

## 1    **Triple-effect correction for Cell Painting data with contrastive and 2    domain-adversarial learning**

3    Chengwei Yan<sup>1,4</sup>, Yu Zhang<sup>2,4</sup>, Jiuxin Feng<sup>1,4</sup>, Heyang Hua<sup>2</sup>, Zhihan Ruan<sup>1</sup>, Zhen Li<sup>3</sup>, Siyu Li<sup>2</sup>,  
4    Chaoyang Yan<sup>1</sup>, Pingjing Li<sup>1</sup>, Jian Liu<sup>1,\*</sup> and Shengquan Chen<sup>2,\*</sup>

5    <sup>1</sup>Centre for Bioinformatics and Intelligent Medicine, Nankai University, Tianjin 300071, China

6    <sup>2</sup> School of Mathematical Sciences and LPMC, Nankai University, Tianjin 300071, China

7    <sup>3</sup> MOE Key Laboratory of Bioinformatics and Bioinformatics Division of BNRIST,  
8    Department of Automation, Tsinghua University, Beijing 100084, China

9    <sup>4</sup> These authors contributed equally

10   \* Corresponding authors: chenshengquan@nankai.edu.cn and jianliu@nankai.edu.cn

11

12

13

### 14    **Abstract**

15   **Cell Painting (CP), as a high-throughput imaging technology, generates extensive cell-  
16   stained imaging data, providing unique morphological insights for biological research.  
17   However, CP data contains three types of technical effects, referred to as triple effects,  
18   including batch effects, gradient-influenced row and column effects (well position effects).  
19   The interaction of various technical effects can obscure true biological signals and  
20   complicate the characterization of CP data, making correction essential for reliable analysis.  
21   Here, we propose cpDistiller, a triple-effect correction method specially designed for CP data,  
22   which leverages a pre-trained segmentation model coupled with a semi-supervised Gaussian  
23   mixture variational autoencoder utilizing contrastive and domain-adversarial learning.  
24   Through extensive qualitative and quantitative experiments across various CP profiles, we  
25   demonstrate that cpDistiller effectively corrects triple effects, especially well position effects,  
26   a challenge that no current methods address, while preserving cellular heterogeneity.  
27   Moreover, cpDistiller effectively captures system-level phenotypic responses to genetic  
28   perturbations and reliably infers gene functions and interactions both when combined  
29   with scRNA-seq data and independently. cpDistiller also excels at identifying gene and  
30   compound targets, which is a critical step in drug discovery and broader biological research.**

31

32 Advanced high-dimensional assay technologies, such as transcriptomics and epigenomics  
33 profiling, offer remarkable depth and breadth in molecular-level biological research<sup>1</sup>. Despite  
34 their strengths, these technologies often focus exclusively on specific molecular changes,  
35 lacking the capability to observe changes at the system level of cell state which involves many  
36 complex and unknown processes. To obtain information at the cellular system level, high-  
37 throughput imaging technologies have been developed to produce useful profiles of cell  
38 phenotypes by imaging stained cells<sup>2–4</sup>. However, these image-based technologies also have  
39 their limitations, as they typically focus on biological processes with known associations or  
40 assumptions, thereby constraining the discovery in existing knowledge<sup>5</sup>. Moreover, traditional  
41 methods that include both high-dimensional assay and image-based technologies are often  
42 constrained by their complexity and high costs. To overcome these issues, the technology,  
43 known as Cell Painting (CP), has been proposed as a solution. Specifically, CP technology  
44 involves staining eight cellular components with six remarkably cheap and easy dyes and  
45 imaging them in five channels on a fluorescence microscope<sup>6</sup>, which is simple to operate and  
46 less costly<sup>7</sup>. Beyond ease of use, CP technology operates on a new paradigm by collecting data  
47 on a large scale without an initial focus on specific hypotheses or known knowledge to reveal  
48 unanticipated biology at play. With these advantages, CP technology has shed light on novel  
49 phenotypes and cellular phenotypic heterogeneity, providing a valuable complement to  
50 genomics<sup>8</sup>. It also has been successfully used to characterize genes<sup>9–11</sup> and compounds<sup>2,11,12</sup> in  
51 several steps of the drug discovery process.

52 As CP technology has developed and been increasingly applied, a large amount of data has  
53 been accumulated<sup>13</sup>. In 2023, scientists established the Joint Undertaking for Morphological  
54 Profiling (JUMP) dataset, a standardized collection of cell-stained images featuring over  
55 116,000 unique compound perturbations and more than 15,000 unique genetic perturbations<sup>14</sup>.  
56 Despite the dataset's substantial analytical potential, technical effects caused by non-biological  
57 factors pose significant challenges in the analysis process. Some of these technical effects,  
58 resulting from variations across different laboratories and different batches in a laboratory have  
59 been observed in the JUMP dataset<sup>15</sup>. In addition to these well-recognized technical effects  
60 observed in most data collection techniques, previous studies have found that CP features  
61 extracted using the conventional tool, CellProfiler<sup>16</sup>, exhibit distinctive well position effects<sup>14</sup>.  
62 Concretely, well position effects arise from the unique design of the CP experiment. In CP  
63 technology, experimental plates are organized into 16 rows and 24 columns, totaling 384 wells,  
64 with each well influenced by both row effects and column effects. We collectively refer to row  
65 effects, column effects, and the effects of different batches as triple effects. The complex and  
66 combined triple effects can lead to deviations from accurate biological profiles, and thus need  
67 to be corrected urgently.

68 However, no methods have been specifically designed to correct triple effects especially  
69 well position effects in CP profiles. Although there exist some batch correction methods  
70 designed for single-cell data, directly applying them to CP data remains challenging for several

71 reasons. First, the characteristics of CP data differ significantly from single-cell data, as CP  
72 data is denser and exhibits lower variability compared to single-cell data<sup>17</sup>. Second, well  
73 position effects in CP data contrast with batch effects in single-cell data, as row or column  
74 effects show a gradient-influenced pattern, where greater differences in row or column numbers  
75 leads to more pronounced effects. Third, the triple effects, especially row and column effects,  
76 are complexly interactive and need to be corrected simultaneously. Some methods, such as  
77 scVI<sup>18</sup>, can correct only one type of technical effects and are constrained to correct multiple  
78 technical effects. Although methods like Harmony<sup>19</sup> model one type of technical effects at a  
79 time and can correct all effects one-by-one, they are unable to simultaneously model triple  
80 effects in CP data.

81 Besides the challenges of correcting triple effects, existing studies that rely on features  
82 extracted by the CellProfiler tool still encounter several unresolved disputes and limitations<sup>14</sup>.  
83 First, it is currently unknown whether the well position effects are inherent to the CP data or  
84 arise from the feature extraction process using CellProfiler. Second, the CellProfiler software,  
85 as a non-end-to-end feature extraction tool, is overly reliant on traditional computer vision  
86 features and requires expert selection, which may overlook certain relevant phenotypic  
87 variation<sup>4</sup>. Third, the CellProfiler software is a non-data-driven feature extraction method,  
88 where the effectiveness of its predefined features is heavily reliant on the characteristics of the  
89 input images and domain knowledge.

90 Here, we demonstrated a one-stop method named cpDistiller for correcting triple effects and  
91 extracting latent patterns in CP data. cpDistiller mainly comprises three modules: the extractor  
92 module for deriving more comprehensive image information, the joint training module for  
93 integrating dual-source features, and the technical correction module for simultaneously  
94 correcting batch, row, and column effects. Specifically, the extractor module, inspired by  
95 transfer learning, employs a pre-trained segmentation model in an end-to-end and data-driven  
96 manner, which is adjusted to extract features from nearly 30 terabytes of raw images. The joint  
97 training module aligns the features extracted by both CellProfiler and the extractor module,  
98 improving the model's ability to better characterize cell-to-cell variation. The technical  
99 correction module employs a semi-supervised Gaussian mixture variational autoencoder  
100 (GMVAE), incorporating contrastive and domain-adversarial learning strategies, to  
101 simultaneously correct technical effects. Based on comprehensive experiments across various  
102 CP profiles, we demonstrated that cpDistiller excelled in both qualitative visualizations and  
103 quantitative metrics, outperforming five baseline methods in single-batch well position effect  
104 correction as well as simultaneous triple-effect correction, all while preserving biological  
105 heterogeneity. Besides, we showcased the extensive capabilities of cpDistiller, including the  
106 ability to integrate more information-rich image features, the support for incremental learning,  
107 and the robustness to various feature selection strategies. Moreover, we emphasized that  
108 cpDistiller effectively captures system-level phenotypic responses to genetic and chemical  
109 perturbations, serving as a powerful tool to complement single-cell RNA sequencing (scRNA-

110 seq) data for uncovering gene functions and relationships. In addition to combination with  
111 scRNA-seq, cpDistiller has the potential to provide unbiased insights into gene associations  
112 independently. Furthermore, by improving the matching of genetic perturbations with their  
113 target genes and enhancing gene-compound similarity assessments, cpDistiller shows strong  
114 promise for accelerating the identification of targets, which is quite valuable in facilitating drug  
115 discovery and various fields of biological research.

## 116 **Results**

### 117 **The overall architecture of cpDistiller**

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

122 cpDistiller processes CP images at a resolution of  $1080 \times 1080$  pixels and extracts  
123 comprehensive features from each well. These features form a matrix, where rows represent  
124 samples (cells, wells or perturbations) and columns correspond to extracted features. We first  
125 extract features from the raw images using CellProfiler, a widely-used but not data-driven  
126 approach, and refer to these features as CellProfiler-based features. Drawing inspiration from  
127 transfer learning, we further develop the extractor module based on an end-to-end pre-trained  
128 segmentation model to automatically extract features, and refer to these features as cpDistiller-  
129 extractor-based features (Methods).

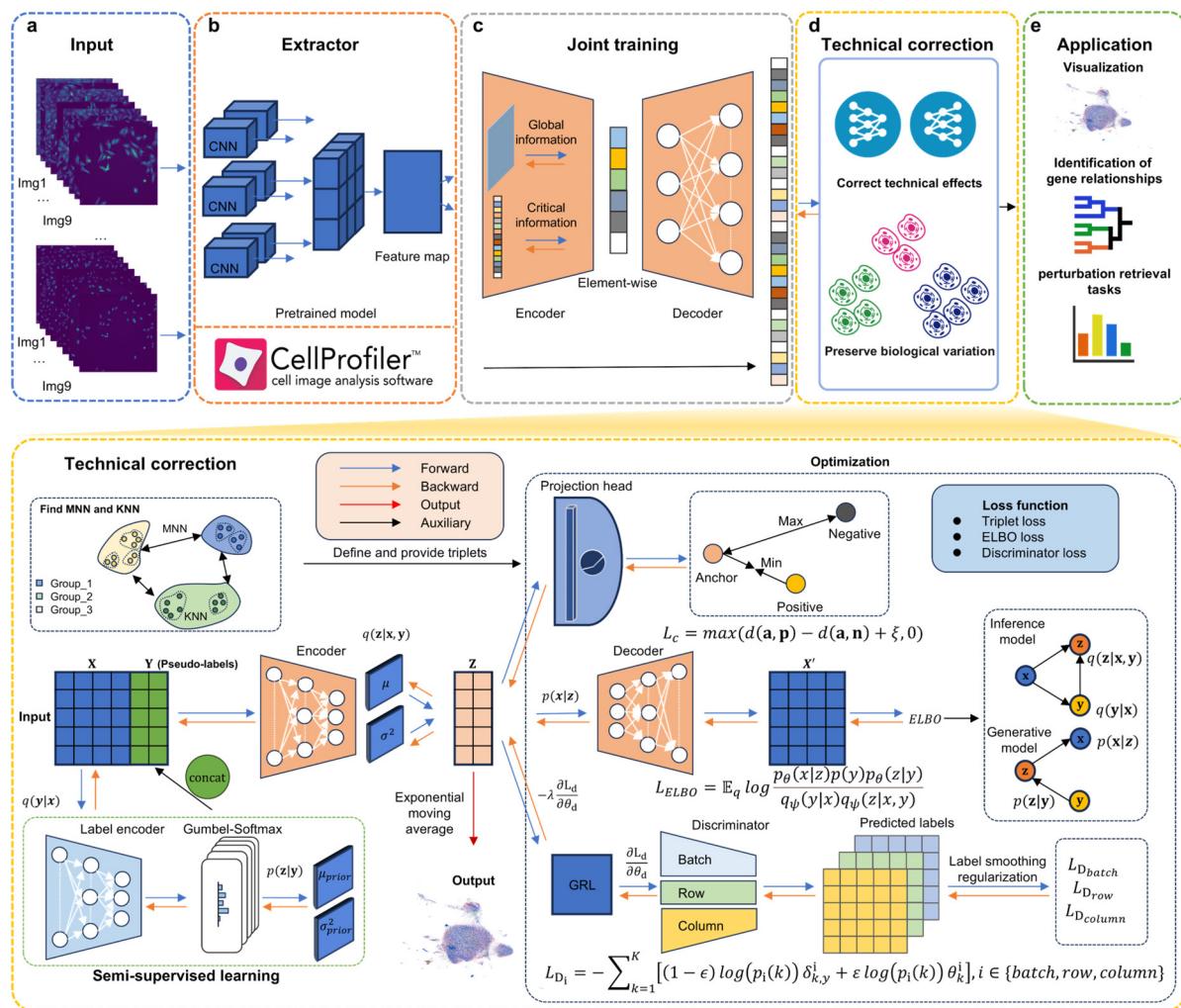
130 Then, the two sets of features are introduced into the joint training module for integration.  
131 The CellProfiler-based features retain unprocessed, while the cpDistiller-extractor-based  
132 features are transformed through an attention mechanism-based encoder, reducing them to a  
133 latent space, and then reconstructed via a decoder. This encoder-decoder structure, applied  
134 exclusively to the cpDistiller-extractor-based features, ensures feature refinement by reducing  
135 noise and enhancing the quality of the representations. The attention mechanism emphasizes  
136 key features across dimensions, facilitating a more effective combination with the CellProfiler-  
137 based features while filtering out irrelevant information (Methods).

