# Beyond Pixel Simulation: Pathology Image Generation via Diagnostic Semantic Tokens and Prototype Control

**Minghao Han**[1,2,*], **YiChen Liu**[3,*], **Yizhou Liu**[1,2], **Zizhi Chen**[1,2]
**Jingqun Tang**[4], **Xuecheng Wu**[5], **Dingkang Yang**[1,2,§], **Lihua Zhang**[1,2,§]

[1]College of Intelligent Robotics and Advanced Manufacturing, Fudan University

[2]Fysics Intelligence Technologies Co., Ltd. (Fysics AI)

[3]School of Intelligent Science and Technology, University of Science and Technology Beijing

[4]ByteDance    [5]School of Computer Science and Technology, Xi'an Jiaotong University

[*]Equal contribution, [§]Corresponding author

## Abstract

In computational pathology, understanding and generation have evolved along disparate paths: advanced understanding models already exhibit diagnostic-level competence, whereas generative models largely simulate pixels. Progress remains hindered by three coupled factors: the scarcity of large, high-quality image–text corpora; the lack of precise, fine-grained semantic control, which forces reliance on non-semantic cues; and terminological heterogeneity, where diverse phrasings for the same diagnostic concept impede reliable text conditioning. We introduce UNIPATH, a semantics-driven pathology image generation framework that leverages mature diagnostic understanding to enable controllable generation. UNIPATH implements Multi-Stream Control: a Raw-Text stream; a High-Level Semantics stream that uses learnable queries to a frozen pathology MLLM to distill paraphrase-robust Diagnostic Semantic Tokens and to expand prompts into diagnosis-aware attribute bundles; and a Prototype stream that affords component-level morphological control via a prototype bank. On the data front, we curate a 2.65M image–text corpus and a finely annotated, high-quality 68K subset to alleviate data scarcity. For a comprehensive assessment, we establish a four-tier evaluation hierarchy tailored to pathology. Extensive experiments demonstrate UNIPATH's SOTA performance, including a Patho-FID of 80.9 (51% better than the second-best) and fine-grained semantic control achieving 98.7% of the real-image. The meticulously curated datasets, complete source code, and pre-trained model weights developed in this study will be made openly accessible to the public.

**Date:** December 25, 2025

**Corresponding:** mhhan22@m.fudan.edu.cn,    dicken@fysics.ai,    lihuazhang@fudan.edu.cn

## 1 Introduction

Gigapixel Whole Slide Images (WSIs) sit at the core of modern cancer diagnostics and are being reshaped by foundation models [38, 40, 42]. Yet pathology AI is evolving along two largely disparate paths: (i) understanding models increasingly capture diagnostic-grade signals across tasks [9, 24, 36, 37], while (ii) generative models predominantly pursue perceptual realism for augmentation with weak diagnostic conditioning and limited semantic control [1, 3, 18, 32, 47]. This gap matters as without diagnosis-aware conditioning, generators tend to optimize for appearance cues rather than pathology-grounded semantics.

In practice, progress in pathology text-to-image generation is impeded by three compounding bottlenecks: **(i) Data scarcity**. Large, high-quality image–text corpora are rare because WSIs are gigapixel-scale and subspecialist annotation is costly, which constrains semantic supervision. **(ii) Lack of precise semantic control**. Because semantic labels are scarce, SOTA methods rely on non-semantic controls (segmentation masks [11, 43] or reference images [47]); when present, text-based control is usually restricted to a single cancer with coarse-grained labels [46]. **(iii) Terminological heterogeneity**. The same diagnostic concept is expressed with institution- and pathologist-specific phrasing. General-purpose text encoders [30, 31] struggle to align these variants to a consistent meaning, making text conditioning unreliable and weakening semantic control.

We tackle these challenges by unifying diagnostic understanding and image generation in a single model. We present UNIPATH, which converts diagnosis-aware semantics into fine-grained, pathology-relevant image generation via **Multi-Stream Control (MSC)** while



**Figure 1** UNIPATH achieves overall leading performance against state-of-the-art baselines across our four-tier evaluation hierarchy.

retaining strong diagnostic understanding. MSC comprises three streams: **Raw-Text Stream (RTS)**, which preserves and forwards the user prompt; **High-Level Semantics Stream (HLS)**, instantiated by learnable queries to a pathology MLLM, distills **Diagnostic Semantic Tokens (DST)** robust to paraphrase and reporting style and expands surface prompts into diagnosis-aware attribute bundles; and **Prototype Stream (PS)**, which conditions the generator on retrieved morphology primitives from a prototype bank, enabling component-level control over glandular architecture, nuclear atypia and other key attributes. As a result, UNIPATH unifies diagnostic understanding with semantically controllable generation within a single model, but fully capitalising on this still hinges on abundant and high-quality data.

To this end, we curate a large-scale corpus balancing quantity and diagnostic richness. Starting from 69,044 WSIs in HISTAI [25], rigorous quality control and diversity filtering yield 1.03M diagnostically rich patches. We generate descriptions with a pathology MLLM and merge the resulting pairs with public data to form a 2.65M image–text corpus. For high-fidelity use, we annotate a 68K refined subset with Gemini-2.5 Pro [10] and GPT-5 [27].

To comprehensively evaluate UNIPATH's controllability and semantic fidelity, we establish a **Four-Tier Evaluation Hierarchy** tailored for controllable generation in pathology, including **Visual Fidelity**, **Text-Image Alignment**, **Fine-grained Semantic Control**, and **Downstream Task Utility**. Extensive experiments based on this hierarchy, summarized in Figure 1, demonstrate that UNIPATH achieves SOTA generative fidelity, highly competitive image–text alignment, and exceptional concept-level controllability, while retaining strong diagnostic understanding. In summary, our main contributions are:

- We present UNIPATH, a unified large multimodal model that couples a pathology understanding module with a controllable generator, enabling semantics-driven pathology image generation while preserving understanding.

- We design a multi-stream control architecture that combines a high-level semantics stream to mitigate terminological heterogeneity and a prototype stream to enable fine-grained, component-level morphological control.

- We curate a large-scale, high-quality corpus of 2.65M image–text pairs and a refined 68K subset with pathology-aware quality control to support training and evaluation.

- We conduct a comprehensive evaluation across fidelity, alignment, controllability, and downstream utility, demonstrating strong image fidelity, robust image–text alignment, and superior concept-level controllability.
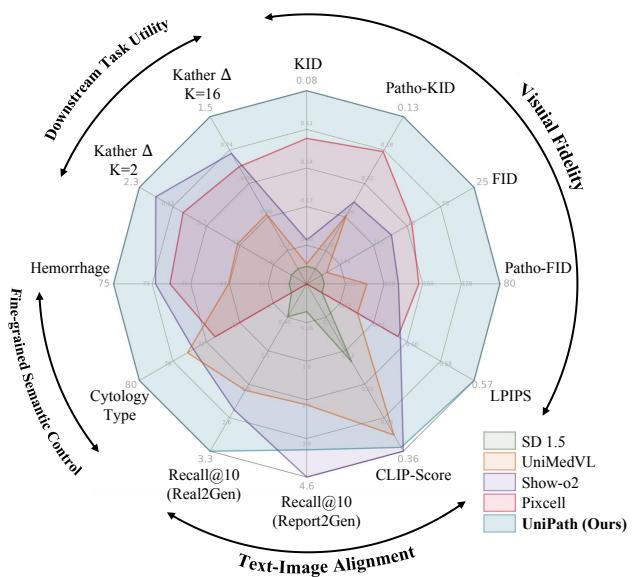
## 2 Related Work

**Foundation Models in Pathology.** Foundation models in pathology are evolving along two disparate paths. On the understanding side, **Pathology MLLMs** have advanced rapidly, approaching diagnostic-grade capability [7, 9, 19, 24, 36, 44]. Progress is driven by high-quality data, such as Quilt-LLaVA's [34] instructional narratives, and by novel architectures like CPath-Omni [37] (unifying multiple tasks) and Patho-R1 [50] (using RL for reasoning). These works confirm that MLLMs capture complex diagnostic concepts, laying the groundwork for leveraging their strong understanding to guide generation. In contrast, **Pathology Generative Models** still rely heavily on non-semantic controls to circumvent core challenges. For instance, ToPoFM [18] imposes topology-guided constraints, while others use masks [3] or reference images [47]. Text-conditioned attempts remain coarse, such as PathLDM [46], which is limited to a single cancer type and provides only coarse-grained textual control. This reliance on general-purpose encoders also fails to normalize pathology phrasing. UNIPATH addresses these gaps via multi-stream control and large-scale, pathology-aware data curation.

