

Triple-effect correction for Cell Painting data with contrastive and domain-adversarial learning

Received: 10 February 2025

Accepted: 14 July 2025

Published online: 25 July 2025

 Check for updates

Chengwei Yan^{1,4}, Yu Zhang^{2,4}, Jiuxin Feng^{1,4}, Heyang Hua², Zhihan Ruan¹, Zhen Li³, Siyu Li², Chaoyang Yan¹, Pingjing Li¹, Jian Liu¹✉ & Shengquan Chen²✉

Cell Painting (CP), as a high-throughput imaging technology, generates extensive cell-stained imaging data, providing unique morphological insights for biological research. However, CP data contains three types of technical effects, referred to as triple effects, including batch effects, gradient-influenced row and column effects (well-position effects). The interaction of various technical effects can obscure true biological signals and complicate the characterization of CP data, making correction essential for reliable analysis. Here, we propose cpDistiller, a triple-effect correction method specially designed for CP data, which leverages a pre-trained segmentation model coupled with a semi-supervised Gaussian mixture variational autoencoder utilizing contrastive and domain-adversarial learning. Through extensive qualitative and quantitative experiments across various CP profiles, we demonstrate that cpDistiller effectively corrects triple effects, especially well-position effects, while preserving cellular heterogeneity. Moreover, cpDistiller effectively captures system-level phenotypic responses to genetic perturbations and reliably infers gene functions and interactions both when combined with scRNA-seq data and independently. cpDistiller also demonstrates promising capability for identifying gene and compound targets, highlighting its potential utility in drug discovery and broader biological research.

Advanced high-dimensional assay technologies, such as transcriptomics and epigenomics profiling, offer remarkable depth and breadth in molecular-level biological research¹. Despite their strengths, these technologies often focus exclusively on specific molecular changes, lacking the capability to observe changes at the system level of cell state, which involves many complex and unknown processes. To obtain information at the cellular system level, high-throughput imaging technologies have been developed to produce useful profiles of cell phenotypes by imaging stained cells^{2–4}. However, these image-based technologies also have their limitations, as they typically focus

on biological processes with known associations or assumptions, thereby constraining the discovery in existing knowledge⁵. Moreover, traditional methods that include both high-dimensional assays and image-based technologies are often constrained by their complexity and high costs. To overcome these issues, the technology, known as Cell Painting (CP), has been proposed as a solution. Specifically, CP technology involves staining eight cellular components with six remarkably cheap and easy dyes and imaging them in five channels on a fluorescence microscope⁶, which is simple to operate and less costly⁷. Beyond ease of use, CP technology operates on a new paradigm by

¹Centre for Bioinformatics and Intelligent Medicine, College of Computer Science, Nankai University, Tianjin, China. ²School of Mathematical Sciences and LPMC, Nankai University, Tianjin, China. ³MOE Key Laboratory of Bioinformatics and Bioinformatics Division of BNRIST, Department of Automation, Tsinghua University, Beijing, China. ⁴These authors contributed equally: Chengwei Yan, Yu Zhang, Jiuxin Feng. ✉ e-mail: jianliu@nankai.edu.cn; chenshengquan@nankai.edu.cn

collecting data on a large scale without an initial focus on specific hypotheses or known knowledge to reveal unanticipated biology at play. With these advantages, CP technology has shed light on novel phenotypes and cellular phenotypic heterogeneity, providing a valuable complement to genomics⁸. It also has been successfully used to characterize genes^{9–11} and compounds^{2,11,12} in several steps of the drug discovery process.

As CP technology has developed and been increasingly applied, a large amount of data has been accumulated¹³. In 2023, scientists established the Joint Undertaking for Morphological Profiling (JUMP) dataset, a standardized collection of cell-stained images featuring over 116,000 unique compound perturbations and more than 15,000 unique genetic perturbations¹⁴. Despite the dataset's substantial analytical potential, technical effects caused by non-biological factors pose significant challenges in the analysis process. Some of these technical effects, resulting from variations across different laboratories, batches within laboratories, and even different microscopes, have been observed in the JUMP dataset^{15,16}. As a representative example, the open reading frame (ORF) overexpression dataset from cpg0016, used here as a unified data source, exhibits batch effects arising from technical variation across experiments¹⁴. In addition to these well-recognized technical effects observed in most data collection techniques, previous studies have found that CP features extracted using the conventional tool, CellProfiler¹⁷, exhibit distinctive well-position effects⁴. Concretely, well-position effects arise from the unique design of the CP experiment. In CP technology, experimental plates are typically organized into 16 rows and 24 columns, totaling 384 wells, with each well influenced by both row effects and column effects. We collectively refer to row effects, column effects, and the effects of different batches as triple effects. The complex and combined triple effects can lead to deviations from accurate biological profiles and thus need to be corrected urgently.

To the best of our knowledge, no methods have been specifically designed to correct triple effects, especially well-position effects in CP profiles. Although there exist some batch correction methods designed for single-cell data, directly applying them to CP data remains challenging for several reasons. First, the characteristics of CP data differ significantly from single-cell data, as CP data is denser and exhibits lower variability compared to single-cell data¹⁸. Second, well-position effects in CP data contrast with batch effects in single-cell data, as row or column effects show a gradient-influenced pattern, where greater differences in row or column numbers lead to more pronounced effects. Third, the triple effects, especially row and column effects, are complexly interactive and need to be corrected simultaneously. Some methods, such as scVI¹⁹, can correct only one type of technical effect and are constrained to correct multiple technical effects. Although methods like Harmony²⁰ model one type of technical effect at a time and can correct all effects one by one, they are unable to simultaneously model triple effects in CP data.

Besides the challenges of correcting triple effects, existing studies that rely on the current standard feature extraction pipeline using the CellProfiler still encounter several unresolved disputes and limitations⁴. First, while well-position effects are primarily introduced at the imaging level, it remains unclear whether different feature extraction paradigms, such as traditional pipelines or pre-trained deep learning models, capture these effects in a consistent manner, particularly those exhibiting gradient-influenced patterns. Second, although CellProfiler is a flexible platform that supports deep learning tools through its plugin ecosystem²¹, the current standard feature extraction pipeline using CellProfiler still predominantly relies on traditional computer vision features and requires expert selection during its feature extraction pipeline, which may overlook certain relevant phenotypic variation⁴. In contrast, pre-trained segmentation models extract features through a different paradigm, where spatial and morphological patterns are learned directly from large-scale image

data, rather than relying on manually engineered processing. This distinct extraction paradigm may capture phenotypic variations that are underrepresented in conventional pipelines, thereby providing complementary information for downstream analysis.

Here, we show a one-stop method named cpDistiller for correcting triple effects and extracting latent patterns in CP data. cpDistiller mainly comprises three modules: the extractor module for deriving more comprehensive image information, the joint training module for integrating dual-source features, and the technical correction module for simultaneously correcting batch, row, and column effects. Specifically, the extractor module, inspired by transfer learning, employs a pre-trained segmentation model in an end-to-end manner, which is adjusted to extract features from nearly 30 terabytes of raw images. The joint training module aligns the features extracted by both CellProfiler and the extractor module, improving the model's ability to better characterize cell-to-cell variation. The technical correction module employs a semi-supervised Gaussian mixture variational autoencoder (GMVAE), incorporating contrastive and domain-adversarial learning strategies, to simultaneously correct technical effects. Based on comprehensive experiments across various CP profiles, we demonstrate that cpDistiller excels in both qualitative visualizations and quantitative metrics, outperforming five baseline methods in single-batch well-position effect correction as well as simultaneous triple-effect correction, all while preserving biological heterogeneity. Besides, we showcase the extensive capabilities of cpDistiller, including the ability to integrate more information-rich image features, the support for incremental learning, and the robustness to various feature selection strategies. Moreover, we emphasize that cpDistiller effectively captures system-level phenotypic responses to genetic and chemical perturbations, serving as a powerful tool to complement single-cell RNA sequencing (scRNA-seq) data for uncovering gene functions and relationships. In addition to the combination with scRNA-seq, cpDistiller has the potential to provide unbiased insights into gene associations independently. Furthermore, by improving the matching of genetic perturbations with their target genes and enhancing gene-compound similarity assessments, cpDistiller shows strong promise for accelerating the identification of targets, which is quite valuable in facilitating drug discovery and various fields of biological research.

Results

The overall architecture of cpDistiller

cpDistiller maps the input data to the low-dimensional embedding space that aims to correct triple effects while capturing true biological signals. Specifically, cpDistiller is composed of three main modules: the extractor module, the joint training module, and the technical correction module (Fig. 1).

cpDistiller processes CP images at a resolution of 1080 × 1080 pixels and extracts comprehensive features from each well. These features form a matrix, where rows represent samples (cells, wells, or perturbations) and columns correspond to extracted features. We first extract features from the raw images using CellProfiler, a widely used approach, and refer to these features as CellProfiler-based features. Drawing inspiration from transfer learning, we further develop the extractor module based on an end-to-end pre-trained segmentation model to automatically extract features, and refer to these features as cpDistiller-extractor-based features (“Methods”).

Then, the two sets of features are introduced into the joint training module for integration. The CellProfiler-based features retain unprocessed, while the cpDistiller-extractor-based features are transformed through an attention mechanism-based encoder, reducing them to a latent space, and then reconstructed via a decoder. This encoder-decoder structure, applied exclusively to the cpDistiller-extractor-based features, ensures feature refinement by reducing noise and enhancing the quality of the representations. The attention

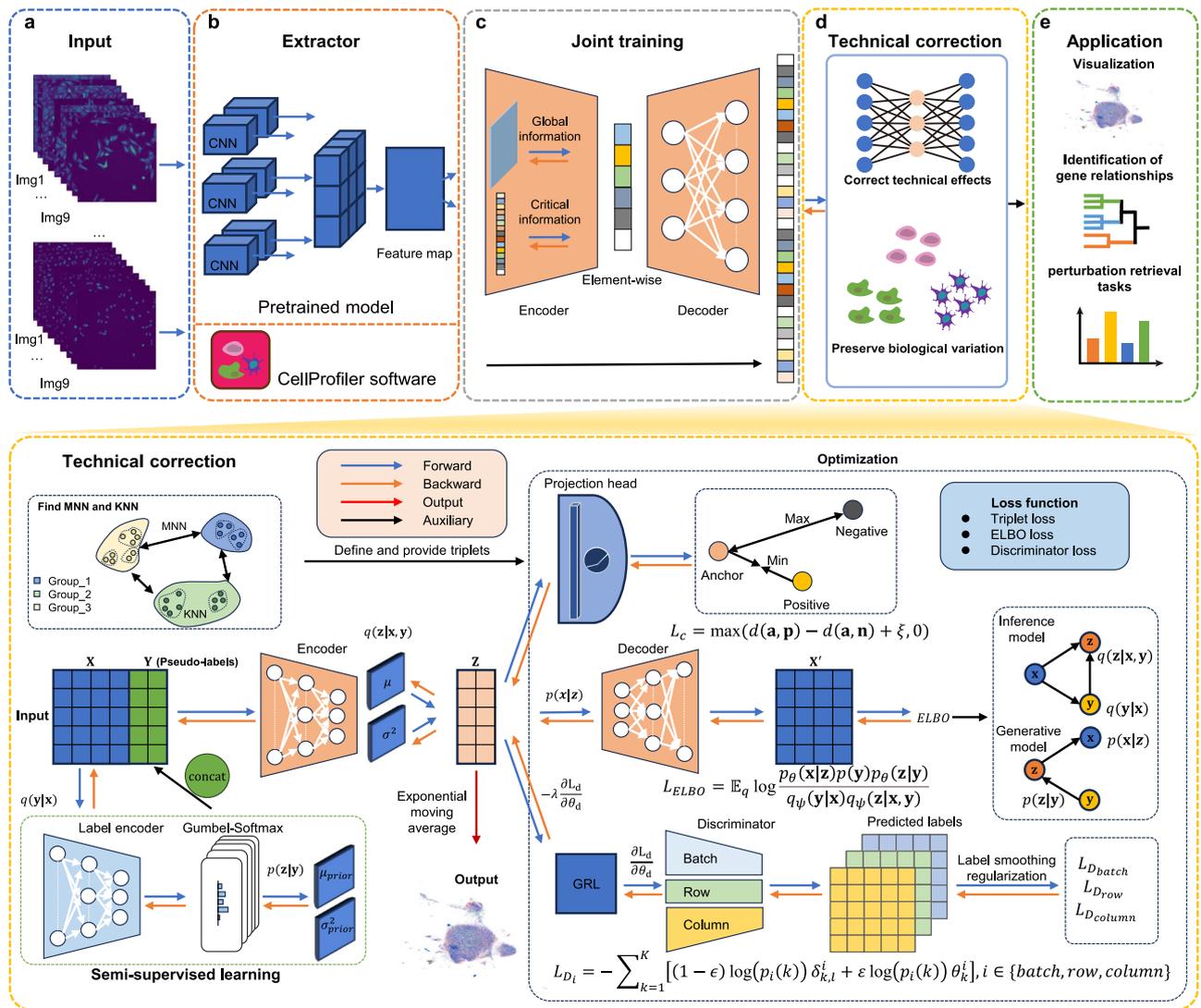


Fig. 1 | Overview of cpDistiller. **a**, **b** cpDistiller takes CP images as input (**a**) and extracts features by CellProfiler and the extractor module (**b**). **c** The joint training module integrates features from both the cpDistiller-extractor module and CellProfiler. The CellProfiler-based features retain unprocessed, while the cpDistiller-extractor-based features undergo the encoder-decoder architecture with an attention mechanism. These refined features are then input into the technical correction module. **d** The technical correction module, built on a GMVAE, infers pseudo-labels and obtains low-dimensional representations. In the low-dimensional space derived from the GMVAE, representations pass through a

projection head before applying the triplet loss to restore more accurate nearest-neighbor relationships. In addition, three discriminators are designed to predict batch, row, and column labels, and the technical correction module uses the GRL to enable adversarial learning to correct triple effects. Moreover, the exponential moving average (EMA) is applied during training, and the averaged parameters are used as the final model parameters. **e** The low-dimensional embeddings obtained by cpDistiller facilitate a range of applications, including data visualization, identification of gene relationships, and various perturbation retrieval tasks.

mechanism emphasizes key features across dimensions, facilitating a more effective combination with the CellProfiler-based features while filtering out irrelevant information (“Methods”).

