

Evidence-based diagnostic reasoning with multi-agent copilot for human pathology

Chengkuan Chen^{1,‡}, Luca L. Weishaupt^{1,2,‡}, Drew F. K. Williamson^{1,2,3,*}, Richard J. Chen^{1,2,3,*}, Tong Ding^{1,5,*}, Bowen Chen^{1,2,*}, Anurag Vaidya^{1,2}, Long Phi Le², Guillaume Jaume^{1,2,3}, Ming Y. Lu^{1,2,3,6+}, and Faisal Mahmood^{1,3,4,7,+}

¹*Department of Pathology, Brigham and Women's Hospital, Harvard Medical School, Boston, MA*

²*Health Sciences and Technology, Harvard-MIT, Cambridge, MA*

³*Department of Pathology, Massachusetts General Hospital, Harvard Medical School, Boston, MA*

⁴*Cancer Program, Broad Institute of Harvard and MIT, Cambridge, MA*

⁵*Harvard John A. Paulson School of Engineering And Applied Sciences, Harvard University, Cambridge, MA*

⁶*Electrical Engineering and Computer Science, Massachusetts Institute of Technology (MIT), Cambridge, MA*

⁷*Harvard Data Science Initiative, Harvard University, Cambridge, MA*

‡ Co-1st authors, * Co-2nd authors, + Co-senior authors

Lead contact: Faisal Mahmood (faisalmahmood@bwh.harvard.edu)

Abstract

Pathology is experiencing rapid digital transformation driven by whole-slide imaging and artificial intelligence (AI)^{1–5}. While deep learning-based computational pathology has achieved notable success^{6–11}, traditional models primarily focus on image analysis without integrating natural language instruction or rich, text-based context. Current multimodal large language models (MLLMs) in computational pathology^{12–17} face limitations, including insufficient training data, inadequate support and evaluation for multi-image understanding, and a lack of autonomous, diagnostic reasoning capabilities. To address these limitations, we introduce PathChat+, a new MLLM specifically designed for human pathology, trained on over 1 million diverse, pathology-specific instruction samples and nearly 5.5 million question answer turns. Extensive evaluations across diverse pathology benchmarks demonstrated that PathChat+ substantially outperforms the prior PathChat copilot, as well as both state-of-the-art (SOTA) general-purpose and other pathology-specific models. Furthermore, we present SlideSeek, a reasoning-enabled multi-agent AI system leveraging PathChat+ to autonomously evaluate gigapixel whole-slide images (WSIs) through iterative, hierarchical diagnostic reasoning, reaching high accuracy on DDxBench, a challenging open-ended differential diagnosis benchmark, while also capable of generating visually grounded, human-interpretable summary reports.

Introduction

Pathology is undergoing a digital transformation, driven by the adoption of whole-slide imaging and advances in artificial intelligence (AI)^{1–5}. Deep learning-based computational pathology^{6–9} has achieved remarkable success in diverse tasks such as cancer detection, subtyping and grading^{10,11,18–22}, biomarker screening^{23–26}, patient outcome prediction^{27–37} and more^{38–48}. However, most computational pathology algorithms have been vision-only, operating on images without leveraging the rich textual knowledge from the pathology literature and clinical context that pathologists rely on^{11,18,19,21}. While early domain-specific visual language models, such as pathology-specialized CLIP-like encoders^{37,49–51}, demonstrated potential in zero-shot classification and retrieval, they were limited in their ability to encode rich clinical knowledge and to handle complex queries based on user instructions. Recently, the rise of foundation LLMs^{52–55} and multimodal generative AI^{56–60} has opened a new frontier in computational pathology, enabling a new generation of machine learning models to not only analyze histology images but also understand unstructured clinical context, medical terminology, and acquire diagnostic reasoning while facilitating intuitive, human-like interactions in natural language. For instance, PathChat¹², a recently introduced multimodal generative AI copilot, achieved high diagnostic accuracy on expert-selected region-of-interests (ROIs) when constrained to a multiple choice setting, and was shown to generate more preferable responses to diverse open-ended questions compared to state-of-the-art visual language models at the time of study.

Despite these advances, the recent development of domain-specific multimodal large language models (MLLMs)^{12–17}, including PathChat, have suffered from limited training data and large-scale evaluation across diverse benchmarks. Furthermore, their utility remains underexplored beyond single image question answering, with most pathology-specialized models currently lacking direct support for multi-image or interleaved image-text understanding. This limitation precludes their application in clinical or research scenarios where the user requires analysis based on contextual information across multiple selected regions of one or more slides. While recent progress in slide-level generative modeling^{16,61–66} has shown some early promises of conversational question answering for diagnostic assessment and reporting, these systems tend to be trained on limited coverage of tissue types and indications, lack the ability to transparently ground their prediction in finegrained morphological features via diagnostic reasoning, and are inherently restricted in scope to tasks in their training distribution due to their static modeling design. Finally, the current MLLM-enabled workflow is largely reactive, answering user-posed questions while requiring manual or pre-defined selection of regions of interest – lacking an explicit mechanism to autonomously drive the diagnostic evaluation process of routine whole slide images (WSIs) via multi-step reasoning^{67,68}. In the broader context of computational pathology, some early works^{69–72} have explored the idea of learning to navigate or zoom for histopathology image analysis, even incorporating large language models as core components of the workflow. However, these systems often lack the intrinsic design to carry out planning and diagnostic reasoning. Moreover, all such approaches currently rely on tailored annotations (at the slide or even ROI level) for disease-specific or task-specific training, limiting their ability to generalize across the wide range of organ systems and disease types encountered in pathology.

AI agent copilot for grounded visual reasoning in human pathology

To address these challenges, we introduce PathChat+, a new state-of-the-art multimodal large language model-based AI copilot for human pathology that builds upon the foundation of the first iteration of PathChat. PathChat+ is trained on over 1 million diverse visual language instruction samples, with comprehensive coverage of all pathology specialties, tissue types and diseases (**Figure 1A**), as well as direct support for multi-image understanding and high-resolution processing for each image. Furthermore, we introduce SlideSeek (see **Architectural design of SlideSeek of Methods** for more details), a reasoning-enabled multi-agent AI system leveraging the core capabilities of PathChat+ to enable autonomous assessment of gigapixel WSIs. Given a task description, SlideSeek methodically, hierarchically, and iteratively investigates different regions, integrates

multi-scale information, and identifies key morphologies required to make accurate diagnoses via diagnostic reasoning analogous to real-world pathology workflow. Its holistic slide-level analysis captures essential contextual cues—such as tissue architecture and peripheral features that may be overlooked in single ROI-based approaches.

In this study, we extensively evaluate the non-agentic foundation model, PathChat+, against state-of-the-art (SOTA) general purpose, instruction-following multimodal large language models (MLLMs) such as Llama⁵⁴, closed-source frontier models such as GPT-4o, as well as specialized medical and pathology models such as Quilt-LLaVA¹⁴. Additionally, we use a comprehensive collection of diverse, challenging cases spanning 52 unique diagnoses to assess the open-ended differential diagnostic capability on WSIs of SlideSeek and other MLLMs capable of multi-image analysis on gigapixel WSIs (See **Evaluation datasets** under **Methods** section).

Results

Large-scale ROI-level evaluation on diverse benchmarks and diseases

We began by assessing the general visual language understanding capabilities of PathChat+, the standalone MLLM, on a diverse collection of benchmarks spanning pathology visual question answering, image classification, and captioning (**Figure 2A**). For visual question answering (VQA), we use the 5 benchmarks from PathMMU⁷³, the current largest collection of domain-expert validated questions spanning multiple organ types and specialties, designed for evaluating MLLMs pathology VQA. Additionally, we report performance on PathBenchQA MCQ, a benchmark of multiple choice questions focused on making a diagnosis from a representative ROI image and relevant clinical context, as described in the original study¹². For image classification, we additionally select 3 diverse pathology image classification benchmarks, BRACS (fine-grained breast neoplasm subtyping), UniToPatho (colorectal polyp assessment), and HiCervix (cervical cytological cell classification) and convert them into a multiple choice format for close-ended model evaluation. Finally, for image captioning, we create the PathBenchQA Caption dataset, where a board-certified pathologist has annotated each image in PathBenchQA with a ground truth caption and use it for assessing the abilities of MLLMs to generate detailed morphological descriptions. We consider a comprehensive list of 10 MLLM models for comparisons, including closed-source frontier models: GPT-4o⁷⁴, Claude 3.5 Sonnet¹, Gemini 2.0 Flash², open-source general purpose models: Qwen2.5-VL⁷⁵, Llama 3.2⁵⁴ and LLaVA-OneVision⁶⁰, and specialized medical models that include pathology training: HuatuoGPT-Vision⁷⁶, LLaVA-Med 1.5¹³, PA-LLaVA¹⁵, and Quilt-LLaVA¹⁴. We also include PathChat 1¹² as a baseline. For each task, owing to the large number of model comparisons, we only plot the best performing model of each model category (i.e. the best performing closed source frontier model, the best open source general model, and the best performing medical-specialized model). For performance of all individual models, we refer readers to **Extended Data Table 4 - 12**.

In VQA, PathChat+ obtains the highest accuracy across all benchmarks (**Figure 2B**), consistently outperforming the best model from each model category as well as the predecessor PathChat 1 model. On the overall combined PathMMU test set, PathChat+ outperforms general purpose frontier models by an average of 13.8% in absolute accuracy score (+15.7% vs. Claude 3.5 Sonnet, +16.8% vs. GPT-4o and +9.0% vs. Gemini 2.0 Flash; $p < 0.001$ for all). When compared to Gemini 2.0 Flash specifically, the overall best performing model on the VQA benchmarks, PathChat+ is anywhere between 5.6% and 17.4% absolute percentage higher across the different sources of PathMMU ($p < 0.001$ for all), and scored 12.4 % higher ($p = 0.011$) on PathBenchQA MCQ. For comparison with domain-specific models, this lead extends even further to between 8.7% and 16.4% higher ($p < 0.001$ for all) on PathMMU and 36.2% ($p < 0.001$) on PathQABench MCQ. In classification, we observe a similar trend, with PathChat+ outperforming the best of all model categories, often by a substan-

