# Supplementary Material: Multistain Pretraining for Slide Representation Learning in Pathology

Guillaume Jaume<sup>1,2,\*</sup>, Anurag Vaidya<sup>1,2</sup>,\*, Andrew Zhang<sup>1,2,\*\*</sup>, Andrew H. Song<sup>1,2,\*\*</sup>, Richard J. Chen<sup>1,2</sup>, Sharifa Sahai<sup>1,2</sup>, Dandan Mo<sup>2</sup>, Emilio Madrigal<sup>2</sup>, Long Phi Le<sup>1,2</sup>, and Faisal Mahmood<sup>1,2</sup>

<sup>1</sup> Harvard Medical School, Cambridge, MA, USA <sup>2</sup> Mass General Brigham, Boston, MA, USA {gjaume, avaidya, faisalmahmood}@bwh.harvard.edu

We provide complementary information on the model architecture and training, additional results, and interpretability examples:

- 1. Sec. 1: Implementation details around pretraining, aggregator architecture, and graph optimal transport loss.
- 2. Sec. 2: Detailed descriptions of datasets used for pretraining and the downstream evaluations.
- 3. Sec. 3: Supervised Multiple Instance Learning (MIL) baselines.
- 4. Sec. 4 and 5: Additional results on downstream breast and kidney tasks.
- 5. Sec. 6: Ablations of loss and aggregator architecture.
- 6. Sec. 7: Additional interpretability examples of breast cancer cases.
- 7. Sec. 8: Limitations of MADELEINE.

### 1 Implementation details

#### 1.1 Contrastive loss

Formally, we define a batch of *B* cases, where each case includes *K* pairs  $(\mathbf{h}_i^{\text{HE}}, \mathbf{h}_i^{s_k})_{k=1}^K$ , where  $s_k$  represents a non-H&E stain. The objective  $\mathcal{L}_{\text{INFONCE}}$  is given by:

$$\mathcal{L}_{\text{INFONCE}} = -\frac{1}{K} \sum_{k=1}^{K} \frac{1}{2B} \left( \sum_{b=1}^{B} \log \frac{\exp\left(\frac{1}{\tau} \left(\boldsymbol{h}_{b}^{\text{HE}}\right)^{\text{T}} \boldsymbol{h}_{b}^{s_{k}}\right)}{\sum_{b'=1}^{B} \exp\left(\frac{1}{\tau} \left(\boldsymbol{h}_{b}^{\text{S}}\right)^{\text{T}} \boldsymbol{h}_{b'}^{s_{k}}\right)} + \sum_{b=1}^{B} \log \frac{\exp\left(\frac{1}{\tau} \left(\boldsymbol{h}_{b}^{s_{k}}\right)^{\text{T}} \boldsymbol{h}_{b}^{\text{HE}}\right)}{\sum_{b'=1}^{B} \exp\left(\frac{1}{\tau} \left(\boldsymbol{h}_{b}^{s_{k}}\right)^{\text{T}} \boldsymbol{h}_{b'}^{\text{HE}}\right)} \right),$$
(1)

where the first and second terms represent the H&E-to- $s_k$  and  $s_k$ -to-H&E contrastive loss, respectively.  $\tau$  represents the Softmax temperature parameter.

<sup>\*</sup> Co-first authorship.

<sup>\*\*</sup> Co-second authorship.

#### 1.2 Multi-head attention architecture

MADELEINE uses a multi-head attention-based Multiple Instance Learning (MIL) architecture. Before applying each head, the patch embeddings are passed through a common pre-attention network consisting of 3 layers with 512 hidden units, layer normalization, GELU activation, and 0.1 dropout. Each attention head comprises a gated-attention network, consisting of a 2-layer MLP with 512 hidden units with Softmax activation and 0.25 dropout. The attention score  $a_{i,j}^m$  for each patch derived from the  $m^{th}$  attention head of M total heads are defined as:

$$a_{i,j}^{m} = \frac{\exp\left(\mathbf{w}_{m}(\tanh(\mathbf{V}_{m}\widetilde{\mathbf{H}}_{i,j}^{\mathrm{HE}}) \odot \operatorname{sigm}(\mathbf{U}_{m}\widetilde{\mathbf{H}}_{i,j}^{\mathrm{HE}})\right)}{\sum_{j'=1}^{N_{\mathrm{HE}}}\exp\left(\mathbf{w}_{m}(\tanh(\mathbf{V}_{m}\widetilde{\mathbf{H}}_{i,j'}^{\mathrm{HE}}) \odot \operatorname{sigm}(\mathbf{U}_{m}\widetilde{\mathbf{H}}_{i,j'}^{\mathrm{HE}})\right)}, \quad \forall m.$$
(2)

The output of each head is concatenated, and a post-attention network consisting of two linear layers with 2048 and 512 units is applied to get a slide embedding for each stain.

#### 1.3 Additional information on the GOT objective

Additional details of the Graph Optimal Transport objectives are as follows,

1) Graph building: Each stain-specific graph is defined by instantiating 256 randomly sampled patches as nodes from the slide (sampling is done as each slide can have > 10,000 patches, making it computationally infeasible to calculate the complete optimal transport-based loss). Then, an edge is built between two nodes (*i.e.*, two patches) if the cosine similarity between their patch embeddings is larger than a threshold. The threshold is dynamically constructed and is set at the lowest similarity value, increasing by 0.1 times the difference between the highest and lowest similarity values.

