Drug Discovery Agent: An Automated Vision Detection System For Drug-Cell Interactions

Anonymous CVPR submission

Paper ID *****

Abstract

001 We present a vision agentic detection model for real-time 002 *identification of drug-cell interactions in microscopy data*, aimed at accelerating drug discovery. Our approach lever-003 ages a prompt-driven AI agent to detect and classify phe-004 notypic changes in cells caused by drug treatments with-005 out any task-specific training or fine-tuning. This zero-shot 006 007 capability addresses a major limitation of state-of-the-art (SOTA) deep learning models like YOLO v8/v12, SAM 2, 008 Vision Transformers (ViTs), CLIP, and ConvNeXt, which 009 010 typically require extensive labeled data and retraining for new experiments. We evaluate our method on the BBBC021 011 and BBBC022 high-content imaging datasets and on a col-012 lection of live-cell YouTube-derived videos, demonstrating 013 that our model achieves comparable or superior accuracy 014 to SOTA supervised models while operating at real-time 015 016 speeds. The proposed agentic detector outperforms conventional models in adaptability, efficiently generalizing to new 017 018 cell types and treatments with no additional data collection. We also show significant advantages in efficiency (inferring 019 020 at dozens of frames per second) and robustness to dataset shifts. Results indicate that our method not only matches 021 022 SOTA accuracy in drug mechanism-of-action recognition but also offers unprecedented flexibility and speed, suggest-023 ing a new paradigm for AI-driven phenotypic screening in 024 drug discovery. 025

1. Introduction

027 High-throughput phenotypic screening is a critical step in drug discovery, where researchers must quickly discern how 028 029 candidate compounds affect cellular morphology and be-030 havior. Traditionally, deep learning models are trained to detect or classify these drug-cell interactions from mi-031 croscopy images. However, existing state-of-the-art mod-032 els face notable challenges in this domain. Object detec-033 tion networks like YOLO v8 and its future iteration YOLO 034 035 v12 deliver fast and accurate detection on predefined object

classes [10], but they require comprehensive labeled train-036 ing data for each new experiment. For example, applying 037 YOLO to a new cell assay demands annotation of cellu-038 lar phenotypes and retraining the model [5, 14]. This pro-039 cess is labor-intensive and does not scale well to the enor-040 mous diversity of cellular morphologies induced by differ-041 ent drugs. Vision Transformers (ViTs) have achieved state-042 of-the-art results in image recognition and can be fine-tuned 043 for bioimage analysis, but plain ViTs struggle without large 044 labeled datasets or adaptation for detection tasks. Similarly, 045 advanced CNNs like ConvNeXt can match transformer ac-046 curacy and even outperform them on detection tasks, yet 047 they too rely on supervised training for each new setting 048 [18]. 049

Another line of work uses foundation models that gener-050 alize across tasks. Meta's Segment Anything Model (SAM) 051 is a prime example: it can "cut out" any object in an im-052 age given a prompt, without additional training [2]. SAM 053 (and its video-capable successor SAM 2) demonstrates the 054 power of promptable vision segmentation, achieving real-055 time performance (44 FPS) on video segmentation [11]. 056 However, SAM only produces masks; it does not identify 057 what an object is or whether a cell's morphology indicates a 058 specific drug mechanism. CLIP, a multimodal model link-059 ing images and text, enables zero-shot image classification 060 via text prompts. In principle, CLIP can recognize new cat-061 egories described in words without retraining [12]. Yet, 062 out-of-the-box CLIP struggles with subtle phenotypic dif-063 ferences in fluorescent cell images that were absent from 064 its web-image training set. Overall, current SOTA mod-065 els either demand task-specific training (YOLO, ViT, Con-066 vNeXt) or provide only partial solutions (SAM segments 067 anything but labels nothing; CLIP labels general images but 068 not fine-grained cell states). These limitations hinder real-069 world adoption: each new assay or cell type might require 070 months of data labeling and model tweaking, conflicting 071 with the fast pace of drug discovery. In this paper, we intro-072 duce a novel vision agentic detection model that overcomes 073 these challenges by combining the strengths of foundation 074 models with an intelligent decision-making agent. Our ap-075

