

# MVH-MIL: Computation-Aware WSI Analysis with Hybrid Hyperbolic-Euclidean Embeddings and Dynamic Sparsity

Anonymous CVPR submission

Paper ID 11826

## Abstract

001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030

Pathological whole slide images (WSIs) are fundamental in pathology analysis, but their massive image scale and subtle morphological features demand highly trained and experienced pathologists. Deep learning techniques have thus been applied to assist pathologists; however, three key challenges remain in WSI analysis: (i) suboptimal geometric representation, as conventional embeddings inadequately model the intrinsic hierarchical tissue ontologies; (ii) a geometric trade-off between capturing global hierarchy context and preserving local discriminative detail, which hampers multi-scale reasoning; and (iii) inefficient, homogeneous computation pipelines that fail to accommodate pathological heterogeneity. Existing methods offer limited improvements, as they rely on a single geometric paradigm or homogeneous computations. We propose MVH-MIL, a novel geometric-aware hierarchical framework for WSI analysis that: (1) introduces hyperbolic embeddings to naturally model the global, tree-like organization of pathological structures; (2) employs a Hybrid Geometric Embedding strategy that synergistically combines hyperbolic representations with Euclidean embeddings; and (3) integrates a Mixture of WSI Experts (MoWE) module that adaptively routes tissue-specific tokens to specialized sub-networks for efficient, expert-driven handling of pathological heterogeneity. Across 7 WSI datasets spanning 7 cancer types, MVH-MIL consistently improves performance over strong state-of-the-art backbones. Extensive experiments demonstrate that MVH-MIL provides a route for a more robust and scalable approach to efficient digital pathology modeling. The project will be open-sourced.

031

## 1. Introduction

032  
033  
034  
035  
036

Accurate cancer diagnosis and subsequent treatment planning hinge on pathological interpretation, the gold standard for identifying disease from microscopic images [11]. With the rise of digital pathology, tissue biopsies are now commonly scanned into gigapixel-scale Whole Slide Images

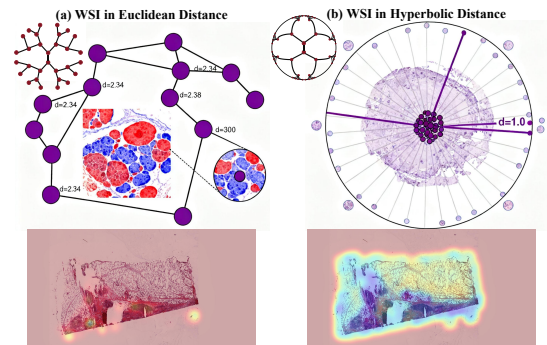


Figure 1. WSI geometric representations. (a) Euclidean embeddings fail to model hierarchical relationships. (b) Hyperbolic geometry naturally encodes the tree-like structure of pathology, from tissue architecture to cellular morphology.

(WSIs) that capture a wealth of multi-scale spatial information. Despite their richness in information, the sheer size of these images (e.g.,  $60,000 \times 80,000$  pixels) presents a major bottleneck. Manual WSI review is an arduous task demanding highly specialized expertise, which introduces the risk of inconsistent and variable diagnostic outcomes [20]. Furthermore, their massive spatial scale, high morphological heterogeneity, and inherent hierarchical organization (cells to entire slides) pose major barriers to automated analysis [9, 12]. To address these challenges, Multiple Instance Learning (MIL) is adopted as a dominant framework for WSI analysis [24, 33]. It tackles gigapixel scale by partitioning WSIs into small tiles (e.g.,  $224 \times 224$  pixels) and using embeddings to compress the features. In MIL, a foundation model (e.g., CONCH [22]) first converts tile latent representations into embeddings, and a slide-level model (e.g., ABMIL [14]) then aggregates these embeddings for slide-level predictions. Though practical, current MIL methods struggle to balance three critical requirements:

037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055

1) **Hierarchical modeling capability** [23]. The intrinsic hierarchical nature of WSIs poses challenges for representation in Euclidean space, which is commonly used by traditional slide-level methods [14, 30]. The Euclidean metric alone is suboptimal for capturing complex spatial

056  
057  
058  
059  
060

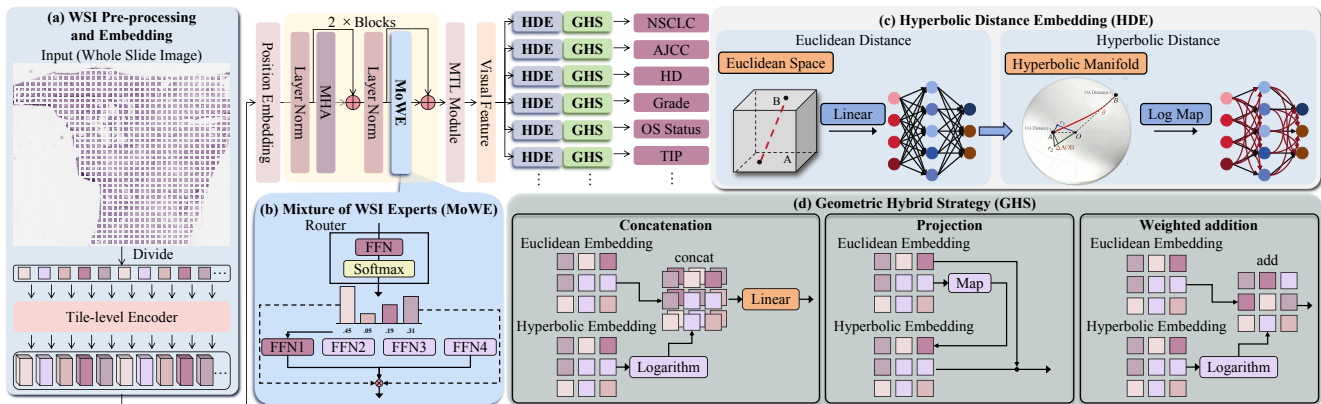


Figure 2. MVH-MIL overview. (a) A WSI is first partitioned into tiles, then extracted into embeddings with a pathological foundation model for the later slide-level modeling. (b) MoWE employs WSI-specialized experts to extract domain-specific features. (c) HDE converts Euclidean representations to hyperbolic distance to capture both hierarchical relationships. (d) GHS combines hyperbolic and Euclidean embeddings to preserve local details and create an integrated, multi-scale representation.

relationships [13] and restricts the capacity of existing approaches to model structural hierarchies, ranging from individual cells to glands and large tissue architecture (Fig. 1).

2) *Efficient homogeneous computation* [34]. Features in pathology images are highly heterogeneous. Diagnostic regions of interest are typically small and dispersed within extensive areas of benign or irrelevant tissues. Existing methods [27, 28] apply uniform computation across entire slides, failing to adapt to this variability and wasting substantial resources on trivial regions.

3) *Multi-scale reasoning* [25]. WSI analysis requires the simultaneous understanding of high-level tissue architecture and fine-grained cellular morphology. The geometric trade-off limits the ability of current methods [21, 34] to effectively capture global hierarchical context while preserving locally discriminative details. It also impedes the integration of information across spatial scales.

To address these requirements, we propose Mixture Visual Hierarchy MIL (MVH-MIL) (Fig. 2), a novel MIL framework for WSI analysis with three key contributions:

1) We introduce **hyperbolic embeddings modeled in the Poincaré ball** [17, 31] leveraging the unique ability of hyperbolic geometry to represent hierarchical tissue structures. 2) We integrate a **dynamic sparse activation mechanism** into a Vision Transformer (ViT) to enable computation-aware feature extraction. It adaptively focuses computational resources on diagnostically relevant regions, improving efficiency while preserving feature discriminability. 3) We employ a **hybrid embedding strategy** that integrates hyperbolic embeddings with Euclidean embeddings, utilizing Euclidean space to preserve local morphological details and achieve a unified representation of multi-scale contextual information.

