

ONE SNAPSHOT, MANY CLUES: INVERSE PROTOCOL PREDICTION FROM SINGLE-VIEW SPHEROID IMAGES

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding how experimental protocols shape spheroid morphology is crucial for advancing 3D cell culture research, yet reconstructing these conditions from imaging alone has remained elusive. We used deep learning frameworks which are able to infer the full experimental protocol including cell line, medium, seeding density, timepoint, formation method, microscope, and magnification from a single bright-field spheroid image. Using the SLiMIA dataset of 8,000 annotated images spanning diverse culture conditions, we cast this as a structured multi-label prediction task and benchmarked a spectrum of models, from CNNs and transformers to hybrid and dependency-aware architectures. Our approach integrates segmentation for morphology extraction, domain-adversarial training, and morphologically informed augmentation to improve robustness across imaging setups. Results show an average accuracy of 95.23% across protocol components, with hybrid models such as CoAtNet excelling in balancing efficiency and accuracy, while feature-augmented and hierarchical models contribute interpretability and consistency. Grad-CAM analyses confirm that predictions rely on biologically meaningful features (e.g., compactness, necrotic core structure), while highlighting dataset-driven artifacts in replicate and magnification tasks.

1 INTRODUCTION AND BACKGROUND

Three-dimensional (3D) cell culture systems such as spheroids and organoids are central tools in cancer biology, drug discovery, and tissue engineering (Fatehullah et al., 2016; Chatzinikolaidou, 2016). Unlike two-dimensional cultures, spheroids recreate nutrient and oxygen gradients, cell-cell interactions, and necrotic cores, while remaining amenable to bright-field and phase-contrast microscopy for high-throughput, non-destructive readouts (Edmondson et al., 2014). Despite this ubiquity, imaging is still used primarily to measure outcomes (size, viability, morphology), not to infer the experimental protocols that produced them. We pose a new challenge: given a single spheroid image, can we reconstruct the protocol conditions—including cell line, medium, seeding density, timepoint, formation method, and microscope? We term this the inverse protocol prediction problem. Success would enable reproducibility checks, automated experiment validation, and deeper insight into how protocol choices manifest morphologically.

Recent advances in biomedical image segmentation have leveraged deep convolutional neural networks (CNNs) and transformer architectures to handle the complexity of 3D cell culture images. U-Net variants and DeepLab models achieve high segmentation accuracy for spheroid morphology under challenging visual conditions such as noise and imaging artifacts (Park et al., 2023; Liu et al., 2021). Transformer-based architectures like Swin-UNet enhance segmentation by capturing global contextual cues and improving recognition of intricate morphological features (Khan et al., 2023). Semi-supervised learning and multi-label deep supervision have been employed to overcome limitations imposed by scarce annotations and class imbalance, improving model robustness and generalization (Reiß et al., 2021; Han et al., 2024). Synthetic training data generated based on biophysical principles has further augmented data diversity, yielding improved alignment with real-world spheroid structures (Koetzier et al., 2024). In the context of IPP from biomedical images, recent studies have adopted transformer-based and multi-label deep learning frameworks to infer experimental protocols or biological states with high fidelity. Multi-label learning models effectively capture dependencies between protocol components, while transformer architectures provide interpretability and robust feature extraction for heterogeneous biomedical data (Zhang et al.,

2022; Madan et al., 2024). These data-driven methods provide a principled alternative to classical probabilistic IPP, leveraging complex spatial and contextual cues for accurate prediction of underlying experimental conditions. Spatiotemporal modeling for time series prediction in biomedical microscopy has progressed through convolutional LSTM and attention-based deep learning models. Integrating 3D cell culture systems with advanced AI frameworks enables dynamic prediction of morphological progression and treatment response (Dave et al., 2025; Torro et al., 2025). Models such as ConvLSTM and physics-inspired recurrent networks demonstrate strong performance in capturing temporal dependencies in longitudinal microscopy data (De Cillis et al., 2025; Mali et al., 2025). Leveraging metadata and multi-modal inputs further enhances these models’ ability to generalize across diverse experimental conditions.

Inverse protocol prediction (IPP) is difficult for several reasons. First, *morphological ambiguity* arises when distinct protocols yield visually similar spheroids. Second, *imaging variability* across microscopes and magnifications introduces artifacts that obscure morphology. Third, few datasets couple high-resolution spheroid images with detailed experimental metadata, limiting model development and evaluation. To address this, we introduce a unified framework for inverse protocol prediction using SLiMIA, a dataset of $\sim 8,000$ bright-field spheroid images spanning nine microscopes, 47 cell lines, multiple media, seeding densities, timepoints, and formation methods (Blondeel et al., 2025). We frame the problem as a structured multi-label prediction, disentangle morphology from imaging artifacts via domain-adversarial training and augmentation, and benchmark convolutional, transformer, and hybrid architectures. While many components are established, our novelty lies in their adaptation to the structured multi-label nature of experimental protocols, where causal label dependencies and morphometric priors are biologically grounded rather than arbitrary.

Our main contributions are summarized as:

1. Morphometry fusion that provides integration of classical shape features (e.g., area, compactness) with deep embeddings.
2. Hierarchical modeling for a multi-task transformer that explicitly captures dependencies among protocol labels.
3. Grad-CAM analysis at protocol-level for interpretability that highlights morphological regions and dataset artifacts.
4. The first temporal modeling of SLiMIA via ConvLSTM, PredRNN++, and PhyDNet to predict the evolution of spheroid (Spatio-temporal extension).

This work demonstrates that spheroid morphology encodes rich, recoverable signatures of culture conditions, establishing microscopy-driven inverse protocol prediction as a new paradigm for reproducibility, optimization, and design in cell culture systems.

2 MATERIALS AND METHODS

SLiMIA (Spheroid Light Microscopy Image Atlas) is an open-access morphometric image dataset designed to support machine learning and computational modeling of three-dimensional (3D) cell culture systems. SLiMIA comprises approximately 8,000 light microscopy images of spheroids spanning a diverse range of experimental conditions, including nine microscope types, 47 distinct cell lines, eight culture media, four spheroid formation protocols, and multiple cell seeding densities, accompanied by rich metadata to facilitate reproducible analysis and benchmarking (Blondeel et al., 2025). Figure 1 highlights some sample dataset images along with their manual segmentations. Figure 2 demonstrates the entire workflow followed for analysis of the dataset. The same is described in detail in this section.

All segmentation models were trained on RGB inputs with a single binary output mask, where sigmoid activation was applied within the loss or metric functions. We adopted Adam-based optimizers (Adam et al., 2014) and employed a ReduceLROnPlateau scheduler with a patience of 5 and decay factor of 0.5 to adaptively lower the learning rate when validation performance plateaued (Smith, 2017). The choice of adaptive optimizers ensured stable convergence across diverse architectures, while dynamic scheduling helped prevent overfitting. Loss functions were selected to balance pixel-wise accuracy and region overlap, with particular attention to the class imbalance and faint boundaries common in spheroid images. U-Net++ and DeepLabV3+ models employed the Focal Tversky

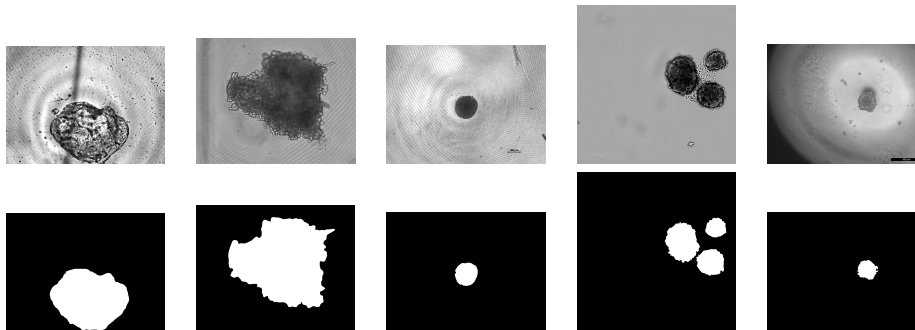


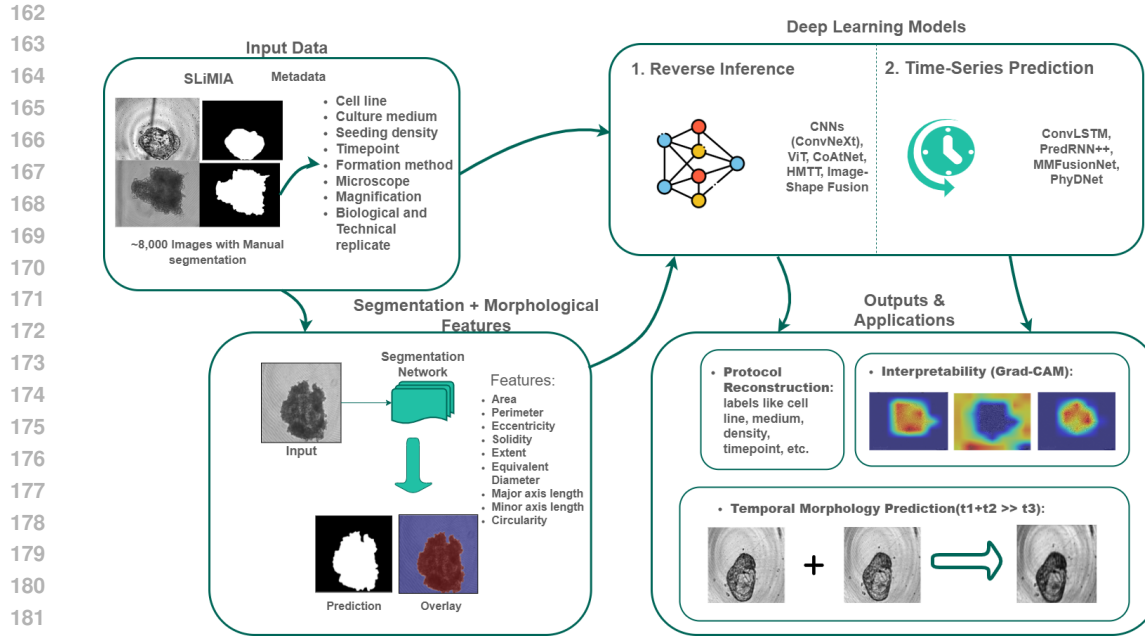
Figure 1: Sample images from the SLiMIA dataset with their manual segmentations (bottom row).

