Bridging Self-Supervision and Mechanism of Action Discovery in Morphological Profiling

Anonymous CVPR submission

Paper ID *****

Abstract

In the quest to interpret complex cellular responses, 001 002 self-supervised learning (SSL) methods have been developed, promising powerful generic representations. How-003 ever, their performance on biologically critical tasks such 004 as mechanism of action (MoA) classification remains lim-005 006 ited. We argue that no single model-no matter how sophisticated or generalisable—can produce representations 007 that are optimal for all downstream tasks, as different 008 009 objectives impose conflicting requirements. To address this, we propose a novel framework called Task-guided 010 Representation exaptation (TRex). In TRex, a generic 011 (possibly self-supervised) model first extracts broad and 012 rich morphological embeddings, which are then refined by a 013 lightweight adaptation network optimised for biological rel-014 evance linked to the specific downstream tasks. This modu-015 lar design enables rapid and resource-efficient transforma-016 017 tion of generic features into biologically meaningful, taskfocused representations — without the need to retrain large-018 scale models. We evaluate TRex on a 20-plate Cell Painting 019 020 dataset spanning two cell lines and show that MoA-based 021 adaptation not only significantly improves MoA classifica-022 tion performance (doubling the mAP), but also enhances compound recognition. Our results highlight the limitations 023 of static, generic representations and demonstrate the util-024 ity of task-aware adaptation for maximising the biological 025 relevance of morphological profiling. 026

027 1. Introduction

Cell Painting is a high-content imaging assay designed for 028 029 morphological profiling, offering a rich, multiplexed readout of cellular responses to chemical and genetic pertur-030 bations [1, 4, 9]. By staining multiple cellular compart-031 ments with a combination of fluorescent dyes, the assay 032 provides valuable insights into cellular structure and func-033 tion. However, to fully exploit the information encoded in 034 035 these images, robust and biologically meaningful methods

for representation learning and analysis are required. Over036the years, approaches to extracting and interpreting fea-037tures have evolved rapidly from handcrafted descriptors [2]038to deep learning-based techniques [8], and more recently,039to self-supervised learning frameworks [6, 7], reducing the040need for extensive annotations and promising improved rep-041resentations.042

The earliest approach to feature extraction from Cell 043 Painting images relied on CellProfiler [2], an open-source 044 image analysis platform that extracts hand-crafted morpho-045 logical features. These features, which include measure-046 ments of shape, texture, intensity, and spatial organization, 047 are computed at the single-cell level and aggregated into 048 well-level feature vectors. Despite wide-scale adoption, 049 this methodology presents several limitations. It typically 050 requires manual parameter tuning for different datasets or 051 tasks, and it is unlikely to capture the complex phenotypic 052 variations present in cellular images. The reliance on pre-053 defined morphological descriptors may limit the discovery 054 of novel patterns in large-scale profiling experiments. 055

To address these limitations, deep learning-based ap-056 proaches have been introduced to learn potentially superior 057 image features directly from data. One such approach is 058 DeepProfiler [8], which uses convolutional neural networks 059 (CNNs) to extract morphological features in a data-driven 060 manner. DeepProfiler employs EfficientNet-B0, initialized 061 with ImageNet weights and fine-tuned on a Cell Paint-062 ing dataset comprising 232 plates and 488 perturbations. 063 Feature extraction is performed using activations from the 064 block6a layer, generating 672-dimensional feature vectors 065 per cell. Training is conducted using a weakly supervised 066 learning (WSL) approach, where classification loss enables 067 the model to learn treatment-specific representations from 068 single-cell images. To mitigate confounding technical vari-069 ation and enhance the biological relevance of extracted fea-070 tures, DeepProfiler applies sphering transform-based batch 071 correction. 072

To further improve the scalability and generalisability of 073 morphological profiling, and to eliminate the need for extensive annotations required for training, researchers turned 075