Experiments on 7 downstream tasks over 10 state-of-the-art (SOTA) methods across 7 WSI datasets, and prove

MVH-MIL’s resilience and flexibility. In summary, by integrating pathology-aware inductive biases and employing a targeted hybrid geometric embedding strategy, our framework effectively meets the core requirements in WSI analysis with superior overall performance.

## 2. Related Works

### 2.1. Multiple Instance Learning

Deep learning for WSI analysis has evolved toward a two-stage MIL paradigm, separating tile-level feature extraction from slide-level aggregation. Early approaches, such as SlideAve (averaging) and SlideMax (max-pooling), adopted simplistic aggregation schemes that either diluted critical signals or amplified noise. To address this, attention-based methods like ABMIL [14] and CLAM [21] became central, using adaptive weighting to prioritize informative tiles. Further refinements like DTFD-MIL [34] and DSMIL [19] used feature distillation and contrastive learning, respectively. However, these attention-driven models failed to solve the spatial context problem, struggling to capture inter-tile relationships. To solve this spatial limitation, TransMIL [27] used transformers to encode tile relationships. This introduced new issues: poor modeling of hierarchical features and a computational bottleneck from quadratic attention. Research then focused on this bottleneck, splitting into two branches. One branch applies large backbones like LongNet in GigaPath [30], remaining memory-intensive and monolithic. The other branch, State Space Models (SSMs), captures long-range dependencies with linear or sub-quadratic complexity. Representatives includes MambaMIL [32] and PathRWKV [2]. While both branches improve efficiency, they converge on the same unresolved problem: their monolithic designs are ill-equipped to model the diverse and heterogeneous feature distributions found in WSIs.

## 2.2. Mixture of Experts in Transformers

The Mixture of Experts (MoE) paradigm has re-emerged as a leading approach for scaling Transformer models beyond the limits of dense computation. It replaces the standard dense Feed-Forward Network (FFN) layers with multiple parallel "expert" FFNs. For each input token, a sparse subset of these experts (e.g., the Top-1 or Top-2) are dynamically selected by a lightweight gating network called router. This conditional computation mechanism decouples the model's total parameter count from its computational cost. It allows for a massive increase in model capacity while keeping the computational cost for inference constant. The resurgence was pioneered by GShard-M4 [18], which first successfully applied MoE to scale Transformers for machine translation. Building on this, the Switch Transformer [8] significantly simplified the architecture by using a Top-1 routing strategy, and introduced an auxiliary load balancing loss to ensure stable training and prevent expert collapse. More recently, models like Mixtral 8x7B [15] have demonstrated the power of Top-K routing. They achieve superior performance by combining the outputs of the two best-scoring experts. MoE has also proven effective beyond language, as Vision MoE [26] successfully adapted it to the Vision Transformer by routing individual image patches, confirming MoE's status as a general-purpose scaling solution.

## 2.3. Hyperbolic Space

Previous MIL methods generally rely on Euclidean geometry. This is considered "geometrically suboptimal" for WSI representation [13], as it fails to capture the intrinsic hierarchical organization of WSIs (Fig. 1a). To better modeling these structures, hyperbolic geometry has been proposed as an effective alternative [13, 17, 31]. Specifically, hyperbolic spaces (e.g., the Poincaré ball [31]) exhibit exponential volume growth, which is well-suited for embedding tree-like, hierarchical data. This matches the features in WSIs (Fig. 1b). Recent works such as HyperPath [13] and HVHM [17] have leveraged hyperbolic neural networks to model these tissue hierarchies, enabling an improved representation of structural organization. However, these methods often face a fundamental geometric trade-off [7, 25]: while hyperbolic space excels at capturing global hierarchy, it distorts fine-grained local and discriminative features that are better preserved in Euclidean space. This trade-off creates the need to jointly model global hierarchy and local details.

## 3. Methods

### 3.1. WSI Pre-processing and Embedding

A WSI  $I$  is first loaded at a specific microns-per-pixel (mpp) resolution and subsequently tessellated into a non-

overlapping grid of tiles  $\{T_{i,j}\}$ , each of size  $T_s$  (Fig. 2a). To ensure the quality of these tiles, a two-step screening process is applied. First, tiles are discarded if their tissue content is below a predefined threshold (e.g.,  $< 70\%$  area coverage). Second, tiles are eliminated if their pixel variance is beneath a quantitative cutoff (e.g.,  $Var(I) < 0.05$ , with  $I$  as the pixel intensity normalized to the  $[0,1]$  range). This ensures only informative tiles proceed to the next stage.

Each valid tile is then converted into a compact, semantic feature vector via a pathological foundation model (e.g., GigaPath [30]). This feature extraction step serves to improve both the training efficiency and the overall performance of the downstream slide-level backbone (e.g., TransMIL [27]).

### 3.2. Mixture of WSI Experts

The slide-level backbone is based on ViT [6], with the major modification on its FFN component. A basic FFN processes the same transformation for each tile embedding  $h_i$ :

$$FFN(h_i) = GELU(h_i W_1 + b_1) W_2 + b_2 \quad (1)$$

This is incompatible with the heterogeneous nature of WSI, ignoring its diverse semantic properties. To address this, we adapt the concept of Mixture-of-Experts (MoE) [26] to computational pathology by proposing the Mixture of WSI Experts (MoWE) (Fig. 2b). This layer replaces the dense FFN with  $K$  specialized expert networks  $\{E_1, E_2, \dots, E_K\}$ , each learning distinct pathological patterns. However, computing all  $K$  experts for every token  $h_i$  (a "dense" MoE) is computationally prohibitive. Instead, MoWE implements a sparse, Top-K routing strategy. A single, learnable routing network  $R$  (implemented as a single linear layer) first computes  $K$  gating logits  $l_i \in \mathbb{R}^K$  for each token  $h_i$ :

$$l_i = R(h_i) = h_i W_g \quad (2)$$

where  $W_g$  is the router's weight matrix. Next, rather than directly implement softmax as MoE, we use a Top-K operation to select the indices  $I$  of the Top-K experts:

$$I, l_{i,k} = \text{Top-K}(l_i, k) \quad k \ll K \quad (3)$$

where  $I = \{j_1, \dots, j_k\}$  is the set of chosen expert indices and  $l_{i,k}$  are their corresponding logits. We calculate the gating weights  $w_i$  by applying softmax only to these  $k$  logits:

$$w_i = \text{Softmax}(l_{i,k}) \quad (4)$$

The final output  $\text{MoWE}(h_i)$  is the sparsely weighted sum of only the  $k$  activated experts:

$$\text{MoWE}(h_i) = \sum_{j \in I} w_{i,j} \cdot E_j(h_i) \quad (5)$$

This sparse approach significantly reduces computational cost, as only  $k$  experts are computed. It also maintains high model capacity with a large  $K$ .

225 Additionally, we introduce the auxiliary loss  $\mathcal{L}_{\text{aux}}$  during  
226 training to prevent expert collapse, where the router over-  
227 whelmingly favors a few experts. It encourages the router to  
228 distribute tokens evenly across all  $K$  experts. This is based  
229 on the fraction of tokens  $f_i$  routed to expert  $i$ , and the mean  
230 routing probability  $P_i$  for that expert across the batch:

$$231 \quad \mathcal{L}_{\text{aux}} = \alpha \cdot \sum_{i=1}^K f_i \cdot P_i \quad (6)$$

232 where  $\alpha$  is a hyperparameter.  $\mathcal{L}_{\text{aux}}$  is added to the main  
233 model loss and optimized jointly.