138 The technical correction module, with a GMVAE at its core, infers pseudo-labels for each  
139 well in a semi-supervised mode using the features integrated by the joint training module.  
140 These pseudo-labels represent the Gaussian components that characterize the underlying  
141 patterns of each well, capturing the differences among these patterns. By combining these  
142 pseudo-labels along with the integrated features, the technical correction module derives the  
143 latent low-dimensional representations for each well. Besides, we employ both contrastive and  
144 domain-adversarial learning strategies. In contrastive learning, we use  $k$ -nearest neighbors  
145 (KNN) and mutual nearest neighbors (MNN) to construct triplets, applying triplet loss<sup>20</sup> to  
146 restore more accurate nearest-neighbor relationships for each well. To facilitate domain-  
147 adversarial learning and avoid stage-wise training in generative adversarial networks (GANs)<sup>21</sup>,

148 we further apply the gradient reversal layer (GRL)<sup>22</sup> in an end-to-end training approach,  
149 making it harder for discriminators to distinguish data sources and thus removing technical  
150 effects. Additionally, for calculating the discriminator loss, we design a soft label distribution  
151 tailored to the gradient-influenced pattern of the CP data, where each entry in the distribution  
152 vector reflects its proximity to the actual data source for a more accurate representation.

153 Currently, no specialized methods are available for correcting triple effects, especially well  
154 position effects in CP data. However, a recent benchmark study has explored the application of  
155 single-cell batch correction methods as a potential solution<sup>15</sup>. In light of this, we compared  
156 cpDistiller with methods that excel at removing batch effects in single-cell data, including  
157 Seurat v5<sup>23</sup>, Harmony<sup>19</sup>, Scanorama<sup>24</sup>, scVI<sup>18</sup> and scDML<sup>25</sup>. Among them, Seurat v5, Harmony,  
158 and Scanorama can iteratively correct triple effects by sequentially incorporating different  
159 technical labels (removal of batch, row, and column effects one after another), as they are based  
160 on low-dimensional representations, which allow for flexible, stepwise corrections. In contrast,  
161 scVI and scDML, directly model the original high-dimensional inputs to obtain low-  
162 dimensional representations, making them less suitable for iterative correction and limiting  
163 them to address only one type of technical effects (one of batch, row, or column effects). Due  
164 to the lack of metrics for evaluating the correction of technical effects in CP profiles, we use  
165 single-cell analysis metrics to assess the capacity of different methods to remove technical  
166 effects while preserving biological variation. Metrics such as average silhouette width (ASW)<sup>26</sup>,  
167 technic average silhouette width (tASW)<sup>27</sup> and graph connectivity<sup>27</sup> are used to measure the  
168 model's ability to remove technical effects, while isolated label silhouette<sup>27</sup>, isolated label F1<sup>27</sup>,  
169 perturbation average silhouette width (pASW)<sup>27</sup>, and normalized mutual information (NMI)<sup>27</sup>  
170 are used to evaluate the characterization of heterogeneity (Methods).

171



**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 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. Additionally, three discriminators are designed to predict batch, row, and column labels, and the technical correction module use 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.

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

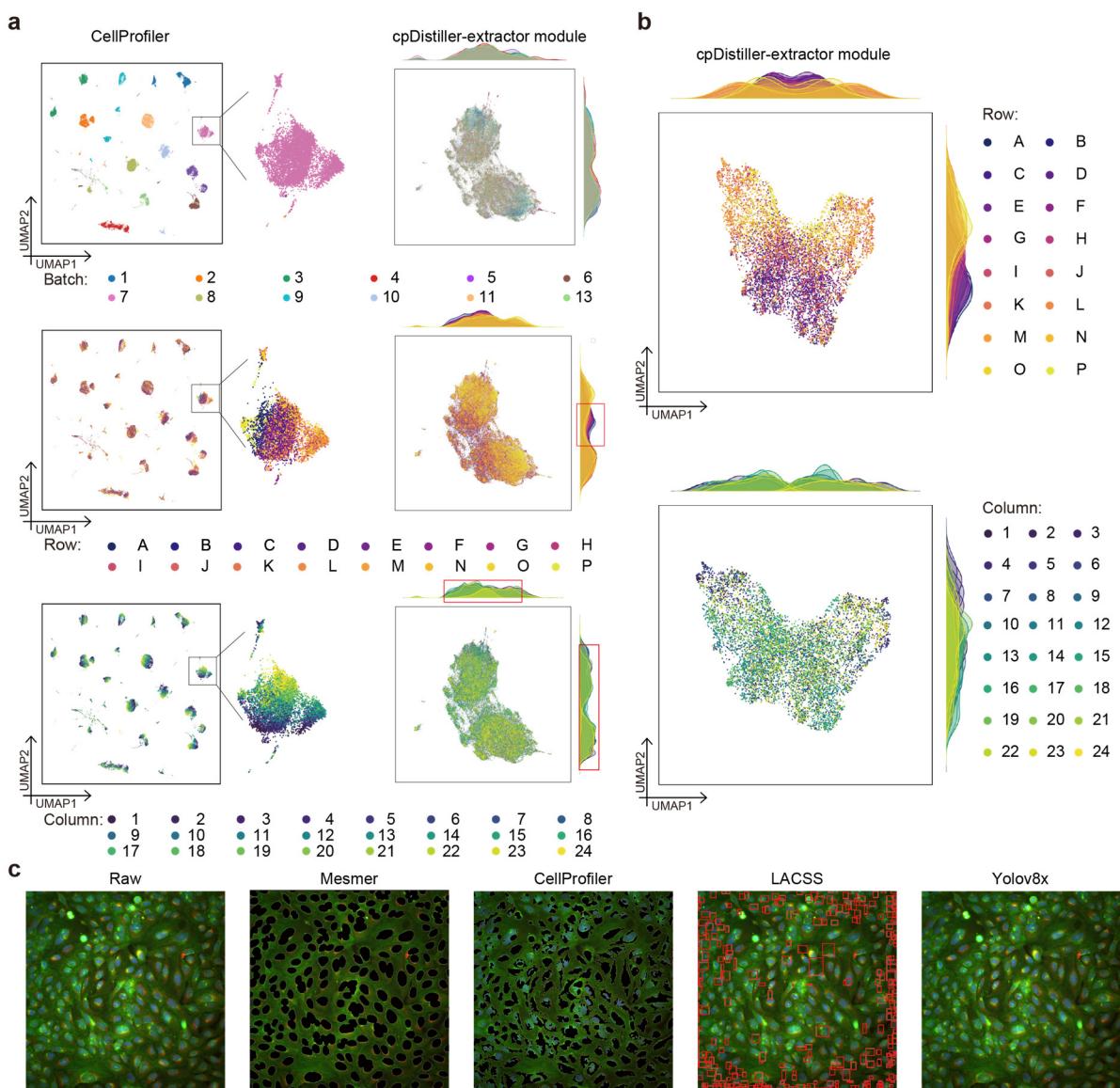
175 Previous studies have identified three types of technical effects in CP data: batch, row, and  
176 column effects<sup>28</sup>. Notably, row and column effects suggest the novel well position effects.  
177 Following the prior study<sup>14</sup>, we employed CellProfiler to obtain 1,446-dimensional features  
178 (CellProfiler-based features) from open reading frame (ORF) overexpression dataset in  
179 cpg0016, and then performed uniform manifold approximation and projection (UMAP)<sup>29</sup> for  
180 these features. From the UMAP visualization, we observed distinct clustering patterns  
181 corresponding to different batches, rows, and columns (Fig. 2a). Batch effects are evident as  
182 patterns from different batches cluster into separate groups, while row and column effects are  
183 displayed as color gradients in the visualization, highlighting the presence of triple effects (Fig.  
184 2a). To investigate whether row and column effects are evident within a single batch, we further  
185 observed the CellProfiler-based features from Batch\_7. The UMAP visualizations revealed  
186 color gradients, indicating a gradient-influenced pattern where more distant rows or columns  
187 exhibit more pronounced effects, a phenomenon consistently observed across 12 batches (Fig.  
188 2a and Supplementary Figs. 1-3).

189 To evaluate whether the triple effects are due to inherent characteristics of CP data or biases  
190 in the feature extraction process with CellProfiler, we developed the cpDistiller-extractor  
191 module, based on the pre-trained segmentation model Mesmer<sup>30</sup>, to extract deep-learning  
192 features (cpDistiller-extractor-based features). As illustrated in Fig. 2a, UMAP visualizations  
193 of the cpDistiller-extractor-based features from 12 batches display slight batch effects and clear,  
194 consistent effects related to row and column (Fig. 2a). We further plotted and observed the  
195 density distributions of the scatter points along the x and y axes in the UMAP visualizations.  
196 These distributions revealed distinct peaks corresponding to different batches, indicating the  
197 presence of batch effects (Fig. 2a). Additionally, we observed that the distributions for rows  
198 and columns were non-overlapping, with row-specific and column-specific peaks highlighted  
199 in red boxes, suggesting the existence of row and column effects (Fig. 2a). Furthermore, UMAP  
200 visualizations of the cpDistiller-extractor-based features from Batch\_7 demonstrated clear  
201 evidence of row and column effects (Fig. 2b), a pattern that was consistently observed across  
202 all batches (Supplementary Figs. 4 and 5).

203 To investigate whether the cpDistiller-extractor module can capture more precise  
204 information regarding cell positions and shapes from CP images, we compared its segmentation  
205 results with those obtained from CellProfiler. As illustrated in Fig. 2c, the Mesmer model,  
206 which serves as the core of cpDistiller-extractor module, outperformed CellProfiler in  
207 segmentation capabilities, particularly in capturing details of cell morphology (Supplementary  
208 Text 1 and Supplementary Fig. 6). Before adopting Mesmer, we also assessed the widely used  
209 LACSS<sup>31</sup> and YOLOv8<sup>32</sup> models as alternatives. However, neither was able to accurately  
210 locate cell positions (Fig. 2c and Supplementary Fig. 7). We further tested the YOLOv8 model  
211 with various parameters, but the results confirmed that it remained ineffective for cell detection

212 and segmentation in CP data (Supplementary Figs. 8-13). Details regarding the pre-trained  
213 model settings can be found in Supplementary Text 2.

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



**Fig. 2 | CellProfiler and cpDistiller-extractor module confirm triple effects.** **a**, The UMAP visualizations of embeddings obtained by the CellProfiler (left) and cpDistiller-extractor module (right), colored by batch, row, and column, respectively. **b**, The UMAP visualizations of embeddings obtained by 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 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 fail to perform cell segmentation from CP images.

