# MEDIAR: Harmony of Data-Centric and Model-Centric for Multi-Modality Microscopy

**Gihun Lee\***
KAIST AI
opcrisis@kaist.ac.kr

**SangMook Kim\***
KAIST AI
sangmook.kim@kaist.ac.kr

**Joonkee Kim\***
KAIST AI
joonkeekim@kaist.ac.kr

**Se-Young Yun**†
KAIST AI
yunseyoung@kaist.ac.kr

## Abstract

Cell segmentation is a fundamental task for computational biology analysis. Identifying the cell instances is often the first step in various downstream biomedical studies. However, many cell segmentation algorithms, including the recently emerging deep learning-based methods, still show limited generality under the multi-modality environment. *Weakly Supervised Cell Segmentation in Multi-modality High-Resolution Microscopy Images*[1] was hosted at NeurIPS 2022 to tackle this problem. We propose MEDIAR, a holistic pipeline for cell instance segmentation under multi-modality in this challenge. MEDIAR harmonizes data-centric and model-centric approaches as the learning and inference strategies, achieving a **0.9067 F1-score** at the validation phase while satisfying the time budget. To facilitate subsequent research, we provide the source code and trained model as open-source: https://github.com/Lee-Gihun/MEDIAR.

## 1 Introduction

Identifying cell organisms in microscopy images is fundamental for various biomedical applications [3, 46, 58, 62, 66]. By partitioning the high-content images into the interested regions, segmenting cell instances is often the first step to extracting meaningful biological signals [4, 16, 21, 33, 40, 56]. As a typical microscopy system generates thousands of images in a session [4, 36], an automated computational approach enables the large-scale simultaneous comprehensive analysis [2, 16, 40].

With the recent advances in deep learning (DL) in a wide range of vision tasks [20, 27, 35, 37, 57], DL methods have been widely adopted in microscopy image analysis [4, 5, 14, 22, 59, 63], showing remarkable success. However, training the deep neural networks often requires a large number of labeled data [37, 40, 43, 55], and learning on the datasets with limited diversity leads the poor generalization of the model [23, 38, 54]. Such an issue is more enlarged in microscopy imaging datasets, where manually annotating cells is highly labor-intensive and time-consuming [18, 22].

*Weakly Supervised Cell Segmentation in Multi-modality High-Resolution Microscopy Images* (CellSeg Challenge) was hosted at NeurIPS 2022 to tackle this problem. By learning from the 1,000 labeled and 1,500+ unlabeled microscopy images, the competition aims to conduct cell instance segmentation for various situations. Although the images consist of various microscopy types, tissue types, and staining types, the metadata for the image (e.g., modality-related annotation) is not provided.

---

*Equal Contribution. † Corresponding Author.
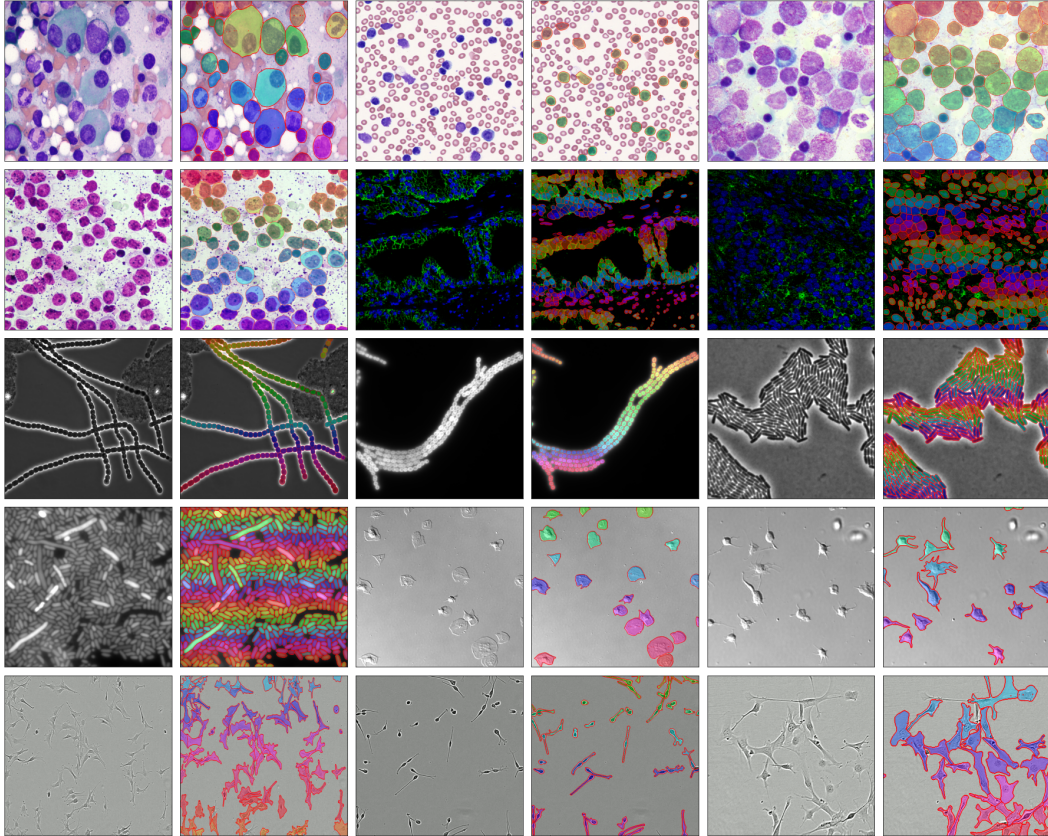[1]https://neurips22-cellseg.grand-challenge.org/

Figure 1: MEDIAR prediction results on *NeurIPS 2022 CellSeg Challenge* validation images. Our proposed method identifies the cell instances evenly well across different modalities.

The solution is evaluated by two criteria as follows:

- **Prediction Performance:** F1-score evaluated at the IoU threshold 0.5 for the true positive.
- **Time Efficiency:** The running time exceeding the time tolerance:

$$\text{Time Tolerance}(H, W) = \begin{cases} 10s & \text{if } H \times W \leqslant 10^6 \\ \frac{H \times W}{10^6} 10s & \text{if } H \times W > 10^6, \end{cases}$$

where $H$ and $W$ each stands for the height and width of the image.

In this paper, we propose MEDIAR, a framework to build a single generalist model for cell instance segmentation by harmonizing data-centric and model-centric approaches. On the data-centric side, MEDIAR starts from an extensive pretraining and replays data to retain the knowledge from the pretraining. Seeking balanced training towards heterogeneous modalities, MEDIAR discovers the latent modality and amplifies lacking modality samples. On the model-centric side, MEDIAR consists of a model structure with two separated heads each for cell recognition and instance distinction. To perform seamless prediction on large-scale images, we propose a stochastic test-time augmentation strategy combined with ensemble prediction. As a result, MEDIAR shows remarkable performance on a variety of microscopy images with multi-modalities while satisfying the time efficiency in the sense of tolerance budget.

To summarize, our main approaches are as follows:

- We suggest an overview of the factors that make the generalization of cell segmentation difficult. Not only the different microscopy technology but also the imaging protocol, cell types, cell shapes, and even the magnification can be the source of the instance-level heterogeneity **(Section 2)**

- We propose **MEDIAR**, a framework to conduct cell segmentation under multi-modality by a single generalist model. By combining data-centric and model-centric approaches, our method performs cell segmentation evenly well across the modalities **(Section 3)**.

