CELLPAINTR: CONTRASTIVE BATCH CORRECTION TRANSFORMER FOR LARGE SCALE CELL PAINTING

Anonymous authors

Paper under double-blind review

Abstract

Cell Painting, a high-content imaging-based profiling method, has emerged as a powerful tool for understanding cellular phenotypes and drug responses. However, batch effects severely constrain the integration and interpretation of data collected across different laboratories and experimental conditions. To mitigate this issue, here we introduce CellPainTR, a novel embedding approach through Tranformer for unified batch correction and representation learning of Cell Painting data, thereby addressing a critical challenge in the field of image-based profiling. Our approach employs a Transformer-like architecture with Hyena operators, positional encoding via morphological-feature-embedding, and a special source context token for batch correction, combined with a multi-stage training process that incorporates masked token prediction and supervised contrastive learning. Experiments on the JUMP Cell Painting dataset demonstrate that CellPainTR significantly outperforms existing approaches such as Combat and Harmony across multiple evaluation metrics, while maintaining strong biological information retention as evidenced by improved clustering metrics and qualitative UMAP visualizations. Moreover, our method effectively reduces the feature space from thousands of dimensions to just 256, addressing the curse of dimensionality while maintaining high performance. These advancements enable more robust integration of multi-source Cell Painting data, potentially accelerating progress in drug discovery and cellular biology research.

029 030 031

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

032

034

1 INTRODUCTION

Cell Painting, a powerful high-content imaging technique, has emerged as a promising tool for 035 biological research and drug discovery. This method generates rich, multidimensional data capturing intricate cellular phenotypes and responses to perturbations. The potential applications of Cell 037 Painting are vast, from identifying novel drug candidates to understanding disease mechanisms at the cellular level (Bray et al., 2016; Gustafsdottir et al., 2013; Ljosa et al., 2013). Unfortunately, fully realizing the potential of Cell Painting is hindered by significant challenges in data integration 040 and interpretation. While recent approaches explore direct learning from cellular images, most es-041 tablished pipelines rely on extracted Cell Painting features (morphological measurements extracted 042 from Cell Painting assays, in our case using CellProfiler software) Ando et al. (2017); Celik et al. 043 (2022); Borowa et al. (2024); Korsunsky et al. (2019) - Typical Cell Painting datasets consist of thousands of such engineered features, leading to the "curse of dimensionality" and computation-044 ally intensive analyses (Caicedo et al., 2017). Moreover, these datasets are prone to batch effects systematic variations unrelated to biology - due to differences in experimental conditions, imaging 046 equipment, and data processing across laboratories (Singh et al., 2017). 047

To address these critical limitations, we propose CellPainTR, a novel foundational model designed
 to learn robust representations of Cell Painting features. By operating directly in the feature space,
 our approach uniquely integrates batch correction and representation learning into a single model.
 Through a three-stage training process, CellPainTR learns to distinguish batch-specific artifacts from
 biologically meaningful features, enabling adaptive balancing of batch effect correction and preservation of underlying biological signals. Leveraging a Transformer-like architecture (Vaswani et al., 2017) and Hyena operators (Poli et al., 2023), the model efficiently handles long-range dependencies



Figure 1: Cell Painting workflow and CellPainTR training. Data acquisition – (a) Cell Model Preparation: Cells of interest are seeded into multi-well plates. (b) Perturbation: Cells undergo var-ious treatments or perturbations (e.g., drug compounds, genetic modifications). (c) Staining: Cells are stained with a mix of fluorescent dyes to highlight different cellular components. (d) Imag-ing: High-content screening using fluorescence microscopy captures detailed images of the stained cells. Morphological feature extraction – (e) Image Analysis: Software like CellProfiler processes the raw images to extract quantitative features, resulting in high-dimensional feature vectors for each cell. CellPainTR training - (f) CellPainTR Processing: The extracted feature vectors are input into CellPainTR, our novel Transformer-based model designed for unified batch correction and representation learning. Finally, CellPainTR produces batch-corrected and dimensionality-reduced representations of the cell profiles, enabling more robust integration and analysis of multi-source Cell Painting data.

and high-dimensionality inherent in Cell Painting data while maintaining direct compatibility with
 established analysis workflows and preserving the interpretability critical to biological research.

We evaluate the effectiveness of CellPainTR using a comprehensive set of metrics, including batch correction indices and compound-specific pattern preservation metrics from the Broad Institute benchmark (Arevalo et al., 2024). Experimental result 1 demonstrate that CellPainTR significantly outperforms existing approaches such as Combat and Harmony (Johnson et al., 2007; Korsunsky et al., 2019) across multiple evaluation metrics. Key results include superior batch correction scores (0.80 with controls, 0.69 without controls) compared to Combat (0.56, 0.37) and Harmony (0.57, 0.40), while maintaining strong biological information retention as evidenced by improved cluster-ing metrics and qualitative UMAP visualizations. The versatility of our foundational model allows it to be fine-tuned for specific tasks and adapted to diverse experimental setups, enhancing its utility across different applications.

This work has significant implications for the broader field of image-based profiling and machine
 learning in biological research. By providing a more robust and adaptable solution to the batch effect
 problem, CellPainTR paves the way for more reliable and scalable analyses of Cell Painting data,
 with the potential to accelerate drug discovery and deepen our understanding of cellular biology.

2 BACKGROUND

Classical Batch Correction Methods. The challenge of batch effects in high-dimensional biolog ical data, particularly in the field of image-based profiling like Cell Painting, has been a significant focus of research in recent years (Arevalo et al., 2024; Ando et al., 2017; Celik et al., 2022; Kraus

et al., 2024; Borowa et al., 2024). Several classical algorithms have been developed to address batch effects in biological data. One prominent example is Combat, which uses a parametric empirical Bayes framework to adjust for batch effects (Johnson et al., 2007). The Combat method models the batch effect as an additive and multiplicative effect on the data, and estimates these effects using an empirical Bayes approach. This can be expressed mathematically as:

$$y_{ij} = \alpha_i + \beta_i x_j + \gamma_b + \delta_b x_j + \epsilon_{ij} \tag{1}$$

where y_{ij} is the observed data, x_j is the covariate of interest, α_i and β_i are the sample-specific intercept and slope, γ_b and δ_b are the batch-specific intercept and slope, and ϵ_{ij} is the residual error.

Another method, Harmony, uses iterative clustering and linear adjustment to correct batch effects in multi-modal single-cell data (Korsunsky et al., 2019). Harmony aims to identify a shared low-dimensional representation across batches by aligning cluster centroids in an iterative fashion. This can be expressed as:

$$x_i = W_b h_i + b_b \tag{2}$$

where z_i is the corrected representation for sample *i*, h_i is the original high-dimensional representation, W_b and b_b are the batch-specific linear transformation parameters.

While these classical methods have demonstrated effectiveness in many scenarios, they may struggle to handle the high dimensionality and complex batch effects present in Cell Painting datasets (Singh et al., 2024).

Self-Attention and the Hyena Operator. The self-attention operator (Vaswani et al., 2017) is a fundamental mechanism of Transformers. Specifically, given a sequence $x \in \mathbb{R}^{L \times D}$ with length L and D features, the self-attention operator A(x) is defined as:

$$A(x) = \sigma(xW_q)(xW_k)^T(xW_v) \tag{3}$$

where $W_q, W_k, W_v \in \mathbb{R}^{D \times D}$ are learnable projection matrices, and σ is a softmax operation. This allows the model to capture pairwise relationships between tokens in the sequence. However, one limitation is that self-attention becomes computationally expensive for long sequences, with a complexity of $\mathcal{O}(L^2)$. To address the computational challenge of self-attention, the Hyena operator (Poli et al., 2023) was introduced as a replacement for self-attention in Transformers. The Hyena operator is characterized by a structured self-attention mechanism that involves long convolutions and element-wise gating:

143

132 133

113

114

122

 $y_t = (h * u)_t = \sum_{\tau=0}^{L-1} h_{t-\tau} u_{\tau}.$ (4)

In standard convolutional architectures, the filter length ℓ is typically constrained by $\ell \ll L$, where L is the input sequence length. This constraint helps control computational costs. However, by parameterizing the filter as a function of the temporal offset τ (i.e., $h_{\tau} = \gamma_{\theta}(\tau)$), we can design extended convolution kernels without a proportional increase in parameters. This technique, known as implicit convolution, enables efficient modeling of long-range dependencies. The implicit convolution mechanism is exemplified by the Hyena operator, which employs a recursive framework incorporating extended convolutions and point-wise modulation:

 $y = x^{N} \cdot (h^{N} * (x^{N-1} \cdot (h^{N-1} * (\cdots x^{1} \cdot (h^{1} * v)))))$ (5)

Here, v denotes the initial input, $\{x^i\}_{i=1}^N$ represent successive transformations, N indicates the recursion depth, * symbolizes convolution, and \cdot denotes Hadamard (element-wise) product.