217 **cpDistiller effectively corrects well position effects and preserves cellular phenotypic  
218 heterogeneity**

219 We first conducted experiments on the cpg0016 ORF profiles derived from the U2OS cell type  
220 in the JUMP dataset to demonstrate the effectiveness of cpDistiller. The cpg0016 ORF profiles  
221 consisted of 12 batches, and we used Batch\_1 as an example. In the raw data, UMAP  
222 visualization reveals a gradient-influenced pattern, with more distant rows or columns  
223 exhibiting more pronounced effects (Fig. 3a,b and results from other methods and batches in  
224 Supplementary Figs. 14-25). We then visualized the low-dimensional embeddings obtained by  
225 cpDistiller and observed that the gradient-influenced pattern was significantly reduced. The  
226 data distribution across different rows was more uniform, and the previously pronounced  
227 column effects, particularly the large differences between low-numbered columns (e.g.,  
228 Column\_1) and high-numbered columns (e.g., Column\_24), were well-mixed and greatly  
229 corrected (Fig. 3a,b). This demonstrated that cpDistiller can effectively correct both row and  
230 column effects. Moreover, overexpression reagents and compounds are theoretically expected  
231 to induce significant differences in cell phenotypes between negative and positive controls.  
232 Specifically, by utilizing UMAP to visualize the impact of various perturbations, cpDistiller  
233 successfully differentiated positive controls, including JCP2022\_037716, JCP2022\_035095,  
234 and JCP2022\_012818 from negative controls across various batches (Supplementary Figs. 14-  
235 25). In addition to analyzing the controls, we also explored the ORF perturbations in treatment.  
236 cpDistiller successfully identified perturbations caused by 12 overexpression reagents in  
237 treatment and preserved their unique stimulatory effects, which were obscured in the raw data  
238 due to well position effects (Fig. 3c). In contrast, UMAP visualizations and hierarchical  
239 clustering results suggest that methods like scDML and scVI, which are limited to correcting  
240 only one type of technical effects, struggle to preserve biological variation while also failing to  
241 effectively correct well position effects (Supplementary Figs. 14-25 and Supplementary Fig.  
242 26a). Although these methods may achieve partial mixing across rows and columns, they still  
243 fail to remove most non-biological noise, such as the clear striping pattern, particularly  
244 noticeable in Batch\_5 and Batch\_8 (Supplementary Figs. 18 and 21). While methods like  
245 Harmony, Seurat v5, and Scanorama can iteratively correct well position effects, they also fall  
246 short in preserving the true biological variation of ORF perturbations (Supplementary Fig. 26a).

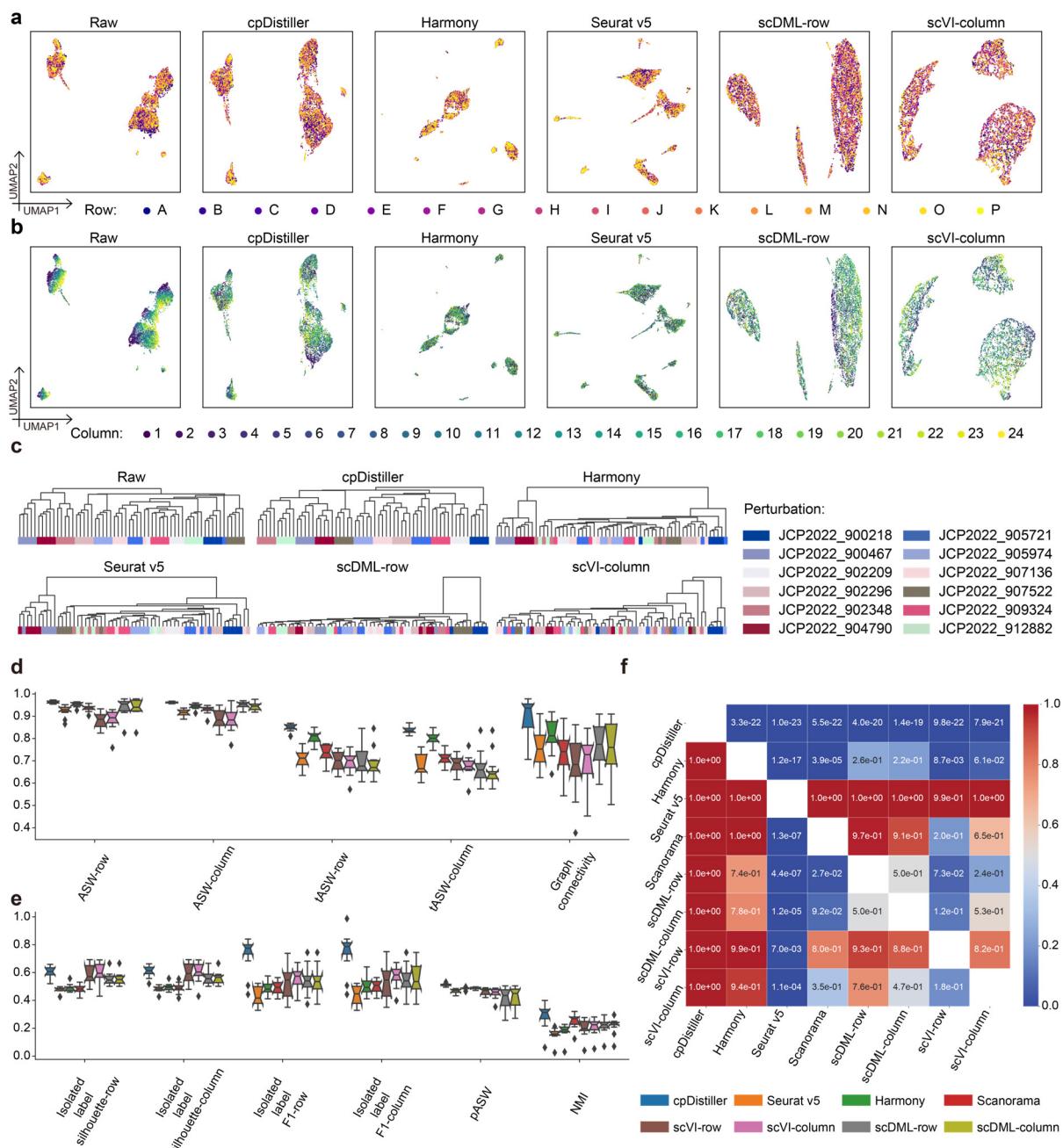
247 To quantitatively demonstrate the advantages of cpDistiller in correcting well position  
248 effects, we further conducted experiments in 12 batches and assessed the performance by ASW,  
249 tASW, and graph connectivity as suggested in refs.<sup>15,26,27</sup>. ASW and tASW were used to assess  
250 the extent of data mixing across different technical labels. Graph connectivity assumed that,  
251 once technical effects were corrected, data with the same biological labels, specifically  
252 perturbation labels, should cluster more tightly. cpDistiller outperformed other baseline  
253 methods in all metrics and consistently demonstrated stable performance across 12 batches (Fig.  
254 3d), indicating superior performance in both uniformly mixing CP profiles and accurately  
255 reflecting cell phenotype differences induced by various perturbations.

256 We further assessed the preservation of biological variation via metrics such as isolated label  
257 scores (including isolated label F1 and isolated label silhouette), pASW, and NMI. Isolated  
258 labels refer to perturbation labels with the least number of the technical effects. Here, the  
259 perturbation labels computed from isolated labels all belonged to positive controls. Higher  
260 isolated label scores reflect a greater impact of positive controls on cell phenotypes compared  
261 to others. cpDistiller excelled in isolated label scores, particularly achieving the highest isolated  
262 label F1 score among all baseline methods (Fig. 3e). Additionally, for the clustering metrics of  
263 pASW and NMI, which were used to assess the clustering quality of replicated experiments  
264 involving the same types of perturbations as suggested in refs.<sup>15,27</sup>, cpDistiller also surpassed  
265 other methods (Fig. 3e). Harmony has demonstrated excellent performance in removing batch  
266 effects for CP data in recent benchmark analysis<sup>15</sup> and also emerges as the most effective  
267 baseline method for correcting row and column effects in our study. To further demonstrate the  
268 significant advantages of cpDistiller, we used one-sided paired Wilcoxon signed-rank tests to  
269 compare the performance of different methods across 12 batches. The *P*-values highlighted  
270 cpDistiller's significant advantages in removing well position effects and preserving biological  
271 variation across batches, outperforming all baseline methods (Fig. 3f).

272 Besides, in our previous experiments with baseline methods like Seurat v5, Harmony, and  
273 Scanorama, which can iteratively remove well position effects, we initially corrected rows  
274 before columns. To investigate whether the performance of these methods was influenced by  
275 the order of correction, we also applied the reverse order, namely correcting columns before  
276 rows. The box plots across 12 batches showed no significant difference in performance when  
277 the same method was applied with different correction orders (Supplementary Fig. 27a).  
278 However, when we used one-sided paired Wilcoxon signed-rank tests for a more detailed  
279 quantification of performance differences, we observed that, although the differences were not  
280 statistically significant, the *P*-values indicated a trend where correcting rows before columns  
281 generally yielded better results than the reverse order for all baseline methods (Supplementary  
282 Fig. 27b).

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

286



**Fig. 3 | cpDistiller can effectively correct well position effects while preserving biological variation for ORF profiles across 12 batches. a,b,** UMAP visualizations of embeddings obtained by different methods in Batch\_1 colored by row (a) and column (b). **c,** Dendograms illustrate the clustering of ORF perturbations induced by 12 reagents in treatment, graphically rendered based on low-dimensional representations generated by different methods. **d,e,** Quantitative evaluation of the performance of different methods on ORF profiles across 12 batches for correcting well position effects while preserving biological variation. The top plot focuses on technical correction capabilities (d), while the bottom plot highlights biological preservation abilities (e). Since the well position effects encompass both row and column effects, the metrics of ASW, tASW, isolated label scores yield two results, 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. **f,** 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.

288 **cpDistiller enables effective and simultaneous correction of triple effects**

289 After integrating ORF profiles from 12 batches and visualizing them using UMAP, we  
290 observed batch effects in addition to well position effects (Fig. 4a). To study biological-process-  
291 related mechanisms of action based on feature similarity<sup>33</sup>, it is necessary to remove batch  
292 effects to ensure reliable comparisons across batches. For competently removing batch effects,  
293 in addition to the discriminators for row and column labels, we also created an additional  
294 discriminator to identify batch labels, encouraging cpDistiller to learn low-dimensional  
295 representations that are indistinguishable with respect to the sources of batch labels, thereby  
296 removing batch effects (Methods). At the same time, we additionally considered KNN intra  
297 batches and MNN inter batches to construct triplets, using the contrastive learning technique  
298 to restore more accurate nearest-neighbor relationships for each well (Methods). We next  
299 verified that cpDistiller can satisfactorily correct triple effects.

300 First, we qualitatively compared the performance of different methods using UMAP  
301 visualizations, which showed that cpDistiller achieved successful mixing of data across batches,  
302 rows, and columns simultaneously (Fig. 4a and Supplementary Fig. 28). Moreover, to  
303 qualitatively assess how well different methods preserve the specificity of perturbations,  
304 compound perturbations shared across 12 batches on the target plates provided a reliable  
305 measure<sup>15</sup>. UMAP visualizations showed that cpDistiller ensured perturbations caused by the  
306 same compounds had a tendency to be clustered (Fig. 4a). To be specific, cpDistiller  
307 successfully identified multiple perturbations in target plates, including JCP2022\_010404,  
308 JCP2022\_000794, and JCP2022\_047545. However, methods that are limited to correcting only  
309 one type of triple effects, such as scDML and scVI, failed to correct well position effects,  
310 resulting in the persistence of striping pattern for rows and columns (Fig. 4a). On the other  
311 hand, methods capable of iteratively correcting triple effects, including Scanorama and Seurat  
312 v5, tended to overcorrect, leading to fail to preserve true biological signals (Supplementary Fig.  
313 28). To further evaluate if cpDistiller effectively maintained biological variation during triple-  
314 effect correction, we reanalyzed the ORF perturbations that were assessed during the correction  
315 of well position effects. As shown in Supplementary Fig. 26b, cpDistiller can still identify these  
316 differential perturbations, whereas other methods cannot.