2)  $\mathcal{L}_{WD}$ : Denoting  $\mathbf{T} \in \mathbb{R}^{N_{\text{HE}} \times N_{s_k}}_+$  as the transport plan, we can minimize the Wasserstein Distance (WD) between distributions  $\hat{p}_{\text{HE}}$  and  $\hat{p}_{s_k}$  by finding the optimal transport plan

$$\mathcal{L}_{\rm WD}(\hat{p}_{\rm HE}, \hat{p}_{s_k}) = \min_{\mathbf{T}} \sum_{j} \sum_{m} \mathbf{T}_{j,m} \cdot C(v_j^{\rm HE}, v_m^{s_k}),$$
(3)

such that  $\sum_{j=1}^{N_{\text{HE}}} \mathbf{T}_{j,m} = 1/N_{s_k}$ ,  $\forall m$  and  $\sum_{m=1}^{N_{s_k}} \mathbf{T}_{j,m} = 1/N_{\text{HE}}$ ,  $\forall j$ . The cost between cross-modal embeddings  $C(v_j^{\text{HE}}, v_m^{s_k})$  is computed with the cosine distance metric.

3)  $\mathcal{L}_{GWD}$ : In addition to the node-matching with  $\mathcal{L}_{WD}$ , we also wish to match the graph topology via comparing the edge distance between stain-specific graphs. Denoting  $\widetilde{\mathbf{T}} \in \mathbb{R}^{N_{\text{HE}} \times N_{s_k}}_+$  as the transport plan as before,

$$\mathcal{L}_{\text{GWD}}(\hat{p}_{\text{HE}}, \hat{p}_{s_k}) = \min_{\mathbf{T}} \sum_{j,j',m,m'} \widetilde{\mathbf{T}}_{j,m} \widetilde{\mathbf{T}}_{j',m'} \cdot C(v_j^{\text{HE}}, v_m^{s_k}, v_{j'}^{\text{HE}}, v_{m'}^{s_k}), \quad (4)$$

such that  $\sum_{j=1}^{N_{\text{HE}}} \widetilde{\mathbf{T}}_{j,m} = 1/N_{s_k}, \forall m \text{ and } \sum_{m=1}^{N_{s_k}} \widetilde{\mathbf{T}}_{j,m} = 1/N_{\text{HE}}, \forall j$ . The cost between the pairs  $(v_j^{\text{HE}}, v_{j'}^{\text{HE}})$  and  $(v_m^{s_k}, v_{m'}^{s_k})$  is given as  $C(v_j^{\text{HE}}, v_m^{s_k}, v_{j'}^{\text{HE}}, v_{m'}^{s_k}) = \|c(v_j^{\text{HE}}, v_{j'}^{\text{HE}}) - c(v_m^{\text{HE}}, v_{m'}^{\text{HE}})\|$ , with  $c(\cdot, \cdot)$  representing the cosine similarity metric.

#### 1.4 MADELEINE pretraining

We pretrained MADELEINE with AdamW optimizer and a batch size of 90 for 120 epochs. The learning rate is linearly ramped up during a 5-epoch warmup from 1e-9 to 1e-4. Then, we employed a cosine scheduler to reach the final learning rate of 1e-8 after 120 epochs. To increase training diversity and simplify batch processing, we sample a fixed and random subset of patches per slide, specifically 2048 patch embeddings. In slides with fewer patches, we perform random oversampling. All training settings are summarized in Appendix Table 1.

Table 1: MADELEINE pretraining and architectural hyperparameters.  $3 \times 24$ GB NVIDIA 3090Ti GPUs were used for training. Batch size refers to the total batch size across all GPUs.

Hyperparameter	Value
Heads	4
Head activation	GELU
Patch embedding dimension	512
Pre-attention hidden dimension	512
Patches sampled (training)	2048
Stain encoding dimension	32
AdamW $\beta$	(0.9, 0.999)
Batch size	90
Warmup epochs	5
Max epochs	120
Learning rate schedule	Cosine
Learning rate (start)	0
Learning rate (post warmup)	1e-4
Learning rate (final)	1e-8
Weight decay	0.01
INFONCE Temperature	0.001
Patches sampled for GOT	256
GOT $\gamma$	0
Automatic mixed precision	bloaft16
Early stopping criteria	SmoothRank [3]

#### 1.5 Early stopping with rank analysis

Following [3], we use the rank as a measure of the quality of the underlying latent space learned during MADELEINE pretraining and save the model weights from the highest rank iteration (no models are saved during the first 20 epochs of training). We compute the rank as the entropy of the d (assuming d < n) L1-normalized singular values of the slide embedding matrix  $H \in \mathbb{R}^{n \times d}$ . Specifically, we have:

$$\operatorname{RankMe}(H) = \exp(-\sum_{k=1}^{d} p_k \log(p_k)) , \qquad (5)$$

$$p_k = \frac{\sigma_k(H)}{|\sigma(H)|_1} + \epsilon \tag{6}$$

where  $\sigma_k$  denotes the k-th singular of H (sorted from large to low), and  $\epsilon$  is small constant set to 1e - 7 for numerical stability.

#### **1.6** Additional information on evaluation

Few-shot evaluation is based on linear probing and prototyping classification.

Linear probing Linear probing is implemented using a logistic regression objective based on sklearn. We use the default sklearn L2 regularization (set to 1.0) with an lbfgs solver. We set the maximum number of training iterations to 10,000 for all experiments.

**Prototyping** We define positive and negative slide "prototypes"  $p^+, p^-$  as the average of k (k=1,5,10,25) slide embeddings using downstream task labels. Subsequently, we measure the similarity between a query slide embedding  $q_i$ and the two prototypes using the L2 distance. We apply this evaluation for morphological and molecular subtyping, and allograft rejection prediction.