**Unified Generation and Understanding.** Unified models in general AI, such as BLIP3o [5], BAGEL [12], and Show-o2 [41], now couple strong understanding with generation to steer synthesis [20, 39]. Transferring this to pathology is non-trivial, as general-purpose models lack the two key prerequisites for this domain: (i) diagnosis-aware semantics to handle terminological heterogeneity, and (ii) component-level morphology control for specific histological structures. UNIPATH adopts this unified design by introducing multi-stream control to address both challenges: an HLS stream provides diagnosis-aware semantics, while a PS stream enables morphological control.

## 3 Data Curation

To address the data scarcity bottleneck outlined in Section 1, we built two key corpora to support the training and evaluation of UNIPATH: (i) a 2.65M-pair large-scale corpus for broad visual-textual alignment, and (ii) a 68K high-quality subset for fine-tuning and evaluation.

### 3.1 Large-Scale Pre-training Corpus

Our 2.65M pre-training corpus consists of two components: 1.62M pairs from public data [36] and 1.03M high-information patches extracted and annotated from 69,044 HISTAI WSIs [25], thereby increasing data diversity.

**Representative Patch Selection.** A dual-strategy pipeline is employed to curate this 1.03M corpus, combining knowledge-guided retrieval and data-driven clustering. First, we use TRIDENT [48] to tile WSIs into $384 \times 384$ patches at $20\times$ magnification. Features from all patches are then extracted using CONCH [23], and two parallel strategies are applied to balance relevance and diversity:

- **Knowledge-Guided Retrieval:** To capture diagnosis-relevant, information-rich patches, we adopt a unified, knowledge-guided retrieval scheme driven by a powerful LLM. For each WSI, LLM processes metadata in parallel: (i) it uses the diagnosis and organ source to generate organ-specific visual feature descriptions; (ii) it processes the original microscopic examination conclusion. Because this conclusion is often long and exceeds CONCH's context limit, LLM splits it into multiple short, complete, and independent passages. We then use CONCH to compute the similarity between those queries and all patches, retrieving the clearest, most typical, and most diagnosis-relevant regions.

- **K-means Cluster Sampling:** In parallel, K-means clustering is performed on all patch features. Samples are then drawn from each cluster to ensure morphological diversity and cover the long-tail distribution. The number of clusters depends on the size of the WSI.

**Description Generation.** Then we de-duplicate patches with UNI2-h [8] visual feature similarity > 0.95. Patch-level descriptions were then generated by PathGen-LLaVA [36] and subsequently summarized by Qwen3-8B [45] into refined, information-dense final descriptions. WSI tiling and filtering ran for a month using one NVIDIA A800 GPU. The subsequent annotation consumed nearly four days on 16 NVIDIA H100 GPUs.

### 3.2 High Quality Refined Subset

To enable high-fidelity fine-tuning, reliable evaluation, and prototype bank construction, a 68K subset was filtered.
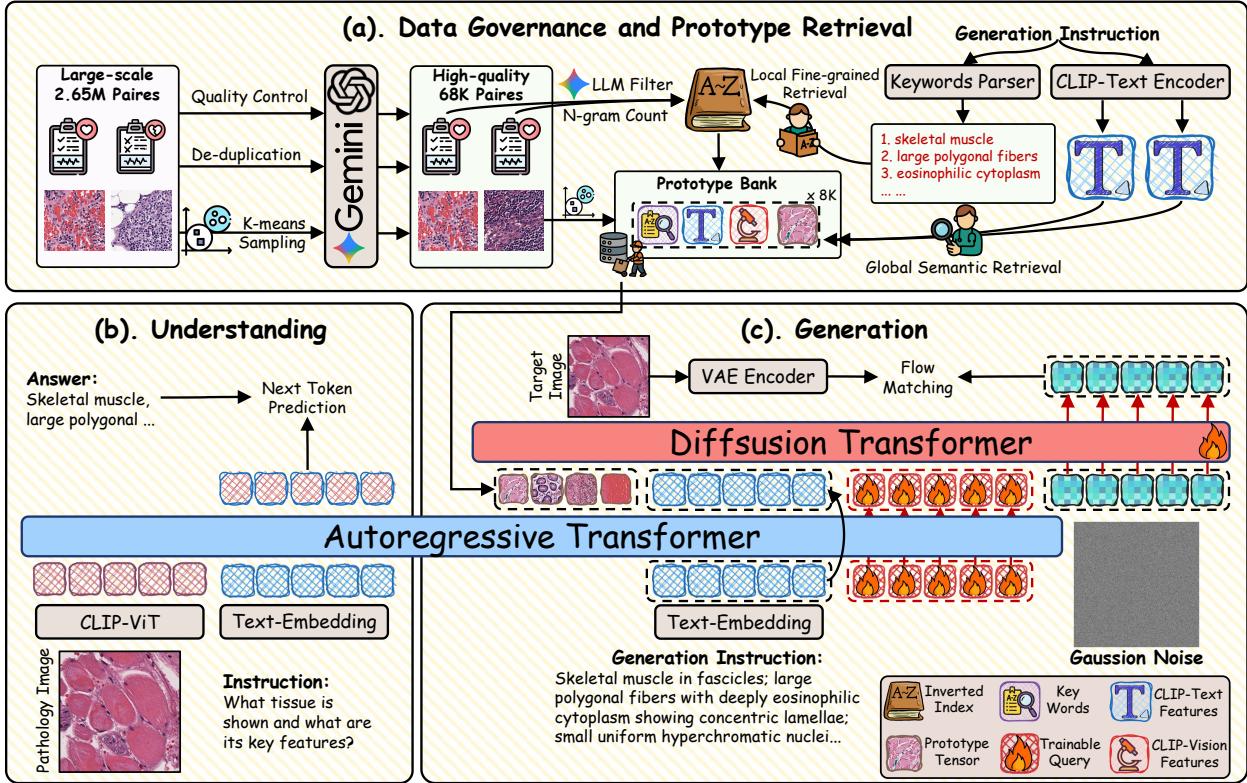
**Figure 2 Overview of UNIPATH**: Unifying pathology understanding and synthesis via Multi-Stream Control. **(a)**: Data governance and prototype retrieval. **(b)**: A frozen pathology MLLM (Understanding) steers a **(c)**: generative DiT (Generation) through MSC fusing raw-text, high-level semantics (emitting diagnostic semantic tokens via learnable queries), and prototype cues.

