

MIPHEI-ViT: MULTIPLEX IMMUNOFLUORESCENCE PREDICTION FROM H&E IMAGES USING ViT FOUNDATION MODELS

A PREPRINT

 **Guillaume Balezo**

Sanofi

Paris, France

guillaume.balezo@minesparis.psl.eu

 **Roger Trullo***

InstaDeep*

Paris, France

r.trullo@instadeep.com

 **Albert Pla Planas**

Sanofi

Barcelona, Spain

albert.plaplanas@sanofi.com

 **Etienne Decencière**

Centre de Morphologie Mathématique

Mines Paris PSL

Fontainebleau, France

etienne.decenciere@minesparis.psl.eu

 **Thomas Walter**

Center for Computational Biology, Mines Paris, PSL,

Institut Curie, INSERM U1331

Paris, France

thomas.walter@minesparis.psl.eu

May 16, 2025

ABSTRACT

Histopathological analysis is a cornerstone of cancer diagnosis, with Hematoxylin and Eosin (H&E) staining routinely acquired for every patient to visualize cell morphology and tissue architecture. On the other hand, multiplex immunofluorescence (mIF) enables more precise cell type identification via proteomic markers, but has yet to achieve widespread clinical adoption due to cost and logistical constraints. To bridge this gap, we introduce MIPHEI (Multiplex Immunofluorescence Prediction from H&E), a U-Net-inspired architecture that integrates state-of-the-art ViT foundation models as encoders to predict mIF signals from H&E images. MIPHEI targets a comprehensive panel of markers spanning nuclear content, immune lineages (T cells, B cells, myeloid), epithelium, stroma, vasculature, and proliferation. We train our model using the publicly available ORION dataset of restained H&E and mIF images from colorectal cancer tissue, and validate it on two independent datasets. MIPHEI achieves accurate cell-type classification from H&E alone, with F1 scores of 0.88 for Pan-CK, 0.57 for CD3e, 0.56 for SMA, 0.36 for CD68, and 0.30 for CD20, substantially outperforming both a state-of-the-art baseline and a random classifier for most markers. Our results indicate that our model effectively captures the complex relationships between nuclear morphologies in their tissue context, as visible in H&E images and molecular markers defining specific cell types. MIPHEI offers a promising step toward enabling cell-type-aware analysis of large-scale H&E datasets, in view of uncovering relationships between spatial cellular organization and patient outcomes.

Keywords Computer Vision · Histopathology · Image Translation · Foundation model · In silico labelling

*Roger Trullo was a Sanofi employee at the time of the work

1 Introduction

The analysis of Hematoxylin and Eosin (H&E)-stained tissue slides is a cornerstone in the diagnosis of many pathologies, including cancer, providing insights into cell types, cellular phenotypes, tissue architecture and their alterations due to disease.

Multiplex immunofluorescence (mIF) Andreou et al. [2022], Tan et al. [2020], Im et al. [2019] imaging is a powerful technique that improves the analysis of tissue sections by providing detailed information beyond what H&E staining can reveal. mIF achieves this by simultaneously visualizing and quantifying multiple protein markers within a single tissue section, utilizing fluorescently labeled antibodies that bind to specific proteins, allowing for the identification of cell types based on marker expression. This capability makes mIF useful across various domains, including cancer biology, immunology, and infectious disease, where understanding spatial cell organization is critical. By combining molecular and spatial information at single-cell resolution, mIF supports detailed characterization of tissue microenvironments and cellular interactions.

The evolution of immunolabeling techniques from traditional immunohistochemistry (IHC) Duraiyan et al. [2012] has led to more advanced mIF techniques like PhenoCycler Black et al. [2021], which enables the detection of up to 100 markers through multiple imaging cycles, each capturing 4 channels. This iterative process can however degrade tissue integrity. Among recent advancements, the Orion scanner introduced by Lin et al. Lin et al. [2023a] allows the simultaneous detection of up to 20 markers in a single cycle, preserving tissue quality while still providing rich multiplexed information. Orion also allows capturing high-quality, restained H&E images from the same tissue section.

While mIF imaging offers several advantages, it also presents important challenges. Preparing and processing mIF slides is time-consuming and labor-intensive, requiring expensive reagents and specialized equipment, which are not always available in all clinical settings, limiting its accessibility. Compared to other techniques for spatial biology, such as Imaging Mass Cytometry (IMC) Mann et al. [2001] and VisiumHD spatial transcriptomics Oliveira et al. [2024], mIF detects fewer molecular markers but offers much higher spatial resolution. Like these technologies, mIF is constrained by high costs, which prevents it from being adopted in clinical practice.

In contrast, H&E slides are routinely generated in clinical practice. Since different cell types are characterized by distinct protein expression patterns, we hypothesize that certain markers can be predicted from cell morphology captured in H&E. With the availability of high-quality datasets obtained from technologies like the ORION scanner, which captures aligned H&E and mIF images with a high number of markers, we can now train AI models to infer mIF data from H&E images. This would allow us to identify those proteins that are predictable from morphological cues and perform these predictions on large retrospective cohorts, including clinical trial data, for which an acquisition of mIF data would not be feasible. Relating the predicted mIF data to outcome or treatment response can then contribute to biomarker discovery and hypothesis generation in oncology.

In this study, we aim to predict the expression of key markers from H&E, covering nuclear content (Hoechst), vasculature (CD31), immune populations (CD45, CD68, CD4, FOXP3, CD8a, CD45RO, CD20, PD-L1, CD3e, CD163), epithelial and stromal structures (E-cadherin, Pan-CK, SMA), and proliferation (Ki67). To achieve this, we introduce **MIPHEI**: Multiplex Immunofluorescence Prediction from H&E-stained Whole Slide Images, a U-Net-inspired architecture integrating state-of-the-art foundation models as encoders. Unlike traditional cell type classification models that rely on manual labels or pseudo-labels, our method is trained directly on aligned mIF data, avoiding biases from predefined annotations. It is also highly modular and does therefore not rely on specific nucleus segmentation and cell classification methods.

The key contributions of this study are as follows:

- **Prediction of Numerous Markers from H&E:** We identify the proteomic markers and cell types predictable from H&E images from a wide range of 15 markers.
- **Rigorous Validation and Reproducibility:** We demonstrate generalization of the method through validation with cell-level metrics for accurate cell type identification on both internal and two external datasets, moving beyond pixel-level evaluation. Our study sets a new benchmark framework for H&E to mIF translation.
- **New State-of-the-art:** Our model outperforms previous models and a random baseline on 15 markers, most of which are indiscernible to pathologists.
- **ViT Foundation Models in U-Net Architecture:** one major methodological novelty relies in the integration of Vision Transformer (ViT) Dosovitskiy et al. [2020] foundation models as the backbone in a U-Net architecture for histopathological image translation, leveraging advanced representation learning.

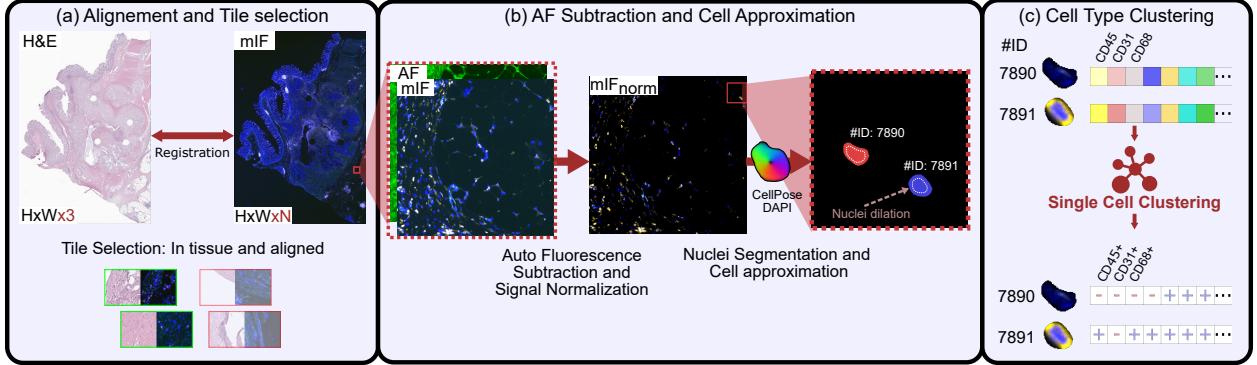


Figure 1: Preprocessing pipeline: (a) H&E and mIF images are aligned using Valis Gatenbee et al. [2023]. Tissue regions are then selected via Otsu thresholding on H&E, and a trained CNN filters misaligned tiles caused by the restaining and acquisition process. (b) Autofluorescence is subtracted from mIF images, followed by DAPI-based nuclei segmentation and nuclei dilation to approximate cell boundaries. (c) Pseudo-labels are generated by computing per-cell mean marker expression and applying GMM clustering to define marker positivity, determining labels (e.g., CD3e⁺).

