

Training state-of-the-art pathology foundation models with orders of magnitude less data

Mikhail Karasikov¹, Joost van Doorn¹, Nicolas Känzig¹, Melis Erdal Cesur², Hugo Mark Horlings², Robert Berke¹, Fei Tang¹, and Sebastian Otálora¹

¹ Kaiko.AI, Zurich, Switzerland

² The Netherlands Cancer Institute, Amsterdam, The Netherlands

{mikhail,sebastian}@kaiko.ai

Abstract. The field of computational pathology has recently seen rapid advances driven by the development of modern vision foundation models (FMs), typically trained on vast collections of pathology images. Recent studies demonstrate that increasing the training data set and model size and integrating domain-specific image processing techniques can significantly enhance the model’s performance on downstream tasks. Building on these insights, our work incorporates several recent modifications to the standard DINOv2 framework from the literature to optimize the training of pathology FMs. We also apply a post-training procedure for fine-tuning models on higher-resolution images to further enrich the information encoded in the embeddings. We present three novel pathology FMs trained on up to two orders of magnitude fewer WSIs than those used to train other state-of-the-art FMs while demonstrating a comparable or superior performance on downstream tasks. Even the model trained on TCGA alone (12k WSIs) outperforms most existing FMs and, on average, matches Virchow2, the second-best FM published to date. This suggests that there still remains a significant potential for further improving the models and algorithms used to train pathology FMs to take full advantage of the vast data collections.

Keywords: Foundation models · Computational pathology · Whole Slide Images.

1 Introduction

Recently, there has been an increased interest in developing vision foundation models for various types of images, including medical imaging data. These FMs generate informative representations that can be used in various downstream tasks such as classification, segmentation, object detection, etc. In particular, the development of foundation models for computational histopathology, commonly referred to in the literature as *pathology FMs*, has rapidly accelerated [19]. This progress has been primarily driven by the ever-increasing amount of unlabeled Whole Slide Image (WSI) data available in public and proprietary sources, the development of more complex model architectures (e.g., ViT [11]), and the steady refinement of training workflows (e.g., DINOv2 [20]).

Most current state-of-the-art pathology FMs are based on either the DINO [9] or DINOv2 [20] self-supervised learning (SSL) algorithms. DINOv2 was developed as an extension of DINO [9] and iBOT [31] to train general-purpose vision FMs that capture both the global context and local structure of the images, using Vision Transformers (ViTs) [11] as an underlying image encoder. Depending on the type of input images, pathology FMs can be designed to produce tile- or slide-level representations. In this work, we only consider tile-level FMs. For training a tile-level FM, the original WSIs are pre-processed to extract smaller regions, *tiles* (also often referred to in the literature as *patches*), typically of size 224×224 pixels.

One of the recent milestones in the development of pathology FMs was UNI [10], a ViT-L16 model trained with DINOv2 on over 100k WSIs from various sources, where the authors set a new standard in the performance of a pathology FM and conducted numerous experiments evaluating it on diverse downstream tasks. Very recently, the same group released a successor model UNI-2 [1] trained on over 200M tiles sampled from over 350k diverse H&E and IHC WSIs. In Kaiko-FM [2], the authors trained relatively performant models solely on TCGA and introduced *online patching*, an efficient technique for sampling WSI tiles of arbitrary size directly during training to reduce the space overhead. H-optimus-0 [22] is a ViT-g14 trained with DINOv2 on 500k WSIs with several hundreds of millions of tiles. Their model is one of the largest in terms of the number of parameters and still remains one of the best-performing published models according to various benchmarks. Another prominent model is Virchow2 [32], which is a ViT-H14 trained on a substantially larger data set comprising 3.1M WSIs. In Hibou [18], the authors trained a family of FMs with 936,441 H&E, 202,464 non-H&E, and 2,676 cytology slides sourced from 306,400 unique cases. Very recently, Atlas [3] was released as a preprint, where the authors trained a new pathology FM and demonstrated outstanding performance on the HEST [15] benchmark and six out of eight downstream tasks from eva [16]. However, without released weights, external evaluation of that model appears impossible.

Training pathology FMs at large scale [32,30,22,18,10,2,18,3] have pushed the frontier by amassing tens of thousands to millions of WSIs from both public and proprietary sources. However, the tendency towards scaling leaves a critical open question: is it crucial to have such a large data set in order to train a pathology FM at the state-of-the-art level, or can similar results be achieved with far fewer WSIs?

In this work, we address the question posed above and present three novel pathology FMs trained on relatively small publicly available data sets and a proprietary set of over 80k WSIs from the Netherlands Cancer Institute (NKI). Despite being trained on orders of magnitude fewer WSIs than most other state-of-the-art models published to date, our models achieve comparable and often higher performance on most downstream tasks. We additionally perform an ablation study to determine the contributions of the individual changes we made to the standard DINOv2 training workflow. Drawing an analogy be-

tween the whole range of models applied to data at different scales (from single molecules, to cells, to tissue samples, and to entire organisms) and depth zones of the ocean, we call our pathology FMs after the middle depth zone in the ocean, bathypelagic (*midnight*) zone. The shared models and the source code and data necessary to reproduce the evaluation experiments are available at <https://github.com/kaiko-ai/midnight>.

2 Methods

Training data We trained our FMs and performed the ablation study on three public collections of WSIs: TCGA, GTEx, and CPTAC, and a proprietary data set NKI-80k. TCGA contains 12k FFPE slides from 32 cancer types collected in different hospitals by the The Cancer Genome Atlas (TCGA) Research Network: <https://www.cancer.gov/tcga>. GTEx contains 25k WSIs across 23 tissue types from 838 donor individuals, collected by the Genotype-Tissue Expression project [25]. CPTAC contains 7.2k WSIs from clinical tumor samples from 13 cohorts collected by the Clinical Proteomic Tumor Analysis Consortium [12]. In addition to those three open-access data sets, we use a proprietary set of 80k WSIs (NKI-80k) from the Netherlands Cancer Institute. These slides have magnifications of 0.25 and 0.5 $\mu\text{m}/\text{px}$, most of them being at 0.25 $\mu\text{m}/\text{px}$. This data set includes mostly FFPE H&E but also Frozen Tissue and immunohistochemistry slides from 10,141 patients and 31 organs.