**Filtering Pipeline.** Our pipeline first clusters the 2.6M samples via K-means (k=128) on their UNI2-h [8] visual features. Subsequently, we compute the Laplacian variance and discard the bottom 50% with the lowest sharpness. To efficiently filter high-quality, representative samples from it, we employ a two-stage automated sampling strategy: first, we perform proportional random sampling within each cluster to obtain a morphologically diverse and manageable candidate pool. This pool then undergoes an automated quality review, where we use a small MLLM (Qwen3-VL 8B [45]) to identify and discard all remaining poor-quality samples (*e.g.*, poor staining and fragmented tissue).

**Re-annotation & Partitioning.** To create high-quality labels, a strict sequential re-annotation process was employed. First, Gemini-2.5 Pro was used to rewrite the descriptions for all candidate samples. Following this, GPT-5 was introduced as an independent reviewer to assess the quality and factual accuracy of each generated description. We only retained pairs that passed the GPT-5 review. From this 68K subset, we first created the 10K Test Set and the 8K Prototype Bank by sampling, prioritizing smaller clusters to cover rare features, and leaving the 50K Fine-tuning Set for Stage 2. Finally, we invited a pathology expert to conduct a spot check: 93.6% of the image-text pairs were rated as "Usable" (detailed in **Appendix E.3**). We also validated that the Prototype Bank and the Test Set are strictly disjoint, ensuring no data leakage (see **Appendix E.5**).

## 4 Methodology

### 4.1 Overview of UNIPATH

As illustrated in Figure 2, the UNIPATH architecture is designed as a unified framework that seamlessly integrates diagnostic-grade pathology understanding with high-fidelity, controllable image generation. The framework operates through the synergy of three core components.

**Understanding Backbone.** We employ a powerful pathology MLLM (Patho-R1 7B [50]) as our understanding

backbone. Keeping its parameters fully frozen preserves the robust semantic understanding, enabling the model to extract stable, consistent diagnostic semantics and effectively overcome terminological heterogeneity.

**Generation Backbone.** We adopt a 0.6B parameter Diffusion Transformer (DiT) [28] derived from PixArt-$\alpha$ [6]. This model is more efficient than larger alternatives like SDXL (2.3B) [29] and Next-DiT (2B) [14], better suited for our domain-specific dataset. We follow the LDM methodology [33], training the DiT in the VAE's latent space [13]. To achieve higher quality, faster convergence, and more efficient inference, we replace the traditional DDPM objective [16] with a Flow Matching objective [21].

**The Multi-Stream Control.** A key innovation of our UNIPATH is the Multi-Stream Control (MSC). The MSC is a trainable module that acts as the interface between the understanding and the generation backbone. Its primary responsibility is to receive and process multi-source control signals originating from the user and the MLLM. It encodes and fuses these signals into a composite conditional sequence, $C_{comp}$. This sequence is then injected into the DiT to steer the generation with fine-grained control.

**Information Flow.** During training, the VAE encodes an image to $z_0$ and the MSC encodes its caption to $C_{comp}$. The DiT learns the conditional flow $z_1 \rightarrow z_0$). During inference, the MSC encodes a prompt to $C_{comp}$, guiding the DiT to generate $z_0$. This unified design grants UNIPATH controllability while retaining its understanding capabilities.

## 4.2 Multi-Stream Control (MSC)

To handle phrasing heterogeneity and expose component-level morphology control, we propose multi-stream control: it normalizes semantics via a high-level semantics stream, preserves literal intent via a raw-text stream, and injects morphology via a prototype stream.

**High-Level Semantics (HLS) Stream.** The HLS stream is designed to both read out phrasing-invariant diagnostic semantics from the frozen pathology MLLM and expand surface prompts into diagnosis-aware attribute bundles, all via a lightweight implementation. We append $N_q$ learnable queries, $Q^{(0)} \in \mathbb{R}^{N_q \times d_c}$, to the end of the prompt embedding sequence $E_{\text{prompt}} \in \mathbb{R}^{L_r \times d_c}$. The combined input sequence $S^{(0)}$ is constructed as:

$$S^{(0)} = [E_{\text{prompt}}; Q^{(0)}], \tag{1}$$

where $[;]$ denotes concatenation along the sequence dimension. This sequence $S^{(0)}$ is processed by the frozen MLLM backbone. After the final layer $L$, we directly slice the hidden states corresponding to the queries, $H_Q^{(L)} = \text{Tail}_{N_q}(S^{(L)})$, which we term the **Diagnostic Semantic Tokens (DST)**. This mechanism forces the queries to distill the high-level diagnostic semantics embedded within the user's prompt, without requiring any updates to the backbone parameters. Finally, these tokens are passed through an LN and an MLP projection to yield the DST condition $C_{\text{DST}}$:

$$C_{\text{DST}} = \text{MLP}_{\text{DST}}(\text{LN}(H_Q^{(L)})) \in \mathbb{R}^{N_q \times d_c}. \tag{2}$$

**Raw-Text Stream (RTS).** The RTS stream serves as a complement to the HLS stream by preserving the user's literal intent and textual diversity. While the HLS stream excels at extracting abstract "consensus semantics," it may also over-smooth specific stylistic or nuanced details present in the original prompt. Therefore, we reuse the input embeddings $E_{\text{prompt}} \in \mathbb{R}^{L_r \times d_c}$ and compute the RTS conditional tokens $C_{\text{RTS}}$ via a dedicated $\text{MLP}_{\text{RTS}}$.

**Prototype Stream (PS).** The PS stream is designed to achieve component-level morphological control. We employ a non-parametric retrieval mechanism, which remains consistent and frozen during both training and inference. First, we process the user's raw text prompt. We use the powerful pathology vision language model CONCH [23] to encode it into an L2-normalized query vector $q \in \mathbb{R}^{d_q}$, which captures the core semantics of the prompt:

$$q = \text{Norm}(\text{CONCH}_{\text{text}}(\text{prompt})) \in \mathbb{R}^{d_q}. \tag{3}$$

To ensure the retrieved prototypes are relevant to the query in terms of both semantic diagnostic concepts and morphological visual appearances, we adopt a hybrid retrieval strategy: (i) **Global Semantic Retrieval**: We use $q$ to perform Top-k cosine similarity retrieval on both the text index $I_{\text{text}}$ and the vision index $I_{\text{vision}}$:

$$U_g = \text{TopK}_{k_t}(I_{\text{text}} \cdot q) \cup \text{TopK}_{k_v}(I_{\text{vision}} \cdot q). \tag{4}$$

5

(ii) **Local Fine-grained Retrieval**: We parse keywords $\mathcal{K}$ from $q$ and query an inverted index $\mathcal{I}$ to recall fine-grained morphological keywords, yielding the set $U_l = \bigcup_{w \in \mathcal{K}} \mathcal{I}(w)$. We then take the union of prototype IDs recalled from both strategies, clip this sequence to a fixed length $K_m$, and retrieve the corresponding features:

$$P = \text{Proto}[\hat{U}] \in \mathbb{R}^{K_m \times d}, \quad \hat{U} = \text{Clip}\left((U_g \cup U_l), K_m\right), \tag{5}$$

where $\text{Proto} \in \mathbb{R}^{M \times d_P}$ is our offline bank storing the UNI2-h [8] features for $M$ prototypes. We also compute the prototype conditional tokens $C_{\text{PS}}$ via $\text{MLP}_{\text{PS}}$. Finally, the three conditional tokens are fused into a single composite sequence via concatenation:

$$C_{comp} = [C_{\text{DST}}; C_{\text{RTS}}; C_{\text{PS}}]. \tag{6}$$

This $C_{comp}$ is injected into the DiT via cross-attention.

