

TeamPath: Building MultiModal Pathology Experts with Reasoning AI Copilots

Anonymous ACL submission

Abstract

Advances in AI have introduced several strong models in computational pathology to usher it into the era of multi-modal diagnosis, analysis, and interpretation. However, the current pathology-specific visual language models still lack capacities in making diagnosis with rigorous reasoning paths as well as handling divergent tasks, and thus challenges of building AI Copilots for real scenarios still exist. Here we introduce TeamPath, an AI system powered by reinforcement learning and router-enhanced solutions based on large-scale histopathology multimodal datasets, to work as a virtual assistant for expert-level disease diagnosis, patch-level information summarization, and cross-modality generation to integrate transcriptomic information for the clinical usage. We also collaborate with pathologists from Yale School of Medicine to demonstrate that TeamPath can assist them in working more efficiently by identifying and correcting expert conclusions and reasoning paths. Overall, TeamPath can flexibly choose the best settings according to the needs, and serve as an innovative and reliable system for information communication across different modalities and experts.

1 Introduction

Pathological diagnosis is a complex yet essential component of clinical decision-making. In this manuscript, we present TeamPath, a framework that augments VLMs with multi-modal reasoning and a task-sensitive routing mechanism, enabling robust performance in several pathology-related tasks. Our approach begins with the careful selection of base models and the design of medical-specific prompts to curate high-quality, reasoning-enriched training data. Through comprehensive analyses, we demonstrate both the necessity of equipping VLMs with reasoning capabilities to address complex pathology tasks and the importance of constructing high-quality datasets for model success.

We further showcase the effectiveness of TeamPath across diverse downstream applications, including multi-modal pathology visual question answering (Pathology VQA) and caption summarization. By leveraging an LLM-driven router, TeamPath dynamically selects the most suitable strategy to meet task requirements, functioning as a reliable and adaptive system. Importantly, we invite pathologists to evaluate the model’s reasoning pathways, thereby validating its practical utility as a medical assistant. Finally, we introduce a new task, known as spatial transcriptomic profiles generation, to assess the cross-modality generative ability of TeamPath. Overall, TeamPath provides a new avenue for integrative analyses that combine molecular and histopathological signatures.

2 Methods

Problem definition. In this manuscript, we aim to construct a pathology-expert-level visual language model $\mathcal{M}()$ which accepts text prompts T and pathology image P as inputs. The outputs of our model follow the instructions and information provided in T and P . To train $\mathcal{M}()$, we collect a dataset $D_p = \{(T_1, P_1), \dots, (T_n, P_n)\}_1^n$ with n items for training, and transfer the trained model to various downstream applications.

Constructing TeamPath as a system. To enhance our system’s multitasking capabilities, we adopted a method commonly used in current basic model development, namely training a language model-based router ($\mathcal{R}()$) according to tasks and requirements. This router accepts questions as input data and outputs the model it selects to solve specific problems. The advantage of this design is to unify the TeamPath as a system for various downstream applications in digital pathology, and select the solution that best meets needs to save costs and improve model capabilities. We mark the best solution settings (one of the following choices:

Reinforcement Learning (RL) (Sheng et al., 2024), Supervised FineTuning (SFT) (Zheng et al., 2024), and test-time verification and correction (TTVC) inspired by test-time scaling (TTS) (Snell et al., 2025)) of each question, and train \mathcal{R} with questions and choices. Here, RL is used for solving questions that require reasoning, and SFT is used for summarization and cross-modality generation. Since AI Copilot needs interactions with physicians, TTVC is used for tasks requiring human-AI collaboration. Current base model of TeamPath is Patho-R1-7B (Zhang et al., 2025), which is selected after carefully comparing it with different LMMs such as Qwen2.5VL-7B (Team, 2024), Qwen2.5VL-3B (Team, 2024), MedVLThinker-7B (Huang et al., 2025), PathGen-LLaVA-13B (Sun et al.), InternVL3-8B (Zhu et al., 2025), and MedGemma-4B (Selligren et al., 2025).

Empowering TeamPath with reasoning capabilities. To enhance the reasoning capabilities of TeamPath for complex pathological analysis, we adopt Group Relative Policy Optimization (GRPO), an efficient reinforcement learning algorithm that forgoes the critic model used in traditional PPO (Shao et al., 2024).

For each pathological query q , GRPO samples a group of G outputs $\{o_1, o_2, \dots, o_G\}$ from the current policy π_θ and optimizes the GRPO objective, where the key novelty lies in the group-relative advantage estimation:

$$\hat{A}_{i,t} = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}. \quad (1)$$

Here, $\mathbf{r} = \{r_1, r_2, \dots, r_G\}$ represents the reward scores for all outputs in the group, obtained from a reward model trained on the quality of pathological reasoning. This group-relative formulation eliminates the need for a separate value function V_ψ required in PPO, significantly reducing computational overhead while maintaining training stability.

For GRPO, the reward of question i is:

$$r_i = r(\hat{y}_i, y_i) = \begin{cases} 1, & \text{is_equivalent}(\hat{y}_i, y_i) \\ 0, & \text{otherwise} \end{cases}, \quad (2)$$

where \hat{y} and y represent the model outputs and observed answers, respectively. `is_equivalent()` is a function used to determine if the answer is correct or not.

In our ablation studies, we also consider introducing open-ended questions to model training; in

that case, we utilize the BLEU score as a reward. The reward for closed-ended samples is the same, but for open-ended sample j , the reward is:

$$r_j = r(\hat{y}_j, y_j) = \text{BLEU}(\hat{y}_j, y_j). \quad (3)$$

The comparative nature of this approach aligns naturally with pathological diagnosis workflows, where medical experts simultaneously evaluate multiple diagnostic hypotheses. By learning from the relative quality of responses within each group, TeamPath develops more nuanced reasoning capabilities for tasks requiring differential diagnosis, evidence synthesis, and step-by-step pathological analysis.

In our ablation studies, we also consider Dynamic sAmpling Policy Optimization (DAPO) as an alternative reinforcement learning algorithm. DAPO removes the KL divergence and adjusts the group-level normalization method.

3 Results

Dataset and Method Overview. The curation of high-quality datasets is increasingly critical for advancing PFMs and VLMs, particularly in the era of multimodal reasoning and summarization. At the same time, careful attention must be paid to preventing data leakage to ensure unbiased evaluation of model performance. To this end, and leveraging prior data collection strategies, we distilled a subset of data from PathGen-1.6M (Sun et al.), which is a large-scale resource comprising nearly 10,000 WSIs and 1.6 million ROIs derived from TCGA data (Weinstein et al., 2013), for the usage in the finetuning stage with reinforcement learning. Reasoning data were constructed using COT templates generated based on the advanced reasoning model o4-mini (OpenAI, 2025), with subsequent quality validation performed by pathologists at Yale School of Medicine. Importantly, this dataset does not overlap with the benchmark testing set used for Pathology VQA evaluation (He et al., 2020b), namely PathMMU (Sun et al., 2024), which contains ROIs paired with questions across five diagnostic categories and represents one of the most advanced evaluation sources. In addition, another subset distilled from PathGen-1.6M was curated as the testing dataset for the ROI summarization task. To assess performance in cross-modality generation, we leveraged HEST-1K (Jaume et al., 2024) and STImage1K4M (Chen et al., 2024), two multi-omic histopathology collections to assess the

prediction of transcriptomic profiles as molecular signatures from ROIs. These two datasets are used to construct training, validation, and test sets. The overall data preprocessing workflow and sample sizes are summarized in Extended Data Figure 1.

The overall process of dataset curation and model training is summarized in Figures 1 (a)-(d). TeamPath emerges as a robust multimodal AI assistant for both disease analysis and modality generation. To refine its reasoning capabilities, we employ Group Relative Policy Optimization (GRPO) (Guo et al., 2025) to finetune the base model (the default setting is Patho-R1-7B), thereby enhancing its ability to perform reasoning over pathology images. With this capacity for structured reasoning, TeamPath demonstrates strong performance in addressing Pathology VQA tasks, as shown in our comprehensive benchmarking analysis. Importantly, the model also maintains high performance on tasks where reasoning is less critical, such as image summarization (known as caption generation) and cross-modality generation. This adaptability enables TeamPath to support task-specific optimization through either reinforcement learning or supervised finetuning. In collaboration with expert pathologists, we further demonstrate that TeamPath can function as a clinical copilot, assisting in applications such as correcting inaccurate conclusions and identifying flawed reasoning paths. Taken together, TeamPath advances both biomedical research and clinical practice in histopathology analysis. Finally, a comparative summary of task- and metric-specific rankings, shown in Figure 1 (e), demonstrates the superior performance of TeamPath across multiple dimensions.