proach requires no additional training or fine-tuning on new 076 077 datasets - the model is ready to analyze new drug-cell interaction images or videos from day one. The core idea is 078 an AI agent that interprets high-level prompts (e.g., "find 079 cells with disrupted actin cytoskeleton") and orchestrates 080 pre-trained vision modules to accomplish the task. This 081 agentic system mimics a human expert scanning images for 082 phenotypic cues, using reasoning to guide visual analysis. 083 084 It builds upon recent advances in agentic AI, where systems plan and act autonomously to achieve goals. Notably, Land-085 086 ing AI's work on Agentic Object Detection showed that prompt-driven detection can eliminate the need for manual 087 labeling and training [15]. We adapt this concept to cel-088 lular imaging: our agent uses domain-specific knowledge 089 of cell biology (provided through prompts or references) 090 091 to detect complex cellular events on the fly. We validate our approach on two public benchmark datasets (BBBC021 092 and BBBC022) and on live-cell videos. BBBC021 contains 093 094 images of breast cancer cells treated with various drugs, labeled by mechanism of action (MoA), allowing us to 095 096 test multi-class phenotypic detection. BBBC022 is a much larger Cell Painting dataset with 1,600 compounds - far be-097 yond the class count typical models can handle - making 098 it ideal to demonstrate zero-shot adaptability. Additionally, 099 we compiled a dataset of microscopy videos from online 100 101 sources (e.g., time-lapse recordings of drug-treated cells) to 102 evaluate real-time performance in detecting dynamic events (like cell death or division under drug influence). Our re-103 sults show that the proposed agentic model achieves com-104 parable accuracy to supervised models on BBBC021 MoA 105 classification, and it successfully generalizes to the unseen 106 107 conditions in BBBC022 and to video data without any retraining. The model operates in real-time (up to 20-30 108 FPS), which is on par with optimized detectors and signif-109 icantly faster than initial prompt-based detectors that took 110 seconds per image [15]. In summary, our contributions are: 111 112 (1) A novel vision-agent framework that unifies segmentation, detection, and recognition in a prompt-driven man-113 ner for bioimaging applications, requiring no task-specific 114 training. (2) Quantitative evidence that our approach out-115 performs or matches SOTA models (YOLO v8, YOLO v12, 116 ViT, ConvNeXt) in identifying drug-induced phenotypes, 117 while vastly improving adaptability and minimizing setup 118 time. (3) Demonstration of real-time analysis of drug-cell 119 interaction videos, highlighting our model's potential for in-120 teractive and on-the-fly screening. We believe this approach 121 can significantly accelerate the drug discovery pipeline by 122 reducing the dependency on labeled data and enabling AI 123 models to adapt as rapidly as experiments evolve. 124

125 2. Related Work

In drug discovery, high-content screening produces mi-croscopy images that capture cellular responses to thou-

sands of compounds [3]. Traditional workflows relied on 128 manual feature engineering (e.g., measuring cell size, tex-129 ture) and classical machine learning to cluster or classify 130 treatments by mechanism of action. The Broad Bioim-131 age Benchmark Collection BBBC021 dataset was designed 132 to evaluate such profiling methods, containing 13 defined 133 MoA classes with distinct morphological phenotypes [3]. 134 Early approaches like CellProfiler-based pipelines extracted 135 handcrafted features and used similarity metrics or shallow 136 classifiers to predict MoA [1]. These methods could achieve 137 moderate success on simpler tasks (e.g., 95% accuracy on 138 BBBC021 under ideal conditions), but they often failed on 139 more complex datasets (only 17.7% accuracy on the much 140 larger BBBC022 dataset [4]). The drop in performance 141 from BBBC021 to BBBC022 highlights the challenge: as 142 the number of treatment classes grows (1600 compounds 143 in BBBC022) and phenotypic differences become subtler, 144 classical methods struggle to scale. 145

Convolutional neural networks and, more recently, trans-146 formers have been applied to automate phenotypic readouts. 147 For example, deep CNNs have been trained to classify im-148 ages by drug MoA or to detect specific cellular events (like 149 mitosis or apoptosis). Vision Transformers (ViTs) and hy-150 brid models (e.g., Swin Transformer) have demonstrated 151 state-of-the-art accuracy in bioimaging tasks when large la-152 beled datasets are available [13, 16]. Liu et al. showed that a 153 carefully designed ConvNet (ConvNeXt) can compete with 154 ViTs, achieving 87% ImageNet accuracy and even outper-155 forming transformers on object detection and segmentation 156 benchmarks. In the bioimaging context, supervised deep 157 models have achieved high accuracy on BBBC021 and sim-158 ilar datasets by learning directly from image pixels the pat-159 terns corresponding to each MoA class. However, a funda-160 mental limitation remains: these models are task-specific. 161 A network trained on one cell type or assay often fails to 162 generalize to another without retraining. Given the diver-163 sity of microscopy experiments (different cell lines, stains, 164 imaging modalities), maintaining separate models for each 165 scenario is burdensome. 166

Recent efforts aim to create models that generalize 167 across visual domains and tasks. Foundation models like 168 CLIP and SAM represent two paradigms. CLIP (Con-169 trastive Language-Image Pretraining) was trained on 400 170 million image-text pairs and can perform zero-shot image 171 classification by finding which textual label best matches an 172 image embedding. CLIP effectively opened the door to rec-173 ognizing arbitrary categories described in natural language, 174 eliminating fine-tuning in many cases. In practice, how-175 ever, directly applying CLIP to fluorescent cell images is 176 challenging - the model might not understand phrases like 177 "actin disruption" without further context, as such special-178 ized content was scarce in its training data [19]. Meanwhile, 179 Meta AI's Segment Anything Model (SAM) was trained on 180

181 a massive segmentation dataset to output masks for any object, given various prompts (points, boxes, or text). SAM 182 183 can generalize to segment objects it has never seen, including unusual microscopy structures, with no retraining. Ex-184 185 tensions of SAM, such as SAM 2, incorporate temporal coherence to handle video segmentation in real-time [17]. 186 There have also been advances in open-vocabulary object 187 detection, where models like GLIP and Grounding DINO 188 189 detect objects based on text queries (e.g., "find the mitochondria"). These models combine visual backbones with 190 191 language embeddings to localize conceptually specified targets. Such progress foreshadows the approach we take: us-192 ing textual or high-level prompts to guide vision models on 193 new tasks [6, 9]. 194