2 Related Works

2.1 Virtual Staining

Image-to-image translation involves transforming images from one domain to another, with Pix2Pix Isola et al. [2017] being one of the most well-known methods in this domain for paired images. Utilizing conditional adversarial networks, these techniques enable high-quality image synthesis, supporting tasks such as style transfer, super-resolution, and domain adaptation, and have driven significant advancements across various applications, as demonstrated by Wang et al. Wang et al. [2018].

In recent years, image-to-image translation techniques have become increasingly popular for life science applications, such as predicting super-resolved microscopy images Wang et al. [2019] or bright-field-like images from holographic images Wu et al. [2019]. In histopathology, such virtual staining techniques have been used to predict immunohistochemistry (IHC) data Sun et al. [2023], DoanNgan et al. [2022], to do virtual multiplexing Wu et al. [2023], Bao et al. [2021], to predict mIF images from IHC images Ghahremani et al. [2022], or both H&E Rivenson et al. [2019], Bai et al. [2023], Zhang et al. [2020], Cao et al. [2023] and mIF Christiansen et al. [2018] from unlabeled autofluorescence images.

In the scope of this study, the most relevant work involves translating H&E images to mIF, which directly addresses the challenge of performing cell type calling from standard histological stains. A notable example is SHIFT Burlingame et al. [2020] based on conditional generative adversarial network (cGAN) to predict markers like DAPI, pan-cytokeratin, and α -smooth muscle actin. Another relevant work is ImmunoAIzer Bian et al. [2021] using a semi-supervised adversarial approach to predict CD3, CD20, PanCK, and DAPI channels from H&E, trained on both aligned and unaligned data. More recently, the same author published HEMIT Bian et al. [2024], releasing a dataset and a method designed for translating H&E to mIHC images targeting DAPI, CD3, and pan-cytokeratin (Pan-CK) markers. The authors propose a hybrid ViT-CNN generator architecture combining CNNs with Swin Transformers Liu et al. [2021].

2.2 ViT on Dense Prediction Tasks

ViTs Dosovitskiy et al. [2020] are now widely used in image classification and often outperform CNNs when pretrained with advanced contrastive self-supervised learning (SSL) methods such as MoCov3 Chen et al. [2021], iBOT Zhou et al. [2021], and DINOV2 Oquab et al. [2023]. These SSL approaches have been adapted to histopathology, enabling recent foundation models to achieve strong and robust performance across diverse downstream tasks. CTransPath Wang et al. [2022] uses MoCov3 Chen et al. [2021] with a Swin Transformer architecture Liu et al. [2021], and was trained on 15.6 million tiles from opensource datasets. UNIV2 Chen et al. [2024] leverages DINOV2 to train a ViT-H/14 model on 200 million tiles from 3.5k proprietary H&E and IHC slides. H-optimus-0 Saillard et al. [2024] is a ViT-G/14 variant with SSL pretraining based on iBOT Zhou et al. [2021] and DINOV2 Oquab et al. [2023], trained on hundreds of millions of tiles from 500,000+ H&E whole slide images. These models are among the state-of-the-art in computational pathology, with DINOV2-based models like UNIV2 and H-optimus-0 outperforming ImageNet-pretrained encoders and CTransPath on unseen datasets.

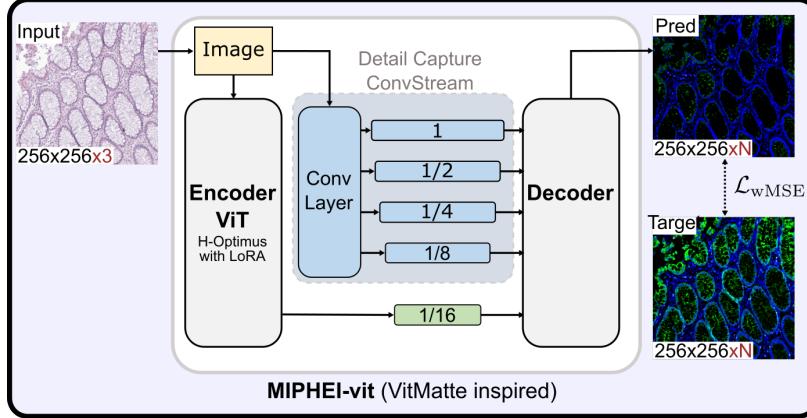


Figure 2: **MIPHEI architecture:** A U-Net-inspired model, based on VitMatte, is trained to predict mIF images from H&E, using the H-optimus-0 ViT foundation model, Tanh activation, and a custom weighted MSE loss to correct for the unbalanced distribution across markers.

While plain ViTs are effective for image-level tasks, they often struggle with dense prediction tasks like image translation due to limited ability to capture fine local details, unlike CNNs which excel at leveraging local continuity and multiscale features. To overcome this, hybrid CNN-ViT architectures have been proposed. CellViT Hörst et al. [2024], for example, uses the UNETR architecture Hatamizadeh et al. [2022], which integrates a ViT encoder into a U-Net-like design, employing convolutional transpose blocks to create hierarchical features, offering a simple adaptation. ViTMatte Yao et al. [2024] combines a plain ViT with a convolutional module for pyramidal feature extraction and detail refinement, using ViT token features only as the bottleneck. More advanced models like ViT-Adapter Chen et al. [2022] and ViT-CoMer Xia et al. [2024] go further by enabling richer multi-scale interactions between ViT and CNN components, but at higher computational cost. While these hybrid methods have mainly been applied to segmentation, we hypothesized they could also benefit dense image translation tasks like mIF prediction, especially when using foundation models trained on large, diverse datasets.

3 Dataset

We present the datasets used in this study, with Figure 3 providing additional details on tissue distribution, domain shifts, and dataset composition. The preprocessed data for ORION and HEMIT datasets is available on Zenodo.

3.1 Orion CRC Dataset (Train/Val/Test)

The ORION colorectal cancer (CRC) dataset Lin et al. [2023b,a] was acquired with a novel system capturing 18-channel immunofluorescence (IF) images alongside H&E staining on the same tissue samples. This dataset includes 41 Whole Slides Images (WSIs) with both H&E and mIF data. The 15 markers (with additional DNA stain) of interest for this study are: Hoechst, CD31, CD45, CD68, CD4, FOXP3, CD8a, CD45RO, CD20, PD-L1, CD3e, CD163, E-cadherin, Pan-CK, SMA and Ki67. The PD-1 channel was not used due to poor signal quality. The H&E images were acquired with an Aperio GT450 microscope (Leica Biosystems) and registered to the IF images in the original study, at a resolution of 0.325 microns per pixel (mpp).

3.2 HEMIT Dataset (Test)

The HEMIT dataset Bian et al. [2024] consists of H&E and multiplex-immunohistochemistry (mIHC) images the same tissue sections, acquired using Mantra system scanner (PerkinElmer, Waltham, MA, USA) with DAPI, CD3, and Pan-CK markers. It includes 5,292 paired 1024x1024-pixel patches aligned at pixel level at 40x magnification.