In our experiments, we found that including GTEx and CPTAC slides in training did not bring substantial improvements (see Results). Thus, we trained our final FMs only on the TCGA FFPE and NKI-80k slides.

Extraction of training tiles We trained our FMs on tiles of size 256×256 cropped from the original WSIs at magnifications of 2, 1, 0.5, and 0.25 $\mu\text{m}/\text{px}$. All tiles were sampled uniformly at random from arbitrary positions of the WSIs with *online patching* [2], with the foreground area threshold set to 40%. Further, to filter out low-informative tiles (e.g., those with mainly adipose tissue), we apply a filter in the HSV color space from [32]. More precisely, a tile is only accepted if $\geq 60\%$ of its pixels have their hue, saturation, and value in ranges [90, 180], [8, 255], and [103, 255], respectively (see examples in Fig. A1-Left).

For all cropped tiles, we apply color augmentations in the Hematoxylin-Eosin-DAB (HED) space [24] (Fig. A1-Right). These augmentations effectively increase the diversity of the training data and help make the FM more robust to various staining methods used in the WSIs.

Self-supervised training with DINOv2 We use the DINOv2 self-distillation framework to train ViT-g14 models with 1.1B parameters (and ViT-B14 with 86M parameters in the ablation experiments) with self-supervised learning. Our algorithm is based on the original DINOv2 algorithm [20] with several modifications. First, as suggested in [32], we use a more stable KDE regularizer [28]

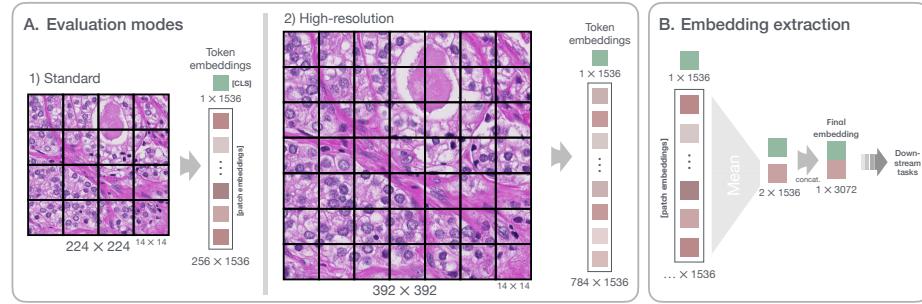


Fig. 1. Schematic representation of the FM evaluation. **Panel A:** Evaluation of a vision transformer FM at standard and high resolution. For high resolution, larger tiles of size 392×392 are cropped into $(392/14)^2 = 784$ patches of the same size 14×14 pixels. (The grids are shown schematically for simplicity. The actual numbers of patches the tiles are cropped into are 256 and 784 instead of 4^2 and 7^2 , as shown in the graph.) **Panel B:** Aggregating token embeddings produced by the ViT into the final CLS+Mean token embedding.

instead of the original KoLeo loss to ensure the diversity of tile embeddings generated by the FM. We start from the checkpoints pre-trained in [20] and train on 32 Nvidia-H100 GPUs with 80 GB memory for 1M iterations with the base learning rate of 3.5×10^{-4} and the learning schedules compressed accordingly, with the batch size of 12 per GPU. (In total, it extracts throughout the training 384×10^6 tiles from the WSIs.) We accumulate gradients over every two training steps, resulting in an effective total batch size of 768.

High-resolution post-training As in [20], after training, we optionally further fine-tune the FM on larger images for 120k iterations to improve its performance, especially on high-resolution images. Similar to reducing the patch size in the underlying vision transformer studied in [6], this technique effectively increases the number of patch tokens generated from every image by the ViT. At the same time, it allows us to start from an FM pre-trained at standard resolution and thereby shorten the training on large images. More precisely, for this post-training, we increase the size of training tiles from 256 to 512 pixels, and accordingly reduce the magnification by 2-fold, to 1, 0.5, 0.25, and $0.125 \mu\text{m}/\text{px}$, to preserve the actual size of the tile regions (512, 256, 128, and $64 \mu\text{m}$). In addition to increasing the resolution, we also scale up the parameters of the DINO transform from 98 and 224 to 168 and 392 for the local and global crop views, respectively. Since the use of larger images increases the memory requirements, we reduce the batch size to 6 per GPU and train on 48 GPUs with accumulating gradients over every four training steps, resulting in the effective total batch size of 1152. The base learning rate in this stage is reduced to 10^{-4} . Consequently, at inference time, each input image is resized to 392×392 before passing it to the FM for the embedding generation (see Fig. 1A). Note that this resizing does

not change the actual region of the WSI contained in the tile. To the best of our knowledge, we are the first to apply this high-resolution fine-tuning procedure for training pathology FMs.

Training three new FMs We trained three new FMs: 1) We trained our first FM on the 12k TCGA WSIs alone, with the training methodology described above. We refer to this model as **Midnight-12k**. 2) For our second model, we applied the same training algorithm on both TCGA and NKI-80k combined. (Each batch was sampled from TCGA or NKI-80k at random with equal probabilities.) We refer to this model as **Midnight-92k**. 3) Finally, we fine-tuned the **Midnight-92k** FM with the high-resolution post-training technique described above, with reduced training schedules, for 120k more iterations. We refer to this model as **Midnight-92k/392**.

Evaluation methodology Our evaluation protocol is based on two open-source benchmarks: eva [16] and HEST [14]. eva includes various tile- and slide-level classification tasks and two tile segmentation tasks, assessing how well the FMs encode tissue morphology in different tissues and cancers. The original data sets from which the downstream tasks in eva were derived are summarized in Table 1. For all tasks in eva, the original tiles are resized to the desired dimensions before passing them to the FMs for embedding generation (e.g., cropping the center squares from the original 700×460 tiles at $1.995 \mu\text{m}/\text{px}$ in BreaKHis and resizing them to 224×224 results in tiles of size 224×224 at $0.97 \mu\text{m}/\text{px}$.) For all tasks except Camelyon16 and Panda, we also disabled early-stopping in eva’s default protocol to ensure that all evaluation runs fully converge.