317 Moreover, we quantitatively evaluated the performance of different methods. cpDistiller  
318 generally exhibited superior performance in correcting triple effects compared to baseline  
319 methods (Fig. 4b). Although Scanorama performed better than cpDistiller in tASW, it scored  
320 lowest in graph connectivity. It indicated that while the data was well mixed across batches,  
321 rows, and columns, it lost the specificity of the perturbations. Additionally, when assessing the  
322 preservation of biological variation through clustering metrics, cpDistiller consistently  
323 performed best across all metrics (Fig. 4b). Besides, we noted the balance between technical  
324 effect removal and biological conservation. We found that all baseline methods struggled to  
325 achieve the equilibrium (Fig. 4c). Specifically, scVI and Harmony both performed comparably  
326 well among the baseline methods. However, scVI tended to favor the preservation of biological

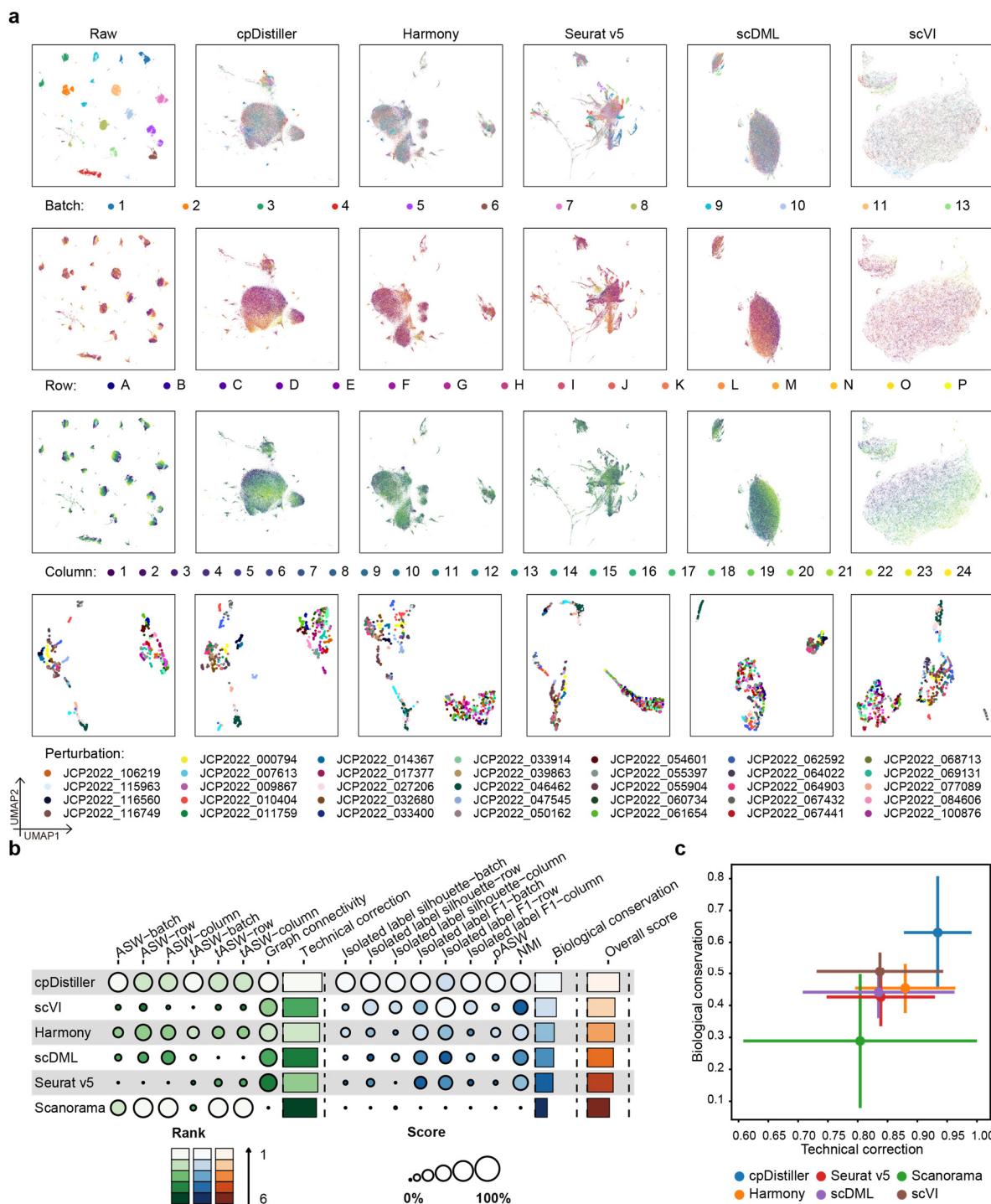
327 variation, while Harmony focused more on technical correction (Fig. 4c). In contrast,  
328 cpDistiller balanced both aspects, offering superior overall performance in both removing triple  
329 effects and preserving cellular phenotypic heterogeneity (Fig. 4c).

330 Furthermore, to comprehensively evaluate the advantages of cpDistiller, we conducted a  
331 series of experiments to demonstrate its ability to leverage more information-rich image  
332 features for enhanced performance, support incremental learning for ongoing studies, and  
333 maintain robustness to feature selection. To investigate the impact of cpDistiller-extractor-  
334 based features on the performance of cpDistiller, we conducted experiments using only  
335 CellProfiler-based features. We found that cpDistiller still outperformed baseline methods in  
336 correcting technical effects while preserving biological variation, even when utilizing only  
337 CellProfiler-based features, as the baseline methods did (Supplementary Text 3 and  
338 Supplementary Fig. 29a-c). Moreover, with leveraging cpDistiller-extractor-based features,  
339 cpDistiller can exploit a broader range of image characteristics from raw images and lead to  
340 further performance improvements (Supplementary Text 4 and Supplementary Fig. 30).

341 Besides, cpDistiller supports incremental learning (Supplementary Text 5). Methods like  
342 Harmony and Scanorama require reprocessing and realigning the entire dataset to integrate new  
343 data, which can be cumbersome, especially when working with large public datasets<sup>15</sup>. These  
344 processes often involve modifying existing representations, which can disrupt the continuity of  
345 prior analytical results. In contrast, cpDistiller can leverage the model parameters learned from  
346 previous tasks, enabling the direct integration of new data without the necessity of realignment  
347 (Supplementary Fig. 31a). These flexibilities are particularly beneficial for ongoing studies,  
348 enabling seamless integration while preserving the integrity of previous analyses.

349 Additionally, cpDistiller demonstrates robustness to feature selection (Supplementary Text  
350 6). We first extracted features with dimensions of 4,752 from fluorescence images and 7,638  
351 when combined with brightfield images by using the CellProfiler software (Methods). We then  
352 utilized these two types of features to carry out two key tasks: correcting well position effects  
353 in single batch and simultaneously correcting triple effects. The results showed that even when  
354 dealing with high-dimensional features that may contain substantial redundant features and  
355 noise, cpDistiller consistently outperformed baseline methods in correcting technical effects  
356 and preserving biological variation (Supplementary Fig. 31b,c).

357 Overall, cpDistiller achieved a satisfactory balance between correcting triple effects and  
358 preserving biological variation, and demonstrated strong capabilities in information-rich  
359 feature extraction, incremental learning, and robust feature selection.



**Fig. 4 | cpDistiller can simultaneously correct triple effects while preserving cellular phenotypic heterogeneity.** **a**, UMAP visualizations of embeddings obtained by different methods across 12 batches in ORF profiles, colored by batch, row, column, and perturbation, respectively. The embeddings are consistent with all visualizations, and the selected 34 perturbations are shown. **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. Due to triple effects, which encompass batch, row and column effects, the metrics of ASW, tASW, isolated label scores yield three results, respectively. **c**, Scatter plot summarizes overall performance on triple-effect correction in ORF profiles. The x-axis represents the technical correction score, while the y-axis shows the biological conservation score. Each point corresponds to the average value per method, with error bars indicating the standard deviation.

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

362 Inferring gene functions and interactions, as a fundamental step in many areas of biological  
363 research, often relies on various types of sequencing data such as scRNA-seq data and  
364 chromatin immunoprecipitation sequencing (ChIP-seq) data<sup>34</sup>. This inferring task is associated  
365 with numerous unknown and highly intricate biological processes, but the sequencing methods  
366 tend to focus on only a limited number of interesting molecular-level changes, making obtaining  
367 thorough and accurate inferences challenging<sup>35</sup>. In contrast, we demonstrated that cpDistiller  
368 offers comprehensive system-level phenotypic characteristics under genetic perturbations and  
369 has the potential to integrate molecular-level RNA data for revealing gene functions and  
370 relationships.

371 To demonstrate the capability of the embeddings learned by cpDistiller in preserving  
372 system-level phenotypic characteristics of perturbation heterogeneity, we applied cpDistiller to  
373 the controls in the ORF data from the cpg0016 dataset. Specifically, we focus on controls,  
374 including the positive controls JCP2022\_037716, JCP2022\_012818, and JCP2022\_035095,  
375 and the negative control JCP2022\_915131. The positive controls are expected to produce  
376 noticeable phenotypic changes, while the negative control should induce minimal changes<sup>14</sup>.  
377 Theoretically, phenotypes under identical genetic perturbation in different wells or different  
378 plates should yield consistent patterns. Therefore, we evaluated the similarity of identical  
379 positive and negative controls duplicated across different wells and plates within a batch, using  
380 embeddings generated by cpDistiller and baseline methods, respectively. Taking Batch\_4 for  
381 example, after performing hierarchical clustering based on cosine similarity (Fig. 5a), we found  
382 that the baseline methods failed to cluster identical treatment sets. In contrast, cpDistiller  
383 successfully grouped duplicated positive and negative controls from different wells, as well as  
384 revealing treatment-specific embeddings. These results illustrated that cpDistiller can  
385 effectively capture phenotypic signals across various genetic perturbations.

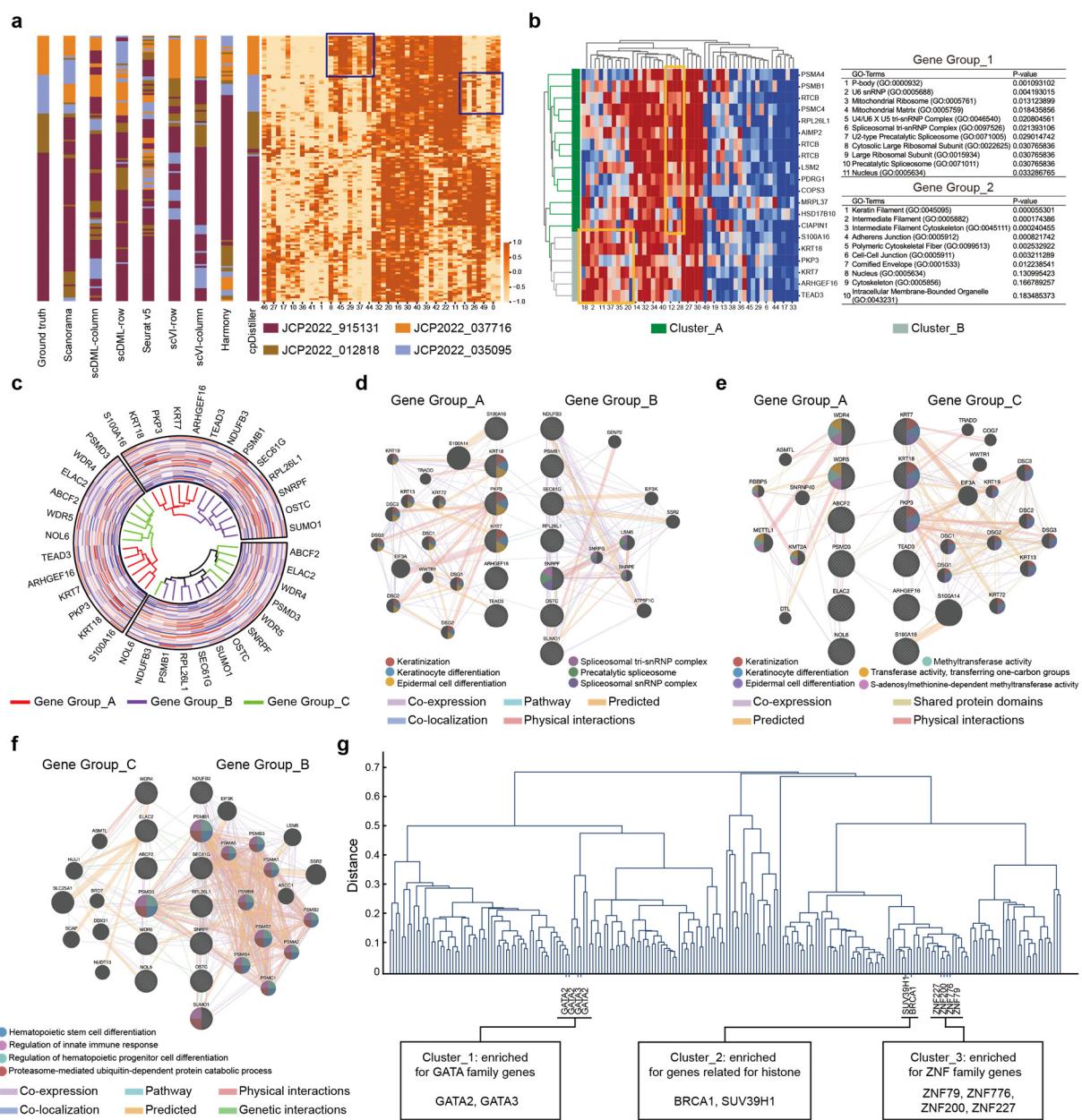
386 To elucidate cpDistiller's capability to integrate with scRNA-seq data for inferring gene  
387 functions, we used the embeddings obtained by cpDistiller to analyze genetic perturbation data  
388 from the ORF dataset. Given that the ORF dataset encompasses large-scale perturbations  
389 targeting 12,602 genes, we first employed the ARCHS4 tool<sup>36</sup> to establish gene groups based  
390 on the similarities of their scRNA-seq data, where genes with highly similar expression patterns  
391 were grouped together (Methods). Next, we calculated the Euclidean distance between the  
392 embeddings from cpDistiller for the individual genes in gene groups and then performed  
393 hierarchical clustering based on these distances. For example, we found gene Group\_1 and  
394 gene Group\_2 are highly separated into distinct clusters (Fig. 5b). Gene Group\_1 is enriched  
395 in Cluster\_A, while gene Group\_2 is enriched in Cluster\_B, indicating that cpDistiller has  
396 learned group-specific embeddings, resulting in significantly distinct morphological  
397 embeddings for genes in Group\_1 and Group\_2. To elucidate that the morphological  
398 embeddings from cpDistiller have the capacity to reveal gene functions, we conducted cellular  
399 component (CC) analyses via gene ontology (GO) enrichment for the genes in Group\_1 and