The technical correction module, with a GMVAE at its core, infers pseudo-labels for each well in a semi-supervised mode using the features integrated by the joint training module. These pseudo-labels represent the Gaussian components that characterize the underlying patterns of each well, capturing the differences among these patterns. By combining these pseudo-labels along with the integrated features, the technical correction module derives the latent low-dimensional representations for each well. We also employ both contrastive and domain-adversarial learning strategies. In contrastive learning, we use *k*-nearest neighbors (KNN) and mutual nearest neighbors (MNN) to construct triplets, applying triplet loss²² to restore more accurate nearest-neighbor relationships for each well. To

facilitate domain-adversarial learning and avoid stage-wise training in generative adversarial networks (GANs)²³, we further apply the gradient reversal layer (GRL)²⁴ in an end-to-end training approach, making it harder for discriminators to distinguish data sources and thus removing technical effects. In addition, for calculating the discriminator loss, we design a soft label distribution tailored to the gradient-influenced pattern of the CP data, where each entry in the distribution vector reflects its proximity to the actual data source for a more accurate representation.

To the best of our knowledge, no specialized methods are available for correcting triple effects, especially well-position effects in CP data. However, a recent benchmark study has explored the application of single-cell batch correction methods as a potential solution¹⁵. In light of this, we compared cpDistiller with methods that excel at removing batch effects in single-cell data, including Seurat v5²⁵, Harmony²⁰, Scanorama²⁶,

scVI¹⁹, and scDML²⁷. Among them, Seurat v5²⁵, Harmony²⁰, and Scanorama²⁶ can iteratively correct triple effects by sequentially incorporating different technical labels (removal of batch, row, and column effects one after another), as they are based on low-dimensional representations, which allow for flexible, stepwise corrections. In contrast, scVI¹⁹ and scDML²⁷ directly model the original high-dimensional inputs to obtain low-dimensional representations, making them less suitable for iterative correction and limiting them to address only one type of technical effect (one of batch, row, or column effects). Although established metrics for evaluating technical correction in CP profiles were previously limited, recent studies have introduced well-motivated and biologically relevant evaluation criteria tailored to CP profiling data²⁸. Therefore, we adopted both these profiling-specific metrics and widely used single-cell analysis metrics to assess the effectiveness of different methods in correcting technical effects while preserving biological variation. Metrics such as average silhouette width (ASW)²⁹, technic average silhouette width (tASW)³⁰ and graph connectivity³⁰ are used to measure the model's ability to remove technical effects, while perturbation average silhouette width (pASW)³⁰, normalized mutual information (NMI)³⁰, phenotypic activity²⁸ and phenotypic consistency²⁸ are used to evaluate the characterization of heterogeneity ("Methods"). Besides, although our study focuses on triple effects, we also evaluate the models' performance in mitigating plate effects, which represent a well-recognized source of systematic variation in CP assays³¹, to enhance the robustness of our assessment ("Methods").

The features derived from CellProfiler and cpDistiller-extractor module confirm triple effects

Previous studies have identified three types of technical effects in CP data: batch, row, and column effects³¹. Notably, row and column effects suggest the special well-position effects. Following the prior study¹⁴, we employed CellProfiler to obtain 1,446-dimensional features (CellProfiler-based features) from the open reading frame (ORF) overexpression dataset in cpG0016, and then performed uniform manifold approximation and projection (UMAP)³² for these features. From the UMAP visualization, we observed distinct clustering patterns corresponding to different batches, rows, and columns (Fig. 2a). Batch effects are evident as patterns from different batches cluster into separate groups, while row and column effects are displayed as color gradients in the visualization, highlighting the presence of triple effects (Fig. 2a). To investigate whether row and column effects are evident within a single batch, we further observed the CellProfiler-based features from Batch_7. The UMAP visualizations revealed color gradients, indicating a gradient-influenced pattern where more distant rows or columns exhibit more pronounced effects, a phenomenon consistently observed across 12 batches (Fig. 2a and Supplementary Figs. 1–3).

To investigate whether different feature extraction paradigms capture triple effects in a consistent manner, especially unique well-position effects, we developed the cpDistiller-extractor module, based on the pre-trained segmentation model Mesmer³³, to extract deep-learning features (cpDistiller-extractor-based features). As illustrated in Fig. 2a, UMAP visualizations of the cpDistiller-extractor-based features from 12 batches display slight batch effects and clear, consistent effects related to row and column (Fig. 2a). We further plotted and observed the density distributions of the scatter points along the x and y axes in the UMAP visualizations. These distributions revealed distinct peaks corresponding to different batches, indicating the presence of batch effects (Fig. 2a). In addition, we observed that the distributions for rows and columns were non-overlapping, with row-specific and column-specific peaks highlighted in red boxes, suggesting the existence of row and column effects (Fig. 2a). Furthermore, UMAP visualizations of the cpDistiller-extractor-based features from Batch_7 demonstrated clear evidence of row and column effects (Fig. 2b), a pattern that was consistently observed across all batches

(Supplementary Figs. 4 and 5). These results indicate that cpDistiller-extractor-based features are also capable of capturing well-position and batch effects. Compared to CellProfiler-based features, the batch effects appeared less pronounced, while row and column effects were still consistently observed. Notably, these deep learning-derived features did not exhibit the typical gradient-influenced patterns observed with CellProfiler-based features. Instead, the well position effects manifested in a more diffuse and less spatially structured form.

To investigate whether the cpDistiller-extractor module can capture relevant information, such as the cell nucleus from CP images, we compared its segmentation results with those obtained from CellProfiler. As illustrated in Fig. 2c, the Mesmer model, which serves as the core of the cpDistiller-extractor module, produced reasonable segmentation results for cell nuclei that were broadly consistent with CellProfiler software (Supplementary Note 1 and Supplementary Fig. 6). Before adopting Mesmer, we also assessed the widely used LACSS³⁴ and YOLOv8³⁵ models as alternatives. However, neither was able to accurately locate cell nucleus (Fig. 2c and Supplementary Fig. 7). We further tested the YOLOv8 model with various parameters, but the results confirmed that it remained ineffective for cell nucleus detection and segmentation in CP data (Supplementary Fig. 7). Details regarding the pre-trained model settings can be found in Supplementary Note 2.

In conclusion, features from both CellProfiler and the cpDistiller-extractor module captured information about the cell nucleus in CP images and revealed the inherent batch, row, and column effects of CP data.

cpDistiller effectively corrects well-position effects and preserves cellular phenotypic heterogeneity

We first conducted experiments on the cpG0016 ORF profiles derived from the U2OS cell type in the JUMP dataset to demonstrate the effectiveness of cpDistiller. The cpG0016 ORF profiles consisted of 12 batches, and we used Batch_1 as an example. In the visualization of Baseline (defined in "Methods" as CellProfiler-based features after standard pre-processing without any correction), UMAP visualization reveals a gradient-influenced pattern (highlighted with a red box), with more distant rows or columns exhibiting more pronounced effects (Fig. 3a, b and results from other methods and batches in Supplementary Figs. 8–19). We then visualized the low-dimensional embeddings obtained by cpDistiller and observed that the gradient-influenced pattern was significantly reduced. The data distribution across different rows was more uniform, and the previously pronounced column effects, particularly the large differences between low-numbered columns (e.g., Column_1) and high-numbered columns (e.g., Column_24), were well-mixed and greatly corrected, as highlighted by the red boxes in Fig. 3a, b. These demonstrated that cpDistiller can effectively correct both row and column effects. Importantly, edge effects are commonly observed in CP experiments and often appear as similar profiles among wells located at the outermost rows and columns, such as Row_A and Row_P or Column_1 and Column_24³¹. However, the results suggested that the observed improvements are not attributable to edge effects, which were minimal in the ORF dataset, but rather reflect cpDistiller's ability to correct the well-position effects beyond those related to outermost wells (Supplementary Note 3 and Supplementary Figs. 20–22). It is noted that cpDistiller reconstructs local neighborhoods based on feature similarity, ensuring that wells are brought closer only when supported by underlying biological resemblance ("Methods"). As highlighted by the blue boxes in Fig. 3b, this integration involves wells from various positions across the columns, rather than being limited to edge wells. Moreover, overexpression reagents and compounds are theoretically expected to induce significant differences in cell phenotypes between negative and positive controls. While simple cell counts can reflect certain phenotypic changes, such as cytotoxicity or cell death^{36,37}, they are insufficient to capture the full spectrum of morphological

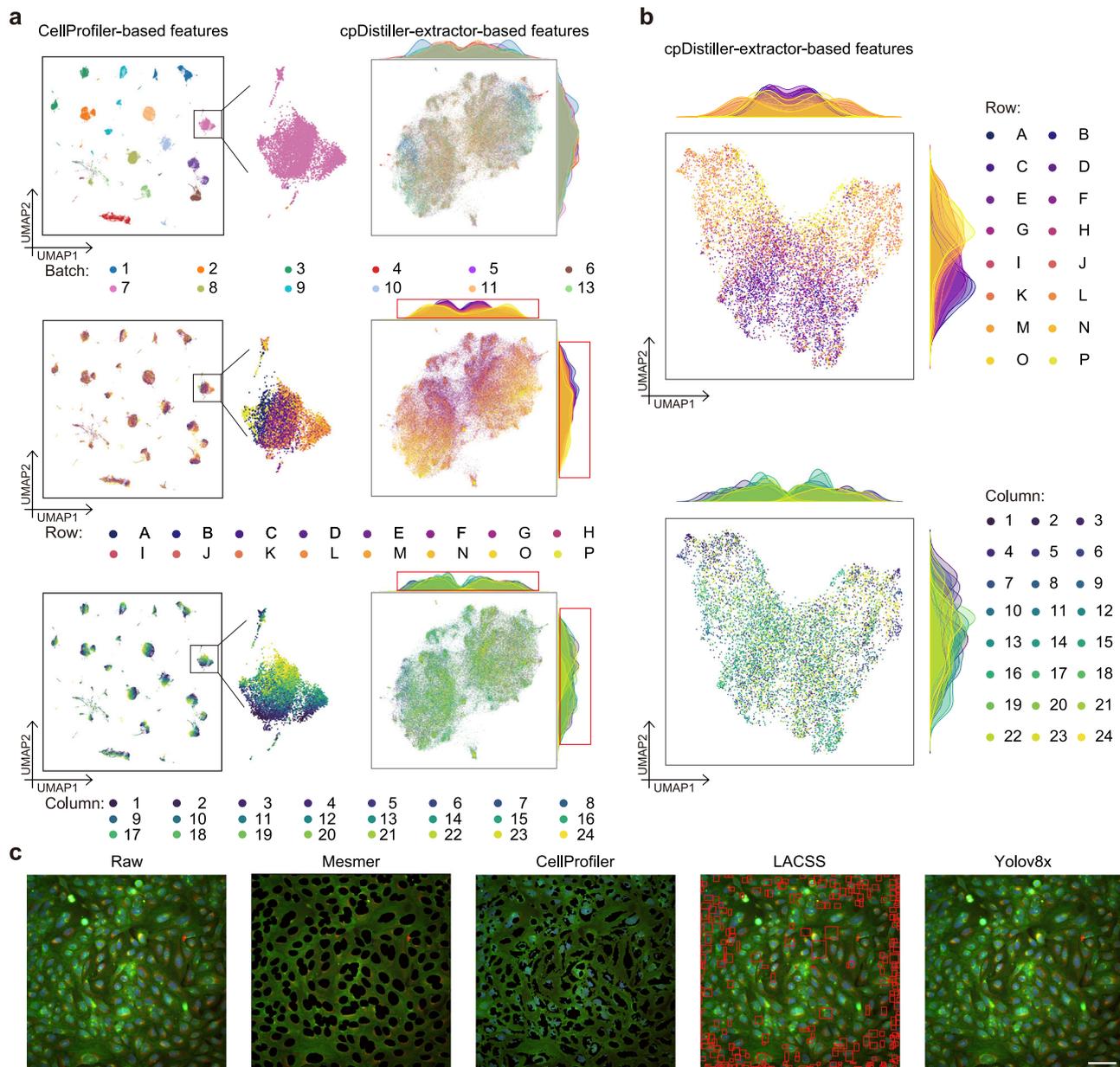


Fig. 2 | CellProfiler and cpDistiller-extractor module confirm triple effects. **a** The UMAP visualizations of CellProfiler-based features (left) and cpDistiller-extractor-based features (right), colored by batch, row, and column, respectively. **b** The UMAP visualizations of embeddings obtained by the cpDistiller-extractor module in Batch_7, colored by row and column, respectively. **c** Images show

segmentation or detection results of different models. The raw image is derived from one of the nine sub-images in a well from the CP images. Segmentation results are shown as black areas for Mesmer and CellProfiler, while detection results in the LACSS are highlighted with red boxes. YOLOv8 fails to produce reliable cell nucleus segmentation results from CP images. The scale bar represents 100 μm .

variation (Supplementary Note 4 and Supplementary Figs. 23–26). By utilizing UMAP to visualize the impact of various perturbations, cpDistiller successfully differentiated positive controls, including JCP2022_037716 (AMG-900), JCP2022_035095 (LY2109761), and JCP2022_012818 (C23H17Cl2N5O4) from negative controls across various batches, highlighting its ability to extract biologically meaningful representations beyond cell count alone (Supplementary Figs. 8–19 and 24). In addition to analyzing the controls, we also explored the ORF perturbations in treatment. cpDistiller successfully identified perturbations caused by 12 reagents in treatment and preserved their unique stimulatory effects, which were obscured in results of Baseline due to well-position effects (Fig. 3c and Supplementary Fig. 27). In contrast, UMAP visualizations and hierarchical clustering results suggest that methods like scDML²⁷ and scVI¹⁹, which are limited

to correcting only one type of technical effects, struggle to preserve biological variation while also failing to effectively correct well-position effects (Supplementary Figs. 8–19 and 27). Although these methods may achieve partial mixing across rows and columns, they still fail to remove most non-biological noise, such as the clear striping pattern, particularly noticeable in Batch_5 and Batch_8 (Supplementary Figs. 12 and 15). While methods like Harmony²⁰, Seurat v5²⁵, and Scanorama²⁶ can iteratively correct well-position effects, they also fall short in preserving the true biological variation of ORF perturbations (Supplementary Fig. 27).