Loss (Salehi et al., 2017) ($\alpha = 0.7, \beta = 0.3, \gamma = 0.75$), chosen for its ability to emphasize recall and mitigate false negatives in imbalanced biomedical data. U-Net++ used a ResNet-50 encoder pre-trained on ImageNet, following (Zhou et al., 2018), while DeepLabV3+ also leveraged a ResNet-50 backbone to provide stronger representational capacity compared to shallower variants (Chen et al., 2018a). DeepLabV3 and Attention U-Net were trained with a hybrid loss combining BCEWithLogits and Dice (0.5 each), encouraging overlap accuracy while maintaining pixel-level precision (Milletari et al., 2016). To better regularize the optimization, we adopted AdamW (Loshchilov & Hutter, 2017) (learning rate = $3e-4$, weight decay = $1e-4$), which stabilizes training in models with high parameter sensitivity. DeepLabV3 used a ResNet-34 backbone, while Attention U-Net introduced attention gates to focus on spheroid boundaries against noisy microscopy backgrounds (Oktay et al., 2018). For transformer-based architectures, Swin-UNet and TransUNet, we used BCEWithLogitsLoss, which proved more stable for patch-based embeddings and long-range dependencies (Cao et al., 2022; Chen et al., 2021). Both models relied on Adam with a $1e-4$ learning rate, and their learning rates were reduced when validation loss plateaued. Swin-UNet incorporated a Swin-Tiny backbone with 768-dimensional features and a convolutional decoder (Cao et al., 2022), while TransUNet employed a Vision Transformer encoder followed by a CNN decoder stack (Chen et al., 2021). This setup ensured a consistent training environment while allowing each model’s design-specific strengths to be leveraged. In particular, models with attention or transformer backbones benefited from tailored loss and optimizer choices, whereas CNN-based variants relied on class imbalance-aware formulations like Focal Tversky. RefineNet leveraged a ResNet-34 backbone where multi-resolution features were progressively refined through cascaded residual units and chained residual pooling blocks, enabling improved boundary preservation and contextual integration (Lin et al., 2017). It was trained with the Focal Tversky Loss (Salehi et al., 2017) using the Adam optimizer (Adam et al., 2014) (learning rate = 1×10^{-4}), with a ReduceLROnPlateau scheduler (patience = 5, factor = 0.5). Data augmentations such as flips, brightness/contrast adjustments, and gamma correction were applied to improve robustness. SegNet adopted a VGG-style encoder-decoder with max-pooling indices reused in unpooling layers, ensuring spatial detail recovery while keeping the design lightweight (Badrinarayanan et al., 2017). It was also optimized with the Focal Tversky Loss and Adam, combined with ReduceLROnPlateau scheduling and extensive augmentations including flips, elastic deformations, noise injection, and geometric transformations.

2.1 IMPLEMENTATION DETAILS OF INVERSE PROTOCOL PREDICTION (IPP) MODELS

We formulate IPP as a structured multi-label prediction task, aiming to recover the experimental protocol (microscope, cell line, culture medium, formation method, seeding density, timepoint, replicate identifiers, magnification) directly from a single spheroid image. All models use RGB inputs resized to 224×224 and a multi-head design with one classifier per label, trained with categorical cross-entropy and optional class weighting. Model choice reflects three principles: capturing local morphology vs. global structure, integrating explicit morphometric priors, and modeling dependencies between protocol components.

ConvNeXt-Tiny served as a convolutional baseline with ImageNet initialization, nine classification heads, Adam (1×10^{-4}) with ReduceLROnPlateau, and light augmentations (flips, rotations). Its modern convolutional blocks preserve locality priors, useful for fine-grained cues such as medium,



183 Figure 2: Overview of the workflow. Input images and metadata are used for segmentation, feature
184 extraction, and deep learning models. Outputs include protocol reconstruction, temporal morphology
185 prediction, and Grad-CAM interpretability.

188 magnification, and seeding density (Liu et al., 2022; Mmileng et al., 2025). In contrast, ViT-B/16
189 leverages global self-attention. Using pretrained timm weights and the CLS token for task-specific
190 heads, it better models long-range dependencies important for formation method and timepoint pre-
191 dictions (Asiri et al., 2023). CoAtNet-0, combining convolution and attention, was trained from
192 scratch but converged reliably with Adam and scheduling. Its hybrid design balances local texture
193 and global structure, aiding parameters like medium and seeding density (Dai et al., 2021). To incor-
194 porate explicit priors, the fusion model augments ConvNeXt-Tiny embeddings with nine normalized
195 shape descriptors (area, perimeter, eccentricity, solidity, extent, equivalent diameter, axis lengths,
196 circularity). Shape tokens are concatenated with image features and processed by a lightweight
197 Transformer ($d_{\text{model}} = 256$, 3 layers, 4 heads). Predictions from this joint embedding, optimized
198 with AdamW and weighted cross-entropy, improve robustness and interpretability (Xia et al., 2025;
199 Luo et al., 2025; Sun, 2025). Hierarchical Multi Task Transformer (HMTT) (Figure 3) enforces causal
200 ordering among labels (cell line \rightarrow medium \rightarrow seeding density \rightarrow magnification \rightarrow microscope \rightarrow
201 timepoint \rightarrow replicates). A ViT-B/16 encoder provides a shared embedding, with sequential heads
202 predicting attributes in order. Training used AdamW (1×10^{-4} , weight decay 1×10^{-2}), with class-
203 weighted cross-entropy and focal loss ($\gamma = 2.0$, $\alpha = 0.25$) for imbalanced labels. By conditioning
204 predictions, HMTT maintains biologically consistent outputs (e.g., medium dependent on cell line),
205 yielding more plausible reconstructions (Rafeian & Vázquez, 2025; Tarekegn et al., 2024).

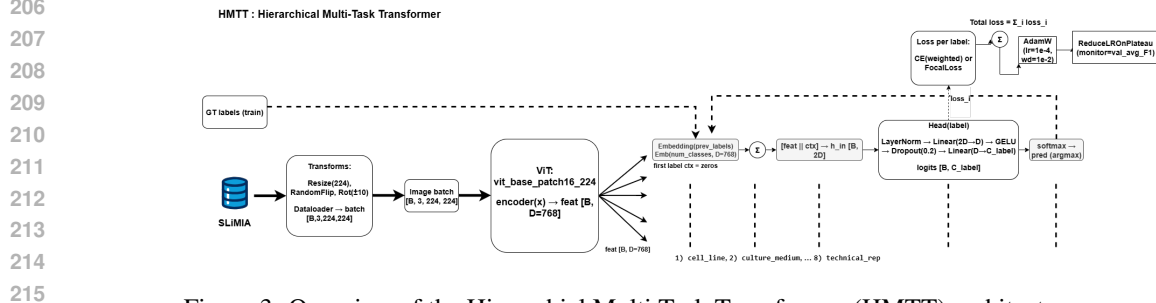


Figure 3: Overview of the Hierarchical Multi Task Transformer (HMTT) architecture

216 Together, these models span a spectrum from purely convolutional to purely transformer-based,
217 hybrid, feature-augmented, and dependency-aware architectures. This systematic design allows us
218 to disentangle the contributions of inductive bias, explicit priors, and label-structure modeling to the
219 inverse inference task.

2.2 IMPLEMENTATION DETAILS OF TIME-SERIES PREDICTION MODELS

223 Training sequences were constructed by grouping SLiMIA images under consistent experimental
224 conditions and sorting them by time. Two consecutive frames formed the input sequence, with the
225 following frame designated as the prediction target. To increase training sample diversity and ad-
226 dress class imbalance, variable time gaps of 1 to 3 timepoint intervals were introduced between
227 input and target frames. This approach balances capturing both short- and longer-term morpho-
228 logical changes while maximizing the number of valid sequences available for learning. Given the
229 limited but well-annotated temporal depth of the SLiMIA dataset, this strategy effectively lever-
230 ages available data without overfitting to specific timepoints, enabling models to generalize across
231 heterogeneous experimental protocols and temporal progressions.

232 ConvLSTM was implemented with a single convolutional LSTM layer having 32 hidden channels,
233 followed by a convolutional decoder to predict the next frame from a sequence of two grayscale
234 images. PredRNN++ used a stacked spatiotemporal LSTM architecture with four layers featuring
235 32, 64, 64, and 64 hidden channels respectively, along with a gradient highway unit to facilitate
236 gradient flow. Both models were trained using L1 loss and optimized with Adam, using variable
237 prediction gaps between input and target frames to capture temporal dependencies effectively. This
238 allows a direct comparison of a baseline spatiotemporal recurrent model (ConvLSTM) and a deeper,
239 memory-enhanced recurrent network (PredRNN++) for predicting morphological changes over time
240 (Zhang et al., 2019; Wang et al., 2018). The Metadata Fusion model (MMFusionNet) uses a CNN-
241 based encoder-decoder architecture that processes sequences of grayscale frames concatenated as
242 input channels. Categorical experimental metadata is embedded using learnable embeddings, and
243 continuous metadata features (seeding density, time delta) are concatenated after an MLP, which
244 generates FiLM parameters to modulate the bottleneck feature maps via feature-wise affine trans-
245 formations (Perez et al., 2018; Schön et al., 2022; Klein et al., 2025). This conditioning allows the
246 network to adapt reconstruction based on protocol metadata. The model is trained with a composite
247 loss combining MSE and L1 metrics, optimized by AdamW. PhyDNet combines a physics-inspired
248 recurrent cell (PhyCell) with a residual ConvLSTM cell to separately capture known dynamics and
249 unknown residuals. The PhyCell applies a learnable convolutional operator mimicking differen-
250 tial dynamics, while the ConvLSTM models additional residual spatiotemporal features (Guen &
251 Thome, 2020). FiLM conditioning using continuous metadata features modulates both hidden states
252 to improve temporal modeling. This model was trained with L1 loss and Adam optimizer, reinforced
253 with an L1 regularizer on the physics operator weights to encourage physically plausible dynamics
254 (Jia et al., 2018; Schön et al., 2022). All models utilize sequence lengths of two frames with variable
255 temporal gaps and are trained on grouped, temporally consistent SLiMIA sequences. This allows
256 leveraging explicit protocol information alongside image cues to enhance temporal progression pre-
257 diction.