400 Group\_2, respectively (Methods). We found that the top three significantly enriched cellular  
401 components for the genes in Group\_1 were all linked with the function of RNA processing,  
402 including 'P-bodies', 'U6 snRNP', and 'Mitochondrial Ribosome'<sup>37-39</sup> (Fig. 5b). By contrast, the  
403 top three significantly enriched cellular components for the genes in Group\_2 were all involved  
404 in the function of forming structural elements, including 'keratin filament', 'intermediate  
405 filament', and 'intermediate filament cytoskeleton'<sup>40</sup>. The CC analysis results revealed that the  
406 genes in Group\_1 and the genes in Group\_2 are associated with distinct functions, which is  
407 consistent with the cluster categorization where genes in Group\_1 and Group\_2 have distinct  
408 embeddings learned by cpDistiller.

409 To demonstrate cpDistiller's effectiveness in elucidating gene relationships when combined  
410 with scRNA-seq data, we further analyzed the embeddings obtained by cpDistiller for genes in  
411 gene groups through multiple types of biological analyses. We conducted hierarchical  
412 clustering analysis for individual genes in gene groups using embeddings obtained by  
413 cpDistiller and found that some gene groups formed distinct clusters. For example, as  
414 illustrated in Fig. 5c, gene Group\_A and Group\_B, as well as gene Group\_A and Group\_C, are  
415 distinguishable into different clusters. However, gene Groups\_B and Group\_C are not easily  
416 separated. These results indicated that although genes within different groups have dissimilar  
417 scRNA-seq expression patterns, they may have indistinguishable phenotypic embeddings  
418 obtained from cpDistiller. To elucidate that the embeddings from cpDistiller can illustrate gene  
419 relationships, we utilized the GeneMANIA tool<sup>41</sup> to construct gene networks to visualize gene  
420 relationships based on large and diverse databases<sup>41</sup> (Methods). We found that the correlation  
421 in the network for genes in Group\_A and Group\_B is minimal (Fig. 5d), as well as the  
422 correlation for genes in Group\_A and Group\_C (Fig. 5e). By contrast, genes in Group\_B and  
423 Group\_C exhibit closer relationships (Fig. 5f). These findings aligned with the clustering  
424 results from cpDistiller, where genes in Group\_A and Group\_C, as well as those in Group\_A  
425 and Group\_B, have significantly different embeddings, while those in Group\_B and Group\_C  
426 share similar embeddings. These results indicate that genes with similar embeddings by  
427 cpDistiller tend to have closer relationships, illuminating that cpDistiller can uncover gene  
428 relationships. Moreover, we used GeneMANIA to predict the significantly enriched gene  
429 functions within different groups and found that functional enrichment results consistently  
430 aligned with the clustering results (Fig. 5d-f). These confirmed that analyzing the embeddings  
431 obtained by cpDistiller has the potential to infer gene functions. In conclusion, these results  
432 demonstrated the phenotypic embeddings obtained by cpDistiller effectively capture gene  
433 functions and relationships, highlighting its potential as a valuable tool for exploring gene  
434 interactions.

435 In addition to combination with scRNA-seq data, we further demonstrated that cpDistiller  
436 alone is also capable of uncovering gene relationships. We first applied cpDistiller to each batch  
437 in the ORF dataset, with each batch containing approximately 2,000 genetic perturbations.  
438 Since many genetic perturbations did not induce significant phenotypic changes, we then

439 screened for active genes that exhibited substantial phenotypic changes compared to controls  
440 in each batch (Methods). We calculated the similarity of the embeddings from cpDistiller for  
441 these active genes, and then performed hierarchical clustering based on the similarity to  
442 identify significantly clustered genes (Methods). As illustrated in Fig. 5g, taking Batch\_1 as an  
443 example, we found that *SUV39H1* and *BRCA1* are grouped together. Both *SUV39H1* and  
444 *BRCA1* exhibit co-expression as well as various connections<sup>42</sup> and they are associated with  
445 condensed chromosomes, histone H3-K9 methylation, and related chromosomal  
446 modifications<sup>43</sup>. Additionally, genes from the zinc finger protein family which share protein  
447 domains<sup>44,45</sup>, including *ZNF227*, *ZNF200*, *ZNF776*, and *ZNF79*, are successfully clustered  
448 together based on the embeddings obtained by cpDistiller. Similarly, *GATA2* and *GATA3*,  
449 members of the *GATA* transcription factor family crucial for cell differentiation and  
450 development, which share protein domains<sup>46</sup> and exhibit co-expression<sup>44,45</sup>, are also found to  
451 be clustered together. These results confirmed that analyzing the embeddings learned by  
452 cpDistiller alone also has the potential to reveal gene relationships.



**Fig. 5 | Analyses of gene functions and gene relationships.** **a**, Dendograms illustrate the clustering of duplicates of controls, graphically rendered based on low-dimensional representations generated by different methods. The additional heatmap visualizes the 50-dimensional embeddings obtained by cpDistiller, highlighting the clustering patterns of controls. **b**, Hierarchical clustering and heatmap of the 50-dimensional embeddings obtained by cpDistiller for gene Group\_1 and gene Group\_2. Detailed views show the cellular compound analysis of GO term analysis for gene Group\_1 and gene Group\_2. **c**, Circular dendograms are performed for genes in different gen groups, graphically rendered based on embedding 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.

455 **cpDistiller can facilitate the search for gene and compound targets**

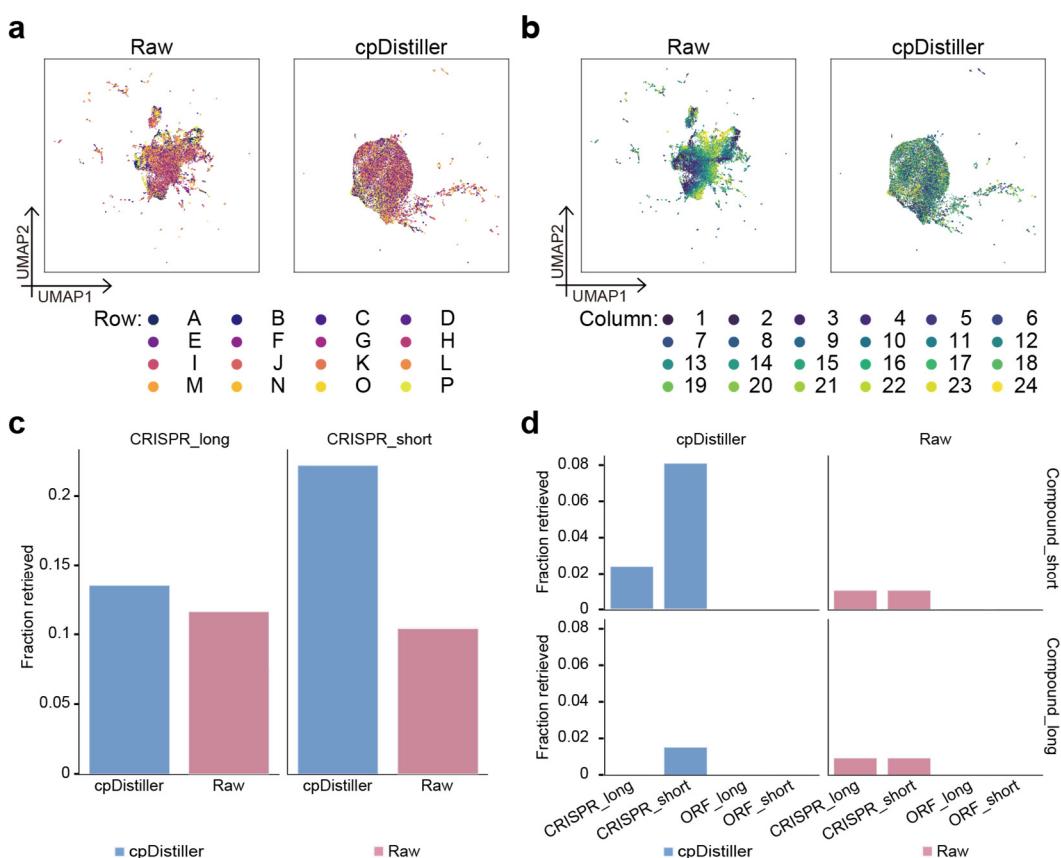
456 The JUMP dataset is primarily generated using U2OS cells, with all ORF experimental data in  
457 the cpg0016 dataset conducted on this cell type. However, since biological research often  
458 extends beyond U2OS cells, the JUMP dataset includes a pilot dataset, cpg0000, conducted on  
459 A549 cells in a single batch. The cpg0000 dataset provides paired genetic and compound  
460 perturbations targeting the same gene, offering substantial potential for uncovering biological  
461 targets<sup>47</sup>. Leveraging this design, researchers have established simulated tasks to retrieve gene  
462 and compound targets using features extracted by CellProfiler<sup>47</sup>. However, as shown in Fig.  
463 6a,b, when we performed UMAP visualization using CellProfiler-based features for A549 cells,  
464 we observed noticeable row and column effects, which could obscure real biological signals.  
465 To address this, we applied cpDistiller to correct both row and column effects. As illustrated in  
466 Fig. 6a,b, the data distributions across different rows and columns were more uniform.

467 We further demonstrated that the cpDistiller-derived embeddings facilitate the identification  
468 of gene and compound targets. Specifically, the cpg0000 dataset includes genetic and  
469 compound perturbations for 160 genes in A549 cells, conducted at both long and short time  
470 points. Each gene is perturbed through one ORF treatment, two gene knockouts by Clustered  
471 Regularly Interspaced Short Palindromic Repeats (CRISPR) guides, and two compound  
472 experiments<sup>47</sup>. Using sister CRISPR guides targeting the same gene, researchers designed  
473 retrieval tasks to simulate gene target identification, calculating the fraction retrieved score to  
474 evaluate retrieval performance<sup>47</sup>. Following their workflows, we calculated the fraction  
475 retrieved scores using the embeddings from cpDistiller and compared the scores to those  
476 generated using CellProfiler-based features (Methods). As shown in Fig. 6c, using CellProfiler-  
477 based features, the fraction retrieved score for retrieving sister CRISPR guides at long time  
478 points (long-time CRISPR retrieval) is 0.123 and the score for retrieving sister CRISPR guides  
479 at short time points (short-time CRISPR retrieval) is 0.11. In contrast, with embeddings  
480 obtained by cpDistiller, the fraction retrieved scores improve to 0.139 for long-time CRISPR  
481 retrieval and 0.235 for short-time CRISPR retrieval. This represented a notable increase,  
482 particularly for short-time retrieval tasks, where cpDistiller achieves over a two-fold increase  
483 compared to CellProfiler-based features. These results highlighted the pivotal role of  
484 cpDistiller in enhancing the accuracy of gene retrieval tasks, making it a valuable tool for  
485 identifying genes involved in similar processes and uncovering critical gene relationships.