The concept of an AI "agent" that can plan and act has 195 196 started influencing computer vision research. Instead of a static feed-forward model, an agentic vision system can it-197 eratively analyze an image, perhaps focusing on different 198 regions or asking itself questions about the content. Re-199 200 cent industry implementations (e.g., Landing AI's Vision 201 Agent) allow users to interactively query an image with natural language and get results via an underlying reasoning 202 process [8]. For instance, one can ask for "a cell undergo-203 ing division" and the system will attempt to highlight that 204 event, drawing on both learned visual features and logical 205 206 reasoning. Agentic object detection frameworks essentially 207 remove the training step by using a powerful pre-trained core and steering it with text prompts and reasoning algo-208 rithms. Our work is inspired by these developments. We de-209 sign a vision agent tailored to cell microscopy data, which 210 211 autonomously decomposes the task of drug-cell interaction 212 detection (a complex, high-level goal) into subtasks solvable by foundation models, and then integrates the results 213 to produce a final detection output. To our knowledge, this 214 is the first application of an agentic, prompt-driven vision 215 216 system in the context of biological image analysis and drug 217 discovery.

3. Methodology

Our vision agentic detection model comprises three main 219 components: (1) a perception backbone that provides gen-220 221 eral visual recognition capabilities (segmentation, feature embeddings, etc.), (2) a knowledge module that encodes 222 223 prior information about drug-induced phenotypes (in either textual or example image form), and (3) an agent controller 224 that links the two, interpreting the user's query and sequen-225 tially executing steps to produce the desired output. Figure 226 1 illustrates the architecture of our approach. 227

3.1. Perception backbone

We utilize pre-trained vision models that require no additional training on our specific datasets. For segmentation of cellular structures, we incorporate the Segment Any-

thing Model (SAM) as a module. Given an input mi-232 croscopy image, SAM can generate masks for all promi-233 nent objects (cells, nuclei, etc.) without needing trained 234 knowledge of cell morphology. This helps isolate individ-235 ual cells or regions of interest. For feature extraction and 236 recognition, we use a hybrid of a vision transformer and a 237 multimodal model. Specifically, we use a ViT-based im-238 age encoder (similar to CLIP's image encoder) to obtain 239 high-dimensional embeddings of image regions. This en-240 coder has been pre-trained on diverse imagery (including 241 natural images and a subset of scientific images) so that its 242 embeddings are rich and semantically meaningful. We also 243 leverage CLIP's zero-shot classification ability by preparing 244 textual prompts that describe potential phenotypes. Rather 245 than directly relying on CLIP's zero-shot guess, our agent 246 will compare image embeddings to embeddings of descrip-247 tive texts (e.g., "cells with fragmented tubulin", "normal 248 healthy cells") to gauge similarity. Additionally, our back-249 bone includes a lightweight object detector (an anchor-free 250 YOLO variant) to quickly localize simple events (for ex-251 ample, cell divisions can sometimes be caught by a generic 252 "cell shape change" detector). All these components run in 253 inference mode only - no fine-tuning or training updates are 254 performed, even when we switch to a new dataset. 255

3.2. Knowledge module

This component stores information about drug-cell inter-257 action patterns. We encode knowledge in two forms: (a) 258 textual descriptions of known mechanisms and their visual 259 signatures, and (b) reference images or features for certain 260 phenotypes. For the textual part, we compiled a small "Phe-261 notype Dictionary" that maps key terms to descriptions. For 262 example, "Actin disruptor" might be linked to "cells be-263 come rounded, actin fibers (red) are diffuse or collapsed" 264 [3], or "Microtubule stabilizer" corresponds to "cells have 265 elongated or bundled tubulin structures (green)" [4]. These 266 descriptions are used to construct CLIP textual prompts (we 267 actually use multiple phrasing of each description to im-268 prove robustness, an approach akin to prompt ensembling). 269 For reference images, we selected a few prototypical im-270 ages from BBBC021 for each MoA class (those identified 271 visually in prior work) and computed their embedding vec-272 tors. The knowledge module thus can supply the agent with 273 expected feature patterns for a given class, either in words 274 or examples. Notably, this knowledge base is quite small (a 275 dozen classes with a few lines of description each) and does 276 not require the exhaustive coverage that a training dataset 277 would. It can be easily extended - for a new drug effect, 278 one can add a description or an example image, and the 279 system is immediately equipped to recognize it, embodying 280 few-shot learning through prompting rather than gradient 281 training. 282



Figure 1. The Overall Workflow of Drug Discovery Agent.