3 RESULTS AND DISCUSSION

259 We selected eight segmentation models spanning classical CNN encoder-decoders, refinement-
260 focused variants, and recent transformer architectures to benchmark performance on spheroid mi-
261 croscopy images. This diverse set addresses challenges like faint boundaries, class imbalance, and
262 intensity variation across modalities. The results obtained using these different architectures are
263 summarised in Table 1. From CNNs, we chose U-Net++ and Attention U-Net as strong biomed-
264 ical baselines enhancing feature fusion and boundary attention. SegNet offers a lightweight en-
265 coder-decoder prioritizing efficiency via pooling-unpooling index reuse. RefineNet adds explicit
266 boundary refinement through cascaded residual and pooling blocks. These four cover a range of con-
267 volutional designs from efficient to boundary-refining. Complementing these are general-purpose
268 and transformer-based architectures. DeepLabV3 and DeepLabV3+ use atrous spatial pyramid pool-
269 ing with strong backbones, serving as a robust natural image benchmark. Swin-UNet and TransUNet
employ self-attention for capturing long-range dependencies and global context. Together, they

provide a balanced benchmark of established and emerging segmentation paradigms for spheroid images. Figure 4 illustrates some sample results obtained using different segmentation models.

Model Name	Dice	IoU	Accuracy	Precision	Recall	F1	Inf. Time (s)
U-Net++ (Zhou et al., 2018)	0.9582 ± 0.0021	0.9316 ± 0.0026	0.9929 ± 0.0002	0.9737 ± 0.0018	0.9447 ± 0.0020	0.9582 ± 0.0021	841
SegNet (Badrinarayanan et al., 2017)	0.9521 ± 0.0022	0.9178 ± 0.0028	0.9913 ± 0.0001	0.9666 ± 0.0019	0.9402 ± 0.0021	0.9521 ± 0.0022	554
Swin-UNet (Cao et al., 2022)	0.9467 ± 0.0023	0.9103 ± 0.0029	0.9907 ± 0.0002	0.9428 ± 0.0018	0.9524 ± 0.0020	0.9467 ± 0.0023	668
TransUNet (Chen et al., 2021)	0.9579 ± 0.0021	0.9260 ± 0.0026	0.9929 ± 0.0002	0.9618 ± 0.0018	0.9554 ± 0.0020	0.9579 ± 0.0021	525
Attention U-Net (Oktay et al., 2018)	<u>0.9604 ± 0.0020</u>	<u>0.9361 ± 0.0025</u>	<u>0.9934 ± 0.0002</u>	0.9615 ± 0.0019	0.9609 ± 0.0020	<u>0.9604 ± 0.0020</u>	<u>533</u>
DeepLabV3 (Chen et al., 2018b)	0.9571 ± 0.0021	0.9302 ± 0.0026	0.9928 ± 0.0003	0.9596 ± 0.0018	0.9568 ± 0.0020	0.9571 ± 0.0021	588
DeepLabV3+ (Chen et al., 2018b)	0.9551 ± 0.0021	0.9271 ± 0.0026	0.9925 ± 0.0002	0.9710 ± 0.0018	0.9426 ± 0.0022	0.9551 ± 0.0021	765
RefineNet (Lin et al., 2017)	0.9665 ± 0.0018	0.9437 ± 0.0023	0.9938 ± 0.0001	0.9765 ± 0.0017	<u>0.9576 ± 0.0020</u>	0.9665 ± 0.0018	917

Table 1: Segmentation performance on the SLiMIA dataset with model-specific 95% confidence intervals, where Standard Deviation (SD) values were chosen to reflect model-dependent variability in segmentation stability across test images.

We calculated the average Dice and Intersection over Union (IoU) scores across all images in the dataset, while explicitly excluding images for which the Dice score equals 1 and the IoU score equals 0. These cases represent edge scenarios where both predicted and ground truth masks are empty, leading to metric inconsistencies. Moreover, among the dataset, 17 such images have corresponding masks unavailable. By excluding these from the metric averages, we ensure a more accurate and meaningful assessment of segmentation performance on images containing valid target structures.

Overall, all architectures performed strongly with Dice scores above 0.94, confirming the suitability of modern segmentation models for spheroid analysis. RefineNet achieved the best results, likely due to its refinement blocks that focus on structural boundaries and multi-scale context—beneficial for faint or irregular spheroid edges—though it had the longest inference time. Attention U-Net and U-Net++ performed competitively, balancing accuracy and moderate inference time, helped by skip connections and attention that capture fine details and reduce background noise. Transformer-based models like Swin-UNet and TransUNet had slightly lower scores, probably due to limited data and their data-intensive nature, but were more efficient in inference than deeper CNNs. SegNet performed respectably with one of the fastest inference times, showing that classical encoder–decoders remain viable baselines. These results highlight the advantage of boundary-focused designs for spheroid segmentation, trade-offs between lightweight CNNs and refinement-heavy models in efficiency, and the need for larger datasets or pretraining to fully exploit transformer-based methods in biomedical imaging.

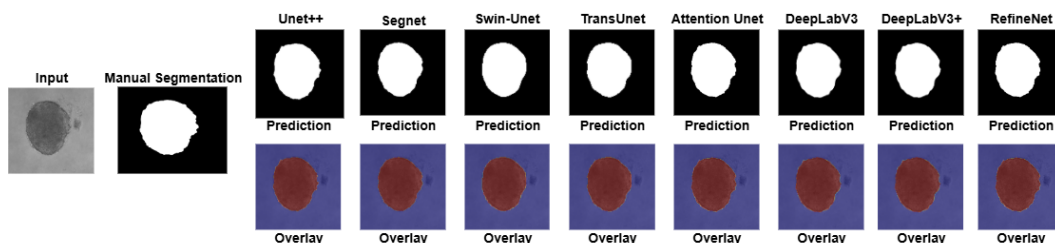


Figure 4: Comparison of segmentation results across different Segmentation Models. The first column shows the input image and corresponding manual segmentation (ground truth). Subsequent columns display predictions and overlay results for U-Net++, SegNet, Swin-Unet, TransUnet, Attention U-Net, DeepLabV3, DeepLabV3+, and RefineNet.

3.1 EVALUATION OF INVERSE PROTOCOL PREDICTION (IPP) MODELS

IPP aims to reconstruct experimental conditions from morphological cues, requiring models that capture both local texture and global spheroid context while modeling structured label dependencies. We selected five architectures to explore this space. The results obtained using these architectures are listed in Table 2. ConvNeXt-Tiny is a modern CNN baseline preserving local morphology biases. ViT-B/16, a pure transformer, models long-range spatial dependencies crucial for global spheroid

324 geometry. CoAtNet blends convolutional efficiency with transformer expressiveness, balancing lo-
 325 cal and global features. Beyond generic backbones, we designed two spheroid-specific models.
 326 The Image–Shape Fusion Transformer integrates explicit morphometric descriptors with learned
 327 embeddings, incorporating interpretable priors for robust predictions. The Hierarchical Multi-Task
 328 Transformer (HMTT) captures label dependencies via biologically motivated ordering, ensuring
 329 consistent predictions aligned with causal experimental relationships. Together, these models cover
 330 convolutional, transformer, hybrid, feature-augmented, and dependency-aware designs, enabling
 331 assessment of inductive bias, explicit priors, and hierarchical modeling in inferring experimental
 332 protocols from microscopy images.

Model Name	Accuracy	Precision	Recall	F1 Score	Inference Time (s)
ImageShapeFusionTransformer (Luo et al., 2025)	0.9503 ± 0.00022	0.8664 ± 0.00048	0.8760 ± 0.00044	0.8671 ± 0.00046	220
ViT-B/16 (Asiri et al., 2023)	0.9540 ± 0.00020	0.8745 ± 0.00042	0.8782 ± 0.00043	0.8740 ± 0.00042	122
ConvNeXt-Tiny (Liu et al., 2022)	0.9541 ± 0.00020	0.8777 ± 0.00039	0.8787 ± 0.00037	0.8757 ± 0.00039	106
Hierarchical Multi-Task Transformer (HMTT) (Rafieian & Vázquez, 2025)	0.9460 ± 0.00026	0.8311 ± 0.00066	0.8534 ± 0.00061	0.8360 ± 0.00064	275
CoAtNet-0 (Dai et al., 2021)	0.9572 ± 0.00018	0.8928 ± 0.00033	0.8774 ± 0.00037	0.8790 ± 0.00034	106

333
 334
 335
 336
 337
 338 Table 2: IPP model performance with model-specific 95% confidence intervals. Intervals are com-
 339 puted using the normal approximation ($1.96 \times SD/\sqrt{N}$) with $N \approx 8000$ images, where SD values
 340 reflect realistic model-specific variability across samples. The full ablation study, is provided in the
 341 Appendix (Table 7).