¹https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf

²<https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/#gemini-2-0-flash>

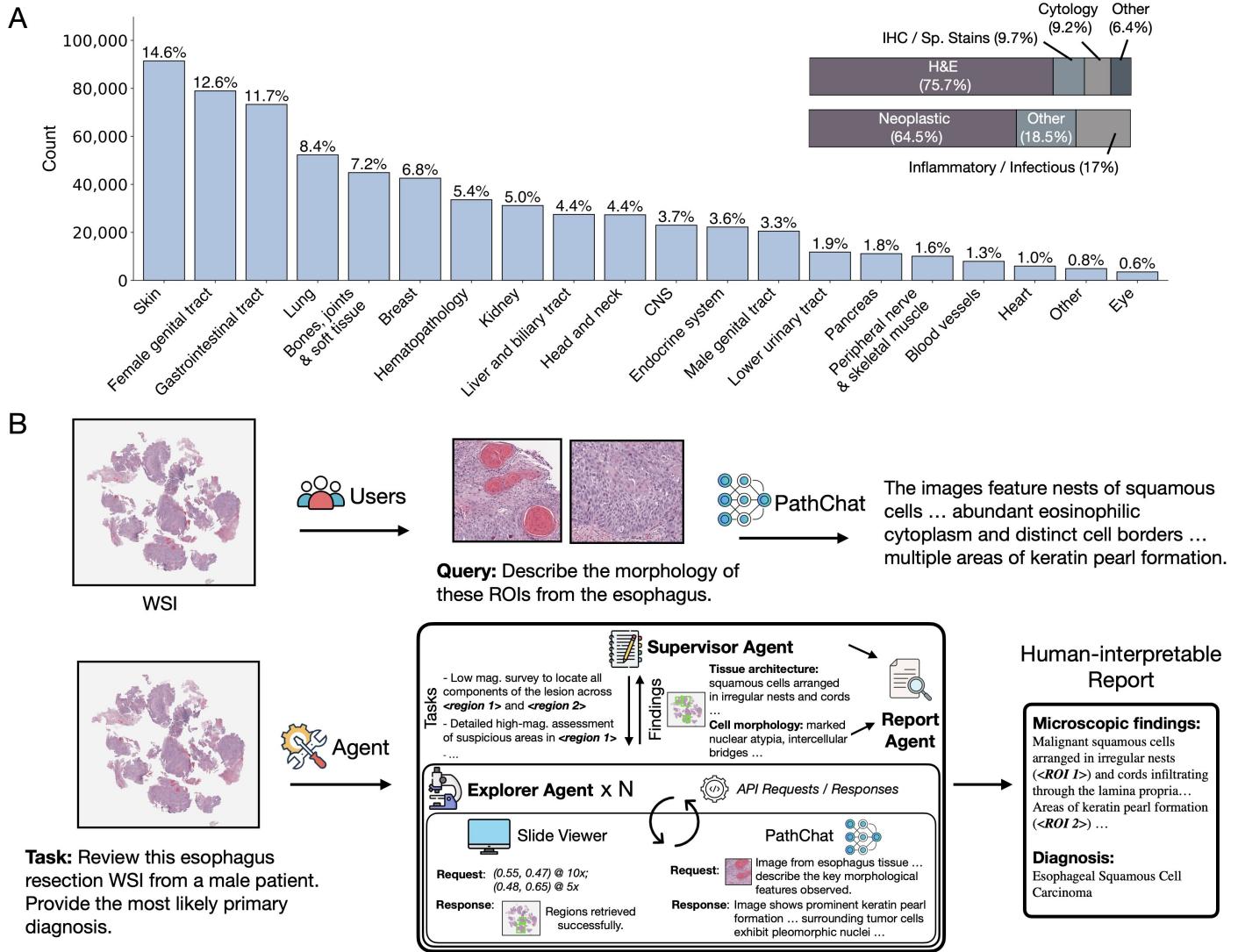


Figure 1: PathChat+ model development and SlideSeek overview. **A.** We curated the largest instruction finetuning dataset to date for human pathology with 1.13 million instructions, 5.49 million question answers turns based on 624 thousand unique images. A distribution of images is provided by tissue site, by staining and preparation (e.g. H&E histology vs. IHC / Special Stains), and by disease category (e.g. Neoplastic vs. inflammatory / infectious disease). Overall, our training set provides comprehensive coverage of diverse tissue, sample preparation and disease types. **B.** Interaction with a conventional chatbot-like multimodal large language model (MLLM) copilot (**top**) typically requires the user to manually select regions of interest (ROIs) from digitized slides (or specific fields of view from a digital microscope) and pose explicit questions. In contrast, our multi-agent-based AI system, SlideSeek (**bottom**), begins with a detailed, high-level task description and autonomously manages the analysis of the slide. A reasoning large language model (LLM) serves as the supervisor agent, continually tracking progress, refining diagnostic plans, and determining additional regions for examination. During each planning iteration, the supervisor instructs a team of specialized pathologist agents, each of which interacts with the slide viewer API and PathChat+ to concurrently analyze specific slide regions and report findings back to the supervisor. This iterative, hierarchical workflow continues until the supervisor agent determines sufficient evidence has been collected to establish a well-supported differential diagnosis. Subsequently, a separate report agent synthesizes the morphological evidence from critical ROIs into an interpretable, visually grounded diagnostic summary report presented to the user.

tial margin (e.g. +31.4% vs. Gemini 2.0 Flash and +43.2% vs. HuatuoGPT-Vision on BRACS, respectively; $p < 0.001$). Lastly, PathChat+ remains the best model on PathBenchQA Caption for generating morphological image captions, achieving the highest METEOR⁷⁷ score when compared to the expert annotated ground truth. Several examples of PathChat+ generated morphological captions can be found in **Extended Data Figure 1**. Overall, we observe that current domain-specific models, while competitive with open-source general purpose MLLMs, tend to under-perform closed-source frontier models, despite their emphasis on domain-specific knowledge acquisition. This underscores the need for curating larger, more comprehensive medical-specific datasets for training, a key distinguishing factor for PathChat+ compared to other models in the domain of human pathology.

Towards open-ended differential diagnosis with multi-image understanding

While current pathology question-answering benchmarks predominantly emphasize single-image interpretation and closed-ended classification within a limited set of predefined categories, pathologists in real-world clinical workflows routinely synthesize information across multiple images to generate open-ended differential diagnoses. To more closely reflect this practice, in this study, we also evaluate the capabilities of PathChat+ and other MLLMs in assisting pathologists by accurately recommending differential diagnoses within an open-ended setting based on multi-image inputs. To this end, we introduce DDxBench, a collection of 150 diagnostic H&E slides covering 52 unique neoplastic indications from all major organ sites (see **Evaluation datasets of Methods** for more details). Each slide is annotated with multiple ROIs from tumor tissues, providing a more realistic diagnostic challenge compared to single-ROI-based evaluations. For every slide, the corresponding question supplies essential contextual information such as tissue site and patient sex, accompanied by multiple ROI images, and prompts the model to generate a primary diagnosis along with two additional possible diagnoses in an open-ended manner to mimic real-world differential diagnosis (**Figure 3A**). For this part of the evaluation, we include all models assessed in the previous section that support multi-image inputs. Given the open-ended nature, to ensure accuracy in evaluation, we ask a board-certified anatomic pathologist to manually label each model's predictions as either correct or incorrect against the assigned ground truth diagnosis.

When considering the accuracy of the primary diagnosis predicted by each model, we find that PathChat+ achieves an overall highest accuracy of 80.0% on this challenging benchmark (**Figure 3B, Extended Data Table 13**), surpassing PathChat 1 (+8%, $p = 0.036$) and the next best model, Gemini-2.0 Flash by a substantial margin (+33.3%, $p < 0.001$). Next, we consider the entire set of differentials (*i.e.* primary diagnosis + two additional predicted possible diagnoses) and whether any one of them matches the groundtruth. This is a helpful measure since occasionally it is difficult to definitively confirm the diagnosis based on morphology alone, and the ability of an AI assistant to propose a reasonable set of candidate possibilities, *i.e.* differential diagnosis, can still help speed up diagnostic workflow by guiding pathologists towards the appropriate next steps (*e.g.* ancillary tests) necessary to reach a final diagnosis, even if the model does not correctly infer the diagnosis directly. In this setting, accuracy of PathChat+ reaches 93.3% (+4.67% against PathChat, $p = 0.118$); +22.7% against the next best model, GPT-4o, $p < 0.001$) suggesting its potential to assist pathologists in accurately identifying key candidate diagnoses for consideration and prioritization.

Evaluation of multi-agent system for interpretable slide assessment

To address the fundamental challenge of efficiently analyzing gigapixel-scale whole slide images (WSIs), we developed SlideSeek, a reasoning-enabled multi-agent MLLM system leveraging PathChat+. Unlike conventional MLLM approaches requiring manual ROI selection, SlideSeek autonomously navigates and evaluates WSIs through a hierarchical, multi-step reasoning process that closely mimics the diagnostic workflow of human pathologists. At SlideSeek's core is a supervisor agent that orchestrates the diagnostic exploration by formulating initial hypotheses based on low-magnification overview, strategically planning regions for detailed examination, and coordinating multiple explorer agents to investigate different tissue regions concur-

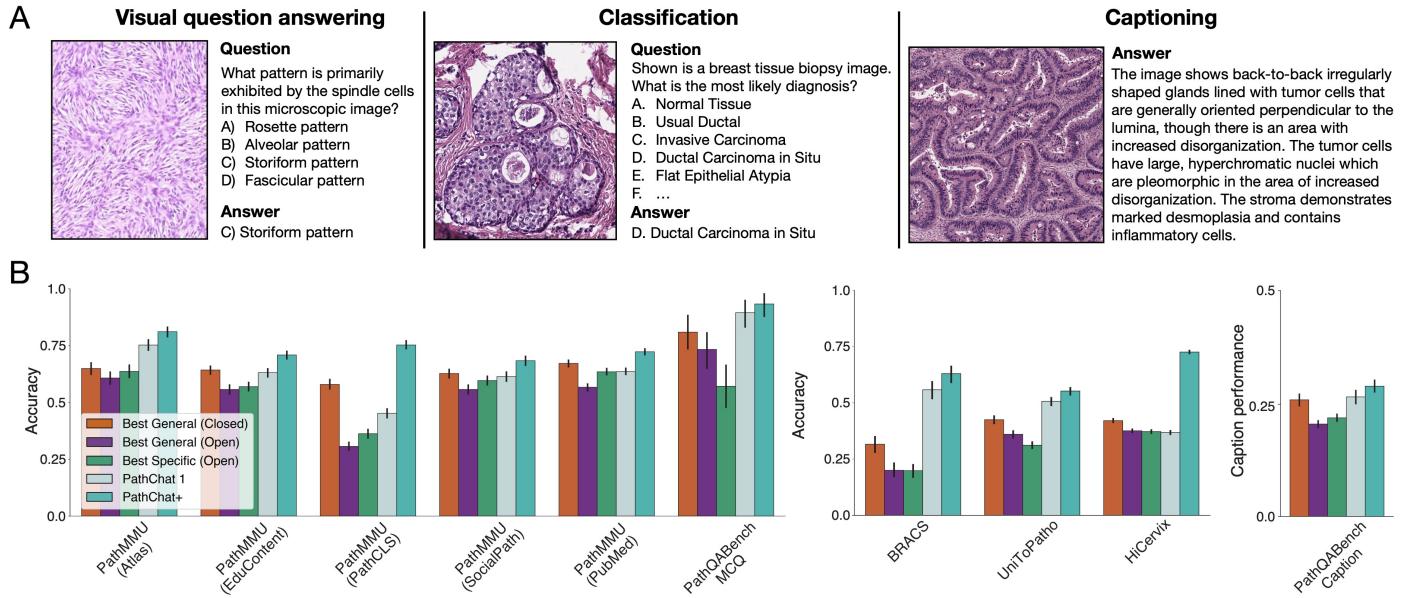


Figure 2: Performance of MLLMs on single ROI benchmarks. **A.** Illustrative example of visual question answering (VQA), image classification and caption respectively. The input always includes a salient histology image ROI and the corresponding instruction for the tasks (see **Evaluation datasets** for the details). **B.** Performance of MLLMs on VQA, classification (accuracy score) and image captioning (METEOR score) tasks. Due to the large number of model comparisons, we only plot the best performing model of each model category (i.e. the best performing closed source frontier model, the best open source general model, and the best performing medical-specialized model). Performance of all individual models are presented in **Extended Data Table 4 - 12**. Error bars represent 95% confidence intervals.

rently. This hierarchical approach enables efficient analysis of gigapixel WSIs by focusing high-magnification examination only on diagnostically relevant regions identified during low-power exploration, similar to how pathologists first scan slides at low power before examining suspicious areas at higher magnification. An example of this multi-agentic exploration is presented in **Figure 4**.

