduction

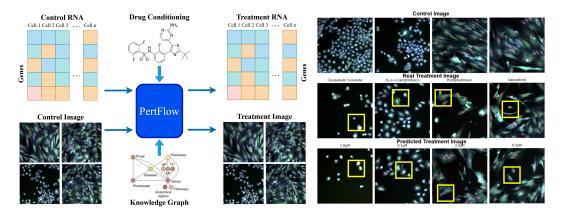


Figure 1: (LEFT) Cross-modal mapping from control RNA-seq and image to treatment RNA-seq and image with drug conditioning through PertFlow. (RIGHT) Comparison of generated treatment vs real treatment images with drug name and concentration. Yellow boxes indicate similar features.

Understanding how drugs alter cellular states is essential for drug discovery, mechanistic understanding, and personalized medicine. Traditional drug response models typically focus on either transcriptomic data or imaging, missing the complex interdependencies between molecular and morphological changes that occur simultaneously in cells. Recent advances in high-throughput profiling now allow paired RNA sequencing and imaging, offering complementary insights: transcriptomics captures molecular mechanisms and gene regulation, while morphology reflects structural and phenotypic changes. These modalities are linked as gene expression can drive morphological transformations, and structural changes can modulate gene activity yet most models treat them in isolation.

Existing methods fall short as transcriptomics-based approaches cannot model morphological effects; image-based models lack molecular interpretability; and cross-modal predictors generate only one modality from another, without joint modeling. Moreover, most studies prioritize genetic over chemical perturbations and analyze rather than predict multi-modal responses. Joint generation of multi-modal responses poses three main challenges: (1) aligning transcriptomic and morphological data across fundamentally different representational spaces; (2) capturing complex drug conditioning involving compound, dose, cell type, and timepoint; and (3) simultaneously predicting discrete gene expression and continuous image data with biological realism and cross-modal consistency.

We introduce **PertFlow** (Figure 1), a novel unified generative framework for jointly predicting treatment gene expression and synthesizing cellular morphology from control conditions, conditioned on drug metadata. Our contributions are: (1) First method to jointly predict transcriptomic and generate morphological responses to chemical perturbations. (2) A shared embedding space integrating control RNA-seq, control images, and drug metadata to model complex dependencies. (3) Multi-token cross-attention to align molecular and morphological features across modalities. We set the benchmark for state-of-the-art performance on the GDPx3 dataset, improving cross-modal alignment and prediction quality over single-modality and diffusion baselines. PertFlow could support downstream applications in virtual drug screening, mechanism discovery, and integrated pharmacological modeling by enabling joint prediction of RNA-seq and image responses to perturbations; a capability, to our knowledge, the first and unique among current methods.

2 RELATED WORKS

Drug-Conditioned and Cross-Modal Modeling: Recent methods predict transcriptional responses to chemical perturbations but remain largely transcriptomics-focused. PRnet Qi et al. (2024) employs a perturbation-conditioned generative model to predict expression changes for novel compounds at bulk and single-cell levels, while TranSiGen Tong et al. (2024) uses self-supervised learning to reconstruct drug-induced profiles from basal expression and compound structure, though limited to denoising and reconstruction. MolGene-E Ohlan et al. (2025) shifts toward inverse design by using contrastive learning to harmonize bulk and single-cell data, but generates molecules from expression rather than predicting responses. Other approaches emphasize repurposing and mechanism discovery without producing new profiles, leaving morphological drug effects unexplored. In parallel, integration of transcriptomic and imaging modalities has emphasized prediction over generation. BLEEP Xie et al. (2023) applies bi-modal contrastive learning to predict spatial gene expression from H&E images, while SCHAF Comiter (2024) is among the few generative models, using GANs to synthesize spatially resolved single-cell omics from histology. TransformerST Zhao et al. (2024) fuses histology with gene expression for super-resolution predictions, prioritizing data enhancement. Multi-modal perturbation frameworks such as Perturb-multi-modal Saunders et al. (2025) and CRISPR ST Binan et al. (2025) integrate imaging and sequencing to study genetic perturbations, but focus on measurement rather than synthesis and mainly on genetic rather than chemical interventions. Fusion-based methods Lu et al. (2024) combine chemical, transcriptomic, and other biological data for prediction and classification, yet cross-modal generative modeling of cellular responses remains unaddressed.

Morphological Profiling and Generative Frameworks: Cellular imaging provides critical insights into drug mechanisms, with Cell Painting Bray et al. (2016) capturing multiplexed phenotypes under perturbations and widely used in virtual screening. Advances in deep learning have enhanced morphological profiling through convolutional models and computer vision Tang et al. (2024), while tools such as CellProfiler McQuin et al. (2018) automate analysis and Cellpose Stringer et al. (2021) improves segmentation. However, generative modeling of cellular images conditioned on perturbations remains limited. Progress in generative frameworks highlights potential for this gap: diffusion models achieve state-of-the-art performance in image, protein, and molecule generation Guo et al. (2024), but suffer from high computational cost and slow sampling. Rectified Flow Liu et al. (2022) offers a more efficient alternative by learning straight-line transport between distributions, reducing sampling steps without loss of quality. This efficiency stems from flow matching Lipman et al. (2022), which linearly interpolates between noise and data, making rectified flow especially suited for large-scale drug screening. While most current methods are single-modal, advances such as Stable Diffusion 3 Esser et al. (2024) demonstrate the feasibility of multi-modal generative modeling, opening opportunities for predictive simulation of cellular responses.

3 METHODS

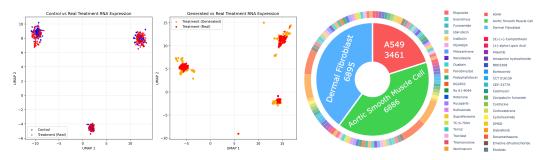


Figure 2: (RIGHT) UMAP representation of control vs real treatment vs generated treatment gene expression data. (LEFT) Distribution of cell lines and compounds in dataset

Problem Formulation. We formalize the drug conditioned multi modal generation problem as learning a mapping from control cellular states to treatment responses across transcriptomic and morphological modalities. Given control gene expression $\mathbf{x}_{\text{rna}}^{\text{ctrl}} \in \mathbf{R}^G$ where G is the number of genes, control cellular images $\mathbf{x}_{\text{img}}^{\text{ctrl}} \in \mathbf{R}^{C \times H \times W}$ with C channels and spatial dimensions $H \times W$, and drug conditioning information $\mathbf{c} = \{c_{\text{compound}}, c_{\text{cell}}, c_{\text{conc}}, c_{\text{time}}\}$ including compound identity, cell line, concentration, and timepoint, our objective is to generate treatment outcomes, where f_{θ} represents our unified generative model parameterized by θ : $\mathbf{x}_{\text{reat}}^{\text{treat}}, \mathbf{x}_{\text{img}}^{\text{treat}} = f_{\theta}(\mathbf{x}_{\text{rna}}^{\text{ctrl}}, \mathbf{c}_{\text{img}}^{\text{ctrl}}, \mathbf{c})$