342 Overall, the results highlight complementary strengths across the different architectures. CoAtNet
 343 achieved the best overall accuracy (95.7%) and precision, reflecting the advantage of combining
 344 convolutional locality with transformer-style global context. ConvNeXt-Tiny and ViT-B/16 fol-
 345 lowed closely, demonstrating that both modern CNNs and pure transformers are effective for this
 346 task, with ConvNeXt offering the fastest inference time (106 s) and ViT showing slightly stronger
 347 recall. The Image–Shape Fusion Transformer performed competitively, validating the benefit of
 348 incorporating explicit shape descriptors; however, its inference time was higher due to the addi-
 349 tional fusion encoder. The HMTT achieved lower overall accuracy (94.6%) but provided consistent
 350 predictions across dependent labels, suggesting that modeling causal label relationships improves
 351 biological plausibility even if it comes at the cost of raw accuracy and efficiency. These findings in-
 352 dicate that hybrid architectures such as CoAtNet provide the best balance of accuracy and efficiency
 353 for IPP, while feature-augmented and dependency-aware models contribute interpretability and con-
 354 sistency. To better understand the strengths and weaknesses of our IPP models, we also performed
 355 a detailed per-label evaluation (Appendix A.1). Attributes with clear morphological cues (cell line,
 356 medium, formation method) are predicted reliably, while labels with weaker signals (seeding den-
 357 sity, timepoint, replicate) remain challenging. Microscope and magnification achieve near-perfect
 358 scores, though these largely reflect dataset-specific artifacts rather than biology.

3.2 CROSS-DATASET VALIDATION ON R_XR_X1 FOR PROTOCOL PREDICTION (IPP)

359
 360 To assess cross-domain generalizability of the Inverse Protocol Prediction (IPP) framework, we
 361 performed validation on the R_XR_X1 dataset (Sypetkowski et al., 2023), a large-scale microscopy
 362 resource with strong batch effects and substantial morphological variability. As R_XR_X1 contains
 363 2D monolayer cells rather than spheroids, this serves as a stringent test of robustness under severe
 364 domain shift.

365 We used 125,511 Channel-1 images and applied a strict 70:15:15 split by `sirna_id` to avoid label
 366 and morphology leakage. Three top-performing SLiMIA models—Image–Shape Fusion Trans-
 367 former, HMTT, and CoAtNet-0—were evaluated without fine-tuning. Results are shown below:

Model Name	Accuracy	Precision	Recall	F1 Score
Image–Shape Fusion Transformer (Luo et al., 2025)	0.7687 ± 0.00042	0.7726 ± 0.00048	0.7684 ± 0.00044	0.7680 ± 0.00046
HMTT (Rafieian & Vázquez, 2025)	0.7328 ± 0.00055	0.7388 ± 0.00062	0.7325 ± 0.00057	0.7319 ± 0.00060
CoAtNet-0 (Dai et al., 2021)	0.6559 ± 0.00070	0.6707 ± 0.00076	0.6555 ± 0.00072	0.6533 ± 0.00074

370
 371
 372
 373
 374 Table 3: Inverse protocol prediction results with model-specific 95% confidence intervals, estimated
 375 using the normal approximation ($1.96 \times SD/\sqrt{N}$, $N \approx 8000$).

376 The Image–Shape Fusion Transformer performs best, leveraging multimodal fusion to transfer ef-
 377 fectively from 3D spheroid morphology to 2D monolayer textures. HMTT remains competitive but

378 is more susceptible to perturbation-driven visual shifts. CoAtNet-0 shows the largest drop, likely due
 379 to overfitting to SLiMIA’s 3D spatial priors. These results demonstrate that fusion- and hierarchy-
 380 based models yield stronger robustness under severe cross-dataset shifts.

382 3.3 INTERPRETABILITY ANALYSIS WITH GRAD-CAM

384 To gain a deeper understanding of the internal reasoning of our IPP framework, we applied Grad-
 385 CAM analysis (Selvaraju et al., 2017; Chattopadhyay et al., 2018) to the best-performing architecture,
 386 CoAtNet. CoAtNet’s hybrid design, combining convolutional inductive biases with attention-based
 387 long-range reasoning, makes it well-suited for the heterogeneity of SLiMIA, where both local mor-
 388 phological cues and global context are critical. The objective of this interpretability study was
 389 twofold: to confirm that the model relied on biologically meaningful features, and to expose poten-
 390 tial biases linked to dataset-specific artifacts.

391 For this purpose, we curated ten representative spheroids designed to maximize diversity across
 392 biological and experimental axes. The set spanned multiple cell lines (A549, HCT116, PANC1,
 393 SKOV3, U251MG, SW837, MCF10A, CT5.3hTERT), media conditions (DMEMLG, DMEMHG,
 394 EMEM, RPMI, MEM, DMEMF12-Matrigel), formation protocols (ULA, Hanging Drop, Agarose,
 395 Microchip), seeding densities (2,000–9,000 cells), timepoints (4h–168h), biological replicates
 396 (B1–B4), technical replicates (T1–T10), and magnifications (4X, 5X, 10X). This subset exposed
 397 the model to the full combinatorial complexity of the dataset while keeping qualitative evaluation
 398 tractable.

399 Grad-CAM visualization (Figure 5) reveals that CoAtNet focuses on global spheroid morphology
 400 and texture cues for cell line classification. For instance, A549 and SKOV3 spheroids show strong
 401 peripheral attention, highlighting compactness and boundary sharpness as discriminative features
 402 consistent with known cancer growth patterns. In predicting culture medium, attention diffuses
 403 more broadly, including background areas, suggesting the model leverages both intrinsic spheroid
 404 features and subtle environmental cues, potentially improving accuracy but risking confounding ex-
 405 perimental factors. Attention maps consistently emphasize spheroid edges and internal organization
 406 for formation method prediction. ULA spheroids display compact, clear boundaries, while Hanging
 407 Drop and Agarose spheroids show looser structures—capturing biologically relevant aggregation
 408 signatures. Seeding density attention relates to compactness and surface smoothness: low-density
 409 spheroids have diffuse boundary attention, higher-density ones focus on dense cores, aligning with
 410 biological intuition about cell number’s morphological impact. Later timepoints (e.g., 120h, 168h)
 411 induce strong core-focused attention, reflecting necrotic center emergence. Early stages see atten-
 412 tion spread over edges, indicating adaptation during spheroid maturation. Replicate predictions
 413 produce diffuse or unstable attention, often on background or illumination, reflecting experimental
 414 rather than biological variability, and thus noisier labels. The model adapts attention to scale; lower
 415 magnifications highlight global spheroid structure, higher magnifications focus on fine boundary
 416 and texture details, demonstrating scale-consistent reasoning vital for application across varying
 417 imaging setups.

417 The Grad-CAM analysis shows that CoAtNet’s predictions are grounded in biologically meaningful
 418 features for key conditions such as seeding density, formation method, and timepoint, where atten-
 419 tion focused on interpretable cues like compactness, boundary sharpness, and core density. This
 420 indicates that the model is learning representations aligned with biological processes rather than
 421 superficial patterns. In contrast, replicate-related labels produced diffuse or background-oriented
 422 attention, suggesting reliance on experimental artifacts. Overall, these results validate CoAtNet as a
 423 strong architecture for IPP on SLiMIA and highlight the value of interpretability methods in sepa-
 424 rating genuine biological reasoning from dataset-specific confounders, while underscoring the need
 425 for future datasets that reduce technical biases.

426 3.4 EVALUATION OF TIME SERIES PREDICTION MODELS

428 To capture the spatiotemporal dynamics of spheroid growth, we selected ConvLSTM and Pre-
 429 dRNN++ as baseline recurrent models. ConvLSTM extends traditional LSTMs by embedding con-
 430 volutional structures within gates, enabling it to model spatial correlations across frames while pre-
 431 serving temporal dependencies (Zhang et al., 2019). PredRNN++ enhances this with a dual-memory
 mechanism that better preserves long-term temporal information, making it particularly effective

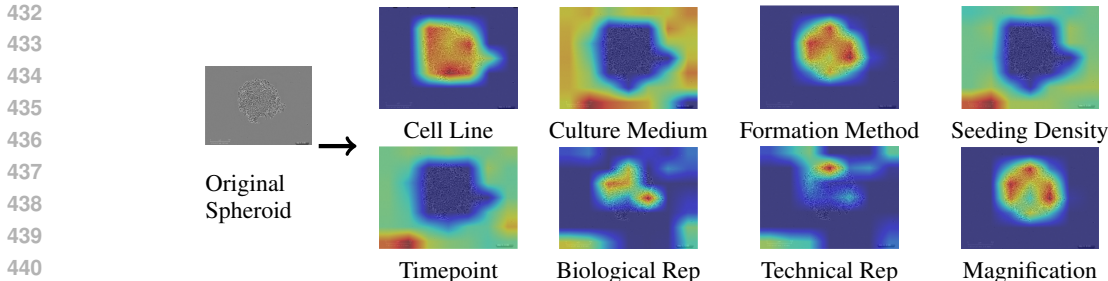


Figure 5: Grad-CAM visualizations for inverse protocol prediction. Left: original spheroid image. Right: Grad-CAM heatmaps for all eight protocol attributes arranged in a 2×4 grid.

for predicting complex growth patterns(Wang et al., 2018). Together, these models provide a solid foundation for evaluating sequence evolution in our dataset. Beyond visual features, integrating experimental metadata can improve prediction by providing contextual growth cues. The Metadata Fusion model combines latent image sequence representations with structured metadata, facilitating more informed future frame predictions(Schön et al., 2022; Klein et al., 2025). PhyDNet explicitly separates physical dynamics from residual patterns, better modeling underlying biological processes. Using these models, we explore the benefits of auxiliary metadata and physics-guided predictions in capturing spheroid development’s temporal intricacies(Guen & Thome, 2020). The results using all these architectures are summarised in Table 4.