TeamPath improves the performance of ROI-level assessment with reasoning ability. The increasing complexity of histopathology image analysis presents significant challenges for developing expert-level VLMs. One particularly demanding setting is Pathology VQA, which requires models to correctly respond to questions grounded in histopathology images. Unlike traditional classification tasks (e.g., disease-state classification or cancer cell identification), Pathology VQA involves a broader and more complex range of scenarios (He et al., 2020a) and demands higher accuracy in answer production. To evaluate model performances under this setting, we employ the recently published PathMMU dataset, which includes VQA pairs spanning five categories, ranging from expert-annotated questions to images from

social media. Importantly, PathMMU is excluded from the training data of all evaluated models to ensure fairness. Reflecting the real-world requirements faced by pathologists, we emphasize the need for high-quality, fine-grained answers that integrate multimodal information and contribute meaningfully at the clinical level. Our baseline comparisons encompass (1) general-domain VLMs, including o4-mini, GPT-4o (Hurst et al., 2024), Qwen2.5VL-3B, Qwen2.5VL-7B (Team, 2024), and InternVL3-8B (Zhu et al., 2025); (2) medical-domain VLMs, including MedGemma-4B (Sellergren et al., 2025) and MedVLThinker-7B (Huang et al., 2025); (3) pathology-specific VLMs, including PathGen-LLaVA-13B (Sun et al.) and Patho-R1-7B (Zhang et al., 2025); and (4) a random-answer baseline. Model performance is assessed by computing accuracy relative to expert-generated answers within PathMMU, enabling a rigorous and fair benchmarking analysis. The comparison of sample size used for training and testing is shown in Extended Data Table 1, and we can see that the amount of images used in testing is large enough to support a general conclusion.

Figures 2 (a)-(c) show our benchmarking results across different categories, including PubMed, SocialPath, Atlas, EduContent, and PathCLS. PathMMU also pre-defines different sample types, and “overall” represents all testing samples in the selected category, “tiny_test” represents testing samples used for expert evaluation, and “test” represents the rest of the samples. We find that TeamPath outperforms all other baseline models, including domain-expert VLMs with similar or larger parameter size, such as Patho-R1-7B and PathGen-LLaVA-13B, in nearly all evaluations. TeamPath also performs better than strong general VLMs, such as o4-mini and GPT-4o, further demonstrating the strength of expert models in addressing medical challenges. Moreover, o4-mini and GPT-4o still perform better than most of the selected baselines, indicating they possess a certain level of understanding of pathological knowledge. Other general VLMs and medical VLMs performed poorly in this task. We further visualize the comprehensive benchmarking analysis, including ranking and accuracy of each method with all samples in Figure 2 (d), which shows that TeamPath also has the lowest rank by considering all categories jointly. Therefore, our experiment results show that introducing reasoning capacities to build pathology-expert VLMs can enhance their ability in making diagnoses, and thus

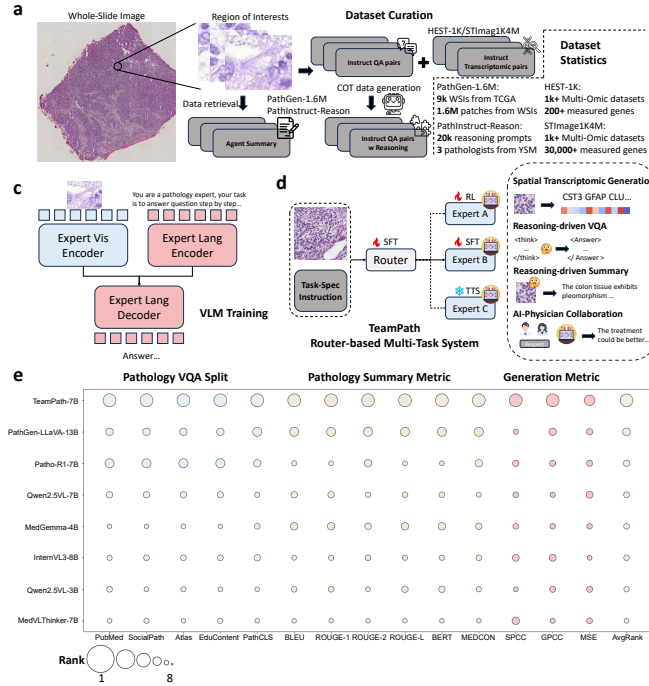


Figure 1: Landscape of TeamPath (a) Steps of dataset curation. We extract image-text pairs from a processed TCGA dataset (PathGen-1.6M). (b) Word cloud visualization of ROI captions (upper) and questions (bottom). (c) The core visual language model architecture of TeamPath. (d) TeamPath as a system with an LLM-enhanced router (with over 80% accuracy in choosing the correct approach) and the corresponding capacities in various downstream applications. The logo fire means that we need to adjust the parameters of models, and the logo snowflake means that we do not change the parameters. (e) Overall ranking list of different methods across tasks and metrics.

TeamPath can serve as a strong performer for the key feature identification and content understanding of ROIs.

To examine the performance of TeamPath specific in disease diagnosis, we utilize GPT-5.2 to extract the type of questions in PathMMU and select diagnosis-related questions to make comparison. According to Extended Data Figures 2 (a)-(b), TeamPath performs better than the second-best baseline Patho-R1-7B as well as random guessing in handling disease-diagnosis-related quires, and the improvement is consistent across most of the disease categories in both tiny group and large group. Therefore, TeamPath can also improve diagnostic accuracy after training.

To obtain a more intuitive understanding of the key contributions of TeamPath following reinforcement learning training, we selected two case studies where TeamPath provided the correct answer while other models failed to make accurate judgments.

Figure 3 highlights the importance of precise morphological criteria in recognizing lipoblasts. While several models incorrectly selected option C, describing large, clear vacuoles displacing the nucleus to the periphery, a hallmark of mature

adipocytes. We found that TeamPath correctly identified option B as the defining feature of lipoblasts (Hisaoka, 2014). Lipoblasts are diagnostically recognized by the presence of moderately sized cytoplasmic fat vacuoles that indent or scallop the nucleus, a distinction that separates them from both mature adipocytes and other stromal features. By emphasizing nuclear indentation rather than displacement, TeamPath demonstrated accurate pathological reasoning aligned with standard diagnostic criteria. This correctness not only underscores the reliability of TeamPath in differentiating subtle histologic features but also highlights the critical nuance needed in distinguishing malignant lipoblastic cells from benign adipocytic processes. Moreover, Extended Data Figure 3 demonstrates that TeamPath correctly identified synaptophysin as the targeted marker in the immunohistochemical stain of section A. The brown, cytoplasmic staining pattern observed is a hallmark of synaptophysin, which is widely used as a marker of neuroendocrine differentiation. While other models misclassified the stain as estrogen receptor or S100 protein, TeamPath distinguished the subtle morphological and staining features that separate synaptophysin from

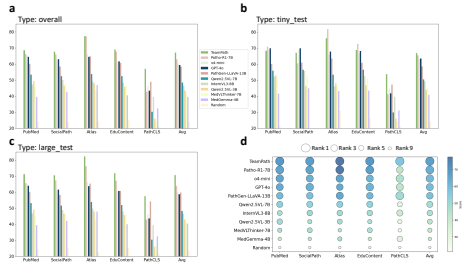


Figure 2: Benchmarking results with PathMMU for the pathology VQA task. We note that since we did not have information about the testing setting of PathGen-LLaVA-13B, we used results reported by the model creators in (Sun et al.). (a) Accuracy across different categories of all selected methods with all samples. (b) Accuracy across different categories of all selected methods with samples from a tiny set. (c) Accuracy across different categories of all selected methods with samples from a large set. (d) Joint visualization with accuracy and ranking information for all selected methods. The darker the bubble color, the higher the model score; The larger the bubble shape, the lower the model ranking.

nuclear markers like estrogen receptor or more diffuse proteins such as S100. This highlights both the accuracy and interpretive strength of TeamPath in immunohistochemistry tasks, particularly in recognizing marker-specific staining patterns and avoiding common pitfalls that lead to misclassification. We also note that previous pathology expert models have obvious shortcomings, such as Patho-R1-7B’s garbled output and PathGen-LLaVA’s lack of interpretable diagnostic outputs. Instead, TeamPath can make correct identification supported by comprehensive explanations, explained in the information provided by the reasoning paths.

We also explored the contributions of different training strategies and highlighted the importance of selecting base models based on a set of ablation studies, discussed in Appendix G and Extended Data Figures 4 (a)-(d), as well as in Appendix H and Extended Data Figure 5 for the data ablation study.

TeamPath acts as a Copilot in the pathologists-AI collaboration system. Beyond demonstrating the capacity of TeamPath in handling VQA sets as a pathology expert, we further explore its potential as an AI-assisted collaborator (Mialon et al., 2023; Liu et al., 2025a). An effective copilot should not only provide accurate responses to user queries but also reduce the effort required to resolve them, thereby saving both time and cost. To this end, we designed an algorithm in TeamPath with test-time verification and correction (TTVC)

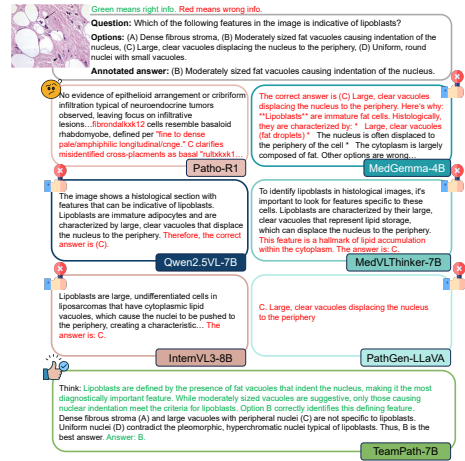


Figure 3: Case study (topic: synaptophysin, which is a precursor cell that develops into an adipocyte (fat cell)) based on the outputs from different models. We highlight the correct information with green text and incorrect information with red text. For the models with errors, we consider two cases. The first case is a wrong answer, and the second case is a confused reasoning path.