Relevance to CellPainTR. The Hyena operator's ability to efficiently model long-range dependencies in high-dimensional data makes it particularly relevant to the CellPainTR model. By incorporating the Hyena operator, CellPainTR can effectively handle the complex feature interactions present in Cell Painting data, which is critical for addressing batch effects while preserving biologically relevant information. Furthermore, the Bidirectional Hyena (Oh et al., 2023) extension, which removes the temporal causality, is particularly well-suited for the Cell Painting domain, where morphological feature interactions are not constrained by sequential order. This non-causal, bidirectional mechanism aligns with the requirements of the CellPainTR model.

Figure 2: CellPainTR architecture. (a) CellPainTR incorporates several novel innovations: a linear adaptor layer combined with Morphological Feature Embeddings, a source context token (SRC), and a bidirectional Hyena operator. (b) Input embedding layer architecture. (c) After multi-step training, CellPainTR can be applied to various downstream tasks, including compound classification, compound retrieval, qualitative analysis, and batch correction. (d) Hyena Operator architecture.

3 CELLPAINTR

162 163

164

165

166 167

168

169

170 171 172

173 174

175

176

177

178

179 180

181 182

183

184

185

186

187

199

209

This section describes our novel feature embedding approach to unified batch correction and representation learning for Cell Painting data. We introduce a Transformer-like architecture that effectively handles the high dimensionality and complex relationships inherent in Cell Painting features while simultaneously addressing batch effects and preserving compound-specific molecular mechanisms of action (MoA) patterns (Arevalo et al., 2024; Moshkov et al., 2022; Chen et al., 2023).

CellPainTR's architecture and training process are specifically designed to address the unique chal-188 lenges presented by Cell Painting data (see Fig. 2(a)). The high dimensionality of the data is tackled 189 through the use of Hyena operators, which efficiently capture long-range dependencies across the 190 extensive feature space. The linear adaptor for morphological feature embedding ensures that the 191 full spectrum of continuous features is preserved, addressing the challenge of information loss of-192 ten associated with discretization methods. The feature context embedding mechanism replaces 193 traditional positional encoding, better capturing the intrinsic relationships between morphological 194 features that are critical in Cell Painting data. Finally, the source-specific token and multi-stage 195 training process directly address the batch effect problem by learning to distinguish between source-196 specific variations and true biological signals. This comprehensive approach enables CellPainTR 197 to simultaneously correct for batch effects and learn biologically meaningful representations, a key requirement for effective analysis of multi-source Cell Painting datasets. More details are as follows. 198

200 3.1 DESIGN ARCHITECTURE

201 Linear Adapter for Morphological Feature Embedding. Cell Painting data consists of thousands 202 of morphological features, each representing a specific aspect of cellular structure and organization. 203 To effectively capture the rich information in this high-dimensional data, we introduce a linear adap-204 tor module to embed the continuous features without any loss of information (Weisbart et al., 2024). 205 Unlike traditional approaches that rely on discrete token representations, our linear adaptor maps the 206 original feature space directly to the model's input embeddings as shown in Figure 2(b). This allows 207 the CellPainTR model to operate on the full spectrum of the morphological features, preserving the 208 nuanced relationships between them. For each feature *i*, the embedding is computed as:

$$E_i = \mathbf{W}_i \cdot C_i + \mathbf{b}_i \quad \text{where } E_i \in \mathbb{R}^{d_{\text{model}}}$$
 (6)

where $\mathbf{W}_i \in \mathbb{R}^{d_{\text{model}}}$ is a learnable weight matrix, $C_i \in \mathbb{R}$ is the input feature value, and $\mathbf{b}_i \in \mathbb{R}^{d_{\text{model}}}$ is a learnable bias vector. This linear adaptor approach ensures that the CellPainTR model can effectively leverage the full contextual information present in the Cell Painting data, laying the foundation for robust batch correction and representation learning (Tromans-Coia & Jamali, 2023).

Feature Context Embedding. Traditional Transformer models use positional encoding to incorporate the sequential nature of the input data. However, in the case of Cell Painting features, the order

of the features does not necessarily reflect any meaningful biological context. To better capture the intrinsic relationships between the morphological features, we introduce a feature context embedding mechanism (Seal et al., 2024), as shown in Figure 2(b). Specifically, the feature context embedding replaces the standard positional encoding by learning an embedding of the feature context. In this approach, each morphological feature is encoded with its own embedding $(M_1, M_2, ..., M_L)$ with $M_i \in \mathbb{R}^{1 \times d_{model}}$, which is then added to the expression embeddings. The complete feature embedding matrix is then constructed by concatenating all feature embeddings:

$$\mathbf{E} = [M_1; M_2; \dots; M_L] \in \mathbb{R}^{L \times d_{\text{model}}}$$
(7)

where *L* is the number of morphological features and d_{model} is the embedding dimension, and $[\cdot; \cdot]$ denotes concatenation along the row direction. This method allows us to provide the CellPainTR model with explicit feature context information. By using the feature context embedding, the Cell-PainTR model can learn to associate the morphological features with their biological context, rather than relying on their position in the input sequence. This enhances the model's ability to capture the complex interdependencies between the features, which is crucial for effective batch correction and representation learning.

Source Context Token. Cell Painting datasets often originate from multiple experimental sources, each with its own unique batch-related characteristics. To explicitly model these source-specific variations, CellPainTR incorporates a special source context token (Weisbart et al., 2024). The source context token S with dimension $S_{dim} = M_{dim}$ is initialized as a learnable parameter, drawn from its own embedding with a vocabulary size K with K the number of source in the dataset; $(S_1, S_2, ..., S_K)$ and is concatenated with the input feature embeddings, as shown in Figure 2(a). We incorporate a learnable source embedding that is concatenated to the feature embeddings:

$$\mathbf{S}_k = \text{Embedding}(k) \quad \text{where } k \in \{1, \dots, K\}, \mathbf{H} = [\mathbf{E}; \mathbf{S}_k] \quad \text{where } \mathbf{S}_k \in \mathbb{R}^{1 \times d_{\text{model}}}$$
(8)

where K is the number of sources in the dataset. During training, the model learns to associate the source context token with the unique batch effects present in each data source. This allows the model to adaptively correct for batch-related biases while preserving the biologically relevant information in the learned representations. The inclusion of the source context token is a key innovation that enables the CellPainTR model to handle the challenges of integrating Cell Painting data from diverse experimental sources, a critical requirement for advancing drug discovery and cellular biology research (Tromans-Coia & Jamali, 2023).

248 249 250

239 240

224

3.2 TRAINING

The CellPainTR model is trained using a multi-stage process that combines unsupervised and supervised learning objectives to achieve unified batch correction and representation learning. By progressively exposing the model to increasingly diverse data contexts, we enable it to learn robust, generalizable representations that preserve compound-specific MoA relationships while mitigating batch-related confounders. More details are as follows.

Channel-Wise Masked Morphology (CWMM). The initial stage of training utilizes Channel-Wise
 Masked Morphology (CWMM), a novel approach inspired by Masked Language Modeling (MLM)
 in natural language processing and Masked Expression Modeling (MEM) in single-cell RNA se quencing analysis (Oh et al., 2023). CWMM is tailored to handle the continuous values of morpho logical features in Cell Painting data while respecting its channel-wise structure (Seal et al., 2024).
 For additional details on CWMM refer to Appendix A.1

262 As shown in Fig. 3(a), in the CWMM task, a subset of input morphological features is randomly 263 masked, with the model tasked to predict these masked values based on the surrounding context. 264 Features are grouped based on both their channel origin and the cellular compartment they describe 265 (e.g., "DNA channel - Nucleus" features form one group, "Mito channel - Cytoplasm" features form 266 another). The masking probability for each training batch is chosen from a range of [0.05, 0.4]267 and is applied uniformly across all feature sets, preserving the biological relationships within each channel-compartment combination. Importantly, only non-zero values are masked and replaced with 268 a [MASK] token, as distinguishing between true and false zero values is not feasible in this context 269 (Way et al., 2021).



