

# Supplemental Material:

## TRIDENT: A Trimodal Cascade Generative Framework for Drug and RNA-Conditioned Cellular Morphology Synthesis

Rui Peng<sup>1,2#</sup> Ziru Liu<sup>4#</sup> Lingyuan Ye<sup>5</sup> Yuxing Lu<sup>1</sup> Boxin Shi<sup>3\*</sup> Jinzhao Wang<sup>1\*</sup>

<sup>1</sup> Department of Big Data and Biomedical AI, College of Future Technology, Peking University

<sup>2</sup> Center for BioMed-X Research, Academy for Advanced Interdisciplinary Studies, Peking University

<sup>3</sup> School of Computer Science, Peking University <sup>4</sup> Yuanpei College, Peking University <sup>5</sup> Tsinghua University

### A. Comparison to MorphDiff

The results of in-distribution comparison against MorphDiff are shown right. Beyond performance, TRIDENT encodes drug information and, in contrast to MorphDiff, does not require post-perturbation RNA during inference. Thus, for a novel drug, cell morphology can be predicted solely from its SMILES string, enabling true virtual screening.

### B. Implementation Details

All models are implemented in PyTorch and trained on a high-performance computing cluster equipped with eight NVIDIA A100 GPUs, each with 80GB of memory. The core of our Morphology Generation Module is a Diffusion Transformer architecture. This transformer is configured with 28 layers, a hidden dimension of 1152, and 16 attention heads. The complete TRIDENT framework is trained end-to-end for a total of 100,000 steps. We employ the AdamW optimizer with a constant learning rate of 1e-4. A global batch size of 32 is used, distributed across the eight GPUs. All Cell Painting images are processed at a resolution of  $512 \times 512$  pixels. The total training process for the final model takes approximately four days.

### C. Dataset pairing consistency

In constructing MorphoGene, we implemented a strict alignment on cell line (MCF-7), perturbation time (24h), and compound dosage: 89% of pairing data have exact dose matches, and the remaining 11% use nearest-neighbor dose matching to minimize dose-response discrepancies.

### D. CellProfiler Feature Construction

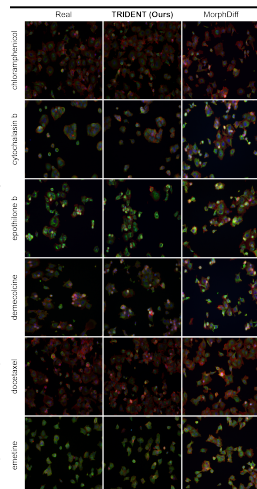
To derive quantitative descriptors of cellular phenotype, we construct a bespoke analysis workflow using CellProfiler (v5.0). This pipeline is engineered to process each Cell Painting image and output a single, comprehensive feature vector summarizing its morphological characteristics.

The workflow’s core is a three-step segmentation process to delineate cellular structures. First, nucleus are identified as primary objects from the blue (DNA) channel. Next, cell boundaries are segmented as secondary objects by propagating outwards from the identified nucleus, using the green channel to define the cell periphery. Finally, the cytoplasm is defined as a tertiary region, calculated by subtracting the nuclear mask from the corresponding cell mask. A quality control step is integrated to discard all objects touching the image border, ensuring that all downstream measurements are derived from complete, intact cells.

Following segmentation, a comprehensive suite of CellProfiler modules is executed to extract measurements from all three compartments (nucleus, cytoplasm, and cells) across all channels. The extracted features include morphological descriptors of size and shape, such as area, perimeter, major and minor axis lengths, eccentricity, and solidity. Furthermore, statistics on pixel intensity distribution like mean, median, standard deviation, median absolute deviation, and quartiles are computed. Finally, the pipeline captures relational metrics, such as spatial relationships between cells, inter-channel signal correlations, and radial intensity distributions.

To generate a single profile for each image, these per-object measurements are aggregated by calculating their mean, median, and standard deviation across all valid cells. This process yields a final, high-dimensional profile of 6,345 distinct morphological features for each image, providing a detailed quantitative fingerprint of the cellular phenotype in response to perturbation.