(Snell et al., 2025; Liu et al., 2025a) and engaged expert pathologists from Yale School of Medicine (YSM) to collaborate with TeamPath in analyzing histopathology images and generating answers on demand. Specifically, we randomly subsampled 10 question–image pairs from each category within the PathMMU “tiny_test” set and examined two capacities: (1) the ability of TeamPath to act as an auto-verifier or auto-corrector for incorrect expert assessments, and (2) the ability of TeamPath to revise and correct reasoning pathways when human experts fail to provide accurate answers. When performing the TTVC, we prepared the inputs as the images, questions, and reasoning from experts, and performed verification and correction based on TeamPath. The overall paradigm for these two tasks is summarized in Figure 4 (a). Through this study, we aim to establish future paradigms of human–AI collaboration in biomedical research and clinical practice, highlighting the role of TeamPath as a reliable and strong copilot.

We jointly compared the expert-provided results with those corrected by TeamPath and visualized the corresponding accuracies in Figure 4 (b). Our analysis shows that TeamPath significantly improves accuracy across all PathMMU categories (p-value = 0.004), demonstrating that its corrective contribution is consistent and robust regardless of the source of pathology ROIs or questions. Notably,

even in categories where expert performance is relatively low, such as PubMed, TeamPath achieves substantial gains. These improvements demonstrate the effectiveness of TeamPath as a corrector, as reflected by the observed accuracy differences. To further illustrate this capability, we conducted a case study (Figure 4 (c)) in which the expert provided an incorrect answer, whereas TeamPath generated the correct response with an improved reasoning path. In this example, the task involved identifying characteristic nuclear features within the image. The expert’s reasoning correctly accounted for cell size but overlooked nucleolar details, leading to an erroneous conclusion. In contrast, TeamPath integrated multiple features, including nuclear size, shape, staining depth, and prior knowledge of the cancer cell line, to eliminate incorrect options and arrive at the correct decision. Moreover, TeamPath was also able to revise flawed reasoning paths when experts could not provide an answer (e.g., “I do not know”), as shown in Extended Data Figure 6. In summary, through collaborative evaluation with pathologists, we demonstrate the capacity of TeamPath to not only fix erroneous answers but also provide explicit reasoning steps, thereby enhancing both the transparency and interpretability of model-assisted pathology analysis.

Here we provided more detailed analyses for understanding the case that TeamPath’s discrepancies arising in correcting the feedback from different pathologists. We first compared the Pearson Correlation Coefficients (PCCs) among the three pathologists’ accuracy. Extended Data Figure 7 (a) shows that the decision-making of our two pathologists is relatively consistent, while their accuracy rates on specific datasets also differ. Therefore, we can ensure a certain degree of diversity in expert selection used for human-AI collaboration experiment. Moreover, we also compared the number of corrected samples made by TeamPath for different pathologists, and Extended Data Figure 7 (b) shows that the number of corrections made by different experts is relatively similar across most datasets, with the sole exception being the SocialPath dataset. Since this dataset is primarily sourced from social media, the complex data pattern can pose challenges to our designed system. Our experiments demonstrate that TeamPath contributes to improving the performance across different pathologists, thereby exhibiting relatively broad applicability.

We have also performed ablation studies for the

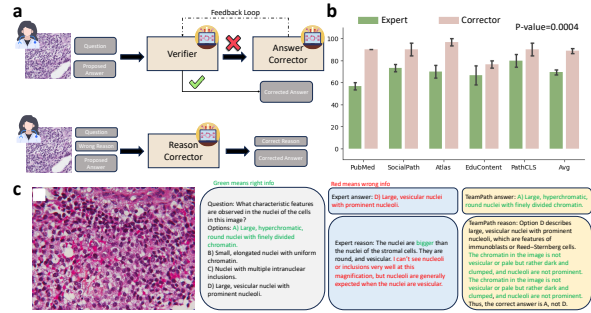


Figure 4: Results of using TeamPath as the answer corrector/reason corrector. TeamPath can work with pathologists together to improve the decision accuracy and provide explainable reasons to support the decision. (a) The illustration of self-verification/correction steps for both answers and reasoning paths. (b) Accuracy before and after correction based on selected samples from PathMMU. We report the average scores and standard deviation across three experts. The test is a one-sided Wilcoxon Rank-sum test. (c) A case study to demonstrate the power of TeamPath as an AI assistant.

verifier with three different choices (using the corrector, o3 (OpenAI, 2025), and o4-mini). Extended Data Figure 8 shows that using o4-mini can achieve the best performance on average, while it can also reduce the cost compared with using o3 or a more advanced model, and thus o4-mini is selected here to perform verification.

TeamPath performs better in summarizing the key information from histopathology images.

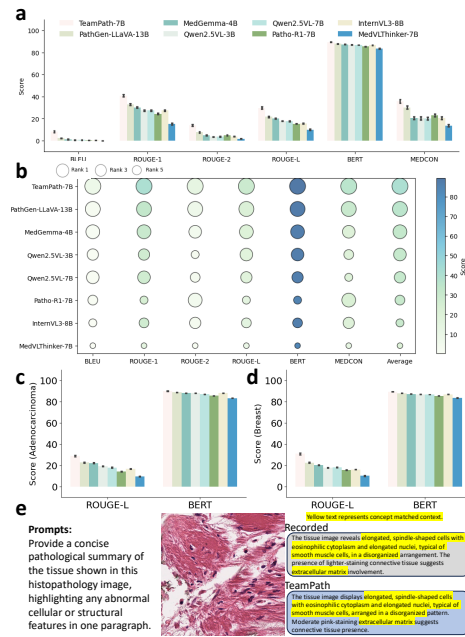
In practical applications, pathology image analysis often extends beyond generating correct answers and reasoning steps to encompass the extraction of important image features for macroscopic or high-level descriptions. To evaluate this capability, we designed experiments aimed at summarizing histopathology information from different ROIs, thereby assessing the capacity of TeamPath to capture and convey high-level image content. For this purpose, we constructed a testing dataset by subsampling 3,000 images and their corresponding captions from PathGen-1.6M. In their setting, the original image caption is enhanced by multi-agent (GPT-4V as backbone) collaboration. These captions were further annotated to include tissue- and disease-state information based on prompting Deepseek-R1 (Guo et al., 2025) to extract the answer from the input caption. To support training, we curated a separate dataset of 50,000 images, ensuring no overlap with the testing set. For benchmarking, we employed the same set of baseline models used in the Pathology VQA experi-

443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473

474 ments. Model performance was evaluated using
 475 multiple similarity metrics between generated sum-
 476 maries and reference captions, including BLEU
 477 (Papineni et al., 2002), ROUGE-1/2/L (Lin, 2004),
 478 BERTScore (Zhang et al.), and MEDCON (Sol-
 479 daini and Goharian, 2016). All metrics were scaled
 480 to a 0–100 range, with higher values indicating
 481 better performance.

482 Figure 5 (a) compares the performance of Team-
 483 Path with other VLMs across all selected metrics on
 484 the testing set. TeamPath consistently outperforms
 485 the baselines across every metric, demonstrating its
 486 strength in generating summaries that align closely
 487 with reference annotations in both content and struc-
 488 ture. To provide a holistic assessment, we further vi-
 489 sualized the aggregated rankings and average scores
 490 of all methods in Figure 5 (b), which highlights
 491 the leading performance of TeamPath across the
 492 joint set of evaluation metrics. Recognizing that
 493 performance may vary by sample source, we also
 494 examined model performance across specific tissue
 495 and disease contexts. We first analyze the produced
 496 categories in the testing dataset and illustrate the
 497 distribution of top 10 categories in Extended Data
 498 Figure 9. Among these labels, we select the top 1
 499 label in each class, including adenocarcinoma and
 500 breast tissue, for further investigation. Figure 5 (c)
 501 reports ROUGE-L and BERT scores for samples
 502 from patients with adenocarcinoma, while Figure 5
 503 (d) shows results for breast tissue samples. In both
 504 cases, TeamPath maintains superior performance
 505 compared with competing baselines. As an illustra-
 506 tive case study, Figure 5 (e) presents an example
 507 output from TeamPath, which accurately captures
 508 key organizational and pathological features—such
 509 as elongated spindle-shaped cells with eosinophilic
 510 cytoplasm and elongated nuclei, characteristic of
 511 smooth muscle cells. By contrast, outputs from
 512 baseline models (Extended Data Figure 10) contain
 513 less precise descriptions and, in some cases, incor-
 514 rect content, further underscoring the advantages
 515 of TeamPath in summarization tasks.

516 Therefore, we conclude that TeamPath demon-
 517 strates as a strong performer in providing the high-
 518 level interpretations with pathology features of as-
 519 signed ROIs.



519 Figure 5: Benchmarking results of the caption summary
 520 task. (a) Performances of different methods for sum-
 521 marizing the caption based on ROI-level information
 522 across all metrics. We report the average scores and
 523 scaled standard deviation ($0.1 \cdot \text{sd}$) with all sam-
 524 ples in the testing set. (b) Joint visualization with
 525 metric scores and ranking information for all selected
 526 methods. The darker the bubble color, the higher the
 527 model score; The larger the bubble shape, the lower
 528 the model ranking. (c) ROUGE-L and BERT scores
 529 based on samples from the selected disease across
 530 all methods. (d) ROUGE-L and BERT scores based
 531 on samples from the selected tissue across all meth-
 532 ods. (e) A case study of caption summary generation
 533 based on TeamPath.