Figure 3: CellPainTR training. (a) Channel-Wise Masked Morphology (CWMM) for selfsupervised learning. (b) Intra-Source Supervised Contrastive Learning: combining CWMM and single source supervised contrastive learning. (c) Inter-Source Supervised Contrastive Learning: combining CWMM and multi source supervised contrastive learning.

Mathematically, the objective function for the CWMM pre-training task is formulated as:

$$\ell_{\text{CWMM}} = \frac{1}{G} \sum_{g=1}^{G} \frac{1}{|M_g|} \sum_{i \in M_g} (F_{g,i} - F'_{g,i})^2$$
(9)

With G the number of feature groups (channel-compartment combinations), M_q represents the set of 290 masked indices for feature group g, $|M_q|$ the number of masked features in group g, $F_{q,i}$ denotes the 291 true value of the *i*-th morphological feature in group g, $F'_{g,i}$ the predicted value for the *i*-th masked feature in group g. Through this pre-training process, CellPainTR acquires generalizable features 292 293 that capture the biological meaning encoded in the morphological data, the contextual relationships 294 between features, and the channel-wise dependencies within the Cell Painting data structure. This 295 comprehensive understanding of the intricate patterns and correlations present in Cell Painting data 296 establishes a robust foundation for the subsequent supervised learning stages, ultimately enhancing 297 the model's capacity for batch correction and representation learning. For additional details, please refer to Appendix A.2 298

299 Intra Source Supervised Learning. Following the unsupervised CWMM pretraining, the model 300 undergoes a fine-tuning phase using a supervised contrastive learning approach within each data 301 source. This intra-source supervised learning stage is designed to encourage the model to learn 302 representations that are both discriminative and invariant to batch-related variations within a given 303 experimental source (plate effects, liquid handling, reagent batches, ...). As shown in Fig. 3(b), 304 during this phase, we allow biological feature metadata to flow to the model during training (more specifically InChIKey via the supervised contrastive objective), while ensuring that each batch con-305 tains data from only a single source. The objective function for this stage combines the CWMM loss 306 with a supervised contrastive loss, weighted equally: 307

$$\ell_{\text{intra}} = \ell_{\text{CWMM}} + \ell_{\text{supcon}} \tag{10}$$

The supervised contrastive loss ℓ_{supcon} is computed directly using the CLS token output of the encoder. This can be expressed as:

$$\ell_{\text{supcon}} = -\sum_{i \in I} \frac{1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(h_i \cdot h_p/\tau)}{\sum_{a \in A(i)} \exp(h_i \cdot h_a/\tau)}$$
(11)

Here, *I* is the set of indices, P(i) is the set of positives for sample *i* (i.e., samples with the same biological label - InChIKey), A(i) is the set of all samples except *i*, h_i and h_p are the normalized CLS token outputs for samples *i* and *p*, respectively, and τ is a temperature parameter. This approach allows the model to learn representations that capture biologically relevant information while remaining robust to source-specific batch effects. For additional details, please refer to Appendix A.3

Inter Source Supervised Learning. The final stage of training involves fine-tuning the model
 using a supervised contrastive learning objective that spans multiple data sources. This inter-source
 supervised learning step enables the model to learn representations that are not only batch-corrected
 but also generalize well across diverse experimental conditions and data sources.

282 283 284

279

280

281

287 288

289

308

311 312

313 314 As shown in Fig. 3(c), the key distinction in this stage is that we now allow different sources to be mixed within the same batch, while source context token specific to each source are optimized during the training. Thanks to the source context token, the supervised contrastive loss now operates across samples from multiple sources, encouraging the model to learn representations that are invariant to source-specific variations (microscope settings, experimental setups, ...) while still capturing biologically meaningful information. For additional details, please refer to Appendix A.4

4 EXPERIMENTS

330 331

332

Our experimental framework was designed to rigorously evaluate the effectiveness of our proposed method for batch correction and representation learning in Cell Painting data. We conducted a series of experiments using a large-scale dataset, applying a multi-stage training process and assessing performance through both qualitative and quantitative measures. Qualitatively, we examined the model's ability to preserve clustering patterns that align with known compound Mechanisms of Action (MoA). Quantitatively, we used established metrics from the Broad Institute benchmark (Arevalo et al., 2024), to measure the preservation of compound-specific cellular responses.

340 Dataset. For our experiments, we utilized the cpg-0016 (Chandrasekaran et al., 2023) Cell Painting 341 dataset from the JUMP consortium. This extensive dataset encompasses a diverse array of biologi-342 cal perturbations, including over 100,000 small molecule compounds and genetic perturbations such 343 as ORF overexpression and CRISPR knockouts targeting thousands of genes. The dataset captures 344 a wide spectrum of morphological features, typically numbering in the thousands per cell, which provide detailed measurements of cell shape, size, texture, and intensity across multiple cellular 345 compartments (Bray et al., 2016). The dataset includes both negative and positive controls, which 346 are crucial for establishing baselines and validating assay performance. Negative controls, such 347 as non-targeting controls or DMSO treatments, provide a reference point for normal cellular mor-348 phology (Caie et al., 2010). Positive controls, consisting of known bioactive compounds or genetic 349 perturbations with well-characterized effects, serve to validate the assay's sensitivity and specificity 350 (Gustafsdottir et al., 2013). 351

Our preprocessing pipeline was designed to address common challenges in Cell Painting data. First, we imputed all infinite and missing values with zero to ensure computational feasibility. Next, we applied MAD (Median Absolute Deviation) normalization, using the negative control of each plate as the baseline. This step helps to mitigate plate-to-plate variations and standardize feature scales. Finally, we employed a clipping strategy, constraining values between the 0.01 and 0.99 quantiles. This step was crucial in managing extreme outliers, which are common artifacts in Cell Painting data and can significantly impact model performance if left unaddressed.

- To evaluate the robustness and generalizability of our model, we experimented with various traintest split strategies. These included hiding certain data source generations from training, excluding specific plates, and random partitioning of the data (Goodfellow et al., 2016). This approach allowed us to choose the best-performing model under different scenarios of data availability and batch effects.
- Trainining. The pretraining phase of our model utilized all compound-related data from the dataset.
 During this phase, the model learned to predict masked morphological features based on the surrounding context, leveraging the bidirectional nature of the Hyena blocks. This approach enabled the model to capture complex relationships within the data effectively, laying the groundwork for subsequent supervised learning stages (Devlin et al., 2018).
- Following pretraining, we moved to a more focused training phase using a curated subset of the 369 data. This subset concentrated on compounds with rich metadata, including control compounds 370 and those with known Mechanisms of Action (MoA) (Schurer et al., 2011). This phase, which took 371 approximately one week on the same hardware, allowed the model to refine its representations based 372 on more specific biological contexts (Goodfellow et al., 2016). During this stage, we introduced the 373 source-specific token, enabling the model to learn and account for source-specific variations in the 374 data (Devlin et al., 2018). This approach was crucial in addressing batch effects while preserving 375 biologically relevant information. 376
- 377 The final stage of our training process involved fine-tuning the model with a strong emphasis on contrastive learning. This phase also lasted approximately one week and used the same curated

Source Mechanism of Action

Figure 4: UMAP visualizations. Comparing Baseline (uncorrected), CellPainTR at its second training step, and the final CellPainTR model. Columns represent methods or training stages, and rows depict data aspects. The top row uses MoA coloring for phenotypic variations, with red for negative controls. The bottom row uses source coloring for batch effects. The Baseline (left) shows strong batch effects and weak spatial constraints in compound clusters, indicating a lack of robust biological signals. CellPainTR's improvement is seen in the middle and right columns. The second step (middle) shows clear source-based clusters with reduced intra-source fragmentation. The final step (right) achieves cohesive batch integration and clear compound-specific patterns, supported by superior Batch Correction (0.76) and Graph Connectivity (0.84) scores, demonstrating effective batch variability mitigation while preserving biological signals.

dataset as the training phase (Schurer et al., 2011). The contrastive learning approach encouraged
the model to learn representations that are invariant to batch effects while still capturing meaningful
biological differences (Chen et al., 2020). By contrasting samples from different sources but with
similar biological properties, the model learned to distinguish between batch-related variations and
true biological signals (He et al., 2020).