### 234 3.3. Hyperbolic Distance Embedding

235 We adopt hyperbolic distance embedding (HDE) (Fig. 2c)  
236 on the output from the slide-level backbone, additionally  
237 modeling the features from a broader scale. Existing MIL  
238 methods [14, 19, 21, 30] often operate in Euclidean space  
239 (Fig. 1a), measuring similarity via the  $L_2$  norm:

$$240 \quad D_{\text{euc}}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (7)$$

241 This geometry excels at modeling local discriminative  
242 patterns but misaligns with the intrinsic, multi-scale hierar-  
243 chy of WSIs. To address this, we adopt hyperbolic space  
244 via the Poincaré ball model  $\mathbb{M}_c^n$ , a principled alternative for  
245 embedding hierarchies. Its distance metric  $D_{\text{hyp}}$  is derived  
246 from Möbius addition ( $\oplus_c$ ):

$$247 \quad D_{\text{hyp}}(\mathbf{x}, \mathbf{y}) = \frac{2}{\sqrt{c}} \operatorname{arctanh}(\sqrt{c} \|\mathbf{x} \oplus_c \mathbf{y}\|) \quad (8)$$

$$248 \quad \mathbf{x} \oplus_c \mathbf{y} = \frac{(1 + 2c\langle \mathbf{x}, \mathbf{y} \rangle + c\|\mathbf{y}\|^2) \mathbf{x} + (1 - c\|\mathbf{x}\|^2) \mathbf{y}}{1 + 2c\langle \mathbf{x}, \mathbf{y} \rangle + c^2\|\mathbf{x}\|^2\|\mathbf{y}\|^2} \quad (9)$$

249 This projection is a form of manifold learning, embed-  
250 ding hierarchical tissue structures from Euclidean space  
251 onto a hyperbolic manifold. It scales exponentially  
252 ( $O(e^{cr})$ ), unlike polynomial Euclidean growth ( $O(r^n)$ ),  
253 preserving their intrinsic geometry.

### 254 3.4. Geometric Hybrid Strategy

255 The global advantage of HDE introduces a local trade-  
256 off. Hyperbolic geometry’s reliance on a conformal factor,  
257  $\lambda_c = \frac{2}{1 - c\|\mathbf{x}\|^2}$ , can destabilize fine-grained local features.  
258 The critical dichotomy of global hierarchical capacity ver-  
259 sus local feature preservation motivates our geometric fu-  
260 sion. This is based on the premise that hyperbolic geometry  
261  $D_{\text{hyp}}$  generalizes Euclidean geometry  $D_{\text{euc}}$  when  $c \rightarrow 0$ :

$$262 \quad \lim_{c \rightarrow 0} D_{\text{hyp}}(\mathbf{x}, \mathbf{y}) = 2\|\mathbf{x} - \mathbf{y}\| \quad (10)$$

263 This demonstrates the mathematical continuity between  
264 the two geometries. We therefore propose three distinctive  
265 fusion strategies:

266 **Concatenation:** With identical  $(z_E, z_H)$  embeddings, fu-  
267 sion can be performed using concatenation  $\oplus$ :

$$z_{\text{concat}} = [z_E \oplus \log_0^c(z_H)], \quad (11)$$

269 this method preserves complete geometric information  
270 while expanding dimensionality to  $2m$ , followed by a lin-  
271 ear layer to compress the channel dimension back to  $m$ .

272 **Projection:** For  $h_i, h_j \in \mathbb{R}^d$ , we compute  $d_E(h_i, h_j) =$   
273  $\|h_i - h_j\|_2$  and their hyperbolic embeddings  $z_{H_i} =$   
274  $\exp_0^c(W_H h_i)$ ,  $z_{H_j} = \exp_0^c(W_H h_j)$  in  $\mathbb{H}_c^m$ . The fused dis-  
275 tance  $d_f$  is:

$$d_f(h_i, h_j) = d_H(z_{H_i}, z_{H_j}) \cdot \left(1 + \gamma \cdot \phi\left(\frac{d_E(h_i, h_j)}{\sigma_E}\right)\right) \quad (12)$$

276 where  $d_H(\cdot, \cdot)$  denotes the Poincaré distance in  $\mathbb{H}_c^m$ ,  $\sigma_E$  is  
277 the dataset’s Euclidean scale, and  $\gamma$  is a learnable intensity.  
278 The projection function  $\phi: \mathbb{R} \rightarrow [0, 1]$  employs a tem-  
279 pered sigmoid activation, with learnable centering parame-  
280 ter  $\mu$  and temperature  $\kappa$  controlling the transition sharpness:  
281

$$\phi(x) = \frac{1}{1 + e^{-\kappa(x - \mu)}}, \quad (13)$$

282 The Poincaré distance exhibits exponential growth with the  
283 embedding radius:  
284

$$d_H(z_{H_i}, z_{H_j}) = \operatorname{arcosh}\left(1 + 2\frac{\|z_{H_i} - z_{H_j}\|^2}{(1 - \|z_{H_i}\|^2)(1 - \|z_{H_j}\|^2)}\right) \quad (14)$$

285 **Weighted addition:** For a tile-level feature  $h \in \mathbb{R}^d$ , the  
286 hybrid embedding  $(z_E, z_H)$  combines  $z_E = W_E h \in$   
287  $\mathbb{E}^m$  and  $z_H = \exp_0^c(W_H h) \in \mathbb{H}_c^m$ , where  $\exp_0^c(\mathbf{v}) =$   
288  $\tanh(\sqrt{c} \frac{\|\mathbf{v}\|}{2}) \frac{\mathbf{v}}{\sqrt{c}\|\mathbf{v}\|}$ . The fused representation is:  
289

$$z = \alpha \cdot z_E + (1 - \alpha) \cdot \log_0^c(z_H), \quad (15)$$

290 where  $\alpha$  is a learnable gate, balancing the geometries via  
291 the inverse log mapping ( $\log_0^c$ ). Element-wise summation  
292 preserves dimensionality for cross-geometric refinement.  
293

294 These strategies represent three distinct fusion philoso-  
295 phies. In our empirical evaluations, weighted addition was  
296 ultimately adopted as the geometric hybrid strategy (GHS)  
297 (Fig. 2d). It avoids deferring the complex integration bur-  
298 den to a subsequent linear layer, as concatenation does. It  
299 also doesn’t limit the fusion’s scope purely to distance com-  
300 putation, as projection does. Instead, this strategy projects  
301 the hyperbolic representation carrying hierarchical context  
302 back to its tangent space via the logarithmic map ( $\log_0^c$ ). It  
303 also performs a direct, learnable interpolation with the Eu-  
304 clidean representation carrying local details. This process

305 yields a unified feature representation that most effectively  
306 coordinates the strengths of both geometries. It aligns per-  
307 fectly with the intrinsic multi-scale nature of WSI analysis.