534 **TeamPath introduces new modalities with a**
 535 **cross-modality generation pipeline.** Building on
 536 our previous research and the existing literature,
 537 we observe that current histopathology image anal-
 538 yses primarily rely on textual and visual interpre-
 539 tations. However, given the breadth of biological
 540 signatures that can contribute to disease modeling
 541 and diagnosis, there is a clear opportunity to de-
 542 sign new pipelines that integrate molecular infor-
 543 mation with histopathology features. Such integra-
 544 tion can enable the generation of new modalities
 545 and provide deeper insights into cellular hetero-
 546 geneity, lineage tracing, and disease mechanisms
 547 (Chen et al., 2025b; Song et al., 2024). Therefore,
 548 we finetune TeamPath using paired histopathology
 549 images and transcriptomic profiles generated with
 550 the Visium technology (Genomics), a platform for
 551 spatial transcriptomics (ST). Each ST spot includes

a histopathology image as background and a corresponding gene expression profile. Inspired by Cell2Sentence (Levine et al., 2024) and Loki (Chen et al., 2025b), we convert gene expression profiles into ranked gene lists, ordering genes from highest to lowest expression. The task is then to generate these “spot sentences” and map them back into the transcriptomic space. For training and evaluation, we use two of the largest public datasets: HEST-1K (invasive ductal carcinoma, IDC) and STImage1K4M (brain tissue). HEST-1K includes a broad range of cancer datasets, whereas STImage1K4M contains samples from both disease and normal tissues, thereby enhancing the modeling of ST data. Baseline models for this task include the same VLMs evaluated in the Pathology VQA setting, supplemented with Cell2Sentence-1B. Model performance is assessed using Spot-level Pearson Correlation Coefficient (SPCC), Gene-level Pearson Correlation Coefficient (GPCC), and mean squared error (MSE). For SPCC and GPCC, higher values indicate better performance, whereas lower MSE values reflect higher accuracy.

Figures 6 (a) and (b) demonstrate that TeamPath outperforms all baseline methods when evaluated by both SPCC and MSE across datasets from different sources, underscoring its ability to generate spot-level gene expression profiles that closely resemble measured results. Extended Data Figures 11 (a) and (b) further confirm TeamPath’s better performance in GPCC, highlighting its capacity to preserve gene-level heterogeneity across spatial spots. To examine the impact of base model selection on cross-modality generation, we finetuned Qwen2.5VL-7B for the same task and compared it with TeamPath. As shown in Extended Data Figures 11 (c) and (d), TeamPath, which was built on a pathology-knowledge-enhanced VLM, outperformed the finetuned Qwen-series model. We also emphasize the importance of task-specific finetuning, supported by the clear performance gap between the unadapted base model and TeamPath in generating high-quality expression profiles. UMAP visualizations of the generated profiles (Figures 6 (c) and (d)) show that outputs from TeamPath are more structured and closely aligned with reference profiles compared to those from the base model. This observation is further validated by cluster-level heatmaps of gene expression patterns in brain (Figure 6 (e)) and IDC (Figure 6 (f)) datasets, where TeamPath more accurately recapitulates the biological signal present in the ground

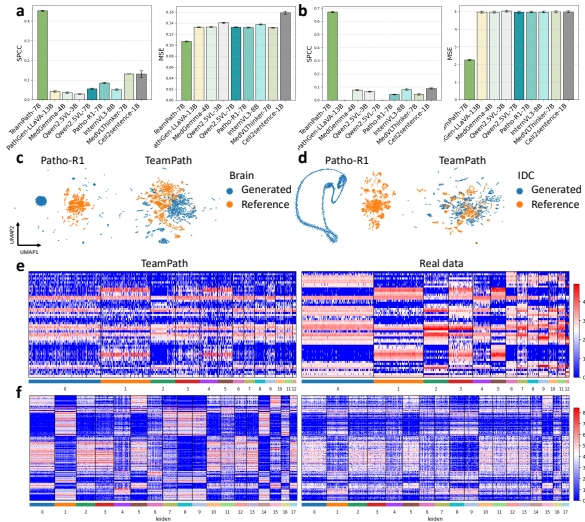


Figure 6: Evaluation of model performances for transcriptomic profile generation. (a) SPCC (higher is better) and MSE (lower is better) scores across different methods for the brain tissue. We report the average scores and scaled standard deviation ($0.1 \cdot sd$) for better visualization. (b) SPCC and MSE scores across different methods for the IDC samples. We report the average scores and scaled standard deviation ($0.1 \cdot sd$) for better visualization. (c) UMAP visualization from the testing set of brain to compare the generated results between Patho-R1 (base) and TeamPath colored by data sources. (d) UMAP visualization from the testing set of IDC to compare the generated results between Patho-R1 (base) and TeamPath colored by data sources. (e) Comparison of expression profiles between generated data and real data based on the brain tissue. (f) Comparison of expression profiles between generated data and real data based on the IDC samples.

truth data. Collectively, these findings demonstrate that the effectiveness of TeamPath in cross-modality generation arises from both the choice of a pathology-informed base model and targeted task-specific finetuning. With these advantages, TeamPath represents a promising approach for generating in-silico or unseen expression profiles directly from histopathology images, thereby providing molecular-level insights into disease phenotypes.

4 Discussion

Here we present TeamPath, an advanced AI copilot to advance research in computational pathology and disease diagnosis by formulating a multi-task AI assistant, which supports various tasks with an automatic router for solution selection. TeamPath demonstrates as a strong method in working with physicians in addressing key pathology challenges.

608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662

References

Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. 2025a. *SFT or RL? an early investigation into training rl-like reasoning large vision-language models*. *Transactions on Machine Learning Research*.

Jiawen Chen, Muqing Zhou, Wenrong Wu, Jinwei Zhang, Yun Li, and Didong Li. 2024. *Stimage-1k4m: A histopathology image-gene expression dataset for spatial transcriptomics*. *Advances in Neural Information Processing Systems*, 37:35796–35823.

Weiqing Chen, Pengzhi Zhang, Tu N Tran, Yiwei Xiao, Shengyu Li, Vrutant V Shah, Hao Cheng, Kristopher W Brannan, Keith Youker, Li Lai, and 1 others. 2025b. *A visual-omics foundation model to bridge histopathology with spatial transcriptomics*. *Nature Methods*, pages 1–15.

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. 2025. *SFT memorizes, RL generalizes: A comparative study of foundation model post-training*. In *Forty-second International Conference on Machine Learning*.

10X Genomics. Visium technology.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. *Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning*. *arXiv preprint arXiv:2501.12948*.

Xuehai He, Zhuo Cai, Wenlan Wei, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. 2020a. *Pathological visual question answering*. *arXiv preprint arXiv:2010.12435*.

Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. 2020b. *Pathvqa: 30000+ questions for medical visual question answering*. *arXiv preprint arXiv:2003.10286*.

Masanori Hisaoka. 2014. *Lipoblast: morphologic features and diagnostic value*. *Journal of UOEH*, 36(2):115–121.

Xiaoke Huang, Juncheng Wu, Hui Liu, Xianfeng Tang, and Yuyin Zhou. 2025. *Medvlthinker: Simple baselines for multimodal medical reasoning*. *arXiv preprint arXiv:2508.02669*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. *Gpt-4o system card*. *arXiv preprint arXiv:2410.21276*.

Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Truong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, Curtis Langlotz, and 1 others. *Radgraph: Extracting clinical entities and relations from radiology reports*.

Guillaume Jaume, Paul Doucet, Andrew H. Song, Ming Y. Lu, Cristina Almagro Pérez, Sophia J Wagner, Anurag Jayant Vaidya, Richard J. Chen, Drew FK Williamson, Ahnong Kim, and Faisal Mahmood. 2024. *HEST-1k: A dataset for spatial transcriptomics and histology image analysis*. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Daniel Levine, Syed A Rizvi, Sacha Lévy, Nazreen Palikkavaliyaveetil, David Zhang, Xingyu Chen, Sina Ghadermarzi, Ruiming Wu, Zihe Zheng, Ivan Vrkic, and 1 others. 2024. *Cell2sentence: Teaching large language models the language of biology*. In *International Conference on Machine Learning*, pages 27299–27325. PMLR.

Chin-Yew Lin. 2004. *Rouge: A package for automatic evaluation of summaries*.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. *Visual instruction tuning*. *Advances in neural information processing systems*, 36:34892–34916.

Tianyu Liu, Simeng Han, Xiao Luo, Hanchen Wang, Pan Lu, Biqing Zhu, Yuge Wang, Keyi Li, Jiapeng Chen, Rihao Qu, and 1 others. 2025a. *Towards artificial intelligence research assistant for expert-involved learning*. *arXiv preprint arXiv:2505.04638*.

Zihe Liu, Jiashun Liu, Yancheng He, Weixun Wang, Jiaheng Liu, Ling Pan, Xinyu Hu, Shaopan Xiong, Ju Huang, Jian Hu, and 1 others. 2025b. *Part i: Tricks or traps? a deep dive into rl for llm reasoning*. *arXiv preprint arXiv:2508.08221*.

Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Melissa Zhao, Aaron K Chow, Kenji Ikemura, Ahnong Kim, Dimitra Pouli, Ankush Patel, and 1 others. 2024. *A multimodal generative ai copilot for human pathology*. *Nature*, 634(8033):466–473.

Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. *Gaia: a benchmark for general ai assistants*. In *The Twelfth International Conference on Learning Representations*.