### 2 Datasets

Overall, our study comprises a total of 23,580 whole slide images (WSIs) from two organs (breast and kidney) and includes eight different immunohistochemistry and special stains. We use 16,281 of these WSIs for pretraining and 7,299 for downstream evaluation. We now detail all the sources of the WSIs used in the study.

Acrobat (multi-stain, pretraining) Acrobat is a multi-stain dataset originally proposed as part of the AutomatiC Registration Of Breast cAncer Tissue MICCAI challenge [12, 13]. It comprises 4,211 whole slide images (WSI) sourced from 1,153 patients diagnosed with primary breast cancer. These WSIs are available at a magnification of  $10 \times$  (equivalent to  $1\mu$ m/px) and show tissue resections, which have been processed using either hematoxylin and eosin (H&E) staining or immunohistochemistry (IHC). For every patient included in the dataset, there is one WSI that has been stained with H&E, along with a minimum of one and a maximum of four WSIs of tissue from the same tumor that has been stained with ER (N=844), PR (N=837), HER2 (N=534), or KI67 (N=843). The collection of slides was digitized at Karolinska Institutet in Stockholm, Sweden, during routine clinical workflows. Data can be downloaded at https://acrobat.grand-challenge.org/data/.

**TCGA Breast (H&E, downstream)** We collected N=1,041 primary cases from the TCGA Breast Invasive Carcinoma (BRCA) cohort, which comprises N=831 Invasive Ductal Carcinoma (IDC) and N=210 Invasive Lobular Carcinoma (ILC). For each case, we downloaded the corresponding disease-specific survival and associated censorship status, subtype (IDC and ILC), and molecular status: ER (N=996; 780 positive, 216 negative), PR (N=993; 678 positive, 315 negative), and HER2 (N=693; 158 positive, 535 negative) from UCSC Xena [4] and cBioPortal [1]. WSIs can be downloaded from Genomics Data Commons (https://portal.gdc.cancer.gov/).

**BCNB (H&E, downstream)** The Cancer Core-Needle Biopsy WSI (BCNB) dataset comprises N=1,058 patients, with a single side associated with each patient [17]. BCNB includes the molecular status of each patient: ER (N=1,058; 831 positive, 227 negative), PR (N=1,058; 790 positive, 268 negative), HER2 (1,058; 277 positive, 781 negative), and KI67 (1,058; 156 positive, 902 negative). The dataset was originally collected from hospital systems in Beijing, China, and is made publicly available at (https://paperswithcode.com/dataset/bcdalnmp).

AIDPATH (H&E, downstream) AIDPATH dataset contains 50 breast cancer WSIs stained with H&E. The dataset additionally provides HER2 expression (positive or negative, where equivocal cases are analyzed with FISH) (7 positive, 41 negative) and KI67 expression (provided as a percentage). We convert the continuous KI67 expression values to a binary task using 50% as a threshold for

IHC status prediction (19 positive, 31 negative). The dataset is made publicly available at https://mitel.dimi.uniud.it/aidpath-db.

**BWH Breast (H&E, downstream)** We collected an invasive breast cancer cohort (N=1,265) from the archives of Hopsital-A, which comprises N=982 IDC and N=283 ILC cases. All cases were primary breast cancers and included resections and biopsies. All slides were scanned at  $20 \times$  or  $40 \times$  magnification. Using patient reports, we additionally collected molecular status: ER (N=874; 613 positive, 261 negative), PR (N=874; 504 positive, 370 negative), and HER2 (N=816; 151 positive, 665 negative).

MGH Breast (Estrogen and Progesterone receptor stains, downstream) We use another private breast cohort for IHC quantification of ER abundance (N=962) and PR abundance (N=1,071). We frame both ER and PR quantification as 3- and 6-class problems. For a detailed breakdown, see Appendix Table 2.

**BWH Kidney (multi-stain, pretraining)** We collected a private renal transplant cohort comprising kidney biopsies from 1,069 renal transplant cases. Each case includes one to three tissue blocks, where each block consists of one to two H&E-stained (N=4,638) and one to two periodic acid-Schiff (PAS) (N=4,630) slides, one Jones-stained slide (N=2,326) and one Trichrome-stained slide (N=2,328). In total, each case includes 6 to 18 slides. We held out 20% of the cohort (210 cases, 463 H&E slides) as an independent test set and used the rest for MADELEINE pretraining. All slides were processed at 20×. We use H&E slides of the held-out cases to screen for Antibody-mediated rejection (AMR, 2 class; 107 positive, 356 negative) and quantify Interstitial Fibrosis and Tubular Atrophy (IFTA, 3 classes; mild : 292, moderate : 104, advanced : 67). As each case includes several H&E slides, we define two sub-tasks: "single-slide" prediction, where we use a single slide per case (N=463 H&E slides), and "all-slides" prediction, where we use all available slides per case (N=305 H&E slides).

Model/Data	$\begin{bmatrix} \mathbf{E} \\ \mathcal{C} = 6 \end{bmatrix}$	$\mathbf{R}$ $\mathcal{C} = 3$	$\mathbf{P}$ $ \mathcal{C} = 6$	$\mathbf{R}$ $\mathcal{C} = 3$
$\begin{array}{c} 0 \\ < 1\% \\ 1 - 10\% \\ 10 - 50\% \\ 50 - 90\% \\ > 90\% \end{array}$	175 160 120 172 176 159	335 292 335	168       169       219       170       175       169	337 389 344

Table 2: MGH Breast label distribution.