## 308 4. Experiment

### 309 4.1. Datasets and Downstream Tasks

310 We validate MVH-MIL on 7 core downstream tasks across  
311 7 WSI datasets, covering diverse diagnostic scenarios in  
312 computational pathology (Fig. 3). The CAMELYON16  
313 [1] dataset with the **Breast Metastasis (BrMet)** task clas-  
314 sifies lymph nodes as normal or tumorous. The TCGA [4]  
315 datasets evaluate multiple diverse diagnostic tasks. Specifi-  
316 cally, TCGA-Lung performs classification of subtypes of  
317 **non-small cell lung cancer (NSCLC)** and predicts the  
318 **AJCC-defined tumor staging (T-Stage)**, a metric reflect-  
319 ing the tumor’s local progression. For **Histological Diagn-  
320 osis (HistoDx)** tasks, TCGA-BRCA classifies breast cancer  
321 subtypes and TCGA-GBM diagnoses brain cancer based  
322 on morphological patterns. Meanwhile, TCGA-BRCA  
323 predicts **HER2 (IHC-HER2)** expression status from WSIs  
324 to assess suitability for targeted therapy. For **Cancer Grad-  
325 ing (Grade)** tasks, TCGA-BLCA and TCGA-CESC per-  
326 form on bladder and cervical tissues respectively, determin-  
327 ing tumor differentiation levels. And TCGA-CESC dataset  
328 with the **Lymphovascular Invasion (LymInv)** task identi-  
329 fies the presence of cancer invasion into lymphatic or vas-  
330 cular channels. The TCGA-LGG dataset with the **T-Stage**  
331 task predicts brain tumors staging.

### 332 4.2. Implementation Details

333 All models were uniformly initialized without pre-training.  
334 The training regimen included a 20-epoch linear warm-  
335 up, followed by 80 epochs optimized via Adam [16], with  
336 the learning rate decayed to  $0.01 \times$  its initial value us-  
337 ing a cosine schedule. We adopted hyperparameter search  
338 over three orders of magnitude ( $1 \times 10^{-4}$ ,  $1 \times 10^{-5}$ , and  
339  $1 \times 10^{-6}$ ), and reported the best performance for each  
340 model. We set  $Batchsize = 4$ , with  $Samp_{Max} = 2,000$   
341 per WSI. Experiments were conducted on 4 NVIDIA RTX  
342 4090 and 2 NVIDIA H100, using Python 3.10.16, PyTorch  
343 2.4.0, and CUDA 12.1. The pipeline is supported by Un-  
344 Puzzle [20].

### 345 4.3. Comparison with SOTA Methods

346 We compared MVH-MIL against 2 baselines: **SlideAve** and  
347 **SlideMax**, and 8 SOTA methods: **ABMIL** [14], **CLAM**  
348 [21], **DSMIL** [19], **TransMIL** [27], **LongNet** [5], **GigaP-  
349 ath** [30], **S4MIL** [10], and **MambaMIL** [32]. The compar-  
350 ison is conducted on 7 challenging pathological tasks across  
351 7 distinct datasets, which cover diverse clinical objectives  
352 (Tab. 1). This comprehensive evaluation demonstrates the  
353 generalizability and robustness of our proposed method.

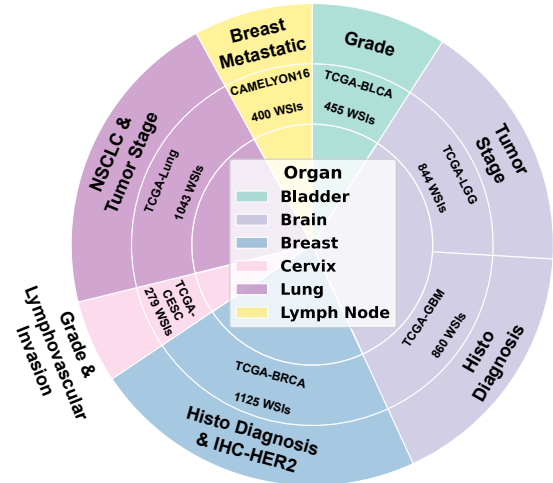


Figure 3. The implemented datasets and downstream tasks.

354 The results clearly highlight the superiority of MVH-  
355 MIL. Our method achieves SOTA performance, rank-  
356 ing first on 7 out of the 7 tasks. Specifically, MVH-  
357 MIL achieves the best scores on CAMELYON16 BrMet  
358 (96.9%), TCGA-BRCA HistoDx (86.4%), TCGA-BRCA  
359 IHC-HER2 (63.4%), TCGA-CESC GRADE (56.1%), and  
360 TCGA-GBM HistoDx (99.3%). It also achieves the joint-  
361 best performance on TCGA-Lung NSCLC (75.7%, tied  
362 with SlideMax) and TCGA-BLCA GRADE (97.7%, tied  
363 with LongNet). Furthermore, for the remaining 3 tasks  
364 where it does not rank first, MVH-MIL consistently se-  
365 cures the second-best performance, as indicated by the underlined  
366 scores for TCGA-Lung T-Stage (50.9%), TCGA-LGG T-  
367 Stage (77.9%), and TCGA-CESC LymInv (70.4%). This  
368 consistent high performance across all evaluated tasks un-  
369 derscores the model’s robustness. Notably, MVH-MIL sig-  
370 nificantly outperformed all other SOTA methods ( $p < 0.05$ )  
371 in the Wilcoxon signed-rank test [29], with the only excep-  
372 tion being S4MIL ( $p = 0.064$ ). This result validates the  
373 statistical significance of our approach.

### 374 4.4. Visualization Analysis

375 To understand where different methods focus, we visualize  
376 Class Activation Maps (CAMs) using Grad-CAM [3]. We  
377 compare the CAMs of three MVH-MIL variants (MVH-  
378 MIL<sub>cat</sub>, MVH-MIL<sub>emb</sub>, MVH-MIL<sub>add</sub>) with Euclidean-  
379 based methods (Fig. 4). The results show a clear distinction  
380 in focus. Most Euclidean-based methods, such as ABMIL  
381 and GigaPath, produce overly smooth and diffuse heatmaps  
382 that create a halo effect over the entire tissue. Others, like  
383 DSMIL and S4MIL, fail to pinpoint specific critical regions,  
384 showing minimal or no meaningful activation.

385 In contrast, all the MVH-MILs generate more inter-  
386 pretable CAMs. They successfully identify multi-scale spa-  
387 tial relationships, as their heatmaps are not one uniform,  
388 monolithic blob. Instead, they highlight both broad con-

| Method   | TCGA-BRCA<br>HistoDx | TCGA-GBM<br>HistoDx | TCGA-Lung<br>T-Stage | TCGA-LGG<br>T-Stage | TCGA-BLCA<br>Grade | TCGA-CESC<br>Grade | TCGA-Lung<br>NSCLC | CAMELYON16<br>BrMet | TCGA-BRCA<br>IHC-HER2 | TCGA-CESC<br>LymInv | P-value |
|----------|----------------------|---------------------|----------------------|---------------------|--------------------|--------------------|--------------------|---------------------|-----------------------|---------------------|---------|
| SlideAve | 81.0                 | 97.2                | 46.3                 | 74.5                | 96.6               | 52.6               | <u>75.1</u>        | 69.6                | 61.2                  | 63.0                | 0.002   |
| SlideMax | 80.1                 | 97.9                | 50.0                 | 70.5                | 95.6               | 47.4               | <b>75.7</b>        | 96.1                | 57.9                  | 59.3                | 0.008   |
| ABMIL    | 83.3                 | 97.9                | <b>51.8</b>          | <u>77.9</u>         | 96.6               | <u>54.4</u>        | 71.6               | 93.7                | <u>62.8</u>           | 66.7                | 0.015   |
| CLAM     | 78.7                 | 94.4                | 45.4                 | 75.8                | 95.5               | <u>54.4</u>        | 72.2               | 73.4                | 54.1                  | 66.7                | 0.002   |
| DSMIL    | 81.9                 | 97.9                | 50.0                 | 75.2                | 96.6               | 50.9               | 73.4               | 86.1                | 59.0                  | 63.0                | 0.002   |
| TransMIL | 81.9                 | <u>98.6</u>         | <u>50.9</u>          | 74.5                | 96.6               | 50.9               | 73.4               | 93.7                | 59.0                  | <u>70.4</u>         | 0.011   |
| LongNet  | 80.5                 | 94.5                | 45.4                 | 75.8                | <b>97.7</b>        | 50.5               | 74.6               | <u>96.2</u>         | 55.2                  | 66.7                | 0.008   |
| GigaPath | <u>84.2</u>          | <u>98.6</u>         | 50.0                 | 75.8                | <u>97.3</u>        | 50.9               | <u>75.1</u>        | 95.3                | <u>62.8</u>           | 66.7                | 0.002   |
| S4MIL    | 81.4                 | <u>98.6</u>         | <b>51.8</b>          | 73.8                | 96.6               | 52.6               | 74.6               | 93.8                | 58.5                  | <b>74.1</b>         | 0.064   |
| MambaMIL | 83.3                 | 97.9                | 50.0                 | <b>79.2</b>         | 95.5               | 50.9               | 71.0               | 95.3                | 59.6                  | 63.0                | 0.006   |
| MVH-MIL  | <b>86.4</b>          | <b>99.3</b>         | <u>50.9</u>          | <u>77.9</u>         | <b>97.7</b>        | <b>56.1</b>        | <b>75.7</b>        | <b>96.9</b>         | <b>63.4</b>           | <u>70.4</u>         | -       |