A key advantage of SlideSeek is its interpretable reasoning chain, where each diagnostic conclusion is explicitly linked to specific morphological findings at precise slide coordinates. During exploration, SlideSeek systematically examines an average of 47.4 ± 40.1 regions across multiple magnification levels (11.9 ± 6.2 high-power fields, 17.9 ± 31.2 medium-power fields, and 17.6 ± 17.9 low-power fields), creating a comprehensive assessment comparable to human pathologists' workflow. Identifying areas of clinical significance at lower magnification and using reasoning to navigate them at increasing magnifications reduces the number of ROIs that have to be processed over conventional computation pathology methods, which usually require examining all ROIs that contain tissue at medium to high magnifications (e.g. DDxBench WSIs contain 1020 ± 783 896×896 pixel tissue ROIs at 20x). A distribution of the number of ROIs required for SlideSeek to make a diagnosis can be found in **Figure 3D**. There was no statistically significant correlation between the number of ROIs examined and performance. Furthermore, SlideSeek demonstrates metacognitive awareness of its diagnostic confidence, achieving 82.7% accuracy for the 127 cases where it expressed high confidence, compared to 65.2% accuracy for the 23 cases where it acknowledged uncertainty, an important capability for responsible clinical deployment. **Extended Data Figure 2** presents two cases in which the agent classifies its own prediction as low or high confidence and provides visually grounded reasoning for this classification. When evaluated on DDxBench, SlideSeek achieved diagnostic accuracy comparable to PathChat+ applied to pre-selected high-magnification ROIs, demonstrating that autonomous multi-agent reasoning can effectively analyze entire slides

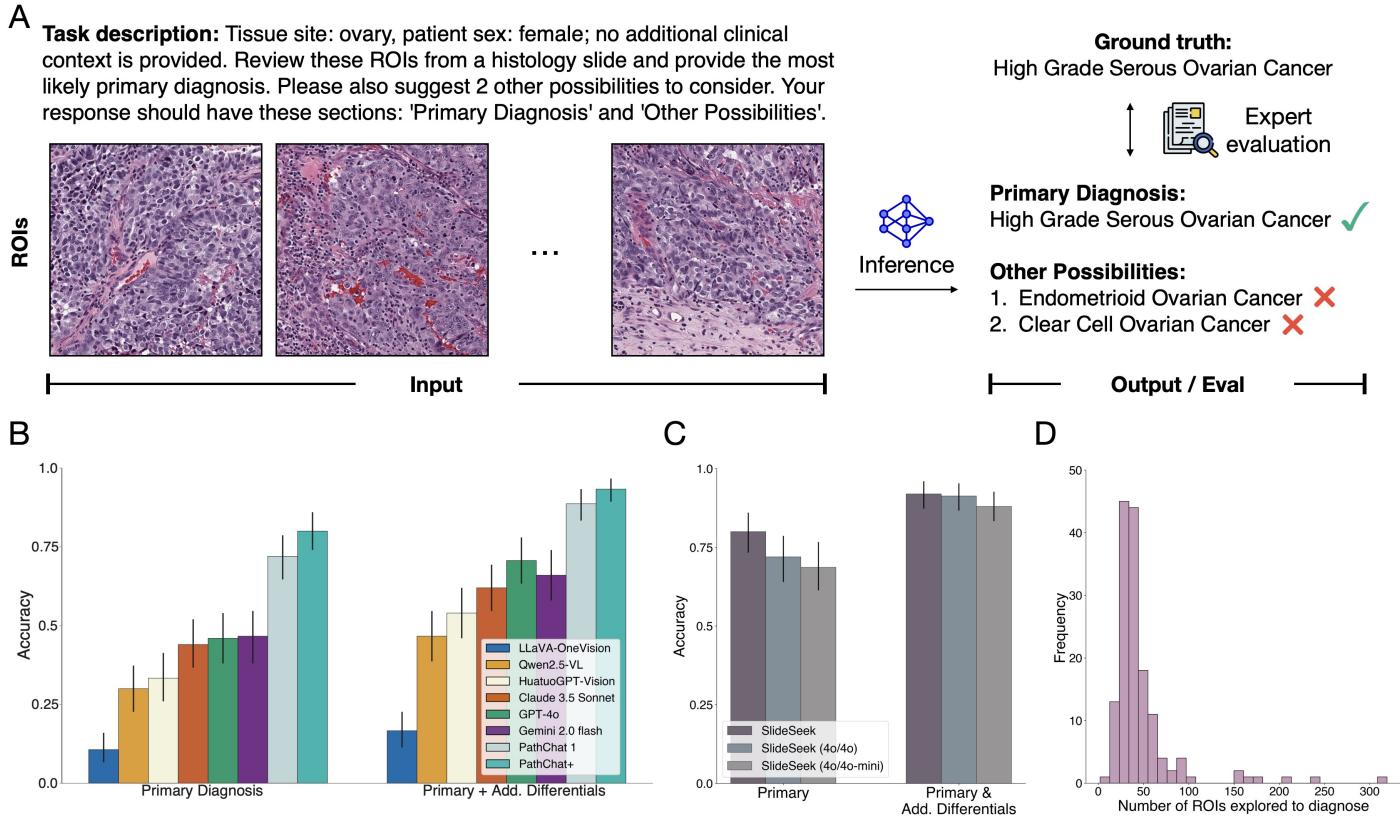


Figure 3: Performance of MLLMs and SlideSeek on open-ended differential diagnosis. **A.** Illustration of task input and evaluation. As shown here, for evaluating MLLMs, the input includes multiple histology image ROIs and the corresponding instruction for the tasks (see **Evaluation datasets** for the details). For SlideSeek, which is capable of directly acting on WSIs, we use the WSI directly as input instead of the ROI. Each model is tasked with generating a primary diagnosis and two additional possible diagnoses in an open-ended manner. Due to the open-ended nature, a board-certified anatomic pathologist manually assessed each model’s predictions against the assigned ground truth diagnosis. Performance is computed for whether the correct disease was listed as the primary diagnosis (Primary) or listed in the differential, which contains the primary and two additional diagnoses (Primary & Add. Differentials). **B.** Performance of MLLMs based on ROIs. The detail scores are presented in **Extended Data Table 13**. Error bars represent 95% confidence intervals. **C.** Performance of SlideSeek DDxBench based on WSIs. Varying the LLM configuration of the SlideSeek model in an ablation study on DDxBench. Swapping the reasoning agent (o1) for a non-reasoning agent (GPT-4o) and subsequently decreasing the parameter count of the explorer agent (GPT-4o to GPT-4o-mini) both cause a drop in performance. **Extended Data Table 14** presents numerical scores. **D.** The number of ROIs explored per case by SlideSeek to make a diagnosis.

without requiring expert pre-selection of diagnostic regions. Furthermore, SlideSeek generates comprehensive, evidence-based reports that explain its diagnostic reasoning through explicit visual grounding, enhancing trust and facilitating human expert verification.

When evaluating the relative contributions of different LLMs for the agent configurations in an ablation study on DDxBench (see **Figure 3C** and **Extended Data Table 14**), we found that the reasoning capabilities of the supervisor agent significantly impact diagnostic performance. Configurations using a high-performance reasoning model (o1) as supervisor outperformed those using an LLM that cannot reason (GPT-4o) by 8.0% in primary diagnosis accuracy, though this difference did not reach statistical significance ($p = 0.141$). Furthermore, the capabilities of explorer agents proved crucial, with configurations using capable explorer agents

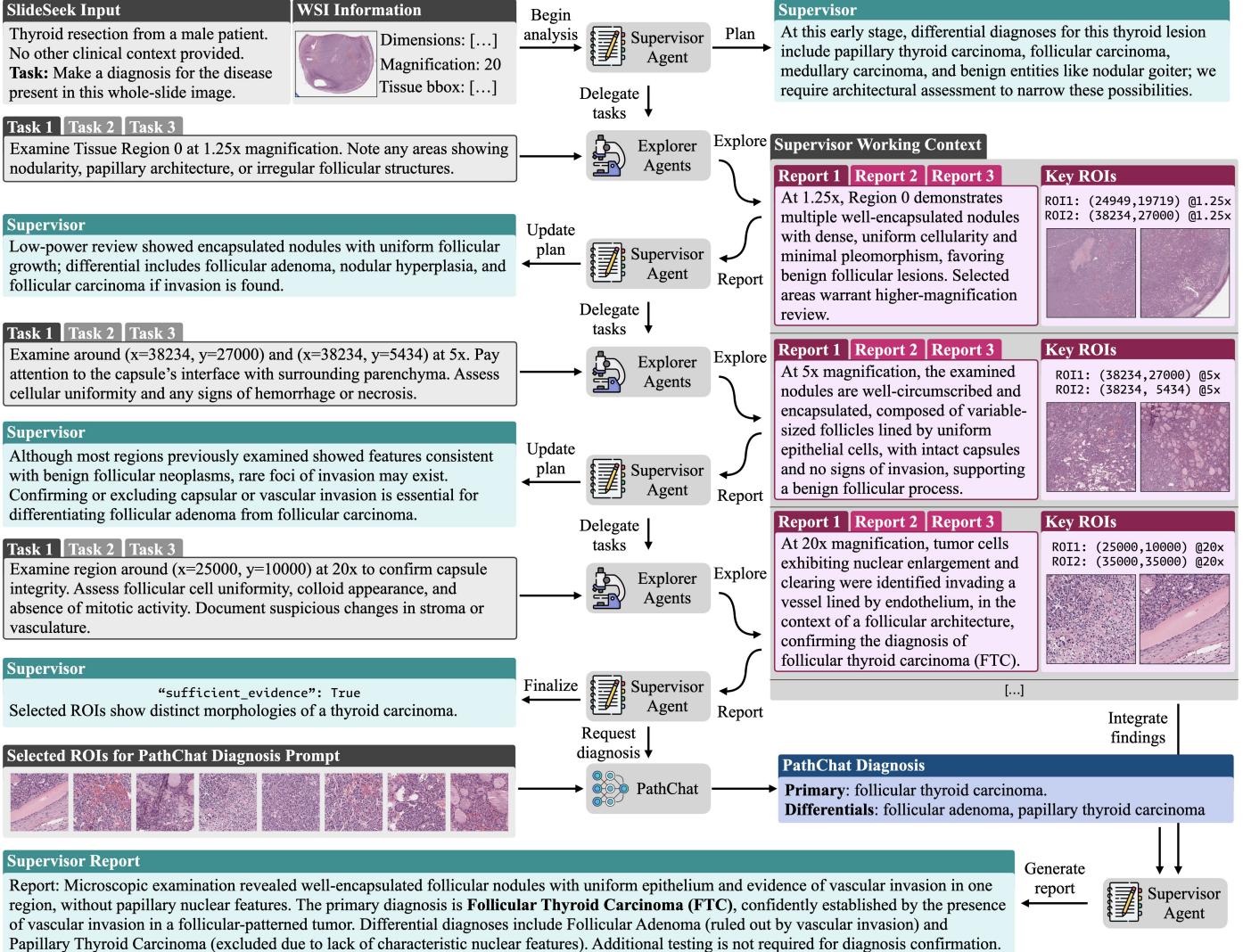


Figure 4: Example of SlideSeek on DDxBench Illustration of the interaction between the supervisor agent and its team of pathologists correctly diagnosing a whole slide image of thyroid tissue containing follicular thyroid carcinoma. The supervisor is provided with a system prompt that provides a low-resolution thumbnail image of the slide, clinical information, and the task to diagnose the disease present in the slide. At each step, the supervisor refines its hypotheses and plan based on the latest observations, creating new tasks for the explorer agents to focus on suspicious regions of tissue at higher magnification. The explorer agents carry out the tasks in parallel by scanning around the slide at the specified magnification level and noting the tissue's morphology for their individual reports to the supervisor through the help of PathChat+. While most of the tissue seems benign at lower magnification, the explorer agents find invasive cancer cells at high magnification. Ten regions of interest (ROIs) are passed to PathChat+ for a primary and differential diagnosis. Finally, the supervisor summarizes the exploration and diagnosis in a report.

(GPT-4o) outperforming those with significantly lower parameter count and inferior representational capacity (GPT-4o-mini) explorers by 11.3% ($p = 0.025$), highlighting the importance of effective microscopic feature recognition at the region level.

Discussion

In this study, we introduced PathChat+, a new state-of-the-art multimodal large language model specifically developed for human pathology. Our model demonstrated significant advances in pathology-specific

visual-language tasks, consistently outperforming existing general-purpose frontier models, as well as domain-specific models, across diverse benchmark datasets, highlighting the importance of domain-specific instruction tuning with large-scale, comprehensive training data.

Beyond single-image assessments, we demonstrated that PathChat+ significantly advances open-ended differential diagnostic capabilities, and can be extended to perform autonomous assessment of gigapixel WSIs effectively via our design of a multi-agent AI system, SlideSeek. SlideSeek autonomously drives diagnostic workflows through iterative multi-scale reasoning and multi-image understanding, aligning closely with real-world clinical practices. Specifically, we observed performance of SlideSeek on DDxBench matching or exceeding PathChat+ alone applied to curated ROIs, underscores the substantial clinical potential of multi-agent reasoning systems in digital pathology.

Despite these promising results, several important limitations and future directions should be noted. First, although our evaluation is extensive, it remains primarily retrospective and benchmark-driven. Prospective validation in clinical workflows is essential for accurately assessing real-world clinical impact. Lastly, while SlideSeek autonomously identifies and analyzes diagnostically relevant regions, promising future research directions involve designing agentic systems for analyzing complex cases with a large number of serial sections and tissue blocks and incorporating multimodal information from electronic health records, genomic data, and additional clinical imaging modalities to further enhance clinical utility.