#### 3 Baselines

#### 3.1 Supervised multiple instance learning baselines

We provide a detailed description of the four multiple-instance learning (MIL) approaches used in the study.

- 1. **ABMIL** [6]: Attention-based multiple instance learning (ABMIL) is a popular MIL architecture. ABMIL operates as follows: first, it assigns patch-level importance scores through a gated-attention mechanism. Attention scores are used to weigh patch embeddings, which are subsequently summed to build a slide representation used for classification.
- 2. TransMIL [10]: Transformer-based multiple instance learning (TransMIL) replaces the gated attention from ABMIL with a low-rank Transformer. TransMIL first squares the sequence of low dimensional representations then applies a Pyramidal Positional Encoding module to encode spatial information and finally uses Nystrom attention [16] to approximate self-attention scores between patches. The CLS token is finally taken as the slide-level representation.
- 3. Information Bottleneck MIL [8]: Information bottlenecks (IB) are used to compress the WSI by removing uninformative instances (patch embeddings). IB aims to find patch instances that minimize the mutual information between the distribution of patches and patch representations. By only keeping such instances, [8] postulate that the most informative patches can be retained, which can then be aggregated into a slide embedding.
- 4. Low-rank MIL [15]: While TransMIL tries to learn slide-level representations by encoding patch correlations, it does not leverage the redundancy in WSIs, which [15] used to propose iterative low-rank attention (ILRA). Each ILRA block consists of two layers: one aims to project the sequence of patch representations to a low-rank space by cross-attending it with a latent matrix, and the second reconstructs the input. Performing max-pooling over the output of k such layers yields a low-rank slide-level representation.

#### 3.2 Slide-level baselines

- 1. **MEAN**: The MEAN baseline is defined by taking the arithmetic average of the patch embedding constituting the slides.
- 2. **HIPT** [2]: Hierarchical Image Pyramid Transformer (HIPT) proposes a 3level slide encoding schema, where each level is independently trained with a Transformer. The first level transforms patches into patch embeddings, which are then aggregated in region embeddings and finally into a slide embedding.
- 3. GigaSSL [7]: GigaSSL, similar to the INTRA baseline, is a method for learning slide representations based on different views of the same slide. It creates different views of a slide by sampling patches and applying augmentations such as random cropping. The different views are then pulled using a contrastive loss to learn the slide representation. Author-provided GigaSSL slide embeddings for TCGA Breast were taken from https://data.mendeley. com/datasets/d573xfd9fg/3.

- 8 G. Jaume et al.
- 4. GigaPath [18]: GigaPath is a concurrent work to ours scaling intra-SSL to large cohorts. It includes its own pan-cancer patch encoder that was pretrained on 171,000+ WSIs (> 30,000 patients) using DINOv2. The slide encoder was trained using a LongNet model with masked auto encoding. We used the official GigaPath demo<sup>3</sup> using (1)  $256 \times 256$ -pixel patching at  $20 \times$  and the latest HuggingFace tile and slide encoders. GigaPath-(MEAN) is defined by taking the average of all patch embeddings, and GigaPath is defined by building a slide embedding using global pooling of all LongNet tokens at the 11th Transformer layer.

# 4 Additional Breast results

Table 3: Few-shot molecular status prediction from H&E in TCGA Breast. Evaluation using Macro-AUC. Standard deviation reported over ten runs. All results for k = 25 training samples per class. MADELEINE refers to INFONCE + GOT. Besides HIPT and GigaSSL, all models use the *same* patch encoder. GigaSSL embeddings for BWH and BCNB cohorts were not available. Best in **bold**, second best is <u>underlined</u>.