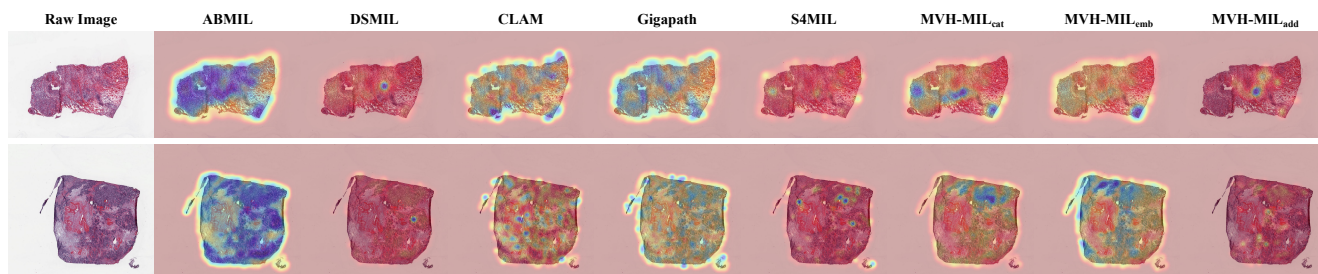
Table 1. The performance comparison with SOTA methods. **Bold**: the best result. Underline: the second best result.

Figure 4. Grad-CAM visualization for different MIL methods. Compared to other Euclidean space-based methods, all three MVH-MIL methods accurately identify the multi-scale spatial relationships among tiles.

389 textual regions and specific, high-attention focuses. This  
 390 structured attention demonstrates a more sophisticated un-  
 391 derstanding of the WSI. Notably, MVH-MIL<sub>add</sub> provides  
 392 the most refined visualization, with clean boundaries and a  
 393 superior capability to seek pathologically significant areas.

## 394 5. Discussion

395 In this section, we evaluate the effectiveness of each pro-  
 396 posed module. We first take out a baseline, SlideViT, based  
 397 on ViT [6]. We modify it to fit as the slide-level backbone.  
 398 Specifically, the patch embedding module is replaced with  
 399 a linear layer to directly take the tile embeddings. The origi-  
 400 nal 1D position embedding is also replaced with a 2D Sin-  
 401 Cos position embedding from GigaPath [30] to embed the  
 402 coordinates of tiles. We then replace its original FFN mod-  
 403 ule with our proposed MoWE module, or individually add  
 404 the GHS module to evaluate their individual contributions.  
 405 These revisions are named SlideMoWE and SlideGHS.

### 406 5.1. Mixture of WSI Experts Analysis

#### 407 5.1.1. Module Ablation

408 We first evaluate the effectiveness of MoWE (Tab. 2). The  
 409 results show that MoWE is a crucial and highly effec-  
 410 tive component for improving model performance. It out-  
 411 performs the SlideViT in 8 out of 7 downstream tasks,  
 412 with notable gains in TCGA-Lung T-Stage (49.1% vs.  
 413 47.7%) and TCGA-LGG T-Stage (77.2% vs. 75.8%). This  
 414 demonstrates that the dynamic expert routing mechanism  
 415 in MoWE is more powerful and adaptive than the static  
 416 FFN. This effectiveness is further validated by the training

and validation curves in Fig. 6. Although the training loss  
 of SlideMoWE is relatively high during comparison, it ex-  
 hibits the most stable validation curve among all compared  
 methods, showing strong generalization and excellent resis-  
 tance to overfitting. Furthermore, the radar chart (Fig. 5i)  
 confirms that SlideMoWE achieves a broad and competi-  
 tive performance profile across diverse tasks. However, the  
 above results also demonstrate the weakness of individually  
 implementing MoWE. SlideMoWE is beaten by SlideViT  
 in 3 of 10 experiments, demonstrating its instability.

#### 5.1.2. Hyperparameter Evaluation

We then compare different hyperparameter settings in  
 MoWE to evaluate their influence on the model perfor-  
 mance (Tab. 3). Specifically, we adopt three sets of num-  
 ber of experts  $N$  ( $N = 4, 8, 16$ ) in MoWE, each with  
 corresponding Top-Ks (Top-K  $< N$ ). Results on three  
 benchmark datasets and downstream tasks demonstrate that  
 Top-K = 2 often yields strong results (with maximum  
 +1.8%, +2.1%, +2.3% on TCGA-Lung, TCGA-BLCA, and  
 CAMELYON16, compared to the lowest score in each  
 task), balancing performance and efficiency. However, its  
 performance is not consistent across different  $N$ . Specifi-  
 cally, the Top-K = 2,  $N = 4$  setting achieves the best re-  
 sult on CAMELYON16 (93.8%), yet it underperforms com-  
 pared to  $N = 8$  and  $N = 16$  settings on TCGA-BLCA  
 (96.6% vs. 97.7%). This may be due to the different WSI  
 processing strategies in different datasets, and different at-  
 tention requirements for different downstream tasks.

On the other hand,  $N$  also shows a complex relationship  
 with performance. A larger  $N$  does not consistently lead to

| Method    | TCGA-BRCA<br>HistoDx | TCGA-GBM<br>HistoDx | TCGA-Lung<br>T-Stage | TCGA-LGG<br>T-Stage | TCGA-BLCA<br>Grade | TCGA-CESC<br>Grade | TCGA-Lung<br>NSCLC | CAMELYON16<br>BrMet | TCGA-BRCA<br>IHC-HER2 | TCGA-CESC<br>LymInv | P-value |
|-----------|----------------------|---------------------|----------------------|---------------------|--------------------|--------------------|--------------------|---------------------|-----------------------|---------------------|---------|
| SlideViT  | 80.5                 | 97.2                | 47.7                 | 75.8                | 97.7               | 47.4               | 74.0               | 93.7                | 56.3                  | 66.7                | 0.008   |
| SlideMoWE | 81.9                 | 97.2                | <u>49.1</u>          | <u>77.2</u>         | 96.6               | 45.6               | <u>74.6</u>        | <u>93.8</u>         | 56.8                  | 63.0                | 0.002   |
| SlideGHS  | <u>83.3</u>          | <u>97.9</u>         | 44.5                 | 75.8                | 96.6               | 45.6               | 74.0               | 92.2                | <u>57.9</u>           | 59.3                | 0.002   |
| MVH-MIL   | <b>86.4</b>          | <b>99.3</b>         | <b>50.9</b>          | <b>77.9</b>         | <b>97.7</b>        | <b>56.1</b>        | <b>75.7</b>        | <b>96.9</b>         | <b>63.4</b>           | <b>70.4</b>         | -       |