It's noted that visual analysis and quantitative metrics do not always align, and both perspectives are essential for a comprehensive evaluation of model performance³⁸. To quantitatively demonstrate the advantages of cpDistiller in correcting well-position effects, we further conducted experiments in 12 batches, evaluating each batch

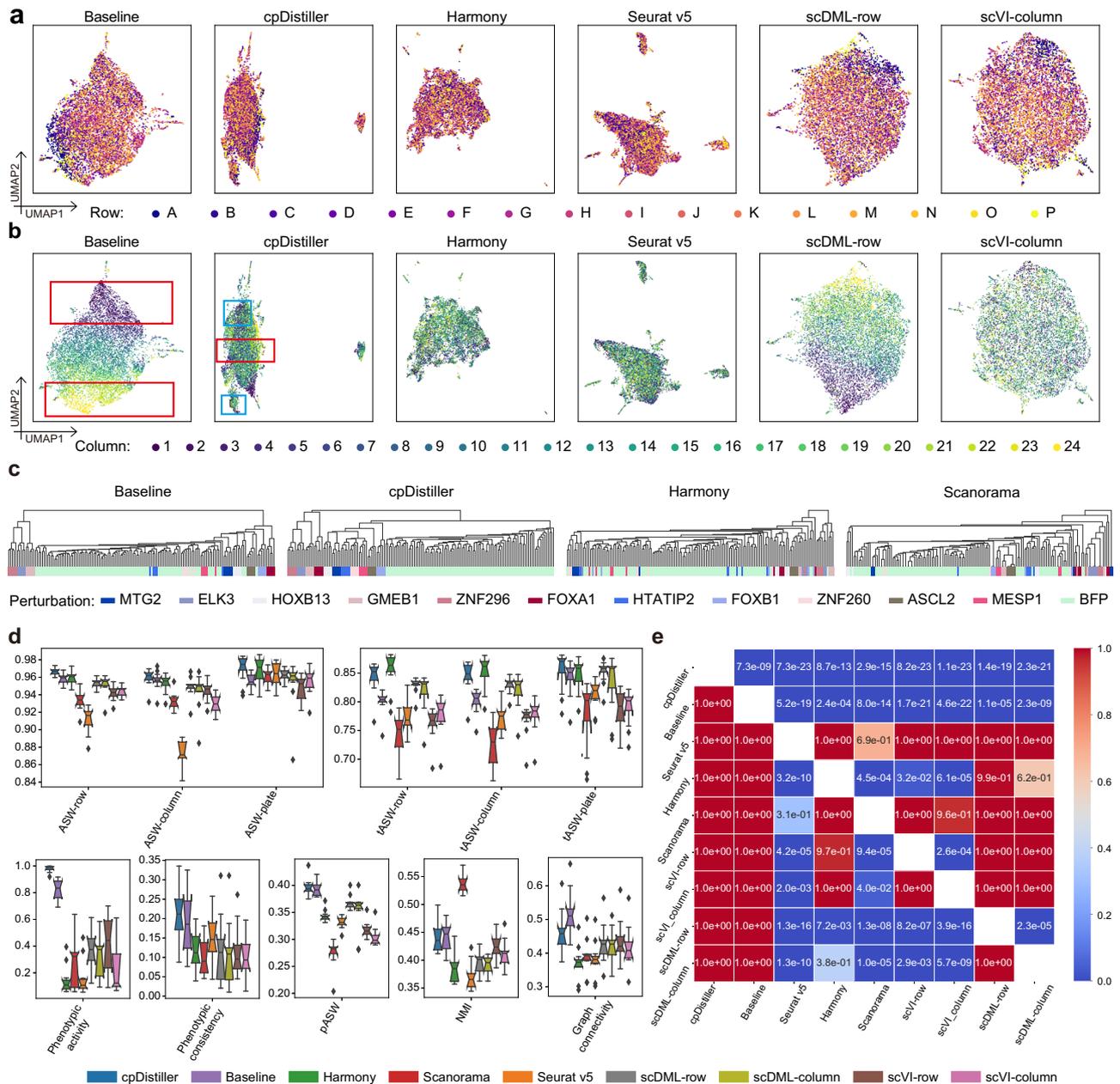


Fig. 3 | cpDistiller can effectively correct well-position effects while preserving biological variation for ORF profiles in 12 batches. **a, b** UMAP visualizations of embeddings obtained by different methods in Batch_1 colored by row (**a**) and column (**b**). Suffixes “-row” and “-column” on scDML and scVI indicate the corrected effects. **c** Dendrograms illustrate the clustering of ORF perturbations induced by 12 reagents in treatment, graphically rendered based on low-dimensional representations generated by different methods. Baseline represents the preprocessed but uncorrected CellProfiler-based features. **d** Quantitative evaluation of the performance of different methods on ORF profiles across 12 batches, where each batch is analyzed independently to assess the correction of well-position effects while preserving biological variation. The ASW and tASW metrics yield three results

based on row labels, column labels, and plate labels, respectively. In the boxplots, center lines indicate the medians, box limits show upper and lower quartiles, whiskers represent the 1.5 × interquartile range, and notches reflect 95% confidence intervals via Gaussian-based asymptotic approximation. Each data point represents one biologically independent batch ($n = 12$), with batches differing in their experimental context and typically containing distinct gene perturbations. **e** Heatmap shows p -values from one-sided paired Wilcoxon signed-rank tests. Each cell in the heatmap reflects the statistical significance of one method’s (row) superiority over another (column), derived from 132 evaluations across 12 batches using 11 metrics ($n = 132$).

independently without cross-batch integration. Performance was then assessed by ASW, tASW, and graph connectivity as suggested in refs. 29,30. ASW and tASW were used to assess the extent of data mixing across different technical labels. Graph connectivity assumed that, once technical effects were corrected, data with the same biological labels, specifically perturbation labels, should cluster more tightly. cpDistiller achieved higher scores than baseline methods in most evaluation criteria and maintained robust performance across batches

(Fig. 3d), indicating superior performance in both uniformly mixing CP profiles and accurately reflecting cell phenotype differences induced by various perturbations.

We further assessed the preservation of biological variation via metrics such as phenotypic activity, phenotypic consistency, pASW, and NMI. Phenotypic activity evaluates whether replicates of a perturbation produce consistent morphological changes that are distinguishable from negative controls, while phenotypic consistency

assesses whether perturbations with shared biological functions yield comparable phenotypic profiles. cpDistiller excelled in phenotypic activity and phenotypic consistency, particularly achieving the highest scores among all baseline methods, indicating its ability to effectively detect perturbation effects and retain biologically meaningful relationships across perturbations (Fig. 3d). In addition, for the clustering metrics of pASW and NMI, which were used to assess the clustering quality of replicated experiments involving the same types of perturbations as suggested in ref. 30, cpDistiller also showed robust and competitive performance compared to other methods (Fig. 3d). Harmony²⁰ has demonstrated excellent performance in removing batch effects for CP data in recent benchmark analysis¹⁵ and also emerges as the most effective baseline method for correcting row and column effects in our study. Nevertheless, it does not perform well in preserving biological heterogeneity, which limits its utility in downstream biological interpretation. Interestingly, Baseline also showed high biological scores. This observation is consistent with prior studies in omics research, which have shown that uncorrected data can sometimes preserve biological heterogeneity^{30,39–41}. Nevertheless, there is broad consensus that proper correction is necessary to reduce technical effects while retaining meaningful signals^{30,39–41}. In particular, we noted that the high scores of Baseline may partially result from shared well-positions across plates, where replicate samples may appear more similar due to shared technical effects, potentially inflating biological consistency scores, a concern also raised in refs. 42,43. These findings highlight the importance of using appropriate correction methods to disentangle technical variation from true biological signals. To further demonstrate the significant advantages of cpDistiller, we used one-sided paired Wilcoxon signed-rank tests to compare the performance of different methods across 12 batches. The *p*-values highlighted cpDistiller's significant advantages in removing well-position effects and preserving biological variation across batches, outperforming all baseline methods (Fig. 3e).

Besides, in our previous experiments with baseline methods like Seurat v5²⁵, Harmony²⁰, and Scanorama²⁶, which can iteratively remove well-position effects, we initially corrected rows before columns. To investigate whether the performance of these methods was influenced by the order of correction, we also applied the reverse order, namely correcting columns before rows. The box plots across 12 batches showed consistent difference in performance when the same method was applied with different correction orders (Supplementary Fig. 28a). However, when we used one-sided paired Wilcoxon signed-rank tests for a more detailed quantification of performance differences, we observed that the differences were statistically significant for Harmony and Scanorama, indicating that correcting rows before columns generally yielded better results than the reverse order. In contrast, Seurat v5 did not exhibit a consistent trend between the two correction orders (Supplementary Fig. 28b).

In summary, cpDistiller demonstrated superior performance on 12 batches of ORF data spanning diverse CP profiles, achieving a satisfactory balance between correcting well-position effects and preserving cellular phenotypic heterogeneity.

cpDistiller enables effective and simultaneous correction of triple effects

After integrating ORF profiles from 12 batches and visualizing them using UMAP, we observed batch effects in addition to well-position effects (Fig. 2a). Although these batch effects became less prominent after applying the median absolute deviation (MAD) normalization (“Methods”), a commonly used preprocessing step in CP studies, we further corrected them to ensure reliable cross-batch comparisons, which is essential for studying biological-process-related mechanisms of action based on feature similarity⁴⁴. For competently removing batch effects, in addition to the discriminators for row and column labels, we also created an additional discriminator to identify batch

labels, encouraging cpDistiller to learn low-dimensional representations that are indistinguishable with respect to the sources of batch labels, thereby removing batch effects (“Methods”). At the same time, we additionally considered KNN intra batches and MNN inter batches to construct triplets, using the contrastive learning technique to restore more accurate nearest-neighbor relationships for each well (“Methods”). We next verified that cpDistiller can satisfactorily correct triple effects.

First, we qualitatively compared the performance of different methods using UMAP visualizations, which showed that cpDistiller, Seurat v5, and Harmony achieved successful mixing of data across batches, rows, and columns simultaneously (Fig. 4a and Supplementary Fig. 26). Moreover, to qualitatively assess how well different methods preserve the specificity of perturbations, compound perturbations shared across 12 batches on the target plates provided a reliable measure¹⁵. UMAP visualizations showed that cpDistiller ensured perturbations caused by the same compounds had a tendency to be clustered (Fig. 4a). To be specific, cpDistiller successfully identified multiple perturbations in target plates, including JCP2022_010404 (SB-203580), JCP2022_000794 (NVP-AEW541), and JCP2022_047545 (SPEBRUTINIB). Although compounds such as SPEBRUTINIB are known to be associated with cytotoxic effects, changes in cell count alone cannot fully capture the morphological variations (Supplementary Fig. 26). The results indicated that cpDistiller's ability to capture richer morphological information allowed it to better differentiate perturbations beyond cell count (Supplementary Note 4 and Supplementary Figs. 23–26). In contrast, methods that are limited to correcting only one type of triple effects, such as scDML²⁷ and scVI¹⁹, failed to correct well-position effects, resulting in the persistence of striping pattern for rows and columns (Fig. 4a). On the other hand, methods capable of iteratively correcting triple effects, including Scanorama²⁶ and Seurat v5²⁵, tended to overcorrect, leading to fail to preserve true biological signals (Supplementary Fig. 26). Given the known limitations of UMAP, such as its sensitivity to initialization and parameter choices, we also explored the use of principal component analysis (PCA) visualizations as a complementary approach to assess the performance. The results indicated that while PCA provided a more straightforward and interpretable layout of major variance directions, its explained variance ratio was relatively low, and it failed to clearly reveal gradient-influenced well-position effects or distinguish perturbations as effectively as UMAP (Supplementary Note 5 and Supplementary Figs. 29, 30).

Moreover, we quantitatively evaluated the performance of different methods. cpDistiller consistently ranked highest for technical correction and biological conservation in overall metrics, demonstrating balanced effectiveness in both aspects (Fig. 4b and Supplementary Table 1). Although Scanorama achieved a high NMI score, this metric tends to favor methods that produce a larger number of clusters⁴⁵, even without well-defined biological separation, which may not reflect true biological structure in datasets with numerous and imbalanced perturbation labels. Besides, we found that all methods yielded low phenotypic consistency scores, which likely stems from our evaluation setup (Supplementary Note 6). To further validate cpDistiller's ability to recognize perturbations with biologically similar mechanisms, we calculated phenotypic consistency scores for the active genes and compounds identified by other methods. The results showed that cpDistiller consistently produced higher phenotypic consistency scores for all methods (Fig. 4c). This suggests that cpDistiller is particularly effective at identifying perturbations with similar biological mechanisms, reinforcing its ability to capture biologically meaningful variation in the CP data.