## 4.3 Prototype Bank Construction

The Prototype Bank is a non-parametric, frozen instance bank of 8K real samples (rather than K-means centroids) to preserve true morphological diversity.

**Retrieval Indices.** We built three index components to support our hybrid retrieval strategy:

- **Dense Indices ($I_{\text{text}}$ & $I_{\text{vision}}$):** We built two L2-normalized dense indices, $I_{\text{text}}$ and $I_{\text{vision}}$ ($\in \mathbb{R}^{M \times d_q}$), by encoding the 8K refined texts and images using the CONCH text and vision encoders, respectively.

- **Inverted Index ($\mathcal{I}$):** We created a pathology vocabulary by extracting Top-5000 N-gram candidates from the 50K subset, and then using Gemini-2.5 Pro to review and refine them based on pathology rules. This vocabulary was then used to parse the 8K prototype bank texts to build the inverted index $\mathcal{I}$.

**Prototype Feature Bank.** The Prototype Feature Bank (Proto) provides the content for injection. It was generated by extracting features from the 8K images using the UNI2-h extractor, resulting in the final matrix $\text{Proto} \in \mathbb{R}^{M \times d_P}$.

## 4.4 Training Strategy

**Flow Matching Objective.** We adopt the Flow Matching (FM) [21] objective to train our DiT backbone. The FM framework learns a vector field to transport samples from a prior distribution (*e.g.*, Gaussian) to a target continuous distribution. Given a ground-truth latent $z_0$ (from VAE encoder) and our composite condition $C_{\text{comp}}$, the training proceeds as follows. At each step, we sample a timestep $t \sim \mathcal{U}(0, 1)$, and noise $z_1 \sim \mathcal{N}(0, 1)$. Following [22], we compute the interpolated latent $z_t$ via linear interpolation:

$$z_t = t z_0 + (1 - t) z_1. \tag{7}$$

The corresponding target velocity vector $v_t$ (*i.e.*, the derivative of $z_t$ with respect to $t$) is analytically given by:

$$v_t = \frac{dz_t}{dt} = z_0 - z_1. \tag{8}$$

Our DiT model, parameterized by $\theta$ and denoted as $v_\theta$, is trained to predict this velocity, conditioned on the corrupted latent $z_t$, the timestep $t$, and our composite condition $C_{\text{comp}}$. The training objective $\mathcal{L}_{\text{Flow}}$ is defined as the L2 loss:

$$\mathcal{L}_{\text{Flow}}(\theta) = \mathbb{E}_{z_0, C_{\text{comp}}, t, z_1} \left[ \left\| v_\theta(z_t, t, C_{\text{comp}}) - v_t \right\|^2 \right]. \tag{9}$$

**Two-Stage Training Strategy.** We employ a two-stage training strategy. Stage 1 (Semantic Alignment) uses the 2.65M large-scale corpus to pre-train the DiT ($v_\theta$) and MSC modules to learn fundamental visual-textual alignment. Stage 2 (High-Quality Fine-tuning) then continues to train the DiT and MSC on the 50K high-quality fine-tuning set using a smaller learning rate, significantly enhancing visual fidelity and fine-grained controllability.

**Table 1** Quantitative comparison of **Visual Fidelity** and **Text-Image Alignment** on our 10K High-Quality Test Set. Patho-FID/KID are computed using the UNI2-h backbone. ★ indicates models fully fine-tuned on our large dataset. Retrieval metrics (Recall/mAP) are reported as Report2Gen (T2I) / Real2Gen (I2I). The best-performing model is **in-bold**, and the second-best-performing model is underlined.

| | Unif. Model | Visual Fidelity ↓ | | | | | Text-Image Alignment ↑ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | FID | KID | Patho-FID | Patho-KID | LPIPS | CLIP-Score | Recall@10 | Recall@50 | mAP@10 | mAP@50 |
| Real Data | - | - | - | - | - | - | 0.372 | 9.60/- | 27.55/- | 3.45/- | 4.21/- |
| *General Text to Image Generation Models* | | | | | | | | | | | |
| SD1.5★ [33] | ✗ | 160.36 | 0.221 | 259.69 | 0.321 | 0.634 | 0.147 | 0.27/0.38 | 1.48/1.33 | 0.23/0.49 | 0.40/0.61 |
| SDXL★ [29] | ✗ | 291.17 | 0.305 | 295.50 | 0.343 | 0.657 | -0.035 | 0.34/0.19 | 1.21/0.95 | 0.52/0.27 | 0.67/0.42 |
| Pixart-$\alpha$★ [6] | ✗ | 270.36 | 0.277 | 336.78 | 0.401 | 0.603 | 0.105 | 0.84/0.64 | 3.40/2.75 | 1.17/0.91 | 1.38/1.15 |
| BLIP3o★ [5] | ✔ | 95.04 | 0.183 | 358.61 | 0.431 | 0.598 | 0.080 | 0.93/1.10 | 3.78/3.68 | 1.23/1.61 | 1.51/1.72 |
| Show-o2★ [41] | ✔ | 98.82 | 0.200 | 184.19 | 0.239 | <u>0.601</u> | **0.357** | **4.64**/<u>2.40</u> | **15.26**/<u>8.15</u> | **5.85**/<u>2.74</u> | **5.50**/2.58 |
| *Pathological / Medical Text to Image Generation Models* | | | | | | | | | | | |
| Pixcell [47] | ✗ | <u>80.74</u> | <u>0.119</u> | <u>163.44</u> | <u>0.177</u> | 0.602 | - | - | - | - | - |
| PathLDM [46] | ✗ | 93.91 | 0.170 | 174.32 | 0.254 | 0.606 | 0.182 | 0.12/0.19 | 0.70/0.66 | 0.19/0.27 | 0.30/0.36 |
| UniMedVL [26] | ✔ | 156.17 | 0.219 | 216.27 | 0.255 | 0.619 | 0.319 | 2.73/1.98 | 8.47/6.67 | 3.86/2.62 | 3.85/<u>2.60</u> |
| **UNIPATH (Ours)** | ✔ | **25.70** | **0.081** | **80.86** | **0.134** | **0.570** | <u>0.348</u> | <u>3.92</u>/**3.30** | <u>14.08</u>/**11.92** | <u>5.66</u>/**4.64** | **5.54**/**4.67** |

# 5 Experiments and Results

## 5.1 Implementation Details

**Model Architecture.** Our DiT backbone (0.6B) comprises 28 Transformer layers, 16 heads, and a hidden dimension of 1152. We use the Stable Diffusion 3 VAE [13] with 8× downsampling. The understanding backbone (Patho-R1 7B) [50] and the VAE are frozen during training. We use 64 learnable queries in the HLS stream, yielding 64 DST.

**Training Details.** In Stage 1, we utilize 2.58M text-image pairs (excluding the 68K subset) and train for 10,000 steps with a global batch size of 512. The learning rate is linearly warmed up to a peak of $1e^{-4}$, followed by cosine annealing. Stage 2 loads the pretrained weights and fine-tunes exclusively on the 50K data for 500 steps with a fixed learning rate of $2e^{-5}$. For the PS Stream, we set $K_m = 16$. All experiments are conducted on 16 NVIDIA H100 GPUs.