Online Methods

PathChat+ dataset curation

For instruction tuning of PathChat+, we curated a new dataset of diverse instruction formats and topics from images and notes from internal case reports, teaching materials, presentations, as well as annotated regions of interests from in-house WSIs. To begin, we remove images of size smaller than 336×336 pixels. We then developed an image-based classifier by finetuning the image encoder of the lightweight foundation model CONCH⁷⁸ on a small set of manually labeled examples to identify pathology images from non-pathology-related images (e.g. photos of lab equipment or people) for quality control and similarly used lightweight local LLMs such as Qwen2-7B⁷⁹ for filtering out low-quality source materials that are overly concise (e.g. image caption that simply reads “This is an H&E stained slide”) or ambiguous (e.g. images associated with text that describes findings from medical literature without clear reference to the image or specific features observed). Similar to previous MLLM works such as LLaVA^{56,60}, for some instruction formats, we use general purpose LLMs^{54,79} to automatically structure the original source text into a question / answer format (e.g. rephrasing an image caption into a natural-sounding response to the question “Can you provide a morphological description for this lung biopsy image.” or converting representative regions paired with a case report into a multiple choice question with the ground truth diagnosis and other mentioned differential diagnoses or morphologically distinct entities as choices). In those scenarios, we design text prompts specifically for each data source and iteratively refine them with the help of pathologists until we observe satisfactory data quality on a randomly audited subset of the data. Finally, similar to previous work¹², we combine non-pathology images with sampled natural images from publicly available data sources such as MS COCO⁸⁰ to construct guardrail examples to instruct the model to refuse queries not related to pathology image analysis with a programmed message: “Sorry I can only assist you with queries related to pathology.” In total, our dataset consists of 1,133,241 instruction examples, with 5.49 million question / answer turns for 624 thousand unique images (median width: 759 pixels, median height: 607 pixels) after excluding 8,034 guardrail images. The instruction formats can be roughly characterized by “conversation” ($n = 238,983$), “description” ($n = 163,342$), “multiple choice” ($n = 78,443$), “free response” ($n = 354,558$), “text-only” ($n = 289,881$) and “guardrail” ($n = 8,034$).

PathChat+ architecture and training

PathChat+ is a multimodal large language model (MLLM) designed and trained to understand both visual (*i.e.* pathology images) and textual inputs (*e.g.*, an instruction prompt to describe the key morphology observed in a histopathology image, and suggest a likely diagnosis). We follow the general architectural design of PathChat 1¹² and popular state-of-the-art MLLMs, consisting of a vision encoder for encoding images from RGB pixels into rich, compact feature representations, a decoder-only, autoregressive LLM for converting tokenized visual and textual representations into meaningful natural language outputs, a multimodal projector for bridging the representation space of the vision encoder and the LLM. We chose CONCH v1.5 as the vision encoder, based on the standard (ViT-L) architecture, as it has been demonstrated to produce robust baseline performance in diverse computational pathology tasks⁸¹, and serves as a strong starting point for further integration and multimodal foundation model development. Similarly, the multimodal projector consists of an attention pooling layer that reduces the dense feature map produced by the vision encoder into a sequence of 128 tokens, followed by a 2-layer MLP with GeLU activation to project them into the same dimension as the embedding layer of the LLM. Finally, for the LLM backbone, we adopt Qwen2.5, which is a widely-used state of the art open-source LLM. Specifically, similar to PathChat 1, we chose the 14 billion parameter instruction-following variant, which offers a good trade-off between computational cost and performance⁷⁹.

Compared to the previous PathChat 1 model, we take inspiration from the AnyRes strategy⁸² for processing high-resolution images by dividing each input image into a grid of 448×448 tiles, first padding to the nearest supported grid size as needed. Each tile is then independently processed by the vision encoder before concatenating into a combined token sequence. If more than one tile is present in the grid, a thumbnail of the original image (resized to 448×448), serving as a low-resolution, global view, is also encoded and its token sequence is prepended to rest of the sequence representing the particular image. We use a maximum grid size of 4 tiles, supporting image sizes of up to 896×896 (represented as a 2×2 grid of tiles + the image thumbnail), while larger images are rescaled down to the nearest grid size. Our strategy therefore represents each image by a variable number of tokens, between 128 tokens (448×448 image or smaller) and 640 tokens (*e.g.* 896×896 image). During both training and inference, images are represented initially by special placeholder token ids in the input token sequence. After passing through the LLM’s embedding layer, these placeholder tokens are then replaced by the actual token representations generated by the vision encoder and multimodal projector for the corresponding images, forming the complete embedded token sequence that is processed by all subsequent layers in the LLM. When there are multiple images in a given instruction, we mark their boundary in the sequence using the newline (“\n”) token.

For model training, we follow the widely-used two stage approach where in the first pretraining stage, we freeze the LLM backbone, and use a subset of one hundred thousand images and annotated captions to train the multimodal adapter, which has been found to be beneficial to downstream performance⁸³ compared to without pretraining. In the second, instruction finetuning stage, both the LLM and projector are unfrozen and are trained to model the likelihood of the groundtruth reference answer turns in the training instruction sequences (no loss is applied to image tokens or text tokens in the question turns). Pretraining was completed using an $8 \times$ NVIDIA A100 80GB GPU node while finetuning was performed using 24 A100 GPUs in a multi-node distributed training setup (see **Extended Data Tables 1 and 2** for detailed hyperparameters used in the training process.)

SlideSeek architecture design

SlideSeek is a multi-agent system capable of autonomously analyzing gigapixel WSIs via planning and self-guided hierarchical exploration. We outline the key components below and illustrate the overall workflow in **Algorithm 1**.

The supervisor agent is central to the SlideSeek architecture, managing high-level strategy, task orchestration, and systematic exploration of WSIs. It receives the initial analytical task description, contextual patient and specimen information (e.g., tissue type, patient sex), and preliminary visual information of the slide. Visual data include slide dimensions, coordinates of bounding boxes surrounding each tissue region, and a low-resolution thumbnail overview. Tissue bounding boxes are identified via the Trident WSI preprocessing library⁸⁴, specifying each region through upper-left and lower-right corner coordinates (e.g., top-left corner: (637, 4467), bottom-right corner: (10612, 13034)).

The supervisor agent receives the following instructions: “You are the supervising pathologist managing a laboratory of skilled pathologists working remotely. Your task is to coordinate a systematic examination of this whole-slide image for diagnostic purposes. Create a detailed exploration plan with specific regions to investigate at appropriate magnifications. Break your plan into logical steps and assign concrete tasks to your pathologists. Review their findings after each step before proceeding to the next step.”

Upon initiation with the above information, the supervisor formulates initial hypotheses about the slide pathology and then develops a comprehensive and structured exploration plan, prioritizing regions for detailed examination, identifying specific pathological features of interest, and designating required magnifications. To operationalize this plan, the supervisor agent defines specific tasks for subordinate explorer agents. Tasks explicitly indicate precise spatial coordinates and extents for examination (e.g., “Examine tissue region #3 from x=1000-2000, y=3000-4000”), required magnification levels (e.g., 1.25x for architectural assessment, 20x for cellular details), specific pathological features to document, and relevant context derived from the current hypotheses. After subordinate agents complete their tasks, the supervisor evaluates their findings, updating hypotheses and modifying the exploration plan as necessary. Justifications for each change are clearly articulated, reflecting insights derived from cumulative observations. The iterative cycle of planning, task assignment, review, and hypothesis updating continues until sufficient evidence is gathered to support a definitive diagnosis. Upon completion, the supervisor signals readiness to finalize the analysis by setting a completion flag, collating key findings into selected ROIs, and forwarding these to a specialized diagnostic AI (PathChat+) for differential diagnosis. We use OpenAI o1⁶⁷ as the supervisor agent due to its state-of-the-art instruction following and advanced reasoning capabilities.

Multiple explorer agents perform the detailed investigation tasks assigned by the supervisor in parallel. Each explorer agent receives clearly defined tasks specifying well-defined tissue regions, recommended magnification levels, and morphological features to analyze. Explorer agents systematically navigate these regions by referencing the thumbnail image of the slide and making navigation requests to a slide-viewer API, which uses the OpenSlide library to retrieve image pixels from the WSI at the specified location and magnification. The requested ROI images are subsequently also forwarded to PathChat+ to identify key morphologies observed in each image, helping the explorer agent identify diagnostically relevant ROIs in reference to the overall task context and hypotheses provided by the supervisor agent. Once exploration of an assigned region is deemed to be complete, explorer agents submit detailed reports back to the supervisor that include key ROIs identified in the region and their findings, allowing iterative refinement of hypotheses and planning until the analysis is fully completed. Since the explorer agents primarily focus on making structured tool call requests to the viewer API, information summarization, and report synthesis without high-level reasoning, we use OpenAI GPT-4o for each explorer agent as it has lower observed latency (due to the lack of test-time scaling) and higher affordability compared to o1.

The iterative analysis between the supervisor agent and explorer agents continues until the supervisor agent determines that all diagnostically relevant regions have been explored and there is sufficient morphological evidence to render a diagnostic recommendation, ending further exploration. At this point, it identifies up to

10 diagnostically relevant ROIs discovered during exploration, and submits them to PathChat+ for differential diagnosis. Finally, a reporting agent is called upon to draft a report summarizing the microscopic findings of the slide, and the recommended primary diagnosis and additional possibilities in the differential. Notably, the report is a humanly interpretable summary of SlideSeek’s output, explicitly linking morphological features to specific ROIs identified by the explorer agents and supervisor.

Algorithm 1 Iterative Multi-Agent Exploration

```

1: Initialize supervisor agent with slide description and context information.
2: Supervisor formulates initial hypotheses.
3: numExplorers ← number of explorer agents
4: analysisComplete ← false
5: while analysisComplete == false do
6:   Supervisor creates exploration plan based on current hypotheses.
7:   Supervisor assigns tasks to explorer agents based on numExplorers.
8:   for  $i \leftarrow 1$  to numExplorers do
9:     Explorer agent  $i$  navigates assigned ROI at specified magnification.
10:    Explorer agent  $i$  documents morphological features and rationale.
11:    Explorer agent  $i$  returns findings to supervisor.
12:  end for
13:  Supervisor reviews explorer agents’ findings.
14:  Supervisor updates hypotheses.
15:  sufficientEvidence ← check if sufficient evidence collected
16:  if sufficientEvidence == true then
17:    analysisComplete ← true
18:  else
19:    Supervisor modifies exploration plan.
20:  end if
21: end while
22: Supervisor collates final ROIs.
23: Supervisor submits findings to diagnostic AI (PathChat+).
24: Report agent drafts summary report with explicit references to ROIs and their supporting morphological
evidence

```

MLLM evaluation

We compare PathChat+ against ten state-of-the-art MLLM baselines which can be categorized as: (1) General purpose closed frontier models (GPT-4o⁷⁴, Claude 3.5 Sonnet³, and Gemini 2.0 Flash⁴), (2) General purpose open models (LLaVA-OneVision⁶⁰, Llama 3.2⁵⁴, and Qwen2.5-VL⁷⁵), and (3) Medical and pathology-specific open models (Quilt-LLaVA¹⁴, PA-LLaVA¹⁵ and HuatuoGPT-Vision⁷⁶, and LLaVA-Med-1.5¹³). We provide a brief description of each category below.

General purpose closed models. For closed weight, general purpose multimodal frontier models, we use their official APIs for evaluation. Specifically, we used GPT-4o (OpenAI), Claude 3.5 Sonnet (Anthropic), and Gemini 2.0 Flash (Google), the flagship general purpose frontier model from each commercial entity available at the time of model evaluation in our study.

³https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf

⁴<https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/#gemini-2-0-flash>

General purpose open models. For general purpose MLLMs with open weights that can be downloaded and deployed locally, we access their weights from their respective, official HuggingFace repositories. We included the following representative models: LLaVA-OneVision, Llama 3.2, and Qwen2.5-VL. We use their default generation parameters and image processors for all tasks.

Medical purpose open model. Similar general purpose open models, we select representative models, biomedical specialized MLLMs with at least some pathology image data in their training, including: LLaVA-Med-1.5, Huatuo-GPT-Vision, Quilt-LLaVA and PA-LLaVA. We access their weights from their respective HuggingFace repositories linked from their official Github pages at the time of study and use their default general parameters and image processors for all of our evaluation tasks. Note we only include LLaVA-Med in the image captioning task as it is not trained on multiple-choice instruction following data¹⁴.

Evaluation datasets

PathMMU⁷³ is an expert-validated pathology benchmark. We use the validation and test splits to evaluate all models. Originally, there are 710 Q&A pairs with 510 images in the validation set and 9,677 Q&A pairs with 7,213 images in the test set. Within the test set, a small subset (test-tiny) comprising 1,156 Q&A pairs was selected to establish expert-level performance, which was assessed by two groups of pathologists. We follow the paper⁷³ to split the data into five categories (i.e., PubMed, SocialPath, EduContent, Atlas and PathCLS). In SocialPath, we removed 47 image links that are no longer available, resulting in 699, 1,156, and 9,618 Q&A pairs in our final validation, test-tiny and test subset. Note that the Atlas subset is not the same as reported in the original paper although the number of Q&A pair is the same⁵. To facilitate evaluation, we append the instruction “Answer with the option’s letter from the given choices directly.” to each question.