3.3 IMMUcan CRC Dataset (Test)

IMMUcan (Integrated iMMUnoprofiling of large adaptive CANcer patient cohorts) Hong et al. [2020] is a European initiative launched in 2019 to advance Tumor Micro Environment (TME) profiling. We use only the CRC cases, comprising 35 registered H&E and mIF slides from two cohorts, which include both advanced and less severe stages.

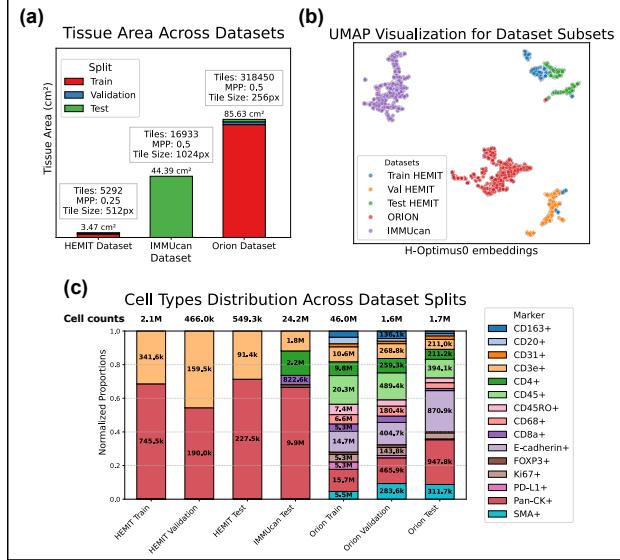


Figure 3: Dataset Overview: (a) Distribution of tissue samples and tiles across datasets and splits. (b) 2D UMAP visualization of H-Optimized-0 embeddings, highlighting domain shifts across our datasets. (c) Normalized cell type distribution across datasets and splits.

H&E slides were scanned with a Hamamatsu NanoZoomer 2.0-HT (0.5 mpp), and mIF images were acquired on a PerkinElmer Vectra Polaris. The mIF panel includes DAPI, CD3, CD8, CD4, FOXP3, and Pan-CK. Unlike HEMIT and ORION, the H&E and mIF slides in IMMUCAN come from consecutive sections, requiring a dedicated evaluation protocol detailed later in Metric Overview. Although currently private, the dataset is expected to become publicly available.

3.4 Data Preprocessing

We designed a preprocessing pipeline tailored to our H&E-to-mIF prediction task, addressing WSI registration, artifact removal, autofluorescence subtraction with normalization, and pseudo-label extraction, as illustrated in Figure 1.

3.4.1 H&E to mIF registration on consecutive cuts

Accurate registration between H&E and mIF images is crucial for training and evaluating our models. For Orion and HEMIT datasets, data was already registered Lin et al. [2023b], Bian et al. [2024]. For the IMMUCAN dataset, we used Valis Gatenbee et al. [2023] for registration of consecutive H&E and mIF slides.

3.4.2 Tile Selection

To reduce artifacts from restaining and acquisition in the ORION dataset, we filtered aligned H&E and mIF tiles using several quality control steps. We used thresholding on an empty mIF channel—without antibody staining but capturing shared noise—to identify artifacts affecting all mIF channels. Poor H&E quality tiles were identified by clustering H-optimus-0 embeddings and discarding clusters associated to obvious artifacts. We also manually annotated misaligned tiles and trained a CNN to detect them automatically from H&E and DAPI images. In total, about 40k tiles (10%) were excluded.

For the IMMUCAN dataset, where H&E and mIF are from consecutive sections, we selected well-aligned tiles to ensure reliable correlation analysis based on three criteria: Pearson correlation of aligned nuclei density maps from 32×32 patches (threshold=0.25 per tile), tissue overlap ($IoU > 0.5$), and tissue percentage ($> 40\%$). We extracted 1024 × 1024 pixel tiles at 20x (0.26 mm²), retaining 17k tiles covering 44.4 cm².

3.4.3 Autofluorescence Subtraction & Data Normalization

Autofluorescence (AF), captured as an independent channel I_{AF} , refers to light naturally emitted by the tissue across channels, introducing noise in other markers. We subtract the AF from each marker channel I_{IF}^c , using λ^c and b^c which were manually adjusted using a Napari Sofroniew et al. [2025] tool we developed:

$$I_{cor}^c = \max(0, I_{IF}^c - \lambda^c \cdot I_{AF} + b^c)$$

Next, each channel is normalized using the 99.9th percentile $q_{0.999}^c$, computed per marker from the distribution of foreground pixel intensities across the training set, and then log-transformed:

$$I_{norm}^c = 255 \cdot \log \left(\frac{\min(I_{cor}^c, q_{0.999}^c)}{q_{0.999}^c} + 1 \right)$$

This logarithmic transformation compresses high-intensity values and reduces the impact of extreme outliers, while normalization ensures a consistent dynamic range across markers. The autofluorescence Napari tool, selected parameters, code, model, and data access instructions are available on our GitHub, within the Sanofi Public organization, under specific license conditions including a limitation to non-commercial use only.

3.4.4 Single-Cell Pseudo-Label Extraction

To establish a ground truth for evaluating our model’s ability to identify marker-positive cells from H&E images, we used a standard cell type calling approach on mIF images Lin et al. [2023a]. We first segmented nuclei from the DAPI channel using Cellpose Stringer et al. [2021] fine-tuned on our data. As a proxy for cell regions, we dilated the nuclear regions by $2 \mu\text{m}$. Single-cell analysis was performed by extracting mean fluorescence intensities, followed by unsupervised gating using a Gaussian Mixture Model (GMM) to distinguish positive from negative cells Zhang et al. [2022], with posterior probabilities estimated from the GMM. Although effective, this approach is sensitive to artifacts, approximate boundaries, and signal spillover. To improve robustness, we implemented a hierarchical gating strategy making sure that the biological marker hierarchy was preserved. For instance, given that CD3 positive cells are also CD45 positive, we kept CD3 positivity only for CD45 positive cells. By applying these rules, we obtained a high confidence annotation of our cells.

For the IMMUCan dataset, we used single-cell data provided by the consortium, which was generated using in-house clustering and nucleus segmentation. Since the same nuclei are not necessarily present in consecutive H&E and mIF sections, we applied Hoverfast Liakopoulos et al. [2024] for nucleus segmentation from H&E images, allowing us to perform the cell-level correlation analysis explained in section 5.2.

4 Methodology

This section outlines our overall approach for predicting mIF images from H&E images. We used the Orion dataset for training, and HEMIT and IMMUCan datasets for testing.

4.1 Model Architecture

Our model utilizes a U-Net generator to perform the image translation task from H&E to mIF signals (Figure 2). The architecture supports various encoder types, including CNNs like ConvNeXt v2 Woo et al. [2023], Swin Transformer Liu et al. [2021], and plain ViTs Dosovitskiy et al. [2020], enabling integration of recent foundation models in histology, such as CTransPath, Univ2, and H-optimus-0 Wang et al. [2022], Chen et al. [2024], Saillard et al. [2024]. CNN-based encoders naturally align with the classical U-Net design, providing pyramidal features. For ViTs, which maintain fixed feature dimensions across layers, we explore two alternatives to generate multi-scale features: (1) convolutional transpose for upsampling token features, as used in UNETR Hatamizadeh et al. [2022] and CellViT Hörist et al. [2024], and (2) a hybrid approach inspired by ViTMatte Yao et al. [2024], where a convolutional stream extracts pyramidal features while ViT features serve as the bottleneck. Unlike original ViTMatte, which concatenates a trimap to the input image, we use only the H&E RGB image and omit convolutional necks and window attention. For simplicity, we still refer to this variant as VitMatte.

Building on these encoded features, the decoder reconstructs outputs using nearest interpolation for upsampling, combined with skip connections, dual 3×3 convolutional layers, batch normalization, and ReLU activation. Multiple output heads predict various mIF signals from the same decoder outputs, with a final Tanh activation.

4.2 Loss Function

We train MIPHEI using a weighted Mean Square Error (MSE) loss to ensure accurate translation from H&E to mIF signals while accounting for varying signal intensities and prevalence across markers. Each marker’s loss is scaled by the inverse of its standard deviation to balance contributions Mai et al. [2021].