- In the Data-Centric view, we provide a learning strategy to balance the latent modalities and retain the knowledge from pretraining data **(Section 4)**. In the Model-Centric view, we provide a model structure to capture cell instances and a corresponding inference strategy to conduct prediction for high-resolution images under the time budget **(Section 5)**.

- We analyze each component's effect in our approach, discuss the key factors of our success, and introduce the open problems. We further release the trained model, which achieved the F1-score **0.9067** on multi-modality microscopy challenge datasets (Figure 1 visualizes the prediction results), to facilitate the subsequent studies **(Section 6)**.

## 2   Difficulties in Multi-Modality Cell Segmentation
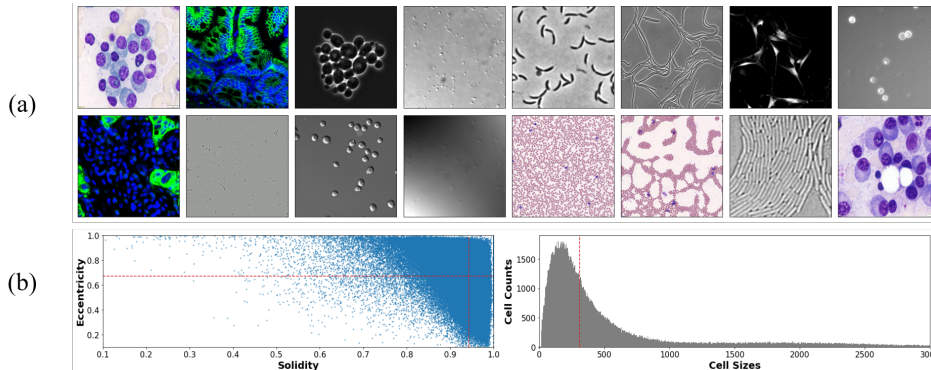


Figure 2: An overview of microscopy images in the CellSeg datasets. (a) images from various modalities. (b) statistics of the cells. The red dotted lines stand for the median value of each measure.

**Microscopy types & Tissue types**   Although the first source of multi-modality in microscopy images is the difference in the microscopy technologies used for the imaging (e.g., brightfield, fluorescent, phase-contrast, and differential interference contrast), the modality may broadly differ by the tissue type. We visualize the examples of different modalities in Figure 2(a). Sometimes the two sources of image-level heterogeneous modality are highly correlated, as a specific microscopy technology can have an advantage in observing particular tissue types.

**Cell shapes and sizes**   Another source of modality originates from cell types, which implies instance-level heterogeneity. To explore the distribution of cell shapes, we use three measures (i) Eccentricity [64] $= \frac{\text{Axis}_{\text{short}}}{\text{Axis}_{\text{long}}}$ (ii) Solidity [64] $= \frac{\text{Area}}{\text{Convex Area}}$, (iii) Cell Size (pixels in each cell object). Note that eccentricity measures the minor and major axis ratio, and solidity measures the object's density using its convex hull. As visualized in Figure 2(b), the cells in the images have various shapes and sizes. Note that the cells belonging to the same image may have different shapes or sizes, depending on their cell phase or the microscopy magnification.

**Annotation Inconsistency**   As the criteria for the cell annotation may vary across the annotators, the heterogeneity also comes at the label-level. For example, the different standards on (i) discarding the cells in the image boundary, (ii) contour shapes of the cell, (iii) cell recognition on the object, and (iv) cell boundary could be the source of annotation inconsistency. Such noise in the data labels often degrades the performance after training [49, 67].

**Others**   Although not directly related to the multi-modality, the cell images are sometimes contaminated or damaged during the staining and microscopy scanning, making cell recognition more challenging. Another critical issue that makes cell instance segmentation difficult is touching cell objects. As the cells are often closely located with other cells without explicit object boundaries, it is often hard to assign pixels with almost similar signals to each object. Those factors also need to be considered carefully in deploying the cell instance segmentation method.

3

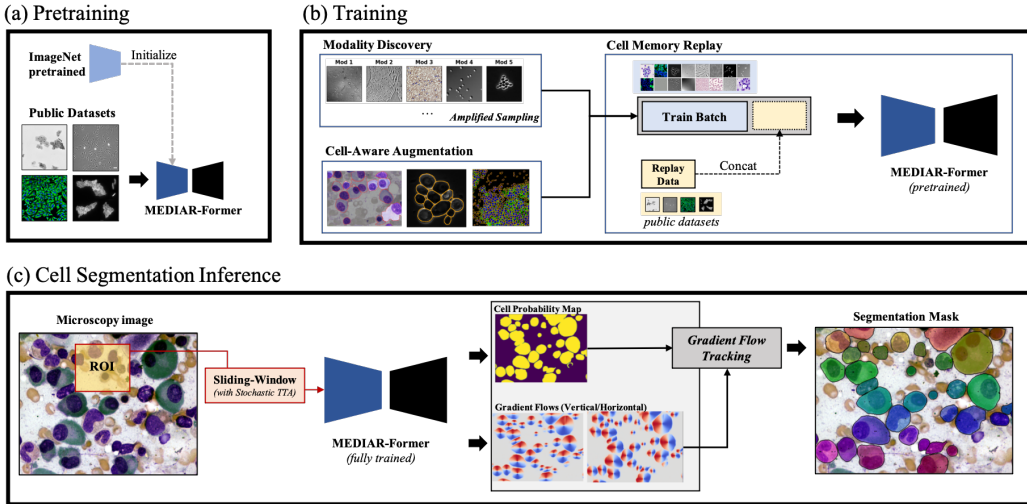# 3   MEDIAR: Harmony of Data-centric and Model-centric Approaches



Figure 3: An overview of MEDIAR framework.

The heterogeneous modalities come from various factors inevitably induce a naive approach to be felt into imbalanced learning towards specific dominant modalities. Our approach consists of two main streams to overcome this challenge: the data-centric approach and the model-centric approach. On the data-centric side, we hypothesize that (i) extensive pretraining helps capture modality-invariant features in the cell, and (ii) learning on the modality-balanced data improves generalization across modalities. On the model-centric side, we hypothesize that (i) we can train the generalist model by reducing interference between different modalities and (ii) merging multiple predictions using ensemble and test-time augmentation results in better generalization.

We illustrate an overview of the pipeline of our approach in Figure 3, which harmonizes the data-centric and model-centric approaches. Each phase in the pipeline ((a) pretraining, (b) training, (c) cell segmentation inference) is closely related to the remarkable performance of our proposed method MEDIAR. Note that although our framework consists of several components in each phase, they are mostly orthogonal. This implies that one can modify a part of our framework for further improvement without performance degradation due to their dependencies. In the following sections, we provide the details of key components in our proposed framework and how they contribute to the performance.