OpenAI. 2025. *Openai o3 and o4-mini system card*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a method for automatic evaluation of machine translation*. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, and 1 others. 2011. *Scikit-learn: Machine learning in python*. *the Journal of machine Learning research*, 12:2825–2830.

663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716

717	Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri,	Jonathan Bright, and 1 others. 2020. Scipy 1.0: funda-	773
718	Atila Kiraly, Madeleine Traverse, Timo Kohlberger,	mental algorithms for scientific computing in python.	774
719	Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau,	<i>Nature methods</i> , 17(3):261–272.	775
720	and 1 others. 2025. Medgemma technical report.		
721	<i>arXiv preprint arXiv:2507.05201</i> .		
722	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Weixun Wang, Shaopan Xiong, Gengru Chen, Wei Gao,	776
723	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Sheng Guo, Yancheng He, Ju Huang, Jiaheng Liu,	777
724	Zhang, YK Li, Yang Wu, and 1 others. 2024.	Zhendong Li, Xiaoyang Li, and 1 others. 2025. Rein-	778
725	Deepseekmath: Pushing the limits of mathematical	forcement learning optimization for large-scale learn-	779
726	reasoning in open language models. <i>arXiv preprint</i>	ing: An efficient and user-friendly scaling library.	780
727	<i>arXiv:2402.03300</i> .	<i>arXiv preprint arXiv:2506.06122</i> .	781
728	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	John N Weinstein, Eric A Collisson, Gordon B Mills,	782
729	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	Kenna R Shaw, Brad A Ozenberger, Kyle Ellrott, Ilya	783
730	Lin, and Chuan Wu. 2024. Hybridflow: A flexible	Shmulevich, Chris Sander, and Joshua M Stuart. 2013.	784
731	and efficient rlhf framework. <i>arXiv preprint arXiv:</i>	The cancer genome atlas pan-cancer analysis project.	785
732	<i>2409.19256</i> .	<i>Nature genetics</i> , 45(10):1113–1120.	786
733	Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Avi-	Wen-wai Yim, Yujuan Fu, Asma Ben Abacha, Neal	787
734	riral Kumar. 2025. Scaling LLM test-time compute	Snider, Thomas Lin, and Meliha Yetisgen. 2023. Acibench: a novel ambient clinical intelligence dataset	788
735	optimally can be more effective than scaling param-	for benchmarking automatic visit note generation.	789
736	eters for reasoning . In <i>The Thirteenth International</i>	<i>Scientific Data</i> , 10(1):586.	790
737	<i>Conference on Learning Representations</i> .		791
738	Luca Soldaini and Nazli Goharian. 2016. Quickumls:	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xi-	792
739	a fast, unsupervised approach for medical concept	aochen Zuo, YuYue, Weinan Dai, Tiantian Fan, Gao-	793
740	extraction.	hong Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin	794
741	Andrew H Song, Mane Williams, Drew FK Williamson,	Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan	795
742	Sarah SL Chow, Guillaume Jaume, Gan Gao, An-	Tong, Chi Zhang, Mofan Zhang, and 17 others. 2025.	796
743	drew Zhang, Bowen Chen, Alexander S Baras, Robert	DAPO: An open-source LLM reinforcement learning	797
744	Serafin, and 1 others. 2024. Analysis of 3d pathol-	system at scale . In <i>The Thirty-ninth Annual Confer-</i>	798
745	ogy samples using weakly supervised ai. <i>Cell</i> ,	<i>ence on Neural Information Processing Systems</i> .	799
746	187(10):2502–2520.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-	800
747	Yuxuan Sun, Hao Wu, Chenglu Zhu, Sunyi Zheng, Qizi	berger, and Yoav Artzi. Bertscore: Evaluating text	801
748	Chen, Kai Zhang, Yunlong Zhang, Dan Wan, Xiaox-	generation with bert.	802
749	iao Lan, Mengyue Zheng, and 1 others. 2024. Path-	Wenchuan Zhang, Penghao Zhang, Jingru Guo, Tao	803
750	mmu: A massive multimodal expert-level benchmark	Cheng, Jie Chen, Shuwan Zhang, Zhang Zhang,	804
751	for understanding and reasoning in pathology. In <i>Euro-</i>	Yuhao Yi, and Hong Bu. 2025. Patho-r1: A multi-	805
752	<i>pean Conference on Computer Vision</i> , pages 56–73.	modal reinforcement learning-based pathology expert	806
753	Springer.	reasoner. <i>arXiv preprint arXiv:2505.11404</i> .	807
754	Yuxuan Sun, Yunlong Zhang, Yixuan Si, Chenglu Zhu,	Wenhao Zhang, Yuexiang Xie, Yuchang Sun, Yanxi	808
755	Kai Zhang, Zhongyi Shui, Jingxiong Li, Xuan Gong,	Chen, Guoyin Wang, Yaliang Li, Bolin Ding, and Jin-	809
756	XINHENG LYU, Tao Lin, and 1 others. Pathgen-1.6	gren Zhou. 2026. On-policy RL meets off-policy	810
757	m: 1.6 million pathology image-text pairs generation	experts: Harmonizing supervised fine-tuning and	811
758	through multi-agent collaboration. In <i>The Thirteenth</i>	reinforcement learning via dynamic weighting . In	812
759	<i>International Conference on Learning Representa-</i>	<i>The Fourteenth International Conference on Learn-</i>	813
760	<i>tions</i> .	<i>ing Representations</i> .	814
761	Qwen Team. 2024. Qwen2 technical report. <i>arXiv</i>	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	815
762	<i>preprint arXiv:2407.10671</i> .	Ye, and Zheyang Luo. 2024. LlamaFactory: Unified	816
763	Dave Van Veen, Cara Van Uden, Louis Blanke-	efficient fine-tuning of 100+ language models . In	817
764	meier, Jean-Benoit Delbrouck, Asad Aali, Christian	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	818
765	Bluethgen, Anuj Pareek, Malgorzata Polacin, Edu-	<i>sociation for Computational Linguistics (Volume 3:</i>	819
766	ardo Pontes Reis, Anna Seehofnerová, and 1 others.	<i>System Demonstrations)</i> , pages 400–410, Bangkok,	820
767	2024. Adapted large language models can outper-	Thailand. Association for Computational Linguistics.	821
768	form medical experts in clinical text summarization.	Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu,	822
769	<i>Nature medicine</i> , 30(4):1134–1142.	Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan,	823
770	Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt	Weijie Su, Jie Shao, and 1 others. 2025. Internvl3:	824
771	Haberland, Tyler Reddy, David Cournapeau, Ev-	Exploring advanced training and test-time recipes	825
772	geni Burovski, Pearu Peterson, Warren Weckesser,	for open-source multimodal models. <i>arXiv preprint</i>	826
		<i>arXiv:2504.10479</i> .	827

A Method Details

Ablation studies of training framework. To demonstrate the efficiency and optimization of our training framework for this system, we have considered several training strategies, including 1. Supervised finetuning (SFT), which collects paired data with images and queries as inputs and answers as outputs; 2. Reinforcement Learning (RL), which utilizes the same input and output data, but we train TeamPath with GRPO (Shao et al., 2024) or DAPO (Yu et al., 2025) to tackle the reasoning capacity of selected base models; 3. SFT+RL, which utilizes the paired data with images, queries, and reasoning paths as inputs and answers as outputs, to train TeamPath with SFT and then with RL. The first step of SFT training ensures the model acquires knowledge in relevant pathological domains, while the second step of RL training enhances the model’s generalization capabilities. Our base models used for ablation studies include Qwen2.5VL-7B and Patho-R1-7B.

Using TeamPath as an AI Copilot to help pathologists. To formalize our method as an AI Assistant, we consider two case studies inspired by the Path VQA experiments with pathologists. We invite the pathologists to answer 50 questions extracted from PathMMU from the five categories, and record their answers as well as reasoning steps. Our first case is a verifier-corrector pipeline, which can detect the incorrect answers made by pathologists and generate the correct answers. Our pipeline utilizes one verifier (a VLM, default as o4-mini) to verify whether the answers and questions proposed by pathologists are correct or not. If it is justified as wrong, we will call the corrector (also a VLM, default as TeamPath used for Pathology VQA) to fix it. Otherwise, the correct answer will be returned. We have a specific threshold to limit the number of epochs in this loop. Our algorithm is summarized in Algorithm 1. We define the success of a self-verification/correction system as follows: if the expert answer is correct, or the expert answer is wrong but the answer produced by this system is correct.

The second case is a reasoning-correction pipeline. Here we have a reasoning corrector, which takes the wrong answers and reasoning paths from pathologists, and generates the correct reasoning path with the correct answer. To improve efficiency, we consider an offline version, that is, by collecting pathologists’ reasons and answers first, we then uti-

lize TeamPath to check or correct these questions. This pipeline can detect the wrong information provided in the reasoning process and generate the correct thinking steps. These two components focus on different aspects and work together as a prototype for building an AI Copilot that can work with pathologists and be deployed in the medical system. We have provided an example of correction in the main text. Our prompts used in these two pipelines are summarized in Appendix F.

Algorithm 1 Verifier-Corrector Pipeline in TeamPath.

Input: Question Q_M , pathology image I_S , human answer O_A , reasoning path O_R , verify prompt T_V , correct prompt T_C , number of iteration N .

Helper Models: Verifier \mathcal{M}_v (An advanced LMM, such as O4-mini), corrector \mathcal{M}_c (An pathology-specific LMM, such as TeamPath with RL finetuning), concatenation function $\cdot||\cdot$.