UniToPatho⁸⁵ is an annotated dataset for colorectal polyps classification and adenomas grading. The dataset contains six type of tissues from colorectal polyps (i.e., Normal, hyperplastic polyp, tubular adenoma_high-grade dysplasia, tubular adenoma_low-grade dysplasia, tubulo-villous adenoma_high-grade dysplasia and tubulo-villous adenoma_low-grade dysplasia). We used 2,799 ROIs from the test split of this dataset. For each ROI, we created a closed-ended multiple-choice question that asks: “Shown is a colorectal tissue biopsy image. What is the most likely diagnosis? Answer with the option’s letter from the given choices directly.”

BRACS⁸⁶ contains six subtypes of lesions and normal tissue sectioned from the breast. We used 570 ROIs from the test split. For each ROI, we create a closed-ended multiple-choice question based on fine-grained labels (i.e., normal tissue, pathological benign, unusual ductal hyperplasia, flat epithelial atypia, atypical ductal hyperplasia, ductal carcinoma in situ, and invasive carcinoma) that asks: “Shown is a breast tissue biopsy image. What is the most likely diagnosis? Answer with the option’s letter from the given choices directly.”

HiCervix⁸⁷ contains twenty-nine classes of cervical cytological cells organized into three hierarchies. We use 8,051 ROIs from the test split. For each ROI, we create a closed-ended multiple-choice question based on the labels in first hierarchy (i.e., atypical squamous cell, negative for intraepithelial lesion or malignancy, atypical glandular cell and organism) that asks: “Shown is an image of cervical cytological cells. Which category do the cells most likely belong to? Answer with the option’s letter from the given choices directly.”

PathQABench contains 150 ROIs paired with clinical context, and asks to select the most likely diagnosis from a set of plausible options, as described in the original study¹² for multiple-choice-style question answering (PathQABench MCQ). To curate PathQABench Caption, we ask a board-certified pathologist to provide the detailed caption per ROI. We append “Answer with the option’s letter from the given choices directly.” and

⁵<https://github.com/PathMMU-Benchmark/PathMMU/blob/main/data/instructions.md#atlas-subset>

“Generate a detailed caption for this image that describes the key morphological features observed. Do not suggest any diagnosis or further testing.” at the end of each question for PathQABench MCQ and PathQABench Caption, respectively.

DDxBench contains 150 H&E stained slides with diverse tissue sites and tumor types, with a total of 52 unique diagnoses (detailed in **Extended Data Table 3**). We used it to curate a multi-image understanding benchmark, with the goal of assessing each MLLM’s ability to make accurate differential diagnoses from representative regions of disease from multi-image understanding. For each slide, the corresponding benchmark question contains 10 representative, tumor ROIs, sampled using a mix of automated tumor classifier labeling and human-in-the-loop review for quality control. The multi-image input, together with the tissue biopsy site and biological sex of the patient is used as input to the model. For evaluating SlideSeek, in order to evaluate its ability to autonomously navigate and assess the slide, the entire WSI, instead of tumor ROIs is directly used as input. For general-purpose models, we explicitly request the JSON schema:{“Primary diagnosis”:str, “Other”:list[str]}, either through structured output if supported by official API or via the system prompt to guide the model response format for evaluation.

Computing hardware and software

All experiments and analyses in the study were performed using Python (version 3.10). For model training, we adapted the official LLaVA-OV⁵⁶ code repository by incorporating our custom vision encoder and multimodal projector implemented using Timm (version 0.9.8). We used up to $24 \times 80\text{GB}$ NVIDIA A100 GPUs configured for multi-node multi-GPU distributed training using PyTorch (version 2.4.1, CUDA 11.8) as the deep learning framework. DeepSpeed (version 0.14.4) was used to enable accelerated training of PathChat+ MLLM. All inference jobs were performed using 24GB NVIDIA 3090 GPUs. Pillow (version 10.2.0) is used for image processing. Matplotlib (version 3.7.1) and Seaborn (version 0.12.2) were used to create plots and figures. SlideSeek was implemented using LangChain (version 0.3.22). All whole slide operations including tissue delineation and slide navigation were implemented using Trident (version 0.1.1)⁸⁴. Other miscellaneous libraries used are listed in the **Reporting Summary**. In total, training the combined system of PathChat+ (including the vision encoder, the multimodal projector and the large language model) took approximately 1,275 A100 GPU hours.

Statistical analysis

We estimated 95% confidence intervals for all reported metrics using nonparametric bootstrapping (1,000 replicates). To assess the statistical significance of observed differences in performance between model pairs, we performed two-sided paired permutation tests (1,000 permutations). The null hypothesis was that no difference exists between the performances of the two models. For each permutation, we randomly swapped pairs of predictions between the two models to generate a new set of performance differences. We computed the p-value as the proportion of permuted differences whose absolute values exceeded the observed difference.

Data availability.

The public data used in our evaluation can be accessed via the official repository: PathMMU, UniToPatho, BRACS, and HiCervix. PathQABench data was curated using a combination of WSIs from TCGA and in-house pathology database at the Brigham and Women’s Hospital. The original TCGA WSIs and associated clinical metadata are available from the NIH genomic data commons (portal.gdc.cancer.gov). The in-house subset was curated with institutional permission through IRB approval for the current study and thus cannot be made publicly available. Similarly, the DDxbench data is also curated from in-house database. All requests for data collected or curated in-house will be evaluated based on institutional and departmental policies to determine whether the data requested is subject to intellectual property or patient privacy obligations. Instruction data was curated from image captions, notes, reports from in-house educational resources and patient data.

Code availability

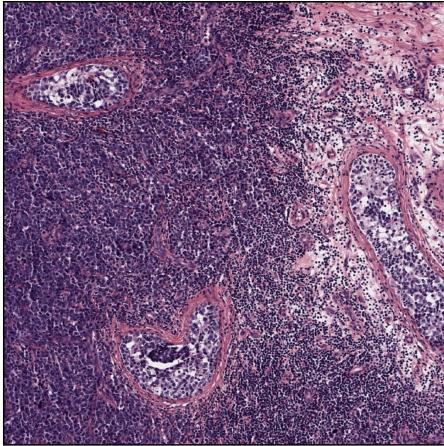
The general open MLLMs used in comparisons can be accessed via their official APIs or repository: OpenAI (model id: gpt-4o-2024-11-20), OpenAI o1 (model id: o1-2024-12-17), OpenAI GPT-4o mini (model id: gpt-4o-mini-2024-07-18), Anthropic (model id: claude-3-5-sonnet-20241022) and Gemini (model id: gemini-2.0-flash-001). Models evaluated in this study can be accessed via their official GitHub repositories: LLaVA-OneVision, Llama 3.2, Qwen2.5-VL, LLaVA-Med, Quilt-LLaVA, PA-LLaVA, and Huatuo-GPT-Vision. Similarly, the weights can be accessed through Hugging Face: LLaVA-OneVision (model id: llava-onevision-qwen2-7b-ov-hf), Llama 3.2 (model id: Llama-3.2-11B-Vision-Instruct), Qwen2.5-VL (model id: Qwen2.5-VL-7B-Instruct), LLaVA-Med (model id: llava-med-v1.5-mistral-7b), Quilt-LLaVA (model id: Quilt-Llava-v1.5-7b), PA-LLaVA (model id: Pathology-LLaVA), and HuatuoGPT-Vision (model id: HuatuoGPT-Vision-7B).

Author contributions

M.Y.L. and F.M. conceived the study and designed the experiments. C.C., M.Y.L., R.J.C., B.C., T.D. and D.F.K.W performed data collection and processing. C.C., L.W. and M.Y.L. performed model development. C.C., L.W., M.Y.L., T.D., D.F.K.W. and R.J.C. performed experimental analysis and interpreted the results. R.J.C., D.F.W.K., A.V., G.J., L.P.L. provided feedback on the analysis. C.C., L.W., M.Y.L., and F.M. prepared the manuscript with input from all co-authors. M.Y.L. and F.M. supervised the research.

Acknowledgements

This work was supported in part by the BWH president's fund, BWH & MGH Pathology. We thank Timothy Janicki, Richard Kenny, Abe Ahmed and the system administration staff at the MGB Enterprise Research Infrastructure & Services (ERIS) Research Computing Core for their dedicated support in providing and maintaining access to NVIDIA A100 computing resources.

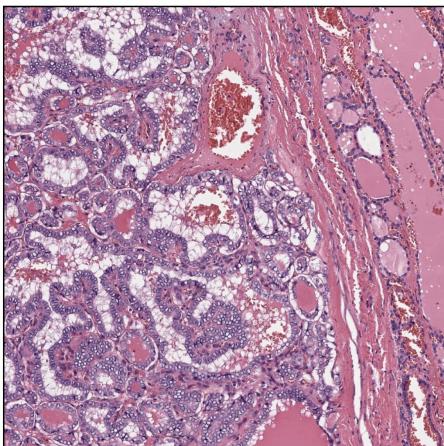


Pathologist: Several testicular tubules surrounded by a sea of lymphocytes though there is some atrophic testis at the right side of the image. The tubules are infiltrated by lymphocytes as well. The lymphocytes are large and atypical with speckled chromatin.

PathChat+: The image shows a dense infiltrate of **atypical lymphoid cells** within the **testicular parenchyma**. These cells are large with irregular nuclear contours and prominent nucleoli. The infiltrate is diffuse and extensive, involving the interstitial spaces and surrounding seminiferous tubules. The **seminiferous tubules** themselves appear compressed and distorted by the infiltrating lymphoid cells. There is a notable absence of normal testicular architecture, with the lymphoid cells replacing the usual cellular components. The background stroma is scant, with minimal fibrous tissue present.

Quilt-Llava: Upon examination of the histopathology image, I observe a lymph node with distinct follicles. These follicles are characterized by a lighter germinal center, which is a typical feature of secondary follicles. The germinal centers appear to be active, indicating a state of immune response.

In the central region of the image, there is a notable presence of tingible body macrophages. These cells are characterized by their cytoplasm filled with nuclear debris, giving them a distinct appearance. The presence of these macrophages suggests an ongoing process of antigen presentation and lymphocyte activation.

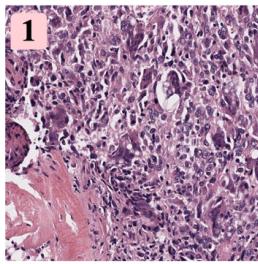


Pathologist: At the right of the image are normal thyroid glands and a thin capsule separating them from the tumor that makes up the bulk of the image. The tumor cells are arranged in papillary structures. They have hypochromatic, cleared nuclei , many of which demonstrate nuclear grooves.

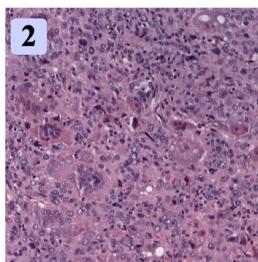
PathChat+: The image shows a section of thyroid tissue with a distinct area of **tumor invasion into the adjacent thyroid parenchyma**. The tumor cells exhibit a **papillary architecture**, characterized by fibrovascular cores lined by neoplastic epithelial cells. These cells have enlarged, overlapping nuclei with a **ground-glass appearance**, **nuclear grooves**, and occasional intranuclear cytoplasmic inclusions. The surrounding thyroid parenchyma displays normal follicular structures with colloid-filled spaces. The interface between the tumor and normal thyroid tissue is irregular, indicating invasive growth.

GPT-4o: The image displays a histological section stained with hematoxylin and eosin (H&E), showcasing key morphological features. There are numerous gland-like structures with prominent lumen, lined by cuboidal to columnar epithelial cells. The nuclei of these cells appear basophilic and are evenly spaced with distinct nucleoli in some areas. Surrounding the glandular structures, there is a pink-staining extracellular matrix indicative of connective tissue. A dilated blood vessel containing erythrocytes is evident in the upper right portion. Additionally, there are regions of fluid-filled spaces and a fibrous, acellular stroma visible in the periphery. The overall architecture suggests extensive cellular organization amid vascular and stromal components.