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

to self-supervised learning (SSL). One such method is Sub-076 077 Cell [6], which utilizes Vision Transformers (ViTs) to ex-078 tract hierarchical and spatially-aware cellular features. Sub-Cell learns from individual cells extracted from the Human 079 080 Protein Atlas (HPA) [10], a near-proteome-wide dataset of high-resolution immunofluorescence images spanning 35 081 different cell lines. The model optimizes a three-component 082 loss function as the basis of feature learning: a reconstruc-083 084 tion loss to ensure general feature extraction, a cell-specific similarity loss to enforce consistency across augmented 085 086 views of the same cell, and a protein-specific localization loss to minimize variation between images stained for the 087 088 same protein across different cell types. Additionally, an attention pooling mechanism suppresses background artifacts 089 and enhances the focus on cellular features, improving ro-090 091 bustness across diverse microscopy tasks without requiring fine-tuning. 092

093 Another SSL based approach, DINO [5], also uses ViTs to learn rich, task-agnostic representations of cellular mor-094 phology without manual annotations or supervision. Devel-095 oped originally for natural images, DINO has proven highly 096 effective for microscopy data, capturing biologically mean-097 098 ingful variation across subcellular, single-cell, and at popu-099 lation levels. The architecture uses a teacher-student framework, where both networks are trained to produce similar 100 feature representations from different augmented views of 101 the same image, encouraging the model to focus on the in-102 103 variant, biologically relevant structures. Unlike DeepProfiler, which relies on weak supervision and convolutional 104 architectures, DINO operates entirely self-supervised and 105 benefits from the global contextual awareness of transform-106 ers. Compared to SubCell, which learns single-cell fea-107 108 tures using contrastive objectives over protein localization, 109 DINO emphasizes semantic consistency across scales and 110 has been shown to uncover hierarchical biological structure, including MoA, cell-cycle stages, and protein local-111 ization patterns. Its general-purpose embeddings make it a 112 strong candidate for further task-specific adaptation, as ex-113 plored in this work through the Task-guided Representation 114 exaptation (TRex) framework. 115

116 More recently, SSL has also been explored for imagelevel feature extraction, eliminating the need for cell seg-117 mentation. In [7], self-supervised vision transformers were 118 trained on a subset of the JUMP Cell Painting dataset [3] to 119 120 extract morphological features directly from full-field images, bypassing segmentation-based workflows. While this 121 segmentation-free approach significantly reduces computa-122 tional costs (by skipping the cell segmentation stage), the 123 biological relevance of the resulting features - particularly 124 in the context of MoA prediction - remains an open ques-125 tion, as demonstrated by the reported MoA performance. 126

In [6], the authors presented a detailed evaluation of compound matching (Cmpd) and mechanism of action

(MoA) prediction on the JUMP Cell Painting chemical per-129 turbation dataset, comparing the performance of CellPro-130 filer, DeepProfiler, DINO, and their own models. The re-131 sults reveal minimal variation between methods, with no-132 tably low performance for MoA prediction across all ap-133 proaches. The Cmpd mAP ranged from 0.27 (DeepProfiler) 134 to 0.39 (SubCell), while the MoA mAP remained consis-135 tently poor, varying from 0.16 (DeepProfiler) to 0.18 (Cell-136 Profiler). Notably, hand-crafted CellProfiler features out-137 performed all machine learning-based methods, including 138 SSL approaches, on the MoA task, suggesting that existing 139 learned representations may not yet capture the biologically 140 relevant features needed for this application. While self-141 supervised features promise to transfer to a wider range of 142 tasks, their ability to effectively capture task-specific bio-143 logical information remains unclear. 144

In this work, we address this limitation by introducing TRex, a novel framework for task-guided representation adaptation. TRex begins with a general-purpose feature extractor—preferably a pre-trained self-supervised model—and adds a lightweight adaptation network that refines these embeddings for a specific downstream task. Unlike traditional fine-tuning, our adaptation module operates on frozen features and requires minimal computation and data, enabling efficient adaptation of powerful SSL backbones.