486 In addition to retrieving sister perturbations, researchers also conducted cross-modality  
487 gene-compound retrieval tasks to simulate the identification of compound targets<sup>47</sup>. Concretely,  
488 they searched for compounds that produced similar effects on cell morphology as the query  
489 gene using CellProfiler-based features. They calculated the fraction retrieved scores to evaluate  
490 retrieval performance. Following their methodology, we calculated the fraction retrieved scores  
491 using embeddings generated by cpDistiller to demonstrate that cpDistiller can improve retrieval  
492 results. As shown in Fig. 6d, when using CellProfiler-based features for compound  
493 perturbations, the fraction retrieved scores are nearly 0.008 for retrieving all compound-

494 CRISPR pairs. In contrast, when using embeddings extracted by cpDistiller for compound  
495 perturbations at both long and short time points, the fraction retrieved scores for retrieving most  
496 compound-CRISPR pairs over a two-fold increase (Fig. 6d). Given the complexity and  
497 importance of gene-compound retrieval tasks, even slight improvements hold substantial  
498 value<sup>47</sup>. These results demonstrated that cpDistiller achieved notable improvements in gene-  
499 compound retrieval tasks compared to uncorrected CellProfiler-based features, highlighting its  
500 potential to advance the discovery and evaluation of novel chemical compounds.

501 In conclusion, cpDistiller satisfactorily corrected well position effects and obtained  
502 embeddings that offered clear advantages in exploring gene and compound targets, providing  
503 valuable insights for biological research and drug discovery.



**Fig. 6 | The performance of cpDistiller in retrieval tasks.** **a,b**, UMAP visualization for the A549 cells using the features extracted by CellProfiler and embeddings learned by cpDistiller, colored by row (**a**) and column (**b**). **c**, Fraction retrieved scores for retrieving sister CRISPR guides in A549 cells using embeddings of cpDistiller and CellProfiler-based features, 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 learned by cpDistiller and CellProfiler-based features, respectively. Short (\_short) and long (\_long) time points mean the experimental conditions.

505 **Discussion**

506 We developed cpDistiller, a method specifically designed to correct triple effects, particularly  
507 well position effects in Cell Painting (CP) data. cpDistiller leverages raw CP images by  
508 integrating high-level features via pre-trained segmentation model with handcrafted features  
509 from traditional methods, followed by a semi-supervised GMVAE utilizing contrastive and  
510 domain-adversarial learning to correct triple effects. We first conducted systematic experiments  
511 to demonstrate that CP data inherently exhibits triple effects, including batch, row, and column  
512 effects. Through comprehensive experiments on multiple batches varying in cell types, plate  
513 designs, perturbation types, we then validated cpDistiller's superior performance in correcting  
514 triple effects and preserving cellular phenotypic heterogeneity. Moreover, we highlighted the  
515 extensive advantages of cpDistiller, including its ability to integrate information-rich details  
516 from raw images, its support for incremental learning, and its consistent robustness across  
517 various feature selection strategies. Additionally, we conducted various downstream  
518 experiments to illustrate the application of cpDistiller in revealing gene functions and gene  
519 relationships, highlighting its ability to uncover gene associations both when combined with  
520 scRNA-seq data and independently. Scientists have found that target-based drug discovery can  
521 be limited in certain situations, and phenotypic drug discovery is sometimes more likely to  
522 succeed<sup>33</sup>. In recent years, research related to CP has expanded significantly, offering a better  
523 perspective on studying drug targets and cellular phenotypic heterogeneity due to its simpler  
524 procedures and lower costs. Therefore, we also explored the potential of cpDistiller in  
525 identifying gene and compound targets, which could enhance drug discovery efforts.

526 While cpDistiller demonstrated excellent performance, there are several areas that could be  
527 explored for future improvement. First, compared to CellProfiler-based features, which has  
528 specific and interpretable for each dimension, the low-dimensional representations learned  
529 from cpDistiller may lack interpretability. Enhancing the interpretability of deep learning  
530 models remains a challenge<sup>33</sup>. Second, we can utilize Segment Anything Model<sup>48</sup> along with  
531 expert annotation and refinement to generate accurate annotations for a subset of CP images.  
532 These annotations will serve as a basis for supervised training of the cpDistiller-extractor  
533 module, allowing it to effectively capture more detailed cellular information. Third, to gain a  
534 more comprehensive understanding of cellular behavior, cpDistiller could be extended to  
535 integrate with multi-omics profiling, connecting morphological and molecular phenotypes at  
536 single-cell resolution, thereby providing a more detailed insight into genetic perturbations and  
537 their effects<sup>49</sup>.

538

539 **Methods**

540 **Overview of cpDistiller**

541 cpDistiller consists of three main modules: the extractor module, the joint training module and  
542 the technical correction module. We first used the cpDistiller-extractor module to extract  
543 features from Cell Painting (CP) images. Then we integrated cpDistiller-extractor-based and  
544 CellProfiler-based features via the joint training module and processed the combined features  
545 through the technical correction module to remove batch, row and column effects.

546 **The extractor module**

547 We develop the extractor module to address the limitations of non-data-driven approaches, such  
548 as CellProfiler, while minimizing biases from human interference. By modeling our extraction  
549 task as the segmentation task in computer vision, we transform raw images into low-  
550 dimensional representations using an end-to-end process. We select Mesmer<sup>30</sup>, pre-trained on  
551 the TissueNet datasets<sup>30</sup>, as the base model for transfer learning due to its strong performance  
552 in cell segmentation tasks on cellular images<sup>50</sup>. By experimenting with various intermediate  
553 layers of Mesmer and considering factors such as overall effectiveness, number of parameters,  
554 and processing speed, we select the backbone of Mesmer and its pre-trained parameters.

555 In the CP assay, each well contains 9 sub-images sized at  $1080 \times 1080$  pixels, captured  
556 across five channels: mitochondria (Mito), nucleus (DNA), nucleoli and cytoplasmic RNA  
557 (RNA), endoplasmic reticulum (ER), and Golgi and plasma membrane and the actin  
558 cytoskeleton (AGP)<sup>14</sup>. To meet Mesmer's dual-channel input requirements and enhance the  
559 visibility of cell contours, we discard the AGP and Mito channels, focusing instead on the DNA  
560 channel and the cell channel created by averaging the ER and RNA channels, as the information  
561 in the AGP and Mito channels is difficult to predict<sup>33</sup>. Besides, to align with the Mesmer's  
562 specifications for pixel density and reduce computational resources, we apply a tiling operation,  
563 adjusting the stride and overlap ratio to crop each large sub-image into multiple small images  
564 with dimensions of  $256 \times 256$  pixels. The number of small images generated from each well's  
565 sub-image is calculated using the following formula:

$$566 N_{image} = (\left\lceil \frac{H - h}{s(1 - o)} \right\rceil + 1) \times (\left\lceil \frac{W - w}{s(1 - o)} \right\rceil + 1)$$

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

$$573 x_{start,i} = i \cdot s \cdot (1 - o)$$

$$574 y_{start,j} = j \cdot s \cdot (1 - o),$$

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

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

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

### 585 The joint training module

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

$$593 \quad \mathbf{Y}_1 = PCA(\mathbf{E}_{pooled}).$$

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

$$596 \quad \mathbf{Y}_2 = reshape(\mathbf{E}_{pooled}) \odot \sigma(Cov1d(AvgPool(reshape(\mathbf{E}_{pooled})))),$$

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

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

$$604 \quad \mathbf{z}_e = \mathbf{Y}_1 \oplus (W_1 \mathbf{Y}_2 + b_1),$$

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

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

611     **The technical correction module**

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

616     **The Gaussian mixture variational autoencoder.** Given the feature matrix  $\mathbf{X}_p \in \mathbb{R}^{n \times I_p}$   
617     integrated by the joint training module, where  $n$  denotes the number of wells and  $I_p$  denotes  
618     the feature dimensions, the GMVAE takes the input  $\mathbf{X}_p$  to obtain low-dimensional  
619     representations  $\mathbf{Z}_p$ . To illustrate the workflow of the GMVAE, we consider a sample  $\mathbf{x}$  from  
620      $\mathbf{X}_p$ . Since CP data encompasses numerous perturbations that may conform to distinct  
621     underlying Gaussian distributions, we use the GMVAE to identify these categorical  
622     distributions (pseudo-labels, denoted as  $\mathbf{y}$ ), which helps the model capture biological variation  
623     and ultimately contribute to biologically meaningful low-dimensional representations (denoted  
624     as  $\mathbf{z}$ ). The categorical distributions are inferred from the posterior distribution  $q_\psi(\mathbf{y}|\mathbf{x})$ , which  
625     follow semi-supervised training pattern. Here,  $q_\psi(\mathbf{y}|\mathbf{x})$  is represented by a feedforward neural  
626     network as:  $q_\psi(\mathbf{y}|\mathbf{x}) = \text{Cat}(\mathbf{y}|\pi_\psi(\mathbf{x}))$  and  $\pi_\psi(\mathbf{x})$  is a probability vector. Since the  
627     categorical distribution cannot be backpropagated in the neural network, we use the Gumbel-  
628     Softmax distribution<sup>51</sup> to facilitate gradient backpropagation, which allows the categorical  
629     distribution to be approximated using a continuous distribution.

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

$$632 \quad \log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\psi(\mathbf{y}, \mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}) + \log p_\theta(\mathbf{z}|\mathbf{y}) - \log q_\psi(\mathbf{z}|\mathbf{x}, \mathbf{y}) - \log q_\psi(\mathbf{y}|\mathbf{x}) + \\ 633 \quad \log p(\mathbf{y})],$$

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

639     The term  $\mathbb{E}_{q_\psi(\mathbf{y}, \mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{z}|\mathbf{y}) - \log q_\psi(\mathbf{z}|\mathbf{x}, \mathbf{y})]$  represents the Kullback-Leibler (KL)  
640     divergence between the variational posterior distribution  $q_\psi(\mathbf{z}|\mathbf{x}, \mathbf{y})$  and the conditional prior  
641     distribution  $p_\theta(\mathbf{z}|\mathbf{y})$ . This divergence ensures that the variational posterior distribution aligns  
642     with the prior, meaning that the learned latent representations  $\mathbf{z}$  conform to the expected

643 Gaussian distribution. Specifically,  $p_\theta(\mathbf{z}|\mathbf{y}) = N(\mathbf{z} \mid \mu_\theta(\mathbf{y}), \sigma_\theta^2(\mathbf{y}))$  represents the prior  
644 Gaussian distribution conditioned on the category  $\mathbf{y}$ , while  $q_\psi(\mathbf{z}|\mathbf{x}, \mathbf{y}) =$   
645  $N(\mathbf{z} \mid \mu_\psi(\mathbf{x}, \mathbf{y}), \sigma_\psi^2(\mathbf{x}, \mathbf{y}))$  represents the variational posterior distribution conditioned on the  
646 input  $\mathbf{x}$  and category  $\mathbf{y}$ .

647  $q_\psi(\mathbf{y}|\mathbf{x})$  provides the probability that  $\mathbf{x}$  originates from various Gaussian distributions,  
648 satisfying  $\sum_{k=1}^K q_\psi(\mathbf{y}|\mathbf{x}) = 1$ , where  $K$  is the predefined number of Gaussian distributions.

649 We use cross-entropy loss to refine the probabilities, guiding them towards high confidence  
650 regions. We treated the prior  $p(\mathbf{y})$  as a constant during loss backpropagation, as it does not  
651 influence the updates to the model's parameters.

652 **The gradient reversal module.** Within the low-dimensional space derived from the GMVAE,  
653 discriminators are employed to identify the source of each well's representation  $\mathbf{z}$ , specifically  
654 the batch, row, and column labels of the corresponding well. To implement adversarial learning  
655 similar to generative adversarial networks (GANs)<sup>21</sup>, different batches, rows, and columns can  
656 be treated as distinct domains. We then employ domain-adversarial learning, specifically the  
657 gradient reversal layer (GRL)<sup>22</sup>, to remove triple effects across these domains. Here,  
658 discriminators are denoted as:  $D_{batch}, D_{row}, D_{column}$ , which are used to predict the batch, row,  
659 and column labels of  $\mathbf{z}$ .

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

668 
$$GRL_\lambda(x) = x,$$

669 
$$\frac{\partial GRL_\lambda(x)}{\partial x} = -\lambda x,$$

670 where  $\lambda$  is a hyperparameter that undergoes a non-linear transformation, varying from 0 to 1.  
671 In the early stages of training, its value is kept small to allow the discriminators to train  
672 sufficiently and develop discriminative capabilities. As training progresses,  $\lambda$  gradually  
673 increases to strengthen adversarial interactions between the encoder part of the GMVAE and  
674 the discriminators. The calculation formula for  $\lambda$  is as follows:

675 
$$\lambda = \frac{2}{1 + \exp(-\gamma \cdot p)} - 1,$$

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

678 
$$\mathbf{z} = E(\mathbf{x}; \theta_g),$$

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

683 
$$D(GRL_\lambda(\mathbf{z}), \theta_{D_i}), i \in \{\text{batch, row, column}\},$$

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

689 Using the column discriminator loss as an example, the loss function is inspired by the label  
690 smoothing cross-entropy loss<sup>52</sup>. This approach is similarly applied to the batch and row  
691 discriminators for avoiding overconfidence:

692 
$$L_{D_{\text{col}}} = - \sum_{k=1}^K [(1 - \epsilon) \log(p_{\text{col}}(k)) \delta_{k,y}^{\text{col}} + \epsilon \log(p_{\text{col}}(k)) \theta_k^{\text{col}}],$$

693 where  $y$  represents the index of the true label, and  $\epsilon$  is a hyperparameter representing the  
694 proportion of soft labels considered when calculating the loss.  $\delta_{k,y}^{\text{col}}$  is a one-hot encoding,  
695 where the position corresponding to the true label is 1, with others set to 0. The probability that  
696 the  $\mathbf{z}$  originates from  $k$ -th column, as predicted by the discriminator, is represented by  
697  $p_{\text{col}}(k)$ .  $\theta_k^{\text{col}}$  represents the soft labels, which indicates the weight assigned to the  $k$ -th  
698 column.