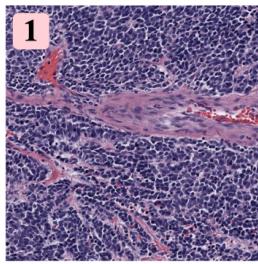
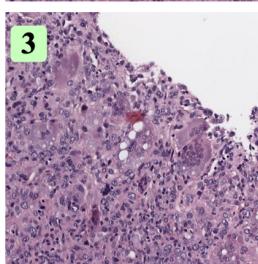
Extended Data Figure 1: Examples of generated captions on PathQABench Caption. In addition to PathChat+, we select the next best performing model for illustrative comparison. In general, PathChat+ is able to capture key diagnostic tissue and nuclear features in the images that are highlighted by the pathologist while other models, despite sometimes being able to describe tissue patterns in the image, often mischaracterize such key features or miss them altogether.



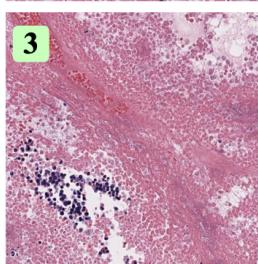
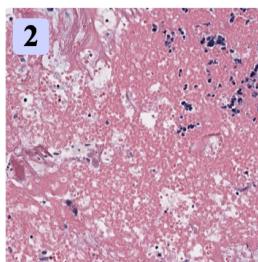
True Dx: Anaplastic Thyroid Cancer
Predicted Dx: Anaplastic Thyroid Cancer
Predicted Confidence: High
Rationale: The integrated high-power and medium-power evaluations across multiple regions strongly support an undifferentiated (anaplastic) thyroid carcinoma. In Region 0, the 20x views from ROIs such as **ROI#1** reveals marked nuclear pleomorphism, hyperchromasia, irregular nuclear contours and disorganized spindle cell and giant cell patterns with areas of necrosis and hemorrhage. Additionally, the high-magnification examination in Tissue Region 1 (**ROI#2**) confirms similar anaplastic features in giant cells, further corroborating a unifying aggressive, high-grade process. Although one high-power view from Tissue Region 1 (**ROI#3**) suggested some reactive features, the overall preponderance of malignant cytological abnormalities strongly aligns with the PathChat prediction. The findings do not favor metastatic carcinoma or thyroid sarcoma as the primary process given the consistent thyroid tissue background in many ROIs, making undifferentiated (anaplastic) thyroid carcinoma the most plausible diagnosis.



True Dx: Anaplastic Thyroid Cancer
Predicted Dx: Anaplastic Thyroid Cancer
Predicted Confidence: High
Rationale: The integrated high-power and medium-power evaluations across multiple regions strongly support an undifferentiated (anaplastic) thyroid carcinoma. In Region 0, the 20x views from ROIs such as **ROI#1** reveals marked nuclear pleomorphism, hyperchromasia, irregular nuclear contours and disorganized spindle cell and giant cell patterns with areas of necrosis and hemorrhage. Additionally, the high-magnification examination in Tissue Region 1 (**ROI#2**) confirms similar anaplastic features in giant cells, further corroborating a unifying aggressive, high-grade process. Although one high-power view from Tissue Region 1 (**ROI#3**) suggested some reactive features, the overall preponderance of malignant cytological abnormalities strongly aligns with the PathChat prediction. The findings do not favor metastatic carcinoma or thyroid sarcoma as the primary process given the consistent thyroid tissue background in many ROIs, making undifferentiated (anaplastic) thyroid carcinoma the most plausible diagnosis.



True Dx: Merkel Cell Carcinoma
Predicted Dx: Small Cell Carcinoma
Predicted Confidence: Low
Rationale: While the regions surveyed at high magnification (for example, region- **ROI#1** and **ROI#2**) consistently show sheets and nests of small, round blue cells with high nuclear-to-cytoplasmic ratios, hyperchromatic nuclei, and focal necrosis, these features are more in keeping with a primary cutaneous neuroendocrine carcinoma, in other words, Merkel cell carcinoma, rather than what is typically termed 'small cell carcinoma of the skin.' Moreover, in **ROI#3**, the tumor architecture raised the possibility of a lymphoid neoplasm, necessitating careful exclusion of cutaneous lymphoma. The overall morphology lacks evidence of glandular or squamous differentiation that might support an aggressive epithelial carcinoma. Thus, I have low confidence that the primary diagnosis should be stated as small cell carcinoma of the skin. More appropriately, the findings lean toward Merkel cell carcinoma as the primary diagnosis with cutaneous lymphoma remaining a consideration, and additional immunohistochemical studies would be required for conclusive classification.



Extended Data Figure 2: Example of SlideSeek classifying a case as high confidence (left) and a case as low confidence (right). The SlideSeek correctly identifies morphological features within the slide that support or refute the predicted diagnosis and supports its argument with specific references to regions of interest (ROIs). *Pathologist's Commentary:* left case: The critic pulls ROIs that appropriately showcase the morphology present in the WSI and which demonstrate morphologies that are consistent with Anaplastic Thyroid Cancer. The image descriptions are accurate. right case: The critic rightly describes how the morphology on display in the ROIs is consistent with Merkel Cell Carcinoma or a cutaneous lymphoma. As small cell carcinoma of the skin is exceedingly rare, the critic is correct in that the diagnosis is much more likely to be one of the aforementioned options.

Hyperparameter	Value
Automatic mixed precision	BF16
DeepSpeed ZeRO	Stage 3
Batch size	128
Learning rate scheduler	Cosine
Warmup ratio	0.03
Peak learning rate	2e-3
AdamW β	(0.9, 0.999)
AdamW ϵ	1e-8
Weight decay	0.
Gradient clipping max. norm	1.0
Training epochs	1
Gradient checkpointing	Yes
TF32	Yes

Extended Data Table 1: Hyperparameters used in PathChat model pretraining. $8 \times 80\text{GB}$ NVIDIA A100 GPUs were used for training. *Batch size* refers to the effective batch size (*i.e.*, the total batch size across GPUs). The learning rate is increased from zero linearly to the *peak learning rate* over the course of *total number batches* \times *warmup ratio* steps and decays back to zero following the *learning rate scheduler*.

Hyperparameter	Value
Automatic mixed precision	BF16
DeepSpeed ZeRO	Stage 3
Batch size per GPU	2
Gradient accumulation steps	2
Learning rate scheduler	Cosine
Warmup ratio	0.03
Peak learning rate	2e-5
AdamW β	(0.9, 0.999)
AdamW ϵ	1e-8
Weight decay	0.
Gradient clipping max. norm	1.0
Training epochs	1
Gradient checkpointing	Yes
TF32	Yes

Extended Data Table 2: Hyperparameters used in PathChat MLLM finetuning. 24 \times 80GB NVIDIA A100 GPUs across 3 compute nodes were used for multi-node distributed training. *Batch size* refers to the total batch size across GPUs. The effective batch size used for optimization is *batch size* \times *gradient accumulation steps*. The learning rate is increased from zero linearly to the *peak learning rate* over the course of *total number batches* \times *warmup ratio* steps and decays back to zero following the *learning rate scheduler*.

Site	Diagnoses
Brain	Anaplastic Astrocytoma, Ependymoma, Glioblastoma, Medulloblastoma, Meningioma, Oligodendrogioma, Pilocytic Astrocytoma, Pituitary Adenoma
Breast	Invasive Breast Carcinoma of No Special Type (NST), Invasive Lobular Carcinoma
Endocrine	Adrenocortical Carcinoma, Anaplastic Thyroid Carcinoma, Follicular Thyroid Carcinoma, Medullary Thyroid Carcinoma, Papillary Thyroid Carcinoma
GI	Colorectal Adenocarcinoma, Esophageal Adenocarcinoma, Esophageal Squamous Cell Carcinoma, Gastric Adenocarcinoma, Gastric Signet Ring Cell Carcinoma, Squamous Cell Carcinoma of the Anus
GYN	Clear Cell Carcinoma of the Ovary, Clear Cell Carcinoma of the Uterus, Endometrioid Carcinoma of the Ovary, Endometrioid Carcinoma of the Uterus, High-Grade Serous Carcinoma of the Ovary, Low-Grade Serous Carcinoma of the Ovary, Mucinous Carcinoma of the Ovary, Serous Borderline Tumor of the Ovary, Serous Carcinoma of the Uterus, Squamous Cell Carcinoma of the Cervix
Lung	Adenocarcinoma of the Lung, Atypical Carcinoid Tumor, Large Cell Neuroendocrine Carcinoma, Small Cell Carcinoma of the Lung, Squamous Cell Carcinoma of the Lung
Male Reproductive Tract	Prostatic Adenocarcinoma, Seminoma
Pancreaticohepatobiliary	Hepatocellular Carcinoma, Pancreatic Adenocarcinoma, Pancreatic Neuroendocrine Tumor (PanNET)
Skin and Connective Tissue	Cutaneous Squamous Cell Carcinoma, Malignant Melanoma, Merkel Cell Carcinoma
Urinary Tract	Adenocarcinoma of the Bladder, Chromophobe Renal Cell Carcinoma, Clear Cell Renal Cell Carcinoma, Nephroblastoma (Wilms Tumor), Papillary Renal Cell Carcinoma, Renal Oncocytoma, Upper Tract Urothelial Carcinoma, Urothelial Carcinoma

Extended Data Table 3: Unique diagnoses by tissue site in DDxBench. GI: Gastrointestinal, GYN: Gynecology.

Model	All-val	All-test-tiny	All-test
Claude 3.5 Sonnet	0.539 (0.501, 0.578)	0.610 (0.583, 0.640)	0.570 (0.560, 0.580)
GPT-4o	0.519 (0.482, 0.557)	0.615 (0.586, 0.643)	0.559 (0.550, 0.569)
Gemini 2.0 Flash	0.605 (0.568, 0.642)	0.676 (0.648, 0.701)	0.638 (0.629, 0.647)
HuatuoGPT-Vision	0.575 (0.536, 0.609)	0.587 (0.561, 0.617)	0.564 (0.554, 0.573)
LLaVA-OneVision	0.428 (0.391, 0.465)	0.489 (0.459, 0.518)	0.452 (0.442, 0.463)
Llama 3.2	0.489 (0.449, 0.528)	0.530 (0.501, 0.558)	0.491 (0.481, 0.501)
PA-LLaVA	0.376 (0.340, 0.415)	0.422 (0.396, 0.451)	0.389 (0.380, 0.399)
PathChat 1	0.579 (0.541, 0.617)	0.640 (0.613, 0.668)	0.609 (0.599, 0.619)
PathChat+	0.692 (0.657, 0.725)	0.750 (0.725, 0.777)	0.728 (0.719, 0.736)
Quilt-LLaVA	0.416 (0.381, 0.454)	0.441 (0.413, 0.470)	0.406 (0.397, 0.416)
Qwen2.5-VL	0.488 (0.449, 0.524)	0.556 (0.529, 0.582)	0.508 (0.498, 0.518)
Expert performance*	-	0.718	-

Extended Data Table 4: Performance on PathMMU multiple-choice questions, validation, test-tiny and test subset. Accuracy is reported on subset of all validation (All-val, $n = 699$), all test tiny (All-test-tiny, $n = 1,150$) and all test (All-test, $n = 9,618$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**. *: The expert performance is quoted from the PathMMU paper⁷³. Note that the quoted expert performance is not directly comparable to all the evaluated model performance numbers since the Atlas subset has been updated by the PathMMU authors and is not the same as reported in the original paper.

Model	PubMed (test-tiny)	PubMed (test-all)
Claude 3.5 Sonnet	0.637 (0.584, 0.698)	0.626 (0.607, 0.642)
GPT-4o	0.644 (0.591, 0.701)	0.591 (0.573, 0.608)
Gemini 2.0 Flash	0.701 (0.648, 0.758)	0.672 (0.654, 0.689)
HuatuoGPT-Vision	0.637 (0.584, 0.690)	0.635 (0.619, 0.653)
LLaVA-OneVision	0.495 (0.438, 0.552)	0.489 (0.471, 0.507)
Llama 3.2	0.580 (0.520, 0.637)	0.535 (0.517, 0.553)
PA-LLaVA	0.427 (0.367, 0.484)	0.401 (0.384, 0.419)
PathChat 1	0.633 (0.580, 0.687)	0.636 (0.620, 0.654)
PathChat+	0.776 (0.730, 0.822)	0.723 (0.707, 0.739)
Quilt-LLaVA	0.438 (0.377, 0.498)	0.416 (0.399, 0.433)
Qwen2.5-VL	0.626 (0.569, 0.683)	0.567 (0.548, 0.584)
Expert performance*	0.729	-

Extended Data Table 5: Performance on PathMMU multiple-choice questions, Pubmed subset. Accuracy is reported on test tiny subset ($n = 281$) and test subset ($n = 3,068$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**. *: The expert performance is quoted from the PathMMU paper⁷³.