Dataset Description. Our study leverages the Ginkgo Data Platform (GDP) series Model & Biologics (2025), a multimodal dataset integrating transcriptomic profiles (GDPx1/GDPx2) and four-channel fluorescence microscopy images (GDPx3) from drug-treated cell cultures. We implemented cross-modal pairing by identifying overlapping compounds and experimental conditions, standardizing metadata (concentration units, cell line nomenclature, temporal alignment), and establishing DMSO controls as baseline references. Transcriptomic preprocessing follows established protocols like total count normalization to 10^6 reads per sample, log1p transformation, and highly variable gene selection (n=8000) using scanpy, to focus on the most informative genomic features. Image preprocessing addresses 16-bit microscopy data through proper intensity scaling (16-bit to [-1,1] range), percentile-based contrast enhancement (1^{st} -99 th percentile) applied per channel, and bilinear interpolation to uniform spatial dimensions.

Drug compounds are represented through multi-modal molecular encodings combining structural and physicochemical information. We extract Morgan and RDKit molecular fingerprints (1024 bits each) from canonical SMILES strings, providing binary structural descriptors capturing substructural patterns and pharmacophoric features. Molecular descriptors include eighteen 2D properties (molecular weight, logP, topological polar surface area, hydrogen bond donors/acceptors, rotatable bonds, aromatic rings, and complexity measures) and five 3D properties when available from SDF structures. For compounds lacking preprocessed molecular data, we implement on-demand SMILES processing with RDKit to ensure comprehensive coverage. Missing molecular information is handled through zero-padding with appropriate masking, while molecular descriptors are normalized using dataset-wide statistics to ensure stable training dynamics across diverse chemical spaces. The dataset allows stratified paired control-treatment comparisons, where drug-treated samples are systematically matched with vehicle DMSO controls from identical cell lines and experimental conditions. This preserves the combinatorial structure of cell line-compound-dose-time relationships across training and validation partitions, ensuring robust model generalization across the full experimental parameter space.

Architecture. PertFlow uses a shared representation learning paradigm with 3 components: (1) individual modality encoders that process control RNA-seq and imaging data, (2) cross-modal attention mechanism that aligns features across modalities, (3) generation heads that produce treatment RNA-seq via direct prediction and treatment images via rectified flow dynamics. Figure 3 shows the architecture pipeline of PertFlow. (TOP) shows the entire pipeline of the architecture with input RNA-seq and image going through their respective encoders. Output from the two encoders then pass through the multi-token cross-modal attention, before entering the shared encoder along with conditioning information which passes through the drug encoder. The transcriptome head uses MSE

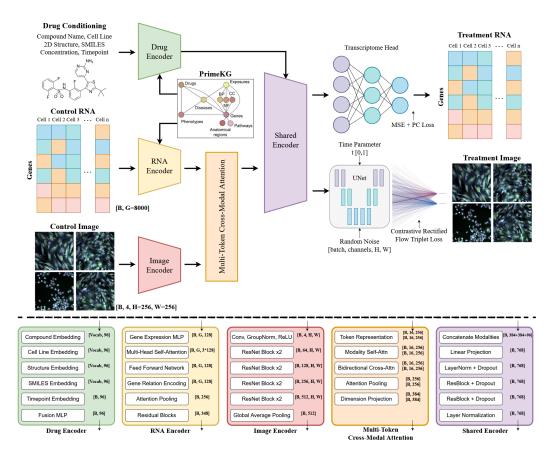


Figure 3: PertFlow architecture for drug conditioned control RNA-image to treatment RNA-image

loss to predict treatment RNA, and the image UNet Huang et al. (2020) with noise and time parameter input uses triplet contrastive loss and rectified flow dynamics to generate treatment cellular image, from the shared embeddings, respectively. (BOTTOM) shows each encoder block and its subcomponents with inputs entering from the top; along with output shapes of each component for each modality.

Our RNA-seq encoder processes gene expression data through multi-layer self-attention Vaswani et al. (2017) to capture gene-gene interactions: $\mathbf{E}_{\text{gene}} = \text{GeneEmbedding}(\mathbf{x}_{\text{rna}}^{\text{ctrl}})$, where each gene expression value is projected to a d_{gene} -dimensional embedding space. We apply L layers of multi-head self-attention (MHA): $\mathbf{A}^{(l)} = \text{MHA}(\mathbf{E}^{(l-1)}, \mathbf{E}^{(l-1)}, \mathbf{E}^{(l-1)})$

$$\mathbf{E}^{(l)} = \text{LayerNorm}(\mathbf{E}^{(l-1)} + \text{FFN}(\text{LayerNorm}(\mathbf{E}^{(l-1)} + \mathbf{A}^{(l)})))$$

The final RNA-seq features are obtained through attention-based pooling:

$$\mathbf{h}_{\text{rna}} = \sum_{i=1}^{G} \alpha_i \mathbf{E}_i^{(L)}, \quad \alpha_i = \frac{\exp(\mathbf{w}^T \tanh(\mathbf{W}_{\text{pool}} \mathbf{E}_i^{(L)}))}{\sum_{j=1}^{G} \exp(\mathbf{w}^T \tanh(\mathbf{W}_{\text{pool}} \mathbf{E}_j^{(L)}))}$$

Control cellular images are processed through a ResNet-style convolutional He et al. (2016) architecture: $\mathbf{h}_{img} = \text{GlobalPool}(\text{ResNet}(\mathbf{x}_{img}^{ctrl}))$. Drug conditioning information combines categorical and continuous variables: $\mathbf{h}_{drug} = \text{Fusion}([\mathbf{e}_{compound}, \mathbf{e}_{cell}, \text{Conc}(c_{conc}), \text{Time}(c_{time})])$, where $\mathbf{e}_{compound}$ and \mathbf{e}_{cell} are learned embeddings for compound and cell line identities.

Knowledge graph integration enhances both molecular and genomic representations through structured biological knowledge from PrimeKG Chandak et al. (2023). For drug embeddings, compounds are mapped to knowledge graph entities capturing molecular interactions, pathways, and pharmacological relationships. The heterogeneous graph neural network processes drug-protein, drug-drug, and protein-protein interactions:

$$\mathbf{h}_{drug}^{kg} = KGDrugEncoder(\mathbf{G}_{drug}, \mathbf{E}_{rel})$$

where G_{drug} represents drug nodes and E_{rel} captures multi-relational edges. Similarly, gene expressions are enhanced with protein interaction networks and pathway information:

 $\mathbf{E}_{RNA}^{kg} = KGGeneEncoder(\mathbf{G}_{gene}, \mathbf{E}_{ppi})$

The knowledge graph embeddings are integrated additively with learned representations:

$$\mathbf{h}_{drug} = \mathbf{h}_{drug} + \alpha_{drug} \mathbf{h}_{drug}^{kg}$$
 and $\mathbf{E}_{RNA} = \mathbf{E}_{RNA} + \alpha_{RNA} \mathbf{E}_{RNA}^{kg}$

where $\alpha_{\text{drug}} = 0.3$ and $\alpha_{\text{RNA}} = 0.3$ are learned weighting factors.