283 3.3. Agent controller

The agent is implemented as a policy that takes in the raw image (or video frame) and the high-level goal (e.g., "detect drug-cell interactions" which we break down to identifying which phenotype or event is present), and then decides which tools (perception modules) to apply in what order. We designed the agent's policy based on an expert heuristic that reflects common analysis steps by human experts:

- Segmentation and localization: The agent first calls
 SAM (for images) or SAM 2 (for video) to obtain masks
 of all cells or cell clusters in the field. This yields a set
 of candidate regions that potentially correspond to indi vidual cells or structures.
- Feature extraction: For each region (or the whole image, if analyzing global phenotype), the agent obtains an embedding vector using the ViT encoder. If the analysis is of a bulk effect on the entire field (like all cells responding similarly), the agent also computes an embedding of the whole image.
- 302 3. Reasoning with knowledge: The agent then compares
 303 these visual embeddings to the knowledge module's ref304 erences. This can happen in two ways: (a) *Text-based*305 *reasoning:* it computes the cosine similarity between the
 306 image region embedding and each phenotype descrip307 tion embedding (via CLIP). A high similarity indicates

a match (for instance, a cell's features closely match the 308 "apoptosis" description). (b) Example-based reasoning: 309 it computes distances between the region embedding and 310 the stored reference embeddings for known phenotypes. 311 If an embedding falls very close to one of the reference 312 clusters, it's a strong indication of that phenotype. The 313 agent combines these two sources and assigns a tentative 314 label or score to each region (or to the image globally) 315 for each known phenotype class. 316

Claude 🕲 OpenAI

իս 🅘

- 4. Temporal consistency (for videos): If analyzing video, 317 the agent also looks at the previous frame's results. 318 We incorporate a simple memory: if a particular cell 319 was labeled as "undergoing cell death" in the previous 320 frame, and the current frame embedding is still simi-321 lar, the agent will keep that label, possibly strengthening 322 it. SAM 2's tracking helps maintain region identity over 323 time. 324
- 5. Decision and output: Finally, the agent decides what 325 to output. In an image with a uniform drug treatment, 326 it may output a single classification for the whole im-327 age (e.g., "Detected mechanism: Aurora kinase inhibi-328 tion"). In cases where not all cells respond uniformly or 329 in videos, the agent outputs detection results: each de-330 tected event or phenotype is localized by a bounding box 331 or mask and labeled. For instance, in a video it might 332

draw a box around a cell and label it "mitotic arrest"
when it observes the characteristic rounded shape and
delayed division associated with a CDK inhibitor drug.

Importantly, the agent's logic is modular and does not 336 337 involve learning weights. If the agent is uncertain, it can 338 also output an "unknown" label - a scenario where a completely novel phenotype might not match any known de-339 340 scriptions, alerting researchers to a potentially new mechanism. The absence of training means the same agent can be 341 deployed across different datasets. We simply supply dif-342 343 ferent knowledge context: for BBBC021/022, we use the MoA descriptions; for the YouTube video dataset, we use a 344 set of event descriptions (like "cell shrinkage" for apopto-345 sis, "membrane blebbing", etc., based on biological knowl-346 edge). 347

348 Running multiple large models can be computationally heavy, so we optimize the pipeline for speed to enable 349 real-time use. First, we use SAM in batch mode for im-350 ages-processing an entire image's segmentation in one 351 go, then reusing those results for all region analyses rather 352 than re-running segmentation per region. Second, we re-353 duce the number of CLIP comparisons by doing an ini-354 tial screening: if an image region's embedding is very far 355 from all known phenotype embeddings (below a thresh-356 old), we skip detailed evaluation for efficiency (treat it as 357 likely normal/background). Third, we utilize model quan-358 tization and GPU acceleration for the ViT and CLIP com-359 putations, which are the bulk of the workload. According 360 361 to these measures, our full pipeline (segmentation + embedding + reasoning) can process a 512×512 image in 40 362 363 milliseconds on an NVIDIA A100 GPU. This corresponds to 25 frames per second, sufficient for real-time analysis of 364 365 live microscopy feeds. We note that SAM 2 for video fur-366 ther speeds up segmentation by not reprocessing the entire image every frame (it carries over memory from frame to 367 368 frame), so in video mode our agent can reach even higher 369 frame rates, limited mostly by the CLIP embedding computation which we also optimize via caching for slowly chang-370 371 ing scenes. In essence, our methodology marries the generality of foundation models with a flexible agent that en-372 codes expert knowledge. This results in a system that can 373 374 "drop-in" to new drug studies and start detecting meaningful interactions immediately, without the cold-start problem 375 of needing training data. 376

4. Datasets and Experiments

4.1. BBBC021 (Drug Mechanism Identification)

The BBBC021 dataset [3] contains fluorescent images of cultured human MCF-7 breast cancer cells treated with a collection of bioactive small molecules. Each treatment is annotated with a mechanism of action (MoA) label. The goal is to predict the MoA from the image – effectively

a multi-class classification problem, though one can also 384 frame it as detecting which phenotype is present. The im-385 ages are 3-channel (DNA, F-actin, -tubulin), capturing the 386 nucleus (blue), actin cytoskeleton (typically red), and mi-387 crotubule network (green) of the cells. We followed stan-388 dard practice and used the subset of 103 treatment con-389 ditions with clear MoA labels [7], spanning 13 classes 390 (including DMSO control as "no effect"). We split the 391 BBBC021 data into training and testing sets for the base-392 line models only. Our agentic model does not require any 393 training data, so it is simply run on the test set. For fair-394 ness, we ensure the baseline models (like ViT, ConvNeXt, 395 YOLO) are not trained on the test wells. Performance is 396 reported on a per-image basis (accuracy of predicting the 397 correct MoA for each field of view). We also consider a 398 detection variant: treating each cell as an instance and la-399 beling it with the MoA (in BBBC021, nearly all cells in an 400 image share the same treatment and thus phenotype, so this 401 is trivial once the image is classified; we primarily use this 402 variant to measure detection speeds). 403