Methods	FID ↓	KID ↓
TRIDENT (Ours)	49.770	0.013
MorphDiff	78.147	0.094



# Equal contribution. \* Corresponding authors.

## E. TRIDENT Algorithm

---

### Algorithm 1 TRIDENT Framework: Training Procedure

---

**Require:** Training data  $\mathcal{D} = \{(G_{pre}^{(i)}, D^{(i)}, I^{(i)}, G_{post}^{(i)})\}$   
**Require:** Pre-trained image VAE ( $\mathcal{E}_{image}, \mathcal{D}_{image}$ )  
**Require:** Diffusion timesteps  $T$ , variance schedule  $\{\beta_t\}_{t=1}^T$   
 (and  $\alpha_t, \bar{\alpha}_t$ )  
**Ensure:** Trained parameters  $\Theta = \{\phi, \psi, \theta, \gamma\}$

- 1: Initialize VAE parameters  $\phi$  (for  $\mathcal{E}_{rna}, \mathcal{E}_{drug}$ ),  $\psi$  (for  $\mathcal{D}_{perturb}$ )
- 2: Initialize Denoising Transformer  $f_\theta$  parameters  $\theta$
- 3: Initialize Time Embedding parameters  $\gamma$  for  $\mathcal{E}_{time}$
- 4: **for** each epoch  $e = 1, \dots, E_{max}$  **do**
- 5:     **for** each batch  $(G_{pre}, D, I, G_{post}) \sim \mathcal{D}$  **do**
- 6:          $\mathbf{X}_{rna} \leftarrow \mathcal{E}_{rna}(G_{pre}), \mathbf{X}_{drug} \leftarrow \mathcal{E}_{drug}(D)$
- 7:          $[\boldsymbol{\mu}_z, \log \boldsymbol{\sigma}_z^2] \leftarrow \mathcal{E}_{perturb}([\mathbf{X}_{rna}, \mathbf{X}_{drug}])$
- 8:          $\boldsymbol{\epsilon}_z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 9:          $\mathbf{z} \leftarrow \boldsymbol{\mu}_z + \boldsymbol{\sigma}_z \odot \boldsymbol{\epsilon}_z$
- 10:          $[\boldsymbol{\mu}_{G_{post}}, \log \boldsymbol{\sigma}_{G_{post}}^2] \leftarrow \mathcal{D}_{perturb}(\mathbf{z})$
- 11:          $\mathcal{L}_{recon} \leftarrow \mathbb{E}_{q_\phi}[-\log p_\psi(G_{post}|\mathbf{z})]$
- 12:          $\mathcal{L}_{KL} \leftarrow \text{DKL}(q_\phi(\mathbf{z}|G_{pre}, D) || p(\mathbf{z}))$
- 13:          $\mathcal{L}_{VAE} \leftarrow \mathcal{L}_{recon} + \mathcal{L}_{KL}$
- 14:          $\mathbf{X}_{image}^0 \leftarrow \mathcal{E}_{image}(I)$
- 15:          $t \sim \mathcal{U}(\{1, \dots, T\})$
- 16:          $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 17:          $\mathbf{X}_{image}^t \leftarrow \sqrt{\bar{\alpha}_t} \mathbf{X}_{image}^0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$
- 18:          $\mathbf{X}_{time} \leftarrow \mathcal{E}_{time}(t)$
- 19:          $\mathbf{X}_{condition} \leftarrow \mathbf{z} + \mathbf{X}_{time}$
- 20:          $\boldsymbol{\epsilon}_\theta \leftarrow f_\theta(\mathbf{X}_{image}^t, \mathbf{X}_{condition})$
- 21:          $\mathcal{L}_{LDM} \leftarrow \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta\|^2$
- 22:          $\mathcal{L}_{TRIDENT} \leftarrow \mathcal{L}_{VAE} + \mathcal{L}_{LDM}$
- 23:         Update parameters  $\phi, \psi, \theta, \gamma$  using  $\nabla \mathcal{L}_{TRIDENT}$
- 24:     **end for**
- 25: **end for**
- 26: **return** Trained parameters  $\Theta = \{\phi, \psi, \theta, \gamma\}$

---



---

### Algorithm 2 TRIDENT Framework: Inference Procedure

---

**Require:** Input  $G_{pre}, D$   
**Require:** Trained parameters  $\Theta = \{\phi, \psi, \theta, \gamma\}$   
**Require:** Pre-trained image VAE ( $\mathcal{E}_{image}, \mathcal{D}_{image}$ )  
**Require:** Diffusion timesteps  $T$  and schedule  $\{\beta_t\}_{t=1}^T$   
**Ensure:** Generated Image  $\hat{I}$