To capture cross-modal RNA-Image dependencies, we use multi-token cross-attention. Each modality is projected to K token representations:

$$\mathbf{T}_{rna} = RNAProj(\mathbf{h}_{rna})$$
 $\mathbf{T}_{img} = ImageProj(\mathbf{h}_{img})$

Each modality goes through a self-attention block, then cross-attention is applied bidirectionally:

$$\mathbf{T}_{ma}^{cross} = MHA(\mathbf{T}_{rna}, \mathbf{T}_{img}, \mathbf{T}_{img}) \qquad \mathbf{T}_{img}^{cross} = MHA(\mathbf{T}_{img}, \mathbf{T}_{rna}, \mathbf{T}_{rna})$$

Enhanced features are obtained through residual connections and attention pooling:

$$\mathbf{h}_{rna}^{enh} = AttentionPool(\mathbf{T}_{rna} + \mathbf{T}_{rna}^{cross}) \qquad \mathbf{h}_{img}^{enh} = AttentionPool(\mathbf{T}_{img} + \mathbf{T}_{img}^{cross})$$

The cross-modal features are combined with drug conditioning to form a unified representation:

$$\mathbf{h}_{shared} = SharedEncoder([\mathbf{h}_{rna}^{enh}, \mathbf{h}_{img}^{enh}, \mathbf{h}_{drug}])$$

This shared representation captures the complex dependencies between molecular states, morphological features, and drug effects necessary for coherent multi-modal generation.

Treatment gene expression is generated through direct prediction from the shared representation:

$$\mathbf{x}_{rna}^{treat} = TranscriptomeHead(\mathbf{h}_{shared})$$

For image generation, we adapt rectified flow dynamics. Given noise $\mathbf{z}_0 \sim \mathcal{N}(0, \mathbf{I})$ and target image $\mathbf{x}_{\text{img}}^{\text{treat}}$, rectified flow defines a linear interpolation path $\mathbf{x}_t = (1-t)\mathbf{z}_0 + t\mathbf{x}_{\text{img}}^{\text{treat}}$, $t \in [0,1]$. The velocity field is defined as $\mathbf{v}_t = \mathbf{x}_{\text{img}}^{\text{treat}} - \mathbf{z}_0$, and our multi-modal-conditioned UNet learns to predict this velocity $\mathbf{v}_\theta(\mathbf{x}_t, t, \mathbf{h}_{\text{shared}}) \approx \mathbf{v}_t$. The UNet incorporates cross-attention layers that attend to image conditioning derived from the shared representation $\mathbf{c}_{\text{img}} = \text{ImageUNet}(\mathbf{h}_{\text{shared}})$.

Our training strategy combines task-specific losses with cross-modal consistency objectives. We use a combination of MSE and auxiliary Pearson Correlation loss for transcriptome prediction:

$$\mathcal{L}_{\text{rna}} = 0.9 \cdot \text{MSE}(\mathbf{x}_{\text{rna}}^{\text{treat}}, \hat{\mathbf{x}}_{\text{rna}}^{\text{treat}}) + 0.1 \cdot \text{PC}(\mathbf{x}_{\text{rna}}^{\text{treat}}, \hat{\mathbf{x}}_{\text{rna}}^{\text{treat}})$$

For rectified flow training, we minimize the velocity prediction error for image generation:

$$\mathcal{L}_{\text{img}} = \mathbf{E}_{t,\mathbf{z}_0,\mathbf{x}_{\text{img}}^{\text{treat}}} \left[\|\mathbf{v}_{\theta}(\mathbf{x}_t,t,\mathbf{h}_{\text{shared}}) - (\mathbf{x}_{\text{img}}^{\text{treat}} - \mathbf{z}_0) \|^2 \right]$$

We implement triplet contrastive consistency Stoica et al. (2025) to ensure well-aligned features produce better predictions than misaligned ones:

$$\mathcal{L}_{triplet} = \mathbf{E}[\max(0, \text{margin} - (\mathcal{L}_{neg} - \mathcal{L}_{pos}))]$$

where \mathcal{L}_{pos} is the prediction error with aligned features and \mathcal{L}_{neg} with misaligned features. The complete training objective combines all losses:

$$\mathcal{L}_{total} = w_{rna}\mathcal{L}_{rna} + w_{img}\mathcal{L}_{img} + w_{triplet}\mathcal{L}_{triplet}$$

where the weights are set to $w_{\rm rna} = 0.5$, $w_{\rm img} = 0.5$, $w_{\rm triplet} = 0.05$.

Treatment RNA-seq is generated through a single forward pass $\mathbf{x}_{\text{rna}}^{\text{treat}} = f_{\theta}(\mathbf{x}_{\text{rna}}^{\text{ctrl}}, \mathbf{x}_{\text{img}}^{\text{ctrl}}, \mathbf{c})$. For high-quality image generation, we use an adaptive DOPRI5 solver Dormand & Prince (1986) that iteratively integrates the learned velocity field:

$$\frac{d\mathbf{x}}{dt} = \mathbf{v}_{\theta}(\mathbf{x}_t, t, \mathbf{h}_{\text{shared}})$$

Starting from noise $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I})$ at t=0, the solver adaptively adjusts step sizes based on error estimation to reach the treatment image at t=1. The adaptive integration ensures both computational efficiency and generation quality. The DOPRI5 method uses a 5th-order Runge-Kutta scheme with embedded 4th-order error estimation for automatic step size control. The step size h is adapted based on the estimated local truncation error to maintain tolerance levels.

Training Parameters. We use 4 attention heads with embedding dimension 128, applying 1 layer of self-attention followed by attention-based pooling. Multi-token representations use K=16 tokens with hidden dimension 256 and 8 attention heads. The rectified flow UNet uses 192 base channels with channel multipliers (1,2,2,2), attention at 16×16 resolution, and cross-attention conditioning at layers 2, 3, 4, and 5. Models are trained with AdamW optimizer ($\beta_1=0.9,\,\beta_2=0.95$), learning rate 10^{-4} with cosine annealing, and automatic mixed precision. Cross-modal consistency weights are gradually increased during training to ensure stable convergence. RNA-seq generation requires a single forward pass, while image generation uses 7-10 DOPRI5 steps with relative tolerance 10^{-3} and absolute tolerance 10^{-4} for high-quality synthesis. The models were trained with an effective batch size of 32 taking 5 hours on 8 H100 NVIDIA GPUs with 80GB VRAM.