## 4   MEDIAR - Data-centric Approaches

### 4.1   Cell-Aware Augmentation

MEDIAR starts from using a intensive combination of data augmentation strategies. With the prevalent augmentation methods, we propose two novel cell-aware augmentations to improve generalization. At first, as the intensity of the cells can differ in the same image in the test time, we cell-wisely randomize the intensity in the image (Cell Intensity Diversification). Second, we excluded the boundary pixels in the label. The boundary exclusion is adopted *only in the pretraining phase*. We provide the combined augmentation policy in Table 1, and visualize the examples in Figure 4.
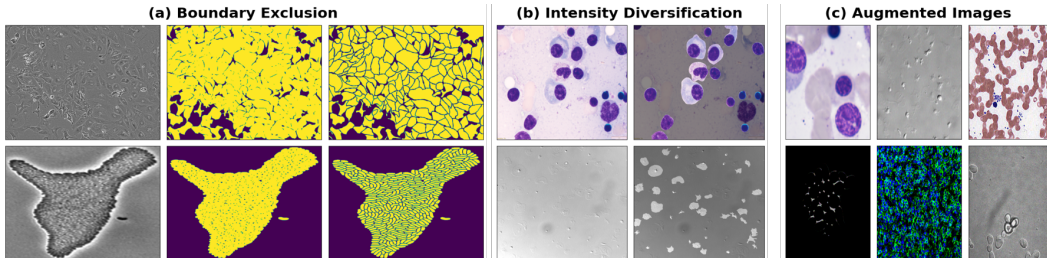


Figure 4: Examples of the proposed policies and augmented images.

Table 1: Augmentation strategies applied in the MEDIAR.

| Strategy | Implementation Details |
|---|---|
| **(P)** Clip (.) | Clip the image pixels into percentile range [0.0, 99,5]. |
| **(P)** Normalization (.) | Normalize the images to $N(\mu = 0, \sigma = 1.0)$ |
| **(P)** Scale Intensity (.) | Scale image pixels into the range [0.0, 1.0]. |
| **(S)** Zoom (0.5) | Zooming to the scale in the range [0.25, 1.5] using nearest interpolation. |
| **(S)** Spatial Crop (1.0) | Cropping images as size (512, 512) at a randomly chosen center. |
| **(S)** Axis Flip (0.5) | Flip the array axis which corresponds to the image channels. |
| **(S)** Rotation (0.5) | Spatially rotate array by 90 degrees (i.e., 90°, 180°, 270°). |
| **(I)** Cell-Aware Intensity (0.25) | The intensity scale of each cell is scaled to the range [1.0, 1.7] |
| **(I)** Gaussian Noise (0.25) | Add Gaussian noise $N(\mu = 0, \sigma = 0.1)$ to the image. |
| **(I)** Contrast Adjustment (0.25) | Change image intensity by factor $\gamma \in [0.0, 2.0]$ |
| **(I)** Gaussian Smoothing (0.25) | Smoothing with Gaussian Filter with $\sigma$=1.0 |
| **(I)** Histogram Shift (0.25) | Non-linear transformation to the intensity histogram with three control points. |
| **(I)** Gaussian Sharpening (0.25) | Sharpening by Gaussian filter with $\sigma$ factor 0.5 & 1.0 with $\alpha \in [10.0, 30.0]$. |
| **(O)** Boundary Exclusion (.) | Map the boundary pixels in the label to the background index. |

* **(P)**: Pre-processing **(S)**: Spatial Augmentation **(I)**: Intensity Augmentation **(O)**: Others
* The value in the parenthesis stands for the probability of each strategy.

## 4.2 Two-phase Pretraining and Fine-tuning

**Pretraining** We use 7,242 labeled images from four public datasets for pretraining: OmniPose [14], CellPose [63], LiveCell [18] and DataScienceBowl-2018 [5]. MEDIAR takes two different phases for the pretraining. In phase 1, the MEDIAR-Former model with encoder parameters pretrained on ImageNet-1k is trained on the public datasets for 80 epochs. In phase 2, the pretrained model is further trained on the joint set of public datasets and train datasets for 60 epochs.

**Fine-tuning** The two pretrained models from phase 1 and phase 2 are fine-tuned with 200 and 25 epochs for each, using the train datasets. We observed that fine-tuned model from phase 1 predicts the modalities appear only in the target datasets. On the other hand, fine-tuning from phase 2 predicts the modalities included in both the public and target datasets. MEDIAR conducts the ensemble prediction using those two models. At phase 2 fine-tuning, we relabel the images, which shows the misaligned prediction between the pretrained model and the phase 1 fine-tuned model to compensate for the possible noisy labels.

## 4.3 Modality Discovery & Amplified Sampling

By observing the datasets, we find that the number of modalities is more than four. Moreover, the distribution of the number of modalities differs, resulting in an imbalanced dataset. This data imbalance may cause degradation of the model's performance for minor modalities.
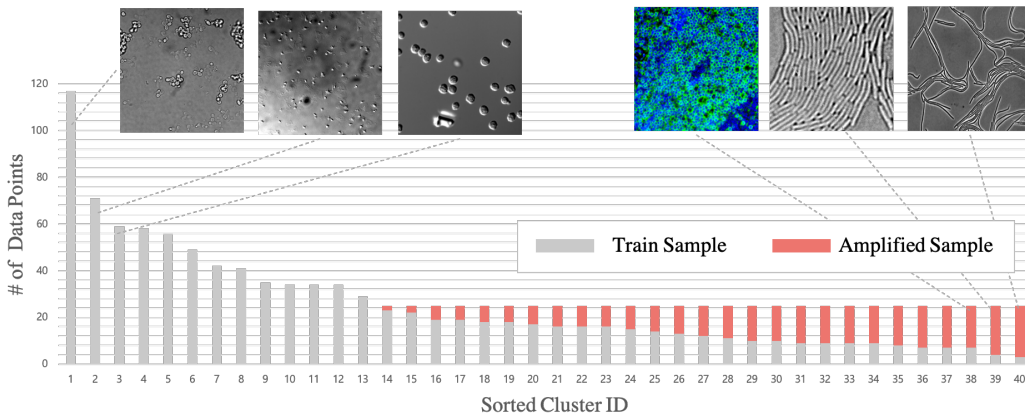


Figure 5: Discovered modalities and amplified samples in the training dataset.

5

We discover the latent modalities and balance their sampling ratio during training to overcome this issue. We group the encoder embeddings from the phase-1 pretrained model via the $k$-means clustering algorithm. We set the number of clusters as 40, which is large enough to sufficiently filter minor modality embeddings. To smooth the sampling ratio towards modalities, we over-sample the minor data samples. The illustration in Figure 5 summarizes our balancing strategy.

## 4.4 Cell Memory Replay

We find that the fine-tuned model performs well in most cases but degrades on some of the modalities in which the pretrained model performs well. We hypothesize that the phenomenon resembles forgetting issue [44, 52] in Continual Learning[47], and memory replay [9, 10] mitigates the problem. We concatenate the data from the public dataset with a small portion to the batch and train with boundary-excluded labels.