Furthermore, to comprehensively evaluate the advantages of cpDistiller, we conducted a series of experiments to demonstrate its utility in multiple conditions, maintain robustness to feature selection, and support incremental learning for ongoing studies. To investigate the utility of cpDistiller, we conducted systematic experiments to

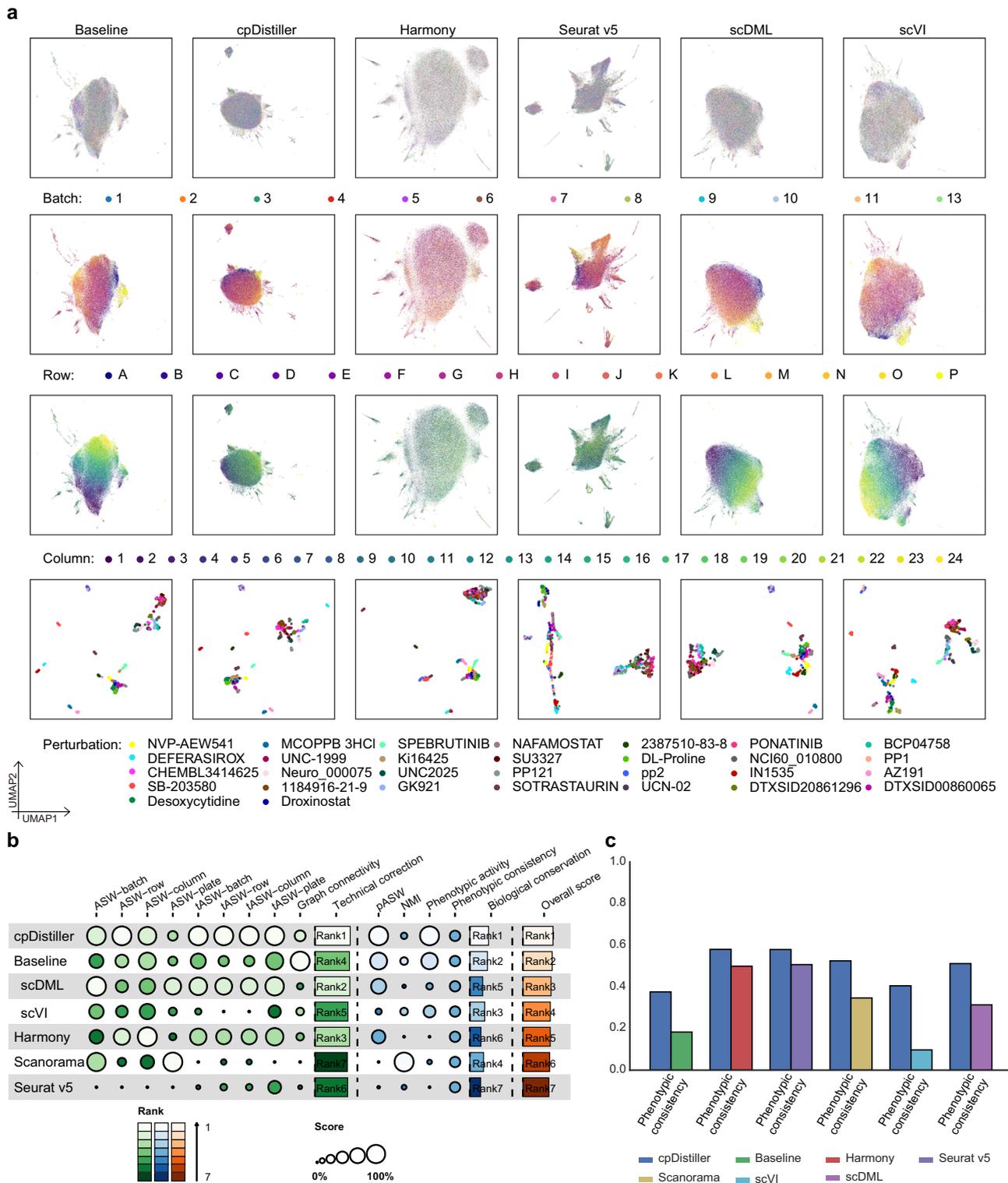


Fig. 4 | cpDistiller can simultaneously correct triple effects while preserving cellular phenotypic heterogeneity. **a** UMAP visualizations of embeddings generated by different methods using the combined ORF profiles from 12 batches, colored by batch, row, column, and perturbation, respectively. The embeddings are consistent with all visualizations, and the selected 30 representative perturbations are shown. Baseline represents the preprocessed but uncorrected CellProfiler-based features. **b** Overview of benchmarking outcomes for different methods. Technical correction and biological conservation scores refer to the average

performance in these two aspects, whereas the overall score represents the aggregate performance across all metrics for different methods. The ASW and tASW metrics yield four results based on batch labels, row labels, column labels, and plate labels, respectively. **c** Bar plots show phenotypic consistency scores after triple-effect correction. Each pair of bars represents scores computed on active perturbation subsets identified by a given method, with one bar for that method and the other for cpDistiller applied to the same subsets.

evaluate the performance of cpDistiller under three different settings: using both cpDistiller-extractor-based features and CellProfiler-based features, using only CellProfiler-based features, and using only cpDistiller-extractor-based features. The results indicate that while cpDistiller-extractor-based features enhance overall performance, using CellProfiler-based features also yields strong performance, making cpDistiller flexible and applicable if only CellProfiler-based features are available (Supplementary Note 7 and Supplementary Fig. 31).

In addition, cpDistiller demonstrates robustness to feature selection (Supplementary Note 8). We first extracted features with dimensions of 4752 from fluorescence images and 7638 when combined with brightfield images by using the CellProfiler software (“Methods”). We then utilized these two types of features to carry out two key tasks: correcting well-position effects in a single batch and simultaneously correcting triple effects. The results showed that even when dealing with high-dimensional features that may contain substantial redundant features and noise, cpDistiller consistently outperformed baseline methods in correcting technical effects and preserving biological variation (Supplementary Fig. 32a, b and Supplementary Tables 2, 3).

Besides, cpDistiller supports incremental learning (Supplementary Note 9). Methods like Harmony²⁰ and Scanorama²⁶ require reprocessing and realigning the entire dataset to integrate new data, which can be cumbersome, especially when working with large public datasets¹⁵. These processes often involve modifying existing representations, which can disrupt the continuity of prior analytical results. In contrast, cpDistiller can leverage the model parameters learned from previous tasks, enabling the direct integration of new data without the necessity of realignment (Supplementary Fig. 32c). These flexibilities are particularly beneficial for ongoing studies, enabling seamless integration while preserving the integrity of previous analyses.

Overall, cpDistiller achieved a satisfactory balance between correcting triple effects and preserving biological variation and demonstrated strong capabilities in the utility, incremental learning, and robust feature selection.

cpDistiller can combine scRNA-seq data to reveal gene functions and relationships

Inferring gene functions and interactions, as a fundamental step in many areas of biological research, often relies on various types of sequencing data, such as scRNA-seq data and chromatin immunoprecipitation sequencing (ChIP-seq) data⁴⁶. This inferring task is associated with numerous unknown and highly intricate biological processes, but the sequencing methods tend to focus on only a limited number of interesting molecular-level changes, making obtaining thorough and accurate inferences challenging⁴⁷. In contrast, we demonstrated that cpDistiller offers comprehensive system-level phenotypic characteristics under genetic perturbations and has the potential to integrate molecular-level RNA data for revealing gene functions and relationships.

To demonstrate the capability of the embeddings learned by cpDistiller in preserving system-level phenotypic characteristics of perturbation heterogeneity, we applied cpDistiller to the controls in the ORF data from the cpg0016 dataset. Specifically, we focus on controls, including the positive controls JCP2022_037716 (AMG-900), JCP2022_012818 (C23H17Cl2N5O4), and JCP2022_035095 (LY2109761), and the negative control JCP2022_915131 (*LacZ*). The positive controls are expected to produce noticeable phenotypic changes, while the negative control should induce minimal changes¹⁴. Theoretically, phenotypes under identical genetic perturbation in different wells or different plates should yield consistent patterns. Therefore, we evaluated the similarity of identical positive and negative controls duplicated across different wells and plates within a batch, using

embeddings generated by cpDistiller and baseline methods, respectively. Taking Batch_4, Batch_5, and Batch_6 for example, hierarchical clustering showed that cpDistiller outperformed baseline methods, including cell count-based clustering, by successfully grouping positive and negative controls and revealing treatment-specific embeddings (Fig. 5a and Supplementary Fig. 25a–c). These results illustrated that cpDistiller can effectively capture phenotypic signals across various perturbations.

To elucidate cpDistiller’s capability to integrate with scRNA-seq data for inferring gene functions, we used the embeddings obtained by cpDistiller to analyze genetic perturbation data from the ORF dataset. Given that the ORF dataset encompasses large-scale perturbations targeting 12,602 genes, we first employed the ARCHS4 tool⁴⁸ to establish gene groups based on the similarities of their scRNA-seq data, where genes with highly similar expression patterns were grouped together (“Methods”). Next, we calculated the Euclidean distance between the embeddings from cpDistiller for the individual genes in gene groups and then performed hierarchical clustering based on these distances. For example, we found gene Group_A and gene Group_B are highly separated into distinct clusters (Fig. 5b). Gene Group_A is enriched in Cluster_A, while gene Group_B is enriched in Cluster_B, indicating that cpDistiller has learned group-specific embeddings, resulting in significantly distinct morphological embeddings for genes in Group_A and Group_B. To elucidate that the morphological embeddings from cpDistiller have the capacity to reveal gene functions, we conducted cellular component (CC) analyses via gene ontology (GO) enrichment for the genes in Group_A and Group_B, respectively (“Methods”). We found that the top significantly enriched cellular components for the genes in Group_A were associated with cell–cell junctions and ubiquitin-mediated protein regulation, including the ‘Cul3-RING ubiquitin ligase complex’, ‘Tight junction’, and ‘Apical junction complex’ (Fig. 5b). These components are functionally linked to maintaining epithelial polarity, cell adhesion, and controlled protein turnover^{49–51}. In contrast, the top significantly enriched cellular components for the genes in Group_B were all involved in mitochondria-related components, including the ‘Mitochondrial inner membrane’, ‘Mitochondrial membrane’, and ‘Mitochondrial intermembrane space’, indicating their potential involvement in mitochondrial organization and energy metabolism^{52,53}. The CC analysis results revealed that the genes in Group_A and the genes in Group_B are associated with distinct functions, which is consistent with the cluster categorization, where genes in Group_A and Group_B have distinct embeddings learned by cpDistiller. Notably, while cpDistiller consistently recovered transcriptomics-defined gene groups as distinct clusters, such a structure was not observed in the embedding of Baseline, where different gene groups often appeared intermixed (Supplementary Fig. 33a). The results demonstrate that cpDistiller not only retains but in fact enhances the biological signals present in CellProfiler-based features, enabling more precise separation of gene groups with distinct functions.

To demonstrate cpDistiller’s effectiveness in elucidating gene relationships when combined with scRNA-seq data, we further analyzed the resulting gene embeddings in gene groups through multiple types of biological analyses. To enable a direct comparison, we conducted hierarchical clustering analysis for individual genes in gene groups using embeddings obtained by both Baseline and cpDistiller and found that they revealed different clustering patterns for several gene groups. As shown in Supplementary Fig. 33b, genes from three groups are intermixed and do not form distinct clusters for Baseline. In contrast, embeddings learned by cpDistiller exhibit a different structure. Gene Group_A and Group_B, as well as gene Group_A and Group_C, are distinguishable into different clusters (Fig. 5c and Supplementary Fig. 33b). However, gene Groups_B and Group_C are not easily separated. These results indicated that although genes within different groups have dissimilar scRNA-seq expression patterns, they

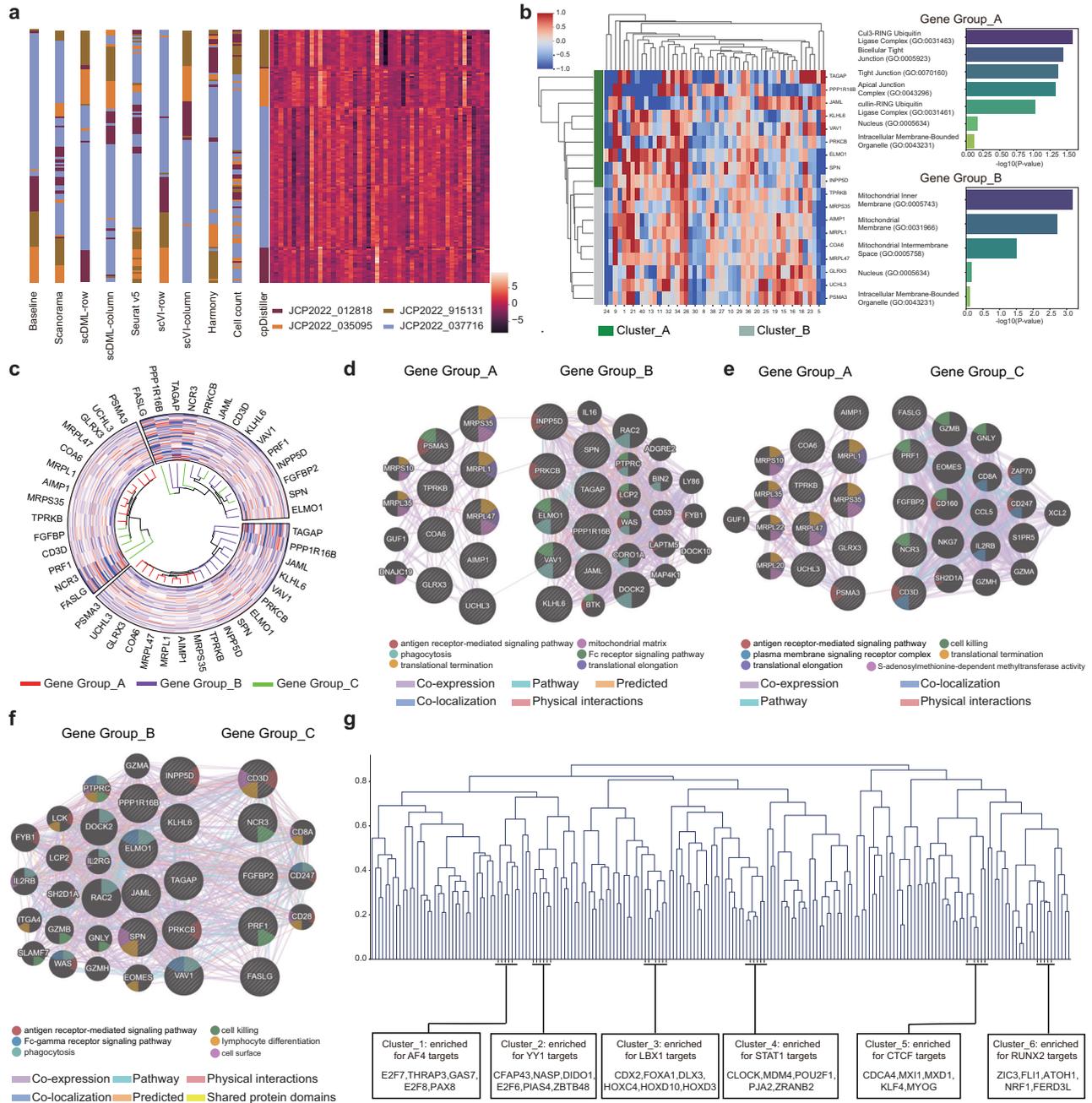


Fig. 5 | Analyses of gene functions and gene relationships. **a** Dendrograms illustrate the clustering of duplicates of controls, graphically rendered based on low-dimensional representations generated by different methods. Suffixes “-row” and “-column” on scDML and scVI indicate the corrected effects. The additional heatmap visualizes the 50-dimensional embeddings obtained by cpDistiller, highlighting the clustering patterns of controls. Baseline represents the preprocessed but uncorrected CellProfiler-based features. **b** Hierarchical clustering and heatmap of the 50-dimensional embeddings obtained by cpDistiller for gene Group A and gene Group B. Detailed views show the cellular compound analysis of GO term analysis for gene Group A and gene Group B. The *p*-values in the GO term analysis

are derived from one-sided Fisher exact tests, as computed by the Enrichr platform. These *p*-values indicate the significance of gene set enrichment under a binomial distribution assumption with independence. **c** Circular dendrograms are performed for genes in different gen groups, graphically rendered based on embeddings obtained by cpDistiller. **d-f** The networks of gene relationships are provided by the GeneMANIA tool for gene Group A and gene Group B (**d**), gene Group A and gene Group C (**e**), and gene Group B and gene Group C (**f**). Each dot represents a gene, with the dot’s color indicating its related enrichment functions from GeneMANIA. **g** Dendrogram of hierarchical clustering for genes using the embeddings learned by cpDistiller.

may have indistinguishable phenotypic embeddings obtained from cpDistiller. To elucidate that the embeddings from cpDistiller can illustrate gene relationships, we utilized the GeneMANIA tool⁵⁴ to construct gene networks to visualize gene relationships based on large and diverse databases⁵⁴ (“Methods”). We found that the correlation in the network for genes in Group A and Group B is minimal (Fig. 5d), as well as the correlation for genes in Group A and Group C (Fig. 5e). By

contrast, genes in Group B and Group C exhibit closer relationships (Fig. 5f). These findings aligned with the clustering results from cpDistiller, where genes in Group A and Group C, as well as those in Group A and Group B, have significantly different embeddings, while those in Group B and Group C share similar embeddings. These results indicate that genes with similar embeddings by cpDistiller tend to have closer relationships, illuminating that cpDistiller can uncover gene

relationships. Moreover, we used GeneMANIA to predict the significantly enriched gene functions within different groups and found that functional enrichment results consistently aligned with the clustering results (Fig. 5d–f). These confirmed that analyzing the embeddings obtained by cpDistiller has the potential to infer gene functions. In conclusion, these results demonstrated the phenotypic embeddings obtained by cpDistiller effectively capture gene functions and relationships, highlighting its potential as a valuable tool for exploring gene interactions.