Output: Corrected outputs O_C

```
1: INIT: initialize all parameters.
2: if  $\mathcal{M}_v(T_V, Q_M||O_R||O_A, I_S)$  is True then
3:    $O_C = O_A$ 
4:   Return  $O_C$ 
5: end if
6: for  $i$  in  $N$  steps do
7:    $O_i, R_i = \mathcal{M}_c(T_C, Q_M||O_R||O_A, I_S)$ 
8:   if  $\mathcal{M}_v(T_V, Q_M||R_i||O_i, I_S)$  is True then
9:      $O_C = O_i$ 
10:    Return  $O_C$ 
11:  else
12:     $O_R = R_i$ 
13:     $O_A = O_i$ 
14:  end if
15: end for
16:  $O_C = O_A$ 
17: Return  $O_C$ 
```

Adapting TeamPath for image summarization and cross-modality generation. To summarize the concepts in histopathology images and further generate image caption, we finetune our base model based on the paired image-caption dataset with the corresponding training set, and we also prepare 3000 samples which are only used for testing, and we utilize Deepseek-R1 (Guo et al., 2025) to extract the disease state and tissue source of the testing samples based on their captions. Our finetuning step follows the setting in Instruction-Tuning imple-

900 mented in Llama-factory. We construct 10 different
901 prompts to ask TeamPath for generating the image
902 captions to reduce the bias of prompt information
903 in the training process.

904 To perform the cross-modality generation task,
905 we also finetune our base model based on the
906 paired image-transcriptomic profile dataset with
907 the corresponding training datasets. We select se-
908 quence data that comes from different batches and
909 resources, but the same tissue/disease, to build a
910 testing dataset. To transfer the information in gene
911 expression space to text space, we first rank the
912 genes of each spot based on their expression pro-
913 files and select the top 100 genes to formulate them
914 in natural language. We then train a linear regres-
915 sor that takes the natural language information as
916 inputs and original gene expression profiles as out-
917 puts based on the training dataset, which finally
918 gives us a method to decode the language informa-
919 tion back to gene expression levels. We finetune
920 our base model with the same approach used in
921 image summarization and also construct 10 differ-
922 ent prompts to ask TeamPath for generating gene
923 expression profiles.

924 The prompts used in this section can also be
925 found in Appendix F.

926 **Evaluations.** In this manuscript, we consider
927 task-specific evaluation (Pedregosa et al., 2011;
928 Virtanen et al., 2020) and follow the settings from
929 previous works with shared tasks.

930 For the evaluation of Path VQA and Human-AI
931 collaboration tasks, we utilize accuracy as a metric.
932 The generated answer should be precisely matched
933 with the provided answer. A higher accuracy repre-
934 sents a better method.

935 For the evaluation of the image caption summa-
936 rization task, we utilize several metrics that can
937 measure the similarity between the generated text
938 and the provided text. These metrics include BLEU,
939 ROUGE-1, ROUGE-2, ROUGE-L, BERT score,
940 and MEDCON (Papineni et al., 2002; Lin, 2004;
941 Zhang et al.; Jain et al.; Yim et al., 2023), supported
942 by a recent publication (Van Veen et al., 2024). We
943 also consider the average score across these metrics.
944 Here are the descriptions:

- 945 • BLEU: The BiLingual Evaluation Understudy
946 (BLEU) score evaluates the quality of gener-
947 ated text by breaking both the generated out-
948 put and the reference text into n-grams, then
949 comparing the overlap between the two sets.
950 The score ranges from 0 to 1 and is typically

951 scaled to a range of 0 to 100, with higher val-
952 ues indicating better model performance.

- 953 • ROUGE: The Recall-Oriented Understudy for
954 Gisting Evaluation (ROUGE) score assesses
955 text quality by computing the F1 score from n-
956 gram overlaps between the generated text and
957 the reference text. In this framework, n-grams
958 from the generated text are treated as predic-
959 tions, while those from the reference text serve
960 as labels. Precision, recall, and the F1 score
961 are calculated using the counts of matching n-
962 grams and their lengths. ROUGE-1 measures
963 unigram overlap, ROUGE-2 measures bigram
964 overlap, and ROUGE-L measures the longest
965 common subsequence. The score ranges from
966 0 to 1 and is typically scaled to a range of 0
967 to 100, with higher values indicating better
968 model performance.
- 969 • BERT: The Bidirectional Encoder Represent-
970 ations from Transformers (BERT) model is
971 pre-trained on large-scale text corpora for lan-
972 guage understanding and excels at producing
973 rich text representations. The BERTScore met-
974 ric leverages this capability by measuring the
975 similarity between embeddings of the gener-
976 ated text and the reference text. The score
977 ranges from 0 to 1 and is typically scaled to a
978 range of 0 to 100, with higher values indicat-
979 ing better model performance.
- 980 • MEDCON: MEDCON limits the recognized
981 concepts and entities to the semantic groups
982 defined in QuickUMLS (Soldaini and Gohar-
983 ian, 2016), including Anatomy, Chemicals,
984 Drugs, Device, Disorders, Genes, Molecu-
985 lar Sequences, Phenomena, and Physiology.
986 These concepts are extracted from both the
987 generated text and the reference text, and the
988 F1 score is calculated based on the overlap
989 between the two sets. The score ranges from
990 0 to 1 and is typically scaled to a range of 0
991 to 100, with higher values indicating better
992 model performance.

993 A higher score of these metrics represents a bet-
994 ter method.

995 For the evaluation of the cross-modality genera-
996 tion task, we consider spot-level Pearson Correla-
997 tion Coefficient (SPCC), gene-level PCC (GPCC),
998 and Mean Squared Error (MSE) as metrics. Higher
999 SPCC and GPCC scores represent a better method,

1000 while a lower MSE score represents a better method.
1001 These metrics are computed between generated
1002 gene expression profiles and observed gene expres-
1003 sion profiles.

1004 **Baselines.** Our baseline methods cover current
1005 state-of-the-art (SOTA) open-source LMMs based
1006 on the open source movement in scientific research
1007 and the powerful influence of open source mod-
1008 els. Moreover, there are a few powerful closed-
1009 source models focusing on digital histopathology.
1010 We apply the access to PathChat (Lu et al., 2024)
1011 but have not received the authorization. These
1012 models include MedGemma-4B, Qwen2.5VL-3B,
1013 Qwen2.5VL-7B, MedVLThinker-7B, InternVL3-
1014 8B, PathGen-LLaVA-13B, and Patho-R1 (7B).
1015 MedGemma-4B is an open-source VLM released
1016 by Google based on finetuning Gemma with
1017 multimodal medical data. Qwen2.5VL-3B and
1018 Qwen2.5VL-7B are open-source VLMs from the
1019 Qwen team, Alibaba Cloud. They are trained
1020 with multimodal data in the general domain.
1021 InternVL3-8B is an open-source VLM released by
1022 OpenGVLab, and it is also trained with multimodal
1023 data from the general domain. For pathology-
1024 specific models, we consider PathGen-LLaVA-
1025 13B (Sun et al.), which is finetuned based on
1026 LLaVA 13B (Liu et al., 2023) with instruction data
1027 from PathGen; as well as Patho-R1, which has a
1028 pathology-specific image encoder and is finetuned
1029 based on Qwen2.5VL-3B with reasoning data.

1030 For the cross-modality generation task, we
1031 also consider a task-specific baseline method,
1032 known as Cell2Sentence (1B) (Levine et al., 2024).
1033 Cell2Sentence is finetuned with instructions and
1034 single-cell transcriptomic data from atlas-level
1035 datasets based on Pythia. This model can generate
1036 cells based on instructions.

1037 B Code and Data Availability

1038 We utilize NCSA, YCRC, and TokyoU HPC plat-
1039 forms to perform experiments. To train Team-
1040 Path, we utilize 32 NVIDIA H100 cores and 8
1041 NVIDIA H200 cores for 24 hours. The CPU
1042 memory upper bound is 80GB. The codes can be
1043 found in [https://github.com/HelloWorldLTY/
1044 TeamPath](https://github.com/HelloWorldLTY/TeamPath), and the license is the MIT license.

1045 The information on the datasets used in this
1046 manuscript can be found in Supplementary File
1047 1. To access TCGA data, an authorized account is
1048 required. To protect personal privacy, we will not
1049 release experts' answers.

C Acknowledgments

We thank Mr. Tong Ding for his suggestion on
model training and task selection.

D Author Contributions

T.L. and W.X. designed this study. T.L., W.X., and
H.Q. ran all the experiments. H.W., P.H., M.D. per-
formed human evaluation. All authors involved in
writing and reviewing. H.Z. supervised this study.

E Institutional Review Board (IRB) Approval.

This project has received approval from Yale IRB,
with project number 2000039055.