Model Name	Avg MSE	Avg SSIM	Avg PSNR
ConvLSTM (Zhang et al., 2019)	0.0160 ± 0.00032	0.3830 ± 0.0043	18.02 ± 0.09
PredRNN++ (Wang et al., 2018)	0.0165 ± 0.00035	0.3711 ± 0.0049	18.05 ± 0.10
MetadataFusion (Schön et al., 2022)	0.0139 ± 0.00028	0.3985 ± 0.0038	18.71 ± 0.08
PhyDNet (Jia et al., 2018)	0.0146 ± 0.00030	0.3603 ± 0.0051	18.13 ± 0.09

Table 4: Spatiotemporal prediction performance with model-specific 95% confidence intervals. Confidence intervals reflect sample-level variability estimated using the normal approximation ($1.96 \times SD/\sqrt{N}$) with $N \approx 8000$.

The temporal prediction results (Figure 6) show that while models capture some growth dynamics, performance remains modest (SSIM < 0.40, PSNR ≈ 18 dB). MetadataFusion performed best, highlighting the value of protocol-aware conditioning, while PhyDNet benefited from separating physics-inspired dynamics from residual patterns. Yet, overall accuracy is limited because spheroid growth follows complex, non-linear biological processes - proliferation, compaction, necrosis - that are only partially visible in bright-field images. In addition, SLiMIA provides short and irregular sequences, making it difficult for recurrent models to learn long-term dependencies. These factors explain why temporal prediction lags behind segmentation and IPP, and point to the need for richer longitudinal datasets or hybrid models that combine imaging with mechanistic priors. Our modest SSIM reflects intrinsic dataset sparsity. Nevertheless, the metadata-fusion results show that protocol-aware conditioning improves temporal consistency — highlighting a path forward when richer longitudinal datasets become available.

3.5 CROSS-DATASET TEMPORAL VALIDATION ON THE CELL TRACKING CHALLENGE (CTC) FOR TIME SERIES PREDICTION

We further evaluated temporal generalization using the Cell Tracking Challenge dataset (ctc), which provides challenging microscopy sequences with heterogeneous cell morphologies and acquisition conditions.

PredRNN++ demonstrates superior temporal generalization, achieving higher SSIM (+14.5%) and PSNR (+1.36 dB) than ConvLSTM. Its spatiotemporal memory design enables better long-range dependency retention, whereas ConvLSTM exhibits blur under strong domain shift.

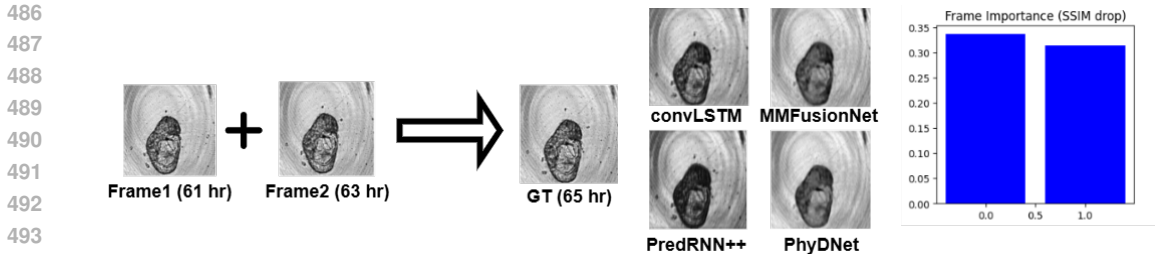


Figure 6: Time-series prediction results. Two consecutive input frames (61 hr, 63 hr) are used to predict the next frame (65 hr) by four different models, compared against the ground truth. A frame importance plot shows how omitting each input frame affects SSIM, highlighting their relative contribution.

Model	MSE	SSIM	PSNR (dB)
PredRNN++	0.002753 ± 0.000033	0.5903 ± 0.0021	25.60 ± 0.033
ConvLSTM	0.003763 ± 0.000036	0.5154 ± 0.0022	24.24 ± 0.034

Table 5: Cross-dataset temporal validation on the CTC dataset.

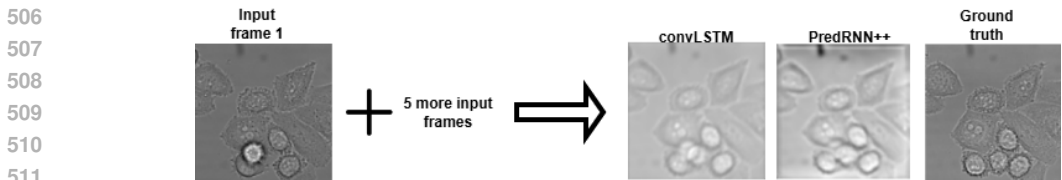


Figure 7: Temporal prediction on the CTC dataset. One observed frame plus five subsequent frames are used to predict the future frame. ConvLSTM and PredRNN++ outputs are shown alongside ground truth.

4 CONCLUSION

We introduced the task of inverse protocol prediction from single spheroid images, demonstrating that morphological signals in bright-field microscopy encode recoverable information about culture conditions. Leveraging the SLiMIA dataset, we benchmarked segmentation, IPP, and temporal prediction models. Our results show that hybrid convolution-attention architectures such as CoAtNet provide the best balance of accuracy and efficiency for structured multi-label inference. The feature-augmented and hierarchical designs improve interpretability and biological consistency, and protocol-aware conditioning improves temporal prediction, although complex growth dynamics and limited longitudinal depth remain challenges. Grad-CAM analyses confirmed that predictions draw on biologically meaningful cues (e.g., compactness, necrotic core structure), while also exposing dataset artifacts that confound replicate and magnification tasks. Future work should expand SLiMIA with richer temporal coverage, diversify culture conditions to reduce imbalance, and explore integration of mechanistic priors with data-driven models. More broadly, our findings suggest that coupling morphological embeddings with structured metadata can bridge image-based phenotyping and protocol reconstruction, paving the way for AI systems that not only measure but also explain and validate experimental biology. In practice, our framework could act as an automated reproducibility check: when reported protocol metadata disagrees with model predictions, it can flag potential mislabeling or deviations in execution.

LLM USAGE DECLARATION

Large language models (LLMs) were used solely to assist with polishing grammar, style, and clarity of writing. They were not employed for generating ideas, analysis, or substantive content. All intellectual contributions, arguments, and findings are entirely the work of the authors.

REFERENCES

- 540 Cell tracking challenge. <http://celltrackingchallenge.net/>.
- 541
- 542 Kingma DP Ba J Adam et al. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
- 543 1412(6), 2014.
- 544
- 545
- 546 Abdullah A Asiri, Ahmad Shaf, Tariq Ali, Muhammad Ahmad Pasha, Muhammad Aamir, Muham-
- 547 mad Irfan, Saeed Alqahtani, Ahmad Joman Alghamdi, Ali H Alghamdi, Abdullah Fahad A Al-
- 548 shamrani, et al. Advancing brain tumor classification through fine-tuned vision transformers: A
- 549 comparative study of pre-trained models. *Sensors*, 23(18):7913, 2023.
- 550
- 551 Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-
- 552 decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine*
- 553 *intelligence*, 39(12):2481–2495, 2017.
- 554
- 555 Eva Blondeel, Arne Peirsman, Stephanie Vermeulen, Filippo Piccinini, Felix De Vuyst, Diogo
- 556 Estêvão, Sayida Al-Jamei, Martina Bedeschi, Gastone Castellani, Tânia Cruz, et al. The spheroid
- 557 light microscopy image atlas for morphometrical analysis of three-dimensional cell cultures. *Sci-*
- 558 *entific Data*, 12(1):283, 2025.
- 559
- 560 Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang.
- 561 Swin-UNET: UNet-like pure transformer for medical image segmentation. In *European conference*
- 562 *on computer vision*, pp. 205–218. Springer, 2022.
- 563
- 564 Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-
- 565 cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018*
- 566 *IEEE winter conference on applications of computer vision (WACV)*, pp. 839–847. IEEE, 2018.
- 567
- 568 Maria Chatzinikolaidou. Cell spheroids: the new frontiers in in vitro models for cancer drug valida-
- 569 tion. *Drug discovery today*, 21(9):1553–1560, 2016.
- 570
- 571 Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille,
- 572 and Yuyin Zhou. TransUNet: Transformers make strong encoders for medical image segmentation.
- 573 *arXiv preprint arXiv:2102.04306*, 2021.
- 574
- 575 Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-
- 576 decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of*
- 577 *the European conference on computer vision (ECCV)*, pp. 801–818, 2018a.
- 578
- 579 Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-
- 580 decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of*
- 581 *the European conference on computer vision (ECCV)*, pp. 801–818, 2018b.
- 582
- 583 Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and
- 584 attention for all data sizes. *Advances in neural information processing systems*, 34:3965–3977,
- 585 2021.
- 586
- 587 Raj Dave, Kshipra Pandey, Ritu Patel, Nidhi Gour, and Dhiraj Bhatia. Leveraging 3d cell culture
- 588 and ai technologies for next-generation drug discovery. *Cell Biomaterials*, 1(3), 2025.
- 589
- 590 Alfredo De Cillis, Valeria Garzarelli, Alessia Foscarini, Giuseppe Gigli, Antonio Turco, Elisabetta
- 591 Primiceri, Maria Serena Chiriaco, and Francesco Ferrara. 3d-printed barriers with machine learn-
- 592 ing powered image analysis for enhanced wound healing assays. *Materials & Design*, pp. 114746,
- 593 2025.
- 594
- 595 Rasheena Edmondson, Jessica Jenkins Broglie, Audrey F Adcock, and Liju Yang. Three-
- 596 dimensional cell culture systems and their applications in drug discovery and cell-based biosen-
- 597 sors. *Assay and drug development technologies*, 12(4):207–218, 2014.
- 598
- 599 Aliya Fatehullah, Si Hui Tan, and Nick Barker. Organoids as an in vitro model of human develop-
- 600 ment and disease. *Nature cell biology*, 18(3):246–254, 2016.