406 4.1 QUALITATIVE EVALUATION RESULTS

Figure 4 illustrates CellPainTR's effectiveness in batch correction and biological signal preserva-tion for Cell Painting data. The visualizations demonstrate the model's ability to account for batch variability stemming from different laboratories, batches and microscopes while maintaining cru-cial biological information. The top row uses compound coloring to represent biological variation, with red indicating the negative control. The bottom row employs source coloring to highlight batch effects. In the compound-colored row (top), CellPainTR demonstrates progressively clearer separa-tion of compounds (each represented by a unique color), especially in the final stage (right). This improved clustering reflects an enhanced preservation of biological information. Simultaneously, the source-colored row (bottom) highlights CellPainTR's effectiveness in mitigating source-specific batch effects, with the final stage showing the most integrated distribution of data points across sources. The improved clustering and reduced overlap in CellPainTR's final stage, compared to both the baseline and its intermediate stage, underscore its superior performance in balancing batch correction with biological signal retention. This visualization demonstrates CellPainTR's superior performance in balancing batch correction with biological signal preservation, addressing a key challenge in integrating high-dimensional microscopy data from diverse sources. For more detailed comparison, please refer to Appendix E

4.2 QUANTITATIVE EVALUATION RESULTS

Table 1 presents a comprehensive comparison of batch correction methods, including our proposed
CellPainTR approach, across various metrics. These metrics can be broadly categorized into: (1)
batch correction measures, which assess the removal of technical variations; (2) compound-specific
pattern preservation indicators (labeled as 'Biological Metrics' in the table), which evaluate how well
the method maintains compound-related effects and known mechanism of action relationships; and
(3) aggregate scores, which average each and both aspects. For consistent interpretation, all metrics
have been normalized to a scale of 0 to 1, where 0 indicates poor performance and 1 represents optimal performance C.2.

Table 1: Performance comparison of batch correction methods

	Batch Correction				Biological Metrics						Aggregate Scores			
Method	Graph Conn.	Sil. Batch	Batch Corr. (control)	Batch Corr. (no control)	Leiden NMI	Leiden ARI	Sil. Label	mAP (control)	mAP (no rep)	Bio Info (control)	Bio Info (no control)	Batch Corr.	Bio Metrics	Overall Score
Baseline	0.86	0.93	0.57	0.40	0.41	0.20	0.50	0.48	0.58	0.91	0.86	0.69	0.56	0.63
Combat Johnson et al. (2007)	0.85	0.93	0.56	0.37	0.39	0.12	0.50	0.48	0.58	0.91	0.85	0.68	0.55	0.61
Harmony Korsunsky et al. (2019)	0.80	0.93	0.57	0.40	0.42	0.24	0.50	0.47	0.58	0.91	0.86	0.68	0.57	0.62
Sphering Kessy et al. (2018)	0.64	0.95	0.70	0.58	0.36	0.35	0.48	0.12	0.23	0.66	0.85	0.72	0.43	0.58
CellPainTR(1)	0.78	0.73	0.78	0.69	0.34	0.26	0.52	0.22	0.32	0.72	0.81	0.75	0.46	0.60
CellPainTR(2)	0.69	0.75	0.71	0.58	0.43	0.15	0.70	0.54	0.63	0.84	0.91	0.68	0.60	0.64
CellPainTR	0.84	0.70	0.80	0.69	0.35	0.17	0.57	0.40	0.54	0.86	0.79	0.76	0.53	0.64

442 To assess the effectiveness of batch correction, we employed several complementary metrics. Graph 443 Connectivity evaluates the preservation of biological relationships across batches by measuring the 444 proportion of k-nearest neighbors that are maintained after batch correction, compared to the orig-445 inal data structure. This metric provides insights into how well the local structure of the data is 446 preserved while removing batch effects, with values closer to 1 indicating better preservation of bio-447 logical relationships. The Silhouette Batch score quantifies the degree of separation between batches 448 by comparing the mean distance between samples from different batches to the mean distance within 449 batches, where values closer to 1 suggest more effective batch integration. Additionally, we introduce a Batch Correction metric calculated as 1 - f1 score, where f1 represents the performance of 450 a classifier trained to detect batch effects in the corrected data. In this context, higher values indi-451 cate superior batch correction as the classifier struggles to distinguish between batches, effectively 452 measuring the degree of batch effect removal. 453

- 454 The preservation of compound-specific patterns is crucial in ensuring that the batch correction process does not compromise the underlying MoA-related information. We quantify this preservation 455 through multiple complementary approaches. Through the Leiden clustering algorithm, we com-456 pute Normalized Mutual Information (NMI) and Adjusted Rand Index (ARI) scores by comparing 457 cluster assignments with known biological labels. These metrics provide different perspectives on 458 biological preservation: NMI values range from 0 to 1, indicating the degree of shared information 459 between cluster assignments and biological labels, while ARI provides a chance-corrected mea-460 sure of agreement, with values above 0.3 indicating strong biological preservation. The Silhouette 461 Label score further strengthens our evaluation by assessing how well samples cluster according to 462 biological conditions rather than batch effects, with values above 0.5 suggesting robust biological 463 signal retention. To evaluate the impact on downstream analysis tasks, we employed mean Average 464 Precision (mAP) scores (see Appendix C.1) for compound retrieval performance in two contexts: 465 with controls (mAP control) and without replicate compounds (mAP no rep). These scores specifically assess the model's ability to identify similar biological conditions across batches, providing a 466 practical measure of biological signal preservation where values above 0.5 indicate strong retention 467 of compound-specific patterns. Aggregate scores provide a holistic view of each method's perfor-468 mance. The Batch Correction score summarizes the overall effectiveness in removing batch effects, 469 while the Biological Metrics score encapsulates the method's ability to preserve biological infor-470 mation. The Overall Score combines these aspects to give a comprehensive assessment of each 471 method's performance. 472
- Our proposed CellPainTR method demonstrates strong performance across these metrics. Notably, 473 it achieves the highest batch correction scores both with and without controls, indicating superior 474 batch integration capabilities. The progressive improvement observed across the three steps of Cell-475 PainTR (denoted as CellPainTR(1), CellPainTR(2), and CellPainTR) highlights the method's ability 476 to balance batch effect removal with biological information preservation. Specifically, CellPainTR 477 consistently outperforms traditional methods like Combat and Harmony in key metrics. While Com-478 bat and Harmony show strengths in certain areas, such as graph connectivity and silhouette batch 479 scores, CellPainTR demonstrates a more balanced performance across all metrics. This is partic-480 ularly evident in the higher batch correction scores and competitive biological metric scores. The 481 final CellPainTR model achieves the highest overall score, tied with CellPainTR(2), suggesting that 482 our approach successfully addresses batch effects while maintaining crucial biological information. This balance is critical in ensuring that batch correction does not come at the cost of losing important 483 biological signals. In summary, the quantitative results presented in Table 1 provide strong evidence 484 for the efficacy of CellPainTR in batch correction tasks. The method's ability to consistently perform 485 well across various metrics, particularly in batch correction and biological information retention, po-

sitions it as a robust solution for addressing batch effects in Cell Painting data analysis. For further
 details on the metrics implementation please refer to Appendix C.

4.3 Ablation Study

To elucidate the contribution of each component in our CellPainTR method, we conducted an ablation study by evaluating the model's performance at different stages of training. The results are
presented in Table 1, where CellPainTR(1) and CellPainTR(2) represent intermediate stages, and
CellPainTR denotes the final model.