## 5 MEDIAR - Model-centric Approaches
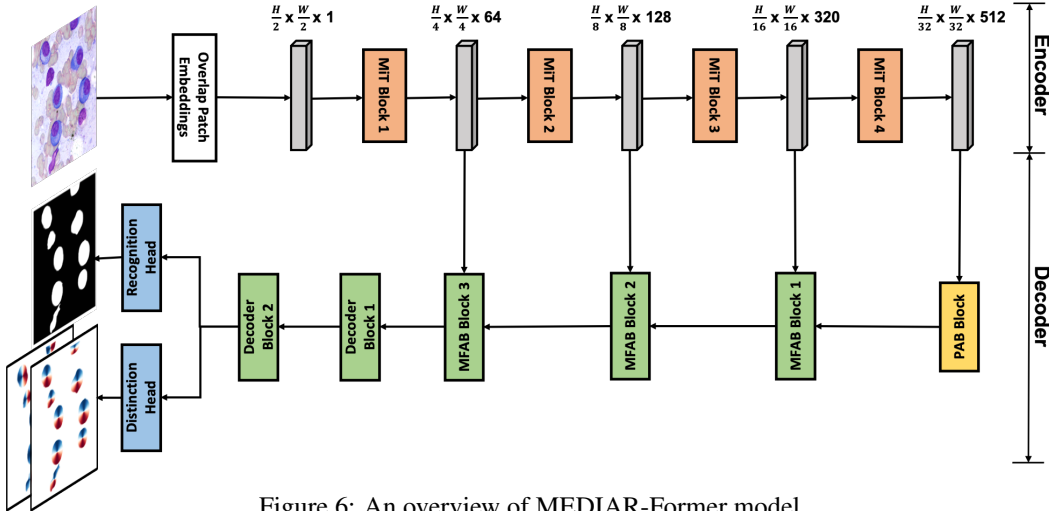
### 5.1 MEDIAR-Former Architecture



Figure 6: An overview of MEDIAR-Former model.

**Encoder & Decoder**    MEDIAR-Former follows the design paradigm of U-Net [57], which allows the hierarchical feature integration from the encoder to the decoder via skip connection. For the encoder and decoder, we adopt SegFormer [68] and MA-Net [19]. The multi-scale features extracted from the encoder are concatenated through skip-connection for the decoder outputs. We use Mish [50] both in the encoder and decoder for better generalization.

The SegFormer encoder consists of multiple MiT Blocks, which include three main components as efficient self-attention, Mix-FFN, and Overlap Patch Mapping. In the efficient self-attention, the conventional self-attention [65] in ViT [17] with the $O(N^2)$ computation complexity is reduced to $O(N^2/R)$, by modifying $K$ in the self attention $\text{Attention}(Q, K, V) = \text{Softmax}QK^T/\sqrt{d_{head}}V$, as $K = \text{Linear}(C \cdot R, C)(\hat{K})$, where $\hat{K} = \text{Reshape}(N/R, C \cdot R)(K)$. In Mix-FFN, 3x3 convolution takes place instead of positional encoding, and Overlap Patch Mapping is recurrently applied to preserve the local continuity between patches. The Decoder MA-Net consists of two key modules, Position-wise Attention Block (PAB) for feature inter-dependencies in global view and Multi-scale Fusion Attention Block (MFAB) for semantic multi-scale feature map fusion.

**Head**    MEDIAR-Former uses two separate heads: *Cell Recognition* (CR) head and *Cell Distinction* (CD) head. Although the prior works use the single-head structure for each outputs [14, 63], semantic prediction for the objects and regression for spatial gradient field interferes with each other when predicted from the same feature space, as prevalent in the Multi-Task Learning [8, 60]. To mitigate this issue, we separate both heads and use two 3x3 Conv-heads with BatchNorm [31].
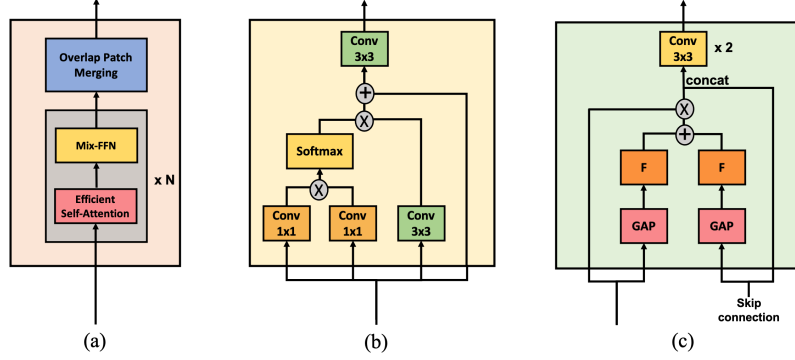
Figure 7: The details of (a) MiT Block, (b) PAB Block, (c) MFAB Block in the Figure 6.

**Learning Objective**   MEDIAR-Former is learned to predict cell binary mask $y^{\text{cell}}$ from CR head ($h_{\text{CR}}$) by the binary cross-entropy loss $\mathcal{L}_{\text{BCE}}$, and cell topological maps $y^{\text{gradient}}$ from CD head ($h_{\text{CD}}$) by mean-square error loss $\mathcal{L}_{\text{MSE}}$. Here, the topological maps are normalized image gradient field generated by pseudo-diffusion as in [63]. The learning objective is as follows:

$$\mathcal{L}(x,y) = \mathcal{L}_{\text{BCE}}(h_{\text{CR}}(f_\theta(x)), y^{\text{cell}}) + \lambda \cdot \mathcal{L}_{\text{MSE}}(h_{\text{CD}}(f_\theta(x)), y^{\text{gradient}}), \qquad (1)$$

where $f_\theta$ is the decoder output given the data $x$. We set $\lambda$ as 0.5 in all experiments.

## 5.2   Gradient Flow Tracking

MEDIAR adopts the gradient flow tracking in CellPose [63]. After filtering the cell candidates, all pixels aggregated into the cell indices by iteratively following the spatial gradient fields. First, the unit vectors are created by normalizing the gradients of each pixel. Second, the mesh grid is generated for the spatial directions, and the values are converted to the smaller starting or ending values of the unit vector directions. Third, the masks are initialized from the peak indices from the mesh grid histogram and extended until convergence. Finally, the error between pseudo-diffusion and gradient field is measured to decide whether to accept it as a cell object. We set the error threshold as 0.4. To improve time efficiency for the whole slide images, MEDIAR conducts the gradient flow tracking in a non-overlapped sliding window manner, with the patch size $2000 \times 2000$.
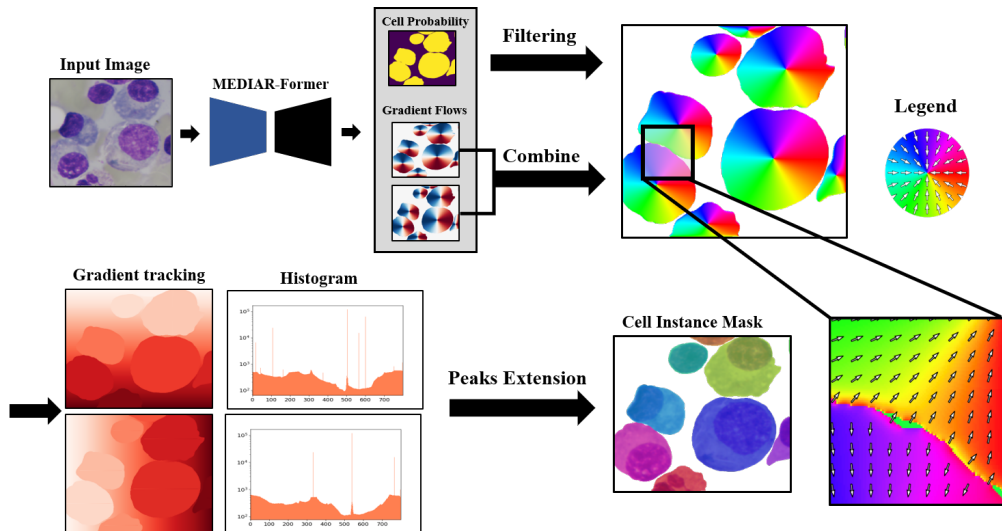


Figure 8: An example of gradient flow tracking for cell instance segmentation.