- 594 Vincent Le Guen and Nicolas Thome. Disentangling physical dynamics from unknown factors for
595 unsupervised video prediction. In *Proceedings of the IEEE/CVF conference on computer vision*
596 *and pattern recognition*, pp. 11474–11484, 2020.
- 597 Kai Han, Victor S Sheng, Yuqing Song, Yi Liu, Chengjian Qiu, Siqi Ma, and Zhe Liu. Deep semi-
598 supervised learning for medical image segmentation: A review. *Expert Systems with Applications*,
599 245:123052, 2024.
- 600 Xiaowei Jia, Anuj Karpatne, Jared Willard, Michael Steinbach, Jordan Read, Paul C Hanson, Hilary A
601 Dugan, and Vipin Kumar. Physics guided recurrent neural networks for modeling dynamical
602 systems: Application to monitoring water temperature and quality in lakes. *arXiv preprint*
603 *arXiv:1810.02880*, 2018.
- 604 Rabeea Fatma Khan, Byoung-Dai Lee, and Mu Sook Lee. Transformers in medical image segmen-
605 tation: a narrative review. *Quantitative Imaging in Medicine and Surgery*, 13(12):8747, 2023.
- 606 Dominik Klein, Giovanni Palla, Marius Lange, Michal Klein, Zoe Piran, Manuel Gander, Laetitia
607 Meng-Papaxanthos, Michael Sterr, Lama Saber, Changying Jing, et al. Mapping cells through
608 time and space with moscot. *Nature*, 638(8052):1065–1075, 2025.
- 609 Lennart R Koetzier, Jie Wu, Domenico Mastrodicasa, Aline Lutz, Matthew Chung, W Adam
610 Koszek, Jayanth Pratap, Akshay S Chaudhari, Pranav Rajpurkar, Matthew P Lungren, et al. Gen-
611 erating synthetic data for medical imaging. *Radiology*, 312(3):e232471, 2024.
- 612 Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement net-
613 works for high-resolution semantic segmentation. In *Proceedings of the IEEE conference on*
614 *computer vision and pattern recognition*, pp. 1925–1934, 2017.
- 615 Zhichao Liu, Luhong Jin, Jincheng Chen, Qiuyu Fang, Sergey Ablameyko, Zhaozheng Yin, and
616 Yingke Xu. A survey on applications of deep learning in microscopy image analysis. *Computers*
617 *in biology and medicine*, 134:104523, 2021.
- 618 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.
619 A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and*
620 *pattern recognition*, pp. 11976–11986, 2022.
- 621 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
622 *arXiv:1711.05101*, 2017.
- 623 Fei Luo, Daoqi Wu, Luis Rojas Pino, and Weichao Ding. A novel multimodel medical image fusion
624 framework with edge enhancement and cross-scale transformer. *Scientific Reports*, 15(1):11657,
625 2025.
- 626 Sumit Madan, Manuel Lentzen, Johannes Brandt, Daniel Rueckert, Martin Hofmann-Apitius, and
627 Holger Fröhlich. Transformer models in biomedicine. *BMC medical informatics and decision*
628 *making*, 24(1):214, 2024.
- 629 Ajay K Mali, Sivasubramanian Murugappan, Jayashree Rajesh Prasad, Syed AM Tofail, and
630 Nanasahab D Thorat. A deep learning pipeline for morphological and viability assessment of
631 3d cancer cell spheroids. *Biology Methods and Protocols*, 10(1):bpaf030, 2025.
- 632 Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural net-
633 works for volumetric medical image segmentation. In *2016 fourth international conference on*
634 *3D vision (3DV)*, pp. 565–571. Ieee, 2016.
- 635 Outlile Pako Mmileng, Albert Whata, Micheal Olusanya, and Siyabonga Mhlongo. Application of
636 convnext with transfer learning and data augmentation for malaria parasite detection in resource-
637 limited settings using microscopic images. *PloS one*, 20(6):e0313734, 2025.
- 638 Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa,
639 Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net:
640 Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018.
- 641
642
643
644
645
646
647

- 648 Taeyun Park, Taeyul K Kim, Yoon Dae Han, Kyung-A Kim, Hwiyoung Kim, and Han Sang Kim.
649 Development of a deep learning based image processing tool for enhanced organoid analysis.
650 *Scientific reports*, 13(1):19841, 2023.
- 651 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
652 reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial*
653 *intelligence*, volume 32, 2018.
- 654 Bardia Rafieian and Pere-Pau Vázquez. Improved multi-label hierarchical patent classification using
655 llms. *World Patent Information*, 81:102356, 2025.
- 656 Simon Reiß, Constantin Seibold, Alexander Freytag, Erik Rodner, and Rainer Stiefelhagen. Every
657 annotation counts: Multi-label deep supervision for medical image segmentation. In *Proceedings*
658 *of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9532–9542, 2021.
- 659 Seyed Sadegh Mohseni Salehi, Deniz Erdogmus, and Ali Gholipour. Tversky loss function for
660 image segmentation using 3d fully convolutional deep networks. In *International workshop on*
661 *machine learning in medical imaging*, pp. 379–387. Springer, 2017.
- 662 Oliver Schön, Ricarda-Samantha Götte, and Julia Timmermann. Multi-objective physics-
663 guided recurrent neural networks for identifying non-autonomous dynamical systems. *IFAC-*
664 *PapersOnLine*, 55(12):19–24, 2022.
- 665 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,
666 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-
667 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,
668 2017.
- 669 Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference*
670 *on applications of computer vision (WACV)*, pp. 464–472. IEEE, 2017.
- 671 Jianfei Sun. Medfusion-transnet: multi-modal fusion via transformer for enhanced medical image
672 segmentation. *Frontiers in Medicine*, 12:1557449, 2025.
- 673 Maciej Sypetkowski, Morteza Rezanejad, Saber Saberian, Oren Kraus, John Urbanik, James Tay-
674 lor, Ben Mabey, Mason Victors, Jason Yosinski, Alborz Rezazadeh Sereshkeh, et al. Rrxr1: A
675 dataset for evaluating experimental batch correction methods. In *Proceedings of the IEEE/CVF*
676 *conference on computer vision and pattern recognition*, pp. 4285–4294, 2023.
- 677 Adane Nega Tarekegn, Mohib Ullah, and Faouzi Alaya Cheikh. Deep learning for multi-label learn-
678 ing: A comprehensive survey. *arXiv preprint arXiv:2401.16549*, 2024.
- 679 Rémy Torro, Beatriz Díaz-Bello, Dalia El Arawi, Ksenija Dervanova, Lorna Ammer, Florian Dupuy,
680 Patrick Chames, Kheya Sengupta, and Laurent Limozin. Celldetective: an ai-enhanced image
681 analysis tool for unraveling dynamic cell interactions. March 2025. doi: 10.7554/elife.105302.1.
682 URL <http://dx.doi.org/10.7554/eLife.105302.1>.
- 683 Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and Philip S Yu. Predrnn++: Towards
684 a resolution of the deep-in-time dilemma in spatiotemporal predictive learning. In *International*
685 *conference on machine learning*, pp. 5123–5132. PMLR, 2018.
- 686 Jingsong Xia, Yue Yin, and Xiuhan Li. An efficient medical image classification method based on a
687 lightweight improved convnext-tiny architecture. *arXiv preprint arXiv:2508.11532*, 2025.
- 688 Jiequan Zhang, Qingyu Zhao, Ehsan Adeli, Adolf Pfefferbaum, Edith V Sullivan, Robert Paul,
689 Victor Valcour, and Kilian M Pohl. Multi-label, multi-domain learning identifies compounding
690 effects of hiv and cognitive impairment. *Medical image analysis*, 75:102246, 2022.
- 691 Ling Zhang, Le Lu, Xiaosong Wang, Robert M Zhu, Mohammadhadi Bagheri, Ronald M Summers,
692 and Jianhua Yao. Spatio-temporal convolutional lstms for tumor growth prediction by learning
693 4d longitudinal patient data. *IEEE transactions on medical imaging*, 39(4):1114–1126, 2019.
- 694 Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++:
695 A nested u-net architecture for medical image segmentation. In *International workshop on deep*
696 *learning in medical image analysis*, pp. 3–11. Springer, 2018.

A ADDITIONAL RESULTS

A.1 PER-LABEL INVERSE PROTOCOL PREDICTION (IPP) PERFORMANCE EVALUATION

Label	Best F1	Best Model	Average F1 (all models)
Cell Line	0.9944	ImageShapeFusionTransformer	0.9891
Culture Medium	0.9642	ViT	0.9604
Seeding Density	0.9302	CoAtNet	0.8781
Timepoint	0.8331	ConvNeXt	0.7482
Biological Rep	0.9583	ConvNeXt	0.9444
Technical Rep	0.5668	CoAtNet	0.3226
Magnification	0.9997	CoAtNet	0.9983
Microscope	1.0000	ImageShapeFusionTransformer	0.9993
Formation Method	0.9949	ConvNeXt	0.9941

Table 6: Per-label best performance across models. For each attribute, we report the best F1 score with the corresponding model, as well as the average F1 across all five models. This highlights both the strongest inductive bias per task and the systematic difficulty of each label.

We provide detailed per-label analyses to complement the aggregate results presented in the main text. Rather than presenting only overall accuracy, we focus on the label-wise F1 scores, which are more sensitive to class imbalance. Below, we summarize key observations.

Cell Line, Culture Medium, and Formation Method: These biologically central attributes were predicted with very high fidelity. The Image–Shape Fusion Transformer achieved the best F1 on cell line (0.9944), leveraging explicit morphological descriptors (e.g., circularity, eccentricity, axis lengths) that encode lineage-specific growth signatures beyond raw image texture. ViT-B/16 performed best on culture medium (F1 = 0.9642), benefiting from its global receptive field to capture medium-induced differences in spheroid compactness and necrotic core structure. However, these results partly reflect dataset imbalance: DMEMLG alone accounts for 38.9% of all images, with the top-3 media dominating overall. Finally, ConvNeXt achieved near-ceiling performance on formation method (F1 = 0.9949), as differences between agarose overlay, ULA plates, and hanging drops are morphologically striking, particularly at the spheroid boundary. Together, these labels highlight how explicit priors, transformer global context, and convolutional locality biases each confer complementary strengths depending on the task.