495 The initial training stage, represented by CellPainTR(1), demonstrated a strong focus on batch effect 496 removal. This stage achieved the highest batch correction score without control (0.69), surpassing 497 both the baseline and traditional methods in this metric. However, this came at the cost of reduced 498 performance in biological metrics, particularly evident in the lower mAP scores (0.22 and 0.32) and 499 decreased graph connectivity (0.78) compared to the baseline (0.86). CellPainTR(2), the intermedi-500 ate training stage, marked a significant shift towards preserving biological information. This stage showed substantial improvements in biological metrics, most notably in Silhouette Label (0.70) and 501 mAP scores (0.54 and 0.63). Moreover, it achieved the highest Leiden NMI (0.43) among all meth-502 ods, suggesting improved cluster separation. These gains in biological signal preservation were 503 accompanied by the best overall biological metrics score (0.60). However, this stage also saw a de-504 crease in batch correction performance compared to CellPainTR(1) and a further reduction in graph 505 connectivity (0.69). The final CellPainTR model emerged as a balanced solution, optimizing both 506 batch correction and biological signal preservation. It achieved the best overall batch correction 507 score (0.76), effectively balancing correction with and without control. The model also recovered 508 graph connectivity (0.84) compared to the intermediate stages while maintaining strong performance 509 in biological information preservation (0.86 for control). Despite these improvements, there was a 510 slight decrease in some biological metrics compared to CellPainTR(2), particularly in Silhouette Label and Leiden NMI. 511

This ablation study reveals a critical trade-off between batch correction and biological signal preservation throughout the training process. The initial stage prioritizes batch effect removal, potentially at the expense of biological signal retention. The intermediate stage then shifts focus to preserving biological information, enhancing clustering and representation quality. Finally, the complete model achieves an optimal balance between these competing objectives. For qualitative comparison of the steps, please refer to Appendix E

518 519

520

489

490

5 CONCLUSION

521 This paper introduces CellPainTR, a novel approach to unified batch correction and representation 522 learning for high-dimensional Cell Painting data. Leveraging a Transformer-like architecture with 523 Hyena operators, CellPainTR addresses the critical challenges of batch effects and dimensional-524 ity reduction in image-based profiling. Extensive experimental results confirmed that CellPainTR 525 successfully addresses batch effects while preserving crucial biological information. However, we 526 acknowledge the limitations of our study, including the restricted comparison to only three other batch correction methods due to difficulties in replicating results from benchmark papers (Arevalo 527 et al., 2024). Additionally, we observed a discrepancy between the qualitative improvements seen in 528 the UMAP visualizations and some of the quantitative metrics, highlighting the complexity of eval-529 uating batch correction methods and the potential limitations of current evaluation approaches (for 530 more details see Appendix D). Despite these challenges, CellPainTR represents a significant step for-531 ward in the analysis of Cell Painting data, offering a powerful approach to integrating multi-source 532 datasets while maintaining compound-specific molecular signature preservation. The method's abil-533 ity to handle high-dimensional data, correct for batch effects, and preserve biological information 534 positions it as a valuable tool for advancing drug discovery and cellular biology research.

535 536

537 ETHICS STATEMENT

538

The utilization of models such as CellPainTR offers significant benefits for advancing our understanding of complex biological systems and potentially improving medical research. However, eth540 ical considerations must guide its use to ensure responsible data handling, avoid biases, and protect 541 individual privacy, underscoring the importance of ethical guidelines and regulations in the applica-542 tion of such models.

Reproducibility Statement

546 We provide detailed implementation information in Section 3.2 and additional details in Appendix A. A comprehensive description of the datasets used in our experiments can be found in Section 4 548 and using the official dataset link: https://github.com/jump-cellpainting/datasets and precision about dataset curation can be found in Appendix B. Our source code is available for access at the following link: https://github.com/CellPainTR/CellPainTR.

551 552 553

554

555

556

557

558

559

562

563

565 566

567

543 544

545

547

549

550

REFERENCES

- D. Michael Ando, Cory Y. McLean, and Marc Berndl. Improving phenotypic measurements in high-content imaging screens. bioRxiv, 2017.
- John Arevalo, Ellen Su, Jessica D. Ewald, Robert van Dijk, Anne E. Carpenter, and Shantanu Singh. Evaluating batch correction methods for image-based cell profiling. *Nature Communications*, 15: 50613, 2024.
- Adriana Borowa et al. Decoding phenotypic screening: A comparative analysis of image represen-561 tations. Computational and Structural Biotechnology Journal, 23:1181–1188, 2024.
 - M.-A. Bray, S. Singh, H. Han, C. T. Davis, B. Borgeson, C. Hartland, and A. E. Carpenter. Cell painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. Nature Protocols, 11:1757–1774, 2016.
 - J. C. Caicedo, S. Cooper, F. Heigwer, S. Warchal, P. Qiu, C. Molnar, and A. E. Carpenter. Dataanalysis strategies for image-based cell profiling. Nature Methods, 14:849-863, 2017.
- 568 Peter D. Caie, Russell E. Walls, Anna Ingleston-Orme, Sandeep Daya, Miles D. Houslay, Richard 569 Eagle, and Neil O. Carragher. High-content phenotypic profiling of drug response signatures 570 across distinct cancer cells. Molecular cancer therapeutics, 9(6):1913–1926, 2010. 571
- Safiye Celik, Jan-Christian Huetter, Sandra Melo, Nathan Lazar, Rahul Mohan, Conor Tillinghast, 572 Tommaso Biancalani, Marta Fay, Berton Earnshaw, and Imran S. Haque. Biological cartography: 573 Building and benchmarking representations of life. In NeurIPS 2022 Workshop on Learning 574 Meaningful Representations of Life, 2022. 575
- 576 Srinivas Niranj Chandrasekaran, Jeanelle Ackerman, Eric Alix, D. Michael Ando, John Arevalo, 577 Melissa Bennion, Nicolas Boisseau, Adriana Borowa, Justin D. Boyd, Laurent Brino, Patrick J. 578 Byrne, Hugo Ceulemans, Carolyn Ch'ng, Beth A. Cimini, Djork-Arne Clevert, Nicole Deflaux, 579 John G. Doench, Thierry Dorval, Regis Doyonnas, Vincenza Dragone, Ola Engkvist, Patrick W. Faloon, Briana Fritchman, Florian Fuchs, Sakshi Garg, Tamara J. Gilbert, David Glazer, David 580 Gnutt, Amy Goodale, Jeremy Grignard, Judith Guenther, Yu Han, Zahra Hanifehlou, Santosh Hariharan, Desiree Hernandez, Shane R. Horman, Gisela Hormel, Michael Huntley, Ilknur Icke, 582 Makiyo Iida, Christina B. Jacob, Steffen Jaensch, Jawahar Khetan, Maria Kost-Alimova, Tomasz 583 Krawiec, Daniel Kuhn, Charles-Hugues Lardeau, Amanda Lembke, Francis Lin, Kevin D. Little, 584 Kenneth R. Lofstrom, Sofia Lotfi, David J. Logan, Yi Luo, Franck Madoux, Paula A. Marin Zap-585 ata, Brittany A. Marion, Glynn Martin, Nicola Jane McCarthy, Lewis Mervin, Lisa Miller, Haseeb 586 Mohamed, Tiziana Monteverde, Elizabeth Mouchet, Barbara Nicke, Arnaud Ogier, Anne-Laure Ong, Marc Osterland, Magdalena Otrocka, Pieter J. Peeters, James Pilling, Stefan Prechtl, Chen 588 Qian, Krzysztof Rataj, David E. Root, Sylvie K. Sakata, Simon Scrace, Hajime Shimizu, David 589 Simon, Peter Sommer, Craig Spruiell, Iffat Sumia, Susanne E. Swalley, Hiroki Terauchi, Aman-590 dine Thibaudeau, Amy Unruh, Jelle Van de Waeter, Michiel Van Dyck, Carlo van Staden, Michał Warchoł, Erin Weisbart, Amélie Weiss, Nicolas Wiest-Daessle, Guy Williams, Shan Yu, Bolek Zapiec, Marek Żyła, Shantanu Singh, and Anne E. Carpenter. Jump cell painting dataset: mor-592 phological impact of 136,000 chemical and genetic perturbations - cpg0016-jump dataset, 2023. https://registry.opendata.aws/cellpainting-gallery/.