In addition to the combination with scRNA-seq data, we further demonstrated that cpDistiller alone is also capable of uncovering gene relationships. We first applied cpDistiller to each batch in the ORF dataset, with each batch containing approximately 2000 genetic perturbations. Since many genetic perturbations did not induce significant phenotypic changes, we then screened for active genes that exhibited substantial phenotypic changes compared to controls in each batch (“Methods”). We then selected the top 200 genes exhibiting the greatest phenotypic divergence from the controls, based on their distances in the embedding space of cpDistiller, and performed hierarchical clustering using the pairwise distances between their embeddings (“Methods”). Beyond GO term analysis, we further explored transcriptional regulation by performing transcription factor (TF) enrichment analysis using the ChEA and Enrichr Submissions TF-Gene Co-occurrence to enhance biological interpretability⁴⁸ (“Methods”). As illustrated in Fig. 5g, taking Batch_1 as an example, gene Cluster_1 were significantly enriched for *AF4* (*AFF1*) targets, suggesting involvement in hematopoietic transcriptional programs⁵⁵. In gene Cluster_2, we observed a strong enrichment of *YY1* targets, linked to transcriptional regulation processes relevant to kidney development⁵⁶. Gene Cluster_3 were notably associated with *LBX1*, *TLX3*, and *VAX1* targets, pointing to the influence of neural transcriptional regulators^{57,58}. For gene Cluster_4, *STAT1* target enrichment was predominant, implicating immune response and cancer-related transcriptional pathways⁵⁹. The gene Cluster_5 exhibited significant enrichment for *CTCF* targets, consistent with roles in chromatin organization and cancer regulation⁶⁰. Finally, the gene Cluster_6 showed prominent enrichment of *RUNX2* targets, indicating potential involvement in bone development and skeletal gene regulation^{61,62}. In contrast to the hierarchical clustering of Baseline (Supplementary Fig. 34a, b), which did not exhibit such well-defined transcriptional associations, these findings collectively demonstrate that cpDistiller embeddings not only preserve but also organize biologically meaningful regulatory programs across distinct gene groups, offering deeper interpretability beyond Baseline. These results also confirmed that analyzing the embeddings learned by cpDistiller alone also has the potential to reveal gene relationships.

cpDistiller shows potential for facilitating gene and compound target discovery

The JUMP dataset is primarily generated using U2OS cells, with all ORF experimental data in the cpG0016 dataset conducted on this cell type. However, since biological research often extends beyond U2OS cells, the JUMP dataset includes a pilot dataset, cpG0000, conducted on A549 cells in a single batch. The cpG0000 dataset provides paired genetic and compound perturbations targeting the same gene, offering substantial potential for uncovering biological targets⁴³. Leveraging this design, researchers have established simulated tasks to retrieve gene and compound targets using features extracted by CellProfiler⁴³. However, as shown in Fig. 6a, b, when we performed UMAP visualization for A549 cells, we observed noticeable row and column effects for the embedding of Baseline, which could obscure real biological signals. To address this, we applied cpDistiller to correct both row and column effects. As illustrated in Fig. 6a, b, the data distributions across different rows and columns were more uniform.

We further demonstrated that the cpDistiller-derived embeddings show potential for facilitating the identification of gene and

compound targets. Specifically, the cpG0000 dataset includes genetic and compound perturbations for 160 genes in A549 cells, conducted at both long and short time points. Each gene is perturbed through one ORF treatment, two gene knockouts by Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) guides, and two compound experiments⁴³. Using sister CRISPR guides targeting the same gene, researchers designed retrieval tasks to simulate gene target identification, calculating the fraction retrieved score to evaluate retrieval performance⁴³. Following their workflows, we calculated the fraction of retrieved scores using the embeddings from cpDistiller and compared the scores with Baseline. This evaluation design, which integrates statistical validation and permutation-based significance testing, ensures a rigorous and biologically meaningful assessment of cpDistiller’s capability to prioritize gene-compound relationships (Supplementary Note 10). As shown in Fig. 6c, Baseline achieved a fraction retrieved score of nearly 0.1 for retrieving sister CRISPR guides at both long time points (long-time CRISPR retrieval) and short time points (short-time CRISPR retrieval). In contrast, cpDistiller leverages the perturbation similarity defined by CellProfiler-based features and refines it through a combination of MNN and KNN relationships, enabling more accurate recovery of true perturbation neighborhoods. As a result, the fraction of retrieved scores improves to 0.240 for long-time CRISPR retrieval and 0.389 for short-time CRISPR retrieval. This represented a notable increase, particularly for short-time retrieval tasks, where cpDistiller achieves over a two-fold increase compared to CellProfiler-based features. These results highlighted cpDistiller holds promise for enhancing the accuracy of gene retrieval tasks, making it a valuable tool for identifying genes involved in similar processes and uncovering critical gene relationships.

In addition to retrieving sister perturbations, researchers also conducted cross-modality gene-compound retrieval tasks to simulate the identification of compound targets⁴³. Concretely, they searched for compounds that produced similar effects on cell morphology as the query gene using CellProfiler-based features. They calculated the fraction retrieved scores to evaluate retrieval performance. Following their methodology (Supplementary Note 10), we calculated the fraction retrieved scores using embeddings generated by cpDistiller to demonstrate that cpDistiller can improve retrieval results. As shown in Fig. 6d, Baseline yields fraction retrieved scores of nearly 0.008 for retrieving all compound-CRISPR pairs. In contrast, when using cpDistiller-derived embeddings for compound perturbations at both long and short time points, the fraction retrieved scores show a notable improvement in retrieving compound-CRISPR pairs (Fig. 6d). Given the complexity and importance of gene-compound retrieval tasks, even slight improvements hold substantial value⁴³. While the current evaluation does not aim to identify definitive targets, the observed performance gains of cpDistiller over Baseline suggest its potential to support the prioritization of biologically plausible gene-compound associations for follow-up validation.

In conclusion, cpDistiller satisfactorily corrected well-position effects and obtained embeddings that hold potential for exploring gene and compound targets, providing valuable insights for biological research and drug discovery.

Discussion

We developed cpDistiller, a method specifically designed to correct triple effects, particularly well-position effects in Cell Painting (CP) data. cpDistiller leverages raw CP images by integrating high-level features via a pre-trained segmentation model with handcrafted features from traditional methods, followed by a semi-supervised GMVAE utilizing contrastive and domain-adversarial learning to correct triple effects. We first conducted systematic experiments to demonstrate that CP data inherently exhibits triple effects, including batch, row, and column effects. Through comprehensive experiments on multiple batches varying in cell types, plate designs, perturbation types, we

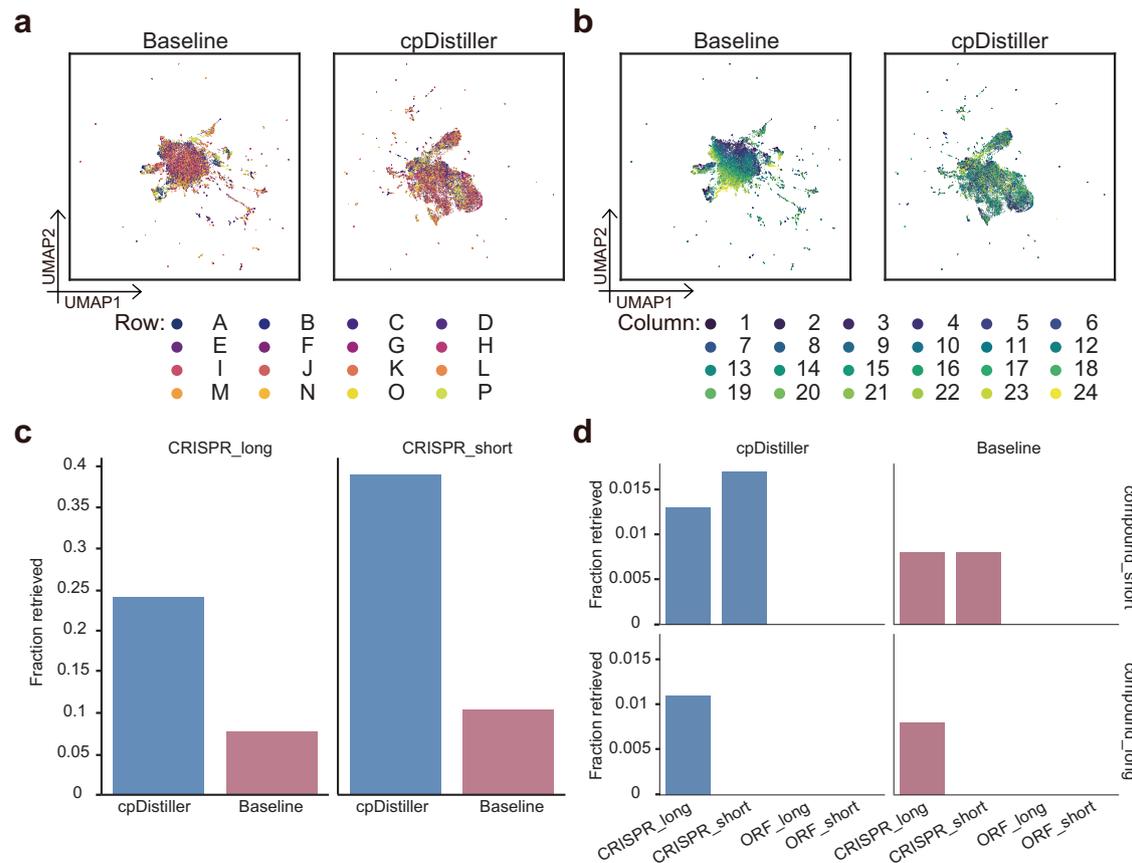


Fig. 6 | The performance of cpDistiller in retrieval tasks. a, b UMAP visualization for the A549 cells using the embedding acquired by Baseline and cpDistiller, colored by row (**a**) and column (**b**). Baseline represents the preprocessed but uncorrected CellProfiler-based features. **c** Fraction retrieved scores for retrieving sister CRISPR guides in A549 cells using embeddings of Baseline and cpDistiller,

respectively. Short (**_short**) and long (**_long**) time points mean the experimental conditions. **d** Fraction retrieved scores for retrieving gene-compound pairs in A549 cells using the embeddings of Baseline and cpDistiller, respectively. Short (**_short**) and long (**_long**) time points mean the experimental conditions.

then validated cpDistiller's superior performance in correcting triple effects and preserving cellular phenotypic heterogeneity. Moreover, we highlighted the extensive advantages of cpDistiller, including its utility in multiple conditions, its support for incremental learning, and its consistent robustness across various feature selection strategies. In addition, we conducted various downstream experiments to illustrate the application of cpDistiller in revealing gene functions and gene relationships, highlighting its ability to uncover gene associations both when combined with scRNA-seq data and independently. Scientists have found that target-based drug discovery can be limited in certain situations, and phenotypic drug discovery is sometimes more likely to succeed⁴⁴. In recent years, research related to CP has expanded significantly, offering a better perspective on studying drug targets and cellular phenotypic heterogeneity due to its simpler procedures and lower costs. Therefore, we also explored whether cpDistiller could reveal biologically meaningful patterns that may assist downstream investigations of gene-compound relationships.

While cpDistiller demonstrated excellent performance, there are several areas that could be explored for future improvement. First, compared to CellProfiler-based features, which have specific and interpretable for each dimension, the low-dimensional representations learned from cpDistiller may lack interpretability. Enhancing the interpretability of deep learning models remains a challenge⁴⁴. Second, we can utilize Segment Anything Model⁵³ along with expert annotation and refinement to generate accurate annotations for a subset of CP images. These annotations will serve as a basis for supervised training of the cpDistiller-

extractor module, allowing it to effectively capture more detailed cellular information. Third, to gain a more comprehensive understanding of cellular behavior, cpDistiller could be extended to integrate with multi-omics profiling, connecting morphological and molecular phenotypes at single-cell resolution, thereby providing a more detailed insight into genetic perturbations and their effects⁶⁴. Fourth, recent efforts have introduced datasets with multiple time points, which could support more detailed modeling of cellular responses⁶⁵. Although these datasets are not yet available at a large scale, several approaches have emerged that combine chemical structure data with CP profiles to enhance biological insight⁶⁵. Such strategies also suggest potential directions for future extensions of cpDistiller. Fifth, among triple effects, our focus on batch effects is limited to technical variation across experiments within cpg0016. However, integrating distinct datasets from different sources remains an open challenge for future work.

Methods

This research does not involve human participants, animal subjects, or the use of biological samples requiring institutional ethics approval. All relevant ethical regulations were followed.