4 EXPERIMENTS

We emphasize, since our method is the first to introduce the problem of multi-modal RNA-Image generation through drug conditioning, we have no previous method to compare to as baseline. To create baselines we trained three diffusion models, along with ablations of knowledge graph module, triplet contrastive objective, and Pearson correlation loss, with RNA only and Image only models. We set the state-of-the-art performance for this problem on the GDPx3 dataset.

Table 1: PertFlow & PertRNA metrics (mean \pm std)

| Model | $\mathbf{MSE}\downarrow$ | $\mathbf{RMSE}\downarrow$ | $\mathbf{MAE}\downarrow$ | Pearson $r \uparrow$ | Spearman $r\uparrow$ |
|---|--|--|--|--|--|
| PertRNA(-PC) PertRNA(-KG) PertRNA | $\begin{array}{c} 0.412 \pm 0.525 \\ 0.360 \pm 1.274 \\ 0.311 \pm 0.956 \end{array}$ | $\begin{array}{c} 0.598 \pm 0.912 \\ 0.475 \pm 0.366 \\ 0.472 \pm 0.271 \end{array}$ | $\begin{array}{c} 0.286 \pm 0.374 \\ 0.112 \pm 0.106 \\ 0.111 \pm 0.025 \end{array}$ | $\begin{array}{c} 0.511 \pm 0.144 \\ 0.770 \pm 0.098 \\ 0.779 \pm 0.081 \end{array}$ | $\begin{array}{c} 0.502 \pm 0.097 \\ 0.791 \pm 0.063 \\ \textbf{0.795} \pm \textbf{0.026} \end{array}$ |
| PertFlow(-KG) PertFlow | | $\begin{array}{c} 0.470 \pm 0.202 \\ \textbf{0.462} \pm \textbf{0.107} \end{array}$ | | | $\begin{array}{c} 0.735 \pm 0.066 \\ 0.792 \pm 0.041 \end{array}$ |

Table 2: PertFlow & PertImage metrics (mean \pm std)

| Model | SSIM ↑ | PSNR ↑ | $\textbf{LPIPS} \rightarrow$ | FID ↓ |
|---|--|--|--|----------------------------|
| $\begin{array}{c} \mathbf{PertDiff}_N \\ \mathbf{PertDiff}_{x0} \\ \mathbf{PertDiff}_V \end{array}$ | $\begin{array}{c} 0.010 \pm 0.003 \\ 0.192 \pm 0.075 \\ 0.194 \pm 0.071 \end{array}$ | 06.62 ± 0.73 10.71 ± 0.54 11.33 ± 0.42 | $\begin{array}{c} 1.087 \pm 0.096 \\ 0.505 \pm 0.045 \\ 0.499 \pm 0.040 \end{array}$ | 246.01 73.63 55.92 |
| PertImage(-triplet) PertImage(-KG) PertImage | $\begin{array}{c} 0.120 \pm 0.102 \\ 0.182 \pm 0.099 \\ 0.187 \pm 0.095 \end{array}$ | 08.22 ± 0.13 11.05 ± 0.58 11.46 ± 0.61 | 0.308 ± 0.178 0.498 ± 0.035 0.509 ± 0.043 | $106.72 \\ 50.38 \\ 46.88$ |
| PertFlow(-KG) PertFlow | 0.206 ± 0.094 0.205 ± 0.097 | | 0.505 ± 0.043 0.511 ± 0.038 | 31.59 24.06 |



Figure 4: Similarity rating

Drug effects on gene expression and cell morphology: We trained PertFlow (control RNA-seq and image to treatment RNA-seq and image), PertRNA (control RNA-seq to treatment RNA-seq), PertImage (control image to treatment image), and their respective ablations omitting the knowledge-graph and contrastive rectified flow objective. PertFlow demonstrates strong performance in predicting treatment gene expression from control conditions in Table 1. Achieving Pearson correlation (0.780 ± 0.264) and Spearman correlation (0.792 ± 0.041) across drug perturbations, measuring with all genes. MSE (0.231 ± 0.708) and MAE (0.110 ± 0.166) indicate robust prediction accuracy for transcriptomic responses. While the PertRNA baseline achieved correlation metrics (Pearson r (0.779 ± 0.081)), Spearman r (0.795 ± 0.026)), PertFlow's joint modeling approach maintains competitive performance while simultaneously generating cellular morphological responses compared to baselines and ablations.

PertFlow also outperformed the PertImage baseline in the image generation task achieving higher SSIM (0.205 \pm 0.097), PSNR (11.66 \pm 0.68), LPIPS (0.511 \pm 0.038), and lower FID (24.06). This demonstrates that incorporating cross-modal information and drug conditioning does not compromise transcriptomic prediction and cellular image generation quality, while enabling the unique capability of joint multi-modal generation. The strong correlation values and lower FID indicates

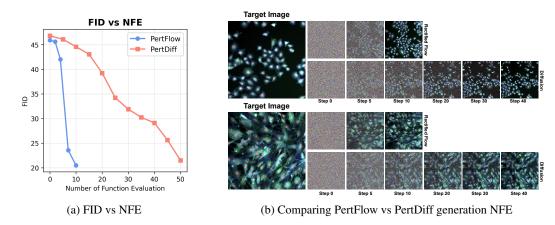


Figure 5: Step-wise comparison between Rectified Flow and Diffusion models duuring inference

that PertFlow successfully captures the complex cellular relationships between drug perturbations and gene expression changes.

We compare our PertFlow model with three PertDiff diffusion variants in Table 2. Diffusion serves as the baseline against our contrastive rectified flow method. The standard DDPM formulation trains the model to predict the injected noise $\epsilon \sim N(0,I)$ from the noisy input x_t . This requires disentangling signal from noise across noise levels, causing instabilities when signal-to-noise ratios are low. The weak performance of PertDiff $_N$ (SSIM: 0.010, FID: 246.01) illustrates these challenges, as the model fails to generate coherent cellular structures from pure noise predictions. The direct x_0 parameterization predicts the clean target x_0 from noisy x_t . While more stable, it forces the model to implicitly learn the full denoising trajectory. PertDiff $_{x_0}$ achieves better but still limited results (SSIM: 0.192, FID: 73.63), reflecting the lack of strong theoretical grounding. The velocity (v) parameterization improves training stability by predicting $v = \alpha_t \epsilon - \sigma_t x_0$, which balances objectives across timesteps, reduces variance, and improves gradient flow. PertDiff $_V$ shows marked improvement (SSIM: 0.194, FID: 55.92), validating this formulation for biological image generation. Compared to baselines and ablations we observe that PertFlow successfully learns more meaningful representations that generalize across different compounds, concentrations, cell lines, and timepoints in the dataset.