594	Tine Chan Simon Kemblish Mahammad Narausi and Caeffron Hinton A simple from and for				
595 596	contrastive learning of visual representations. <i>arXiv preprint arXiv:2002.05709</i> , 2020.				
597	Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrev E. Hinton. Claire: contrastive				
598	learning-based batch correction framework for single-cell rna sequencing. <i>Bioinformatics</i> , 3				
599	btad099, 2023.				
600					
601	Jacob Devlin, Ming-wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training o				
602	ordirectional transformers for language understanding. arxiv preprint arxiv:1810.04805, 201				
603	Ian Goodfellow, Yoshua Bengio, and Aaron Courville. <i>Deep learning</i> , volume 1. MIT Press, 2016.				
604	S. M. Gustafsdottir, V. Ljosa, K. L. Sokolnicki, J. A. Wilson, D. Walpita, M. M. Kemp, and A. E.				
605	Carpenter. Multiplex cytological profiling assay to measure diverse cellular states. <i>PLoS ONE</i> , 8:				
606	e80999, 2013.				
607	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsu-				
608	pervised visual representation learning. Proceedings of the IEEE/CVF Conference on Computer				
609	Vision and Pattern Recognition, pp. 9729–9738, 2020.				
610					
611 612	W. Evan Johnson, Cheng Li, and Ariel Rabinovic. Adjusting batch effects in microarray expression data using empirical bayes methods. <i>Biostatistics</i> , 8:118–127, 2007.				
613	Agence Kassy Alay Lawin and Kashinian Steinman Ontimal whitening and decompletion. The				
614	Agnan Kessy, Alex Lewin, and Korolinian Strimmer. Optimal writening and decorrelation. American Statistician $72(4):300, 314, 2018$				
615	American Statistician, 72(+).507–514, 2010.				
616	Ilya Korsunsky, Nathan Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yu				
617	Baglaenko, Michael Brenner, Po-Ru Loh, and Soumya Raychaudhuri. Fast, sensitive and ac-				
618	curate integration of single-cell data with harmony. <i>Nature Methods</i> , 16:1289–1296, 2019.				
619	O Kraus K Kenvon-Dean S Saberian M Fallah P McLean I Leung V Sharma A Khan				
620	J. Balakrishnan, S. Celik, and D. Beaini. Masked autoencoders for microscopy are scalable learn-				
621	ers of cellular biology. In Proceedings of the IEEE/CVF Conference on Computer Visio				
622	Pattern Recognition, pp. 11757–11768, 2024.				
623	V Liosa V I. Sakalniaki and A. E. Carnantar. Annotated high throughout microscopy image sate				
624 625	for validation. <i>Nature Methods</i> , 9:637, 2013.				
626	Ilva Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Con				
627	ence on Learning Representations, 2019. URL https://openreview.net/forum?i				
628	Bkg6RiCqY7.				
629	Nilita Machkay Michael Domhaldt Santiago Danait Matthew Smith Claim McQuin Allan Cood				
630	man Rebecca Senft Vu Han Mehrtash Rabadi Peter Horvath Reth A Cimini Anne E Carpen				
631	ter Shantanu Singh and Juan C. Caicedo Learning representations for image-based profiling of				
632	perturbations. <i>bioRxiv</i> , 10:503783, 2022.				
633	I man have been a second s				
634	Gyutaek Oh, Baekgyu Choi, Inkyung Jung, and Jong Chul Ye. schyena: Foundation model for				
635	full-length single-cell rna-seq analysis in brain. <i>bioRxiv</i> , 2023.				
636	Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y. Fu, Tri Dao, Stephen Baccus, Yoshua				
637	Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional				
638	language models. arXiv preprint arXiv:2302.10866, 2023.				
039	General C. Caluman Harron D. Manuard, Dishard D. Carid, Association A. Challet and N. 197, C. et al.				
640	Bioassay ontology appotations facilitate cross analysis of diverse high throughput screening date				
640	sets <i>Journal of hiomolecular screening</i> 16(4):415–426 2011				
642	5-6. Vommun of Diomorcomun screening, 10(7),715-720, 2011.				
643	Srijit Seal, Jordi Carreras-Puigvert, Shantanu Singh, Anne E. Carpenter, Ola Spjuth, and Andreas				
644	Bender. From pixels to phenotypes: Integrating image-based profiling with cell health data as				
646	biomorph features improves interpretability. <i>Molecular Biology of the Cell</i> , 35:mr2, 1–13, 2024				
647	S. Singh, MA. Bray, T. R. Jones, and A. E. Carpenter. Pipeline for illumination correction of images for high-throughput microscopy. <i>Nature Protocols</i> , 12:1709–1725, 2017.				

- S. Singh, M.-A. Bray, T. R. Jones, and A. E. Carpenter. Evaluating batch correction methods for image-based cell profiling. *bioRxiv*, 15:558001, 2024.
- 651 Nathan Tromans-Coia and Ali Jamali. Jump-moa plate data. *Preprint*, 2023.
 - A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30: 5998–6008, 2017.
- 656 Gregory P. Way et al. Cell painting gallery: an open resource for image-based profiling. *Nature*, 2021.
 - Nathan Weisbart et al. Cell painting gallery. AWS Registry of Open Data, 2024.
 - A EXPERIMENTAL DETAILS
- For all downstream tasks, we fine-tuned the model using the AdamW optimizer (Loshchilov &
 Hutter, 2019). The fine-tuning process was executed using one NVIDIA Quadro RTX 4000.
- 666 A.1 CHANNEL-WISE MASKED MORPHOLOGY (CWMM)

Channel-Wise Masked Morphology (CWMM) is a novel training objective designed for Cell Painting data analysis. Inspired by Masked Language Modeling (MLM) (Devlin et al., 2018) in natural
language processing and Masked Expression Modeling (MEM) in single-cell RNA sequencing (Oh
et al., 2023), CWMM adapts these concepts to address the continuous values and channel-wise structure of morphological features. In CWMM, a subset of input morphological features is randomly
and uniformaly masked with respect to each channel during training. The model then predicts these
masked values based on surrounding context. Key aspects of the implementation include:

- Masking probability: Randomly selected from a range of 0.05 to 0.4 for each training batch.
- 675 676 677

679

680

685

686 687

688

689

690

652

653

654

655

658

659 660 661

662

665

667

- Uniform application: Masking is applied equally across all five channels in Cell Painting data.
- 678
- Non-zero value focus: Only non-zero values are masked and replaced with a [MASK] token.

The CellPainTR model, using bidirectional Hyena blocks, predicts the true values of masked features. This process requires consideration of both preceding and following features, capturing complex relationships within the data. The CWMM objective function is formulated as in Section 3.2 Eq. (9).

- CWMM offers several advantages:
 - Encourages learning of generalizable features capturing biological meaning in morphological data.
 - Facilitates understanding of contextual relationships and channel-wise dependencies.
 - Establishes a foundation for subsequent supervised learning, enhancing capabilities in batch correction and representation learning.
- 691 692 693

By adapting proven techniques to Cell Painting data's specific challenges, CWMM provides a powerful tool for extracting insights from complex morphological datasets. This approach has the potential to advance our understanding of cellular phenotypes and their relationships to underlying biological processes, benefiting fields from drug discovery to basic cell biology research.

- 698 A.2 PRETRAINING STEP 1
- 699
- 700 The pretraining phase of our model utilized all compound-related data from the dataset. We employed a configuration with 3 recurrences and 4 Hyena layers, training for 3 epochs. The model was optimized using AdamW (Devlin et al., 2018) with a learning rate of 1e-4 and a batch size of 16 for

approximately two weeks. The model also had access to the source context token, as described in
 Section 3.1, but had no constraint or direct objective regarding this supplementary context. During
 this step, all compound data was used, including negative controls, positive controls, and unknown
 compounds. This approach enabled the model to effectively capture complex relationships within
 the data, facilitating the learning of lower-dimensional representations.

707 708

709

A.3 PRETRAINING - STEP 2

710 The second training phase focused on compounds with rich metadata, including controls and those with known Mechanisms of Action (MoA). We reduced the learning rate to 1e-5, increased the 711 batch size to 32, and ran this phase for approximately one week. This allowed the model to refine its 712 representations based on more specific biological contexts. During this stage, we introduced intra-713 source supervised contrastive learning alongside a reconstruction constraint (Channel-Wise Masked 714 Morphology). To implement this, we restricted each batch to a single source, ensuring positive and 715 negative pairs were intra-source. The first epoch served as a warm-up, allowing backpropagation 716 only through the learnable source context token. Subsequently, we unfroze the rest of the model for 717 the remaining epochs. This dual approach enabled the model to account for source-specific varia-718 tions while preserving biological signals, effectively addressing intra-source batch effects without 719 losing relevant information.

720 721

722

A.4 PRETRAINING - STEP 3

The final stage of our training process involved fine-tuning the model with a strong emphasis on contrastive learning. We maintained the learning rate at 1e-5 but increased the batch size to 64.
This phase, lasting approximately one week, used the same curated dataset as the second training phase. Unlike the previous step, we shuffled the data to include multiple sources in each batch, compelling the model to learn source-invariant representations. This approach aimed to preserve biological signals while mitigating inter-source batch effects.