- 1:  $\mathbf{X}_{rna} \leftarrow \mathcal{E}_{rna}(G_{pre})$
- 2:  $\mathbf{X}_{drug} \leftarrow \mathcal{E}_{drug}(D)$
- 3:  $[\boldsymbol{\mu}_z, \log \boldsymbol{\sigma}_z^2] \leftarrow \mathcal{E}_{perturb}([\mathbf{X}_{rna}, \mathbf{X}_{drug}])$
- 4:  $\boldsymbol{\epsilon}_z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5:  $\mathbf{z} \leftarrow \boldsymbol{\mu}_z + \boldsymbol{\sigma}_z \odot \boldsymbol{\epsilon}_z$
- 6:  $\hat{\mathbf{X}}_{image}^T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 7: **for**  $t = T$  **down to** 1 **do**
- 8:      $\mathbf{X}_{time} \leftarrow \mathcal{E}_{time}(t)$
- 9:      $\mathbf{X}_{condition} \leftarrow \mathbf{z} + \mathbf{X}_{time}$
- 10:      $\boldsymbol{\epsilon}_\theta \leftarrow f_\theta(\hat{\mathbf{X}}_{image}^t, \mathbf{X}_{condition})$
- 11:      $\boldsymbol{\epsilon}' \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\boldsymbol{\epsilon}' \leftarrow \mathbf{0}$
- 12:      $\hat{\mathbf{X}}_{image}^{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( \hat{\mathbf{X}}_{image}^t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta \right) + \sigma_t \boldsymbol{\epsilon}'$
- 13: **end for**
- 14:  $\hat{I} \leftarrow \mathcal{D}_{image}(\hat{\mathbf{X}}_{image}^0)$
- 15: **return**  $\hat{I}$