We show that our approach yields substantial gains in both MoA classification and compound recognition, demonstrating that task-aware adaptation is key to unlocking the full potential of self-supervised representations for biological discovery. Our specific contributions are as follows:

- 1. We propose a novel, lightweight TRex architecture for adapting SSL-derived morphological features to biological prediction tasks, achieving substantial performance gains.
- 2. We design an efficient adaptation module composed of four fully connected layers with batch normalization, GELU activation, dropout regularization, and residual connections in selected layers to preserve signal flow and reduce overfitting.
- 3. We systematically evaluate the TRex framework and show that it can effectively double MoA classification performance over prior art, improving mAP from 0.15 to 0.32.
- 4. We evaluate three state-of-the art representations within TRex and show that SSL-derived features offer good performance, with DINO achieving the highest accuracy on MoA prediction tasks.
- 5. We demonstrate that MoA-based training leads to more biologically relevant representations, enhancing not only MoA classification but also compound recognition including generalisation to unseen compounds.
 178
 179
 180
 181

2. Introducing TRex 182

We propose a new two-stage architecture for efficient mor-183 phological profiling. In the first stage, self-supervised learn-184 ing is used to generate a broad set of biologically meaning-185 186 ful features. In the second stage, a learnt task-specific feature transformation is applied to adapt this self-supervised 187 188 feature space into a task-oriented representation optimised for specific downstream objective. We refer to our approach 189 as Task-Guided Representation exaptation (TRex). We 190 adopted the term exaptation from evolutionary biology be-191 192 cause it aligns with how our method repurposes and adapts generic self-supervised features for specific tasks. 193

The motivation behind this approach stems from the fact 194 that it is not possible for a single generic feature representa-195 tion to be equally effective across all morphological profil-196 ing tasks. Take for instance, two compound perturbations, 197 A and B, which coincide in mechanism of action (MoA). 198 It therefore follows that their feature representations should 199 200 aim to be distinct for compound identification (where the goal is to differentiate them), yet simultaneously identical 201 for MoA classification (where their shared functional effect 202 should be captured). Clearly, a static, task-agnostic rep-203 resentation cannot reconcile these conflicting requirements 204 205 without additional adaptation. This limitation is evident in the poor MoA prediction performance reported in [6] across 206 various methods, which plateau around 0.17, highlighting 207 208 the need for a more tailored, task-aware approach.

In our TRex framework, shown in Figure 1, the first 209 stage (TRex Feature Extractor) derives a comprehensive 210 211 pool of generic single-cell embeddings, most likely via a self-supervised approach or even by aggregating multiple 212 213 representations from complementary methods. The second stage (TRex Adaptation Module) applies a lightweight 214 transformation network to adapt these features into a task-215 specific final representation. Since the first stage is compu-216 217 tationally intensive, it is trained only once. In contrast, the second-stage adaptation is fast, efficient, and can be per-218 219 formed even with limited resources. This makes our approach not only biologically meaningful - by aligning rep-220 resentations with task-specific biological distinctions such 221 as compound mechanisms of action or phenotypic similar-222 223 ity- but also computationally practical, allowing for rapid 224 adaptation to new tasks and datasets. 225

Let:

- $f_{SSL} : X \to Z$ be a self-supervised feature extractor, 226 where X represents the input 5-channel microscopy im-227 ages, and $Z \in \mathbb{R}^d$ is the high-dimensional feature repre-228 sentation learned via self-supervised pretraining. 229
- $t_{\text{TASK}} : Z \to Y$ be the task-guided feature transfor-230 **mation network**, where Y is the final representation and 231 TASK is the task at hand, for example the perturbation 232 233 or MoA classification.
- 234 For example, in MoA classification task, the goal of our

second-stage network is to learn a transformation t_{MoA} such 235 that the extracted features Z are mapped to an optimized, 236 lower-dimensional space Y that retains biologically mean-237 ingful information linked to MoA, thus improving MoA 238 classification accuracy. Unlike end-to-end learning, our 239 framework operates on precomputed feature embeddings, 240 significantly reducing computational complexity and data 241 requirements. 242