F Prompt list

The prompt used for the Pathology VQA is:

**Your task: 1. Think through the question
step by step, enclose your reasoning process in
<think>...</think> tags. 2. Then provide the cor-
rect single-letter choice (A, B, C, D,...) inside
<answer>...</answer> tags. 3. No extra informa-
tion or text outside of these tags.**

The prompt used for the self-verifier is:

**You are an expert in pathology. You are given
a QUESTION and a PROPOSED SOLUTION.
Your job is to: 1. Break down each component
of the proposed solution. 2. Think step by step
to verify if the proposed solution is correct given
the question and the figure. 3. Write a line of the
form "The proposed solution is correct" or "The
proposed solution is incorrect" at the end of your
response based on your analysis. QUESTION:
question. PROPOSED SOLUTION: solution.**

The prompt used for the self-corrector is:

**You are also given a question and an analysis
for the question. Your job is to outline your
step-by-step thought process for deriving a cor-
rect solution and also write down the correct
solution. Using this format: <think> Your step-
by-step reasoning of the question and solution
</think><answer> Your final answer </answer>
Question: question Solution: out_verifier.**

The prompt used for the reason corrector is:

**You are given QUESTION, REASON, and SO-
LUTION. Your task is to correct the REASON
and SOLUTION. QUESTION: question. SO-
LUTION: solution. REASON: reason The REA-
SON is WRONG. Your solution:**

1096 The prompts used for training TeamPath for caption
1097 summary include:

1098 **Provide a concise pathological summary of the**
1099 **tissue shown in this histopathology image, high-**
1100 **lighting any abnormal cellular or structural fea-**
1101 **tures in one paragraph. Based on the visual char-**
1102 **acteristics in this image, summarize the likely**
1103 **histological diagnosis and key indicators lead-**
1104 **ing to it in one paragraph. Describe the main**
1105 **histopathological patterns visible in this image**
1106 **and summarize what they suggest about the tis-**
1107 **sue state in one paragraph. Summarize the key**
1108 **morphological findings in this histopathology**
1109 **image, including any signs of malignancy, in-**
1110 **flammation, or necrosis in one paragraph. Gener-**
1111 **ate a pathology report-style summary based**
1112 **solely on this histological section, mentioning**
1113 **tissue type, grade, and diagnostic clues in one**
1114 **paragraph. Briefly summarize the clinical im-**
1115 **plications of the abnormalities visible in this**
1116 **histopathology image in one paragraph. From**
1117 **this histopathology image, extract and summa-**
1118 **rize the most diagnostically relevant features**
1119 **in one paragraph. Identify and summarize any**
1120 **histopathological hallmarks (e.g., mitotic figures,**
1121 **glandular formation, stromal invasion) present**
1122 **in the image in one paragraph. Write a sum-**
1123 **mary suitable for a pathology trainee explain-**
1124 **ing what this histopathology image represents**
1125 **and why in one paragraph. Provide an expert-**
1126 **level summary of the pathological findings in**
1127 **this histopathology image, including your confi-**
1128 **dence in the assessment in one paragraph.**

1129 The prompts used for training TeamPath for cross-
1130 modality generation (using IDC as an example)
1131 include:

1132 **Generate a list of 100 genes in order of de-**
1133 **scending expression from one spot shown in**
1134 **the histopathology image in IDC disease. Cell**
1135 **sentence:, Produce a list of 100 gene names in**
1136 **descending order of expression which represent**
1137 **the expressed genes from one spot shown in**
1138 **the histopathology image in IDC disease. Cell**
1139 **sentence:, Create a ranked list of 100 genes**
1140 **in decreasing order of expression from one**
1141 **spot shown in the histopathology image in**
1142 **IDC disease. Cell sentence:, List the top 100**
1143 **expressed genes from one spot shown in the**
1144 **histopathology image in IDC disease. Cell**
1145 **sentence:, Identify the highest expressed 100**
1146 **genes in decreasing order of expression from**

one spot shown in the histopathology image in
IDC disease. Cell sentence:, Enumerate a list
of 100 genes in descending order of expression
from one spot shown in the histopathology
image in IDC disease. Cell sentence:, Compile
a descending order list of 100 expressed genes
from one spot shown in the histopathology
image in IDC disease. Cell sentence:, Present
a sequence of 100 genes ordered by decreasing
expression level from one spot shown in the
histopathology image in IDC disease. Cell
sentence:, Generate an ordered list of 100 genes
by decreasing expression level from one spot
shown in the histopathology image in IDC
disease. Cell sentence:, Assemble a list of 100
genes from highest to lowest expression from
one spot shown in the histopathology image in
IDC disease. Cell sentence:

1166 **G SFT vs RL: Comparison for training** 1167 **strategies.**

1168 How to train a mature reasoning model has always
1169 been a controversial topic, as there exist several dif-
1170 ferent strategies and their performances and ranks
1171 might vary under different task settings or experi-
1172 ment settings (Chu et al., 2025; Wang et al., 2025;
1173 Liu et al., 2025b). Meanwhile, this kind of discus-
1174 sion has not been investigated in training a large
1175 reasoning model for histopathology analysis, and
1176 thus, we consider several different approaches to
1177 provide an empirical analysis to select and interpret
1178 the best combination, which might inspire future
1179 directions or different researchers.

1180 We conduct these experiments based on differ-
1181 ent base models as well as training strategies. Our
1182 base models include Qwen2.5VL-7B, which does
1183 not contain domain-specific knowledge and Patho-
1184 R1-7B, which contains domain-specific knowl-
1185 edge. We also consider different training strate-
1186 gies, including RL (GRPO), RL (DAPO), SFT, and
1187 SFT+RL (GRPO). Extended Data Figures 4 (a)
1188 and (b) show that GRPO can achieve a more obvi-
1189 ous score improvement while DAPO cannot make
1190 an improvement, which implies that a mixture of
1191 tricks does not contribute to training a multi-modal
1192 pathology reasoning model. Extended Data Figures
1193 4 (c) and (d) show that it is important to select a
1194 base model with domain knowledge for training,
1195 and performing SFT+RL or direct SFT settings
1196 does not benefit this task. Therefore, our optimal

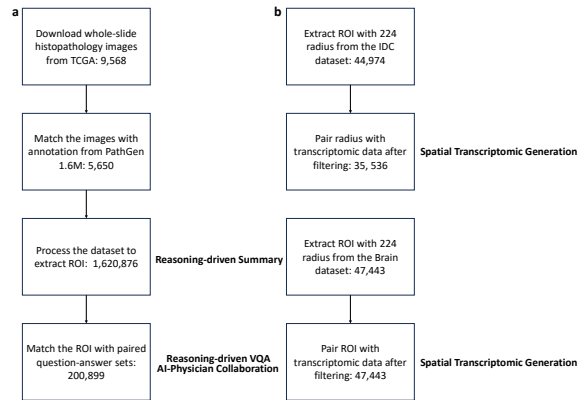
choice to build TeamPath for pathology VQA is Patho-R1-7B+RL (GRPO).

In conclusion, our results showcase that it is important to select a good model with domain knowledge to perform training, and for LMMs that can already possess domain knowledge, directly training them using RL policies can enhance their generalization capabilities without requiring specialized SFT training (also known as cold-start training). Our important findings also align with relevant research across different fields (Chen et al., 2025a; Zhang et al., 2026), which indirectly validates the reliability of our conclusions.

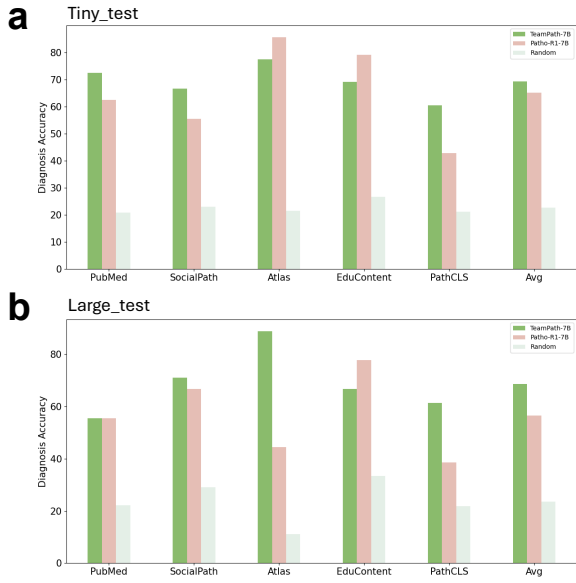
H Studies for training data ablation.

We also investigate if incorporating more diverse data could help TeamPath for generating a better reasoning path or not. To examine it, we compare the results between only using multi-choice VQA (mc VQA) data and using both multi-choice VQA and open-ended VQA (full VQA) data. Based on our evaluation results shown in Extended Data Figure 5 with validation from the PathMMU dataset, using full VQA data cannot boost TeamPath’s performances in generalizing the results for solving questions in the PathMMU dataset. Therefore, our optimal setting only takes mc VQA for RL training. Details of the reward design for these two different types of data are explained in the Method section.

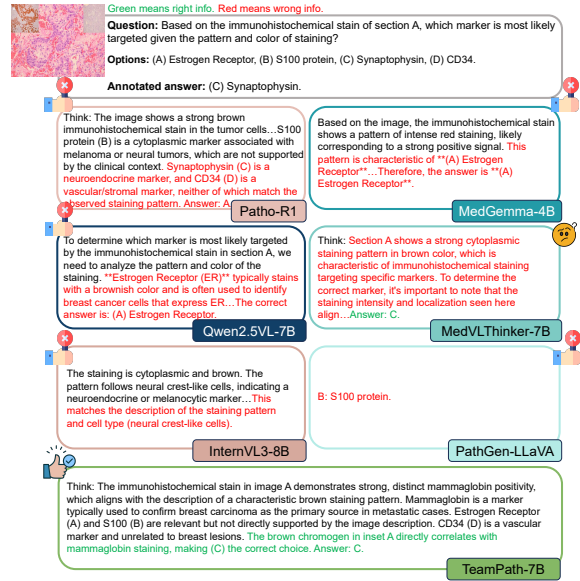
I Supplementary figures



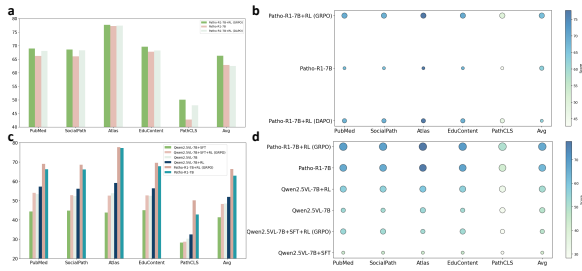
Extended Data Fig. 1: A flowchart of data-preprocessing used to train TeamPath.