7

## 5.3 Ensemble Prediction with TTA

To conduct inference on large-size input images, MEDIAR uses sliding-window inference [61] with the overlap size between the adjacent patches as 0.6. To predict the different views on the image, MEDIAR uses Test-Time Augmentation (TTA). Each image is horizontally or vertically flipped, and the outputs are summed up. During each prediction, MEDIAR generates an importance map from the Gaussian Filter ($\sigma = 0.125$) for each patch. Multiplying the importance map for each patch prediction output prevents recognizing the same cell at the patch boundary as multiple cells. The ensemble final prediction uses the two fine-tuned models from each pretraining phase. The modalities that the phase 1 and phase 2 models predict well are slightly different.

## 6 Experiments

### 6.1 Experimental Setups

**Implementation details**    We use AdamW [42] optimizer with an initial learning rate of 5e-5 in pretraining, both in the first and second phases, and 2e-5 in fine-tuning. The learning rate is decayed using cosine scheduler [41] using 100 interval without restarts. The code is implemented using PyTorch [53] and MONAI library [13] with some modifications. The base model structure and ImageNet1K-pretrained parameters are from PyTorch segmentation package [29]. We use 2 A5000 GPU cards but without the Multi-GPU training. We use mixed precision training [48] as FP-16, which reduces memory usage during training without performance degradation. We further specify the development environment in Table 2 and training hyperparameters in Table 3. More details are provided in the released source code.

Table 2: Development environments and requirements.

| Environment | Specification |
|---|---|
| System | Ubuntu 18.04.5 LTS |
| CPU | AMD EPYC 7543 32-Core Processor CPU@2.26GHz |
| RAM | 500GB; 3.125MT/s |
| GPU (number and type) | NVIDIA A5000 (24GB) 2ea |
| CUDA version | 11.7 |
| Programming language | Python 3.9 |
| Deep learning framework | Pytorch [53] (v1.12, with torchvision v0.13.1) |
| Code dependencies | MONAI [13] (v0.9.0), Segmentation Models [29] (v0.3.0) |
| Specific dependencies | ttach (v0.0.3) [30] for Test-Time Augmentation |

Table 3: MEDIAR training protocols for pretraining and fine-tuning. The epochs in the parenthesis are for the phase 2 model. Note that we include a public data sample in the fine-tuning batch.

| Learning Setups | Pretraining | Fine-tuning |
|---|---|---|
| Initialization (Encoder) | Imagenet-1K [15] | from **Pretraining** |
| Initialization (Decoder & Head) | He normal init | from **Pretraining** |
| Batch size | 9 | 9 (with memory) |
| Total epochs | 80 (60) | 200 (25) |
| Optimizer | AdamW [42] | AdamW [42] |
| Initial learning rate (lr) | 5e-5 | 2e-5 |
| Lr decay schedule | Cosine [41] (100 interval) | Cosine [41] (100 interval) |
| Loss function | MSE, BCE | MSE, BCE |
| Training time | 72 hours | 48 hours |
| Number of model parameters | 121.31 M | 121.31 M |
| Number of flops | 204.26 G | 204.26 G |
| $CO_2eq$ [1] | 15.105g | 9.876g |

&ast; Parameters counter for pytorch models: https://github.com/sovrasov/flops-counter.pytorch
&ast; Flops counter for pytorch models: https://github.com/facebookresearch/fvcore
&ast; Carbon tracker for deel learning models [1]: Emission: https://github.com/lfwa/carbontracker/

8

**Public Datasets Usage** For the pretraining and fine-tuning, we gathered 7,242 labeled data from four public datasets as follows:

- **OmniPose** [14]: contains mixtures of 14 bacterial species. We only use 611 bacterial cell microscopy images and discard 118 worm images.
- **CellPose** [63] includes Cytoplasm, cellular microscopy, fluorescent cells images. We used 551 images by discarding 58 non-microscopy images. We convert all images as gray-scale.
- **LiveCell** [18]: is a large-scale dataset with 5,239 images containing 1,686,352 individual cells annotated by trained crowdsources from 8 distinct cell types.
- **DataScienceBowl 2018** [5]: 841 images contain 37,333 cells from 22 cell types, 15 image resolutions, and five visually similar groups.

Each of the datasets contains instance-level cell mask labels. We jointly combine all the collected datasets. Example images from the collected datasets are povided in Figure 9.
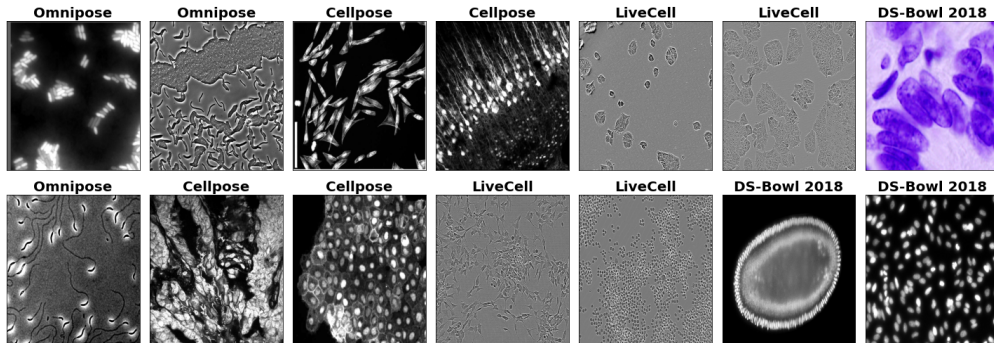


Figure 9: Example images from the collected public datasets.

## 6.2 Best Model Selection Standard

As the ground-truth label for the validation dataset is not provided, we use two different measures to select the best model checkpoint as follows:

- **F1-score** (*with* hold-out): Randomly select 10% of train samples as a hold-out set, and select the model that shows the best hold-out F1-score.
- **Cell Count** (*without* hold-out): Use all data as train samples, and count the number of predicted cells on validation datasets.

The cell count measure is based on the observation that MEDIAR has its strength in sensitivity by avoiding false-positive predictions. This is because MEDIAR discards the cells above the error threshold. We plot the change of the two measures in Figure 10.
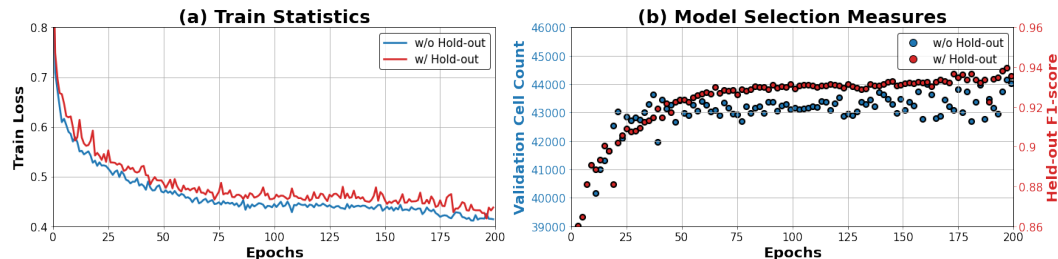


Figure 10: Change of selection measures during fine-tuning from the phase 1 pretrained model.