2.1. TRex design details

The input features are extracted in the stage one model 244 (TRex Feature Extractor), which, in our case, is from one 245 of the three pre-trained models: DeepProfiler, DINO, or 246 Subcell. The key properties of these models are sum-247 marised in Table 1. The DeepProfiler model [8], based on 248 EfficientNet-B0, was trained in a weakly supervised man-249 ner on cellular images from the BBBC and LINCS datasets, 250 generating a 672-dimensional embedding per cell. The 251 DINO method [5], a vision transformer (ViT), was trained 252 in a self-supervised manner on the same datasets and pro-253 duces a 768-dimensional embedding per cell. Both mod-254 els were trained on 224×224 images containing five fluo-255 rescence channels (DNA, ER, RNA, Golgi/Actin, and Mi-256 tochondria). The SubCell model, also based on ViT, was 257 trained in a self-supervised manner on the Human Protein 258 Atlas (HPA) dataset, which consists of three-channel im-259 ages (ER, DNA, Protein). The model generates a 4608-260 dimensional embedding per cell. The HPA training set in-261 cluded 35 cell lines, including both U2OS and A549. All 262 three models have seen U2OS and A549 cells during train-263 ing, but SubCell was additionally trained on 33 other cell 264 lines, which in theory may enable it to produce more gener-265 alisable features across diverse cellular contexts. 266

Method	Dataset	Training	Cell Types	Dim.
DeepProfiler	L+B	Weak-sup.	U2OS,A549	672
DINO	L+B	Self-sup.	U2OS,A549	768
SubCell	HPA	Self-sup.	Various	4608

Table 1. Summary of feature extraction methods tested in our framework. L+B denotes a combination of LINCS and BBBC datasets.

To make these diverse and generic embeddings task-267 relevant, we employ the TRex adaptation network, which 268 transforms the output of each feature extractor into a biolog-269 ically meaningful and task-optimised representation suit-270 able for downstream prediction. This network is lightweight 271 by design, enabling efficient adaptation regardless of the 272 upstream embedding source. It consists of four fully con-273 nected layers, each followed by batch normalization, GELU 274 activation, and dropout regularization to ensure stable gradi-275 ent propagation, improved non-linearity, and reduced over-276

310

317

318

319

320

331

332



Figure 1. Overview of the TRex framework. Single-cell embeddings are extracted from five-channel Cell Painting images using pre-trained DeepProfiler, DINO, or SubCell models (STAGE 1). These embeddings are passed to the TRex adaptation module, which is trained using focal loss for the mechanism of action (MoA) classification task (STAGE 2). During evaluation/inference, features are extracted from the adaptation module's third layer and used for two tasks: (1) replicate retrieval — identifying wells treated with the same compound across experimental batches, and (2) MoA identification — detecting compounds that share the same mechanism of action.

fitting.

277

278 The first layer of the proposed model processes the embeddings from one of these pre-trained models through a 279 fully connected transformation (1024-dimensional output), 280 281 stabilizing activations with batch normalization, introduc-282 ing non-linearity with GELU, and applying dropout for 283 regularization. The second layer projects the transformed features into a 512-dimensional latent space, retaining the 284 285 same activation and normalization mechanisms. To enhance gradient flow and preserve learned representations, a resid-286 ual connection is incorporated in the third layer, allowing 287 288 the input to be directly added to the output before further transformation. The fourth layer (256-dimensional output) 289 290 refines the feature representation before passing it to the final classification layer, which is optimised for a specific 291 downstream task. In our experiments, this final layer was 292 trained either for compound classification (a single-label 293 task) or for Mechanism of Action (MoA) prediction, for-294 295 mulated as a multi-label classification problem to account 296 for cases where a single cell may be associated with multiple MoAs. More broadly, the TRex framework is com-297 patible with other downstream objectives, depending on the 298 299 structure of the final classification head and the chosen optimisation criterion. 300

301 2.2. Key Design Principles of TRex

The TRex framework is built around several core principles
that distinguish it from prior approaches to morphological
profiling:

Feature Refinement Rather than Direct End-to-End
 Training – Rather than training directly from raw im ages, TRex operates on high-quality embeddings pro duced by large self-supervised models. This approach

preserves expressive morphological information while reducing redundancy and computational cost..

- Task-Specific Optimisation The adaptation module is explicitly trained using a task-specific loss (e.g., MoA classification), ensuring that the learned representation is aligned with biologically meaningful distinctions rather than generic morphological similarity.
 Multi-Label Adaptability – TRex naturally supports
- 3. Multi-Label Adaptability TRex naturally supports multi-label classification, enabling the representation of complex phenotypes such as compounds with multiple mechanisms of action a common scenario in biological data.
- 4. Computational Efficiency and Modularity By freez-321 ing the upstream feature extractor and training only the 322 lightweight adaptation module, TRex achieves efficient 323 task adaptation with minimal labelled data and compute 324 resources. In contrast to full SSL training, which re-325 quires multi-GPU server infrastructure and days of com-326 pute time, the TRex adaptation module can be trained 327 in under an hour on a single consumer-grade GPU (e.g., 328 NVIDIA RTX 3090). This makes it practical for scaling 329 across new tasks, cell lines, or experimental conditions. 330