HEST includes nine tile-level tasks that evaluate how well the FM can predict gene expression from histology images. Each task is a regression of the FM’s embeddings of the 224×224 tiles at $0.5 \mu\text{m}/\text{px}$ to normalized transcript counts

Table 1. Data used in the evaluated downstream tasks. All tiles in all tasks are resized to 224×224 (or other respective dimensions) before passing to the FM for computing embeddings. (*) For slide-level tasks Camelyon16 and Panda, the values in columns ‘Tile size’ and ‘Magnification’ represent the tiles cropped from the original WSIs before resizing them to the target dimensions.

Task name	# images	Tile size	Magnification	Organ	Metric
BACH [21]	400	1536x2048	0.42 $\mu\text{m}/\text{px}$ (20x)	Breast	Bal. acc.
BRACS [7]	4,539	variable	0.25 $\mu\text{m}/\text{px}$ (40x)	Breast	Bal. acc.
BreaKHis [23]	7,909	700x460	1.995 $\mu\text{m}/\text{px}$ (40x)	Breast	Bal. acc.
CRC-100K [17]	107,180	224x224	0.5 $\mu\text{m}/\text{px}$ (20x)	Colorectal	Bal. acc.
Gleason TMA [4]	21,496	750x750	0.23 $\mu\text{m}/\text{px}$ (40x)	Prostate	Bal. acc.
MHIST [29]	3,152	224x224	1.25 $\mu\text{m}/\text{px}$ (8x)	Colorectal	Bal. acc.
PatchCamelyon [26]	327,680	96x96	1 $\mu\text{m}/\text{px}$ (10x)	Breast	Bal. acc.
Camelyon16* [5]	399 WSIs	224x224	0.25 $\mu\text{m}/\text{px}$ (40x)	Breast	Bal. acc.
Panda* [8]	1909 WSIs	448x448	0.25 $\mu\text{m}/\text{px}$ (40x)	Prostate	Bal. acc.
CoNSeP [13]	41	1000x1000	0.25 $\mu\text{m}/\text{px}$ (40x)	Colorectal	Dice score
MoNuSAC [27]	294	variable	0.25 $\mu\text{m}/\text{px}$ (40x)	Various	Dice score
HEST (all) [15]	236,495	224x224	0.5 $\mu\text{m}/\text{px}$ (20x)	Various	Pearson ρ

Table 2. Performance metrics for all evaluated FMs on the data sets from Table 1, and their average. pc10 is a tile-level classification task derived from PatchCamelyon (pc) where the training set is reduced to just ten random tiles per class (20 tiles in total). We report balanced accuracy for the classification tasks, dice score (no background) for semantic segmentation (cnsp, mnsc), and the average Pearson correlation for the nine HEST regression tasks. All classification tasks use CLS+Mean embeddings (the concatenation of the CLS token and the mean of all patch tokens in the output of the ViT). For all results, see Extended Table A1.

Model name	#WSIs	pc10	bach	brcs	bkhs	crc	glsn	mhst	pc	c16	pnd	cnsp	mnsc	HEST	Avg.
Midnight-92k/392	92k	.900	.904	.646	.802	.966	.807	.828	.951	.868	.651	.662	.708	.415	.778
UNI-2	350k	.885	.924	.651	.863	.970	.777	.829	.951	.873	.666	.626	.644	.431	.776
Midnight-92k	92k	.882	.889	.615	.793	.967	.823	.831	.948	.872	.643	.629	.656	.425	.767
Virchow2	3.1M	.835	.890	.633	.818	.966	.791	.865	.938	.860	.646	.640	.674	.403	.766
Midnight-12k	12k	.803	.907	.639	.840	.967	.790	.815	.931	.869	.656	.625	.664	.412	.763
Kaiko-B8	29k	.799	.876	.641	.842	.960	.761	.830	.920	.836	.650	.644	.686	.391	.757
tgcga-100M	12k	.789	.873	.619	.814	.968	.798	.808	.928	.870	.675	.622	.656	.415	.757
H-Optimus-0	500k	.831	.752	.620	.813	.962	.769	.850	.943	.847	.672	.644	.687	.425	.755
Prov_GigaPath	171k	.853	.794	.626	.846	.959	.727	.831	.944	.812	.657	.628	.688	.405	.752
Hibou-L	1.1M	.825	.792	.643	.767	.954	.766	.850	.949	.852	.654	.646	.668	.397	.751
UNI	100k	.833	.797	.613	.808	.954	.759	.841	.937	.854	.662	.627	.662	.391	.749
Phikon	12k	.826	.744	.579	.715	.946	.743	.824	.919	.822	.648	.624	.644	.377	.724
Phikon-v2	60k	.756	.737	.607	.725	.953	.753	.796	.900	.807	.634	.626	.645	.391	.718
Lunit	36k	.763	.785	.627	.759	.943	.758	.785	.905	.759	.604	.600	.630	.362	.714
vitg14 (nat. img.)	0	.721	.724	.578	.783	.943	.740	.855	.881	.500	.509	.565	.614	.351	.674
vitg14 (initial)	0	.652	.474	.413	.425	.754	.459	.578	.763	.526	.304	.462	.432	.166	.493

of the top 50 highly variable genes, measured at the respective positions of the tiles. Performance in HEST is measured by the Pearson correlation coefficient between the predicted and actual gene expression, computed across all patients.

3 Results and Discussion

Reaching state-of-the-art performance with less data We evaluated the performance of our FMs and several other state-of-the-art FMs on the downstream tasks described above. For every model, we evaluated both the CLS token and the CLS+Mean token embeddings (the concatenation of the CLS token and the mean of all $(\text{image_size}/\text{patch_size})^2$ patch tokens in the vision transformer, see Fig. 1B). For HEST, we only report the aggregate average of Pearson correlations. It can be seen (Table 2) that even our model **Midnight-12k** trained on just 12k WSIs is superior to most other existing FMs, and is only marginally different from Virchow2 despite being trained on 258× fewer WSIs (12k vs. 3.1M).