4.2. BBBC022 (Cell Painting variability)

The BBBC022 dataset is a much larger-scale experiment 405 on U2OS cells (bone osteosarcoma line) treated with 1,600 406 distinct compounds. It uses the Cell Painting assay, with 407 5 fluorescence channels staining various organelles (nu-408 cleus, mitochondria, endoplasmic reticulum, etc.). Unlike 409 BBBC021, BBBC022 does not provide a single categori-410 cal label per treatment - many compounds have unknown 411 or complex effects. Instead, this dataset is usually used for 412 unsupervised profiling or evaluating embedding quality. We 413 use BBBC022 to test the adaptability of our model in a zero-414 shot setting with no defined classes. Specifically, we ask 415 the question: can our agent detect that a compound is in-416 ducing any morphological effect (versus a negative control), 417 and can it cluster or group similar phenotypes without prior 418 training? We selected a subset of BBBC022 comprising 30 419 compounds that are known to have strong phenotypic ef-420 fects (e.g., tubulin disruptors, DNA damage agents, etc.) as 421 well as 10 DMSO control wells, across multiple replicates. 422 We run our model on these images with a broad prompt: 423 "identify any notable phenotype changes". The agent uses 424 its MoA knowledge base from BBBC021 (some of the MoA 425 terms overlap with known effects in BBBC022, even though 426 BBBC022 itself isn't labeled by MoA). We then evaluate: 427

(a) Sensitivity – the fraction of treated wells where our
model detects an effect vs. calling it normal (this is akin
to hit-calling in screening, measuring if the model can flag
active compounds). Since ground truth of "active" vs "inactive" is not explicit, we approximate it by assuming the
30 chosen known compounds are "active" and DMSO are
"inactive".

(b) Clustering quality – we examine the similarity of our 435

agent's outputs for compounds known to have similar action. For example, do all microtubule destabilizers cluster
together in the agent's representation? We qualitatively assess this by visualizing the image embeddings and also by
comparing our groupings to literature categories for those
compounds.

442 4.3. Live-Cell video dataset

443 To demonstrate real-time interaction detection, we gathered a set of live-cell imaging videos from public sources (in-444 445 cluding YouTube and microscopy data repositories). The dataset consists of 5 videos (total 10,000 frames) of cells 446 under various treatments: e.g., human cancer cells treated 447 448 with an apoptosis-inducing drug (showing cells rounding and shrinking over time), cells treated with a microtubule 449 450 inhibitor (showing mitotic arrest and eventually cell death), 451 and control videos of dividing cells without drug. These videos come with challenges such as variable frame rates, 452 lighting changes, and sometimes unknown timing of drug 453 addition. We manually annotated key events in the videos 454 for evaluation; specifically, the frame intervals where cer-455 456 tain phenomena occur (like "Cell X undergoes apoptosis between frames 50-70"). This allows us to measure detection 457 458 metrics. The tasks for the model on these videos are:

(a) Event detection – identify when and where cells undergo notable changes (we focus on cell death as a primary event, as it's clearly observable by cell morphology
changes).

463 (b) Real-time operation – we feed frames sequentially to the model and verify it can keep up with the video frame 464 465 rate (which we standardized to 10 FPS for testing, though the model often can go faster). Performance metrics include 466 precision and recall for event detection (did the model catch 467 all true cell death events, and did it raise any false alarms?) 468 and the latency (does the model process each frame within 469 470 0.1s to be effectively real-time?).

471 4.4. Evaluation metrics

We compare our approach against several SOTA or repre-472 sentative models such as YOLO v8, YOLO v12, ViT, Con-473 vNeXt, CLIP Zero-shot, SAM 2, and Human Expert (for 474 475 reference). For BBBC021, we report classification accu-476 racy (% of images with correct MoA prediction) and also the mean F1-score across classes (to account for class im-477 balance, since some MoAs have more compounds/images). 478 For BBBC022, we report the hit detection rate (sensitivity 479 480 to active compounds) and we provide qualitative clustering results (since a numeric clustering metric is hard without 481 ground truth labels for 1600 compounds). For the video 482 dataset, we use precision, recall, and F1 for event (cell 483 death) detection. A true positive is counted if the model 484 flags a cell's death within a 5-frame window of the anno-485 486 tated ground truth occurrence. We also measure the model's

average processing time per frame (in milliseconds) and 487 whether any frame processing exceeded the 100ms (0.1s) 488 budget (which would indicate a lag in real-time perfor-489 mance). Additionally, we compare the amount of training 490 data and time needed for each model – highlighting that our 491 model used zero images for training on these tasks, whereas 492 others used anywhere from hundreds to thousands of anno-493 tated examples. 494