Model	SocialPath (test-tiny)	SocialPath (test-all)
Claude 3.5 Sonnet	0.681 (0.624, 0.738)	0.614 (0.591, 0.638)
GPT-4o	0.664 (0.603, 0.725)	0.578 (0.556, 0.599)
Gemini 2.0 Flash	0.681 (0.620, 0.738)	0.628 (0.605, 0.649)
HuatuoGPT-Vision	0.576 (0.511, 0.638)	0.596 (0.575, 0.619)
LLaVA-OneVision	0.489 (0.419, 0.555)	0.489 (0.464, 0.513)
Llama 3.2	0.520 (0.454, 0.585)	0.521 (0.499, 0.545)
PA-LLaVA	0.454 (0.393, 0.520)	0.410 (0.385, 0.433)
PathChat 1	0.638 (0.576, 0.699)	0.614 (0.591, 0.637)
PathChat+	0.690 (0.629, 0.747)	0.684 (0.660, 0.706)
Quilt-LLaVA	0.441 (0.376, 0.507)	0.454 (0.433, 0.478)
Qwen2.5-VL	0.616 (0.550, 0.673)	0.557 (0.536, 0.580)
Expert performance*	0.715	-

Extended Data Table 6: Performance on PathMMU multiple-choice questions, SocialPath subset. Accuracy is reported on test tiny subset ($n = 229$) and test subset ($n = 1,796$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation** in **Methods**. *: The human expert performance is quoted from the PathMMU study⁷³.

Model	EduContent (test-tiny)	EduContent (test-all)
Claude 3.5 Sonnet	0.667 (0.608, 0.722)	0.599 (0.577, 0.618)
GPT-4o	0.659 (0.600, 0.718)	0.577 (0.555, 0.598)
Gemini 2.0 Flash	0.745 (0.694, 0.796)	0.643 (0.622, 0.663)
HuatuoGPT-Vision	0.620 (0.557, 0.682)	0.570 (0.548, 0.591)
LLaVA-OneVision	0.573 (0.514, 0.631)	0.512 (0.489, 0.535)
Llama 3.2	0.565 (0.506, 0.627)	0.507 (0.488, 0.530)
PA-LLaVA	0.424 (0.361, 0.490)	0.421 (0.399, 0.444)
PathChat 1	0.678 (0.616, 0.737)	0.631 (0.609, 0.651)
PathChat+	0.722 (0.663, 0.780)	0.709 (0.689, 0.728)
Quilt-LLaVA	0.498 (0.431, 0.557)	0.437 (0.413, 0.460)
Qwen2.5-VL	0.612 (0.553, 0.671)	0.557 (0.535, 0.580)
Expert performance*	0.690	-

Extended Data Table 7: Performance on PathMMU multiple-choice questions, EduContent subset. Accuracy is reported on test tiny subset ($n = 255$) and test subset ($n = 1,938$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation** in **Methods**. *: The expert performance is quoted from the PathMMU paper⁷³.

Model	Atlas (test-tiny)	Atlas (test-all)
Claude 3.5 Sonnet	0.615 (0.553, 0.683)	0.607 (0.575, 0.637)
GPT-4o	0.635 (0.562, 0.702)	0.592 (0.558, 0.623)
Gemini 2.0 Flash	0.630 (0.562, 0.692)	0.649 (0.621, 0.677)
HuatuoGPT-Vision	0.678 (0.611, 0.745)	0.637 (0.607, 0.667)
LLaVA-OneVision	0.601 (0.538, 0.668)	0.541 (0.511, 0.573)
Llama 3.2	0.615 (0.553, 0.683)	0.608 (0.578, 0.637)
PA-LLaVA	0.476 (0.409, 0.548)	0.441 (0.409, 0.473)
PathChat 1	0.788 (0.731, 0.841)	0.753 (0.727, 0.779)
PathChat+	0.817 (0.764, 0.870)	0.811 (0.786, 0.834)
Quilt-LLaVA	0.495 (0.428, 0.567)	0.476 (0.443, 0.506)
Qwen2.5-VL	0.553 (0.486, 0.620)	0.572 (0.541, 0.601)
Expert performance*	0.683	-

Extended Data Table 8: Performance on PathMMU multiple-choice questions, Atlas subset. Accuracy is reported on test tiny subset ($n = 208$) and test subset ($n = 1,007$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**. *: The expert performance is quoted from the PathMMU study⁷³. Note that the quoted expert performance is not directly comparable to all the evaluated model performance numbers since the Atlas subset has been updated by the PathMMU authors and is not the same as reported in the original paper.

Model	PathCLS (test-tiny)	PathCLS (test-all)
Claude 3.5 Sonnet	0.390 (0.316, 0.463)	0.383 (0.361, 0.403)
GPT-4o	0.418 (0.345, 0.492)	0.451 (0.431, 0.474)
Gemini 2.0 Flash	0.582 (0.514, 0.650)	0.579 (0.557, 0.604)
HuatuoGPT-Vision	0.367 (0.299, 0.441)	0.363 (0.341, 0.384)
LLaVA-OneVision	0.226 (0.164, 0.288)	0.239 (0.220, 0.261)
Llama	0.311 (0.249, 0.379)	0.307 (0.287, 0.328)
PA-LLaVA	0.305 (0.237, 0.379)	0.286 (0.266, 0.307)
PathChat 1	0.424 (0.356, 0.497)	0.452 (0.430, 0.475)
PathChat+	0.751 (0.689, 0.814)	0.753 (0.734, 0.774)
Quilt-LLaVA	0.299 (0.232, 0.367)	0.271 (0.252, 0.293)
Qwen2.5-VL	0.288 (0.220, 0.362)	0.270 (0.249, 0.289)
Expert performance	0.789	-

Extended Data Table 9: Performance on PathMMU multiple-choice questions, PathCLS subset. Accuracy is reported on test tiny subset ($n = 177$) and test subset ($n = 1,809$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**. *: The expert performance is quoted from the PathMMU study⁷³.

Model	BRACS	UniToPatho	HiCervix
Claude 3.5 Sonnet	0.239 (0.204, 0.272)	0.291 (0.274, 0.311)	0.285 (0.276, 0.296)
GPT-4o	0.232 (0.196, 0.265)	0.279 (0.261, 0.296)	0.351 (0.340, 0.361)
Gemini 2.0 Flash	0.316 (0.277, 0.353)	0.425 (0.405, 0.444)	0.421 (0.410, 0.432)
HuatuoGPT-Vision	0.198 (0.167, 0.228)	0.311 (0.294, 0.329)	0.371 (0.360, 0.381)
LLaVA-OneVision	0.147 (0.118, 0.175)	0.360 (0.341, 0.378)	0.376 (0.366, 0.386)
Llama 3.2	0.200 (0.170, 0.235)	0.212 (0.196, 0.229)	0.293 (0.282, 0.302)
PA-LLaVA	0.161 (0.132, 0.191)	0.242 (0.226, 0.260)	0.330 (0.320, 0.340)
PathChat 1	0.558 (0.516, 0.596)	0.506 (0.485, 0.526)	0.368 (0.356, 0.379)
PathChat+	0.630 (0.588, 0.665)	0.552 (0.532, 0.570)	0.725 (0.716, 0.736)
Quilt-LLaVA	0.160 (0.132, 0.189)	0.207 (0.191, 0.222)	0.372 (0.362, 0.382)
Qwen2.5-VL	0.165 (0.135, 0.195)	0.354 (0.336, 0.372)	0.258 (0.249, 0.268)

Extended Data Table 10: Performance on multiple-choice questions of BRACS, UniToPatho and HiCervix. Accuracy is reported on BRACS ($n = 570$), UniToPatho ($n = 2,399$) and HiCervix ($n = 8,051$). 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**.

Model	PathQABench MCQ
Claude 3.5 Sonnet	0.790 (0.714, 0.867)
GPT-4o	0.752 (0.667, 0.829)
Gemini 2.0 Flash	0.810 (0.733, 0.886)
HuatuoGPT-Vision	0.571 (0.476, 0.667)
LLaVA-OneVision	0.257 (0.171, 0.343)
Llama 3.2	0.733 (0.647, 0.810)
PA-LLaVA	0.324 (0.238, 0.410)
PathChat 1	0.895 (0.829, 0.952)
PathChat+	0.933 (0.876, 0.981)
Quilt-LLaVA	0.333 (0.248, 0.429)
Qwen2.5-VL	0.619 (0.524, 0.705)

Extended Data Table 11: Performance on multiple-choice questions of PathQABench MCQ. Accuracy is reported on PathQABench MCQ ($n = 105$) with 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation in Methods**.

Model	PathQABench Caption
Claude 3.5 Sonnet	0.227 (0.217, 0.237)
GPT-4o	0.242 (0.232, 0.254)
Gemini 2.0 Flash	0.264 (0.249, 0.278)
HuatuoGPT-Vision	0.223 (0.214, 0.232)
LLaVA-Med	0.190 (0.181, 0.199)
LLaVA-OneVision	0.116 (0.105, 0.127)
Llama 3.2	0.180 (0.173, 0.187)
PA-LLaVA	0.005 (0.002, 0.009)
PathChat 1	0.270 (0.253, 0.286)
PathChat+	0.294 (0.280, 0.309)
Quilt-LLaVA	0.187 (0.179, 0.195)
Qwen2.5-VL	0.209 (0.199, 0.218)

Extended Data Table 12: Performance on PathQABench Caption for image captioning. METEOR score⁷⁷ is reported on PathQABench Caption ($n = 105$) with 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation** in **Methods**.

Model	Primary Diagnosis	Primary + Add. Differentials
Claude 3.5 Sonnet	0.440 (0.367, 0.520)	0.620 (0.547,0.694)
GPT-4o	0.460 (0.380, 0.540)	0.707 (0.633,0.780)
Gemini 2.0 Flash	0.467 (0.380, 0.547)	0.660 (0.580,0.740)
HuatuoGPT-Vision	0.333 (0.260, 0.413)	0.540 (0.460,0.620)
LLaVA-OneVision	0.107 (0.067, 0.160)	0.167 (0.113,0.227)
PathChat 1	0.720 (0.647, 0.787)	0.887 (0.833,0.933)
PathChat+	0.800 (0.740, 0.860)	0.933 (0.893,0.967)
Qwen2.5-VL	0.300 (0.227, 0.373)	0.467 (0.387,0.547)

Extended Data Table 13: Performance on DDxBench. Primary Diagnosis and Primary + Add. Differentials accuracy are reported on DDxBench ($n = 150$) with 95% confidence intervals from bootstrapping are included in parentheses. For more details see **MLLM evaluation** in **Methods**.

Model	Primary Diagnosis	Primary + Add. Differentials
SlideSeek	0.800 (0.733, 0.860)	0.920 (0.873, 0.960)
SlideSeek w. 4o supervisor	0.720 (0.647, 0.793)	0.913 (0.867, 0.953)
SlideSeek w. 4o supervisor + 4o-mini explorer	0.687 (0.613, 0.760)	0.880 (0.827, 0.933)

Extended Data Table 14: Performance of SlideSeek on DDxBench. Primary Diagnosis and Primary + Add. Differentials accuracy are reported on DDxBench ($n = 150$) with 95% confidence intervals from bootstrapping are included in parentheses. For more details see **Evaluation of multi-agent system for interpretable slide assessment** in **Methods**.