The **Midnight-92k** model trained on the TCGA and NKI-80k WSIs (92k WSIs in total) slightly surpasses Virchow2, and is just 0.009 behind UNI-2 (see Table 2) despite being trained on 4× fewer WSIs (92k vs. 350k). Note that UNI-2 [1] was released in Jan. 2025 as a successor of UNI [10]. Despite UNI-2 having used significantly more data for training, our models demonstrate a comparable and sometimes superior performance on the considered downstream

tasks. Unfortunately, we could not add Atlas [3] to our evaluation because we did not have the model weights for it.

Finally, our post-trained model **Midnight-92k/392** demonstrated a superior average accuracy to all other models in the benchmark and surpassed UNI-2 with an average margin of 0.002. For this evaluation, all the images were resized from their original size specified in Table 1 to 392×392 , instead of resizing them to 224×224 as for all other models. The results (Table 2) indicate that the high-resolution post-training improved the base **Midnight-92k** model especially significantly on the segmentation metrics, CoNSeP and MoNuSAC. Notably, on the PCam (10 shots) task, which is derived from PatchCamelyon by reducing the training set to just ten random tiles per class (20 in total) for every evaluation run and averaging the test accuracy over 50 training runs, this model achieves the balanced accuracy of 0.90 and surpasses all other evaluated models. However, the performance on the Camelyon16 and HEST tasks has significantly degraded, which needs further investigation. Overall, this still resulted in the absolute best-performing model among all the evaluated models.

To ensure a fair comparison, we additionally evaluated UNI [10] on all downstream tasks with resized images. (In [10], they mention a fine-tuning performed with larger 512×512 images but without the details about their procedure.) However, the performance of UNI only degraded on images resized to 512×512 (see row ‘UNI/512’ in Extended Table A1), e.g., 0.89 on PCam, which suggests that their fine-tuning procedure was of a different nature than ours. We also evaluated UNI-2 in the same way as **Midnight-92k/392**, on 392×392 images; however, that also yielded slightly lower performance (see row ‘UNI-2/392’ in Extended Table A1), which again suggests that UNI-2 does not benefit from evaluating on larger images at higher resolution.

Notably, all evaluated pathology FMs surpass the baseline ViT-g14 model trained on natural images (‘vitg14 (nat. img.)’ in Table 2) with a large margin, which highlights the importance of developing domain-specific pathology FMs. However, all pathology FMs in our benchmark perform relatively poorly on MHIST, where the baseline model trained on natural images is the second-best model. This suggests that there still remains a potential for improvement, which sets a particular aim for future work.

Ablation experiments To measure the effect of the adaptations we made to the baseline training workflow, we performed an ablation study, where we trained several smaller ViT-B14 models (Table 3). These training runs were done for 500k iterations with 4 GPUs, a batch size of 64 per GPU, and accumulating gradients over every three training steps, resulting in an effective total batch size of 768.

The first four experiments evaluated the importance of the HSV filter, the KDE regularizer, and the HED augmentations. Without replacing the default KoLeo regularizer with KDE and without the HSV filter, our training did not converge, thus, we report ‘n/a’ in the first row of Table 3. After adding the HSV filter, the training converged to an average accuracy of 0.704 on eva, which was still far from 0.753, obtained with the final config. The HED color augmentations

Table 3. Performance of the models in the ablation study. Each row corresponds to a single ViT-B14 model trained with the modifications specified by check marks.

TCGA	NKI-80k	CPTAC	GTEX	HSV	KDE	HED	eva	HEST
✓							n/a	n/a
✓				✓			0.704	0.367
✓				✓	✓		0.754	0.367
✓				✓	✓	✓	0.753	0.376
✓				✓	✓	✓	0.759	0.374
✓			✓	✓	✓	✓	0.744	0.380
✓				✓	✓	✓	0.750	0.363
✓				✓	✓	✓	0.742	0.368
✓				✓	✓	✓	0.742	0.375
✓				✓	✓	✓	0.750	0.368
✓				✓	✓	✓	0.765	0.373
✓				✓	✓	✓	0.768	0.375

improved the performance on HEST but did not have a large effect on eva. However, we applied it anyway to help make the FMs more robust to different stainings.

Next, we added each of NKI-80k, CPTAC, and GTEX, to evaluate their contribution when added to TCGA. (The tiles were sampled from both data sets with equal probability.) The results (rows 5, 6, and 7 in Table 3) show that all three had a rather small impact, with NKI-80k bringing the highest average gain in accuracy. We also ran data ablations relative to the baseline run on all four data sets: TCGA, GTEX, CPTAC, and NKI-80k (the five last rows of Table 3). Here, we excluded each data set at a time and trained a ViT-B14 on the remaining three data sets. The results show that removing GTEX and CPTAC only marginally affected the FM’s final performance, while removing TCGA and NKI-80k resulted in a higher loss.

Last, to check whether we could get an FM with a comparable performance to that of **Midnight-12k** even with less data, we trained the large ViT-g14 model on just 10% of the TCGA slides ($\sim 1k$ WSIs). The resulting performance was far lower, e.g., only 0.9 on PCAM after 500k iterations. We also trained ViT-g14 on 100M distinct tiles randomly sampled from all the 12k TCGA slides, which resulted in a slightly lower performance than that of **Midnight-12k** (row ‘**tcga-100M**’ in Extended Table A1).

Image segmentation Identifying different cell types can be essential not only for making an accurate diagnosis but also for understanding tumor behavior by analyzing the cellular composition of the micro-environment. In addition to systematically evaluating the FMs’ capability to segment and classify cells in images on the CoNSeP and MoNuSAC tasks (evaluated in eva), here, we selected two images from the CoNSeP data set for a clear visual demonstration and performed the standard semantic segmentation procedure implemented in eva for four models: ViT-g14 (natural images), Lunit, Virchow2, and our models **Midnight-12k** and **Midnight-92k/392**.