Overview of cpDistiller

cpDistiller consists of three main modules: the extractor module, the joint training module, and the technical correction module. We first used the cpDistiller-extractor module to extract features from Cell

Painting (CP) images. Then we integrated cpDistiller-extractor-based and CellProfiler-based features via the joint training module and processed the combined features through the technical correction module to remove batch, row, and column effects.

The extractor module

We develop the extractor module to leverage the distinct feature extraction paradigm of pre-trained segmentation models, where spatial and morphological patterns are learned directly from large-scale image data. Unlike traditional pipeline using CellProfiler, this approach minimizes biases from manually engineered features and captures phenotypic variations that may be underrepresented in a conventional pipeline. By modeling our extraction task as the segmentation task in computer vision, we transform raw images into low-dimensional representations using an end-to-end process. We select Mesmer³³, pre-trained on the TissueNet datasets³³, as the base model due to its strong performance in cell segmentation tasks on cellular images⁶⁶. By experimenting with various intermediate layers of Mesmer and considering factors such as overall effectiveness, number of parameters, and processing speed, we select the backbone of Mesmer and its pre-trained parameters.

In the CP assay, each well contains 9 sub-images sized at 1080×1080 pixels, captured across five channels: mitochondria (Mito), nucleus (DNA), nucleoli and cytoplasmic RNA (RNA), endoplasmic reticulum (ER), and Golgi and plasma membrane and the actin cytoskeleton (AGP)⁴⁴. To meet Mesmer's dual-channel input requirements and enhance the visibility of cell contours, we discard the AGP and Mito channels, focusing instead on the DNA channel and composite cytoplasmic regions created by averaging the ER and RNA channels, as the information in the AGP and Mito channels is difficult to predict⁴⁴. Besides, to align with Mesmer's specifications for pixel density and reduce computational resources, we apply a tiling operation, adjusting the stride and overlap ratio to crop each large sub-image into multiple small images with dimensions of 256×256 pixels. The number of small images generated from each well's sub-image is calculated using the following formula:

$$N_{image} = \left(\left\lceil \frac{H-h}{s(1-o)} \right\rceil + 1 \right) \times \left(\left\lceil \frac{W-w}{s(1-o)} \right\rceil + 1 \right), \quad (1)$$

where H and W represent the height and width of each large sub-image, while h and w denote the height and width of the cropped images. s denotes the sliding stride, which refers to the number of pixels moved per step, while o represents the overlap ratio, indicating the degree of overlap between adjacent cropped images. To calculate the starting position of each crop, we determine the top-left corner coordinates $(x_{start,i}, y_{start,j})$ for each small image based on its index (i, j) :

$$\begin{aligned} x_{start,i} &= i \cdot s \cdot (1 - o) \\ y_{start,j} &= j \cdot s \cdot (1 - o), \end{aligned} \quad (2)$$

where i and j represent the crop index along the height and width. Specifically, i ranges from 0 to $\left\lceil \frac{H-h}{s(1-o)} \right\rceil$ and j ranges from 0 to $\left\lceil \frac{W-w}{s(1-o)} \right\rceil$.

After processing each cropped image through Mesmer, we reverse the tiling operation to reconstruct the large feature map, then apply a 2D max pooling operation followed by flattening to produce the embedding for each sub-image. To aggregate the embeddings of the nine sub-images per well, we sequentially concatenate them to form a single overall embedding for each well.

Finally, we obtain the feature matrix $\mathbf{E}_e \in \mathbb{R}^{n \times I_e}$ extracted by the extractor module, where n denotes the number of wells and I_e denotes the feature dimensions. We also refer to these as cpDistiller-extractor-based features.

The joint training module

We design the joint training module to integrate CellProfiler-based features and cpDistiller-extractor-based features. To reduce potential noise and redundancy in the high-dimensional cpDistiller-extractor-based features, we initially use average pooling to reduce the feature dimensionality and smooth out irrelevant variation, resulting in the feature matrix \mathbf{E}_{pooled} for further processing. Besides, we further employ two approaches to extract valuable information. For the first part, we approximate principal component analysis (PCA) using a linear layer to obtain the low-dimensional representation for each well, which we refer to as critical information:

$$\mathbf{Y}_1 = \text{PCA}(\mathbf{E}_{pooled}). \quad (3)$$

In the second part, we reshape \mathbf{E}_{pooled} back into 2D feature maps and pass them through an attention module to get the global information \mathbf{Y}_2 :

$$\mathbf{Y}_2 = \text{Reshape}(\mathbf{E}_{pooled}) \odot \sigma(\text{Cov1d}(\text{AvgPool}(\text{Reshape}(\mathbf{E}_{pooled}))))), \quad (4)$$

where the AvgPool represents a 2D average pooling operation, while Cov1d indicates a 1D convolution to capture inter-channel dependencies. σ denotes the Sigmoid function, which is used to produce a channel-wise attention map. \odot represents element-wise multiplication between the attention map and the reshaped \mathbf{E}_{pooled} .

To fusing critical and global information in the low-dimensional space, we apply element-wise addition as the encoding process for \mathbf{E}_{pooled} . The latent representations \mathbf{z}_e for \mathbf{E}_{pooled} are computed as:

$$\mathbf{z}_e = \mathbf{Y}_1 \oplus (W_1 \mathbf{Y}_2 + b_1), \quad (5)$$

where \oplus represents element-wise addition. W_1 and b_1 denote the weight and bias parameters of the encoder. Once the latent representations \mathbf{z}_e are obtained, the decoder reconstructs the data back to the same dimensionality as \mathbf{E}_{pooled} .

Ultimately, we combine the CellProfiler-based features with the cpDistiller-extractor-based features transformed by the attention mechanism-based encoder-decoder architecture and feed the combined features into the technical correction module for further refinement.

The technical correction module

The technical correction module removes triple effects from CP data and generates low-dimensional representations that maintain biological variation. Specifically, it consists of three parts: the Gaussian mixture variational autoencoder (GMVAE) as the core component, along with the contrastive learning module and the gradient reversal module.

The Gaussian mixture variational autoencoder. Given the feature matrix $\mathbf{X}_p \in \mathbb{R}^{n \times I_p}$ integrated by the joint training module, where n denotes the number of wells and I_p denotes the feature dimensions, the GMVAE takes the input \mathbf{X}_p to obtain low-dimensional representations \mathbf{Z}_p . To illustrate the workflow of the GMVAE, we consider a sample \mathbf{x} from \mathbf{X}_p . Since CP data encompasses numerous perturbations that may conform to distinct underlying Gaussian distributions, we use the GMVAE to identify these categorical distributions (pseudo-labels, denoted as \mathbf{y}), which helps the model capture biological variation and ultimately contribute to biologically meaningful low-dimensional representations (denoted as \mathbf{z}). The categorical distributions are inferred from the posterior distribution $q_\psi(\mathbf{y}|\mathbf{x})$, which follows a semi-supervised training pattern. Here, $q_\psi(\mathbf{y}|\mathbf{x})$ is represented by a feedforward neural network as: $q_\psi(\mathbf{y}|\mathbf{x}) = \text{Cat}(\mathbf{y}|\boldsymbol{\pi}_\psi(\mathbf{x}))$ and $\boldsymbol{\pi}_\psi(\mathbf{x})$

is a probability vector. Since the categorical distribution cannot be backpropagated in the neural network, we use the Gumbel-Softmax distribution⁶⁷ to facilitate gradient backpropagation, which allows the categorical distribution to be approximated using a continuous distribution.

In GMVAE, the objective is to optimize the Evidence Lower Bound (ELBO), which is expressed as follows:

$$\text{Log}p_{\theta}(\mathbf{x}) \geq E_{q_{\psi}(\mathbf{y}, \mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}|\mathbf{y}) - \log q_{\psi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) - \log q_{\psi}(\mathbf{y}|\mathbf{x}) + \log p(\mathbf{y})], \quad (6)$$

where these components can be broadly divided into three main optimization directions. The first optimization direction focuses on the reconstruction loss, represented by the expectation $E_{q_{\psi}(\mathbf{y}, \mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]$. This loss is measured using mean squared error (MSE), which ensures that the low-dimensional representations \mathbf{z} effectively capture the key information of the original input \mathbf{x} .

The term $E_{q_{\psi}(\mathbf{y}, \mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}|\mathbf{y}) - \log q_{\psi}(\mathbf{z}|\mathbf{x}, \mathbf{y})]$ represents the Kullback-Leibler (KL) divergence between the variational posterior distribution $q_{\psi}(\mathbf{z}|\mathbf{x}, \mathbf{y})$ and the conditional prior distribution $p_{\theta}(\mathbf{z}|\mathbf{y})$. This divergence ensures that the variational posterior distribution aligns with the prior, meaning that the learned latent representations \mathbf{z} conform to the expected Gaussian distribution. Specifically, $p_{\theta}(\mathbf{z}|\mathbf{y}) = N(\mathbf{z}|\mu_{\theta}(\mathbf{y}), \sigma_{\theta}^2(\mathbf{y}))$ represents the prior Gaussian distribution conditioned on the category \mathbf{y} , while $q_{\psi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) = N(\mathbf{z}|\mu_{\psi}(\mathbf{x}, \mathbf{y}), \sigma_{\psi}^2(\mathbf{x}, \mathbf{y}))$ represents the variational posterior distribution conditioned on the input \mathbf{x} and category \mathbf{y} .

$q_{\psi}(\mathbf{y}|\mathbf{x})$ provides the probability that \mathbf{x} originates from various Gaussian distributions, satisfying $\sum_{n=1}^N q_{\psi}(\mathbf{y}|\mathbf{x}) = 1$, where N is the pre-defined number of Gaussian distributions. We use cross-entropy loss to refine the probabilities, guiding them towards high-confidence regions. We treated the prior $p(\mathbf{y})$ as a constant during loss back-propagation, as it does not influence the updates to the model's parameters.

The gradient reversal module. Within the low-dimensional space derived from the GMVAE, discriminators are employed to identify the source of each well's representation \mathbf{z} , specifically the batch, row, and column labels of the corresponding well. To implement adversarial learning similar to generative adversarial networks (GANs)²³, different batches, rows, and columns can be treated as distinct domains. We then employ domain-adversarial learning, specifically the gradient reversal layer (GRL)²⁴, to remove triple effects across these domains. Here, discriminators are denoted as: $D_{batch}, D_{row}, D_{column}$, which are used to predict the batch, row, and column labels of \mathbf{z} .

Specifically, the GRL reverses the gradient during back-propagation, causing the parameters of the GMVAE's encoder, which acts as the generator, to be updated in the opposite direction of the discriminators, thereby achieving the adversarial objective. The GRL can be described as a pseudo-function: $\text{GRL}_{\lambda}(\mathbf{x})$. During the forward pass, the GRL functions as an identity operation, leaving the input parameters unchanged. However, during the backward pass, it scales the gradients from the following layers by $-\lambda$ before sending them back to the preceding layers. The forward and backward passes are described by the following two equations:

$$\begin{aligned} \text{GRL}_{\lambda}(\mathbf{x}) &= \mathbf{x} \\ \frac{\partial \text{GRL}_{\lambda}(\mathbf{x})}{\partial \mathbf{x}} &= -\lambda \mathbf{x}, \end{aligned} \quad (7)$$

where λ is a hyperparameter that undergoes a non-linear transformation, varying from 0 to 1. In the early stages of training, its value is kept small to allow the discriminators to train sufficiently and develop discriminative capabilities. As training progresses, λ gradually

increases to strengthen adversarial interactions between the encoder part of the GMVAE and the discriminators. The calculation formula for λ is as follows:

$$\lambda = \frac{2}{1 + \exp(-\gamma \cdot p)} - 1, \quad (8)$$

where γ is a hyperparameter, and p represents the percentage of the total iteration progress during training. To be specific, we denote the encoder part of the GMVAE as follows:

$$\mathbf{z} = E(\mathbf{x}; \theta_g), \quad (9)$$

where θ_g denotes the learnable parameters of the GMVAE's encoder. Subsequently, the low-dimensional representations \mathbf{z} are passed through the GRL and discriminators, followed by the Softmax function to obtain the probability distribution for batch, row, and column predictions. We further describe the discriminators in detail as:

$$D(\text{GRL}_{\lambda}(\mathbf{z}), \theta_{D_i}), \quad i \in \{batch, row, column\}, \quad (10)$$

where $\theta_{D_{batch}}, \theta_{D_{row}}$ and $\theta_{D_{column}}$ denote the learnable parameters of the batch, row, and column discriminators, which are updated to minimize the discriminator loss. Meanwhile, θ_g is updated through the GRL to maximize the discriminator loss, ensuring that the discriminators cannot distinguish the source of the low-dimensional representations, thereby obtaining representations free of technical effects.

Using the column discriminator loss as an example, the loss function is inspired by the label smoothing cross-entropy loss⁶⁸. This approach is similarly applied to the batch and row discriminators for avoiding overconfidence:

$$L_{D_{column}} = - \sum_{k=1}^K [(1 - \epsilon) \log(p_{column}(k)) \sigma_{k,l}^{column} + \epsilon \log(p_{column}(k)) \theta_k^{column}], \quad (11)$$

where l represents the index of the true label, and ϵ is a hyperparameter representing the proportion of soft labels considered when calculating the loss. $\sigma_{k,l}^{column}$ is a one-hot encoding, where the position corresponding to the true label is 1, with others set to 0. The probability that the \mathbf{z} originates from k -th column, as predicted by the discriminator, is represented by $p_{column}(k)$. θ_k^{column} represents the soft labels, which indicate the weight assigned to the k -th column.

To be specific, standard cross-entropy loss focuses solely on optimizing the predicted probability of the true label, which can lead to overconfidence and neglect the uncertainty in the model's predictions. Label smoothing cross-entropy loss addresses this limitation by redistributing a portion of the probability mass from the true label to the other classes, applying a uniform distribution across them to reduce overconfidence. However, this approach assigns equal importance to all non-true labels, which may not be appropriate for tasks where the predicted probabilities for non-true labels carry varying degrees of significance. In the context of our study, UMAP visualizations of the raw CP data show a gradient-influenced pattern that more distant column indexes exhibit more apparent column effects. This suggests that predicted probabilities closer to the true label index carry more meaning. To better capture this, we define a soft label distribution θ_k^{column} , that reflects the differences among predicted probabilities, assigning varying significance to them based on their distance from the true label, while still avoiding overconfidence:

$$\theta_k^{column} = \begin{cases} \alpha, & \text{if } k=l \\ q^{|k-l|}, & \text{otherwise} \end{cases} \quad (12)$$

where the hyperparameter α represents the weight assigned to the true label, and q is the common ratio that determines the weights for the other labels based on their distance from the true label l . The value of q is determined by solving the higher-degree equation:

$$\begin{aligned} f(q) &= \alpha - \alpha q^{l+1} - \alpha q^{n-l} - 1 + q + \alpha q \\ q : f(q) &= 0, \quad \text{s.t. } 0 < q < 1, \end{aligned} \quad (13)$$

where n represents the number of categories for the technical effects (in this case, the number of column labels). The solution to this equation ensures that the soft label distribution is normalized, meaning that the sum of all probabilities equals 1, while also reflecting a geometric decay in significance as the distance from the true label index increases.