3. Experimental Evaluation

3.1. Datasets

Our training and evaluation were performed on a selec-
tion of 20 plates from the JUMP-CP Pilot dataset, repre-
senting two distinct human cell lines: U2OS, which con-
sists of human bone osteosarcoma epithelial cells, and
A549, which consists of human alveolar basal epithelial
cells. To ensure a representative distribution of morpholog-
ical variations, separate subsets of plates were designated333
336

340 for training and testing for each cell line. For the A549 341 cell line, the training set included plates BR00116991, BR00116992, BR00117050, BR00117052, BR00117054, 342 and BR00117055, while the testing set comprised plates 343 BR00116993, BR00116994, BR00117008, BR00117009, 344 BR00117051, and BR00117053. For the U2OS cell 345 line, the training set consisted of plates BR00116995, 346 BR00117013, and BR00117025, while the test set in-347 348 cluded plates BR00117011, BR00117012, BR00117024, and BR00117026. 349

There were 303 different perturbations present, representing 77 unique usable modes of action. We followed the MoA definitions and evaluation methodology defined in [6].

353 The images of selected plates were pre-processed and segmented into individual cells using DeepProfiler. Default 354 parameter values were used for segmentation to maintain 355 generability and consistency with prior studies. Following 356 segmentation, the data set consisted of a total of 1,079,388 357 358 cells for training and 1,144,682 cells for testing in the U2OS 359 cell line, and 3,179,429 cells for training and 2,626,695 cells for testing in the A549 cell line. These single-cell 360 images were subsequently used for feature extraction and 361 downstream performance analyses in our study. 362

363 3.2. Learning Framework

To address class imbalance in the dataset, the model is 364 trained using focal loss, which dynamically adjusts the loss 365 contribution to emphasize hard-to-classify samples. The 366 training process runs for 100 epochs, utilizing the AdamW 367 368 optimizer with an initial learning rate of 1e-3, which is reduced by a factor of 0.5 if the validation mAP does not im-369 370 prove for five consecutive epochs. Early stopping is applied if no improvement is observed for 10 consecutive epochs, 371 ensuring efficient convergence while preventing overfitting. 372

373 For evaluation of the JUMP-CP test dataset, we extract the output from layer 3 of our model, generating 374 a 512-dimensional embedding that serves as the learned 375 feature representation. This representation supports two 376 key downstream tasks: replicate retrieval, which identifies 377 wells treated with the same compound across experimental 378 batches, and MoA identification, which detects compounds 379 that share the same mechanism of action (MoA). The archi-380 381 tecture design of the proposed method is presented in Figure 382 1.

383 To generate well-level profiles, single-cell embeddings from our model are first aggregated into Field of View 384 (FOV) profiles using mean aggregation, followed by a sec-385 386 ond mean aggregation step to obtain well-level representa-387 tions. Post-processing includes Principal Component Anal-388 ysis (PCA) for dimensionality reduction, followed by Median Absolute Deviation (MAD)-robust standardization to 389 390 enhance data robustness.

Well-level profiles were used to evaluate replicate re-

391

trieval. For mechanism of action (MoA) prediction, well-392level profiles from the same compound treatment were av-393eraged across plates to generate a consensus profile for each394compound. Mean average precision (mAP) was calculated395using cosine similarity between well-level and consensus396profiles, for replicate retrieval and MoA identification, respectively.398

3.3. Experimental Evaluation

We performed a series of experiments to evaluate the ef-400 fectiveness of the proposed TRex framework. First, we as-401 sessed the impact of the training objective by comparing 402 the models trained for compound classification versus the 403 mechanism of action classification. Next, we applied TRex 404 to three different base representations: one obtained via 405 weak supervision (DeepProfiler) and two derived from dis-406 tinct self-supervised learning strategies (DINO, SubCell). 407 This analysis provides insight into whether self-supervised 408 representations may be more effective for downstream bio-409 logical tasks. Finally, we investigated the generalizability of 410 the method by training on datasets restricted to a single cell 411 type and evaluating its performance on data from an unseen 412 cell type. 413