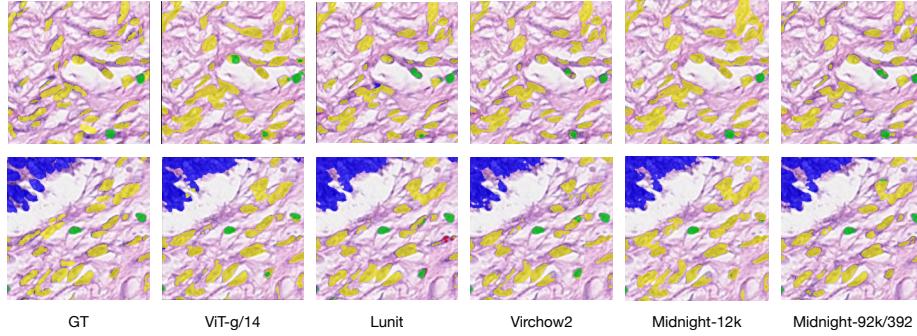


Fig. 2. Examples of segmentation performed with different FMs on two tiles from the CoNSeP data set: ViT-g14 (natural images), Lunit, Virchow2, and our models **Midnight-12k** and **Midnight-92k/392**. Ground truth is shown on the left side — green: inflammatory, blue: epithelial, yellow: spindle-shaped nuclei.

All pathology FMs produce segmentations (Fig. 2) that are noticeably better than the baseline model, which again highlights the importance of training the FMs on pathology images. Notably, our models produce segmentations that are comparable in quality to those from Virchow2, despite being trained on 34–258× fewer WSIs (12–92k vs. 3.1M).

Tile-level classification with FM for slide-level segmentation Metastasis to regional lymph nodes is an early sign of malignant spread. Thus, detecting lymph node metastasis is crucial in many cancer types, as it upstages the disease and impacts both clinical outcomes and treatment strategies. To demonstrate how tile-level pathology FMs perform on slide-level tasks, we trained tile-level

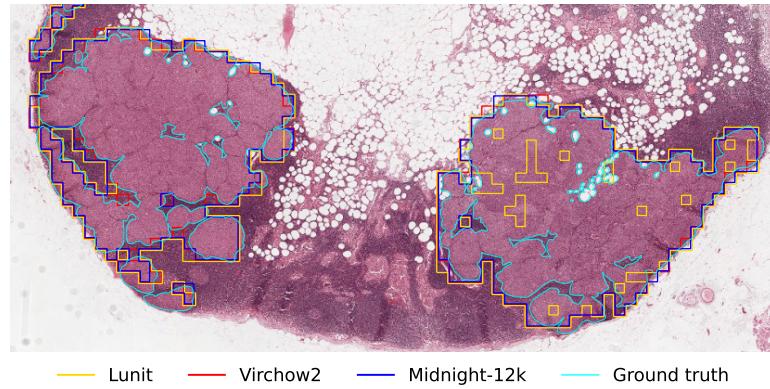


Fig. 3. Tile predictions for Lunit, Virchow2, and **Midnight-12k** and ground truth annotations for slide `test_040` from Camelyon16.

downstream classifiers to detect lymph node metastases in breast cancer on the Camelyon16 data with three different FMs: Lunit, Virchow2, and **Midnight-12k**. We applied these classifiers on a randomly picked slide `test_040` (Fig. 3). It can be seen that the predictions for **Midnight-12k** and Virchow2 are nearly identical and close to the expert annotations. At the same time, the weaker FM Lunit generates far less accurate predictions, which clearly shows the practical importance of using higher-performing FMs to get qualitatively better results on downstream tasks.

4 Conclusion

We have presented three new pathology FMs trained on up to two orders of magnitude fewer WSIs than some state-of-the-art models, yet achieving comparable or superior performance on downstream tasks. We have shown that even with a relatively basic setup, it is possible to train a high-performing pathology FM with far fewer WSIs than one may have previously thought necessary.

We make our **Midnight-12k** model trained solely on TCGA (which currently ranks third among all publicly available models) open for download from <https://huggingface.co/kaiko-ai/midnight> under the MIT license to encourage further research and reproducibility.

With our main results, we do not mean to imply that small data sets are always sufficient to reach state-of-the-art performance. On the contrary, our strong results with magnitudes fewer WSIs show that foundation model training in pathology remains far from saturation. In other words, we believe that there is still significant unrealized potential in today’s algorithms — potential that can be tapped at truly large scale. As pathology AI continues to evolve, we believe our work makes a solid contribution to the collective efforts of devising better pathology FMs. This brings us another step closer to achieving real impact in clinics and lowering the burden on pathologists while ultimately improving the quality of provided healthcare.

Acknowledgements

We thank the Core Facility for Molecular Pathology and Biobanking of the Netherlands Cancer Institute for providing the NKI-80k digitized slides with approval from the Institutional review board (IRBdm22-188). We also thank Prof. Dr. Lodewyk Wessels for his support throughout this project. The GTEx data used for the analyses described in this manuscript were obtained from the GTEx Portal <https://gtexportal.org> on October 1, 2024. The TCGA slides used in this work were generated by the TCGA Research Network: <https://www.cancer.gov/tcga>. We thank Nils Eckstein, Moritz Platscher, Thomas Hufener, and others at kaiko.ai for the helpful discussions and their valuable feedback on this manuscript. We also thank Dang Nguyen, Gianmaria Genetlici, and the rest of the Infrastructure Team at kaiko.ai for maintaining and optimizing the computing resources used in this work.