	Madel/Data	<b>– – –</b>		(木)	D	CND	(*)	D	WH	(*)
	model/Data	FD	DD		ם. קים	DD	(I) HEDal	ED E	DD	U) HED9
		Ŀn	гn	11En2	En	гn	IIEn2	ĽŊ	гn	HER2
	ABMIL [6]	82.7	72.8	62.4	81.4	75.9	64.5	65.4	62.6	57.0
		$\pm$ 3.1	$\pm 2.5$	$\pm 3.3$	$\pm 4.5$	$\pm 2.9$	$\pm$ 3.4	$\pm$ 3.5	$\pm$ 2.6	$\pm 2.4$
	TransMIL [10]	75.1	63.5	55.0	68.7	63.4	56.0	54.8	54.3	51.8
E		$\pm 5.5$	± 7.1	± 2.9	$\pm 7.0$	± 8.2	$\pm 4.1$	± 4.2	± 3.0	± 2.7
Z	IB-MIL [8]	81.6	73.1	<u>62.6</u>	81.0	76.7	64.3	65.8	61.7	55.2
	IT DA [12]	$\pm 2.8$	± 2.5	$\pm 2.6$	$\pm 2.8$	± 2.2	$\pm 4.7$	$\pm 3.1$	± 2.5	$\pm 3.7$
	ILRA [15]	82.3	74.0	62.2	79.0	74.2	63.6	63.2	59.9	55.1
		± 2.7	± 2.4	$\pm 2.6$	± 7.2	± 1.9	± 3.9	± 3.9	± 4.7	$\pm 3.5$
	Mean	79.4	68.5	59.6	78.5	72.4	64.7	62.8	60.7	55.8
	_	$\pm$ 4.7	$\pm$ 4.2	$\pm$ 3.1	$\pm 3.2$	$\pm$ 3.8	$\pm$ 4.3	$\pm$ 2.8	$\pm$ 3.6	$\pm$ 4.3
	Intra	80.1	70.2	60.0	79.6	75.1	67.1	64.0	60.9	55.9
	UIDT [c]	$\pm 3.3$	± 3.2	$\pm 2.5$	± 2.7	± 2.7	± 3.7	± 2.2	± 3.1	$\pm 3.4$
	$HIPT_{CLS-4k}$ [2]	74.1	63.9	61.9	65.3	62.2	54.0	60.3	57.0	52.9
be	a: aar [=]	$\pm 3.0$	± 2.9	$\pm 3.7$	3.7	4.4	3.0	$\pm 3.6$	$\pm 3.6$	$\pm 3.7$
ro	GigaSSL [7]	77.7	69.4	59.9	-	-	-	_	-	-
L L	Cim Deth (Marry) [10]	$\pm 3.1$	± 2.8	$\pm 3.2$	70.0	71.0	- - 7			-
ea	GigaPath-(MEAN) [18]	11.1	08.1	58.3	10.2	(1.0	03.7	04.6	59.3	54.4
Ę.	Circo Doth [19]	$\pm 4.8$	± 2.8	± 3.4	$\pm 3.9$	± 3.9	± 3.8	± 2.2	± 3.1	± 3.4
Η	GigaPath [18]	10.0	00.7	57.4	14.1	08.9	01.8	03.0	08.2	53.2
	MADELEINE	± 5.0	T 2.1	$\pm 3.4$ 61.3	± 3.7	± 4.0	± 3.7	± 2.4	± 3.2	± 3.4
	MADELEINE	+ 2.1	+ 3.4	+ 2.4	+ 1 3	+ 2.3	+31	+ 1.8	+ 3.2	+ 3.8
	MADELEINE-SE	84.6	74.6	62.5	81.8	76.9	69.8	68.5	64.3	59.7
	MADELEINE DE	+ 2.2	+ 2.8	+ 2.6	+ 2.0	+ 2.8	+ 2.3	+ 2.0	+ 2.2	+ 3.9
	) (		<u></u>		<u></u>	= =	05.5		E 0. F	<u> </u>
	MEAN	79.6	70.8	61.1	76.9	74.8	65.5	62.5	58.7	54.3
	Intro	$\pm 4.0$	± 2.0	± 3.4	± 2.2	± 2.7	± 2.9	$\pm 2.8$	± 3.9	± 4.2
	mira	19.9	(1.0	00.3	11.1	(4.)	00.8	03.3	09.4	54.2
	HIPT or a m [9]	± 2.8	$\pm 1.4$ 64.0	$\pm 2.4$ 60.7	$\pm 1.8$	± 2.1	± 3.2	± 2.7	± 3.8	± 4.4
	1111 1 CLS-4k [2]	+ 2 0	+ 2 1	+ 3.2	+2.8	+4 1	+2.3	+ 2.0	+ 26	+ 27
ng	GigaSSL [7]	80.0	71.5	$\pm 3.3$ 62.8						± 2.1
į	C1800011 [1]	+ 2.4	+ 2.0	+ 2.4	_	_	_	_	_	_
et,	GigaPath-(MEAN) [18]	71.7	66.1	56.4	69.7	68.3	61.6	62.0	58.8	53.8
ot	ongai ann (millinn) [10]	5.6	3.8	2.6	$\pm 3.2$	± 5.1	± 3.7	± 2.7	± 2.8	$\pm 3.8$
F	GigaPath [18]	70.5	63.9	56.4	68.5	66.5	59.4	60.4	57.4	52.5
		5.4	2.7	1.9	$\pm 3.3$	± 4.8	± 4.0	± 2.9	$\pm 2.4$	± 3.4
	Madeleine	85.1	76.4	62.6	83.0	80.7	68.5	68.2	65.7	57.1
		$\pm 1.4$	± 1.2	± 2.9	± 1.6	$\pm 1.8$	± 2.2	$\pm 2.7$	$\pm 3.8$	± 3.2
	Madeleine-SE	83.3	74.9	62.9	80.5	77.3	69.8	67.2	64.7	56.9
		$\pm 1.5$	± 1.1	$\pm$ 3.1	$\pm 1.6$	$\pm 1.5$	$\pm 2.0$	$\pm 2.6$	$\pm$ 3.2	$\pm$ 3.8

Table 4: IHC quantification. We quantify the abundance of estrogen (ER) (N=962) and progesterone (PR) (N=1,071) receptor expression in 3-class and 6-class scenarios using IHC. We compare MADELEINE fine-tuned against MADELEINE architecture trained from scratch and MEAN. Results using 5-fold cross-validation with k=25 examples per class and evaluated using macro-AUC. Best in **bold**, second best is <u>underlined</u>.

Model/Data	ER	(†)	PR	. (†)
·	C = 3	C = 6	C = 3	C = 6
MEAN (linear probe)	$74.6 \pm 1.9$	$69.5\pm1.2$	$73.1 \pm 1.8$	$69.1\pm1.0$
ABMIL (Random)	$82.1 \pm 2.0$	$\underline{83.4} \pm 1.3$	$83.8 \pm 1.4$	$\underline{83.9} \pm 1.4$
ABMIL (FineTune)	$89.6 \pm 1.3$	$86.0\pm0.8$	$89.4 \pm 0.9$	$\textbf{85.5}\pm0.9$

Table 5: Survival outcome prediction in TCGA Breast. Models are trained using site-stratified 5-fold cross-validation. Evaluation using Concordance index (c-index). Besides HIPT, GigaSSL and GigaPath, all models use the *same* patch encoder. Best in **bold**, second best is <u>underlined</u>.