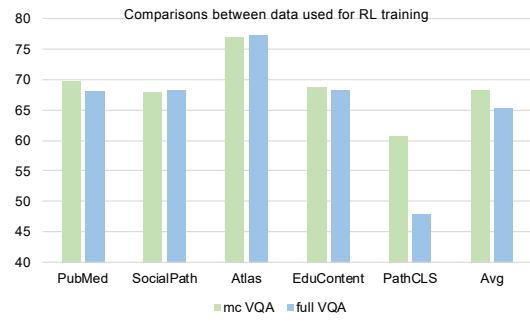
Extended Data Fig. 2: Accuracy of diagnosis-related questions in PathMMU. (a) Results reported based on the tiny_test dataset. (b) Results reported based on the large_test dataset.



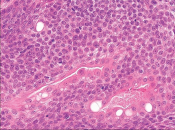
Extended Data Fig. 3: Case study (topic: lipoblast, which is a protein found in the presynaptic vesicles of neurons and neuroendocrine cells that plays a role in synaptic transmission) based on the outputs from different models. We highlight the correct information with green text and incorrect information with red text. For the models with errors, we consider two cases. The first case is a wrong answer, the second case is a confused reasoning path.



Extended Data Fig. 4: Training strategies optimization with different settings. The results are evaluated based on the PathMMU dataset. (a) Accuracy across different categories based on the base model and different RL strategies. (b) Accuracy and rank across different categories based on the base model and different RL strategies. (c) Accuracy across different categories based on different base models and different RL/SFT strategies. (d) Accuracy and rank across different categories based on different base models and different RL/SFT strategies.



Extended Data Fig. 5: Comparisons for the PathMMU VQA question set with different training data. The metric is accuracy.



Green means right info

Question: Analyzing the tissue section, what can be inferred about the tubular structures with respect to their cellular lining?

Options: A) The cells lack special features and are consistent with simple squamous epithelium. B) The cells possess distinct brush borders consistent with renal proximal tubules. C) The cells exhibit goblet cells typical of the intestinal lining. D) The cells have cuticular borders suggesting a level of specialization.

Red means wrong info

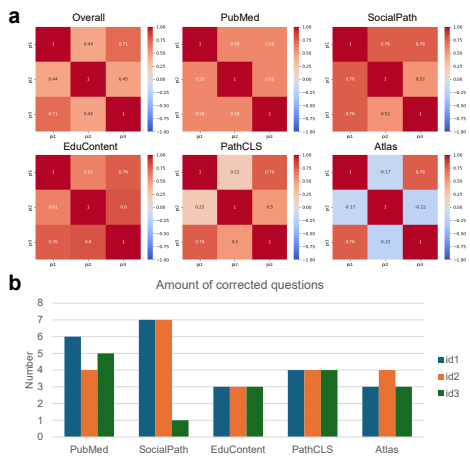
Expert answer: I do not know.

Expert reason: I actually don't know the answer to this question. The cells have brush borders (so they are not squamous), but many cells have brush borders, such as the intestinal epithelium. Goblet cells can be seen in any endoderm-derived epithelium, for example, the bronchus. I don't know what cuticular border means.

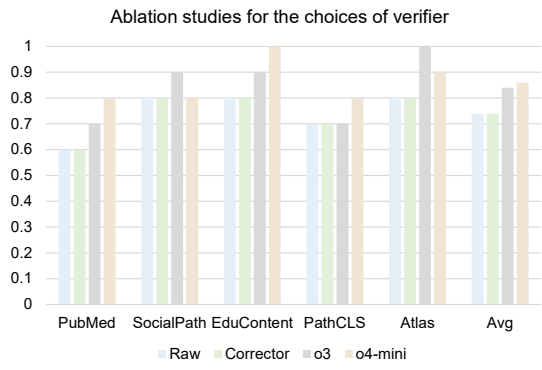
TeamPath answer: C) The cells have cuticular borders suggesting a level of specialization.

TeamPath reason: The tubular structures in the image are not renal proximal tubules. While the cells are well-oriented and dome-shaped, the luminal border lacks the brush border of microvilli seen in proximal tubules. Instead, the smooth, homogeneous cuticular border aligns with transitional epithelium, which is specialized for stretching. The absence of goblet cells (mucinous) and single-layered squamous epithelium (simple squamous) further excludes other options.

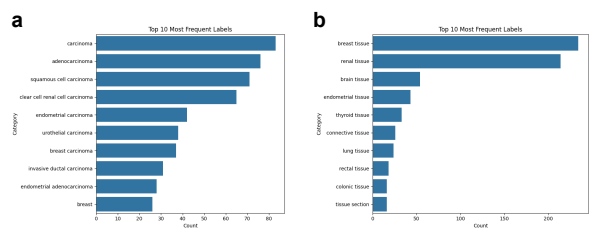
Extended Data Fig. 6: An example of using TeamPath to correct the reasoning path from pathologists.



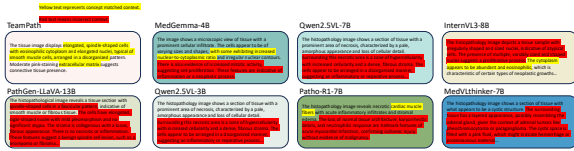
Extended Data Fig. 7: Results of expert feedback and TeamPath-corrected information. (a) The PCCs based on the accuracy of pathologists across different question categories in PathMMU. (b) The number of corrected samples for each pathologist made by TeamPath.



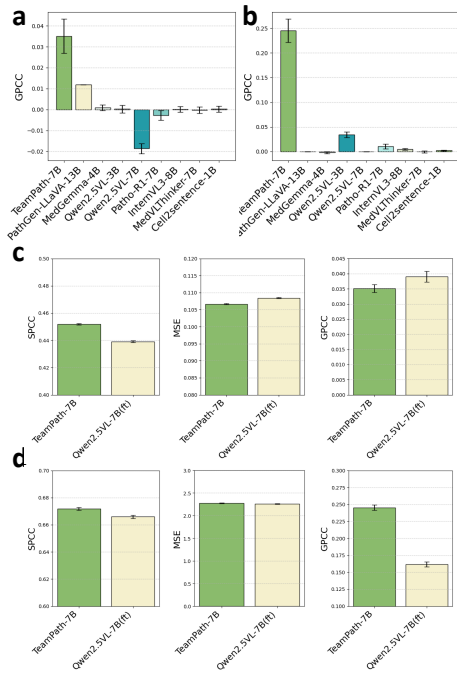
Extended Data Fig. 8: Ablation studies of using different models as verifiers.



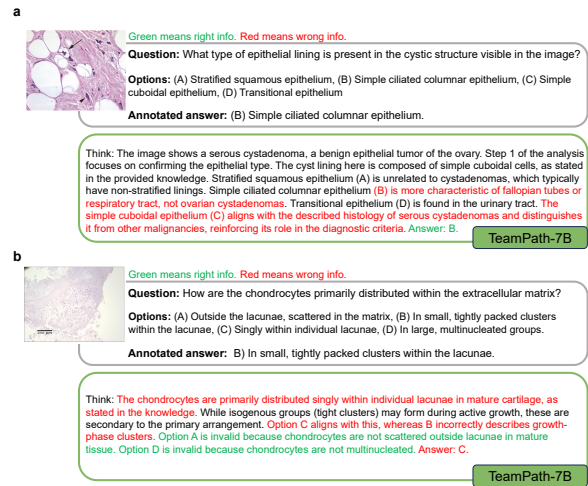
Extended Data Fig. 9: Distribution of image categories used for the testing of summarization performance. (a) Top 10 disease categories and their number. (b) Top 10 tissue categories and their number.



Extended Data Fig. 10: Examples of model outputs for caption summary tasks. We highlight both correct and incorrect information.



Extended Data Fig. 11: Extended analyses of spatial transcriptomic generation. (a) GPCC across different methods based on the brain tissue. (b) GPCC across different methods based on the IDC sample. (c) Comparison between finetuned Qwen2.5VL-7B and TeamPath based on the brain tissue. (d) Comparison between finetuned Qwen2.5VL-7B and TeamPath based on the IDC sample.



Extended Data Fig. 12: Case studies of TeamPath output for pathology VQA with (a) an incorrect reasoning path but a correct answer and (b) an incorrect reasoning leads to an incorrect answer.

J Supplementary Tables

Extended Data Tab. 1: Dataset statistics for training and testing sets.

Data type	Number of samples
Training_size (only MCA)	14,288
PubMed_test_tiny	281
PubMed_test	2,787
SocialPath_test_tiny	216
SocialPath_test	968
Atlas_test_tiny	208
Atlas_test	799
EduContent_test_tiny	255
EduContent_test	1,683
PathCLS_test	1,632
PathCLS_test_tiny	177
Test_total	9,006