## 6.3 Cell Instance Segmentation Performance

The performance is evaluated via the F1-score at the IoU threshold of 0.5 for true positives. We use two models that each are fine-tuned from phase1 and phase2 pretraining. The validation F1-score by different inference strategies and models are in Table 4, and learning curves for each model is plotted in Figure 11. We use window size (512x512) with an overlap of 0.6 and a stacked error threshold of 0.4 for each cell in the inference.

9

Combined with stochastic TTA and ensemble, our MEIDAR achieves F1-score **0.9067** on validation datasets with 101 images. The prediction results on validation images are provided in Figure 1. The results shows that MEDIAR predicts cells suprisingly well in various modalities. We emphasize that our MEDIAR perfectly satisfies the time limit for each image, and relaxing the time constraints would further improve the performance.

Table 4: MEDIAR validation F1-score by different models and inference strategies.

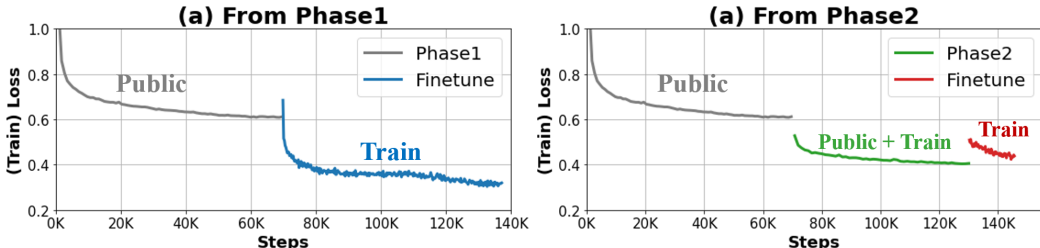| Inference Strategy | Prediction Model | | Implementation Details |
|---|---|---|---|
| | From Phase 1 | From Phase 2 | |
| MEDIAR | 0.9000 | 0.9011 | Sliding window & Gradient tracking |
| + Stochastic TTA | 0.9027 | 0.9048 | Flipping & Importance map |
| + Ensemble | **0.9067** | | Predict using the both models |



Figure 11: MEDIAR learning curves corresponding to Table 4. Note that the learning samples for each curve are different. The text above each curve stands for the samples used in the training.

## 6.4 Ablation Study

**Model Structure** To examine the model structure's effect, we compare MEDIAR Former's performance to the model structures from prior works: U-Net [57], Swin Unetr [24]. For a fair comparison, we train each model from scratch for 200 epochs without pretraining. The initial learning rate is set as 5e-5, and other hyperparameters are the same as the fine-tuning protocol specified in Table 3. The results are in Table 5, which shows that MEDIAR Former significantly outperforms the others.

Table 5: MEDIAR performance on different model structures.

| Model Structure | Model Component | (Train) Loss | (Valid) F1-score |
|---|---|---|---|
| U-Net [57] | | 0.9592 | 0.5473 |
| UNetr [25] | None | 0.6777 | 0.6034 |
| Swin Unetr [24] | | 0.6776 | 0.6320 |
| MEDIAR Former | Base Structure | 0.6206 | 0.6503 |
| | + Encoder Initialization | 0.4507 | 0.8292 |
| | + Decoder Scale-up | 0.4268 | 0.8347 |
| | + Head Separation | 0.4144 | **0.8424** |

**Data-centric Components** In Table 6, we provide the effect of each data-centric approach with the details. Although the labeling consistency does not improve solely, we observe that it predicts some modalities are more robust. We combine the labeling consistency approach when fine-tuning from phase 2, which is used for the ensemble prediction. Note that all the experiments in Table 6 use a combined augmentation strategy in Table 1.

Table 6: Effect of data-centric components with implementation details.

| Data-centric Component | (Valid) F1-score | Implementation Details |
|---|---|---|
| MEDIAR (From Phase1) | 0.8801 | Fine-tune 200 epochs from phase 1 model |
| + Cell-Aware Augmentation | 0.8881 | Apply intensity diversification augmentation |
| + Amplified Sampling | 0.8921 | Balanced sampling for discovered 40 modalities |
| + Cell Memory Replay | **0.9000** | Exclude boundary for public memory data |
| + Labeling Consistency | 0.8979 | Relabeling train datasets by phase2 model |

**Cell Memory Replay** We further investigate the effect of varying memory ratios for the fine-tuning phase. As suggested in Table 7, replaying only one memory data per batch is enough for preserving the knowledge from the pretraining phase, showing the sweet spot for the prediction performance.

Table 7: MEDIAR performance by the varying memory ratio.

| Pretrained Model | Train : Memory | (Valid) F1-score | Implementation Details |
|---|---|---|---|
| | $9 : 0$ | 0.8801 | |
| **Phase 1** | $8 : 1$ | **0.9000** | Fine-tuning for 200 epochs |
| | $7 : 2$ | 0.8881 | |
| | $9 : 0$ | 0.8876 | |
| **Phase 2** | $8 : 1$ | **0.9011** | Fine-tuning for 25 epochs |
| | $7 : 2$ | 0.8849 | |

## 6.5 Time Efficiency

We measure the time cost of MEDIAR prediction using the environment specified in Table 2 with a single A5000 GPU card. The results are plotted in Figure 12. As suggested in Figure 12(a), MEDIAR conducts most images in less than 1sec, and this depends on the image size. Note that although the complexity is not linear as we use the transformer structure in our encoder, its complexity $O(N^2/R)$, which is relaxed from $O(N^2)$. We also measure the total prediction time cost on 101 images in validation data, which includes one WSI with the size of $8415 \times 10496$. As in Figure 12(b) and Figure 12(c), the total prediction time considerably increased by using the TTA strategy, but it only slightly increases the WSI prediction. Even with the TTA and ensemble, we emphasize that MEDIAR satisfies the time budget for all images.
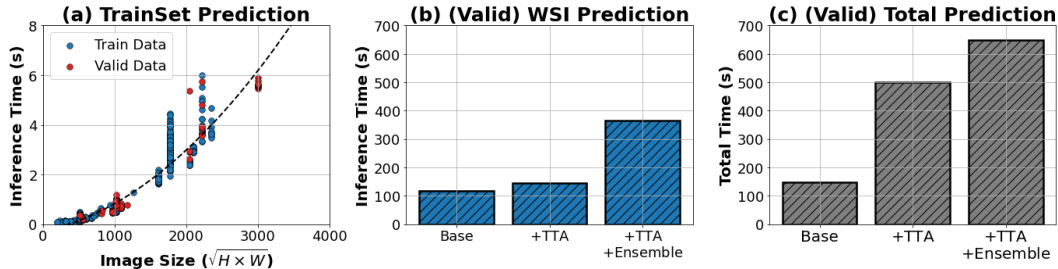


Figure 12: The time cost of the MEDIAR prediction on train datasets (1,000 images) and validation datasets (101 images). The validation datasets include one WSI.