	Model/Data	<b>TCGA Breast</b> $(\uparrow)$
MIL	ABMIL [6] TransMIL [10] IB-MIL [8] ILRA [15]	$\begin{array}{c} 0.669 \pm 0.073 \\ \underline{0.697} \pm 0.046 \\ 0.612 \pm 0.088 \\ 0.657 \pm 0.067 \end{array}$
Slide level	MEAN INTRA HIPT <sub>CLS-4k</sub> [2] GigaSSL [7] GigaPath-(MEAN) [18] GigaPath [18] MADELEINE MADELEINE-SE	$\begin{array}{c} 0.687 \pm 0.079 \\ 0.692 \pm 0.069 \\ 0.547 \pm 0.078 \\ 0.530 \pm 0.038 \\ 0.587 \pm 0.091 \\ 0.521 \pm 0.083 \\ \textbf{0.715} \pm 0.041 \\ 0.696 \pm 0.073 \end{array}$

**Table 6: Molecular subtyping from H&E.** Detection of HER2 and KI67 status (binary) from H&E in AIDPATH and BCNB datasets. Results of MADELEINE and MADELEINE-SE obtained using linear probing. "SL" stands for slide level. Results using 5-fold stratified cross-validation evaluated using macro-AUC. Best in **bold**, second best is <u>underlined</u>.

	Model/Data	$\begin{array}{c} \textbf{HER2} (\uparrow) \\ \textbf{AIDPATH} \end{array}$	KI6 AIDPATH	7 (↑) BCNB
MIL	ABMIL [6] TransMIL [10] IB-MIL [8] ILRA [15]	$\begin{array}{c} 81.1\pm8.9\\ 46.4\pm10.7\\ 73.2\pm11.1\\ 76.1\pm7.8\end{array}$	$\begin{array}{c} \underline{89.2} \pm 7.7 \\ 65.1 \pm 19.9 \\ 87.7 \pm 6.1 \\ 84.9 \pm 4.9 \end{array}$	$\begin{array}{c} \underline{81.9} \pm 3.7 \\ 74.9 \pm 10.1 \\ 81.6 \pm 3.5 \\ 78.8 \pm 3.6 \end{array}$
$\mathbf{SL}$	Mean Intra Madeleine Madeleine-SE	$\begin{array}{c} 77.4  \pm  20.5 \\ \underline{85.8}  \pm  17.4 \\ 81.5  \pm  9.9 \\ \textbf{92.5}  \pm  7.2 \end{array}$	$\begin{array}{c} 80.2\pm2.8\\ 80.2\pm5.9\\ \textbf{91.3}\pm4.9\\ 83.0\pm8.6\end{array}$	$\begin{array}{c} 79.6 \pm 4.0 \\ 80.9 \pm 3.6 \\ 81.4 \pm 4.2 \\ \textbf{82.0} \pm 3.6 \end{array}$

# 5 Additional kidney results

Table 7: Kidney rejection tasks. Linear probe and prototyping for k = 50 reported. HIPT and GigaSSL are not available for non-cancer datasets. Best in **bold**, second best is <u>underlined</u>.

	Model/Data	IF	<b>FA</b> (†)	AN	$\mathbf{IR}(\uparrow)$
		Slide	Patient	Slide	Patient
	ABMIL [6]	74.5	78.2	69.6	71.2
		$\pm 1.9$	$\pm 4.9$	$\pm 3.7$	$\pm$ 5.8
	TransMIL [10]	57.5	60.4	55.9	55.4
Η		$\pm 3.7$	$\pm$ 6.6	$\pm 5.4$	$\pm$ 10.5
Σ	IB-MIL [8]	73.0	80.0	67.4	69.6
		$\pm$ 3.7	$\pm$ 5.1	$\pm 4.3$	$\pm$ 8.5
	ILRA [15]	70.9	77.5	62.4	63.5
		$\pm 5.0$	$\pm$ 4.5	$\pm$ 6.8	$\pm~10.4$
	Mean	73.6	79.3	67.8	70.9
ē		$\pm 2.6$	$\pm$ 2.1	$\pm 3.2$	$\pm$ 3.2
ġ.	Intra	74.5	80.2	67.9	72.0
đ		$\pm 2.2$	$\pm 1.8$	$\pm 3.1$	$\pm$ 3.2
ar	Madeleine	<u>75.3</u>	80.9	71.2	73.8
ine		$\pm 2.5$	$\pm$ 2.3	2.9	$\pm$ 3.2
Ξ	MADELEINE-SE	76.1	82.4	<u>70.0</u>	74.2
		$\pm$ 2.0	$\pm$ 1.8	$\pm$ 2.9	$\pm$ 3.5
	Mean	70.2	75.0	63.8	67.5
60		$\pm 2.4$	$\pm$ 2.8	$\pm 6.0$	$\pm$ 8.1
in	Intra	70.6	75.7	63.4	66.5
K		$\pm 2.5$	$\pm$ 3.2	$\pm 4.3$	$\pm$ 5.4
fo	Madeleine	72.1	77.0	66.7	71.2
ľ		$\pm 2.4$	$\pm$ 3.0	$\pm 3.7$	$\pm$ 4.3
д	MADELEINE-SE	73.1	79.8	65.5	70.4
		$\pm 2.4$	$\pm$ 2.0	$\pm 4.1$	$\pm$ 5.4

## 6 Additional ablations

Table 8: Ablation study of MADELEINE feature extractor. Survival was evaluated using c-index and site-stratified 5-fold cross-validation. Subtyping and molecular status prediction were evaluated using macro-AUC and prototyping evaluation (k=25) repeated five times with fixed seed across baselines. Standard deviation reported over the 5 runs. MADELEINE refers to pretraining on breast cancer using INFONCE + GOT without stain encoding. CONCH is the patch encoder of the Vision+Language model proposed in [9]. Best in **bold**, second best is <u>underlined</u>.