Let $L_{\text{MSE},j}$ denote the MSE loss for the j^{th} marker, σ_j denote the standard deviation, M be the total number of markers, and λ the weight of the reconstruction loss. The weighted MSE loss is then defined as:

$$\mathcal{L}_{\text{wMSE}} = \frac{\lambda}{M} \sum_j \frac{1}{\sigma_j^2} L_{\text{MSE},j}$$

4.3 Foundation Model Training Strategies

While CNN-based encoders within the U-Net architecture are fully fine-tuned in all experiments, we explore two efficient training strategies for large ViT-based foundation model encoders: (1) decoder-only training with a frozen foundation encoder, leveraging robust pretrained features while reducing trainable parameters; (2) Low-Rank Adaptation (LoRA) Hu et al. [2022] with $\text{rank} = 8$ and $\alpha = 1$, which adapts a pretrained ViT encoder by adding trainable low-rank matrices to the Query and Value projections of the attention layers.

5 Experiments

Here we present the experimental workflow that we set up to assess the performance and robustness of our image translation model from H&E to mIF.

5.1 Training Setup

5.1.1 Data Configuration

We split the Orion dataset by slide into training (37), validation (2), and test (2) sets. For HEMIT, we used its original train-validation-test splits, while for IMMUCan, all tiles were treated as test set.

All models are trained on H&E-mIF images from ORION extracted as 256x256 pixel tiles at 0.5 mpp. For normalization, target data is scaled to the range $[-0.9, 0.9]$ to prevent saturation at extreme values in Tanh activation. Input H&E RGB data is normalized using the mean and standard deviation of the foundation model employed.

To enhance generalization, augmentations include spatial transformations, such as horizontal and vertical flips and coarse dropout (zeroing out a random box in both input and target), along with color augmentations like stain augmentation, random brightness and contrast adjustments, gaussian blur, and gaussian noise. These augmentations simulate variations in stain intensity, color distributions, and common imaging artifacts. To improve robustness, we trained a CycleGAN Zhu et al. [2017] following Runz et al. Runz et al. [2021] to translate Orion H&E images toward the IMMUCan style, capturing more complex domain shift, and used both original and precomputed translated tiles as augmentation during training.

5.1.2 Training Configuration

We adopt similar Pix2Pix hyperparameters, given its role as a classical image translation model. We use the Adam optimizer with a learning rate of 2×10^{-4} , kept constant in the first half of training before linearly decaying to zero. A linear warmup phase is applied to the generator’s learning rate for the first 400 iterations to improve stability and adaptation, as recommended for ViTs by Dosovitskiy et al. Dosovitskiy et al. [2020]. On top of data augmentations, regularization includes weight decay (1×10^{-5}), gradient clipping (max norm 1), and dropout (0.1). Non-pretrained weights are initialized using a normal distribution with a gain of 0.02, and biases are set to zero. All training are conducted with a batch size of 16. All experiments were performed on a single A100 GPU.

5.2 Metric Overview

To evaluate our image translation model from H&E to mIF, we used both standard pixel-level and cell-level metrics. At the pixel level, we employed PSNR and SSIM. One of the most important aspect of our model is to capture cell-type information, which is not adequately measured by these metrics. We therefore designed a cell-level fidelity metric.