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

710 assigning varying significance to them based on their distance from the true label, while still  
 711 avoiding overconfidence:

$$\theta_k^{col} = \begin{cases} \alpha, & \text{if } k = y \\ q^{|k-y|}, & \text{otherwise} \end{cases}$$

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

$$716 \quad f(q) = \alpha - \alpha q^{y+1} - \alpha q^{n-y} - 1 + q + \alpha q$$

$$717 \quad q : f(q) = 0, s.t. 0 < q < 1,$$

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

722 **The contrastive learning module.** We first establish nearest neighbor relationships by  
 723 utilizing  $k$ -nearest neighbors (KNN) intra technical effects and mutual nearest neighbors  
 724 (MNN) inter technical effects, based on CellProfiler-based features and cpDistiller-extractor-  
 725 based features processed with average pooling, using cosine distance as the similarity metric.  
 726 This approach is applied to batch, row, and column effects, respectively, constructing the  
 727 adjacency matrix that captures nearest neighbor relationships between wells for each effect.  
 728 We then leverage the relationships identified across multiple technical effects to form triplets.  
 729 We next use triplet loss<sup>20</sup> to restore more accurate nearest-neighbor relationships for each well.

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

$$736 \quad (\mathbf{a}, \mathbf{p}, \mathbf{n}) \Leftrightarrow \begin{cases} (\mathbf{a}, \mathbf{p}) \in S_{knn_{technical\ effects}} \cup S_{mnn_{technical\ effects}} \\ (\mathbf{a}, \mathbf{n}) \notin S_{knn_{technical\ effects}} \cup S_{mnn_{technical\ effects}} \end{cases}$$

737 where  $S_{knn_{technical\ effects}}$  represents the set of all pairs of data points that originate from the  
 738 same type of technical effects and are  $k$ -nearest neighbors within the same category.  
 739  $S_{mnn_{technical\ effects}}$  represents the set of all pairs of data points that originate from the same  
 740 type of technical effects but are mutual nearest neighbors in different categories.

741 The triplet loss function is not directly applied to the latent space of the GMVAE. Instead,  
 742 the representations  $\mathbf{z}$  need to pass through a nonlinear projection head, as previous research

743 has demonstrated that adding such a layer can improve the quality of the learned  
744 representations<sup>53</sup>. This architecture can be represented as  $ph(\mathbf{z})$ :

745 
$$ph(\mathbf{z}) = \text{LeakyReLU}(W_{ph}\mathbf{z} + b_{ph}),$$

746 where the parameters  $W_{ph}$  and  $b_{ph}$  represent the learnable weights and biases of a linear  
747 layer.

748 Then we use the triplet loss to remove triple effects:

749 
$$L_c = \max(d(ph(\mathbf{a}), ph(\mathbf{p})) - d(ph(\mathbf{a}), ph(\mathbf{n})) + \xi, 0),$$

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

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

753 
$$\text{Loss} = w_{ELBO}L_{ELBO} + w_{dis}(L_{D_{row}} + L_{D_{col}}) + w_{con}L_{c'},$$

754 where  $L_{c'}$  represents the triplet loss calculated by considering well position effects.

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

757 
$$\text{Loss} = w_{ELBO}L_{ELBO} + w_{dis}(L_{D_{batch}} + L_{D_{row}} + L_{D_{col}}) + w_{con}L_c,$$

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

762 The parameters  $\theta'_t$  of the final trained model are obtained through exponential moving  
763 average (EMA)<sup>54</sup>, which can be given by:

764 
$$\theta'_t = \alpha_{ema}\theta'_{t-1} + (1 - \alpha_{ema})\theta_t,$$

765 where  $\theta_t$  represents the weighted model parameters obtained at round  $t$ , and  $\alpha_{ema}$  is a  
766 hyperparameter. All training hyperparameters are available in the training details.

## 767 Training details

768 For the extractor module, the overlap ratio  $o$  is set to 0.25 and the sliding stride  $s$  is set to  
769 256. Max pooling with kernel size and step size both set to 16, is applied to merge feature maps  
770 from nine sub-images, yielding the output dimension  $I_e$  to 11,664, corresponding to the 108  
771  $\times$  108 feature map. For the joint training module, average pooling with kernel size and step size  
772 both set to 9, is applied to smooth out irrelevant variation. For the technical correction module,  
773 the hidden dimensions of the encoder and decoder are set to 512, while the hidden dimensions  
774 of discriminators are set to 128. The latent space dimensionality of cpDistiller is set to 50. The  
775 technical correction module uses LeakyReLU with default parameters as the activation  
776 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 7,638 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  $\gamma$  is set to 10 in GRL. The soft label hyperparameter  $\alpha$  is set to 0.75, and the hyperparameter  $\epsilon$  in label smoothing cross-entropy loss is set to 0.1. For the contrastive learning, the  $\xi$  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  $\alpha_{ema}$  set to  $1 - \frac{5}{epochs}$ . In the loss functions, the weights of  $w_{ELBO}$ ,  $w_{dis}$  and  $w_{con}$  are set to 1,  $\frac{I_p}{35}$  and  $\frac{I_p}{35}$ , respectively. The experimental environment includes two 24GB Nvidia 4090 graphics cards and 96 Intel(R) Xeon(R) Gold 5318N CPUs @ 2.10GHz.

## 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 ARCHS4<sup>36</sup> 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 significant<sup>55</sup>, 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 Enrichr<sup>56</sup> platform to conduct GO term enrichment analysis<sup>57</sup> 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.

Gene relationship analysis: we used the GeneMANIA tool<sup>41</sup> 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.<sup>14</sup>, 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

814 between negative treatments and set the 95th percentile of these distances as the threshold for  
815 identifying active genes. Overexpression treatments were classified as active if their Euclidean  
816 distance from the negative treatments exceeded this threshold. In Batch\_1, 226 genes were  
817 identified as active genes.

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

826 The calculation of fraction retrieved in gene and compound retrieval tasks: we calculated  
827 and compared the fraction retrieved scores for gene-gene and gene-compound simulated  
828 retrieval tasks using cpDistiller-derived representations and CellProfiler-based features in the  
829 cpg0000 dataset, following the workflow in ref.<sup>47</sup>.

830 **Pre-processing for CellProfiler-based features.** We utilized features pre-extracted with  
831 CellProfiler from the prior study<sup>14</sup>, which included up to 7,638 features for images with both  
832 brightfield and fluorescence images and 4,752 features for only fluorescence images, after  
833 removing features containing NaN values. Additionally, we followed the feature selection  
834 process described in the study<sup>14</sup> to select 1,446 features. Due to data quality issues reported in  
835 the original dataset<sup>14</sup>, we excluded data from Batch\_12 and the BR00123528A plate. For pre-  
836 computed feature transformation, we followed the Pycytominer<sup>58</sup> and applied its default  
837 operations, including z-score normalization on the entire dataset.

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

#### 847 **Evaluation metrics**

848 To assess the effectiveness of different methods in removing technical effects, we used three  
849 technical correction metrics: average silhouette width (ASW)<sup>26</sup>, technic average silhouette  
850 width (tASW)<sup>27</sup>, and graph connectivity<sup>27</sup>. To evaluate biological preservation, four metrics  
851 were used: isolated label F1<sup>27</sup>, isolated label silhouette<sup>27</sup>, perturbation average silhouette width

852 (pASW)<sup>27</sup>, and normalized mutual information (NMI)<sup>27</sup>. For a more intuitive and clear  
853 comparison, we normalized all metrics to a range of 0 to 1, where higher values indicated better  
854 results. Metrics involving perturbation labels, excluding ASW and tASW, which reflected the  
855 clustering of same perturbations and the dispersion of different ones, were calculated on the  
856 selected subset. The reasons for these calculation approaches are based on two main challenges.  
857 First, the perturbation labels in the JUMP dataset are not well-defined, making it difficult to  
858 accurately cluster similar perturbations. Specifically, the ORF perturbations induced by  
859 different reagents may target functionally similar genes, making some perturbations inherently  
860 difficult to distinguish. Besides, with over 10,000 unique perturbation labels in ORF data, it is  
861 difficult to separate functionally overlapping or densely labeled perturbations. Second,  
862 category imbalance and numerous categories pose a challenge for accurately calculating  
863 clustering metrics. To be specific, perturbations in controls replicated substantially more  
864 frequently than those in the treatment, often by several to dozens of times. Metrics like pASW,  
865 which do not account for category imbalance, result in biased scores. Metrics like NMI  
866 typically require unsupervised algorithms, such as Leiden<sup>59</sup> to generate predicted labels.  
867 However, the numerous and less well-defined perturbation labels in the JUMP dataset pose  
868 additional challenges, leading to inaccuracies in the metric calculations.

869 To address the issues mentioned above, we focused on certain perturbations. For correcting  
870 well position effects, we used controls, including five positive controls (JCP2022\_037716,  
871 JCP2022\_035095, JCP2022\_050797, JCP2022\_012818, and JCP2022\_915132) and four  
872 negative controls (JCP2022\_915128, JCP2022\_915129, JCP2022\_915130, and  
873 JCP2022\_915131). For removing triple effects, we selected perturbations induced by 34  
874 different compounds that showed better clustering on target plates. First, these selections  
875 reduce the impact of unwell-defined perturbation labels, allowing for more reliable clustering  
876 of similar perturbations. In particular, positive and negative controls theoretically should have  
877 significant or negligible effects on cell phenotypes, resulting in marked differences in their  
878 impacts<sup>14</sup>. Furthermore, different positive controls are expected to exhibit distinct effects on  
879 cell phenotypes. Besides, for target plates, perturbations are shared across different batches,  
880 which can be used for aligning data<sup>15</sup>. Additionally, perturbations caused by the same  
881 compounds in target plates that clustered well in the raw data should remain well-clustered  
882 after technical correction. Otherwise, it may indicate overcorrection, as true biological  
883 differences should persist even after correction. Second, these selections help address category  
884 imbalance and reduce categories, allowing for more accurate calculation of clustering metrics.  
885 Specifically, the number of positive and negative controls are generally consistent, and each  
886 perturbation in target plates was duplicated approximately 20 times, ensuring a consistent  
887 number of tests.

888 The specific calculation of above metrics, which follow the approach of previous works by  
889 scArches<sup>26</sup> and scIB<sup>27</sup> for analyzing single-cell data, are further detailed in Supplementary Text  
890 7.

891 **Baseline methods**

892 To evaluate technical correction in CP data, we compared several widely used single-cell  
893 analysis methods: Seurat v5 (v5.0.1)<sup>23</sup>, Harmony (v0.0.9)<sup>19</sup>, Scanorama (v1.7.4)<sup>24</sup>, scVI  
894 (v0.14.6)<sup>18</sup>, and scDML (v0.0.1)<sup>25</sup>. UMAP plots, used to compare embeddings generated by  
895 different methods, were created with the following parameters:  $n\_neighbors=15$ ,  $min\_dist=0.1$ ,  
896 and  $random\_state=9000$ . For removing well position effects, data points within each batch  
897 were standardized using z-score. When correcting triple effects, z-score standardization was  
898 applied within each batch, followed by integration across batches and an additional round of  
899 normalization to mitigate batch effects. All methods, including cpDistiller, were applied to CP  
900 data that had undergone the above preprocessing steps to ensure consistency.

901 For Harmony, Seurat v5 and Scanorama, which can correct multiple technical effects, the  
902 workflow involved sequentially passing the row and column labels if considering removing  
903 well position effects. When aiming to remove triple effects, these methods required first passing  
904 the batch labels, followed sequentially by the row and column labels. For methods that are  
905 limited to correcting only one type of technical effects, if considering well position effects, we  
906 selected row and column effects separately as correction targets, resulting in two results,  
907 respectively. Considering triple effects, we focused solely on batch effects as they represent the  
908 most significant influence. Specific implementation details are as follows:

909 Seurat v5: we followed the example pipeline provided by the Seurat v5 for integrative  
910 analysis. The labels for technical effects were passed sequentially into the *split* function. We  
911 skipped the *NormalizeData* function and *FindVariableFeatures* function, as these are  
912 specifically tailored for scRNA-seq data. During the correction phase, the *CCA* method was  
913 employed to remove technical effects. All parameters were used with default settings.