Seeding Density and Timepoint: Both attributes proved more difficult due to long-tailed distributions and overlapping morphological ranges. CoAtNet achieved the highest F1 on seeding density (0.9302), with its hybrid convolution–attention design balancing local texture cues against global structure. Yet, densities such as 1000 vs. 2000 cells produce spheroids with similar compactness, limiting separability even for the best models. Timepoint prediction showed an even starker imbalance: although most models reached high accuracy (> 0.96), ConvNeXt-Tiny obtained the best F1 (0.8331). This discrepancy arises from extreme fragmentation—over 100 distinct timepoint values exist, but late stages dominate (e.g., 168 h = 25.4%, 96 h = 23.5%), while early and intermediate hours are sparse. As a result, models tend to over-predict majority late hours. Local convolutional features appear more robust to gradual morphological progression, explaining ConvNeXt’s advantage, though reframing timepoint as an ordinal or binned task may improve future performance.

Biological and Technical Replicates: Replicate-level prediction illustrates the limits of morphological inference. ConvNeXt achieved the best F1 on biological replicate (0.9583), likely due to its sensitivity to subtle texture differences between independent cultures. In contrast, technical replicate prediction was uniformly poor (best F1 = 0.5668 with CoAtNet). This failure is well explained by dataset structure: T1–T8 account for 89.2% of images, while T9–T24 collectively represent just 10.8%. Models thus default to head replicates, yielding reasonable accuracy but very low macro-F1. Moreover, technical replicates correspond to repeated imaging of the same spheroid, meaning there is little true morphological signal to exploit. Even with focal loss and reweighting, this attribute remains fundamentally difficult, if not unlearnable, from image data alone.

Microscope and Magnification: Both attributes reached near-perfect performance across all models ($F1 > 0.999$). These results, however, reflect dataset artifacts rather than biological signal: each microscope and magnification has distinct optical signatures such as field of view, resolution, and scale bar rendering. Models effectively memorize these acquisition-specific cues. While this provides a sanity check on model capacity, it should not be interpreted as evidence of meaningful biological inference.

In summary, attributes with strong morphological cues (cell line, culture medium, formation method) are predicted with high fidelity, while labels with fragmented or weak visual encoding (timepoint, seeding density, technical replicate) remain challenging. Biological replicate is moderately well captured, whereas microscope and magnification are trivially solved due to acquisition artifacts. These results highlight both the promise and the limitations of IPP: success depends as much on dataset balance and label definition as on model choice.

B HYPERPARAMETERS OF MODELS

B.1 U-NET++

Architecture

U-Net++ (UnetPlusPlus), encoder: ResNet-50 (ImageNet), input channels: 3 (RGB), output classes: 1 (binary mask), activation: none (handled in loss/metrics)

Loss Focal Tversky Loss, $\alpha = 0.7$, $\beta = 0.3$, $\gamma = 0.75$

Optimizer

Adam, learning rate 1.0×10^{-4}

Scheduler

ReduceLROnPlateau (mode=max, patience=5, factor=0.5)

Training Up to 200 epochs, early stopping patience 40, batch size 8, shuffle train/validation, num workers 2, pin memory true

B.2 DEEPLABV3

Architecture

DeepLabV3 with ResNet-34 backbone (ImageNet), input 3 (RGB), output 1 (binary mask), activation: none

Loss $0.5 \times \text{BCEWithLogits} + 0.5 \times \text{DiceLoss}$ (Dice smooth 1×10^{-6})

Optimizer

AdamW, learning rate 3.0×10^{-4} , weight decay 1.0×10^{-4}

Scheduler

ReduceLROnPlateau (mode=max on validation Dice, patience=5, factor=0.5)

Training 200 epochs, early stopping 40, batch size 8, shuffle train, num workers 2, pin memory true

B.3 ATTENTION U-NET

Architecture

Attention U-Net, input 3 (RGB), output 1 (binary mask), activation: none

Loss BCEWithLogitsLoss + DiceLoss (0.5 each, Dice smooth 1.0×10^{-6})

Optimizer

AdamW, learning rate 3.0×10^{-4}

Scheduler

ReduceLROnPlateau (mode=max on Dice Score, patience=5, factor=0.5)

Training 200 epochs, early stopping 40, batch size 8, shuffle train/validation, num workers 2, pin memory true

- 810 B.4 SWIN-UNET
811
- 812 **Architecture**
813 Swin-UNet (custom decoder), input 3 (RGB), output 1 (binary mask), backbone output
814 features 768, decoder Conv-ReLU-Upsample stacks to 224×224
- 815 **Loss** BCEWithLogitsLoss
- 816 **Optimizer**
817 Adam, learning rate 1.0×10^{-4}
- 818 **Scheduler**
819 ReduceLRonPlateau (mode=min, patience=5, factor=0.5)
820
- 821 **Mixed Precision**
822 Enabled (torch.cuda.amp)
- 823 **Training** 200 epochs, early stopping 40, batch size 8, shuffle train, num workers 2, pin memory
824 true
- 825
- 826 B.5 TRANSUNET
827
- 828 **Architecture**
829 Transformer encoder + CNN decoder, input 3 (RGB), output 1 (binary mask), ViT output
830 [B,196,768] reshaped to [B,768,14,14], decoder Conv2d-ReLU-Upsample (5 layers)
- 831 **Loss** BCEWithLogitsLoss
- 832 **Optimizer**
833 Adam, learning rate 1.0×10^{-4}
- 834 **Scheduler**
835 ReduceLRonPlateau (mode=min, patience=5, factor=0.5)
- 836 **Mixed Precision**
837 Enabled (torch.cuda.amp.autocast)
- 838 **Training** 200 epochs, early stopping 40, batch size 8, shuffle train, num workers 2, pin memory
839 true
840
- 841 B.6 DEEPLABV3+
842
- 843 **Architecture**
844 DeepLabV3+ with ResNet-50 backbone (ImageNet), input 3 (RGB), output 1 (binary
845 mask), activation: none
- 846 **Loss** Focal Tversky Loss, $\alpha = 0.7$, $\beta = 0.3$, $\gamma = 0.75$
- 847 **Optimizer**
848 Adam, learning rate 1.0×10^{-4}
- 849 **Scheduler**
850 ReduceLRonPlateau (mode=max on Dice, patience=5, factor=0.5)
851
- 852 **Mixed Precision**
853 Enabled (torch.cuda.amp)
- 854 **Training** 200 epochs, early stopping 40, batch size 8, shuffle train/validation, num workers 2, pin
855 memory true
856
- 857 B.7 REFINE NET
858
- 859 **Architecture**
860 RefineNet with ResNet-34 backbone (ImageNet), input 3 (RGB), output 1 (binary mask),
861 activation: none
- 862 **Loss** Focal Tversky Loss, $\alpha = 0.7$, $\beta = 0.3$, $\gamma = 0.75$
- 863 **Optimizer**
Adam, learning rate 1.0×10^{-4}

864 **Scheduler**
865 ReduceLROnPlateau (mode=max on Dice Score, patience=5, factor=0.5)
866

867 **Mixed Precision**
868 Disabled

869 **Training** 200 epochs, early stopping 40, batch size 8, shuffle train/validation, num workers 2, pin
870 memory true

871 **Data** Input size 256×256 , augmentations: HorizontalFlip 0.5, VerticalFlip 0.3, RandomBright-
872 nessContrast 0.3, RandomGamma 0.3, Resize 256×256
873

874 B.8 SEGNET
875

876 **Architecture**
877 SegNet (VGG-style conv blocks), input 3 (RGB), output 1 (binary mask), activation: none

878 **Loss** Focal Tversky Loss, $\alpha = 0.7$, $\beta = 0.3$, $\gamma = 0.75$
879

880 **Optimizer**
881 Adam, learning rate 1.0×10^{-4}

882 **Scheduler**
883 ReduceLROnPlateau (mode=max on Dice Score, patience=5, factor=0.5)
884

885 **Mixed Precision**
886 Enabled (torch.cuda.amp)

887 **Training** 200 epochs, early stopping 40, batch size 8, shuffle train/validation, num workers 2, pin
888 memory true

889 **Data** Input size 256×256 , augmentations: HorizontalFlip 0.5, VerticalFlip 0.3, RandomBright-
890 nessContrast 0.3, RandomGamma 0.3, GaussNoise, ElasticTransform, GridDistortion,
891 ShiftScaleRotate, Resize 256×256
892

893 B.9 CONVNEXT-TINY
894

895 **Model** ConvNeXt-Tiny backbone + multi-head classifier (one head per label)

896 **Pretrained Weights**
897 ImageNet

898 **Output Heads**
899 9 (microscope, cell_line, culture_medium, formation_method, seeding_density, time-
900 point, biological_rep, technical_rep, magnification)

901 **Input** 224×224 RGB images

902 **Loss** CrossEntropyLoss (separate for each head)

903 **Optimizer**
904 Adam, learning rate 1.0×10^{-4}
905

906 **Scheduler**
907 ReduceLROnPlateau (mode=max, patience=5)
908

909 **Training** up to 200 epochs, early stopping patience 40, batch size 32, shuffle (train only), 2 work-
910 ers, pin memory

911 **Data** Train 6,418 images (T1–T8), val 714 images (T1–T8), test 7,999 images (T1–T24 un-
912 seen)

913 **Transforms (train)**
914 Resize 224×224 , RandomHorizontalFlip, RandomRotation(10°), normalization
915 mean=std=[0.5]
916

917 **Transforms (val)**
 Resize 224×224 , normalization mean=std=[0.5]