Figure 4 shows similarity rating by expert pathologist evaluation (10-point scale of morphology, detail, and plausibility), which confirmed PertFlow's quality against real treatment samples: median scores were 7.0 and 8.0. Step-wise comparison (Figure 5) highlights the difference in inference dynamics. Our rectified flow generates recognizable structures by NFE 10 showing nuclear boundaries and cytoplasmic organization with near-final morphology. Diffusion requires more steps to reach similar organization. Both methods preserve multi-channel fluorescence distributions, but rectified flow avoids the oscillatory intermediates observed in diffusion, maintaining smoother trajectories.

Recovering drug-induced phenotype and morphology: Figure 1 shows the comparison of generated treatment vs real treatment images with drug name and concentration. The yellow boxes indicate similar cellular features due to drug perturbations in real and generated images. From left to right the drugs have the following effect on the cellular morphology: (1) Cevipabulin is a microtubule-destabilizing agent that binds to tubulin Yang et al. (2021), disrupting microtubule dynamics, which leads to mitotic arrest and apoptosis in cancer cells. It shows anti-proliferative effects by inhibiting microtubule polymerization. (2) S-Camptothecin and its stereoisomers inhibit DNA topoisomerase Hansch & Verma (2007), causing DNA damage during replication. This leads to DNA double-strand breaks, S-phase cell cycle arrest, and apoptosis, especially in rapidly dividing cells. (3) Podophyllotoxin binds to tubulin and inhibits microtubule assembly Desbene & Giorgi-Renault (2002), resulting in mitotic arrest at metaphase and subsequent apoptosis. It also serves as a precursor for etoposide, a topoisomerase II inhibitor. (4) Dabrafenib selectively inhibits mutant BRAF kinase (commonly V600E mutation) Planchard et al. (2022), blocking MAPK/ERK signaling pathway, leading to decreased tumor cell proliferation and inducing apoptosis in BRAF-mutated cancer cells. Figure 6 shows more examples of generated treatment images. Note that the 5th exam-

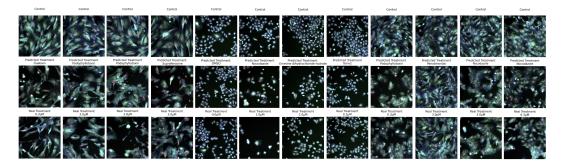


Figure 6: Examples of generated treatment vs real treatment images with drug name and concentr.

ple from the left with DMSO control drug has no effect on the cellular morphology and our model correctly predicts that. Further examples are available in the appendix.

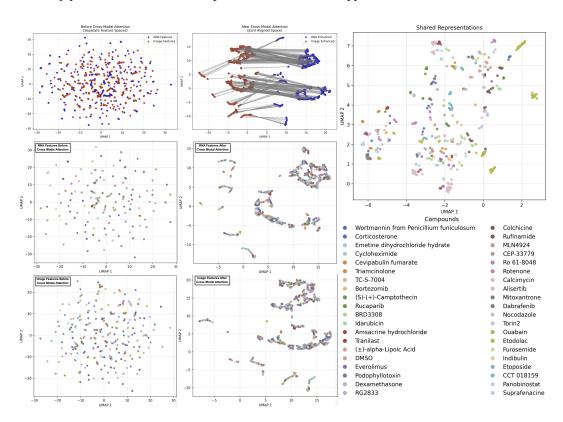


Figure 7: UMAP of RNA and image embeddings before and after cross modal attention block

Figure 2 UMAPs McInnes et al. (2018) demonstrate PertFlow's ability to generate biologically coherent treatment responses across both modalities. The left panel shows clear separation between control and treatment clusters, indicating that drug perturbations induce distinct transcriptomic signatures well-captured in the embedding space (overlap results from control data sampled with treatment as negative control). The right panel shows generated treatment samples clustering closely with real treatment samples, particularly in the upper region corresponding to strong drug responses, indicating PertFlow successfully learns directional gene expression changes induced by drug perturbations. The tight clustering of generated samples with real counterparts for high-response conditions, demonstrates that cross-modal attention effectively leverages imaging information to improve RNA-seq prediction accuracy, producing treatment profiles that are biologically realistic and consistent with experimental observations.

Figure 7 UMAPs illustrates the effect of cross-modal attention on feature alignment and organization. The top panels show the transformation from separate feature spaces (left) to a joint aligned space (middle), where RNA-seq features (blue) and image features (red) demonstrating successful cross-modal correspondence learning. The middle and bottom panels reveal that cross-modal attention dramatically improves feature organization within each modality. Before attention, both RNA-seq features and image features exhibit scattered, unstructured distributions in their respective embedding spaces. After cross-modal attention, both modalities show highly structured, clustered organization, indicating that the attention mechanism aligns features across modalities and also enhances the internal structure and discriminative power of each individual feature space, resulting in more meaningful and interpretable representations for both RNA-seq and imaging data.

UMAP of shared representation space (right) demonstrates the successful integration of RNA, imaging, and drug modalities. Organization based on compounds suggest that the cross-modal attention mechanism effectively combines information from both gene expression and imaging data with drug conditioning into a unified latent space where samples with similar biological states cluster together regardless of their original modality. This enables the model to leverage multi modal information for improved downstream tasks such as treatment response prediction and generation.

5 DISCUSSION

PertFlow represents a foundational step toward unified modeling of multi-modal cellular drug responses, bridging the molecular (transcriptomic) and phenotypic (morphological) effects of chemical perturbations. Unlike previous approaches that treat these modalities in isolation or only predict one from the other, PertFlow achieves simultaneous, drug-conditioned generation of both gene expression profiles and cellular morphology. The integration of control transcriptomic and imaging data into a shared embedding space, combined with rectified flow dynamics, enables biologically consistent synthesis of treatment outcomes. Despite this progress, several challenges remain. First, generalization to unseen cell lines or novel compounds is limited by the scarcity of paired multi-modal datasets with shared metadata. While PertFlow can still infer morphological changes from control data alone, future work should explore integrating chemical structure representations or compound-target interaction graphs to enhance out-of-distribution performance. Second, while PertFlow enables modality translation, aligning embeddings across modalities may inadvertently entangle task-relevant factors. Disentangling causal latent factors remains an open question for cross-modal modeling. Third, the in vitro context of our experiments may not capture drug effects requiring complex microenvironmental interactions, such as immune modulation. Extending Pert-Flow to model cell-cell communication or tissue-level organization could enhance translational utility. Overall, PertFlow sets the stage for future cross-modal generative modeling in drug discovery, offering a unified framework for understanding how molecular mechanisms manifest as observable phenotypes under pharmacological perturbation.