Model/Data		BWH Subtyping (↑)	TCGA PR (†)	BCNB EB. (↑)	Avg
		Subtyping (I)		<b>Die</b> (1)	
CTransPath+Mean	68.6	81.1	65.0	67.7	70.6
	$\pm 4.0$	$\pm 3.9$	$\pm 1.5$	$\pm 2.1$	
CONCH+Mean	68.7	86.2	70.8	76.9	75.6
	$\pm 7.9$	$\pm$ 7.9	$\pm$ 3.0	$\pm$ 2.0	
CTRANSPATH+MADELEINE	65.4	83.1	66.9	68.6	71.0
	$\pm 6.5$	$\pm 4.6$	$\pm 1.6$	$\pm 3.5$	
CONCH+MADELEINE	71.5	94.9	76.4	83.0	81.5
	$\pm$ 4.1	$\pm$ 0.9	$\pm$ 1.9	$\pm$ 1.6	

Table 9: Ablation study of MADELEINE architecture. Survival was evaluated using c-index and site-stratified 5-fold cross-validation. Subtyping and molecular status prediction were evaluated using macro-AUC and prototyping evaluation (k=25) repeated five times with fixed seed across baselines. Standard deviation reported over the 10 runs. "MH" refers to multi-head attention, "SH" to single-head attention, and "SE" to stain encoding. MADELEINE refers to pretraining on breast cancer using IN-FONCE + GOT. Best in **bold**, second best is <u>underlined</u>.

Model/Data	$\begin{array}{c c} \mathbf{TCGA} \\ \mathbf{Survival} \ (\uparrow) \end{array}$	BWH Subtyping (†)	$\begin{array}{c} \mathbf{TCGA} \\ \mathbf{PR} \ (\uparrow) \end{array}$	$\begin{array}{c} \mathbf{BCNB} \\ \mathbf{ER} \ (\uparrow) \end{array}$	Avg
MADELEINE-SH	70.1	91.8	75.0	80.5	79.4
	$\pm$ 7.1	$\pm 1.7$	$\pm 1.4$	$\pm 2.5$	
MADELEINE w. TransMIL	55.7	90.8	75.6	82.1	76.5
	$\pm 7.8$	$\pm 1.7$	$\pm 1.3$	$\pm 1.3$	
Madeleine w. SE	69.6	95.8	74.9	80.5	80.2
	± 7.3	$\pm 0.8$	$\pm 1.1$	$\pm$ 1.6	
Madeleine-MH	71.5	94.9	76.4	83.0	81.5
	$\pm$ 4.1	$\pm 0.9$	$\pm$ 1.2	$\pm$ 1.6	

## 7 Additional interpretability examples



Fig. 1: Additional heatmap examples obtained with MADELEINE A. Attention weights of multi-headed (frozen) ABMIL slide encoder pretrained with MADELEINE overlaid on three randomly chosen samples for TCGA Breast cohort. We show all heads and the average of heads. B. Attention weights of a single head (frozen) ABMIL slide encoder pretrained with MADELEINE overlaid on three randomly chosen samples for TCGA Breast cohort. Multi-headed ABMIL trained with MADELEINE can focus on different morphologies, whereas single-headed ABMIL focuses only on tumor morphology.

### 8 Limitations

MADELEINE is a multimodal pre-training strategy for slide representation learning. It operates under the assumption that representation learning of H&E images can be guided by other stains (immunohistochemistry and special stains). This premise is directly inspired by the standard practice in clinical settings, where H&E staining is routinely performed as the *gold standard* procedure, along with complementary stains. Though this approach is principled, we highlight some limitations of our study and this methodology more broadly.

**Data scaling** Clinical practice is complex and ever evolving. Every year, new IHC and special stains become available, some of which are integrated in the workflow and can be used on a case-by-case basis. In breast cancer, our study focuses on four IHC stains (the most common ones), whereas many more can be employed, such as Epidermal Growth Factor Receptor (EGFR), P53, and E-Cadherin. As each stain offers a different view of a biomarker, increasing the number of stains would make the training signal richer and the resulting representation potentially better.

Lack of large public datasets Acrobat is the only *large-scale* public dataset with H&E and IHC stains. Therefore, without relying on proprietary data, such method cannot be scaled to more stains and other types of cancer. While the NADT-Prostate [14] cohort includes H&E and IHC, it remains limited by its size; for example, 14/18 stains provided have less than 100 examples, preventing efficient pre-training in prostate adenocarcinoma. In addition, TCGA includes known limitations such as site-specific biases [5] and demographic biases [11]. Despite these limitations, TCGA remains the largest public resource for cancer prognostication and survival analyses.

**Model scaling** MADELEINE is trained using a combination of a global objective using contrastive learning and a local objective using graph optimal transport (GOT). Using our current hardware  $(3 \times 3090 \text{ GPUs})$ , we are limited by the maximum batch size for contrastive learning, even using efficient parallelization and bfloat16 quantization. In addition, computing GOT is computationally expensive, with significant memory requirements. Because of this constraint, we must use 256 patch embeddings (or tokens) per stain for computing GOT. Scaling to more tokens would allow finer-grained matching between stains. Local alignment through GOT also requires morphological overlap between tissue sections used for different stains.