References

1. Madabhushi, A. & Lee, G. Image analysis and machine learning in digital pathology: Challenges and opportunities. *Medical image analysis* **33**, 170–175 (2016).
2. Bera, K., Schalper, K. A., Rimm, D. L., Velcheti, V. & Madabhushi, A. Artificial intelligence in digital pathology—new tools for diagnosis and precision oncology. *Nature reviews Clinical oncology* **16**, 703–715 (2019).
3. Heinz, C. N., Echle, A., Foersch, S., Bychkov, A. & Kather, J. N. The future of artificial intelligence in digital pathology—results of a survey across stakeholder groups. *Histopathology* **80**, 1121–1127 (2022).
4. Cui, M. & Zhang, D. Y. Artificial intelligence and computational pathology. *Laboratory Investigation* **101**, 412–422 (2021).
5. Abels, E. *et al.* Computational pathology definitions, best practices, and recommendations for regulatory guidance: a white paper from the digital pathology association. *The Journal of pathology* **249**, 286–294 (2019).
6. Song, A. H. *et al.* Artificial intelligence for digital and computational pathology. *Nature Reviews Bioengineering* 1–20 (2023).
7. Shmatko, A., Ghaffari Laleh, N., Gerstung, M. & Kather, J. N. Artificial intelligence in histopathology: enhancing cancer research and clinical oncology. *Nature Cancer* **3**, 1026–1038 (2022).
8. Lipkova, J. *et al.* Artificial intelligence for multimodal data integration in oncology. *Cancer cell* **40**, 1095–1110 (2022).
9. Waqas, A. *et al.* Revolutionizing digital pathology with the power of generative artificial intelligence and foundation models. *Laboratory Investigation* 100255 (2023).
10. Coudray, N. *et al.* Classification and mutation prediction from non–small cell lung cancer histopathology images using deep learning. *Nature medicine* **24**, 1559–1567 (2018).
11. Campanella, G. *et al.* Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine* **25**, 1301–1309 (2019).
12. Lu, M. Y. *et al.* A multimodal generative ai copilot for human pathology. *Nature* **634**, 466–473 (2024).
13. Li, C. *et al.* Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems* **36**, 28541–28564 (2023).
14. Seyfioglu, M. S., Ikezogwo, W. O., Ghezloo, F., Krishna, R. & Shapiro, L. Quilt-llava: Visual instruction tuning by extracting localized narratives from open-source histopathology videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13183–13192 (2024).
15. Dai, D. *et al.* Pa-llava: A large language-vision assistant for human pathology image understanding. In *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 3138–3143 (IEEE, 2024).
16. Sun, Y. *et al.* Cpath-omni: A unified multimodal foundation model for patch and whole-slide image analysis in computational pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2025). URL https://openaccess.thecvf.com/content/CVPR2025/html/Sun_CPath-Omni_A_Unified_Multimodal_Foundation_Model_for_Patch_and_Whole_CVPR_2025_paper.html.
17. Zhang, W. *et al.* Patho-r1: A multimodal reinforcement learning-based pathology expert reasoner. *arXiv preprint arXiv:2505.11404* (2025). URL <https://arxiv.org/abs/2505.11404>.

18. Lu, M. Y. *et al.* Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering* **5**, 555–570 (2021).
19. Bulten, W. *et al.* Automated deep-learning system for gleason grading of prostate cancer using biopsies: a diagnostic study. *The Lancet Oncology* **21**, 233–241 (2020).
20. Nagpal, K. *et al.* Development and validation of a deep learning algorithm for improving gleason scoring of prostate cancer. *NPJ digital medicine* **2**, 48 (2019).
21. Huang, S.-C. *et al.* Deep neural network trained on gigapixel images improves lymph node metastasis detection in clinical settings. *Nature communications* **13**, 3347 (2022).
22. Ehteshami Bejnordi, B. *et al.* Diagnostic Assessment of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer. *JAMA* **318**, 2199–2210 (2017).
23. Kather, J. N. *et al.* Pan-cancer image-based detection of clinically actionable genetic alterations. *Nature cancer* **1**, 789–799 (2020).
24. Fu, Y. *et al.* Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis. *Nature cancer* **1**, 800–810 (2020).
25. Saldanha, O. L. *et al.* Self-supervised attention-based deep learning for pan-cancer mutation prediction from histopathology. *NPJ Precision Oncology* **7**, 35 (2023).
26. Wagner, S. J. *et al.* Transformer-based biomarker prediction from colorectal cancer histology: A large-scale multicentric study. *Cancer Cell* **41**, 1650–1661 (2023).
27. Beck, A. H. *et al.* Systematic analysis of breast cancer morphology uncovers stromal features associated with survival. *Science translational medicine* **3** (2011).
28. Mobadersany, P. *et al.* Predicting cancer outcomes from histology and genomics using convolutional networks. *Proceedings of the National Academy of Sciences* **115**, E2970–E2979 (2018).
29. Chen, R. J. *et al.* Pan-cancer integrative histology-genomic analysis via multimodal deep learning. *Cancer Cell* **40**, 865–878 (2022).
30. Lee, Y., Park, J., Oh, S. *et al.* Derivation of prognostic contextual histopathological features from whole-slide images of tumours via graph deep learning. *Nat. Biomed. Eng* (2022).
31. Courtiol, P. *et al.* Deep learning-based classification of mesothelioma improves prediction of patient outcome. *Nature medicine* **25**, 1519–1525 (2019).
32. Lu, C. *et al.* A prognostic model for overall survival of patients with early-stage non-small cell lung cancer: a multicentre, retrospective study. *The Lancet Digital Health* **2**, e594–e606 (2020).
33. Amgad, M. *et al.* A population-level digital histologic biomarker for enhanced prognosis of invasive breast cancer. *Nature Medicine* **1**–13 (2023).
34. Boehm, K. M. *et al.* Multimodal data integration using machine learning improves risk stratification of high-grade serous ovarian cancer. *Nature cancer* **3**, 723–733 (2022).
35. Sammut, S.-J. *et al.* Multi-omic machine learning predictor of breast cancer therapy response. *Nature* **601**, 623–629 (2022).
36. Vanguri, R. S. *et al.* Multimodal integration of radiology, pathology and genomics for prediction of response to pd-(l) 1 blockade in patients with non-small cell lung cancer. *Nature cancer* **3**, 1151–1164 (2022).

37. Huang, Z. *et al.* Artificial intelligence reveals features associated with breast cancer neoadjuvant chemotherapy responses from multi-stain histopathologic images. *NPJ Precision Oncology* **7**, 14 (2023).
38. Lu, M. Y. *et al.* Ai-based pathology predicts origins for cancers of unknown primary. *Nature* **594**, 106–110 (2021).
39. Zhu, L. *et al.* An accurate prediction of the origin for bone metastatic cancer using deep learning on digital pathological images. *EBioMedicine* **87** (2023).
40. Chen, C. *et al.* Fast and scalable search of whole-slide images via self-supervised deep learning. *Nature Biomedical Engineering* **6**, 1420–1434 (2022).
41. Kalra, S. *et al.* Yottixel—an image search engine for large archives of histopathology whole slide images. *Medical Image Analysis* **65**, 101757 (2020).
42. Hegde, N. *et al.* Similar image search for histopathology: Smily. *NPJ digital medicine* **2**, 56 (2019).
43. Wang, X. *et al.* Retccl: clustering-guided contrastive learning for whole-slide image retrieval. *Medical image analysis* **83**, 102645 (2023).
44. Yala, A. *et al.* Optimizing risk-based breast cancer screening policies with reinforcement learning. *Nature medicine* **28**, 136–143 (2022).
45. Zhou, Y. *et al.* Multi-site cross-organ calibrated deep learning (muscl): Automated diagnosis of non-melanoma skin cancer. *Medical image analysis* **84**, 102702 (2023).
46. Laleh, N. G. *et al.* Benchmarking weakly-supervised deep learning pipelines for whole slide classification in computational pathology. *Medical image analysis* **79** (2022).
47. Graham, S. *et al.* Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical image analysis* **58**, 101563 (2019).
48. Graham, S. *et al.* One model is all you need: multi-task learning enables simultaneous histology image segmentation and classification. *Medical Image Analysis* **83**, 102685 (2023).
49. Lu, M. Y. *et al.* Visual language pretrained multiple instance zero-shot transfer for histopathology images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 19764–19775 (2023).
50. Ikezogwo, W. *et al.* Quilt-1m: One million image-text pairs for histopathology. *Advances in neural information processing systems* **36**, 37995–38017 (2023).
51. Zhang, S. *et al.* A multimodal biomedical foundation model trained from fifteen million image–text pairs. *NEJM AI* **2**, A1oa2400640 (2025).
52. Ouyang, L. *et al.* Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **35**, 27730–27744 (2022).
53. Brown, T. *et al.* Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020).
54. Dubey, A. *et al.* The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
55. Team, G. *et al.* Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118* (2024).
56. Liu, H., Li, C., Wu, Q. & Lee, Y. J. Visual instruction tuning. In *NeurIPS* (2023).

57. Alayrac, J.-B. *et al.* Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems* **35**, 23716–23736 (2022).
58. Yang, Z. *et al.* The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421* **9**, 1 (2023).
59. Wang, P. *et al.* Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191* (2024).
60. Li, B. *et al.* Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326* (2024).
61. Chen, Y. *et al.* Slidechat: A large vision-language assistant for whole-slide pathology image understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 5134–5143 (2025).
62. Gao, Z. *et al.* Alpaca: Adapting llama for pathology context analysis to enable slide-level question answering. *medRxiv* 2025–04 (2025).
63. Tran, M. *et al.* Generating dermatopathology reports from gigapixel whole slide images with histogpt. *Nature Communications* **16**, 1–17 (2025).
64. Shaikovski, G. *et al.* Prism: A multi-modal generative foundation model for slide-level histopathology. *arXiv preprint arXiv:2405.10254* (2024).
65. Ding, T. *et al.* Multimodal whole slide foundation model for pathology. *arXiv preprint arXiv:2411.19666* (2024).
66. Ahmed, F. *et al.* Pathalign: A vision-language model for whole slide images in histopathology. *arXiv preprint arXiv:2406.19578* (2024).
67. Jaech, A. *et al.* Openai o1 system card. *arXiv preprint arXiv:2412.16720* (2024).
68. Guo, D. *et al.* Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
69. Thandiackal, K. *et al.* Differentiable zooming for multiple instance learning on whole-slide images. In *European Conference on Computer Vision*, 699–715 (Springer, 2022).
70. Ghezloo, F. *et al.* Pathfinder: A multi-modal multi-agent system for medical diagnostic decision-making applied to histopathology. *arXiv preprint arXiv:2502.08916* (2025).
71. Buzzard, Z., Hemker, K., Simidjievski, N. & Jamnik, M. Paths: A hierarchical transformer for efficient whole slide image analysis. *arXiv preprint arXiv:2411.18225* (2024).
72. Sun, Y. *et al.* Cpathagent: An agent-based foundation model for interpretable high-resolution pathology image analysis mimicking pathologists’ diagnostic logic. *arXiv preprint arXiv:2505.20510* (2025). URL <https://arxiv.org/abs/2505.20510>.
73. Sun, Y. *et al.* Pathmmu: A massive multimodal expert-level benchmark for understanding and reasoning in pathology. In *European Conference on Computer Vision*, 56–73 (Springer, 2024).
74. Hurst, A. *et al.* Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).
75. Bai, S. *et al.* Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923* (2025).
76. Chen, J. *et al.* Huatuogpt-vision, towards injecting medical visual knowledge into multimodal llms at scale (2024). URL <https://arxiv.org/abs/2406.19280>.
77. Banerjee, S. & Lavie, A. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 65–72 (2005).

78. Lu, M. Y. *et al.* Towards a visual-language foundation model for computational pathology. *arXiv preprint arXiv:2307.12914* (2023).
79. Yang, A. *et al.* Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115* (2024).
80. Lin, T.-Y. *et al.* Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, 740–755 (Springer, 2014).
81. Vaidya, A. *et al.* Molecular-driven foundation model for oncologic pathology. *arXiv preprint arXiv:2501.16652* (2025).
82. Liu, H. *et al.* Llava-next: Improved reasoning, ocr, and world knowledge, january 2024. URL <https://llavavl.github.io/blog/2024-01-30-llava-next> **1** (2024).
83. Tong, P. *et al.* Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *Advances in Neural Information Processing Systems* **37**, 87310–87356 (2024).
84. Zhang, A., Jaume, G., Vaidya, A., Ding, T. & Mahmood, F. Accelerating data processing and benchmarking of ai models for pathology. *arXiv preprint arXiv:2502.06750* (2025).
85. Barbano, C. A. *et al.* Unitopatho, a labeled histopathological dataset for colorectal polyps classification and adenoma dysplasia grading. In *2021 IEEE International Conference on Image Processing (ICIP)*, 76–80 (2021).
86. Brancati, N. *et al.* Bracs: A dataset for breast carcinoma subtyping in h&e histology images. *Database* **2022**, baac093 (2022).
87. Cai, D. *et al.* Hicervix: An extensive hierarchical dataset and benchmark for cervical cytology classification. *IEEE Transactions on Medical Imaging* 1–1 (2024).