Training on compound recognition. In the first set 414 of experiments, we trained TRex on the task of compound 415 recognition (Table 2). Performance was evaluated on both 416 compound recognition and MoA classification tasks, using 417 datasets for individual cell lines as well as the combined 418 dataset. We note that TRex offers significant improvements 419 in compound recognition across all representations, provid-420 ing gains of +0.09 for A549 and +0.05 for U2OS. The 421 gain for the combined dataset is 0.07. We observe that 422 self-supervised representations offer some limited perfor-423 mance gains over DeepProfiler, and these gains are main-424 tained when using TRex. The best performing combination, 425 TRex + SubCell, achieved compound recognition of 0.34, 426 compared to 0.29 for the SubCell representation alone. We 427 also note that training on compound recognition results in 428 only marginal improvements for MoA classification, with 429 gains between 0.01 and 0.02. The best MoA result is 0.17, 430 which may not be sufficient for many applications. This 431 suggests that compound-based training does not promote 432 biologically relevant representations and is suboptimal for 433 MoA tasks, and the mAP for MoA remains disappointingly 434 low. 435

Training with the MoA objective. In contrast, train-436 ing with the MoA objective leads to a substantial improve-437 ment in MoA classification performance, effectively more 438 than doubling the mAP score from 0.15 to 0.32 for TRex + 439 DINO (Table 3). We note that DINO offered the best base 440 representation for distillation by TRex, outperforming TRex 441 + SubCell by 0.05 on the combined cell dataset. Interest-442 ingly, this approach also results in a notable improvement 443

Method	Cmpd A549	MoA A549	Cmpd U2OS	MoA U2OS	Cmpd Both	MoA Both
DeepProfiler	0.35	0.15	0.24	0.15	0.22	0.14
DINO	0.38	0.16	0.27	0.16	0.26	0.15
SubCell	0.36	0.15	0.29	0.16	0.27	0.14
TRex + DeepProfiler	0.44	0.15	0.30	0.16	0.29	0.16
TRex + DINO	0.47	0.18	0.32	0.16	0.33	0.17
TRex + SubCell	0.45	0.17	0.34	0.17	0.34	0.17

Table 2. Compound recognition (Cmpd) and MoA classification performance (mAP) for models trained with compound supervision on both A549 and U2OS cell lines.

Table 3. Compound recognition (Cmpd) and MoA classification performance (mAP) for models trained with MoA supervision on both A549 and U2OS cell lines.

Method	Cmpd A549	MoA A549	Cmpd U2OS	MoA U2OS	Cmpd Both	MoA Both
DeepProfiler	0.35	0.15	0.24	0.15	0.22	0.14
DINO	0.38	0.16	0.27	0.16	0.26	0.15
SubCell	0.36	0.15	0.29	0.16	0.27	0.14
TRex + DeepProfiler	0.42	0.27	0.28	0.25	0.29	0.26
TRex + DINO	0.44	0.34	0.30	0.28	0.30	0.32
TRex + SubCell	0.42	0.28	0.32	0.25	0.31	0.27

in compound classification performance for TRex across 444 445 all representations. The improvement was most signifi-446 cant for DeepProfiler (0.07), followed by +0.04 for the self-supervised representations. These findings suggest that 447 training with an MoA objective using the proposed multi-448 label formulation provides a more generalizable solution, as 449 it enhances both MoA classification and compound recogni-450 451 tion. This indicates that when tackling tasks related to MoA or compound classification, MoA-driven training should be 452 preferred, as it leads to greater performance gains in both 453 objectives. 454

Generalisation across compounds and cell lines. To 455 evaluate the generalisation properties of the TRex frame-456 work, we conducted experiments in two settings: general-457 isation to unseen chemical compounds and generalisation 458 to an unseen cell line. Training with the MoA objective 459 provides a natural starting point, as 111 compounds are not 460 461 used in TRex MoA training due to the lack of usable MoA annotations. Table 4 shows the mAP for the compounds 462 that were unseen during the adaptation stage. As shown, 463 TRex improves performance even on these previously un-464 seen compounds, increasing mAP from 0.24 to 0.28.