## 5.2 Evaluation Protocols

**Four-Tier Evaluation Hierarch.** Traditional metrics are insufficient for evaluating pathology generative models, as they fail to capture diagnostic relevance or controllability. We thus first performed a preliminary validation to confirm the model's understanding, then established a four-tier evaluation scheme. **Tier 1: Visual Fidelity.** We assess perceptual quality and distributional fidelity using FID [15] and KID [4] (standard and pathology variants), as well as LPIPS [49]. **Tier 2: Text-Image Alignment.** We assess semantic consistency using clip-score, retrieval metrics, and MLLM, as well as human judges. **Tier 3: Fine-grained Semantic Control.** We assess fine-grained semantic control via the "Train-on-Synth, Test-on-Real" paradigm, comparing the real-test-set performance of a Gen2Real classifier against a Real2Real classifier. **Tier 4: Downstream Task Utility.** We evaluate data augmentation utility by comparing the real-test-set performance gain of a baseline classifier (real) versus an augmented classifier (real + synthetic).

**Baselines.** To comprehensively evaluate UNIPATH's generative capabilities, we selected SOTA models from both general and specific domains as baselines. **General T2I** SOTA: We selected SD1.5 [33], SDXL [29], PixArt-$\alpha$ [6], BLIP3o [5], and Show-o2 [41]. As these general-purpose models have minimal exposure to pathology data, we fully fine-tuned them on our dataset using the same training strategy. **Pathology-Specific** SOTA: We selected PixCell [47], PathLDM [46], and UniMedVL [26], which are leading generative models designed for this domain. Notably, PixCell is natively image-conditioned. To adapt PixCell for T2I comparison , we used CONCH to retrieve a prompt-relevant image from our prototype bank to serve as its generation condition. We note that this adaptation, while necessary, is a convenience variant and may not represent the optimal text-conditioning performance of the PixCell architecture.

**Understanding Capability.** For validation, we evaluate UNIPATH on the PathMMU benchmark [35]. The results (see **Appendix F.1**) show that UNIPATH achieves SOTA performance among open-source models and approaches

**(a)** GPT-5 as judge.
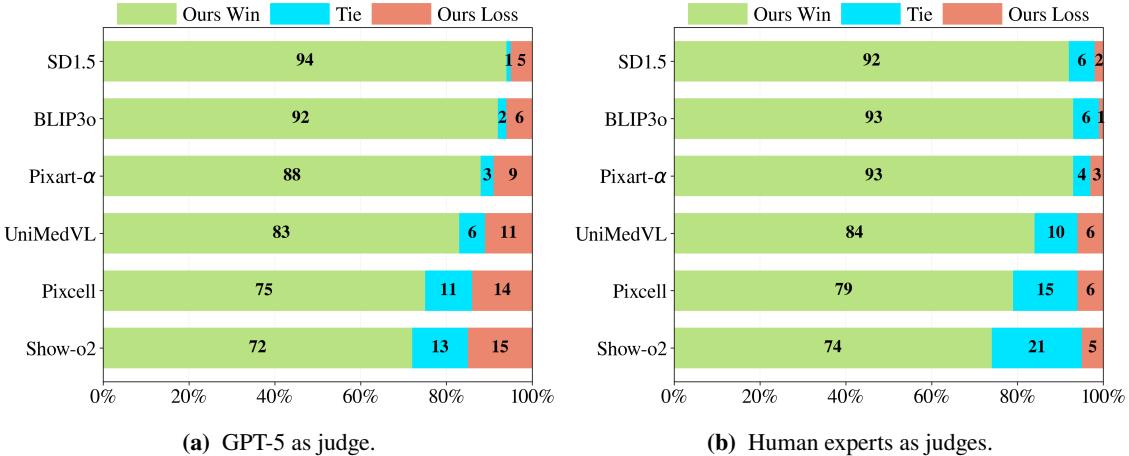
**(b)** Human experts as judges.

**Figure 3** MLLM and human expert evaluation results.

closed-source models. This performance confirms the model's ability to extract robust, phrasing-invariant semantics, which are critical to controllable generation.

## 5.3 Visual Fidelity

We assess visual fidelity in Table 1 (left), covering standard (FID/KID) and pathology (Patho-FID/KID with UNI2-h [8] backbone) metrics, as well as LPIPS. UNIPATH attains state-of-the-art results across all metrics. The gap is substantial: on the challenging Patho-FID, UNIPATH achieves a 50.5% relative reduction compared with the second-best model. It also yields the best scores on FID, KID, Patho-KID, and LPIPS, indicating that the images generated by UNIPATH align best with the real data distribution in both general and pathology feature spaces. We also show results using Virchow2 [51] and MUSK [40] extractors in **Appendix F.2**, where performance remains strong.

| Models | Cytology Type (4-classes) | | Hemorrhage (2-classes) | |
|---|---|---|---|---|
| | Wgt. F1 | Wgt. AUC | Wgt. F1 | Wgt. AUC |
| Real Data-Image | 83.43 | 87.15 | 78.13 | 80.11 |
| Real Data-Text | 95.22 | 98.59 | 83.09 | 88.29 |
| SD1.5★ [33] | 71.26 | 71.34 | 64.41 | 62.88 |
| SDXL★ [29] | 69.41 | 42.31 | 62.25 | 56.78 |
| Pixart-$\alpha$★ [6] | 70.65 | 72.27 | 67.67 | 68.53 |
| BLIP3o★ [5] | 67.31 | 68.78 | 68.57 | 65.38 |
| Show-o2★ [41] | 77.82 | 66.89 | 72.44 | 71.72 |
| Pixcell [47] | 76.37 | 83.22 | 71.57 | 71.59 |
| PathLDM [46] | 69.57 | 63.61 | 62.07 | 47.58 |
| UniMedVL [26] | 78.21 | 72.43 | 68.05 | 65.36 |
| **UNIPATH (Ours)** | 81.43 | 84.05 | 74.96 | 75.27 |
| **UNIPATH-*Aug*.** | **81.49** | **85.29** | **77.02** | **79.03** |

**Table 2** Results on **Fine-grained Control**.

## 5.4 Text-Image Alignment

We next evaluate the models' performance on text-image alignment. Table 1 (right) reports the quantitative results based on a CONCH-based clip-score and retrieval metrics.

**CLIP-Score.** On the CLIP-Score metric, our UNIPATH achieves a high score of 0.348, beating all non-unified models. It trails only the unified model Show-o2 (0.357), confirming that a strong understanding aids semantic generation. The 0.009 CLIP-Score gap to Show-o2 is explained by its SigLIP-distilled semantic tokens (see **Appendix A**); other alignment and downstream metrics favour UNIPATH.

**Retrieval Metrics.** On Report2Gen (T2I retrieval), Show-o2 achieves the best results, while UNIPATH is highly competitive with near-tied mAP. On Real2Gen (I2I retrieval)—which measures feature-space distance to real images—UNIPATH reaches clear SOTA, outperforming all methods by a large margin. Thus, although Show-o2 has a slight T2I edge, UNIPATH generates images that are closer to real pathology in feature space. This trend matches the LPIPS results and shows that UNIPATH is not only semantically aligned but also produces morphology and visual structures most faithful to
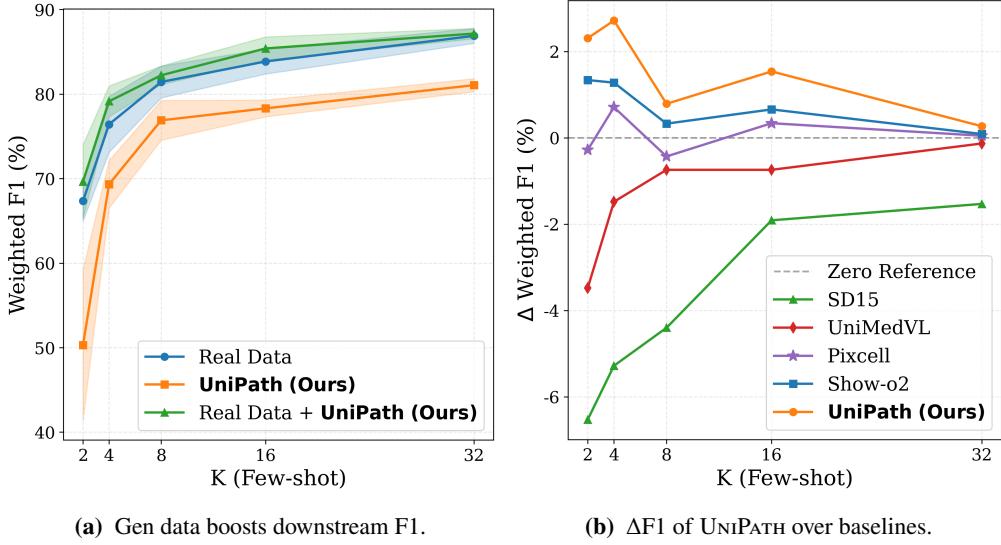
**(a)** Gen data boosts downstream F1.