The contrastive learning module. We first establish nearest neighbor relationships by utilizing k -nearest neighbors (KNN) intra technical effects and mutual nearest neighbors (MNN) inter technical effects, based on CellProfiler-based features and cpDistiller-extractor-based features, using cosine distance as the similarity metric. This approach is applied to batch, row, and column effects, respectively. The intersection of the nearest neighbors is used to construct the adjacency matrix that captures nearest neighbor relationships between wells. We then leverage the relationships identified across multiple technical effects to form triplets and use triplet loss²² to restore more accurate nearest-neighbor relationships for each well.

Specifically, we need to select a triplet $(\mathbf{a}, \mathbf{p}, \mathbf{n})$ to act as the anchor, positive, and negative samples. Each data point serves as the anchor in turn. For each anchor, data points with nearest neighbor relationships to that anchor are considered positive, while those without such relationships are considered negative. Specifically, (\mathbf{a}, \mathbf{p}) pairs are either in the k -nearest neighbors set $S_{\text{knn,technical effects}}$ or in the mutual nearest neighbors set $S_{\text{mnn,technical effects}}$. Conversely, (\mathbf{a}, \mathbf{n}) pairs are not found in either of the two sets. The specific formula is as follows:

$$(\mathbf{a}, \mathbf{p}, \mathbf{n}) \iff \begin{cases} (\mathbf{a}, \mathbf{p}) \in S_{\text{knn,technical effects}} \cup S_{\text{mnn,technical effects}}, \\ (\mathbf{a}, \mathbf{n}) \notin S_{\text{knn,technical effects}} \cup S_{\text{mnn,technical effects}} \end{cases}, \quad (14)$$

where $S_{\text{knn,technical effects}}$ represents the set of all pairs of data points that originate from the same type of technical effects and are k -nearest neighbors within the same category. $S_{\text{mnn,technical effects}}$ represents the set of all pairs of data points that originate from the same type of technical effects but are mutual nearest neighbors in different categories. To accommodate the distinct role of negative controls, we incorporate an additional sampling strategy during triplet construction. When the anchor is not a negative control, negative samples are preferentially drawn from negative controls that are not part of the KNN or MNN sets related to the anchor. In the absence of such candidates, other samples without nearest neighbor relationships to the anchor are considered. Conversely, when a negative control serves as the anchor, positive samples are preferentially selected from other negative controls.

The triplet loss function is not directly applied to the latent space of the GMVAE. Instead, the representations \mathbf{z} need to pass through a nonlinear projection head, as previous research has demonstrated that adding such a layer can improve the quality of the learned representations⁶⁹. This architecture can be represented as $ph(\mathbf{z})$:

$$ph(\mathbf{z}) = \text{LeakyReLU}(W_{\text{ph}}\mathbf{z} + \mathbf{b}_{\text{ph}}), \quad (15)$$

where the parameters W_{ph} and \mathbf{b}_{ph} represent the learnable weights and biases of a linear layer.

Then we use the triplet loss to remove triple effects:

$$L_c = \max(d(\text{ph}(\mathbf{a}), \text{ph}(\mathbf{p})) - d(\text{ph}(\mathbf{a}), \text{ph}(\mathbf{n})) + \xi, 0), \quad (16)$$

where $d(\text{ph}(\mathbf{a}), \text{ph}(\mathbf{p}))$ represents the Euclidean distance between the anchor and positive samples after passing through the projection head, and the ξ is a hyperparameter.

If we aim to correct row and column effects, the overall loss can be written as follows:

$$\text{Loss} = w_{\text{ELBO}}L_{\text{ELBO}} + w_{\text{dis}}(L_{D_{\text{row}}} + L_{D_{\text{column}}}) + w_{\text{con}}L_c, \quad (17)$$

where L_c represents the triplet loss calculated by considering well-position effects.

If we need to correct batch, row, and column effects, the overall loss can be expressed as follows:

$$\text{Loss} = w_{\text{ELBO}}L_{\text{ELBO}} + w_{\text{dis}}(L_{D_{\text{batch}}} + L_{D_{\text{row}}} + L_{D_{\text{column}}}) + w_{\text{con}}L_c, \quad (18)$$

where L_c represents the triplet loss calculated by considering triple effects. In the above loss functions, the weights of w_{ELBO} , w_{dis} , and w_{con} are the weighting coefficients assigned to different components of the loss function. These coefficients control the relative importance of the ELBO, discriminator losses, and the triplet loss in the overall optimization process.

The parameters θ'_t of the final trained model are obtained through exponential moving average (EMA)⁷⁰, which can be given by:

$$\theta'_t = \alpha_{\text{ema}}\theta'_{t-1} + (1 - \alpha_{\text{ema}})\theta_t, \quad (19)$$

where θ_t represents the weighted model parameters obtained at round t , and α_{ema} is a hyperparameter. All training hyperparameters are available in the training details.

Training details

For the extractor module, the overlap ratio o is set to 0.25, and the sliding stride s is set to 256. Max pooling with kernel size and step size both set to 16 is applied to merge feature maps from nine sub-images, yielding the output dimension I_e to 11,664, corresponding to the 108×108 feature map. For the joint training module, average pooling with kernel size and step size both set to 9, is applied to smooth out irrelevant variation. For the technical correction module, the hidden dimensions of the encoder and decoder are set to 512, while the hidden dimensions of the discriminators are set to 128. The latent space dimensionality of cpDistiller is set to 50. The technical correction module uses LeakyReLU with default parameters as the activation function throughout, except for the variance inference component, which utilizes the Softplus activation function. The projection head consists of a linear layer with an output dimension of 50 and a LeakyReLU activation function. The optimizer used is AdamW, with the learning rate set to $3e-3$ for the discriminators and $1e-3$ for the other parts. Although we optimized the adversarial learning, mode collapse is still a common issue. When considering 7638 CellProfiler-based features, the default learning rate is unsuitable for training and can lead to collapse. Therefore, we adjusted the initial learning rate to half of the default value. The γ is set to 10 in GRL. The soft label hyperparameter α is set to 0.75, and the hyperparameter ϵ in label smoothing cross-entropy loss is set to 0.1. For the contrastive learning, the ξ is set to 10, and the nearest neighbor hyperparameters for MNN and KNN are set to 5 and 10, respectively. The default number of epochs is 50, with α_{ema} set to $1 - \frac{5}{\text{epochs}}$. In the loss functions, the weights of w_{ELBO} , w_{dis} , and w_{con} are set to $1, \frac{1}{35}$, and $\frac{1}{35}$, respectively. The experimental environment includes two 24GB Nvidia 4090 graphics cards and 96 Intel(R) Xeon(R) Gold 5318 N CPUs @ 2.10 GHz.

Implementation details of downstream analyses

The establishment of gene groups based on scRNA-seq data: we used the gget tool, a Python package available at <https://github.com/pachterlab/gget>, which enables efficient querying of the top

100 similar genes of the input gene calculated by the ARCHS4⁴⁸ based on scRNA-seq data. Concretely, we input the 12,602 genes experimented in the cpg0016 dataset and retrieved the top 100 most similar genes for each input gene, forming 12,602 gene sets. Since a similarity score between 0.6 and 0.8 was considered significant⁷¹, we selected 0.6 as the threshold to include more genes and got more applicable gene sets. Finally, we intersected each gene set with the genes present in the cpg0016 dataset to obtain the final gene groups of similar genes.

Gene Ontology (GO) term analysis: We utilized the Enrichr⁷² platform to conduct GO term enrichment analysis⁷³ for each gene group, assigning relevant GO terms to the genes. We specifically focused on Cellular Component (CC) analysis to identify the organelles associated with each gene and ranked the results by their statistical significance.

ChEA enrichment analysis: We utilized the Enrichr⁷² platform to perform ChEA analysis for each gene group, aiming to identify transcription factors (TFs) that potentially regulate the genes within each cluster. The results were ranked by statistical significance based on combined scores, highlighting key TFs enriched in each cluster.

Enrichr Submissions TF-Gene Co-occurrence analysis: we used the Enrichr⁷² platform to perform TF-Gene Co-occurrence analysis based on community-submitted datasets, aiming to identify transcription factors that frequently co-occur with each gene group. Enriched TFs were ranked by combined scores to highlight potential regulatory patterns.

Gene relationship analysis: We used the GeneMANIA tool⁵⁴ to evaluate the relationship of genes in different groups. We submitted the genes in gene groups to GeneMANIA and used the default network selections. Specifically, GeneMANIA first identifies genes that are functionally similar or share properties with the submitted genes, then builds the network connecting both the submitted and similar genes. It then displays the gene network regarding co-expression networks, physical interaction, genetic interaction, co-localization, pathways, and predicted and shared protein domain information.

Screening for active genes: following the approach established by ref. 14, we calculated the average of repetitions for each perturbation at the same position across five different plates to obtain the average representation for each perturbation. Considering that negative treatments generally do not induce significant phenotypic changes, we calculated the Euclidean distances between negative treatments and set the 95th percentile of these distances as the threshold for identifying active genes. Overexpression treatments were classified as active if their Euclidean distance from the negative treatments exceeded this threshold.

Ranking the significance of the clusters in single-link hierarchical clustering: We applied a perturbation approach to rank the significance of clusters in single-link hierarchical clustering. For each cluster, we randomly selected the same number of data points as contained in that cluster. We then calculated the minimum clustering distance among these randomly selected points. This perturbation process was repeated 1000 times for each cluster, generating 1000 distance scores. By comparing the actual cluster distance with the distribution of the 1000 perturbed distance scores, we ranked the clusters and got the ranking of the actual cluster, which provides an indication of the significance level of data point aggregation within each cluster.

The calculation of fraction retrieved in gene and compound retrieval tasks: we calculated and compared the fraction retrieved scores for gene-gene and gene-compound simulated retrieval tasks using cpDistiller-derived representations in the cpg0000 dataset, following the workflow in ref. 43.

Pre-processing for CellProfiler-based features. We utilized features pre-extracted with CellProfiler from the prior study¹⁴, which included up to 7,638 features for images with both brightfield and fluorescence

images and 4752 features for only fluorescence images, after removing features containing NaN values. In addition, we followed the feature selection process described in the study¹⁴ to select 1,446 features. Due to data quality issues reported in the original dataset¹⁴, we excluded data from Batch_12 and the BR00123528A plate. For pre-computed feature transformation, we applied a plate-wise normalization strategy using the median absolute deviation (MAD) method, using negative control wells to calibrate feature distributions within each plate⁷⁴ (Supplementary Note 11 and Supplementary Figs. 35 and 36). Subsequently, to ensure comparability across different features, we further standardized the data using z-score normalization.

Data pre-processing for Cell Painting images. To match the dual-channel input format required by Mesmer of the extractor module, we extracted data from three channels: DNA, RNA, and ER. The DNA channel primarily pertains to the nucleus, while the RNA and ER channels are related to the cytoplasm. After computing these with their respective illumination files, we combined them into two separate channels. The processed data was then stored in NPZ format to prepare the image data for the extractor module. For image pre-processing, we utilized the procedure from Mesmer, which included using Contrast Limited Adaptive Histogram Equalization (CLAHE) to enhance local contrast, followed by logarithmic smoothing of the data in the first channel.

Evaluation metrics

To assess the effectiveness of different methods in removing technical effects, we used three technical correction metrics: average silhouette width (ASW)²⁹, technic average silhouette width (tASW)³⁰, and graph connectivity³⁰. To evaluate biological preservation, four metrics were used: perturbation average silhouette width (pASW)³⁰, normalized mutual information (NMI)³⁰, phenotypic activity²⁸, and phenotypic consistency²⁸. The specific calculation of the above metrics, following the approach of previous works including scArches²⁹, scIB³⁰, and copairs²⁸ for single-cell and CP data analysis, is further detailed in Supplementary Note 12.

Baseline methods

To evaluate technical correction in CP data, we compared several widely used single-cell analysis methods: Seurat v5 (v5.0.1)²⁵, Harmony (v0.0.9)²⁰, Scanorama (v1.7.4)²⁶, scVI (v0.14.6)¹⁹, and scDML (v0.0.1)²⁷. UMAP plots, used to compare embeddings generated by different methods, were created with the following parameters: “n_neighbors = 15”, “min_dist = 0.1”, and “random_state = 9000”. We performed plate-wise normalization for pre-computed feature transformation using the median absolute deviation method, using negative controls to calibrate feature distributions within each plate. To further ensure comparability across features, we subsequently standardized the data using z-score normalization. All methods, including cpDistiller, were applied to CP data that had undergone the above preprocessing steps to ensure consistency. To provide a comparison with uncorrected data, we also included Baseline, which represents CellProfiler-based features after the same preprocessing but without any further correction.

For Harmony²⁰, Seurat v5²⁵, and Scanorama²⁶, which can correct multiple technical effects, the workflow involved sequentially passing the row and column labels if considering removing well-position effects. When aiming to remove triple effects, these methods require first passing the batch labels, followed sequentially by the row and column labels. For methods that are limited to correcting only one type of technical effects, if considering well-position effects, we selected row and column effects separately as correction targets, resulting in two results, respectively. Considering triple effects, we focused solely on batch effects as they represent the most important influence. Specific implementation details are as follows:

Seurat v5²⁵: We followed the example pipeline provided by Seurat v5²⁵ for integrative analysis. The labels for technical effects were passed sequentially into the “split” function. We skipped the “NormalizeData” function and “FindVariableFeatures” function, as these are specifically tailored for scRNA-seq data. During the correction phase, the “CCA” method was employed to remove technical effects. All parameters were used with default settings.

Harmony²⁰: Dimensionality reduction was performed using “scanpy.tl.pca” to the default size of 50. The labels for technical effects were passed sequentially into the “scanpy.external.pp.harmony_integrate” function. All parameters were used with default settings.