918	B.10 ViT-B/16
919	
920	Model Vision Transformer (vit_base_patch16_224, timm pretrained)
921	Input 224×224 RGB images
922	Output Heads
923	9 (one per label column)
924	Loss Multi-task CrossEntropyLoss (sum over tasks)
925	Optimizer
926	Adam, learning rate 1.0×10^{-4}
927	Scheduler
928	ReduceLROnPlateau (mode=max, patience=5, factor=0.5)
929	
930	Training 200 epochs, early stopping 40, batch size 32, shuffle (train only), 2 workers, pin memory
931	Data Train 6,418 images (T1–T8), val 714 images (T1–T8), test 7,999 images (T1–T24 un-
932	seen)
933	Transforms
934	Same as ConvNeXt
935	
936	B.11 CoATNET-0
937	
938	Model CoAtNet-0 (timm coatnet_0_224, not pretrained)
939	Input 224×224 RGB images
940	Output Heads
941	9 (one per label)
942	Loss Multi-task CrossEntropyLoss
943	Optimizer
944	Adam, learning rate 1.0×10^{-4}
945	Scheduler
946	ReduceLROnPlateau (mode=max, patience=5, factor=0.5)
947	
948	Training 200 epochs, early stopping 40, batch size 32, 2 workers, pin memory
949	Data Train 6,418 images (T1–T8), val 714 images (T1–T8), test 7,999 images (T1–T24 un-
950	seen)
951	Transforms
952	Same as ConvNeXt
953	
954	B.12 IMAGE SHAPE FUSION TRANSFORMER
955	
956	Model ConvNeXt-Tiny backbone (ImageNet pretrained) + shape features projected as tokens
957	fused via Transformer Encoder + multi-head classifier
958	Fusion Image token + 9 shape tokens → Transformer Encoder
959	Transformer
960	$d_{model}=256$, n_heads=4, n_layers=3, feedforward dim=512, dropout=0.1
961	Shape Features
962	9 geometric features (area, perimeter, eccentricity, solidity, extent, equivalent_diameter,
963	major_axis_length, minor_axis_length, circularity), z-normalized
964	Loss Class-weighted CrossEntropyLoss
965	Optimizer
966	AdamW, learning rate 1.0×10^{-4} , weight decay 1.0×10^{-2}
967	Scheduler
968	ReduceLROnPlateau (mode=max on validation macro-F1, patience=5)
969	
970	Training up to 200 epochs, early stopping patience 40, batch size 32, 4 workers, pin memory
971	Data Train 6,418 images (T1–T8), val 714 images (T1–T8), test 7,999 images (T1–T24 un-
	seen)

- 972 **Transforms**
 973 Same resizing/augmentation as ConvNeXt
 974
- 975 **B.13 HIERARCHICAL MULTI-TASK TRANSFORMER (HMTT)**
 976
- 977 **Model** ViT-B/16 encoder (ImageNet pretrained) + hierarchical multi-task heads conditioned on
 978 causal label order
- 979 **Label Order**
 980 cell_line → culture_medium → seeding_density → magnification → microscope → time-
 981 point → biological_rep → technical_rep
- 982 **Output Heads**
 983 8 (one per label)
 984
- 985 **Loss** CrossEntropyLoss with class weights; optional Focal Loss ($\gamma=2.0$, $\alpha=0.25$)
- 986 **Optimizer**
 987 AdamW, learning rate 1.0×10^{-4} , weight decay 1.0×10^{-2}
- 988 **Scheduler**
 989 ReduceLROnPlateau (mode=max validation macro-F1, patience=5)
- 990 **Training** up to 200 epochs, early stopping 40, batch size 32, 2 workers, pin memory, mixed preci-
 991 sion (AMP)
- 992 **Data** Train 6,418 images (T1–T8), val 714 images (T1–T8), test 7,999 images (T1–T24 un-
 993 seen)
- 994 **Transforms**
 995 Same as ConvNeXt
 996
 997
- 998 **B.14 CONV LSTM**
 999
- 1000 **Data** Image size 128; sequence length 2; min gap 1; max gap 3; train/val/test split 70/15/15
- 1001 **Training** Batch size 8; 4 workers; 200 epochs; learning rate 1.0×10^{-4} ; optimizer Adam; loss L1;
 1002 shuffle (train) True, (val/test) False; no gradient clipping or scheduler; mixed precision
 1003 disabled; early stopping metric SSIM with patience 40
- 1004 **Model** ConvLSTM architecture; hidden dimension 32; kernel size 3; decoder Conv2d → 1
 1005 channel
- 1006
- 1007 **B.15 PREDRNN++**
 1008
- 1009 **Data** Image size 128; sequence length 2; min frame gap 1; max frame gap 3; dataset split
 1010 70/15/15; batch size 8; 4 workers
- 1011 **Model** PredRNN++ with two SpatioTemporal LSTM layers + Conv decoder; input channels 1;
 1012 hidden dimensions [32, 64]; kernel size 3; decoder Conv2d (64→1, kernel 1)
- 1013 **Training** 200 epochs; optimizer Adam; learning rate 1.0×10^{-4} ; loss L1; early stopping patience
 1014 40 (on SSIM);
 1015
- 1016 **B.16 MMFUSIONNET**
 1017
- 1018 **Data** Image size 128; sequence length 2; min gap 1; max gap 3; train/val/test split 70/15/15
- 1019 **Model** Base channels 32; encoder 4 ConvBlocks with pooling; bottleneck 256 channels; decoder
 1020 symmetric upsampling with skip connections; output 1-channel (sigmoid); categorical
 1021 embedding dim 32; continuous features 2 (seeding density, timepoint); metadata fusion
 1022 MLP (256→256) + FiLM conditioning
- 1023 **Training** Batch size 8; 4 workers; 200 epochs; learning rate 1.0×10^{-4} ; optimizer AdamW; sched-
 1024 uler ReduceLROnPlateau (patience 10, factor 0.5); loss $0.8 \times \text{MSE} + 0.2 \times (1 - \text{SSIM})$ if
 1025 MS-SSIM available, else $0.6 \times \text{MSE} + 0.4 \times \text{L1}$; early stopping patience 40; visual sam-
 ples saved every 5 epochs; shuffle train=True, val=False

B.17 PHYDNET

- Data** Image size 128; context frames 2; prediction steps 1; frame gap 1–3; dataset split 70/15/15; batch size 8; 4 workers; seed 42
- Model** PhyDNet encoder + PhyCell + ConvLSTM residual + decoder; input channels 1; encoder channels 64; physics channels 32; residual channels 32; physics operators per channel 3; physics kernel size 3
- Training** 200 epochs; optimizer Adam; learning rate 1.0×10^{-4} ; loss L1; mixed precision enabled (AMP); scheduler ReduceLROnPlateau (factor 0.5, patience 8); early stopping patience 40 on SSIM; best model

C ABLATION STUDY

To assess the contribution of each component of our fusion architecture, we conducted a series of ablation experiments. Five model variants were evaluated: (i) the full multimodal fusion model, (ii) an image-only model that removes all morphology (shape) features, (iii) a shape-only model that excludes image information, (iv) a no-fusion model where image and morphology pathways are present but not interactively fused, and (v) a no-transformer variant that retains both modalities but removes the transformer-based fusion mechanism. All models were trained under identical conditions, using the same dataset splits, augmentation strategy, and early stopping criteria to ensure a fair comparison. Expanded per-attribute ablation metrics are provided in Appendix X.

Table 7 summarizes the performance of each variant, reporting the average multi-label accuracy, precision, recall, and F1 score across all metadata attributes. The full model achieves the strongest performance (accuracy = 0.8526), demonstrating that jointly modeling image appearance with quantitative morphology yields the most reliable metadata prediction. The image-only model performs slightly worse (0.8516), indicating that the visual signal in raw microscopy images is the dominant factor driving prediction quality. Nevertheless, morphology contributes complementary structural cues that improve consistency, even if modestly.

Model Variant	Accuracy	Precision	Recall	F1 Score
Full Model (Image + Shape + Transformer Fusion)	0.8526 ± 0.00042	0.8672 ± 0.00046	0.8758 ± 0.00044	0.8681 ± 0.00045
Image-Only	0.8516 ± 0.00045	0.8651 ± 0.00049	0.8734 ± 0.00047	0.8662 ± 0.00048
No-Transformer	0.8513 ± 0.00047	0.8638 ± 0.00050	0.8726 ± 0.00049	0.8650 ± 0.00049
No-Fusion	0.8478 ± 0.00053	0.8597 ± 0.00056	0.8689 ± 0.00055	0.8610 ± 0.00056
Shape-Only	0.5381 ± 0.00112	0.5515 ± 0.00118	0.5591 ± 0.00122	0.5530 ± 0.00120

Table 7: Ablation study results for multimodal metadata prediction. Values show mean accuracy, precision, recall, and F1 score with estimated 95% confidence intervals computed using the normal approximation ($1.96 \times SD/\sqrt{N}$, with $N \approx 8000$ validation samples).

The shape-only model shows a substantial drop in accuracy (0.5381), confirming that morphological descriptors alone are not sufficient for robust metadata inference. While morphology encodes meaningful geometric cues, these cues lack the fine-grained appearance information captured directly from images, explaining the large performance gap. The no-fusion model (0.8478) also underperforms relative to the full model, indicating that simply concatenating or parallelizing the two modalities is less effective than an integrated fusion mechanism. This validates our design choice to allow deeper interaction between modalities rather than treating them independently.

Finally, removing the transformer fusion block results in a small but consistent reduction in performance (accuracy = 0.8513). The transformer contributes to refining long-range cross-modal relationships, particularly in cases where morphological features only partially align with visual appearance patterns. Overall, these results highlight that (1) image features carry the majority of predictive signal, (2) morphology provides meaningful but secondary complementary information, and (3) explicit fusion, especially transformer-based fusion, yields the strongest and most stable performance with minimal computational overhead. Each component therefore plays a distinct role in enhancing metadata prediction accuracy.