---

## F. Additional Comparison Results

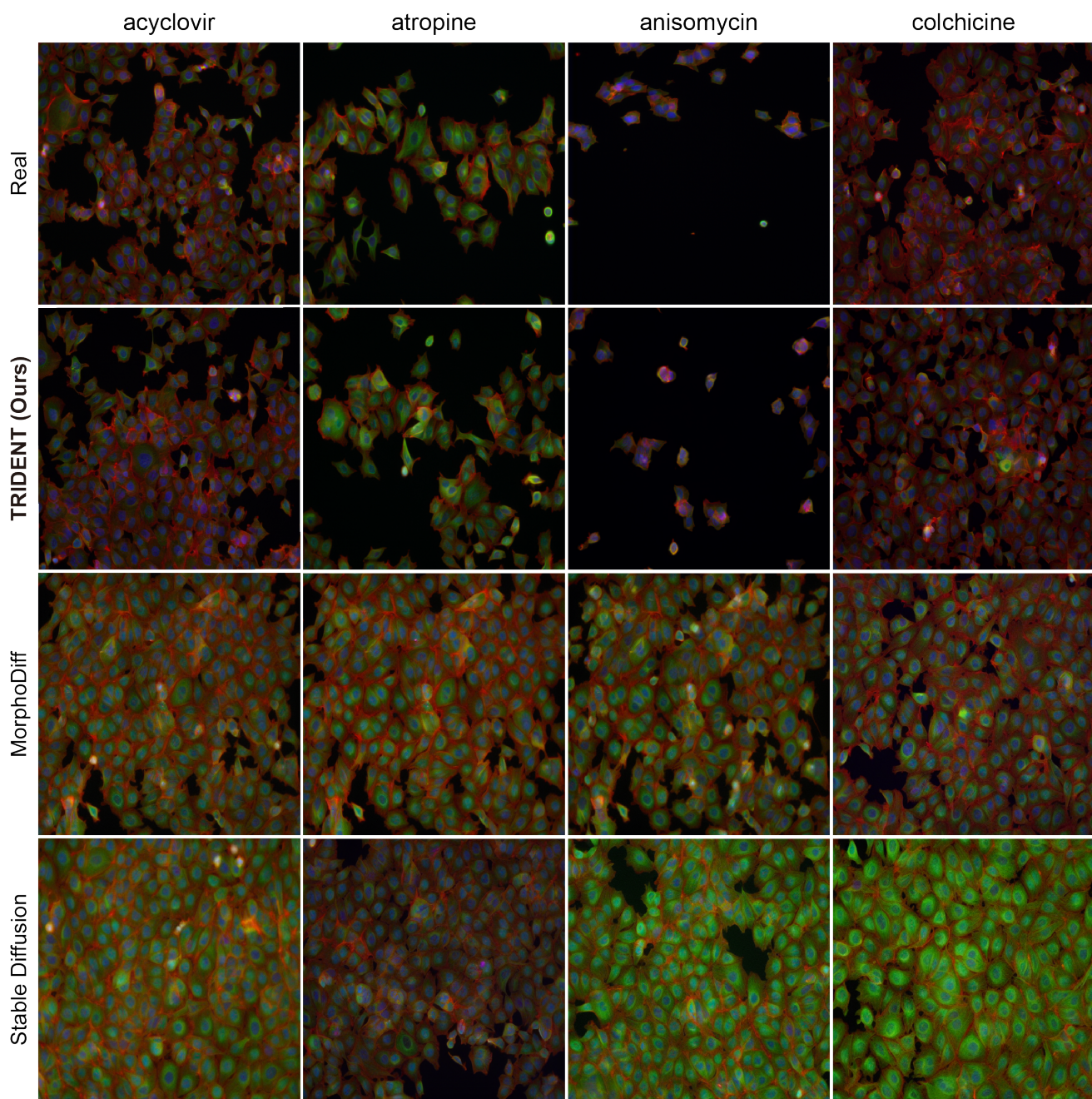


Figure A1. Additional visual comparison of generated cellular morphologies. Ground-truth images (Row 1) are compared to outputs from TRIDENT (Row 2), MorphoDiff (Row 3), and Stable Diffusion (Row 4).

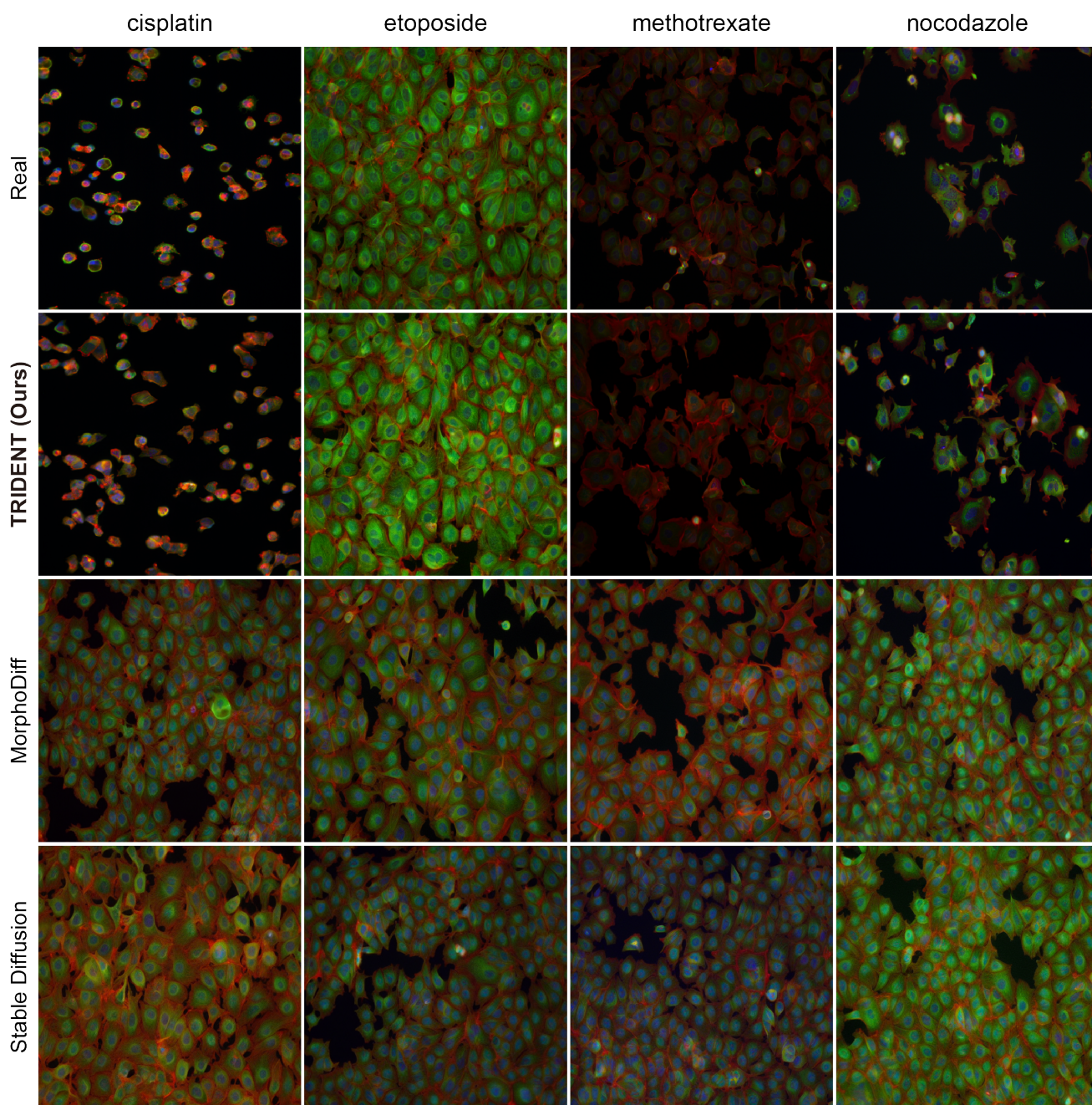


Figure A2. Additional visual comparison of generated cellular morphologies. Ground-truth images (Row 1) are compared to outputs from TRIDENT (Row 2), MorphoDiff (Row 3), and Stable Diffusion (Row 4).