References

1. MahmoodLab/UNI2-h (revision d517a8d) (2025), <https://huggingface.co/MahmoodLab/UNI2-h>, accessed on 11 Feb 2025
2. kaiko.ai, Aben, N., de Jong, E.D., Gatopoulos, I., Känzig, N., Karasikov, M., Lagré, A., Moser, R., van Doorn, J., Tang, F.: Towards large-scale training of pathology foundation models (2024), <https://arxiv.org/abs/2404.15217>
3. Alber, M., Tietz, S., Dippel, J., Milbich, T., Lesort, T., Korfiatis, P., Krügener, M., Cancer, B.P., Shah, N., Möllers, A., et al.: Atlas: A novel pathology foundation model by mayo clinic, charit\`e, and aignostics. arXiv preprint arXiv:2501.05409 (2025)
4. Arvaniti, E., et al.: Automated gleason grading of prostate cancer tissue microarrays via deep learning. *Scientific reports* **8**(1), 12054 (2018)
5. Bejnordi, B., et al.: Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama* **318**(22), 2199–2210 (2017)
6. Beyer, L., Izmailov, P., Kolesnikov, A., Caron, M., Kornblith, S., Zhai, X., Minderer, M., Tschanne, M., Alabdulmohsin, I., Pavetic, F.: Flexivit: One model for all patch sizes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14496–14506 (2023)
7. Brancati, N., Anniciello, A.M., Pati, P., Riccio, D., Scognamiglio, G., Jaume, G., De Pietro, G., Di Bonito, M., Foncubierta, A., Botti, G., et al.: Bracs: A dataset for breast carcinoma subtyping in h&e histology images. *Database* **2022**, baac093 (2022)
8. Bulten, W., Kartasalo, K., Chen, P.H.C., Ström, P., Pinckaers, H., Nagpal, K., Cai, Y., Steiner, D.F., Van Boven, H., Vink, R., et al.: Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge. *Nature medicine* **28**(1), 154–163 (2022)
9. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9650–9660 (2021)
10. Chen, R.J., Ding, T., Lu, M.Y., Williamson, D.F., Jaume, G., Chen, B., Zhang, A., Shao, D., Song, A.H., Shaban, M., et al.: Towards a general-purpose foundation model for computational pathology. *Nature Medicine* (2024)
11. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
12. Edwards, N.J., Oberti, M., Thangudu, R.R., Cai, S., McGarvey, P.B., Jacob, S., Madhavan, S., Ketchum, K.A.: The cptac data portal: a resource for cancer proteomics research. *Journal of proteome research* **14**(6), 2707–2713 (2015)
13. Graham, S., Vu, Q.D., Raza, S.E.A., Azam, A., Tsang, Y.W., Kwak, J.T., Rajpoot, N.: Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical image analysis* **58**, 101563 (2019)
14. Jaume, G., Doucet, P., Song, A.H., Lu, M.Y., Almagro-Perez, C., Wagner, S.J., Vaidya, A.J., Chen, R.J., Williamson, D.F.K., Kim, A., Mahmood, F.: Hest-1k: A dataset for spatial transcriptomics and histology image analysis. In: Advances in Neural Information Processing Systems (Dec 2024)
15. Jaume, G., Doucet, P., Song, A.H., Lu, M.Y., Almagro-Pérez, C., Wagner, S.J., Vaidya, A.J., Chen, R.J., Williamson, D.F., Kim, A., et al.: Hest-1k: A

- dataset for spatial transcriptomics and histology image analysis. arXiv preprint arXiv:2406.16192 (2024)
16. kaiko.ai, Gatopoulos, I., Käenzig, N., Moser, R., Otálora, S.: eva: Evaluation framework for pathology foundation models. In: Medical Imaging with Deep Learning (2024), <https://openreview.net/forum?id=FNBBQ0Pj18N>
 17. Kather, J.N., Halama, N., Marx, A.: 100,000 histological images of human colorectal cancer and healthy tissue (Apr 2018). <https://doi.org/10.5281/zenodo.1214456>, <https://doi.org/10.5281/zenodo.1214456>
 18. Nechoev, D., Pchelnikov, A., Ivanova, E.: Hibou: A family of foundational vision transformers for pathology. arXiv preprint arXiv:2406.05074 (2024)
 19. Ochi, M., Komura, D., Ishikawa, S.: Pathology foundation models. JMA journal **8**(1), 121–130 (2025)
 20. Oquab, M., et al.: Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193 (2023)
 21. Polónia, A., Eloy, C., Aguiar, P.: Bach dataset : Grand challenge on breast cancer histology images (May 2019). <https://doi.org/10.5281/zenodo.3632035>, <https://doi.org/10.5281/zenodo.3632035>
 22. Saillard, C., Jenatton, R., Llinares-López, F., Mariet, Z., Cahané, D., Durand, E., Vert, J.P.: H-optimus-0 (2024), <https://github.com/bioptimus/releases/tree/main/models/h-optimus/v0>
 23. Spanhol, F.A., Oliveira, L.S., Petitjean, C., Heutte, L.: A dataset for breast cancer histopathological image classification. Ieee transactions on biomedical engineering **63**(7), 1455–1462 (2015)
 24. Tellez, D., et al.: Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology. Medical image analysis **58**, 101544 (2019)
 25. The GTEx Consortium, et al.: The genotype-tissue expression (gtex) pilot analysis: Multitissue gene regulation in humans. Science **348**(6235), 648–660 (2015). <https://doi.org/10.1126/science.1262110>, <https://www.science.org/doi/abs/10.1126/science.1262110>
 26. Veeling, B., Linmans, J., Cohen, T., Welling, M.: Rotation equivariant cnns for digital pathology (Sep 2018). https://doi.org/10.1007/978-3-030-00934-2_24, https://doi.org/10.1007/978-3-030-00934-2_24
 27. Verma, R., Kumar, N., Patil, A., Kurian, N.C., Rane, S., Graham, S., Vu, Q.D., Zwager, M., Raza, S.E.A., Rajpoot, N., et al.: Monusac2020: A multi-organ nuclei segmentation and classification challenge. IEEE Transactions on Medical Imaging **40**(12), 3413–3423 (2021)
 28. Wang, T., Isola, P.: Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In: International conference on machine learning. pp. 9929–9939. PMLR (2020)
 29. Wei, J., Suriawinata, A., Ren, B., Liu, X., Lisovsky, M., Vaickus, L., Brown, C., Baker, M., Tomita, N., Torresani, L., et al.: A petri dish for histopathology image analysis. In: International Conference on Artificial Intelligence in Medicine. pp. 11–24. Springer (2021)
 30. Xu, H., Usuyama, N., Bagga, J., Zhang, S., Rao, R., Naumann, T., Wong, C., Gero, Z., González, J., Gu, Y., et al.: A whole-slide foundation model for digital pathology from real-world data. Nature pp. 1–8 (2024)
 31. Zhou, J., Wei, C., Wang, H., Shen, W., Xie, C., Yuille, A., Kong, T.: ibot: Image bert pre-training with online tokenizer. arXiv preprint arXiv:2111.07832 (2021)