**(b)** ΔF1 of UNIPATH over baselines.

**Figure 4** Few-shot classification on Kather-CRC-2016.

the true pathological appearance.

**As-Judge Evaluation.** This strong alignment is confirmed by two pairwise evaluations (Figure 3). Against the strongest baseline, UNIPATH was preferred by GPT-5 in 72% of cases. A panel of three human pathologists confirmed this, preferring UNIPATH in 74% of cases. This convergence of both MLLM and human expert preference underscores a clear advantage in nuanced, human-aligned semantic understanding. See details in **Appendix F.4**.

## 5.5 Fine-grained Semantic Control

We assess semantic fidelity for fine-grained concepts with a "Train-on-Synth, Test-on-Real" (G2R) protocol, since FID does not ensure learnable morphology; specifically, we split the 10K test set 60/20/20 (train/val/test), train a linear probe classifier on frozen CONCH features, and evaluate on the real test split, reporting F1 and AUC.

**Quantitative Analysis.** Table 2 summarizes the results. Real Data–Image serves as the Real2Real (R2R) ceiling, while the strong Real Data–Text baseline indicates that key diagnostic factors are encoded in the prompts. Among all G2R models, UNIPATH achieves performance most comparable to the R2R benchmark. For cytology type, its F1 score is within two points of the R2R result, and for hemorrhage, it surpasses the next-best model by a clear margin.

**Closing the Gap with Augmentation.** We also evaluated augmenting with 5 images per prompt (UNIPATH-Aug), which further boosts performance, reaching 98.7% (hemorrhage) and 97.9% (cytology type) of the real-image AUC. This nearly closes the R2R gap, indicating that the synthetic features are morphologically precise enough to be learned and generalized to real images.

## 5.6 Downstream Task Utility

Finally, we assess UNIPATH's practical utility as a data augmentation tool using a few-shot setting.

**Experimental Setup.** We adopted a 7-class Kather-CRC-2016 [17] classification task (background filtered) under a strict few-shot setup. To synthesize relevant data, Gemini-2.5 Pro generated rich, descriptive prompts from representative images for each class. UNIPATH and all baselines used these prompts for augmentation. We report the mean F1-Weighted score over 5 random seeds on the real test set.

**Quantitative Analysis.** The results in Figure 4 illustrate the superior utility of data generated by UNIPATH. Figure 4 (Left) shows that 'Real Data + UNIPATH' significantly outperforms the "Real Data" baseline across all K-shot settings, which confirms its augmentation effectiveness. Figure 4 (Right) shows that among SOTA methods, only UNIPATH and
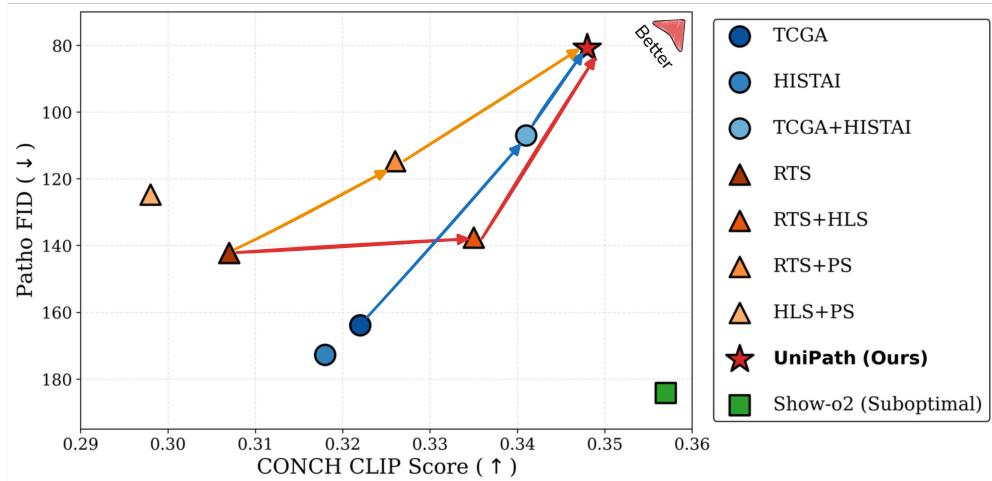
**Figure 5** Dataset and component ablation results.

Show-o2 deliver positive F1 gains at every K, and UNIPATH shows a clear lead, especially in the extreme few-shot, with larger gains than Show-o2 and PixCell. Other models yield negative gains, which indicates that their synthetic data degrades performance.

## 5.7 Ablation Studies

We ablate our (i) data curation strategy and (ii) Multi-Stream Control (MSC) architecture. Figure 5 plots all results on Visual Fidelity vs. Text-Image Alignment.

**Ablation on Data Contribution.** The blue path in Figure 5 shows our data ablation. Starting from the TCGA (1.62M) baseline, adding the 1.03M HISTAI dataset significantly improves both metrics. The final 50K fine-tuning dataset pushes performance further to achieve the best overall results, proving the value of our two-stage data strategy.

**Ablation on MSC Architecture.** The orange and red paths in Figure 5 ablate the MSC architecture. Starting from the poor-performing RTS-only baseline, adding the HLS stream primarily boosts alignment, while adding the PS stream improves both metrics. The poor alignment of the HLS+PS variant proves that all streams are indispensable. Our UNIPATH leverages all three streams to achieve the optimal balance, firmly occupying the "Better" region. We also provide inference-time ablations in **Appendix F.5**, analyzing the sensitivity of key MSC hyperparameters.

## 6 Qualitative Analysis and Controllability

We provide key qualitative visual evidence in Figure 6 to intuitively demonstrate UNIPATH's SOTA performance and the effectiveness of our MSC architecture in addressing core challenges, such as terminological heterogeneity. More diverse examples are provided in **Appendix G**.

**SOTA Qualitative Comparison.** Figure 6 (top) compares UNIPATH against SOTA baselines on complex prompts. Baseline models often exhibit concept dropping, introduce artifacts, or render erroneous morphologies when faced with prompts containing multiple pathological features. In contrast, UNIPATH is the only model capable of accurately and high-fidelity reproducing all specified pathological features. This qualitative result visually confirms our SOTA Visual Fidelity and highlights the CLIP-Score's limitations. While Show-o2 scored slightly higher, this figure proves UNIPATH has superior practical semantic alignment by rendering fine-grained concepts that the metric missed.