Table 2. Ablations on the proposed MoWE and GHS modules. **Bold**: the best result. Underline: the second best result.

| Top-K | TCGA-Lung<br>NSCLC |             |             | TCGA-BLCA<br>Grade |             |             | CAMELYON16<br>BrMet |      |      |
|-------|--------------------|-------------|-------------|--------------------|-------------|-------------|---------------------|------|------|
|       | N=4                | N=8         | N=16        | N=4                | N=8         | N=16        | N=4                 | N=8  | N=16 |
| 2     | <b>74.6</b>        | <b>74.6</b> | <b>74.6</b> | 96.6               | <b>97.7</b> | <b>97.7</b> | <b>93.8</b>         | 91.5 | 91.5 |
| 4     | -                  | <b>74.6</b> | 72.8        | -                  | <b>97.7</b> | 96.6        | -                   | 93.0 | 91.5 |
| 8     | -                  | -           | 72.8        | -                  | -           | 95.6        | -                   | -    | 93.0 |

Table 3. Ablations on different number of experts and Top-K.

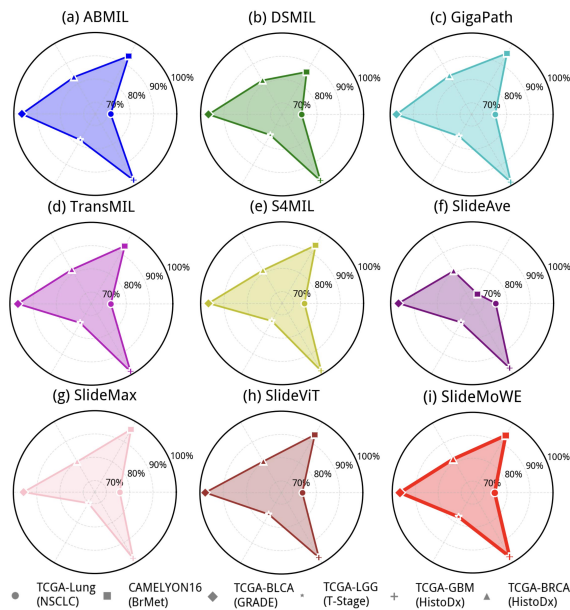


Figure 5. Performance comparison radar charts.

447 better results. For instance, when Top-K = 2, increasing  
 448  $N$  from 4 to 8 improves the score on TCGA-BLCA (from  
 449 96.6% to 97.7%), but significantly degrades it on CAME-  
 450 LYON16 (from 93.8% to 91.5%). For TCGA-Lung, the per-  
 451 formance remains identical (74.6%) regardless of  $N$ . Fur-  
 452 thermore, when Top-K = 4, increasing  $N$  from 8 to 16  
 453 results in a performance drop across all three datasets. This  
 454 suggests that simply adding more experts may not be an ef-  
 455 fective strategy and can even lead to inefficient expert spe-  
 456 cialization, especially when a larger  $K$  is also used.

## 457 5.2. Hyperbolic Distance Embedding Analysis

458 Manifold curvature  $hyp_c$  is the critical component that  
 459 determines the geometric structure in the Poincaré ball  
 460 model. To investigate how  $hyp_c$  influences the perfor-  
 461 mance of the MVH-MIL variants, we conduct experiments  
 462 using 4 different values for this hyperparameter ( $hyp_c \in$

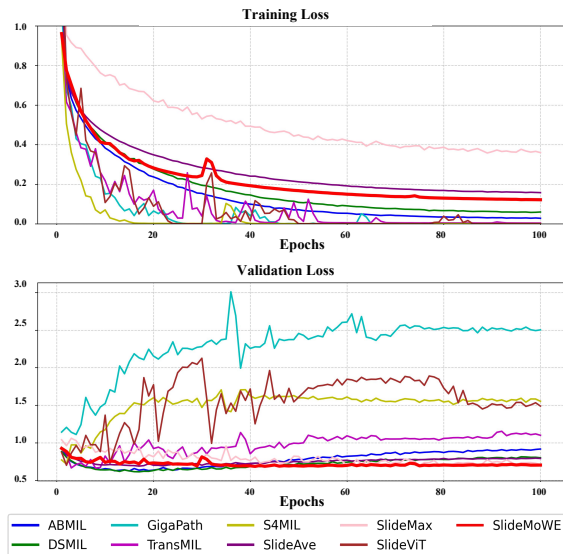


Figure 6. Training and validation loss comparison.

{0.05, 0.075, 0.1, 0.2}). Notably, as  $hyp_c \rightarrow 0$ , the hy-  
 463 perbolic space flattens and collapses to Euclidean space;  
 464 a larger  $hyp_c$  preserves more hyperbolic features in contrast.  
 465

466 Results reveal that the optimal  $hyp_c$  value is not univer-  
 467 sal but is tied to dataset structural properties and variant-  
 468 specific inductive biases (Tab. 4). Specifically,  $hyp_c =$   
 469 {0.05, 0.075, 0.1, 0.2} achieve 8, 14, 14, and 9 opti-  
 470 mal performances on 30 experiments, respectively. Though  
 471 dataset and variant-specific, moderate  $hyp_c$  values ( $hyp_c =$   
 472 0.075, 0.1) tend to show better generalizability.  $hyp_c = 0.1$   
 473 achieves the best overall performance, with a slight average  
 474 performance boost of 0.2% compared with  $hyp_c = 0.075$   
 475 (75.19% vs. 74.99%). We eventually adopt  $hyp_c = 0.1$   
 476 based on its superior overall performance.

## 477 5.3. Geometric Hybrid Strategy Analysis

### 478 5.3.1. Module Ablation

479 The fusion strategy is critical for deciding the combination  
 480 structure of the different geometry distances. We evaluate  
 481 the effectiveness of GHS, compare the performance of three  
 482 fusion strategies, and further evaluate the generalizability of  
 483 GHS on other methods. We first evaluate the effectiveness  
 484 of the GHS module (Tab. 2). The results indicate that im-  
 485 plementing GHS individually provides mixed and often un-  
 486 stable results. While it outperforms the SlideViT baseline  
 487 in 3 out of 10 tasks, with notable gains in TCGA-BRCA