32. Zimmermann, E., Vorontsov, E., Viret, J., Casson, A., Zelechowski, M., Shaikovski, G., Tenenholz, N., Hall, J., Klimstra, D., Yousfi, R., et al.: Virchow2: Scaling self-supervised mixed magnification models in pathology. arXiv preprint arXiv:2408.00738 (2024)

A Appendix

A.1 Supplementary Figures

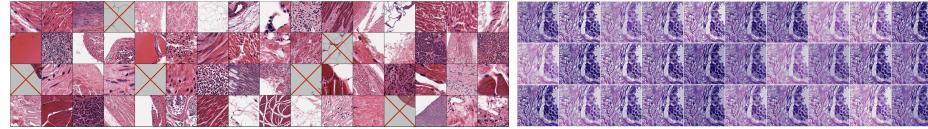


Fig. A1. Left: Random tiles cropped from TCGA FFPE slides at $0.5\mu\text{m}/\text{px}$ that passed the HSV filter and those filtered out. **Right:** Random HED augmentations applied to a single tile sampled from the BACH data set.

A.2 Extended Tables

Table A1. Performance of all evaluated FMs on the evaluated downstream tasks and their average. PCam (10 shots) is a tile-level classification task derived from PCam where the training set is reduced to just ten random tiles per class (20 tiles in total). We report balanced accuracy for the classification tasks, dice score (no background) for semantic segmentation (CoNSeP and MoNuSAC), and the average Pearson correlation for the nine HESt regression tasks. All models were evaluated in two variants on the classification tasks, using CLS+Mean token embeddings (default) and CLS only. (For the latter, the respective rows are highlighted in light grey.) The models are sorted by the average performance with the CLS+Mean tokens, and the top three values for each task are highlighted in bold. (The evaluations with the CLS tokens, rows in grey, are not taken into account for this sorting.) Each subscript for the tasks in eva represents the standard deviation of the sample mean computed over multiple runs of the respective downstream task. (*) For HESt, the subscripts represent standard deviations of the predicted gene expressions, averaged across the nine HESt tasks. (**) For the average (last column), the subscripts represent standard deviations averaged across all 13 evaluated downstream tasks.