B DATASET

730 731

729

732 As supplementary context, drawn from the official dataset paper (Chandrasekaran et al., 2023) and 733 the associated GitHub page, we provide additional details on dataset curation. As described in Ap-734 pendix A, in the first step, we utilize all data related to compound perturbations. The full dataset includes Open Reading Frame (ORF) perturbations, Clustered Regularly Interspaced Short Palin-735 dromic Repeats (CRISPR) perturbations, and compound perturbations. Using the associated meta-736 data files, we specifically select only the compound-related perturbations. In the second and third 737 steps, we further curate the data by using the metadata to select compounds with a referenced Mech-738 anism of Action (MoA), or those identified as positive controls or positive compound pairs. 739

- 740 741 C BASELINE METHOD
- 742 743

744

745 746

747 748

749

750 751

752

754

All methods used in this study follow the implementation provided by the scib Python library, as

described in the benchmark paper by (Arevalo et al., 2024).

C.1 METRICS

Mean Average Precision (mAP) Evaluation:

- 1. Each sample in the dataset serves as a query
- 2. Positive elements: Other biological replicates sharing the same compound (identified by InChIKey)
 - 3. Negative elements for mAP (control): Control wells from the same plate
- 4. Negative elements for mAP (no rep): Wells from the same plate treated with different compounds

5. Higher scores indicate better preservation of compound-specific biological effects after batch

The mean average precision (mAP) metric is also implemented in accordance with the benchmarkpaper (Arevalo et al., 2024).

For each query, there are M - 1 positive samples, which share the same compound, and N negative samples, consisting of profiles from the same plate with different compounds or negative controls. A ranked list is generated using cosine similarity between the query profile and the (M - 1) + Nprofiles. The average precision (AP) of the i^{th} query is defined as:

756

758

767 768 769 $AP_i = \sum_{k=1}^{(M-1)+N} (R_k - R_{k-1})P_k$ (12)

where:

770 771 $R_k = \frac{TP_k}{M-1}$ is recall at rank k

772
773
774
$$P_k = \frac{TP_k}{k}$$
 is precision at rank

 TP_k is the number of positive elements retrieved up to rank k

Finally, the mean average precision (mAP) for a compound is the average of the AP values across all replicates:

779

775 776

 $mAP = \frac{1}{M} \sum_{i=1}^{M} AP_i \tag{13}$

784 C.2 EVALUATION METRICS: MATHEMATICAL FORMULATIONS AND LIMITATIONS

ROLE OF CONTROLS IN METRIC CALCULATION

785 786

800

802

803 804

805

C.2.1

⁷⁸⁷ In our evaluation framework, metrics are reported in two variants: "with controls" and "without controls". This distinction relates directly to the negative controls present in each plate, which are

controls". This distinction relates directly to the negative controls present in each plate, which are also used as baseline in the Median Absolute Deviation normalization process.

With Controls Metrics calculated "with controls" include negative control wells in the evaluation.
 Since these controls serve as the normalization baseline, they represent the most standardized condition across plates. Performance metrics in this setting typically show better results as these samples are inherently more aligned across batches due to their role in the normalization process.

Without Controls The "without controls" variant excludes negative control wells, focusing exclusively on compound-treated samples. This provides a more stringent and realistic assessment of batch correction performance, as it evaluates the method's effectiveness on samples that weren't used in the normalization process.

- 801 Interpretation Considerations
 - 1. Better performance in "with controls" metrics is expected due to the standardization process
 - 2. "Without controls" metrics provide a more challenging benchmark for real-world application
- 806
 807
 808
 3. The gap between these variants can indicate how well the batch correction generalizes beyond standardized samples
- 4. Methods showing consistent performance across both variants suggest robust batch correction capabilities

0.2.2 DA	ATCH CORRECTION METRICS	
Graph Cor	nnectivity Mathematical formulation:	
	$GC = \frac{ \{(i,j) \in kNN_{before} \cap kNN_{after}\} }{ k < \pi}$	(14)
where kNN	$k \times n$ N _{before} and kNN_{after} represent the k-nearest neighbors before and	l after batch correc-
ion, n is the	he number of samples, and k is the number of neighbors (we use $k=15$	5).
Limitations	s:	
1. Se	ensitive to the choice of k	
2. Ma	lay not capture global structure changes	
3. Ca	an be biased by density differences between batches	
Batch Corn with differe	crection and Biological Information Scores For both metrics, we ent objectives and interpretations:	use a classifier but
	$BC = 1 - F1_{batch\ classifier}$	(15)
	$BI = F1_{compound_classifier}$	(16)
where $F1_{bb}$	<i>batch_classifier</i> is the F1 score of the classifier trained to predict l	patch identifier and to different MoAs).
For a given	classification task, the F1 score is calculated as:	
<u>8</u>	$_$ $precision \times recall$	
	$F1 = 2 imes rac{1}{precision + recall}$	(17)
The rational	ale behind these metrics is:	
1. Fo the	or Batch Correction (BC): Higher values (closer to 1) indicate better the classifier fails to distinguish between batches	batch correction as
2. Fo of dif	or Biological Information (BI): Higher values (closer to 1) indicate f biological signals as the classifier successfully distinguishes betwe fferent MoAs	better preservation en compounds with
Limitations	S:	
1. Ma	lay not capture very subtle batch effects or biological signals	
2. Re	equires sufficient samples per batch/compound for reliable estimation	ı
3. F1	1 score can be affected by the distribution of classes in the dataset	
Alternative n available provide dire ame classif	e metrics: While established metrics like kBET exist in the field, to e implementations at the time of development led to our current app rect insight into both batch effect removal and biological signal pre- ification framework.	echnical limitations broach. Our metrics eservation using the
С.2.3 Ви	IOLOGICAL METRICS	
Leiden Clu	ustering Metrics (NMI & ARI) For two clusterings U and V:	
	NMI(U V) = - 2 imes I(U,V)	(18)
	H(U) + H(U) + H(V)	(10)
where I(U,V	V) is the mutual information and H(U), H(V) are the entropies.	
	$ARI = \frac{\mathbf{KI} - E[\mathbf{KI}]}{\max(\mathbf{PI}) - E[\mathbf{PI}]}$	(19)
	$\max(\mathbf{KI}) = E[\mathbf{KI}]$	
where RI is	s the Rand Index and E[R]] is its expected value	

864	1. Both metrics are sensitive to the number of clusters
865	2 ARI can be pessimistic with many small clusters
866	2. Pasulta demond on Leiden algorithm normators
867	5. Results depend on Leiden argorithin parameters
868 869	Mean Average Precision (mAP) For a query q:
870 871	$AP(q) = \sum_{k=1}^{n} P(k) \times rel(k) $ ⁽²⁰⁾
872	$\kappa=1$ where $\mathbf{P}(k)$ is precision at cutoff k rel(k) is relevance of k-th item
873	
074 875	Limitations:
876	1. Sensitive to the number of relevant items
877	2. May not capture complex biological relationships
878	3 Can be biased by compound representation imbalance
879	5. Can be blased by compound representation inibiliance
880	C.2.4 EDGE CASES AND COMMON FAILURE MODES
881	1 Graph Connectivity may fail when:
882	(a) Detah offacts are confounded with higherical signals
883	(a) Batch effects are confounded with biological signals (b) Date her very different scales perces betches
884	(b) Data has very different scales across batches
885	(c) Extreme outliers are present
886	2. Biological metrics may be unreliable when:
887	(a) Very few replicates are available
888	(b) Compound effects are subtle
889	(c) Multiple mechanisms of action overlap
801	3. General considerations:
802	(a) All metrics assume sufficient sample size per batch
893	(b) Metrics may be less reliable with highly imbalanced batches
894	(c) Strong batch effects might mask biological signals in all metrics
895	
896	C.2.5 AGGREGATED SCORES CALCULATION
897	The aggregated scores presented in Table 1 are calculated as weighted averages of their respective
898	constituent metrics:
899	
900	Batch Correction Score The Batch Correction aggregate score is calculated as the arithmetic
901	mean of the batch-related metrics:
902	
903	$BC_{max} = \frac{GC + S_{batch} + BC_{control} + BC_{no_control}}{(21)}$
904	$\frac{1}{4}$
905	where GC is the Graph Connectivity, S_{batch} is the Silhouette Batch score, $BC_{control}$ is the Batch
900	Correction score with controls, and $BC_{no_control}$ is the Batch Correction score without controls.
908	
909	Biological Metrics Score The Biological Metrics aggregate score combines all biology-related
	matriagi

metrics:

913

914

$$Bio_{score} = \frac{NMI + ARI + S_{label} + mAP_{control} + mAP_{no_rep} + BI_{control} + BI_{no_control}}{7}$$
(22)

915 where NMI is the Normalized Mutual Information from Leiden clustering, ARI is the Ad-916 justed Rand Index from Leiden clustering, S_{label} is the Silhouette Label score, $mAP_{control}$ and 917 $mAP_{no.rep}$ are the Mean Average Precision scores with controls and without replicates respectively, and $BI_{control}$ and $BI_{no.control}$ are the Biological Information scores with and without controls. 918
 919
 919
 919
 920
 Overall Score The Overall Score is calculated as the arithmetic mean of the Batch Correction and Biological Metrics scores:

$$Overall_{score} = \frac{BC_{score} + Bio_{score}}{2}$$
(23)

This balanced approach ensures that both batch correction effectiveness and biological signal preservation contribute equally to the final evaluation of method performance.