## 6.6 Using unlabeled images

We examine some standard semi-supervised learning approaches to exploit 1,500+ unlabeled images in train datasets. (i) *Consistency Regularization*: We add a consistency loss [32, 51] term to match the prediction on the clean image and the distorted images to the model's consistency on unlabeled images. For the distortion, we use various combinations of the augmentations in Table 1. (ii) *Reconstruction Error* We add a reconstruction loss [11, 69] term by using additional head module. The model learns to reconstruct the unlabeled images. We test two reconstruction errors: one for only the pixels corresponding to the cell masks, and the other is the entire image reconstruction. (iii) *Pseudo Labeling*: Using the pretrained model and fine-tuned model, we assign pseudo label [28, 70] to the unlabeled images to use as like labeled images. Unfortunately, all the above methods could not improve performance, thereby only the 1,000 labeled images are used in MEDIAR. Although not included in our MEDIAR, we expect that self-supervised learning approaches [6, 7, 12, 26, 34, 39] can be a promising alternative direction for using the unlabeled images.

Table 8: Approaches for using unlabeled images.

| Methods | Baseline | Semi-Supervised Approach | | |
|---|---|---|---|---|
| | | Consistency Loss | Reconstruction Loss | Pseudo Labeling |
| **(Valid) F1-score** | **0.8801** | 0.8720 | 0.8798 | 0.8655 |

## 6.7 Failure cases of MEDIAR prediction

Although our MEDIAR performs cell instance segmentation surprisingly well in various situations, it suffers from capturing cells in a few cases. We provide the failure cases in Figure 13 with the categorized failure types. At first, when the cell regions in microscopy image are distorted, MEDIAR sometimes captures the organisms in the contaminated area as cells (*Contaminated*) or drops the cells in the region (*Missing Boundary, Blurred Staining*). Second, when the cell consists of only a few pixels (*Extremely small*), or the shape has an extraordinarily irregular structure (*Irregular Structure*), the cells are not recognized on occasion. As humans can capture cells even in those cases using their prior knowledge, we expect that integration of prior knowledge may further improve the robustness of prediction. On the other hand, some ambiguous objects are captured as cells (*Ambiguous*). Those cases may result from just a particular type of noise from the staining process or specific cell phenomena (e.g., apoptosis [45]), which depends on the cell recognition criteria.
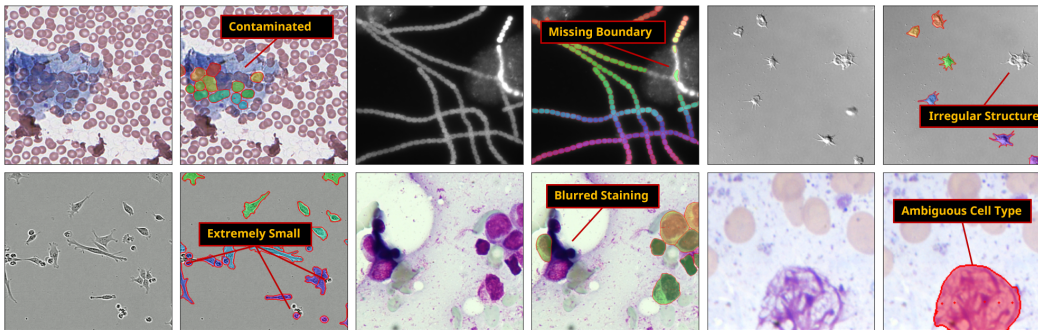


Figure 13: MEDIAR failure cases across different modalities. We zoomed in on the images for clear visualization and categorized the failure types. Note that the magnifications are different.

## 7 Conclusion

This study investigates the difficulties in cell segmentation on multi-modality microscopy images and proposes a robust and generalizable algorithm, MEDIAR, to conduct cell instance segmentation in various situations. To overcome the modality heterogeneity, we harmonize data-centric and model-centric approaches. On the data-centric side, we suggest a strategy to balance modalities during training and a comprehensive pretraining strategy with replaying their knowledge in fine-tuning. On the model-centric side, we propose a model structure to recognize cell regions and identify each cell object with a corresponding efficient inference method. Our MEDIAR shows remarkable success in various microscopy images and identifies the cell instances well across different modalities.

**Broader Impact** We believe that automated analysis of microscopy images is a crucial first step for many bio-medical applications. Providing the trained model with open-source code release may facilitate the advance of biomedical research. However, despite the remarkable success of MEDIAR in microscopy images, it sometimes needs to improve in recognizing cells depending on the imaging quality or when the cells in the image have inconsistent sizes or shapes. Although tuning the model weights using additional datasets can be the solution, the bio-medical participators should consider this problem before deploying the method. Furthermore, as the MEDIAR framework does not use unlabeled datasets, how to properly incorporate approaches for unlabeled datasets would be a promising extension for MEDIAR.

## Acknowledgement

# References

[1] Lasse F Wolff Anthony, Benjamin Kanding, and Raghavendra Selvan. Carbontracker: Tracking and predicting the carbon footprint of training deep learning models. *arXiv preprint arXiv:2007.03051*, 2020.

[2] Théo Aspert, Didier Hentsch, and Gilles Charvin. Detecdiv, a generalist deep-learning platform for automated cell division tracking and survival analysis. *Elife*, 11:e79519, 2022.

[3] Michael Boutros, Florian Heigwer, and Christina Laufer. Microscopy-based high-content screening. *Cell*, 163(6):1314–1325, 2015.

[4] Juan C Caicedo, Sam Cooper, Florian Heigwer, Scott Warchal, Peng Qiu, Csaba Molnar, Aliaksei S Vasilevich, Joseph D Barry, Harmanjit Singh Bansal, Oren Kraus, et al. Data-analysis strategies for image-based cell profiling. *Nature methods*, 14(9):849–863, 2017.

[5] Juan C Caicedo, Allen Goodman, Kyle W Karhohs, Beth A Cimini, Jeanelle Ackerman, Marzieh Haghighi, CherKeng Heng, Tim Becker, Minh Doan, Claire McQuin, et al. Nucleus segmentation across imaging experiments: the 2018 data science bowl. *Nature methods*, 16(12):1247–1253, 2019.

[6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020.

[7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.

[8] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.

[9] Arslan Chaudhry, Albert Gordo, Puneet Kumar Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. *arXiv preprint arXiv:2002.08165*, 2(7), 2020.

[10] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc'Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.

[11] Shuai Chen, Gerda Bortsova, Antonio García-Uceda Juárez, Gijs van Tulder, and Marleen de Bruijne. Multi-task attention-based semi-supervised learning for medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 457–465. Springer, 2019.

[12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

[13] MONAI Consortium. Monai: Medical open network for ai, September 2022. If you use this software, please cite it using these metadata.

[14] Kevin J Cutler, Carsen Stringer, Teresa W Lo, Luca Rappez, Nicholas Stroustrup, S Brook Peterson, Paul A Wiggins, and Joseph D Mougous. Omnipose: a high-precision morphology-independent solution for bacterial cell segmentation. *Nature Methods*, pages 1–11, 2022.

[15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[16] Shujian Deng, Xin Zhang, Wen Yan, Eric I Chang, Yubo Fan, Maode Lai, Yan Xu, et al. Deep learning in digital pathology image analysis: a survey. *Frontiers of medicine*, 14(4):470–487, 2020.