| Method                 | $hyp_c$ | TCGA-BRCA<br>HistoDx | TCGA-GBM<br>HistoDx | TCGA-Lung<br>T-Stage | TCGA-LGG<br>T-Stage | TCGA-BLCA<br>Grade | TCGA-CESC<br>Grade | TCGA-Lung<br>NSCLC | CAMELYON16<br>BrMet | TCGA-BRCA<br>IHC-HER2 | TCGA-CESC<br>LymInv |
|------------------------|---------|----------------------|---------------------|----------------------|---------------------|--------------------|--------------------|--------------------|---------------------|-----------------------|---------------------|
| MVH-MIL <sub>cat</sub> | 0.05    | <b>81.9</b>          | 97.9                | 47.2                 | <b>75.8</b>         | 95.5               | 50.9               | 71.0               | 93.8                | <b>61.7</b>           | 71.4                |
|                        | 0.075   | <b>81.9</b>          | <b>99.3</b>         | 46.8                 | 73.2                | 96.6               | 47.4               | <b>73.4</b>        | 93.8                | 58.5                  | 70.4                |
|                        | 0.1     | 81.4                 | 96.5                | <b>53.7</b>          | 71.1                | 95.5               | 45.6               | 71.0               | <b>95.3</b>         | 57.9                  | <b>74.1</b>         |
|                        | 0.2     | 81.4                 | 98.6                | 50.9                 | 72.5                | <b>97.7</b>        | <b>56.1</b>        | 71.0               | 94.6                | 60.7                  | <b>74.1</b>         |
| MVH-MIL <sub>emb</sub> | 0.05    | 81.9                 | 97.9                | 48.2                 | 72.5                | 95.5               | 47.4               | 72.8               | 89.1                | 58.5                  | 60.4                |
|                        | 0.075   | 82.4                 | <b>98.6</b>         | <b>49.1</b>          | 73.2                | <b>97.7</b>        | <b>49.1</b>        | <b>74.6</b>        | <b>95.3</b>         | <b>61.2</b>           | 66.7                |
|                        | 0.1     | <b>82.8</b>          | <b>98.6</b>         | <b>49.1</b>          | <b>79.2</b>         | 96.6               | <b>49.1</b>        | <b>74.6</b>        | <b>95.3</b>         | <b>61.2</b>           | 66.7                |
|                        | 0.2     | 80.5                 | <b>98.6</b>         | 48.2                 | 74.5                | 95.5               | <b>49.1</b>        | <b>74.6</b>        | <b>95.3</b>         | 56.8                  | <b>70.4</b>         |
| MVH-MIL <sub>add</sub> | 0.05    | <b>86.4</b>          | 97.9                | 47.7                 | 73.2                | <b>97.7</b>        | <b>56.1</b>        | <b>75.7</b>        | 92.2                | <b>63.4</b>           | 63.0                |
|                        | 0.075   | 83.3                 | <b>99.3</b>         | 50.0                 | 75.2                | <b>97.7</b>        | <b>56.1</b>        | 72.2               | 96.1                | 60.1                  | <b>70.4</b>         |
|                        | 0.1     | 82.4                 | 98.6                | <b>50.9</b>          | <b>77.9</b>         | 96.6               | <b>56.1</b>        | 71.6               | <b>96.9</b>         | 59.0                  | <b>70.4</b>         |
|                        | 0.2     | 81.4                 | 97.9                | 50.0                 | 76.5                | 96.6               | 54.4               | 71.6               | <b>96.9</b>         | 61.7                  | 66.7                |

Table 4. Ablations on different fusion methods with different  $hyp_c$ .

488 (83.3% vs. 80.5%) and TCGA-GBM (97.9% vs. 97.2%),  
 489 it is outperformed by SlideViT in 5 of the 10 tasks. This  
 490 demonstrates that hierarchical modeling can be beneficial  
 491 for specific tasks but is not a universally stable solution on  
 492 its own, underscoring the necessity of combining GHS with  
 493 an adaptive component like MoWE to effectively balance  
 494 hierarchical and Euclidean feature representations.

### 495 5.3.2. Strategy Comparison

496 Different fusion strategies also show distinctive perfor-  
 497 mance (Tab. 4). Specifically, the overall performance of  
 498 MVH-MIL<sub>cat</sub> is moderate, securing 5 best scores across the  
 499 ten experiments, but its unpredictable preference suggests a  
 500 potential incompatibility with the Poincaré ball model. Be-  
 501 sides, MVH-MIL<sub>emb</sub> only achieves one best score, showing  
 502 the poorest performance of the three fusion methods. Cor-  
 503 responding to our analysis, MVH-MIL<sub>add</sub> achieves six best  
 504 scores across the ten experiments. This result demonstrates  
 505 its clear superiority, aligning perfectly with the theoretical  
 506 underpinnings of the Poincaré ball model. Using Möbius  
 507 addition ( $\oplus_c$ ) is the native, geometrically-aware operation  
 508 for vector combination within hyperbolic space. This oper-  
 509 ation is natively compatible with the Poincaré distance met-  
 510 ric used for optimization.

### 511 5.3.3. Generalization Evaluation

512 We further evaluate the generalizability of GHS as a plug-  
 513 and-play module by applying it to seven MIL models  
 514 (Tab. 5). The results demonstrate that GHS can significantly  
 515 boost the performance of several backbones, although its  
 516 impact is task and method dependent. On CAMELYON16,  
 517 GHS provides dramatic performance gains for specific  
 518 methods. For example, it improves CLAM from 73.4%  
 519 to 90.7% (+17.3%) and DSMIL from 86.1% to 94.6%  
 520 (+8.5%). Similarly, on TCGA-Lung, GHS consistently  
 521 enhances performance for most methods, notably lifting  
 522 SlideAve from 75.1% to 78.1% (+3.0%) and ABMIL from  
 523 71.6% to 74.6% (+3.0%). On TCGA-BLCA, where per-  
 524 formance is already high, the improvements are more mod-  
 525 est but still present, as seen with CLAM (95.5% to 97.7%)  
 526 and ABMIL (96.6% to 97.7%). However, we also ob-  
 527 serve cases where GHS does not provide a benefit, or even

| Method   | TCGA-Lung<br>NSCLC |             | TCGA-BLCA<br>Grade |             | CAMELYON16<br>BrMet |             |
|----------|--------------------|-------------|--------------------|-------------|---------------------|-------------|
|          | o/GHS              | w/GHS       | o/GHS              | w/GHS       | o/GHS               | w/GHS       |
| SlideAve | 75.1               | <b>78.1</b> | <b>96.6</b>        | <b>96.6</b> | <b>69.6</b>         | 66.7        |
| SlideMax | 75.7               | <b>77.5</b> | <b>95.6</b>        | 94.3        | <b>96.1</b>         | 75.2        |
| ABMIL    | 71.6               | <b>74.6</b> | 96.6               | <b>97.7</b> | 93.7                | <b>94.6</b> |
| CLAM     | <b>72.2</b>        | <b>72.2</b> | 95.5               | <b>97.7</b> | 73.4                | <b>90.7</b> |
| DSMIL    | 73.4               | <b>76.3</b> | <b>96.6</b>        | <b>96.6</b> | 86.1                | <b>94.6</b> |
| GigaPath | 75.1               | <b>76.3</b> | 97.3               | <b>97.7</b> | <b>95.3</b>         | 89.9        |
| S4MIL    | <b>74.6</b>        | 73.4        | <b>96.6</b>        | <b>96.6</b> | <b>93.8</b>         | 92.0        |

Table 5. Ablations on the mainstream methods with/without GHS. o/GHS is the original method. w/GHS implements GHS.

degrades performance, such as with S4MIL on TCGA-  
 Lung (74.6% to 73.4%) and SlideMax on CAMELYON16  
 (96.1% to 75.2%). This variability suggests that GHS is  
 most impactful for methods and tasks where the original  
 backbone struggles with slide heterogeneity, and its feature-  
 constraining nature may not be universally beneficial for all  
 architectures, requiring a compatible structure design.

## 535 6. Conclusion

WSI analysis confronts three key challenges: suboptimal  
 tissue hierarchy representation, global-local feature trade-  
 off, and tissue heterogeneity, which our MVH-MIL frame-  
 work addresses. For hierarchical modeling, we leverage  
 hyperbolic geometry, its exponential distance growth effi-  
 ciently encodes tree-like pathological structures, reducing  
 Euclidean space information loss. A hybrid geometric em-  
 bedding module balances scales by fusing Euclidean (lo-  
 cal pattern-preserving) and hyperbolic (global hierarchy-  
 capturing) embeddings via weighted addition. To tackle  
 heterogeneity, the MoWE module replaces uniform FFNs  
 with K pathology-specialized experts, using a learnable  
 router to dynamically assign tokens to optimal experts  
 for adaptive computation. MVH-MIL achieves SOTA on  
 7 downstream tasks across 7 datasets, outperforming 10  
 SOTA methods. Furthermore, the proposed GHS module  
 and MoWE module are model agnostic and could poten-  
 tially be adapted to improve other deep learning models,  
 effectively tackling hierarchical and heterogeneous data.