Scanorama²⁶: Dimensionality reduction was performed using “scanpy.tl.pca” to the default size of 50. The labels for technical effects were passed sequentially into the “scanpy.external.pp.scanorama_integrate” function. All parameters were used with default settings.

scVI¹⁹: We followed the example tutorial provided by the scVI¹⁹ on GitHub for processing scRNA-seq data. We did not use the “scanpy.pp.highly_variable_genes” function, which is specifically tailored for scRNA-seq data. In the preprocessing stage, we followed the preprocessing operations of previous work to process each feature as follows: $\hat{x}_i = x_i - \min(x) + 1$ ¹⁵. Other operations were used with default settings. We additionally considered the use of “scanpy.pp.highly_variable_genes” during preprocessing, and conducted comparative analyses to assess its potential impact (Supplementary Note 13 and Supplementary Fig. 37).

scDML²⁷: We followed the example tutorial on GitHub provided by scDML²⁷. We did not use the “scanpy.pp.log1p” function and the “scanpy.pp.highly_variable_genes” function, which are specifically tailored for scRNA-seq data. Other operations were used with default settings. In addition, we conducted analyses to assess the impact of applying “scanpy.pp.highly_variable_genes” on the performance of analyses (Supplementary Note 13 and Supplementary Fig. 37).

Statistics and reproducibility

No statistical method was used to predetermine sample size. Due to data quality issues reported in the original dataset¹⁴, we excluded data from Batch_12 and the BR00123528A plate. Randomization and blinding, which are important considerations in biological experiments, were not applicable in our study, as it is a computational analysis based on publicly available CP datasets.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The processed Cell Painting data are available at <https://registry.opendata.aws/cellpainting-gallery/>. The metadata files containing compound annotations and genetic perturbation information are available at <https://github.com/jump-cellpainting/JUMP-MOA> and <https://github.com/jump-cellpainting/datasets>. Source data are provided in this paper.

Code availability

The cpDistiller software, along with detailed documentation and tutorials, is freely available at <https://github.com/BioX-NKU/cpDistiller>⁷⁵.

References

- Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of next-generation sequencing technologies. *Nat. Rev. Genet.* **17**, 333–351 (2016).
- Perlman, Z. E. et al. Multidimensional drug profiling by automated microscopy. *Science* **306**, 1194–1198 (2004).
- Wawer, M. J. et al. Toward performance-diverse small-molecule libraries for cell-based phenotypic screening using multiplexed high-dimensional profiling. *Proc. Natl. Acad. Sci. USA* **111**, 10911–10916 (2014).
- Moshkov, N. et al. Learning representations for image-based profiling of perturbations. *Nat. Commun.* **15**, 1594 (2024).
- Fredin et al. Cell Painting-based bioactivity prediction boosts high-throughput screening hit-rates and compound diversity. *Nat. Commun.* **15**, 3470 (2024).
- Bray, M.-A. et al. Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nat. Protoc.* **11**, 1757–1774 (2016).
- Chandrasekaran, S. N., Ceulemans, H., Boyd, J. D. & Carpenter, A. E. Image-based profiling for drug discovery: due for a machine-learning upgrade? *Nat. Rev. Drug Discov.* **20**, 145–159 (2021).
- Jo, Y. et al. Label-free multiplexed microtomography of endogenous subcellular dynamics using generalizable deep learning. *Nat. Cell Biol.* **23**, 1329–1337 (2021).
- Liberati, P., Snijder, B. & Pelkmans, L. A hierarchical map of regulatory genetic interactions in membrane trafficking. *Cell* **157**, 1473–1487 (2014).
- Collinet, C. et al. Systems survey of endocytosis by multiparametric image analysis. *Nature* **464**, 243–249 (2010).
- Carpenter, A. E. & Singh, S. Bringing computation to biology by bridging the last mile. *Nat. Cell Biol.* **26**, 5–7 (2024).
- Loo, L.-H., Wu, L. F. & Altschuler, S. J. Image-based multivariate profiling of drug responses from single cells. *Nat. Methods* **4**, 445–453 (2007).
- Weisbart, E. et al. Cell Painting Gallery: an open resource for image-based profiling. *Nat. Methods* **21**, 1775–1777 (2024).
- Chandrasekaran, S. N. et al. JUMP Cell Painting dataset: morphological impact of 136,000 chemical and genetic perturbations. Preprint at <https://doi.org/10.1101/2023.03.23.534023> (2023).
- Arevalo, J. et al. Evaluating batch correction methods for image-based cell profiling. *Nat. Commun.* **15**, 6516 (2024).
- Tromans-Coia, C. et al. Assessing the performance of the Cell Painting assay across different imaging systems. *Cytometry A* **103**, 915–926 (2023).
- Carpenter, A. E. et al. CellProfiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biol.* **7**, 1–11 (2006).
- Bouland, G. A., Mahfouz, A. & Reinders, M. J. Consequences and opportunities arising due to sparser single-cell RNA-seq datasets. *Genome Biol.* **24**, 86 (2023).
- Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
- Weisbart, E. et al. Cellprofiler Plugins—an easy image analysis platform integration for containers and python tools. *J. Microsc.* **296**, 227–234 (2024).
- Schroff, F., Kalenichenko, D. & Philbin, J. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* 815–823 (CVPR, 2015).
- Salimans, T. et al. Improved techniques for training gans. In *Advance in Neural Information Processing Systems* (NeurIPS, 2016).
- Ganin, Y. & Lempitsky, V. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning* (ICML, 2015).
- Hao, Y. et al. Dictionary learning for integrative, multimodal and scalable single-cell analysis. *Nat. Biotechnol.* **42**, 293–304 (2024).
- Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nat. Biotechnol.* **37**, 685–691 (2019).

27. Yu, X., Xu, X., Zhang, J. & Li, X. Batch alignment of single-cell transcriptomics data using deep metric learning. *Nat. Commun.* **14**, 960 (2023).
28. Kalinin, A. A. et al. A versatile information retrieval framework for evaluating profile strength and similarity. *Nat. Commun.* **16**, 1–17 (2025).
29. Lotfollahi, M. et al. Mapping single-cell data to reference atlases by transfer learning. *Nat. Biotechnol.* **40**, 121–130 (2022).
30. Luecken, M. D. et al. Benchmarking atlas-level data integration in single-cell genomics. *Nat. Methods* **19**, 41–50 (2022).
31. Caicedo, J. C. et al. Data-analysis strategies for image-based cell profiling. *Nat. Methods* **14**, 849–863 (2017).
32. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. Preprint at <https://doi.org/10.48550/arXiv.1802.03426> (2018).
33. Greenwald, N. F. et al. Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nat. Biotechnol.* **40**, 555–565 (2022).
34. Shrestha, P., Kuang, N. & Yu, J. Efficient end-to-end learning for cell segmentation with machine generated weak annotations. *Commun. Biol.* **6**, 232 (2023).
35. Jocher, G., Chaurasia, A. & Qiu, J. Ultralytics yolov8. <https://github.com/ultralytics/ultralytics> (2023).
36. Dahlin, J. L. et al. Reference compounds for characterizing cellular injury in high-content cellular morphology assays. *Nat. Commun.* **14**, 1364 (2023).
37. Seal, S. et al. Small molecule bioactivity benchmarks are often well-predicted by counting cells. Preprint at <https://doi.org/10.1101/2025.04.27.650853> (2025).
38. Tran, H. T. N. et al. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol.* **21**, 1–32 (2020).
39. Haghverdi, L., Lun, A. T., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
40. Chazarra-Gil, R., van Dongen, S., Kiselev, V. Y. & Hemberg, M. Flexible comparison of batch correction methods for single-cell RNA-seq using BatchBench. *Nucleic Acids Res.* **49**, e42 (2021).
41. Song, Y., Miao, Z., Brazma, A. & Papatheodorou, I. Benchmarking strategies for cross-species integration of single-cell RNA sequencing data. *Nat. Commun.* **14**, 6495 (2023).
42. Cimini, B. A. et al. Optimizing the Cell Painting assay for image-based profiling. *Nat. Protoc.* **18**, 1981–2013 (2023).
43. Chandrasekaran, S. N. et al. Three million images and morphological profiles of cells treated with matched chemical and genetic perturbations. *Nat. Methods* **21**, 1114–1121 (2024).
44. Seal, S. et al. Cell Painting: a decade of discovery and innovation in cellular imaging. *Nat. Methods* **1**, 15 (2024).
45. Mahmoudi, A. & Jemielniak, D. Proof of biased behavior of Normalized Mutual Information. *Sci. Rep.* **14**, 9021 (2024).
46. Wang, K. C. & Chang, H. Y. Epigenomics: technologies and applications. *Circ. Res.* **122**, 1191–1199 (2018).
47. Theodoris, C. V. et al. Transfer learning enables predictions in network biology. *Nature* **618**, 616–624 (2023).
48. Lachmann, A. et al. Massive mining of publicly available RNA-seq data from human and mouse. *Nat. Commun.* **9**, 1366 (2018).
49. Strutt, H., Searle, E., Thomas-MacArthur, V., Brookfield, R. & Strutt, D. A Cul-3-BTB ubiquitylation pathway regulates junctional levels and asymmetry of core planar polarity proteins. *Development* **140**, 1693–1702 (2013).
50. Takahashi, S. et al. The E3 ubiquitin ligase LNX1p80 promotes the removal of claudins from tight junctions in MDCK cells. *J. Cell Sci.* **122**, 985–994 (2009).
51. Van Campenhout, C. A. et al. Dlg3 trafficking and apical tight junction formation is regulated by nedd4 and nedd4-2 e3 ubiquitin ligases. *Dev. Cell* **21**, 479–491 (2011).
52. Vogel, F., Bornhövd, C., Neupert, W. & Reichert, A. S. Dynamic subcompartmentalization of the mitochondrial inner membrane. *J. Cell Biol.* **175**, 237–247 (2006).
53. Maly, D. J. Exploring the intermembrane space. *ACS Chem. Biol.* **2**, 213–216 (2007).
54. Warde-Farley, D. et al. The GeneMANIA prediction server: biological network integration for gene prioritization and predicting gene function. *Nucleic Acids Res.* **38**, W214–W220 (2010).
55. Harman, J. R. et al. A KMT2A-AFF1 gene regulatory network highlights the role of core transcription factors and reveals the regulatory logic of key downstream target genes. *Genome Res.* **31**, 1159–1173 (2021).
56. Patel, S. R. & Dressler, G. R. Expression of Pax2 in the intermediate mesoderm is regulated by YY1. *Dev. Biol.* **267**, 505–516 (2004).
57. Cheng, L. et al. Lbx1 and Tlx3 are opposing switches in determining GABAergic versus glutamatergic transmitter phenotypes. *Nat. Neurosci.* **8**, 1510–1515 (2005).
58. Hallonet, M., Hollemann, T., Pieler, T. & Gruss, P. Vax1, a novel homeobox-containing gene, directs development of the basal forebrain and visual system. *Genes Dev.* **13**, 3106–3114 (1999).
59. Weng, H., Yuan, S., Huang, Q., Zeng, X.-T. & Wang, X.-H. STAT1 is a key gene in a gene regulatory network related to immune phenotypes in bladder cancer: An integrative analysis of multi-omics data. *J. Cell. Mol. Med.* **25**, 3258–3271 (2021).
60. Fang, C. et al. Cancer-specific CTCF binding facilitates oncogenic transcriptional dysregulation. *Genome Biol.* **21**, 1–30 (2020).
61. Liu, T. M. & Lee, E. H. Transcriptional regulatory cascades in Runx2-dependent bone development. *Tissue Eng. Part B Rev.* **19**, 254–263 (2013).
62. Komori, T. Regulation of bone development and extracellular matrix protein genes by RUNX2. *Cell Tissue Res.* **339**, 189–195 (2010).
63. Kirillov, A. et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4015–4026 (ICCV, 2023).
64. Tang, Q. et al. Morphological profiling for drug discovery in the era of deep learning. *Brief. Bioinform.* **25**, <https://doi.org/10.1093/bib/bbae284> (2024).
65. Ewald, J. D. et al. Cell Painting for cytotoxicity and mode-of-action analysis in primary human hepatocytes. Preprint at <https://doi.org/10.1101/2025.01.22.634152> (2025).
66. Amitay, Y. et al. CellSighter: a neural network to classify cells in highly multiplexed images. *Nat. Commun.* **14**, 4302 (2023).
67. Jang, E., Gu, S. & Poole, B. Categorical Reparameterization with Gumbel-Softmax. In *International Conference on Learning Representations (ICLR, 2017)*.
68. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* 2818–2826 (CVPR, 2016).
69. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning* 1597–1607 (ICML, 2020).
70. Tarvainen, A. & Valpola, H. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advance in Neural Information Processing Systems (NeurIPS, 2017)*.
71. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **102**, 15545–15550 (2005).
72. Kuleshov, M. V. et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
73. Ashburner, M. et al. Gene ontology: tool for the unification of biology. *Nat. Genet.* **25**, 25–29 (2000).

74. Serrano, E. et al. Reproducible image-based profiling with Pycyto-miner. *Nat. Methods* **22**, 677–680 (2025).
75. Yan, C. et al. Triple-effect correction for Cell Painting data with contrastive and domain-adversarial learning. *Zenodo*, <https://doi.org/10.5281/zenodo.16361441> (2025).

Acknowledgements

This work was supported by the National Key Research and Development Program of China, grant no. 2020YFA0908700 (J.L.), 2020YFA0908702 (J.L.), and 2024YFA1307703 (S.C.), the National Natural Science Foundation of China grants no. 62473212 (S.C.), 62203236 (S.C.), and 62272246 (J.L.), and the Young Elite Scientists Sponsorship Program by CAST grant no. 2023QNRC001 (S.C.).

Author contributions

S.C. and J.L. conceived the study and supervised the project. C.W.Y., Y.Z., J.F., J.L., and S.C. designed, implemented, and validated cpDistiller. H.H., Z.R., Z.L., S.L., C.Y.Y., and P.L. helped analyze the results. C.W.Y., Y.Z., J.F., J.L., and S.C. wrote the manuscript with inputs from all the authors.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-025-62193-z>.

Correspondence and requests for materials should be addressed to Jian Liu or Shengquan Chen.

Peer review information *Nature Communications* thanks Srijit Seal and the other anonymous reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

Reprints and permissions information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025