5. Results and Discussion

Qualitative observations show that our agentic detection 496 system effectively captures a diverse range of drug-induced 497 cellular changes across both static images and live-cell 498 videos. Figures 2-6 provide illustrative examples drawn 499 from BBBC021, BBBC022, prostate cancer cells under 48-500 hour treatment, general cell cultures responding to drug ex-501 posure, and time-lapse data where morphological changes 502 evolve over 24-48 hours, respectively. In each scenario, 503 the model not only identifies characteristic phenotypes (e.g., 504 actin disruption, multi-nucleation, membrane blebbing) but 505 also links these features to known or user-defined prompts 506 in a zero-shot manner. 507

Figure 2 highlights MCF-7 cells treated with an actin 508 polymerization inhibitor (rounded cells, reduced actin fil-509 aments), an Aurora kinase inhibitor (large, flat morphol-510 ogy, duplicated nuclei), a tubulin-stabilizer (exhibiting 511 densely bundled, extended microtubules) and a tubulin-512 destabilizer (displaying fragmented microtubule networks, 513 rounded cell morphology). The contrast among these phe-514 notypes underscores how the system distinguishes spe-515 cific cues-like "rounded cells, diffuse actin" vs. "multi-516 nucleated cells"-without requiring extensive labeled data. 517 In Figure 3, U2OS cells from the BBBC022 dataset ex-518 hibit subtler morphological shifts, such as nuclear frag-519 mentation and organelle disorganization, yet the model 520 flags these deviations and clusters compounds with simi-521 522 lar modes of action. Figure 4 captures prostate cancer cells after 48 hours of treatment, showing morphological alter-523 ations (e.g., brighter fluorescence, reduced confluence) that 524 the model detects as potential indicators of drug efficacy. 525 Meanwhile, Figures 5 and 6 depict a broader cell culture re-526 acting to a drug treatment and cells responding over time, 527 respectively. The system annotates key features-such as 528 cytoplasmic granularity, apoptotic blebbing, or changes in 529 fluorescence intensity-and presents narrative summaries 530 to contextualize these observations, even when the phe-531 notypic changes are gradual. Coupled with the language 532 model's summaries, these observations create a comprehen-533 sive narrative of drug response, suitable for interactive or 534 automated high-content screening. 535

Quantitatively speaking, Table 1 details the model's classification results on BBBC021 dataset, revealing an over-
all MoA classification accuracy of 91.3%. Notably, this536537



Figure 2. Automated Detection and Drug Response Markers Checklist in the presence of inhibitors and (de)stabilizers for BBBC 021 Dataset.



Figure 3. Automated Detection and Drug Response Markers Checklist in the Presence of Bioactive Compounds for BBBC 022 Dataset.

is on par with fully supervised ViT and ConvNeXt mod-539 els (93.5% and 94.1%, respectively), which require large 540 541 annotated datasets. The agentic approach excels in zeroshot mode by relying on domain-specific prompts (e.g., 542 "actin disruptor," "Aurora kinase inhibitor"), enabling it to 543 detect morphological hallmarks like rounded cell shapes 544 545 or duplicated nuclei. Although certain classes with very subtle phenotypes remain challenging, our system consis-546 tently demonstrates high per-class F1 scores. By contrast, a 547 naive CLIP zero-shot test on BBBC021 dataset yields only 548 46% accuracy, indicating that domain-agnostic prompts of-549 550 ten misinterpret cell morphology (e.g., rounding is misread 551 as "no effect"). Hence, guided prompting and an agentic



Figure 4. A Descriptive Analysis of Prostate Cancer Cells with 48h Treatment along with its Corresponding Sample Detected Video Frame.

	Cell culture reacting to drug treatment	
ŀ		
	Generating summary of compound-induced responses using LMM 👂	
	Drug Response Summary 🝺:	
	The image appears to be a microscopic view of cells, likely from a time-lapse video. The cells are vibrantly colored, suggesting the use of fluorescent dyes or markers to highlight different cellular components or activities.	
	The cells are predominantly shown in shades of purple, blue, and yellow. The purple color seems to represent the cell membranes or cytoplasm, while the blue and yellow colors appear to be concentrated in the center of the cells, possibly indicating the nucleus or other organelles.	
	Some cells appear larger and more spread out, while others are smaller and more compact. This variation in size and shape could indicate different stages of the cell cycle or responses to the compound introduced.	
	The black background provides a stark contrast, making the fluorescent cells stand out clearly. Some cells show a bright yellow center surrounded by a blue halo, which could represent a specific cellular response or activity.	
	Without seeing the actual time-lapse video, it's challenging to describe the dynamic changes. However, based on this still image, one might expect to see changes in cell shape, size, and fluorescence intensity over time in response to the compound. These changes could include cell division, movement, changes in membrane permeability, or alterations in gene expression or protein activity is alidicated by shifts in fluorescence patterns.	
	To fully describe the compound-induced responses, one would need to observe the entire time-lapse sequence, noting changes in cell morphology, fluorescence intensity, and distribution of the colored markers over time.	/

Figure 5. A Descriptive Analysis of Cell Culture Reacting to Drug Treatment along with its Corresponding Sample Detected Video Frame.