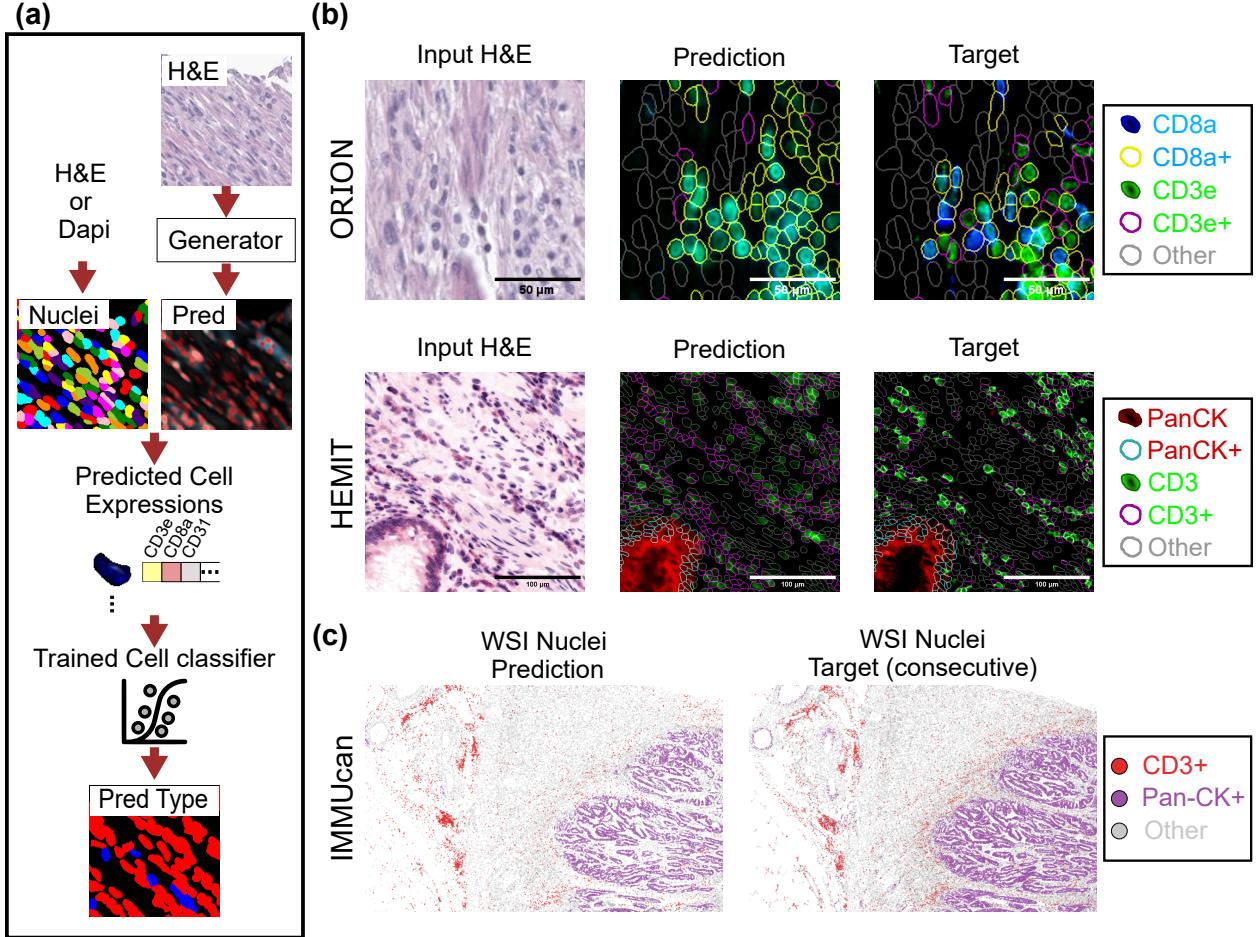


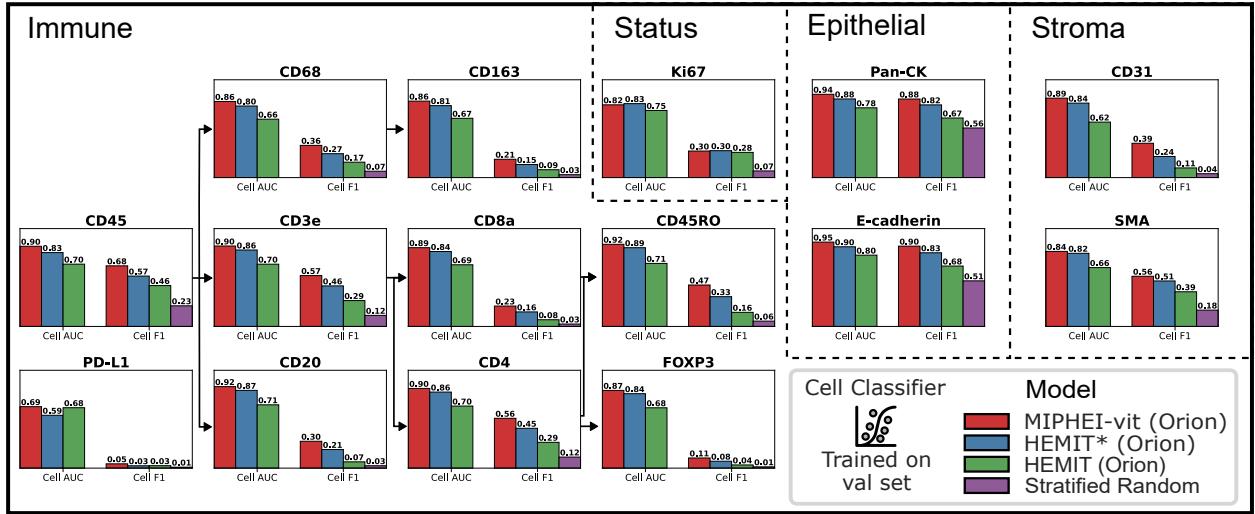
Figure 4: Prediction pipeline and prediction visualization: (a) **Inference pipeline:** mIF images are first generated using a trained U-Net model. Predicted single-cell data are then extracted by averaging predicted mIF signals within each nucleus, using nuclei masks from an external segmentation model. Finally a cell classifier, trained on validation set cells, predicts cell types. (b) **Prediction examples from our best model:** Predicted mIF images and cell types (shown as colored cell boundaries) are compared to target mIF images and annotated cell types from the same stained tissue section in the Orion (CD3e, CD8a) and HEMIT (Pan-CK, CD3) datasets. (c) **IMMUcan large-area visualization:** Nuclei predictions on H&E alongside clustered nuclei from the corresponding consecutive mIF sections.

For this, we extract predicted single-cell data by averaging predicted marker expression within each cell, and classify them using a logistic regression trained on validation set cells. The classifier is then applied to the test set and predictions are compared to pseudo-labels as described in subsection 3.4.4. We use a classifier instead of direct gating on predicted mIF channels to improve robustness to individual channel errors and account for calibration differences between predicted and target fluorescence intensities.

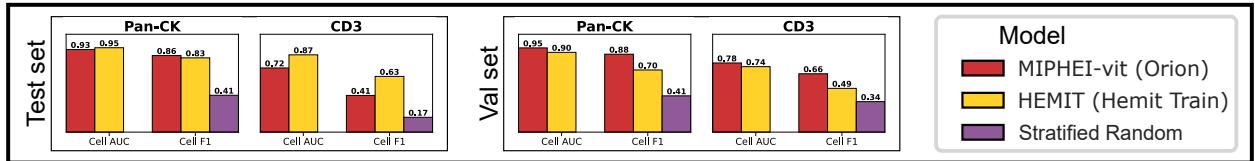
We assess the discriminative power of our model’s predictions by computing Cell AUC (Receiver Operating Characteristic Area Under the Curve), using the predicted mean intensity per nucleus as a score compared to pseudo-labels. To account for class imbalance and better assess model performance, we report F1 score using the cell classifier predictions, as most markers have a lower frequency of positive cells.

We further analyze our model by computing cell count correlations on IMMUcan for CD3, CD8, CD4, FOXP3, and Pan-CK. Since consecutive tissue sections in the IMMUcan dataset are biologically similar but not identical, direct pixel-level alignment is not feasible. However, the minimal spatial separation between sections preserves overall tissue architecture and cellular distribution, allowing for meaningful region-level comparisons. We compute the correlation between predicted and target positive cell counts across registered tiles, using predictions from the logistic regression-based single-cell classification. Following preprocessing in D.1, we compute Pearson correlation coefficients between

(a) Orion Test Set Evaluation (Same sections)



(b) HEMIT Dataset Evaluation (Same sections)



(c) IMMUcan Evaluation (Consecutive sections)

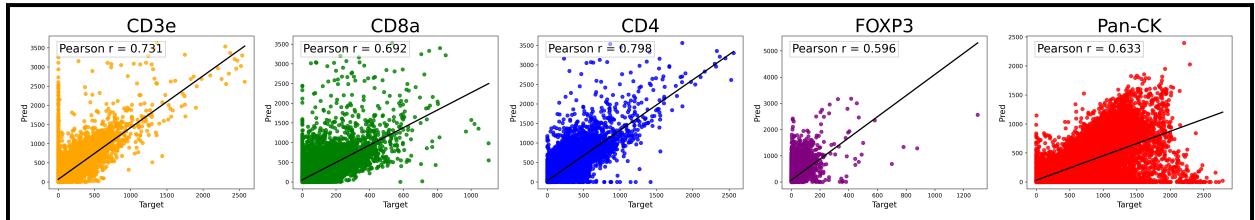


Figure 5: Evaluation analysis: (a) **Orion test set:** Performance comparison across all Orion markers between MIPHEI, HEMIT*, HEMIT (all trained on Orion train set) and a stratified random model predicting classes based on cell type proportions. Markers are grouped by general functions, with hierarchical relationships indicated by arrows. Cell classification model is a logistic regression model trained on Orion validation cells. (b) **HEMIT dataset:** Comparison between MIPHEI (trained on Orion), the HEMIT model (trained on HEMIT train set) and a stratified random model. Cell classification model is a logistic regression: trained on 5% of training cells for MIPHEI, and all available training cells for the HEMIT model. (c) **IMMUcan** (consecutive sections): Pearson correlation and regression plots (with linear fits) between predicted and pseudo-labeled cell type counts from 17k tiles using MIPHEI.

tile-level cell counts (0.26 mm^2 area per tile) for each marker. High correlations, even if imperfect, indicate effective cell identification.

6 Results

This section presents an ablation study to identify critical model components, comparisons with the state-of-the-art HEMIT method, and a detailed analysis of marker-specific performance. We selected HEMIT for benchmarking as it addresses a similar problem, provides public code and data for direct comparison, and outperforms baselines like Pix2Pix used in other related work such as SHIFT.

6.1 Ablation studies

To identify critical components and configurations significantly impacting performance, we trained all models for 15 epochs ($\sim 276k$ iterations) with batch size 16 by default.

Table 1: Ablation Study: Performance Metrics Across Model Configurations on ORION test set

Configuration	PSNR	SSIM	Cell AUC	Cell F1
Impact of GAN Discriminator (UNETR H-optimus-0 LoRA)				
Generator only	27.86	0.840	0.868	0.431
Pix2Pix (GAN)	27.00	0.830	0.817	0.410
Impact of Foundation Model Encoder (UNETR with LoRA)				
CTransPath	26.96	0.837	0.812	0.351
Univ2	27.73	0.838	0.862	0.424
H-optimus-0 (selected)	27.86	0.840	0.868	0.431
Encoder Finetuning Method (UNETR H-optimus-0)				
Frozen encoder	27.44	0.836	0.857	0.413
LoRA (selected)	27.86	0.840	0.868	0.431
Impact of Architecture				
U-Net-ResNet50	26.54	0.832	0.812	0.344
U-Net-ConvNeXtv2 Large	27.40	0.839	0.840	0.379
HEMIT*	27.08	0.837	0.831	0.360
UNETR	27.86	0.840	0.868	0.431
MIPHEI (Best)	27.78	0.837	0.876	0.438

6.1.1 Impact of Discriminator

We evaluated the impact of adding a discriminator in a Pix2Pix-like setup, using UNETR as the generator and the standard patch-based discriminator from Pix2Pix. We found that the discriminator slightly reduces pixel-level metrics, and consistently cell-type classification accuracy. We hypothesize the model favors realism over fidelity (Table 1), thus leading to more realistic looking images, but not providing more accurate predictions. For this reason, we excluded the discriminator from further experiments.