Table 4. Compound recognition performance (mAP) for compounds unseen during TRex adaptation, using MoA-based training

Method	Cmpd (unseen)
DINO	0.24
TRex + DINO	0.28

For the second setting, we trained the adaptation module 465 with MoA objective using data from only one of the two 466 available cell lines (A549 or U2OS), and evaluated perfor-467 mance on both (Table 5). While TRex provides a clear per-468 formance boost on the cell line it was trained on, we observe 469 no significant improvement on the unseen cell line. This 470 outcome is not surprising: since training was performed on 471 a single cell line, the adaptation module had no exposure 472 to examples from other cellular contexts and therefore no 473 basis for learning cell-line-invariant MoA features. In other 474 words, we would not expect generalisation to emerge from 475 a single-cell-line training setup. TRex is intended to spe-476 cialise general-purpose features for a specific downstream 477 task and cellular context, optimising performance where it 478 is needed most. 479

Table 5. Compound recognition (Cmpd) and MoA performance (mAP) for models trained on a single cell line and evaluated on both seen and unseen cell lines. TRex was trained with MoA supervision.

Method	Train Cell	Eval Cell	Cmpd	MoA
DINO	-	A549	0.38	0.16
TRex + DINO	A549	A549	0.44	0.36
TRex + DINO	U2OS	A549	0.36	0.16
DINO	_	U2OS	0.27	0.16
TRex + DINO	U2OS	U2OS	0.30	0.31
TRex + DINO	A549	U2OS	0.25	0.17

533

534

535

536

537

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

In future work, we plan to investigate whether training
TRex on multiple cell lines can support better generalisation
to unseen cellular contexts by encouraging the adaptation
module to learn cell-line-invariant biological features.

484 4. Conclusions

485 In this work, we addressed the limitations of self-supervised learning (SSL) for morphological profiling by introducing a 486 two-stage learning framework that refines generic SSL fea-487 tures for task-specific objectives. While SSL provides scal-488 able feature extraction, our results confirm that its represen-489 490 tations are suboptimal for important tasks such as mecha-491 nism of action (MoA) classification. This is because a single, task-agnostic model cannot effectively satisfy conflict-492 493 ing requirements across different end tasks.

By introducing a lightweight transformation network 494 that adapts SSL features to specific predictive tasks, TRex 495 enables efficient, task-aware optimization without requir-496 ing the retraining of large-scale SSL models. Our exper-497 iments on a dataset of 20 Cell Painting plates from two 498 cell lines demonstrate that training for MoA prediction in 499 a multi-label setting significantly improves both MoA clas-500 501 sification and compound recognition, highlighting the benefits of learning task-specific refinements over purely static 502 503 SSL representations.

Because TRex is computationally efficient and can be 504 505 trained with fewer labeled examples, it provides a practical and scalable solution for large-scale drug discovery ap-506 plications. More broadly, our findings suggest that self-507 508 supervised feature representations should not be treated as 509 fixed but rather as a foundation for adaptive transforma-510 tions that better align with specific biological objectives. 511 Future work could explore extending TRex to other high-512 content imaging assays and integrating additional domain-513 specific constraints to further improve biological inter-514 pretability.

515 References

- 516 [1] Mark-Anthony Bray, Shantanu Singh, Han Han, Chad517 wick T Davis, Blake Borgeson, Cathy Hartland, Maria Kost518 Alimova, Sigrun M Gustafsdottir, Christopher C Gibson, and
 519 Anne E Carpenter. Cell painting, a high-content image-based
 520 assay for morphological profiling using multiplexed fluores521 cent dyes. *Nature Protocols*, 11(9):1757–1774, 2016. 1
- [2] Anne E. Carpenter, Thouis R. Jones, Michael R. Lamprecht,
 Colin Clarke, In Han Kang, Ola Friman, David A. Guertin,
 Jason H. Chang, Ruth A. Lindquist, Jason Moffat, Polina
 Golland, and David M. Sabatini. Cellprofiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biology*, 7(10):R100, 2006. 1
- 528 [3] Srinivas Niranj Chandrasekaran, Beth A Cimini, Amy
 529 Goodale, Lisa Miller, Maria Kost-Alimova, Nasim Jamali,
 530 John G Doench, Briana Fritchman, Adam Skepner, Michelle
 531 Melanson, Alexandr A Kalinin, John Arevalo, Marzieh