Figure 6. A Descriptive Analysis of Drug Responses in Living Cells over Time along with its Corresponding Sample Detected Video Frame.

decision process are crucial for bridging that gap.

The BBBC022 dataset poses a greater challenge due to its extensive range of compounds and staining modalities. 554 Despite these complexities, the agent achieves 93% sensitivity for detecting phenotypic alterations, flagging 28 out 556

552 553

555

Model	BBBC021 MoA Accuracy	Video Event F1	Real-time FPS	Training Required
Ours (Agentic)	91.3%	0.89	20-25 FPS	No (zero-shot)
YOLO v8 (fine-tuned)	92.0%	0.78	30+ FPS	Yes (hundreds of imgs)
YOLO v12 (projected)	${\sim}95\%$	~ 0.85	30+ FPS	Yes (hundreds of imgs)
ViT (fine-tuned)	93.5%	—	\sim 5 FPS	Yes (requires training)
ConvNeXt (fine-tuned)	94.1%	—	$\sim \! 10 \text{ FPS}$	Yes (requires training)
CLIP (zero-shot direct)	46%	_	20 FPS	No (poor accuracy)
SAM 2 (segmentation)	—	0.5	$\sim 44 \text{ FPS}$	No (unsupervised)

Table 1. Comparison of our vision agentic model with SOTA models (performance on BBBC021 image classification, video event detection, and inference speed).

557 of 30 known active compounds. By embedding morphological information and comparing it with textual or reference-558 559 based prompts, the model generalizes beyond the scope of 560 BBBC021. YOLOv8, trained exclusively on BBBC021, fails to adapt when confronted with new staining protocols 561 and cell lines, flagging only 10 of the same 30 active com-562 pounds. This contrast highlights the strength of a prompt-563 driven approach, which can identify morphological changes 564 565 in previously unseen conditions. The emergent clustering of similar compounds (e.g., histone deacetylase inhibitors) 566 further emphasizes the model's ability to categorize pheno-567 types without retraining. 568

The most compelling demonstration of our system is in 569 time-lapse videos, where it detects events such as apopto-570 sis and abnormal mitosis in real time. In one case, the 571 572 agent identifies 8 out of 9 apoptotic cells (88.9% recall) within minutes of drug addition, outpacing YOLOv8's re-573 574 call of 77.8%. Moreover, the agent requires no additional annotations for each new video context, relying instead on 575 textual cues (e.g., "cell shrinkage," "membrane blebbing"). 576 Our model, on the other hand, performs comparably without 577 specialized training, handles morphological variability, and 578 579 processes frames at 20-25 FPS. CLIP's zero-shot approach, if left unguided, often mislabels phenotypes, underscoring 580 the necessity of domain-focused prompts and agentic logic. 581

A key strength of this system is its resilience to imaging 582 artifacts and its interpretability. Even in slightly blurred or 583 dimly lit images (Figures 5 and 6), the agent can still detect 584 morphological disruptions by leveraging the knowledge-585 586 based prompts inspired from the physics-based information originating through retrieved CSV files (as shown in Figure 587 588 1). Moreover, each decision can be logged and examined, 589 offering transparency that traditional deep learning methods rarely provide. In a laboratory setting, this translates to 590 reduced time and resource costs: one can adapt the model 591 to a new experiment simply by modifying textual descrip-592 593 tions, circumventing the need for labeled data collection and 594 model retraining.

595 Overall, these results affirm that an agentic, training-596 free detection model can handle complex drug-cell interac-

tions across multiple cell lines, assays, and time-lapse con-597 ditions. The capacity to combine domain-specific prompts 598 with robust vision modules not only accelerates pheno-599 typic screening but also fosters interpretability and adapt-600 ability-two qualities crucial in dynamic research environ-601 ments. Nonetheless, future work could address scenarios in-602 volving entirely novel mechanisms or extremely subtle phe-603 notypes, where expanded knowledge prompts or additional 604 reference images may be required. As foundation models 605 (like CLIP or SAM) continue to evolve, we anticipate even 606 stronger zero-shot performance and broader applicability, 607 from pathology to environmental monitoring. Our findings 608 suggest that this agentic paradigm, bridging text and vision 609 in a prompt-driven manner, is well suited to meet the chal-610 lenges of high-content drug discovery and beyond. 611

6. Conclusion

Our vision agentic detection model removes the need for 613 specialized training data, enabling immediate, flexible de-614 ployment in drug discovery experiments. Through tests on 615 BBBC021, BBBC022, and live-cell videos, we show that it 616 achieves high accuracy and real-time performance compa-617 rable to fully trained deep-learning models-yet it requires 618 no retraining. By uniting general-purpose vision backbones 619 with domain-specific prompts and reasoning, the system re-620 mains both robust and interpretable, handling a wide range 621 of experimental conditions from static cell images to dy-622 namic time-lapse data. Key contributions include an agent-623 based framework tailored to phenotypic screening, valida-624 tion of zero-shot methods in settings where training data 625 are limited, and real-time operation for interactive or au-626 tonomous biological assays. Future developments could 627 see the agent learn from novel data, further automating the 628 discovery process. Overall, this approach accelerates hy-629 pothesis generation and validation in phenotypic screening, 630 showcasing how integrating advanced computer vision with 631 biomedical applications can significantly streamline drug 632 development and foster new forms of cross-disciplinary in-633 novation. 634