6 Conclusion

We introduced PertFlow, the first unified generative framework for jointly modeling transcriptomic and morphological drug responses using cross-modal attention and rectified flow dynamics. By aligning control RNA-seq and image features through a shared embedding space conditioned on drug metadata, PertFlow enables simultaneous prediction of treatment gene expression and synthesis of cellular morphology. Extensive evaluation on the GDPx3 dataset demonstrates strong cross-modal consistency, biologically realistic image generation, and competitive transcriptomic prediction performance, outperforming single-modality and diffusion baselines. UMAP analysis of cross-modal attention and shared embeddings further strengthen our hypothesis. Our results highlight PertFlow's potential for virtual drug screening, mechanistic hypothesis generation, and multimodal perturbation analysis, paving the way for more integrative and interpretable approaches to pharmacological modeling.

REFERENCES

Loïc Binan, Aiping Jiang, Serwah A Danquah, Vera Valakh, Brooke Simonton, Jon Bezney, Robert T Manguso, Kathleen B Yates, Ralda Nehme, Brian Cleary, et al. Simultaneous crispr

- screening and spatial transcriptomics reveal intracellular, intercellular, and functional transcriptional circuits. *Cell*, 188(8):2141–2158, 2025.
 - 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.
 - Payal Chandak, Kexin Huang, and Marinka Zitnik. Building a knowledge graph to enable precision medicine. *Scientific Data*, 10(1):67, 2023. URL https://doi.org/10.1038/s41597-023-01960-3.
 - Charles Comiter. *Inference of single cell profiles from histology stains with the Single-Cell omics from Histology Analysis Framework (SCHAF)*. Massachusetts Institute of Technology, 2024.
 - Stephanie Desbene and Sylviane Giorgi-Renault. Drugs that inhibit tubulin polymerization: the particular case of podophyllotoxin and analogues. *Current Medicinal Chemistry-Anti-Cancer Agents*, 2(1):71–90, 2002.
 - John R Dormand and Peter J Prince. Runge-kutta triples. *Computers & Mathematics with Applications*, 12(9):1007–1017, 1986.
 - Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
 - Zhiye Guo, Jian Liu, Yanli Wang, Mengrui Chen, Duolin Wang, Dong Xu, and Jianlin Cheng. Diffusion models in bioinformatics and computational biology. *Nature reviews bioengineering*, 2 (2):136–154, 2024.
 - Corwin Hansch and Rajeshwar P Verma. 20-(s)-camptothecin analogues as dna topoisomerase i inhibitors: a qsar study. *ChemMedChem: Chemistry Enabling Drug Discovery*, 2(12):1807–1813, 2007.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
 - Huimin Huang, Lanfen Lin, Ruofeng Tong, Hongjie Hu, Qiaowei Zhang, Yutaro Iwamoto, Xianhua Han, Yen-Wei Chen, and Jian Wu. Unet 3+: A full-scale connected unet for medical image segmentation. In *ICASSP 2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 1055–1059. Ieee, 2020.
 - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
 - Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
 - Xiaohua Lu, Liangxu Xie, Lei Xu, Rongzhi Mao, Xiaojun Xu, and Shan Chang. Multimodal fused deep learning for drug property prediction: Integrating chemical language and molecular graph. *Computational and Structural Biotechnology Journal*, 23:1666–1679, 2024.
 - Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv* preprint arXiv:1802.03426, 2018.
 - Claire McQuin, Allen Goodman, Vasiliy Chernyshev, Lee Kamentsky, Beth A Cimini, Kyle W Karhohs, Minh Doan, Liya Ding, Susanne M Rafelski, Derek Thirstrup, et al. Cellprofiler 3.0: Next-generation image processing for biology. *PLoS biology*, 16(7):e2005970, 2018.
 - API Model and Enzymes Biologics. Ginkgo datapoints: Data generation for ai model training. 2025.

- Rahul Ohlan, Raswanth Murugan, Li Xie, Mohammedsadeq Mottaqi, Shuo Zhang, and Lei Xie. Molgene-e: Inverse molecular design to modulate single cell transcriptomics. *bioRxiv*, 2025.
 - David Planchard, Benjamin Besse, Harry JM Groen, Sayed MS Hashemi, Julien Mazieres, Tae Min Kim, Elisabeth Quoix, Pierre-Jean Souquet, Fabrice Barlesi, Christina Baik, et al. Phase 2 study of dabrafenib plus trametinib in patients with braf v600e-mutant metastatic nsclc: updated 5-year survival rates and genomic analysis. *Journal of Thoracic Oncology*, 17(1):103–115, 2022.
 - Xiaoning Qi, Lianhe Zhao, Chenyu Tian, Yueyue Li, Zhen-Lin Chen, Peipei Huo, Runsheng Chen, Xiaodong Liu, Baoping Wan, Shengyong Yang, et al. Predicting transcriptional responses to novel chemical perturbations using deep generative model for drug discovery. *Nature Communications*, 15(1):9256, 2024.
 - Reuben A Saunders, William E Allen, Xingjie Pan, Jaspreet Sandhu, Jiaqi Lu, Thomas K Lau, Karina Smolyar, Zuri A Sullivan, Catherine Dulac, Jonathan S Weissman, et al. Perturb-multimodal: A platform for pooled genetic screens with imaging and sequencing in intact mammalian tissue. *Cell*, 2025.
 - George Stoica, Vivek Ramanujan, Xiang Fan, Ali Farhadi, Ranjay Krishna, and Judy Hoffman. Contrastive flow matching. *arXiv preprint arXiv:2506.05350*, 2025.
 - Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a generalist algorithm for cellular segmentation. *Nature methods*, 18(1):100–106, 2021.
 - Qiaosi Tang, Ranjala Ratnayake, Gustavo Seabra, Zhe Jiang, Ruogu Fang, Lina Cui, Yousong Ding, Tamer Kahveci, Jiang Bian, Chenglong Li, et al. Morphological profiling for drug discovery in the era of deep learning. *Briefings in Bioinformatics*, 25(4), 2024.
 - Xiaochu Tong, Ning Qu, Xiangtai Kong, Shengkun Ni, Jingyi Zhou, Kun Wang, Lehan Zhang, Yiming Wen, Jiangshan Shi, Sulin Zhang, et al. Deep representation learning of chemical-induced transcriptional profile for phenotype-based drug discovery. *Nature Communications*, 15(1):5378, 2024.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in neural information processing systems, pp. 5998–6008, 2017.
 - Ronald Xie, Kuan Pang, Sai Chung, Catia Perciani, Sonya MacParland, Bo Wang, and Gary Bader. Spatially resolved gene expression prediction from histology images via bi-modal contrastive learning. *Advances in Neural Information Processing Systems*, 36:70626–70637, 2023.
 - J Yang, Y Yu, Y Li, W Yan, H Ye, L Niu, M Tang, Z Wang, Z Yang, H Pei, et al. Cevipabulin-tubulin complex reveals a novel agent binding site on α -tubulin with tubulin degradation effect. sci adv 7: eabg4168, 2021.
 - Chongyue Zhao, Zhongli Xu, Xinjun Wang, Shiyue Tao, William A MacDonald, Kun He, Amanda C Poholek, Kong Chen, Heng Huang, and Wei Chen. Innovative super-resolution in spatial transcriptomics: a transformer model exploiting histology images and spatial gene expression. *Briefings in Bioinformatics*, 25(2):bbae052, 2024.