#### References

- Cerami, E., Gao, J., Dogrusoz, U., Gross, B.E., Sumer, S.O., Aksoy, B.A., Jacobsen, A., Byrne, C.J., Heuer, M.L., Larsson, E., et al.: The cbio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. Cancer discovery 2(5), 401–404 (2012) 5
- Chen, R.J., Chen, C., Li, Y., Chen, T.Y., Trister, A.D., Krishnan, R.G., Mahmood, F.: Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16144–16155 (2022) 7, 9, 10
- Garrido, Q., Balestriero, R., Najman, L., Lecun, Y.: Rankme: Assessing the downstream performance of pretrained self-supervised representations by their rank. In: International conference on machine learning (10 2022) 3, 4
- Goldman, M.J., Craft, B., Hastie, M., Repečka, K., McDade, F., Kamath, A., Banerjee, A., Luo, Y., Rogers, D., Brooks, A.N., Zhu, J., Haussler, D.: Visualizing and interpreting cancer genomics data via the xena platform. Nature biotechnology 38(6), 675–678 (2020) 5
- Howard, F.M., Dolezal, J., Kochanny, S., Schulte, J., Chen, H., Heij, L., Huo, D., Nanda, R., Olopade, O.I., Kather, J.N., et al.: The impact of site-specific digital histology signatures on deep learning model accuracy and bias. Nature communications 12(1), 4423 (2021) 15
- Ilse, M., Tomczak, J., Welling, M.: Attention-based deep multiple instance learning. In: International conference on machine learning. pp. 2127–2136. PMLR (2018) 7, 9, 10, 11
- Lazard, T., Lerousseau, M., Decencière, E., Walter, T.: Giga-ssl: Self-supervised learning for gigapixel images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4304–4313 (2023) 7, 9, 10
- Li, H., Zhu, C., Zhang, Y., Sun, Y., Shui, Z., Kuang, W., Zheng, S., Yang, L.: Taskspecific fine-tuning via variational information bottleneck for weakly-supervised pathology whole slide image classification. In: 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (03 2023) 7, 9, 10, 11
- Lu, M., Chen, B., Williamson, D., Chen, R., Liang, I., Ding, T., Jaume, G., Odintsov, I., Zhang, A., Le, L., Gerber, G., Parwani, A., Mahmood, F.: Towards a visual-language foundation model for computational pathology. Nature Medicine (2024) 12
- Shao, Z., Bian, H., Chen, Y., Wang, Y., Zhang, J., Ji, X., et al.: Transmil: Transformer based correlated multiple instance learning for whole slide image classification. Advances in Neural Information Processing Systems 34, 2136–2147 (2021) 7, 9, 10, 11
- Vaidya, A., Chen, R.J., Williamson, D.F., Song, A.H., Jaume, G., Yang, Y., Hartvigsen, T., Dyer, E.C., Lu, M.Y., Lipkova, J., et al.: Demographic bias in misdiagnosis by computational pathology models. Nature Medicine **30**(4), 1174– 1190 (2024) 15
- Weitz, P., Valkonen, M., Solorzano, L., et al.: A multi-stain breast cancer histological whole-slide-image data set from routine diagnostics. Sci Data 10, 562 (2023)
  5
- Weitz, P., Valkonen, M., Solorzano, L., Carr, C., Kartasalo, K., Boissin, C., Koivukoski, S., Kuusela, A., Rasic, D., Feng, Y., et al.: A multi-stain breast cancer histological whole-slide-image data set from routine diagnostics. Scientific Data 10(1), 562 (2023) 5

- 14. Wilkinson, S., Ye, H., Karzai, F., Harmon, S.A., Terrigino, N.T., VanderWeele, D.J., Bright, J.R., Atway, R., Trostel, S.Y., Carrabba, N.V., Whitlock, N.C., Walker, S.M., Lis, R.T., Sater, H.A., Capaldo, B.J., Madan, R.A., Gulley, J.L., Chun, G., Merino, M.J., Pinto, P.A., Salles, D.C., Kaur, H.B., Lotan, T.L., Venzon, D.J., Choyke, P.L., Turkbey, B., Dahut, W.L., Sowalsky, A.G.: Nascent prostate cancer heterogeneity drives evolution and resistance to intense hormonal therapy (sep 2020) 15
- Xiang, J., Zhang, J.: Exploring low-rank property in multiple instance learning for whole slide image classification. In: The Eleventh International Conference on Learning Representations (2022) 7, 9, 10, 11
- Xiong, Y., Zeng, Z., Chakraborty, R., Tan, M., Fung, G., Li, Y., Singh, V.: Nyströmformer: A nyström-based algorithm for approximating self-attention. In: Proceedings of the AAAI Conference on Artificial Intelligence (2021) 7
- Xu, F., Zhu, C., Tang, W., Wang, Y., Zhang, Y., Li, J., Jiang, H., Shi, Z., Liu, J., Jin, M.: Predicting axillary lymph node metastasis in early breast cancer using deep learning on primary tumor biopsy slides. Frontiers in oncology 11, 759007 (2021) 5
- Xu, H., Usuyama, N., Bagga, J., Zhang, S., Rao, R., Naumann, T., Wong, C., Gero, Z., González, J., Gu, Y., Xu, Y., Wei, M., Wang, W., Ma, S., Wei, F., Yang, J., Li, C., Gao, J., Rosemon, J., Bower, T., Lee, S., Weerasinghe, R., Wright, B.J., Robicsek, A., Piening, B., Bifulco, C., Wang, S., Poon, H.: A whole-slide foundation model for digital pathology from real-world data. Nature (2024) 8, 9, 10