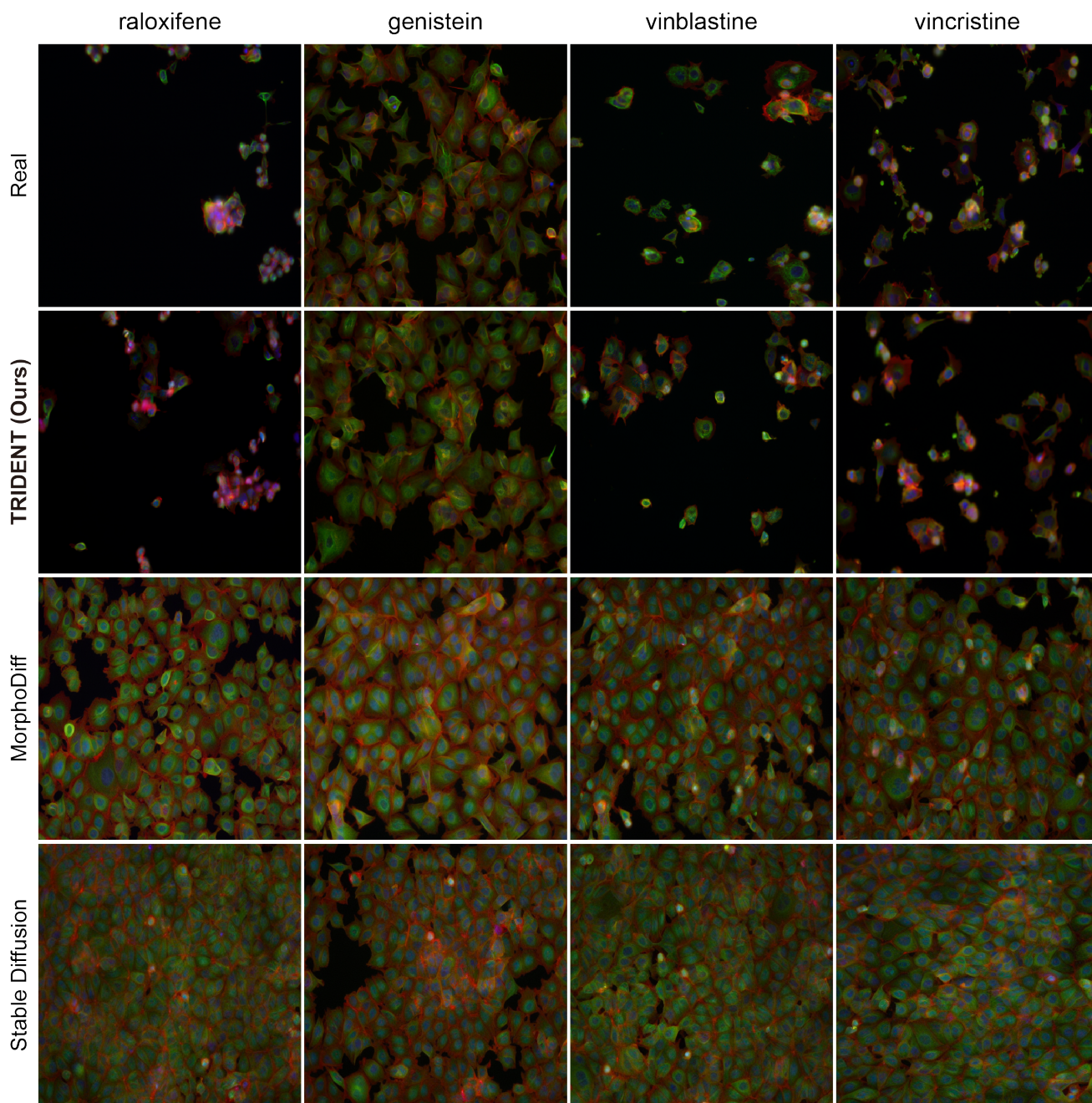


Figure A3. Additional visual comparison of generated cellular morphologies. Ground-truth images (Row 1) are compared to outputs from TRIDENT (Row 2), MorphoDiff (Row 3), and Stable Diffusion (Row 4).