A APPENDIX

A.1 ETHICS STATEMENT

We used large language models solely for manuscript proofreading and grammar checking, with no involvement in code or content generation.

A.2 REPRODUCIBILITY STATEMENT

We will release the code and pretrained model weights upon acceptance.

A.3 IDENTIFYING GENES CHANGING MORPHOLOGIC PHENOTYPES:

We evaluated gene contributions to morphological recovery using gradient-based feature importance with respect to the flow-matching loss during inference. Through cross-modal embedding co-registration, the model correctly identified gene modules linked to treatment-induced morphology changes for example, apoptosis pathway genes activated in A549 cells under camptothecin or etoposide, and reduced activation of cell-cycle modules under proliferation inhibitors compared to negative controls. To map pathways affected by drug treatments, we extracted gene importance scores from the model's latent representations for each sample (Figure 8). Scores capture the model's learned associations between gene expression and drug-induced transcriptional changes. For each drug–cell line pair, we averaged scores across samples, selected the top 200 genes, and performed gene set enrichment analysis with GSEApy using the MSigDB Hallmark 2020 collection (adjusted p $_{\rm i}$ 0.25). This pipeline systematically linked drug-specific transcriptional signatures to biological processes, revealing both universal stress responses (e.g., EMT, TNF-alpha/NF-kB signaling) and cell line–specific activations. Enrichment results were visualized with scanpy-style dotplots, where dot size reflects gene overlap and color intensity indicates significance, enabling clear comparison of pathway activation across compounds and cell types.

Gene set enrichment analysis revealed distinct cell line-specific responses to pharmaceutical compounds, with notable differences in pathway activation between A549 lung cancer cells, human aortic smooth muscle cells (HASMC), and human dermal fibroblasts. A549 cells consistently showed limited pathway enrichment, with most drugs activating only TNF-alpha signaling via NF-kB and occasionally inflammatory response pathways. In contrast, both HASMC and dermal fibroblasts demonstrated robust, multi-pathway responses to the same compounds, suggesting that A549 cells may have inherent resistance mechanisms or altered sensitivity to drug-induced transcriptional changes. Epithelial Mesenchymal Transition (EMT) emerged as the most consistently enriched pathway across drug-cell line combinations, appearing as the top-ranked pathway in nearly all HASMC and fibroblast treatments. This universal EMT activation suggests that pharmaceutical stress triggers fundamental cellular reprogramming programs associated with cell plasticity and survival. TNF- α signaling via NF- κ B represented the second most common response, activated across all three cell lines, indicating that drug treatment consistently triggers inflammatory stress response cascades regardless of the specific compound mechanism of action.

Beyond the universal stress signatures, each cell type exhibited specialized pathway responses reflecting their distinct biological functions. HASMC consistently activated vascular-specific pathways including angiogenesis, coagulation, and hypoxia response, which aligns with their role in vascular homeostasis and their sensitivity to oxygen and hemodynamic stress. Human dermal fibroblasts uniquely enriched for tissue remodeling pathways including myogenesis, adipogenesis, and UV response, consistent with their role in tissue repair and their exposure to environmental stressors. Notably, fibroblasts also showed strong enrichment for interferon gamma response pathways, suggesting heightened immune surveillance capabilities compared to the other cell types. Our multimodal model's gene importance scoring successfully captured biologically relevant drugtarget relationships, as evidenced by cell-type-specific pathway enrichment patterns. The A549 lung cancer cell line's dominant TNF- α /NF- κ B activation across diverse compounds aligns with established literature showing that KRAS-mutant lung adenocarcinomas exhibit heightened inflammatory signaling dependency. This universal inflammatory response contrasts with the diverse EMT and angiogenesis signatures observed in primary cells, recapitulating known differences between transformed cancer cells and stromal cells. Mechanistic validation was further demonstrated through drug-specific responses: DNA-damaging agents activated p53 pathways in TP53-intact A549 cells,

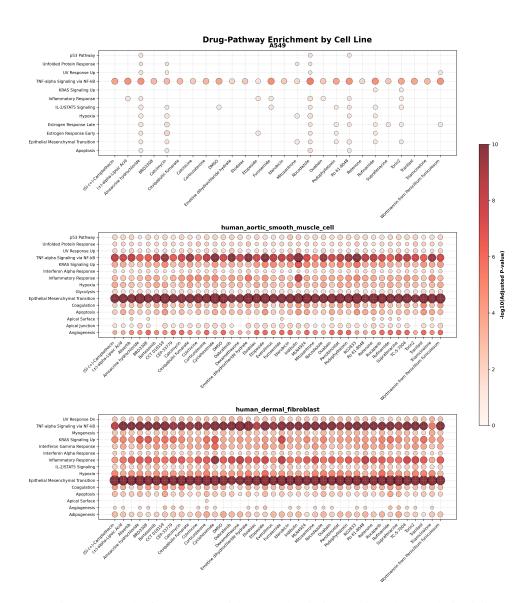


Figure 8: Dotplot for geneset enrichment anlaysis from activated genes during inference

while dabrafenib triggered compensatory KRAS signaling in non-mutant cells - mirroring clinical resistance mechanisms.

The model's biological accuracy extends to morphological predictions, as evidenced by paired Cell Painting data showing dramatic cell number decreases in samples treated with apoptosis-inducing drugs like nocodazole, consistent with the observed enrichment of apoptosis pathways in our transcriptional analysis. These concordant transcriptional and morphological responses demonstrate that our multimodal architecture captures functionally relevant biological relationships rather than mere correlative patterns, validating the gene importance scores as mechanistically informative features for drug response prediction. The enrichment patterns largely validated known drug mechanisms of action. DNA-damaging agents such as etoposide, mitoxantrone, and camptothecin consistently activated p53 pathway and apoptosis responses in responsive cell lines. Anti-inflammatory compounds including dexamethasone and corticosterone showed expected modulation of inflammatory response and IL-2/STAT5 signaling pathways. Targeted inhibitors like dabrafenib demonstrated compensatory KRAS signaling activation, consistent with known resistance mechanisms in cancer cells. However, the limited response in A549 cells to many compounds suggests potential resistance mechanisms that may be clinically relevant for lung cancer treatment strategies.