6.1.2 Model Architecture

We next compared four U-Net-inspired architectures for our mIF prediction model: two convolution-based encoders pre-trained on ImageNet—ResNet50 (37M parameters) and ConvNeXt v2 Large (203M parameters)—and two U-Net variants that both use the H-optimus-0 foundation model as encoder but differ in their ViT integration strategies—UNETR (1.17B parameters, 36M trainable via LoRA) and ViTMatte (1.14B parameters, 6.6M trainable via LoRA). U-Net with a ResNet50 encoder showed lower performance, likely due to its smaller size and lack of domain-specific pretraining. In contrast, U-Net-ConvNeXt achieved strong results, benefiting from greater model capacity and attention-like mechanisms. Transformer-based variants using H-optimus-0 encoder significantly outperformed convolutional ones, emphasizing the advantage of employing pretrained foundation models for the encoder component in our task. UNETR and ViTMatte achieved similar performance, but ViTMatte converged much faster, reaching UNETR’s 10-epoch performance after the first one only. This is likely due to its convolutional ‘detail capture’ stream, which provides intermediate skip features that complement ViT representations and aid spatial reconstruction more effectively than UNETR’s token-based upsampling. Overall, ViTMatte offers the best trade-off between performance and training efficiency.

6.1.3 Foundation Model Benchmark

We evaluated three foundation-model encoders available for our image translation task (UNETR with LoRA): CTransPath, Univ2, and H-optimus-0. CTransPath was significantly outperformed by Univ2 and H-optimus-0, which provided similar results with a slight advantage for H-optimus-0, which we therefore selected.

6.1.4 Finetuning Strategies

Finally, we evaluated the finetuning strategies described in subsection 4.3 using UNETR with the H-optimus-0 encoder. Our results show that LoRA improves performance over a frozen encoder, achieving a higher test cell AUC (0.431 vs. 0.413). Based on this, we adopt LoRA as our fine-tuning strategy.

Based on our ablation study, the best-performing configuration for MIPHEI is ViTMatte with the H-optimus-0 encoder, adapted using LoRA, and trained without a GAN strategy, providing an optimal balance between performance, efficiency, and convergence speed.

6.2 Comparison with State of the Art (HEMIT model)

To benchmark our proposed model, MIPHEI, against the state-of-the-art HEMIT, we evaluate cross-dataset generalization using five models: (1) **MIPHEI**, trained on ORION, (2) the official **HEMIT** checkpoint trained on the HEMIT dataset, (3) **HEMIT*** trained on ORION using the HEMIT-architecture, but with the MIPHEI training scheme (weighted MSE loss, no discriminator), (4) **HEMIT-Orion**, trained on ORION using the original HEMIT pipeline, and (5) a **stratified random model** as a baseline, assigning cell type probabilities based on their distribution in the test datasets, with metrics averaged over 100 runs to account for variability in random sampling. Since most cell subtypes are not identifiable from H&E alone by a pathologist, the task is challenging, and the random baseline provides a reference to show that our model outperforms chance-level expectations. Evaluation is conducted on the ORION test set (in-domain for models trained on ORION), on the HEMIT validation and test sets and the IMMUCan dataset, both serving as external out-of-domain benchmarks.

6.2.1 Evaluation on ORION test set

Table 2: Overall Cell-Level Performances on Test Sets

Cell Metric	MIPHEI (Orion)	HEMIT (HEMIT)	HEMIT* (Orion)	HEMIT (Orion)	Random
ORION Dataset (Test Set)					
AUC	0.876	-	0.831	0.701	-
F1	0.438	-	0.360	0.253	0.140
HEMIT Dataset (Average Validation & Test Sets)					
AUC	0.844	0.863	0.764	0.598	-
F1	0.701	0.663	0.471	0.481	0.333
IMMUCan Dataset					
Pearson	0.690	-	0.667	0.422	-

The results of performance evaluation of all tested models is shown in Table 2, Figure 5). HEMIT-Orion, trained with the original HEMIT pipeline, outperformed the random model, but showed the lowest performance across all markers except PD-L1, for which all methods performed poorly. We hypothesize that this was due to the use of the L1 loss, which may be less effective for markers with low intensity or high background proportion, as suggested by the stronger results of HEMIT* trained with our setup. On the other hand, MIPHEI achieved the best overall performance, surpassing both HEMIT* and HEMIT-Orion across nearly all markers. The only exception was Ki67, where HEMIT* performed slightly better, achieving a Cell F1 score 0.3% higher than MIPHEI.

These results confirm that MIPHEI provides the most robust predictions within its training domain. The improved performance of HEMIT* over HEMIT-Orion further demonstrates the effectiveness of our training pipeline, even when applied to an alternative architecture.

6.2.2 Evaluation on HEMIT dataset

The HEMIT dataset serves as an external test set for assessing the robustness of our approach. We compare MIPHEI (train: ORION), with the HEMIT model (train: HEMIT, released checkpoint). As Figure 3.B suggests, the HEMIT validation and test sets come from different domains, but both are contained within the training domain, with the test domain being more strongly represented. Hence, we report cell-level performances separately for each. We use logistic regression as the cell classifier, trained on 5% of the training cells for MIPHEI and all training cells for HEMIT. As MIPHEI predicts 15 ORION markers while HEMIT predicts only Pan-CK and CD3, each model’s classifier is trained on the full set of markers predicted by its respective model.

On the HEMIT test set, MIPHEI slightly outperforms HEMIT by +3% F1 score for Pan-CK (Figure 5), but HEMIT significantly outperforms MIPHEI on CD3, with a +22% Cell F1 score. However, on the HEMIT validation set, MIPHEI surpasses HEMIT for both markers, achieving +18% F1 for Pan-CK and +17% for CD3. This is a strong result, as MIPHEI was evaluated on an external dataset with only minimal adaptation of the auxiliary cell classification stage, using just 5% of the training data. In contrast, the HEMIT model was trained on data that included images from both the validation and test domains of the HEMIT dataset (Figure 3.C). MIPHEI demonstrates stronger generalization, likely due to its foundation model encoder, enabling better adaptation to unseen domains.

6.2.3 Evaluation on IMMUCAN dataset

We further evaluate model generalization on the IMMUCAN dataset by performing a correlation analysis between predicted cell type counts from H&E and pseudo-label counts from consecutive mIF sections. Pearson correlation is computed over 0.26 mm^2 regions for MIPHEI, HEMIT*, and HEMIT. Each model uses its own logistic regression cell classifier, trained on ORION validation cells for MIPHEI and HEMIT*, and on all HEMIT training cells for HEMIT.

Across all markers, MIPHEI achieves the highest Pearson correlation, outperforming both HEMIT* and HEMIT, further confirming its superiority in cross-domain generalization (Figure 5).

6.3 Marker-Level Analysis

6.3.1 In-domain performance analysis

We have demonstrated that MIPHEI consistently outperforms random predictions for all markers, confirming its ability to capture meaningful morphological cues from H&E to estimate protein expression. However, the predictability varies across markers and cell types (Figure 5.A).

Markers with the highest performance include epithelial markers E-cadherin (F1 0.903) and Pan-CK (F1 0.884), which label epithelial cells forming well-defined clusters and glandular structures in H&E. The immune marker CD45 (F1 0.681) also performs well. While it mainly identifies lymphocytes, which are easily recognizable in H&E, it also includes cells like monocytes, which are harder to spot, making prediction a bit more challenging.

Markers with moderate performance included SMA (F1 0.564), which labels smooth muscle cells, and CD31 (F1 0.386), which marks endothelial cells. Their more subtle morphological features may have contributed to lower accuracy: smooth muscle cells have a spindle-shaped appearance but can be confused with fibroblasts, while endothelial cells are typically sparse and located within vessel structures, making them harder to distinguish. Immune subtype markers CD3e, CD45RO, CD4, CD8a, and CD20 (F1 0.229–0.572) showed moderate performance, with broader T-cell markers like CD3e achieving better performance than more specific ones like CD8a, which are harder to identify. While lymphocytes are visible in H&E, their subtypes remain indistinguishable to pathologists, highlighting our model’s value on this difficult task. Macrophage markers CD68 (F1 0.362) and CD163 (F1 0.206) faced similar challenges, as their heterogeneity complicates identification, with CD163 being a macrophage subtype of CD68, further complicating prediction. Markers with the lowest performance are FOXP3 (F1 0.114) and PD-L1 (F1 0.048), both challenging for different reasons. FOXP3 marks rare regulatory T-cells, a highly specific CD4+ subtype, making it one of the most difficult immune markers to predict. PD-L1, a functional marker, lacks a clear association with a single cell type and shows irregular expression across various cells, making probably prediction from H&E nearly impossible.