**MSC: Robustness to Paraphrasing.** Figure 6 (bottom) visually validates our MSC's ability to resolve terminological heterogeneity. As shown in the samples, UNIPATH generates morphologically consistent images for "Original" and "Variant" captions despite different phrasing. This demonstrates the efficacy of UNIPATH's multi-stream design: the HLS
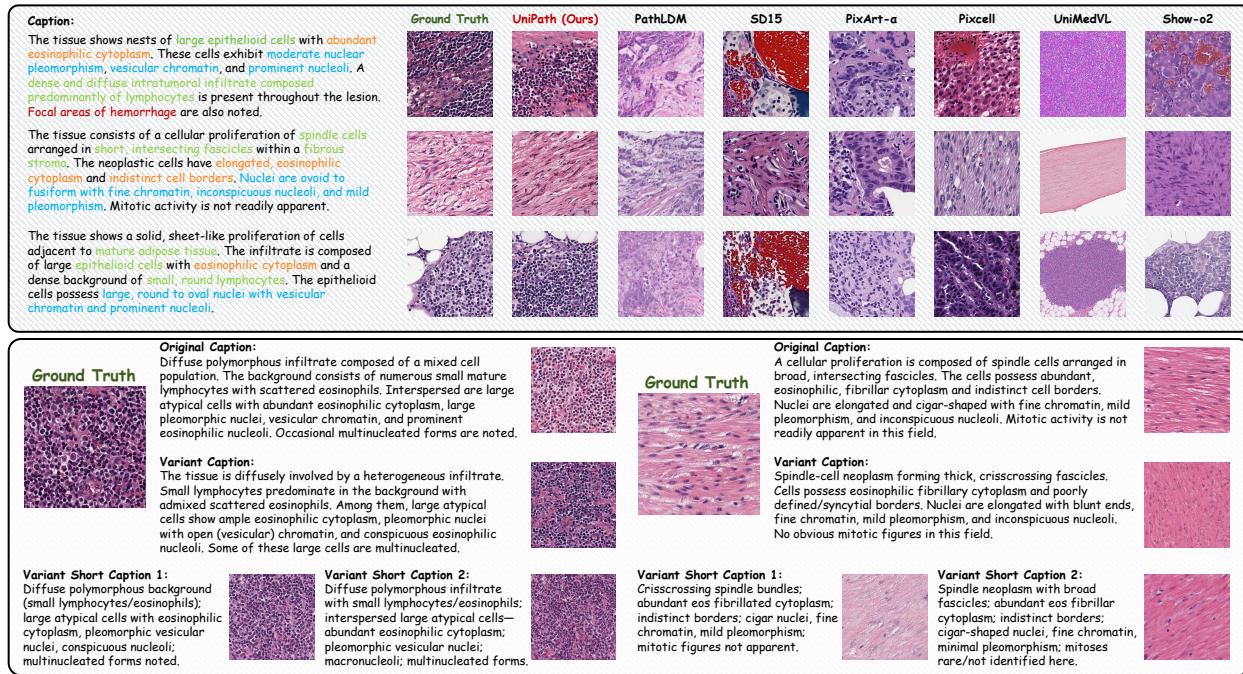
**Figure 6  Qualitative examples of UNiPATH's generation and controllability. (Top)** Comparison of UNiPATH with SOTA baselines on complex prompts. Prompt text is color-coded by pathological concept: (Tissue/Cell Type, Nuclear Features, Cytoplasm, Hemorrhage). **(Bottom)** Robustness of UNiPATH to synonymous prompts, validating the MSC's handling of terminological heterogeneity.

stream distills heterogeneous prompts into the same diagnostic semantic tokens, maintaining stable semantic guidance and achieving robust control.

# 7  Conclusion

In this paper, we introduce UNiPATH, a novel unified pathology model coupling MLLM understanding with a DiT generator. Our core contribution is the Multi-Stream Control architecture, which uses a High-Level Semantic stream and a Prototype stream to tackle terminological heterogeneity and component-level morphological control simultaneously. To support and validate UNiPATH, we curated a large-scale corpus, a 68K high-fidelity subset, and established a four-tier evaluation protocol. Extensive experiments show UNiPATH demonstrates leading performance in both automated metrics and human expert evaluations.

**Broader Impacts and Future Work.**  As a controllable, unified pathology foundational model, UNiPATH offers significant potential in data augmentation by generating high-fidelity, customized synthetic images; in research by enabling the systematic exploration of morphological features; and in education as an interactive training tool. For future work, we plan to extend UNiPATH to support higher-resolution image generation and explore its capabilities in pathological image editing. We provide further discussions of the limitations and ethical issues in the **Appendix A & B**.

# References

[1] Saghir Alfasly, Wataru Uegami, MD Hoq, Ghazal Alabtah, and HR Tizhoosh. Semantic and visual crop-guided diffusion models for heterogeneous tissue synthesis in histopathology. *arXiv preprint arXiv:2509.17847*, 2025.

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[3] Mahesh Bhosale, Abdul Wasi, Yuanhao Zhai, Yunjie Tian, Samuel Border, Nan Xi, Pinaki Sarder, Junsong Yuan, David Doermann, and Xuan Gong. Pathdiff: Histopathology image synthesis with unpaired text and mask conditions. *arXiv preprint arXiv:2506.23440*, 2025.

[4] Miko laj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018.

[5] Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025.

[6] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. PixArt-α: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.

[7] Pingyi Chen, Chenglu Zhu, Sunyi Zheng, Honglin Li, and Lin Yang. Wsi-vqa: Interpreting whole slide images by generative visual question answering. In *European Conference on Computer Vision*, pages 401–417. Springer, 2024.

[8] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature medicine*, 30(3):850–862, 2024.

[9] Ying Chen, Guoan Wang, Yuanfeng Ji, Yanjun Li, Jin Ye, Tianbin Li, Ming Hu, Rongshan Yu, Yu Qiao, and Junjun He. Slidechat: A large vision-language assistant for whole-slide pathology image understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5134–5143, 2025.

[10] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.

[11] Sachin Kumar Danisetty, Alexandros Graikos, Srikar Yellapragada, and Dimitris Samaras. Pathsegdiff: Pathology segmentation using diffusion model representations. In *MICCAI Workshop on Deep Generative Models*, pages 141–150. Springer, 2025.

[12] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.

[13] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.

[14] Peng Gao, Le Zhuo, Dongyang Liu, Ruoyi Du, Xu Luo, Longtian Qiu, Yuhang Zhang, Chen Lin, Rongjie Huang, Shijie Geng, et al. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers. *arXiv preprint arXiv:2405.05945*, 2024.

[15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

[16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[17] Jakob Nikolas Kather, Cleo-Aron Weis, Francesco Bianconi, Susanne M Melchers, Lothar R Schad, Timo Gaiser, Alexander Marx, and Frank Gerrit Zöllner. Multi-class texture analysis in colorectal cancer histology. *Scientific reports*, 6(1):1–11, 2016.

[18] Jingxiong Li, Chenglu Zhu, Sunyi Zheng, Pingyi Chen, Yuxuan Sun, Honglin Li, and Lin Yang. Topofm: Topology-guided pathology foundation model for high-resolution pathology image synthesis with cellular-level control. *IEEE Transactions on Medical Imaging*, 2025.

[19] Yuci Liang, Xinheng Lyu, Wenting Chen, Meidan Ding, Jipeng Zhang, Xiangjian He, Song Wu, Xiaohan Xing, Sen Yang, Xiyue Wang, et al. Wsi-llava: A multimodal large language model for whole slide image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22718–22727, 2025.