555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610

## References

- [1] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermesen, Quirine F Manson, Maschenka Balkenhol, 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. 5
- [2] Sicheng Chen, Tianyi Zhang, Dankai Liao, Dandan Li, Low Chang Han, Yanqin Jiang, Yueming Jin, and Shangqing Lyu. Pathrwkv: Enabling whole slide prediction with recurrent-transformer. *arXiv preprint arXiv:2503.03199*, 2025. 2
- [3] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, pages 1251–1258, 2017. 5
- [4] Antonio Colaprico, Tiago C Silva, Catharina Olsen, Luciano Garofano, Claudia Cava, Davide Garolini, Thais S Sabedot, Tathiane M Malta, Stefano M Pagnotta, Isabella Castiglioni, et al. Tcgabiobio: an r/bioconductor package for integrative analysis of tcga data. *Nucleic acids research*, 44(8):e71–e71, 2016. 5
- [5] Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shao-han Huang, Wenhui Wang, Nanning Zheng, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint arXiv:2307.02486*, 2023. 5
- [6] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3, 6
- [7] Aleksandr Ermolov, Leyla Mirvakhabova, Valentin Khrukov, Nicu Sebe, and Ivan Oseledets. Hyperbolic vision transformers: Combining improvements in metric learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7409–7419, 2022. 3
- [8] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022. 3
- [9] Ruiwei Feng, Xuechen Liu, Jintai Chen, Danny Z Chen, Honghao Gao, and Jian Wu. A deep learning approach for colonoscopy pathology wsi analysis: accurate segmentation and classification. *IEEE Journal of Biomedical and Health Informatics*, 25(10):3700–3708, 2020. 1
- [10] Leo Fillioux, Joseph Boyd, Maria Vakalopoulou, Paul-Henry Cournède, and Stergios Christodoulidis. Structured state space models for multiple instance learning in digital pathology. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 594–604. Springer, 2023. 5
- [11] Matthew G Hanna, Anil Parwani, and Sahussapont Joseph Sirintrapun. Whole slide imaging: technology and applications. *Advances in Anatomic Pathology*, 27(4):251–259, 2020. 1
- [12] Petr Holub, Heimo Müller, Tomáš Běl, Luca Pireddu, Markus Plass, Fabian Prasser, Irene Schlünder, Kurt Zatloukal, Rudolf Nenutil, and Tomáš Brázdil. Privacy risks of whole-slide image sharing in digital pathology. *Nature Communications*, 14(1):2577, 2023. 1
- [13] Peixiang Huang, Yanyan Huang, Weiqin Zhao, Junjun He, and Lequan Yu. Hyperpath: Knowledge-guided hyperbolic semantic hierarchy modeling for wsi analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 262–272. Springer, 2025. 2, 3
- [14] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International conference on machine learning*, pages 2127–2136. PMLR, 2018. 1, 2, 4, 5
- [15] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 3
- [16] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [17] Hyeonjun Kwon, Jinhyun Jang, Jin Kim, Kwonyoung Kim, and Kwanghoon Sohn. Improving visual recognition with hyperbolic visual hierarchy mapping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17364–17374, 2024. 2, 3
- [18] Dmitry Lepikhin, HyounJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020. 3
- [19] Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14318–14328, 2021. 2, 4, 5
- [20] Dankai Liao, Sicheng Chen, Nuwa Xi, Qiaochu Xue, Jieyu Li, Lingxuan Hou, Zeyu Liu, Chang Han Low, Yufeng Wu, Yiling Liu, et al. Unpuzzle: A unified framework for pathology image analysis. *arXiv preprint arXiv:2503.03152*, 2025. 1, 5
- [21] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*, 5(6):555–570, 2021. 2, 4, 5
- [22] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guillaume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. A visual-language foundation model for computational pathology. *Nature medicine*, 30(3):863–874, 2024. 1
- [23] Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. *Advances in neural information processing systems*, 30, 2017. 1
- [24] Linhao Qu, Manning Wang, Zhijian Song, et al. Bi-directional weakly supervised knowledge distillation for whole slide image classification. *Advances in Neural Information Processing Systems*, 35:15368–15381, 2022. 1

- 668 [25] Sameera Ramasinghe, Violetta Shevchenko, Gil Avraham,  
669 and Ajanthan Thalaisyasingam. Accept the modality gap:  
670 An exploration in the hyperbolic space. In *Proceedings of*  
671 *the IEEE/CVF Conference on Computer Vision and Pattern*  
672 *Recognition*, pages 27263–27272, 2024. 2, 3
- 673 [26] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim  
674 Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel  
675 Keysers, and Neil Houlsby. Scaling vision with sparse mix-  
676 ture of experts. *Advances in Neural Information Processing*  
677 *Systems*, 34:8583–8595, 2021. 3
- 678 [27] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian  
679 Zhang, Xiangyang Ji, et al. Transmil: Transformer based  
680 correlated multiple instance learning for whole slide image  
681 classification. *Advances in neural information processing*  
682 *systems*, 34:2136–2147, 2021. 2, 3, 5
- 683 [28] Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing  
684 Zhang, Junzhou Huang, Wei Yang, and Xiao Han. Transpath:  
685 Transformer-based self-supervised learning for histopatho-  
686 logical image classification. In *International Conference on*  
687 *Medical Image Computing and Computer-Assisted Intervention*,  
688 pages 186–195. Springer, 2021. 2
- 689 [29] Frank Wilcoxon. Individual comparisons by ranking meth-  
690 ods. *Biometrics bulletin*, 1(6):80–83, 1945. 5
- 691 [30] Hanwen Xu, Naoto Usuyama, Jaspreet Bagga, Sheng Zhang,  
692 Rajesh Rao, Tristan Naumann, Cliff Wong, Zelalem Gero,  
693 Javier González, and Yu Gu. A whole-slide foundation  
694 model for digital pathology from real-world data. *Nature*,  
695 630(8015):22, 2024. 1, 2, 3, 4, 5, 6
- 696 [31] Jiexi Yan, Lei Luo, Cheng Deng, and Heng Huang. Un-  
697 supervised hyperbolic metric learning. In *Proceedings of*  
698 *the IEEE/CVF conference on computer vision and pattern*  
699 *recognition*, pages 12465–12474, 2021. 2, 3
- 700 [32] Shu Yang, Yihui Wang, and Hao Chen. Mambamil: En-  
701 hancing long sequence modeling with sequence reordering  
702 in computational pathology. In *International conference on*  
703 *medical image computing and computer-assisted interven-*  
704 *tion*, pages 296–306. Springer, 2024. 2, 5
- 705 [33] Jiawen Yao, Xinliang Zhu, Jitendra Jonnagaddala, Nicholas  
706 Hawkins, and Junzhou Huang. Whole slide images based  
707 cancer survival prediction using attention guided deep multi-  
708 ple instance learning networks. *Medical image analysis*, 65:  
709 101789, 2020. 1
- 710 [34] Hongrun Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao,  
711 Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. Dtf-  
712 mil: Double-tier feature distillation multiple instance learn-  
713 ing for histopathology whole slide image classification. In  
714 *Proceedings of the IEEE/CVF conference on computer vi-*  
715 *sion and pattern recognition*, pages 18802–18812, 2022. 2