C.3 DATA PREPROCESSING

For the baseline data preprocessing shown in Table 1 (Baseline), we follow the steps outlined in the benchmark paper (Arevalo et al., 2024):

- 1. Variation Filtering: Features with low variance are filtered using the absolute coefficient of variation, with a threshold of $C_{\text{var}} < 1e^{-3}$.
- 2. Median Absolute Deviation: For each well, feature values are normalized using the median \bar{X} and the absolute deviation $\tilde{\sigma}$, both calculated from control wells.
- 3. **Rank-based Inverse Normal Transformation (INT)**: Feature values are transformed based on their rank within each plate, following Blom's formula:

$$Y_i = \Phi^{-1} \left(\frac{r_i - c}{N - 2c + 1} \right)$$

where r_i is the rank of sample *i*, *N* is the number of samples, and $c = \frac{3}{8}$.

4. **Feature Selection**: Features are selected using a correlation threshold, excluding features with a correlation above 0.9 with other features.

D DISCUSSION

D.1 FEATURE SPACE PROCESSING IN BIOLOGICAL REPRESENTATION LEARNING

Our approach to representation learning in Cell Painting data distinguishes itself through a deliberate focus on engineered feature space processing. Unlike approaches that directly operate on raw image data, our method leverages CellProfiler-generated features, offering several critical advantages:

- **Biological Interpretability:** By maintaining direct mapping to established cellular measurements, we ensure that learned representations retain meaningful biological context. Each transformed feature can be traced back to its original morphological or biochemical interpretation, a crucial requirement for biological research.
- **Methodological Innovation:** Our approach bridges traditional feature-based analysis with modern representation learning techniques. We demonstrate that sophisticated machine learning transformations can be achieved while preserving the semantic meaning of individual cellular features.
- **Practical Workflow Integration:** The method supports immediate adoption within existing Cell Painting protocols, addressing a critical need in drug discovery and biological research workflows.
- D.2 BATCH CORRECTION AND REPRESENTATION LEARNING CHALLENGES

Batch correction in biological data presents unique challenges that our approach systematically addresses:

- 1. Removing technical variations while preserving biological signals
 - 2. Maintaining interpretability of cellular measurements
 - 3. Enabling cross-experimental comparability

972 973	4. Handling high-dimensional, complex biological datasets
974	
975	batch specific artifacts from biologically magningful features
976	bach-specific artifacts from biologicarly meaningful readures.
977	D 2 LIMITATIONS AND MODEL CENEDALIZADU ITV
978	D.5 LIMITATIONS AND MODEL GENERALIZABILITY
979	The current implementation of CellPainTR introduces both significant advances and notable limita-
980	tions:
981	Current Constraints:
982	
983	• Fixed pretrained source context tokens limit direct processing of unseen data sources
984	• Computational intensity of the training process
985	Reliance on CellProfiler-generated features
986	Renalice on common generated reactions
987 988 980	Adaptation Strategies: We propose a straightforward extension mechanism for handling new data sources:
990	1. Extend the source embedding layer
991	2. Fine-tune new source context embeddings
992	3 Follow the established training procedure
993	5. Follow the established training procedure
994 995	D.4 FUTURE RESEARCH DIRECTIONS
996 997	The methodological framework we introduce opens several promising avenues for future research:
998	• Developing more dynamic source context token adaptation mechanisms
999	• Exploring end-to-end learning approaches for cellular imaging
1000	Exploring one to one tourning approaches for contrain maging
1001	• Expanding applicability across diverse biological datasets
1002	• Investigating transfer learning strategies for cellular phenotype representation
1003 1004	D.5 BROADER IMPACT
1005	
1006 1007 1008 1009	Beyond technical innovations, our work contributes to a broader scientific objective: making ad- vanced machine learning techniques more accessible and meaningful in biological research. By maintaining interpretability and providing a robust methodological framework, we offer a critical stepping stone for more transparent and impactful computational approaches in cellular analysis.
1010 1011 1012	The approach demonstrates that sophisticated representation learning can be achieved within con- strained, interpretable feature spaces, challenging the prevailing notion that advanced machine learn- ing requires complete abstraction from domain-specific measurements.
1013 1014	D.6 PRACTICAL CONSIDERATIONS FOR MODEL STAGES AND COMPUTATIONAL EFFICIENCY
1015 1016 1017 1018	The proposed three-stage approach—CellPainTR(1), CellPainTR(2), and CellPainTR—offers flex- ibility depending on the requirements of the analysis. Each stage reflects a trade-off between com- putational demands, batch correction efficacy, and biological signal preservation.
1019 1020 1021 1022	CellPainTR(1): This stage serves as a foundational representation learning model, achieving a good starting point for biological signal extraction. However, it retains pronounced batch effects between data sources, making it more suitable for exploratory analyses or scenarios where inter-source batch effects are less critical (Figure 7 6 5, Tab 1).
1023 1024 1025	CellPainTR(2): Optimized for single-source datasets, CellPainTR(2) effectively addresses intra- source batch effects while retaining biological signals with high fidelity. This model is particularly useful for studies focused on single-source data, such as analyzing cellular responses within a spe- cific experimental batch or environment (Figure (Figure 7 6 5, Tab 1).

CellPainTR: The final stage represents the most balanced model for multi-source datasets, excelling in inter-source batch correction. While there is a minor loss in biological signal fidelity compared to CellPainTR(2), this trade-off enables robust integration and analysis across large-scale, multi-source data, making it well-suited for meta-analyses and cross-experiment studies (Figure 7 6 5, Tab 1).

These distinctions provide practical guidance for researchers on selecting the appropriate model stage based on their experimental setup and goals. The additional visualizations E included in the appendix further illustrate these trade-offs, emphasizing the nuanced improvements achieved by the final stage.

 Computational Efficiency and Reproducibility: The training times reported in this study reflect an unoptimized setup, involving a single Quadro 4000 GPU and no advanced strategies such as distributed training or mixed precision. The reported times are provided for transparency and to guide users who may wish to replicate the experiments using the current public codebase.

1039 To reduce computational time, several straightforward optimizations can be implemented:

- Utilizing modern GPUs with higher processing power.
- Incorporating distributed training techniques for parallel computation.
- Leveraging mixed precision to reduce memory requirements and training time.

It is important to note that the primary focus of this work is the methodological design and its effectiveness in addressing batch correction and biological signal retention, rather than computational efficiency. Nonetheless, these optimizations can significantly lower the training cost without altering the methodology.

By providing these clarifications, we aim to facilitate the reproducibility of our results while ensuring
that researchers can adapt our approach efficiently to their specific computational resources and
research goals.

E VISUALIZATION



Figure 5: Complete PCA Comparison. In this figure, we compare all methods listed in Table 1.
Although the dimensionality reduction method used is PCA, we observe clear structures emerging from the CellPainTR(2) method, indicating strong signal preservation and effective optimization within the learned representation.



Figure 6: **Complete t-SNE Comparison**. In this figure, we compare all methods listed in Table 1. For all steps of CellPainTR, we observe a better organization in the t-SNE, suggesting strong retention of biological signals. Additionally, for CellPainTR(1) and CellPainTR, the microscope and source batch effects are largely corrected. On the other hand, CellPainTR(2) shows strong fragmentation, which is expected due to the intra-source nature of its objective. Interestingly, CellPainTR(2) organizes its manifold by compound regions. Since CellPainTR(2) does not have access to different source compounds simultaneously, this emergent organization is noteworthy.



Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. In this figure, we compare all methods listed in Table Figure 7: **Complete UMAP Comparison**. Sugression of the baseline methods, we observe a strong fragmentation and CellPainTR(1) and CellPainTR manage this fragmentation effectively. However, due to the training nature of CellPainTR(2), the UMAP also shows significant fragmentation. Despite this, the intra-source clusters are more cohesive compared to the baseline methods, indicating the method's success in reducing batch effects within individual sources.