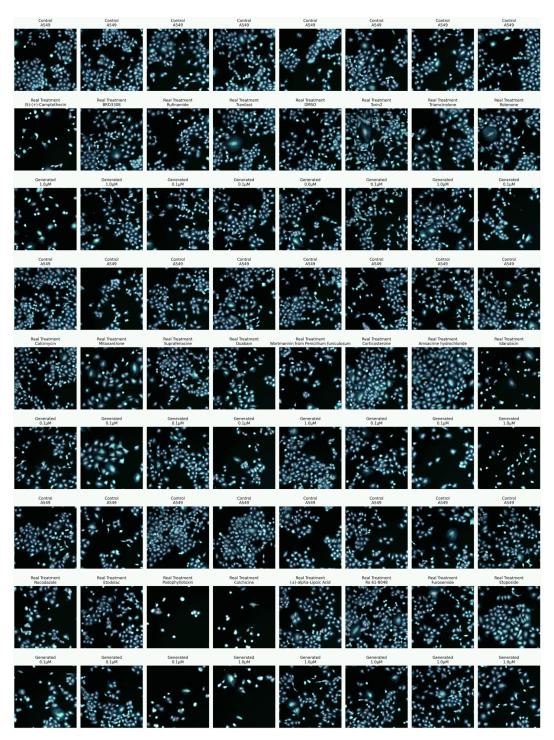


Figure 9: Cell line A549

These findings demonstrate that drug response profiling through transcriptional analysis can reveal both universal cellular stress responses and cell type-specific vulnerabilities, providing insights into both drug mechanism of action and potential therapeutic resistance patterns across different tissue contexts.

A.4 FURTHER EXAMPLES OF GENERATED IMAGES

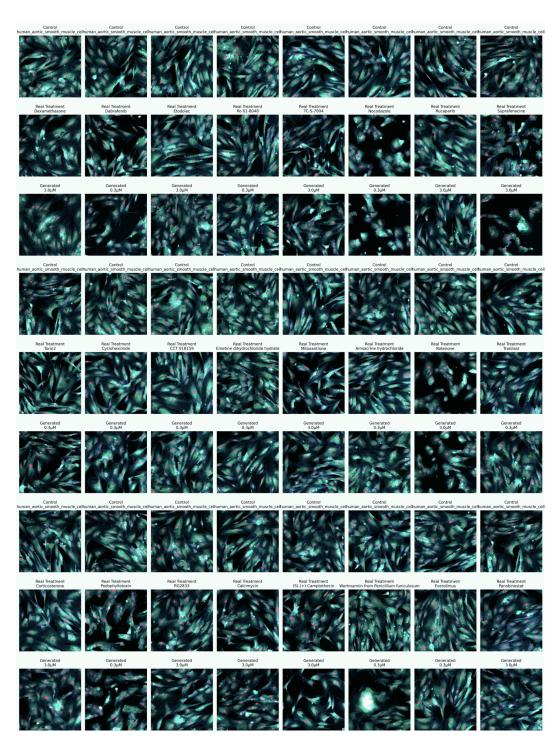


Figure 10: Cell line Aortic Smooth Muscle Cell

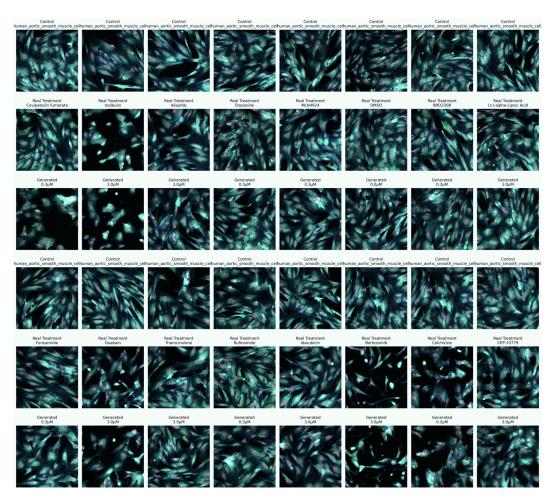


Figure 11: Cell line Aortic Smooth Muscle Cell

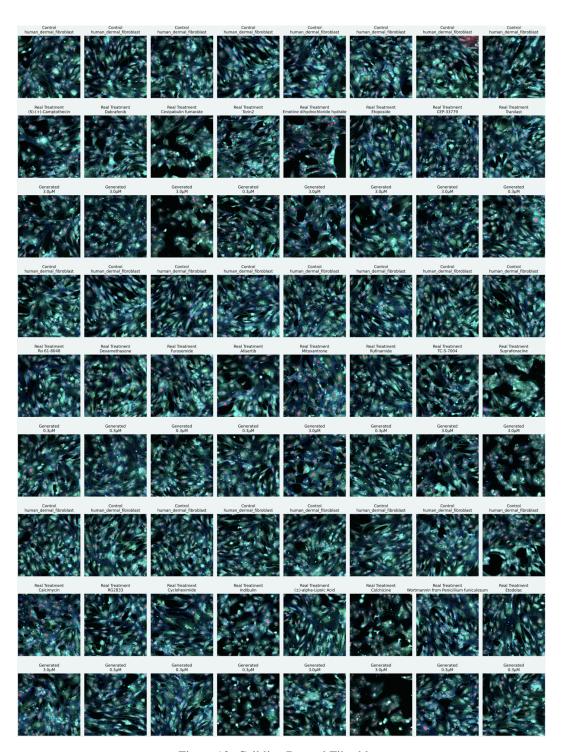


Figure 12: Cell line Dermal Fibroblast

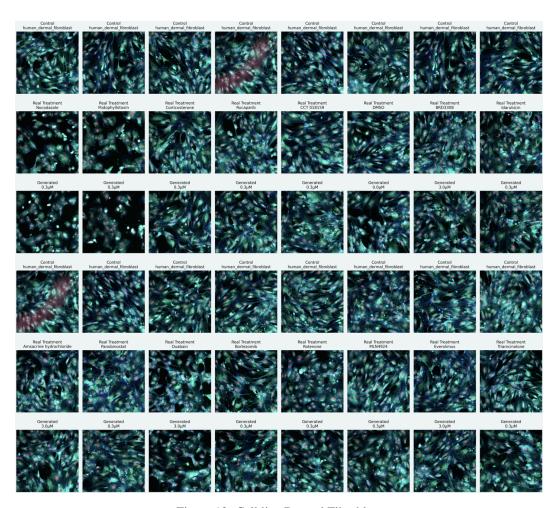


Figure 13: Cell line Dermal Fibroblast