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[18] Christoffer Edlund, Timothy R Jackson, Nabeel Khalid, Nicola Bevan, Timothy Dale, Andreas Dengel, Sheraz Ahmed, Johan Trygg, and Rickard Sjögren. Livecell—a large-scale dataset for label-free live cell segmentation. *Nature methods*, 18(9):1038–1045, 2021.

[19] Tongle Fan, Guanglei Wang, Yan Li, and Hongrui Wang. Ma-net: A multi-scale attention network for liver and tumor segmentation. *IEEE Access*, 8:179656–179665, 2020.

[20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[21] Noah F Greenwald, Geneva Miller, Erick Moen, Alex Kong, Adam Kagel, Thomas Dougherty, Christine Camacho Fullaway, Brianna J McIntosh, Ke Xuan Leow, Morgan Sarah Schwartz, et al. Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nature biotechnology*, 40(4):555–565, 2022.

[22] Noah F Greenwald, Geneva Miller, Erick Moen, Alex Kong, Adam Kagel, Thomas Dougherty, Christine Camacho Fullaway, Brianna J McIntosh, Ke Xuan Leow, Morgan Sarah Schwartz, et al. Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nature biotechnology*, 40(4):555–565, 2022.

[23] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*, 2020.

[24] Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R Roth, and Daguang Xu. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In *International MICCAI Brainlesion Workshop*, pages 272–284. Springer, 2022.

[25] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 574–584, 2022.

[26] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.

[27] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[28] Ruifei He, Jihan Yang, and Xiaojuan Qi. Re-distributing biased pseudo labels for semi-supervised semantic segmentation: A baseline investigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6930–6940, 2021.

[29] Pavel Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019.

[30] Pavel Iakubovskii. Image test time augmentation with pytorch. https://github.com/qubvel/ttach, 2020.

[31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.

[32] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. *Advances in neural information processing systems*, 32, 2019.

[33] Leeat Keren, Marc Bosse, Diana Marquez, Roshan Angoshtari, Samir Jain, Sushama Varma, Soo-Ryum Yang, Allison Kurian, David Van Valen, Robert West, et al. A structured tumor-immune microenvironment in triple negative breast cancer revealed by multiplexed ion beam imaging. *Cell*, 174(6):1373–1387, 2018.

[34] Sungnyun Kim, Gihun Lee, Sangmin Bae, and Se-Young Yun. Mixco: Mix-up contrastive learning for visual representation. *arXiv preprint arXiv:2010.06300*, 2020.

[35] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

[36] Romain F Laine, Ignacio Arganda-Carreras, Ricardo Henriques, and Guillaume Jacquemet. Avoiding a replication crisis in deep-learning-based bioimage analysis. *Nature methods*, 18(10):1136–1144, 2021.

[37] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

[38] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017.

[39] Zhaowen Li, Zhiyang Chen, Fan Yang, Wei Li, Yousong Zhu, Chaoyang Zhao, Rui Deng, Liwei Wu, Rui Zhao, Ming Tang, et al. Mst: Masked self-supervised transformer for visual representation. *Advances in Neural Information Processing Systems*, 34:13165–13176, 2021.

[40] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.

[41] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.

[42] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[43] Mufti Mahmud, M Shamim Kaiser, T Martin McGinnity, and Amir Hussain. Deep learning in mining biological data. *Cognitive computation*, 13(1):1–33, 2021.

[44] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989.

[45] Pascal Meier, Andrew Finch, and Gerard Evan. Apoptosis in development. *Nature*, 407(6805):796–801, 2000.

[46] Erik Meijering. Cell segmentation: 50 years down the road [life sciences]. *IEEE signal processing magazine*, 29(5):140–145, 2012.

[47] Martial Mermillod, Aurélia Bugaiska, and Patrick Bonin. The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers in psychology*, 4:504, 2013.

[48] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.

[49] Shaobo Min, Xuejin Chen, Zheng-Jun Zha, Feng Wu, and Yongdong Zhang. A two-stream mutual attention network for semi-supervised biomedical segmentation with noisy labels. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4578–4585, 2019.

[50] Diganta Misra. Mish: A self regularized non-monotonic neural activation function. *arXiv preprint arXiv:1908.08681*, 4(2):10–48550, 2019.

[51] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020.

[52] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.

[53] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[54] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019.

[55] Samira Pouyanfar, Saad Sadiq, Yilin Yan, Haiman Tian, Yudong Tao, Maria Presa Reyes, Mei-Ling Shyu, Shu-Ching Chen, and Sundaraja S Iyengar. A survey on deep learning: Algorithms, techniques, and applications. *ACM Computing Surveys (CSUR)*, 51(5):1–36, 2018.

[56] Aditya Pratapa, Michael Doron, and Juan C Caicedo. Image-based cell phenotyping with deep learning. *Current Opinion in Chemical Biology*, 65:9–17, 2021.

[57] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[58] Johannes Schindelin, Ignacio Arganda-Carreras, Erwin Frise, Verena Kaynig, Mark Longair, Tobias Pietzsch, Stephan Preibisch, Curtis Rueden, Stephan Saalfeld, Benjamin Schmid, et al. Fiji: an open-source platform for biological-image analysis. *Nature methods*, 9(7):676–682, 2012.

[59] Uwe Schmidt, Martin Weigert, Coleman Broaddus, and Gene Myers. Cell detection with star-convex polygons. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 265–273. Springer, 2018.

[60] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018.

[61] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.

[62] Christoph Sommer, Christoph Straehle, Ullrich Koethe, and Fred A Hamprecht. Ilastik: Interactive learning and segmentation toolkit. In *2011 IEEE international symposium on biomedical imaging: From nano to macro*, pages 230–233. IEEE, 2011.

[63] Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a generalist algorithm for cellular segmentation. *Nature methods*, 18(1):100–106, 2021.

[64] Earl William Swokowski. *Calculus with analytic geometry*. Taylor & Francis, 1979.

[65] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[66] Lucas von Chamier, Romain F Laine, Johanna Jukkala, Christoph Spahn, Daniel Krentzel, Elias Nehme, Martina Lerche, Sara Hernández-Pérez, Pieta K Mattila, Eleni Karinou, et al. Democratising deep learning for microscopy with zerocostdl4mic. *Nature communications*, 12(1):1–18, 2021.

[67] Guotai Wang, Xinglong Liu, Chaoping Li, Zhiyong Xu, Jiugen Ruan, Haifeng Zhu, Tao Meng, Kang Li, Ning Huang, and Shaoting Zhang. A noise-robust framework for automatic segmentation of covid-19 pneumonia lesions from ct images. *IEEE Transactions on Medical Imaging*, 39(8):2653–2663, 2020.

[68] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021.

[69] Jiafan Zhuang, Zilei Wang, and Yuan Gao. Semi-supervised video semantic segmentation with inter-frame feature reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3263–3271, 2022.

[70] Yuliang Zou, Zizhao Zhang, Han Zhang, Chun-Liang Li, Xiao Bian, Jia-Bin Huang, and Tomas Pfister. Pseudoseg: Designing pseudo labels for semantic segmentation. *arXiv preprint arXiv:2010.09713*, 2020.