In summary, structural markers and broad immune markers seem predictable with high accuracy. Markers defining specific immune subtypes are harder to predict, and some functional markers are not predictable with an accuracy high enough to be applicable in practice.

6.3.2 Out-of-domain performance analysis

Our external validation on HEMIT and IMMUCAN confirms that MIPHEI, trained on ORION, generalizes well to datasets with domain shifts in mIF technology and H&E appearance, requiring minimal adaptation.

On the HEMIT dataset, analysis is limited by the small number of available markers, but our model performs well with only minimal adaptation of the cell classifier using 5% of the training cells.

On the IMMUCAN dataset, we observe strong correlation for CD4 (0.80) CD3e (0.73) and CD8 (0.69), and moderate correlation for FOXP3 (0.60) and Pan-CK (0.63). We also observe overdetection for CD3e (i.e. cells with CD3e prediction and measured 0 expression) and underdetection of Pan-CK (cells with predicted 0-expression and positive according to the measurement). Manual inspection showed that this was mainly due to low-quality tiles with high auto-fluorescence.

7 Discussion

In this study, we present MIPHEI, a method trained to predict 16 mIF channels from standard H&E slides. For this, we proposed a U-Net with a ViT-based foundation model as encoder and demonstrated that this architecture outperforms previously proposed methods. Moreover, we showed that MIPHEI generalizes well across datasets. We attribute this to the robust and transferable encodings learned by the foundation model, which was exposed to millions of tiles during pretraining. Our findings highlight the value of leveraging foundation models for mIF prediction, and potentially for other image translation tasks in histology.

We also introduced a validation strategy focused on single-cell metrics, which we believe are the most relevant for this task. Pixel-wise accuracy may be unreliable, as H&E images are usually not informative about cytoplasmic boundaries. Instead, the biologically meaningful information lies in protein expression levels within individual cells, protein positivity, or the resulting cell type. We reflect this in our evaluation framework, which we provide as part of this study.

Accurate, domain-robust mIF prediction opens the door to a range of applications. While direct clinical deployment for diagnostics may remain challenging, MIPHEI proves to be a powerful tool for mining large retrospective cohorts—enabling the identification of cell types without relying on manual annotation. Ultimately, this can allow to identify associations between specific cell populations, their spatial arrangements, and clinically relevant outcomes such as survival or treatment response. MIPHEI thus holds promise for hypothesis generation and exploratory analysis. Furthermore, extending this approach to predict clinically relevant scores, such as the Immunoscore, represents a compelling direction for future work.

Our study is not free of limitations. Although the number of training tiles was substantial, they were derived from only 41 slides. Access to larger and more diverse datasets would likely improve the model’s robustness to domain shifts and biological variability. Additionally, we observed that the cell classification model still requires fine-tuning on a small set of labeled cells when applied to out-of-domain data. As such, if a domain shift is anticipated, some mIF data will still be needed for calibration.

8 Conclusion

In this paper, we present MIPHEI, a deep learning framework that predicts mIF images from H&E using ViT foundation models as encoders within a U-Net architecture. MIPHEI outperforms state-of-the-art models on both internal and external datasets, demonstrating strong generalization across staining protocols and imaging conditions. We evaluated MIPHEI across 15 protein markers and associated cell types. It achieves high accuracy for epithelial (Pan-CK, E-cadherin) and broad immune markers (CD45, CD3e), performs moderately on more specific immune subtypes (CD8a, CD4, CD45RO, CD68, CD163) and stromal markers (CD31, α -SMA), and struggles with FOXP3 and PD-L1 due to their complexity and small number of positive cells.

Acknowledgment

This work uses IMMUcan data funded by IMI2 JU (grant 821558) with support from Horizon 2020 and EFPIA.

Funding

This work was sponsored by ANRT and Sanofi. Furthermore, T. Walter acknowledges funding from Agence Nationale de la Recherche under the France 2030 program, with the reference number ANR-24-EXCI-0004

Conflict of Interest

G. Balezo, R. Trullo, A. Pla Planas are/or were Sanofi employees and may hold shares and/or stock options in the company. E. Decencière and T. Walter have nothing to disclose.

References

- Chrysafis Andreou, Ralph Weissleder, and Moritz F Kircher. Multiplexed imaging in oncology. *Nature Biomedical Engineering*, 6(5):527–540, 2022.
- Wei Chang Colin Tan, Sanjna Nilesh Nerurkar, Hai Yun Cai, Harry Ho Man Ng, Duoduo Wu, Yu Ting Felicia Wee, Jeffrey Chun Tatt Lim, Joe Yeong, and Tony Kiat Hon Lim. Overview of multiplex immunohistochemistry/immunofluorescence techniques in the era of cancer immunotherapy. *Cancer Communications*, 40(4):135–153, 2020.
- Kyuseok Im, Sergey Mareninov, M Fernando Palma Diaz, and William H Yong. *An introduction to performing immunofluorescence staining*. Springer, 2019.
- Jeyapradha Duraiyan, Rajeshwar Govindarajan, Karunakaran Kaliyappan, and Murugesan Palanisamy. Applications of immunohistochemistry. *Journal of Pharmacy and Bioallied Sciences*, 4(Suppl 2):S307–S309, 2012.
- Sarah Black, Darci Phillips, John W Hickey, Julia Kennedy-Darling, Vishal G Venkataraaman, Nikolay Samusik, Yury Goltsev, Christian M Schürch, and Garry P Nolan. Codex multiplexed tissue imaging with dna-conjugated antibodies. *Nature protocols*, 16(8):3802–3835, 2021.
- Jia-Ren Lin, Yu-An Chen, Daniel Campton, Jeremy Cooper, Shannon Coy, Clarence Yapp, Juliann B Tefft, Erin McCarty, Keith L Ligon, Scott J Rodig, et al. High-plex immunofluorescence imaging and traditional histology of the same tissue section for discovering image-based biomarkers. *Nature cancer*, 4(7):1036–1052, 2023a.
- Matthias Mann, Ronald C Hendrickson, and Akhilesh Pandey. Analysis of proteins and proteomes by mass spectrometry. *Annual review of biochemistry*, 70(1):437–473, 2001.
- Michelli F Oliveira, Juan P Romero, Meij Chung, Stephen Williams, Andrew D Gottscho, Anushka Gupta, Susan E Pilipauskas, Syrus Mohabbat, Nandhini Raman, David Sukovich, et al. Characterization of immune cell populations in the tumor microenvironment of colorectal cancer using high definition spatial profiling. *BioRxiv*, pages 2024–06, 2024.
- 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.
- Chandler D Gatenbee, Ann-Marie Baker, Sandhya Prabhakaran, Ottlie Swinyard, Robbert JC Slobbos, Gunjan Mandal, Eoghan Mulholland, Noemi Andor, Andriy Marusyk, Simon Leedham, et al. Virtual alignment of pathology image series for multi-gigapixel whole slide images. *Nature communications*, 14(1):4502, 2023.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018.
- Hongda Wang, Yair Rivenson, Yiyin Jin, Zhensong Wei, Ronald Gao, Harun Günaydin, Laurent A Bentolila, Comert Kural, and Aydogan Ozcan. Deep learning enables cross-modality super-resolution in fluorescence microscopy. *Nature methods*, 16(1):103–110, 2019.
- Yichen Wu, Yilin Luo, Gunvant Chaudhari, Yair Rivenson, Ayfer Calis, Kevin De Haan, and Aydogan Ozcan. Bright-field holography: cross-modality deep learning enables snapshot 3d imaging with bright-field contrast using a single hologram. *Light: Science & Applications*, 8(1):25, 2019.
- Kexin Sun, Zhineng Chen, Gongwei Wang, Jun Liu, Xiongjun Ye, and Yu-Gang Jiang. Bi-directional feature fusion generative adversarial network for ultra-high resolution pathological image virtual re-staining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3904–3913, 2023.
- B DoanNgan, DarrowMorgan Angus, LeeHan Sung, et al. Label-free virtual her2 immunohistochemical staining of breast tissue using deep learning. *BME frontiers*, 2022.
- Eric Wu, Alexandro E Trevino, Zhenqin Wu, Kyle Swanson, Honesty J Kim, H Blaize D’Angio, Ryan Preska, Aaron E Chiou, Gregory W Charville, Piero Dalerba, et al. 7-up: Generating in silico codex from a small set of immunofluorescence markers. *PNAS nexus*, 2(6):pgad171, 2023.
- Shunxing Bao, Yucheng Tang, Ho Hin Lee, Riqiang Gao, Sophie Chiron, Ilwoo Lyu, Lori A Coburn, Keith T Wilson, Joseph T Roland, Bennett A Landman, et al. Random multi-channel image synthesis for multiplexed immunofluorescence imaging. In *MICCAI Workshop on Computational Pathology*, pages 36–46. PMLR, 2021.