914 Harmony: dimensionality reduction was performed using *scanpy.tl.pca* to the default size of  
915 50. The labels for technical effects were passed sequentially into the  
916 *scanpy.external.pp.harmony\_integrate* function. All parameters were used with default settings.

917 Scanorama: dimensionality reduction was performed using *scanpy.tl.pca* to the default size  
918 of 50. The labels for technical effects were passed sequentially into the  
919 *scanpy.external.pp.scanorama\_integrate* function. All parameters were used with default  
920 settings.

921 scVI: we followed the example tutorial provided by the scVI on GitHub for processing  
922 scRNA-seq data. We did not use the *scanpy.pp.highly\_variable\_genes* function, which is  
923 specifically tailored for scRNA-seq data. In the preprocessing stage, we followed the  
924 preprocessing operations of previous work to process each feature as follows:  $\hat{x}_i = x_i -$   
925  $\min(x) + 1^{15}$ . Other operations were used with default settings.

926 scDML: we followed the example tutorial on GitHub provided by scDML. We did not use  
927 the *scanpy.pp.log1p* function and the *scanpy.pp.highly\_variable\_genes* function, which are  
928 specifically tailored for scRNA-seq data. Other operations were used with default settings.

929 **Data availability**

930 We downloaded well-level images and the features extracted by CellProfiler, following the  
931 recommended workflow<sup>14</sup>. For the JUMP datasets, raw stained images and the features  
932 extracted by CellProfiler are accessible via AWS at <https://registry.opendata.aws/cellpainting-gallery/>.

934 **Code availability**

935 The cpDistiller software, along with detailed documentation and tutorials, is freely available at  
936 <https://github.com/Cell-Painting/cpDistiller>.

937 **Acknowledgements**

938 This work was supported by the National Key Research and Development Program of China  
939 grants no. 2020YFA0908700 (J.L.), 2020YFA0908702 (J.L.) and 2024YFA1307703 (S.C.),  
940 the National Natural Science Foundation of China grants no. 62473212 (S.C.), 62203236 (S.C.)  
941 and 62272246 (J.L.), and the Young Elite Scientists Sponsorship Program by CAST grant no.  
942 2023QNRC001 (S.C.).

943 **Competing Interests Statement**

944 The authors declare no competing interests.

945 **References**

- 946 1. Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of next-  
947 generation sequencing technologies. *Nat. Rev. Genet.* **17**, 333–351 (2016).
- 948 2. Perlman, Z. E. *et al.* Multidimensional drug profiling by automated microscopy. *Science* **306**, 1194–1198 (2004).
- 949 3. Wawer, M. J. *et al.* Toward performance-diverse small-molecule libraries for cell-based  
950 phenotypic screening using multiplexed high-dimensional profiling. *Proc. Natl. Acad. Sci.* **111**, 10911–10916 (2014).
- 951 4. Moshkov, N. *et al.* Learning representations for image-based profiling of perturbations. *Nat. Commun.* **15**, 1594 (2024).
- 952 5. Fredin Haslum, J. *et al.* Cell Painting-based bioactivity prediction boosts high-  
953 throughput screening hit-rates and compound diversity. *Nat. Commun.* **15**, 3470 (2024).
- 954 6. Bray, M.-A. *et al.* Cell Painting, a high-content image-based assay for morphological  
955 profiling using multiplexed fluorescent dyes. *Nat. Protoc.* **11**, 1757–1774 (2016).
- 956 7. Chandrasekaran, S. N., Ceulemans, H., Boyd, J. D. & Carpenter, A. E. Image-based  
957 profiling for drug discovery: due for a machine-learning upgrade? *Nat. Rev. Drug Discov.* **20**, 145–159 (2021).
- 958 8. Jo, Y. *et al.* Label-free multiplexed microtomography of endogenous subcellular  
959 dynamics using generalizable deep learning. *Nat. Cell Biol.* **23**, 1329–1337 (2021).
- 960 9. Liberali, P., Snijder, B. & Pelkmans, L. A hierarchical map of regulatory genetic  
961 interactions in membrane trafficking. *Cell* **157**, 1473–1487 (2014).
- 962 10. Collinet, C. *et al.* Systems survey of endocytosis by multiparametric image analysis. *Nature* **464**, 243–249 (2010).
- 963 11. Carpenter, A. E. & Singh, S. Bringing computation to biology by bridging the last mile. *Nat. Cell Biol.* **26**, 5–7 (2024).

- 970 12. Loo, L.-H., Wu, L. F. & Altschuler, S. J. Image-based multivariate profiling of drug  
971 responses from single cells. *Nat. Methods* **4**, 445–453 (2007).
- 972 13. Weisbart, E. *et al.* Cell Painting Gallery: an open resource for image-based profiling.  
973 *Nat. Methods* **21**, 1775–1777 (2024).
- 974 14. Chandrasekaran, S. N. *et al.* JUMP Cell Painting dataset: morphological impact of  
975 136,000 chemical and genetic perturbations. Preprint at bioRxiv  
976 <https://doi.org/10.1101/2023.03.23.534023> (2023).
- 977 15. Arevalo, J. *et al.* Evaluating batch correction methods for image-based cell profiling.  
978 *Nat. Commun.* **15**, 6516 (2024).
- 979 16. Carpenter, A. E. *et al.* CellProfiler: image analysis software for identifying and  
980 quantifying cell phenotypes. *Genome Biol.* **7**, 1–11 (2006).
- 981 17. Boulard, G. A., Mahfouz, A. & Reinders, M. J. Consequences and opportunities arising  
982 due to sparser single-cell RNA-seq datasets. *Genome Biol.* **24**, 86 (2023).
- 983 18. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling  
984 for single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- 985 19. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with  
986 Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
- 987 20. Schroff, F., Kalenichenko, D. & Philbin, J. Facenet: A unified embedding for face  
988 recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision  
989 and Pattern Recognition* (CVPR, 2015).
- 990 21. Salimans, T. *et al.* Improved techniques for training gans. In *Advance in Neural  
991 Information Processing Systems* (NeurIPS, 2016).
- 992 22. Ganin, Y. & Lempitsky, V. Unsupervised domain adaptation by backpropagation. In  
993 *International Conference on Machine Learning* (ICML, 2015).
- 994 23. Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell  
995 analysis. *Nat. Biotechnol.* **42**, 293–304 (2024).
- 996 24. Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell  
997 transcriptomes using Scanorama. *Nat. Biotechnol.* **37**, 685–691 (2019).
- 998 25. Yu, X., Xu, X., Zhang, J. & Li, X. Batch alignment of single-cell transcriptomics data  
999 using deep metric learning. *Nat. Commun.* **14**, 960 (2023).
- 1000 26. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning.  
1001 *Nat. Biotechnol.* **40**, 121–130 (2022).
- 1002 27. Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics.  
1003 *Nat. Methods* **19**, 41–50 (2022).
- 1004 28. Caicedo, J. C. *et al.* Data-analysis strategies for image-based cell profiling. *Nat.  
1005 Methods* **14**, 849–863 (2017).
- 1006 29. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and  
1007 projection for dimension reduction. Preprint at arXiv  
1008 <https://doi.org/10.48550/arXiv.1802.03426> (2018).
- 1009 30. Greenwald, N. F. *et al.* Whole-cell segmentation of tissue images with human-level  
1010 performance using large-scale data annotation and deep learning. *Nat. Biotechnol.* **40**,  
1011 555–565 (2022).
- 1012 31. Shrestha, P., Kuang, N. & Yu, J. Efficient end-to-end learning for cell segmentation  
1013 with machine generated weak annotations. *Commun. Biol.* **6**, 232 (2023).
- 1014 32. Jocher, G., Chaurasia, A. & Qiu, J. Ultralytics yolov8.  
1015 <https://github.com/ultralytics/ultralytics> (2023).
- 1016 33. Seal, S. *et al.* Cell Painting: a decade of discovery and innovation in cellular imaging.  
1017 *Nat. Methods* 1–15 (2024).
- 1018 34. Wang, K. C. & Chang, H. Y. Epigenomics: technologies and applications. *Circ. Res.*  
1019 **122**, 1191–1199 (2018).

- 1020 35. Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature*  
1021 **618**, 616–624 (2023).
- 1022 36. Lachmann, A. *et al.* Massive mining of publicly available RNA-seq data from human  
1023 and mouse. *Nat. Commun.* **9**, 1366 (2018).
- 1024 37. Parker, R. & Sheth, U. P bodies and the control of mRNA translation and degradation.  
1025 *Mol. Cell* **25**, 635–646 (2007).
- 1026 38. Will, C. & Lührmann, R. Spliceosome structure and function. *Cold Spring Harb  
1027 Perspect Biol* **3**: a003707. (2011).
- 1028 39. Waltz, F. & Giegé, P. Striking diversity of mitochondria-specific translation processes  
1029 across eukaryotes. *Trends Biochem. Sci.* **45**, 149–162 (2020).
- 1030 40. Fuchs, E. & Weber, K. Intermediate filaments: structure, dynamics, function and  
1031 disease. *Annu. Rev. Biochem.* **63**, 345–382 (1994).
- 1032 41. Warde-Farley, D. *et al.* The GeneMANIA prediction server: biological network  
1033 integration for gene prioritization and predicting gene function. *Nucleic Acids Res.* **38**,  
1034 W214–W220 (2010).
- 1035 42. Ross, D. T. & Perou, C. M. A comparison of gene expression signatures from breast  
1036 tumors and breast tissue derived cell lines. *Dis. Markers* **17**, 99–109 (2001).
- 1037 43. Wheeler, D. L. *et al.* Database resources of the national center for biotechnology  
1038 information. *Nucleic Acids Res.* **36**, D13–D21 (2007).
- 1039 44. Paysan-Lafosse, T. *et al.* InterPro in 2022. *Nucleic Acids Res.* **51**, D418–D427 (2023).
- 1040 45. Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**,  
1041 D412–D419 (2021).
- 1042 46. Mallon, B. S. *et al.* StemCellDB: the human pluripotent stem cell database at the  
1043 National Institutes of Health. *Stem Cell Res.* **10**, 57–66 (2013).
- 1044 47. Chandrasekaran, S. N. *et al.* Three million images and morphological profiles of cells  
1045 treated with matched chemical and genetic perturbations. *Nat. Methods* **1–8** (2024).
- 1046 48. Kirillov, A. *et al.* Segment anything. In *Proceedings of the IEEE/CVF International  
1047 Conference on Computer Vision* (ICCV, 2023).
- 1048 49. Tang, Q. *et al.* Morphological profiling for drug discovery in the era of deep learning.  
1049 *Brief. Bioinform.* **25**, (2024).
- 1050 50. Amitay, Y. *et al.* CellSighter: a neural network to classify cells in highly multiplexed  
1051 images. *Nat. Commun.* **14**, 4302 (2023).
- 1052 51. Jang, E., Gu, S. & Poole, B. Categorical reparameterization with gumbel-softmax.  
1053 Preprint at arXiv <https://doi.org/10.48550/arXiv.1611.01144> (2016).
- 1054 52. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the inception  
1055 architecture for computer vision. In *Proceedings of the IEEE Conference on Computer  
1056 Vision and Pattern Recognition* (CVPR, 2016).
- 1057 53. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A simple framework for contrastive  
1058 learning of visual representations. In *International Conference on Machine Learning*  
1059 (ICML, 2020).
- 1060 54. Tarvainen, A. & Valpola, H. Mean teachers are better role models: Weight-averaged  
1061 consistency targets improve semi-supervised deep learning results. In *Advance in  
1062 Neural Information Processing Systems* (NeurIPS, 2017).
- 1063 55. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for  
1064 interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci.* **102**, 15545–15550  
1065 (2005).
- 1066 56. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web  
1067 server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
- 1068 57. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. *Nat. Genet.* **25**,  
1069 25–29 (2000).

- 1070 58. Serrano, E. *et al.* Reproducible image-based profiling with Pycytominer. Preprint at  
1071 *arXiv* <https://doi.org/10.48550/arXiv.2311.13417> (2023).
- 1072 59. Traag, V. A., Waltman, L. & Van Eck, N. J. From Louvain to Leiden: guaranteeing well-  
1073 connected communities. *Sci. Rep.* **9**, 1–12 (2019).
- 1074