Model name	# WSIs	PCam (10 shots)	BACH	BRACS	BreastLis	CRC-100K	Gleason	MHIST	PCam	Cam16	Panda	CoNSeP	MoNuSAC	HESt	Avg. (**)	
Midnight-92k [392]	92k	.900 ±.005	.904±.002	.646 ±.001	.802±.000	.966 ±.000	.807 ±.001	.828±.000	.951 ±.000	.868±.002	.651±.002	.662 ±.001	.708 ±.002	.415±.038	.778 ±.004	
Midnight-92k [392] [CLS]	92k	.900±.005	.906±.001	.642±.001	.850±.000	.964±.000	.809±.001	.825±.000	.951±.000	.891±.002	.633±.002	.663±.001	.707±.002	.384±.044	.774±.005	
UNI-2	350k	.858 ±.005	.924 ±.002	.651 ±.004	.863 ±.002	.970 ±.000	.777 ±.003	.829 ±.001	.951 ±.000	.873 ±.002	.666 ±.003	.626 ±.000	.644 ±.002	.414 ±.037	.776 ±.005	
UNI-2 [CLS]	350k	.857±.005	.914±.002	.661±.002	.860±.001	.965±.000	.810±.001	.867±.000	.949±.001	.868±.002	.659±.002	.628±.001	.644±.002	.424±.041	.771±.005	
UNI-2 [392]	350k	.827 ±.006	.928 ±.002	.667 ±.003	.810 ±.001	.967 ±.000	.787 ±.003	.850 ±.000	.929 ±.001	.882 ±.004	.664 ±.003	.629 ±.001	.659 ±.001	.407 ±.043	.767 ±.005	
UNI-2 [392] [CLS]	350k	.821±.006	.917±.003	.663±.003	.829±.001	.965±.000	.791±.002	.849±.001	.927±.001	.888±.003	.653±.004	.629±.001	.643±.003	.425±.032	.767±.004	
Midnight-94k	92k	.852 ±.005	.889 ±.003	.615 ±.005	.793 ±.001	.851 ±.001	.948 ±.000	.851 ±.001	.871 ±.004	.643 ±.000	.629 ±.000	.656 ±.003	.629 ±.002	.656 ±.002	.392 ±.039	.761 ±.005
Midnight-92k [CLS]	92k	.876±.006	.896±.002	.616±.004	.789±.001	.966±.000	.820±.002	.811±.000	.950±.000	.861±.001	.861±.000	.861±.001	.861±.001	.656±.002	.398±.039	.761±.005
Virchow2	3.1M	.855±.006	.890±.004	.633±.003	.818±.001	.966±.000	.791±.004	.865 ±.001	.938±.001	.860±.001	.846±.001	.640±.001	.640±.001	.674±.002	.403±.046	.766±.006
Virchow2 [CLS]	3.1M	.851±.007	.884±.004	.624±.003	.823±.004	.966±.000	.778±.006	.861±.000	.966±.001	.865±.001	.865±.001	.656±.004	.676±.002	.676±.004	.386±.040	.766±.006
Midnight-12k	12k	.803 ±.007	.907 ±.001	.639 ±.002	.840 ±.000	.967 ±.000	.790 ±.001	.815 ±.000	.931 ±.000	.869 ±.003	.656 ±.003	.625 ±.001	.664 ±.001	.412 ±.036	.763 ±.004	
Midnight-12k [CLS]	12k	.791±.007	.904±.001	.644±.001	.841±.000	.966±.000	.801±.001	.807±.001	.930±.000	.850±.004	.663±.003	.626±.001	.663±.001	.395±.040	.760±.005	
Kaike-B8	29k	.799 ±.007	.876 ±.004	.641 ±.004	.842 ±.003	.960 ±.000	.761 ±.006	.850 ±.000	.920 ±.000	.836 ±.001	.650 ±.003	.644 ±.000	.686 ±.001	.391 ±.048	.757 ±.005	
Kaike-B8 [CLS]	29k	.806±.006	.872±.002	.617±.005	.825±.004	.961±.000	.748±.004	.828±.001	.917±.001	.831±.006	.642±.003	.643±.001	.686±.001	.373±.052	.748±.007	
tgc-a-100M	12k	.789 ±.007	.873 ±.002	.619 ±.002	.814 ±.001	.968 ±.001	.774 ±.001	.808 ±.001	.928 ±.001	.870 ±.004	.675 ±.003	.622 ±.000	.656 ±.002	.415 ±.038	.757 ±.005	
tgc-a-100M [CLS]	12k	.774±.007	.864±.001	.615±.001	.779±.000	.967 ±.000	.799±.000	.792 ±.000	.927±.000	.852±.007	.667±.005	.622 ±.001	.656 ±.001	.396±.041	.747±.005	
H-Optimus-0	500k	.83±.007	.752±.004	.620±.006	.813±.003	.969±.000	.754±.004	.842±.001	.943±.001	.893±.003	.672±.003	.644 ±.002	.644 ±.002	.685±.003	.755±.006	
H-Optimus-0 [CLS]	500k	.824±.007	.757±.003	.615±.004	.808±.004	.956±.002	.771±.003	.842±.003	.942±.001	.888±.004	.670±.002	.644±.004	.644±.004	.751±.006	.751±.006	
Prov. GigaPath	171k	.853±.006	.794±.002	.620±.002	.846 ±.001	.955±.000	.751±.004	.727±.004	.851±.000	.944±.000	.812±.004	.657±.003	.628±.001	.658±.001	.405±.045	.752±.005
Prov. GigaPath [CLS]	171k	.782±.006	.766±.002	.616±.002	.817±.002	.951±.000	.766±.001	.850 ±.001	.942±.001	.871±.006	.860±.003	.660±.003	.626±.002	.687±.002	.373±.047	.743±.006
Hibon-L	1.1M	.804 ±.007	.811 ±.002	.637 ±.002	.740 ±.001	.933 ±.001	.763 ±.001	.829 ±.000	.952 ±.001	.823 ±.003	.634 ±.005	.646 ±.001	.668 ±.001	.397 ±.052	.751 ±.006	
UNI	100k	.833 ±.006	.797 ±.006	.613 ±.003	.808 ±.004	.954 ±.001	.759 ±.007	.841 ±.001	.937 ±.002	.854 ±.003	.662 ±.004	.627 ±.001	.662 ±.003	.391 ±.049	.749 ±.007	
UNI [CLS]	100k	.815±.007	.794±.002	.610±.005	.794±.009	.952±.000	.757±.005	.840±.001	.939±.001	.852±.006	.655±.004	.620±.001	.659±.004	.386±.051	.740±.007	
UNI/512	100k	.755 ±.007	.891 ±.004	.624 ±.002	.751 ±.004	.951 ±.001	.757 ±.007	.851 ±.001	.890 ±.001	.811 ±.005	.647 ±.002	.620 ±.001	.671 ±.003	.380 ±.045	.737±.007	
UNI/512 [CLS]	100k	.785±.006	.792±.002	.643±.005	.767±.002	.954±.001	.766±.001	.850 ±.001	.949±.001	.852±.003	.654±.003	.626±.002	.668±.001	.393±.047	.743±.006	
Phikon	12k	.804 ±.007	.811 ±.006	.637 ±.007	.740 ±.001	.933 ±.002	.763 ±.001	.829 ±.001	.952 ±.001	.823 ±.003	.634 ±.005	.645 ±.001	.668 ±.001	.388 ±.048	.741 ±.007	
Phikon [CLS]	12k	.820±.007	.735±.004	.568 ±.015	.713±.004	.942±.002	.729±.003	.804±.001	.923±.000	.809±.008	.644±.004	.623±.001	.644±.003	.387±.050	.747±.007	
Phikon-v2	60k	.756 ±.007	.737 ±.004	.607 ±.003	.728 ±.002	.953 ±.000	.753 ±.001	.796 ±.000	.900 ±.000	.807 ±.003	.634 ±.003	.626 ±.000	.645 ±.003	.391 ±.045	.718±.006	
Phikon-v2 [CLS]	60k	.741±.008	.734±.003	.600±.003	.716±.003	.939±.001	.751±.004	.784±.003	.893±.000	.803±.003	.631±.002	.626±.001	.645±.003	.375±.041	.711±.005	
Lunit	36k	.737 ±.008	.877 ±.008	.612 ±.001	.732 ±.004	.950 ±.000	.754 ±.001	.785 ±.004	.865 ±.001	.814 ±.006	.654 ±.002	.621 ±.001	.655 ±.002	.360 ±.003	.714±.008	
Lunit [CLS]	36k	.753±.008	.782±.008	.614±.004	.750±.007	.938±.001	.747±.003	.779±.001	.901±.001	.730±.006	.610±.004	.599±.001	.629±.001	.353±.054	.707±.008	
vigt4 (nat. img.)	0	.721±.006	.724±.006	.578±.002	.783±.001	.943±.000	.740±.002	.855 ±.001	.881±.001	.509±.010	.509±.006	.565±.001	.614±.001	.351±.042	.674±.006	
vigt4 (nat. img.) [CLS]	0	.719±.005	.728±.005	.583±.002	.832±.001	.935±.000	.744±.002	.862±.001	.874±.001	.507±.004	.382±.016	.564±.001	.614±.001	.342±.043	.668±.006	
vigt4 (initial)	0	.652±.006	.474±.005	.413±.007	.425±.003	.754±.003	.459±.003	.549±.003	.578±.004	.763±.001	.526±.008	.304±.003	.462±.001	.432±.005	.466±.010	
vigt4 (initial) [CLS]	0	.649±.006	.473±.002	.411±.004	.427±.008	.748±.003	.464±.003	.569±.000	.755±.002	.506±.006	.308±.002	.461±.003	.428±.004	.455±.010	.455±.010	