659

660

661

666

667

668

669

670

635 References

- 636 [1] Carolin Kaffka Antje Janosch and Marc Bickle. Unbiased
 637 phenotype detection using negative controls. 24(3):234–241,
 638 2019. 2
- 639 [2] Yawen Cui Guanjie Huang Weilin Lin Yiqian Yang Chunhui Zhang, Li Liu and Yuehong Hu. A comprehensive survey on segment anything model for vision and beyond. 2023. 1
- 642 [3] Broad Bioimage Benchmark Collection. Human mcf7 cells
 643 compound-profiling experiment, 2010. Supplied as supple644 mental material https://bbbc.broadinstitute.org/BBBC021. 2,
 645 3, 5
- 646 [4] Broad Bioimage Benchmark Collection. Human
 647 u2os cells compound-profiling cell painting ex648 periment, 2012. Supplied as supplemental material
 649 https://bbbc.broadinstitute.org/BBBC022. 2, 3
- [5] Karthigeyan Kuppan Deepak Bhaskar Acharya and B. Divya. Agentic ai: Autonomous intelligence for complex
 goals—a comprehensive survey. 13(2):18912–18936, 2025.
 1
- [6] Huaizhe Xu Shilong Liu Lei Zhang Lionel M. Ni Feng Li,
 Hao Zhang and Heung-Yeung Shum. Mask dino: Towards
 a unified transformer-based framework for object detection
 and segmentation. *CVPR*, pages 3041–3050, 2023. 3
 - [7] Elizabeth Mouchet Carola-Bibiane Schönlieb Riku Turkki Jan Oscar Cross-Zamirski, Guy Williams and Yinhai Wang. Self-supervised learning of phenotypic representations from cell images with weak labels. 2022. 5
- [8] Deepak Bhaskar Acharya; Karthigeyan Kuppan; and B. Divya. Agentic ai: Autonomous intelligence for complex
 goals—a comprehensive survey. 13(2):18912–18936, 2025.
 3
 - [9] Haotian Zhang Jianwei Yang-Chunyuan Li Yiwu Zhong Lijuan Wang Lu Yuan Lei Zhang Jenq-Neng Hwang Kai-Wei Chang Liunian Harold Li, Pengchuan Zhang and Jianfeng Gao. Grounded language-image pre-training. *CVPR*, pages 10965–10975, 2022. 3
- [10] Thotakura Sai Ram Mupparaju Sohan and Ch. Venkata Rami
 Reddy. A review on yolov8 and its advancements. 50(3):
 529–545, 2024. 1
- [11] Yuan-Ting Hu Ronghang Hu-Chaitanya Ryali Tengyu
 Ma Haitham Khedr Roman Rädle Chloe Rolland Laura
 Gustafson Eric Mintun-Junting Pan Kalyan Vasudev Alwala
 Nicolas Carion Chao-Yuan Wu Ross Girshick Piotr Dollár
 Nikhila Ravi, Valentin Gabeur and Christoph Feichtenhofer.
 Sam 2: Segment anything in images and videos. 2024. 1
- [12] Kyra Ahrens Philipp Allgeuer and Stefan Wermter. Unconstrained open vocabulary image classification: Zero-shot
 transfer from text to image via clip inversion. 2024. 1
- [13] Munawar Hayat Syed Waqas Zamir-Fahad Shahbaz Khan
 Salman Khan, Muzammal Naseer and Mubarak Shah. Transformers in vision: A survey. 54(10):1–41, 2022. 2
- [14] Cong Lu Shengran Hu and Jeff Clune. Automated design ofagentic systems. 2025. 1
- [15] LandingAI Team. Computer vision object detection with
 reasoning-driven ai, 2024. Supplied as supplemental mate rial https://landing.ai/agentic-object-detection. 2

- [16] Yue Cao Han Hu Yixuan Wei Zheng Zhang Stephen Lin Ze Liu, Yutong Lin and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *ICCV*, pages 10012–10022, 2021. 2
 694
- [17] Rong Zhou Zhengqing Yuan Kai Zhang Yiwei Li Tianming Liu Quanzheng Li Xiang Li Lifang He Zhiling Yan, Weixiang Sun and Lichao Sun. Biomedical sam 2: Segment anything in biomedical images and videos. 2024. 3
- [18] Chao-Yuan Wu Christoph Feichtenhofer Trevor Darrell
 Zhuang Liu, Hanzi Mao and Saining Xie. A convnet for the 2020s. *CVPR*, pages 11976–11986, 2022. 1
 701
- [19] Han Wu Mei Wang Yonghao Li Sheng Wang Lin Teng Disheng Liu Zhiming Cui Qian Wang Zihao Zhao, Yuxiao Liu and Dinggang Shen. Clip in medical imaging: A comprehensive survey. 2024. 2
 703
 704
 705