- Parmida Ghahremani, Yanyun Li, Arie Kaufman, Rami Vanguri, Noah Greenwald, Michael Angelo, Travis J Hollmann, and Saad Nadeem. Deep learning-inferred multiplex immunofluorescence for immunohistochemical image quantification. *Nature machine intelligence*, 4(4):401–412, 2022.
- Yair Rivenson, Hongda Wang, ZhenSong Wei, Kevin de Haan, Yibo Zhang, Yichen Wu, Harun Günaydin, Jonathan E Zuckerman, Thomas Chong, Anthony E Sisk, et al. Virtual histological staining of unlabelled tissue-autofluorescence images via deep learning. *Nature biomedical engineering*, 3(6):466–477, 2019.
- Bijie Bai, Xilin Yang, Yuzhu Li, Yijie Zhang, Nir Pillar, and Aydogan Ozcan. Deep learning-enabled virtual histological staining of biological samples. *Light: Science & Applications*, 12(1):57, 2023.
- Yijie Zhang, Kevin de Haan, Yair Rivenson, Jingxi Li, Apostolos Delis, and Aydogan Ozcan. Digital synthesis of histological stains using micro-structured and multiplexed virtual staining of label-free tissue. *Light: Science & Applications*, 9(1):78, 2020.
- Rui Cao, Scott D Nelson, Samuel Davis, Yu Liang, Yilin Luo, Yide Zhang, Brooke Crawford, and Lihong V Wang. Label-free intraoperative histology of bone tissue via deep-learning-assisted ultraviolet photoacoustic microscopy. *Nature biomedical engineering*, 7(2):124–134, 2023.
- Eric M Christiansen, Samuel J Yang, D Michael Ando, Ashkan Javaherian, Gaia Skibinski, Scott Lipnick, Elliot Mount, Alison O’neil, Kevan Shah, Alicia K Lee, et al. In silico labeling: predicting fluorescent labels in unlabeled images. *Cell*, 173(3):792–803, 2018.
- Erik A Burlingame, Mary McDonnell, Geoffrey F Schau, Guillaume Thibault, Christian Lanciault, Terry Morgan, Brett E Johnson, Christopher Corless, Joe W Gray, and Young Hwan Chang. Shift: speedy histological-to-immunofluorescent translation of a tumor signature enabled by deep learning. *Scientific reports*, 10(1):17507, 2020.
- Chang Bian, Yu Wang, Zhihao Lu, Yu An, Hanfan Wang, Lingxin Kong, Yang Du, and Jie Tian. Immunoaizer: A deep learning-based computational framework to characterize cell distribution and gene mutation in tumor microenvironment. *Cancers*, 13(7):1659, 2021.
- Chang Bian, Beth Phillips, Tim Cootes, and Martin Fergie. Hemit: H&e to multiplex-immunohistochemistry image translation with dual-branch pix2pix generator. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 184–197. Springer, 2024.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9640–9649, 2021.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical image analysis*, 81: 102559, 2022.
- Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*, 30(3):850–862, 2024.
- Charlie Saillard, Rodolphe Jenatton, Felipe Llinares-López, Zelda Mariet, David Cahané, Eric Durand, and Jean-Philippe Vert. H-optimus-0, 2024. URL <https://github.com/bioptimus/releases/tree/main/models/h-optimus/v0>.
- Fabian Hörst, Moritz Rempe, Lukas Heine, Constantin Seibold, Julius Keyl, Giulia Baldini, Selma Ugurel, Jens Siveke, Barbara Grünwald, Jan Egger, et al. Cellvit: Vision transformers for precise cell segmentation and classification. *Medical Image Analysis*, 94:103143, 2024.
- 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.
- Jingfeng Yao, Xinggang Wang, Shusheng Yang, and Baoyuan Wang. Vitmatte: Boosting image matting with pre-trained plain vision transformers. *Information Fusion*, 103:102091, 2024.

- Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022.
- Chunlong Xia, Xinliang Wang, Feng Lv, Xin Hao, and Yifeng Shi. Vit-comer: Vision transformer with convolutional multi-scale feature interaction for dense predictions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5493–5502, 2024.
- Jerry Lin, Yu-An Chen, Daniel Campton, Jeremy Cooper, Shannon Coy, Clarence Yapp, Juliann B. Tefft, Erin McCarty, Keith Ligon, Scott J. Rodig, Steven Reese, Tad George, Sandro Santagata, and Peter K. Sorger. labsyspharm/orion-crc [data set], 2023b. URL <https://doi.org/10.5281/zenodo.7637988>.
- Henoch S Hong, Robin Liechti, Christoph Reinhard, Marie M Morfouace, and IMMUCan consortium. Immucan: Broad cellular and molecular profiling of the human tumor microenvironment. *Cancer Research*, 80(16_Supplement): 1699–1699, 2020.
- N. Sofroniew, T. Lambert, G. Bokota, J. Nunez-Iglesias, P. Sobolewski, A. Sweet, L. Gaifas, K. Evans, A. Burt, D. Doncila Pop, K. Yamauchi, M. Weber Mendonça, L. Liu, G. Buckley, W.-M. Vierdag, T. Monko, L. Royer, A. Can Solak, K. I. S. Harrington, R. Zhao, et al. napari: a multi-dimensional image viewer for python (v0.6.1rc2), 2025. URL <https://doi.org/10.5281/zenodo.15421167>.
- Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a generalist algorithm for cellular segmentation. *Nature methods*, 18(1):100–106, 2021.
- Weiruo Zhang, Irene Li, Nathan E Reticker-Flynn, Zinaida Good, Serena Chang, Nikolay Samusik, Saumyaa Saumyaa, Yuanyuan Li, Xin Zhou, Rachel Liang, et al. Identification of cell types in multiplexed *in situ* images by combining protein expression and spatial information using celesta. *Nature methods*, 19(6):759–769, 2022.
- Petros Liakopoulos, Julien Massonet, Jonatan Bonjour, Medya Tekes Mizrakli, Simon Graham, Michel A Cuendet, Amanda H Seipel, Olivier Michelin, Doron Merkler, and Andrew Janowczyk. Hoverfast: an accurate, high-throughput, clinically deployable nuclear segmentation tool for brightfield digital pathology images. *arXiv preprint arXiv:2405.14028*, 2024.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16133–16142, 2023.
- Vincent Mai, Waleed Khamies, and Liam Paull. Batch inverse-variance weighting: Deep heteroscedastic regression. *arXiv preprint arXiv:2107.04497*, 2021.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.
- Marlen Runz, Daniel Rusche, Stefan Schmidt, Martin R Weihrauch, Jürgen Hesser, and Cleo-Aron Weis. Normalization of he-stained histological images using cycle consistent generative adversarial networks. *Diagnostic Pathology*, 16: 1–10, 2021.