Haghighi, Juan C Caicedo, Daniel Kuhn, Desiree Hernandez, James Berstler, Hamdah Shafqat-Abbasi, David E Root, Susanne E Swalley, Sakshi Garg, Shantanu Singh, and Anne E Carpenter. Three million images and morphological profiles of cells treated with matched chemical and genetic perturbations. *Nature Methods*, 21(6):1114–1121, 2024. 2

- [4] Beth A Cimini, Srinivas Niranj Chandrasekaran, Maria Kost-538 Alimova, Lisa Miller, Amy Goodale, Briana Fritchman, 539 Patrick Byrne, Sakshi Garg, Nasim Jamali, David J Lo-540 gan, John B Concannon, Charles-Hugues Lardeau, Elizabeth 541 Mouchet, Shantanu Singh, Hamdah Shafqat Abbasi, Peter 542 Aspesi, Justin D Boyd, Tamara Gilbert, David Gnutt, San-543 tosh Hariharan, Desiree Hernandez, Gisela Hormel, Karolina 544 Juhani, Michelle Melanson, Lewis H Mervin, Tiziana Mon-545 teverde, James E Pilling, Adam Skepner, Susanne E Swalley, 546 Anita Vrcic, Erin Weisbart, Guy Williams, Shan Yu, Bolek 547 Zapiec, and Anne E Carpenter. Optimizing the cell painting 548 assay for image-based profiling. Nature Protocols, 18(7): 549 1981–2013, 2023. 1 550
- [5] Michal Doron, Thibault Moutakanni, Zihan S. Chen, et al. Unbiased single-cell morphology with self-supervised vision transformers. *bioRxiv*, 2023. Preprint. 2, 3
- [6] Ankit Gupta, Zoe Wefers, Konstantin Kahnert, Jan N. Hansen, Will Leineweber, Anthony Cesnik, Dan Lu, Ulrika Axelsson, Frederic Ballllosera Navarro, Theofanis Karaletsos, and Emma Lundberg. Subcell: Vision foundation models for microscopy capture single-cell biology. *bioRxiv*, 2024. 1, 2, 3, 5
- [7] Vladislav Kim, Nikolaos Adaloglou, Marc Osterland, Flavio M. Morelli, Marah Halawa, Tim König, David Gnutt, and Paula A. Marin Zapata. Self-supervision advances morphological profiling by unlocking powerful image representations. *Scientific Reports*, 15(1):4876, 2025. 1, 2
- [8] Nikita Moshkov, Malte Bornholdt, Sylvain Benoit, et al. Learning representations for image-based profiling of perturbations. *Nature Communications*, 15(1):1594, 2024. 1, 3
- [9] Srijit Seal, Maria-Anna Trapotsi, Ola Spjuth, Shantanu Singh, Jordi Carreras-Puigvert, Nigel Greene, Andreas Bender, and Anne E. Carpenter. Cell painting: a decade of discovery and innovation in cellular imaging. *Nature Methods*, 22(2):254, 2025. 1
- [10] Mathias Uhlén, Linn Fagerberg, Björn M. Hallström, Cecilia Lindskog, Per Oksvold, Adil Mardinoglu, Åsa Sivertsson, Caroline Kampf, Evelina Sjöstedt, Anna Asplund, IngMarie Olsson, Karolina Edlund, Emma Lundberg, Sanjay Navani, Cristina Al-Khalili Szigyarto, Jacob Odeberg, Dijana Djureinovic, Jenny Ottosson Takanen, Sophia Hober, Tove Alm, Per-Henrik Edqvist, Holger Berling, Hanna Tegel, Jan Mulder, Johan Rockberg, Peter Nilsson, Jochen M. Schwenk, Marica Hamsten, Kalle von Feilitzen, Mattias Forsberg, Lukas Persson, Fredric Johansson, Martin Zwahlen, Gunnar von Heijne, Jens Nielsen, and Fredrik Pontén. Tissue-based map of the human proteome. *Science*, 347(6220):1260419, 2015. 2