[20] Bin Lin, Zongjian Li, Xinhua Cheng, Yuwei Niu, Yang Ye, Xianyi He, Shenghai Yuan, Wangbo Yu, Shaodong Wang, Yunyang Ge, et al. Uniworld: High-resolution semantic encoders for unified visual understanding and generation. *arXiv preprint arXiv:2506.03147*, 2025.

[21] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

[22] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.

[23] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guillaume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. A visual-language foundation model for computational pathology. *Nature medicine*, 30(3): 863–874, 2024.

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

[25] Dmitry Nechaev, Alexey Pchelnikov, and Ekaterina Ivanova. Histai: An open-source, large-scale whole slide image dataset for computational pathology, 2025. URL https://arxiv.org/abs/2505.12120.

[26] Junzhi Ning, Wei Li, Cheng Tang, Jiashi Lin, Chenglong Ma, Chaoyang Zhang, Jiyao Liu, Ying Chen, Shujian Gao, Lihao Liu, et al. Unimedvl: Unifying medical multimodal understanding and generation through observation-knowledge-analysis. *arXiv preprint arXiv:2510.15710*, 2025.

[27] OpenAI. Introducing GPT-5, August 2025. URL https://openai.com/index/introducing-gpt-5/. Accessed: 2025-11-03.

[28] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.

[29] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.

[31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21 (140):1–67, 2020.

[32] Ekaterina Redekop, Mara Pleasure, Vedrana Ivezic, Zichen Wang, Kimberly Flores, Anthony Sisk, William Speier, and Corey Arnold. Prototype-guided diffusion for digital pathology: Achieving foundation model performance with minimal clinical data. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5187–5195, 2025.

[33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.

[34] Mehmet Saygin Seyfioglu, Wisdom O Ikezogwo, Fatemeh Ghezloo, Ranjay Krishna, and Linda Shapiro. 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*, pages 13183–13192, 2024.

[35] Yuxuan Sun, Hao Wu, Chenglu Zhu, Sunyi Zheng, Qizi Chen, Kai Zhang, Yunlong Zhang, Dan Wan, Xiaoxiao Lan, Mengyue Zheng, et al. Pathmmu: A massive multimodal expert-level benchmark for understanding and reasoning in pathology. In *European Conference on Computer Vision*, pages 56–73. Springer, 2024.

[36] Yuxuan Sun, Yunlong Zhang, Yixuan Si, Chenglu Zhu, Zhongyi Shui, Kai Zhang, Jingxiong Li, Xingheng Lyu, Tao Lin, and Lin Yang. Pathgen-1.6 m: 1.6 million pathology image-text pairs generation through multi-agent collaboration. *arXiv preprint arXiv:2407.00203*, 2024.

[37] Yuxuan Sun, Yixuan Si, Chenglu Zhu, Xuan Gong, Kai Zhang, Pingyi Chen, Ye Zhang, Zhongyi Shui, Tao Lin, and Lin Yang. Cpath-omni: A unified multimodal foundation model for patch and whole slide image analysis in computational pathology. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 10360–10371, 2025.

[38] Eugene Vorontsov, Alican Bozkurt, Adam Casson, George Shaikovski, Michal Zelechowski, Kristen Severson, Eric Zimmermann, James Hall, Neil Tenenholtz, Nicolo Fusi, et al. A foundation model for clinical-grade computational pathology and rare cancers detection. *Nature medicine*, 30(10):2924–2935, 2024.

[39] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025.

[40] Jinxi Xiang, Xiyue Wang, Xiaoming Zhang, Yinghua Xi, Feyisope Eweje, Yijiang Chen, Yuchen Li, Colin Bergstrom, Matthew Gopaulchan, Ted Kim, et al. A vision–language foundation model for precision oncology. *Nature*, 638(8051):769–778, 2025.

[41] Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou. Show-o2: Improved native unified multimodal models. *arXiv preprint arXiv:2506.15564*, 2025.

[42] Hanwen Xu, Naoto Usuyama, Jaspreet Bagga, Sheng Zhang, Rajesh Rao, Tristan Naumann, Cliff Wong, Zelalem Gero, Javier González, Yu Gu, et al. A whole-slide foundation model for digital pathology from real-world data. *Nature*, 630(8015):181–188, 2024.

[43] Meilong Xu, Saumya Gupta, Xiaoling Hu, Chen Li, Shahira Abousamra, Dimitris Samaras, Prateek Prasanna, and Chao Chen. Topocellgen: Generating histopathology cell topology with a diffusion model. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20979–20989, 2025.

[44] Zhe Xu, Ziyi Liu, Junlin Hou, Jiabo Ma, Cheng Jin, Yihui Wang, Zhixuan Chen, Zhengyu Zhang, Fuxiang Huang, Zhengrui Guo, et al. A versatile pathology co-pilot via reasoning enhanced multimodal large language model. *arXiv preprint arXiv:2507.17303*, 2025.

[45] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

[46] Srikar Yellapragada, Alexandros Graikos, Prateek Prasanna, Tahsin Kurc, Joel Saltz, and Dimitris Samaras. Pathldm: Text conditioned latent diffusion model for histopathology. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5182–5191, 2024.

[47] Srikar Yellapragada, Alexandros Graikos, Zilinghan Li, Kostas Triaridis, Varun Belagali, Saarthak Kapse, Tarak Nath Nandi, Ravi K Madduri, Prateek Prasanna, Tahsin Kurc, et al. Pixcell: A generative foundation model for digital histopathology images. *arXiv preprint arXiv:2506.05127*, 2025.

[48] Andrew Zhang, Guillaume Jaume, Anurag Vaidya, Tong Ding, and Faisal Mahmood. Accelerating data processing and benchmarking of ai models for pathology. *arXiv preprint arXiv:2502.06750*, 2025.

[49] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

[50] Wenchuan Zhang, Penghao Zhang, Jingru Guo, Tao Cheng, Jie Chen, Shuwan Zhang, Zhang Zhang, Yuhao Yi, and Hong Bu. Patho-r1: A multimodal reinforcement learning-based pathology expert reasoner. *arXiv preprint arXiv:2505.11404*, 2025.

[51] Eric Zimmermann, Eugene Vorontsov, Julian Viret, Adam Casson, Michal Zelechowski, George Shaikovski, Neil Tenenholtz, James Hall, David Klimstra, Razik Yousfi, et al. Virchow2: Scaling self-supervised mixed magnification models in pathology. *arXiv preprint arXiv:2408.00738*, 2024.

# Appendix

## Appendix Contents

## A   Limitations

**Text–Image Alignment Performance.** As shown in Table 1, UNIPATH is SOTA on every downstream metric, on Real2Gen retrieval, and in the MLLM-as-Judge comparison, yet trails Show-o2 on the single CONCH CLIP-Score (0.348 vs. 0.357). We argue that this gap reflects a bias in the evaluator–model paradigm, rather than a true semantic alignment weakness.

- **Show-o2's semantic tokens are contrastive-distilled**. During training, Show-o2 loads SigLIP weights into its Semantic Layers and minimizes a distillation loss, forcing those tokens to mimic SigLIP patch features. Although generation itself uses flow matching, the resulting high-level tokens remain "SigLIP-style" at inference.

- **The evaluator is also contrastive.** CONCH [23] is a CLIP-family model whose similarity metric is the same cosine space SigLIP is trained in. A model whose internal tokens are pre-aligned to this space naturally receives a higher CLIP score.

- **UNIPATH uses MLLM-derived pathology semantics without SigLIP distillation.** Our HLS stream extracts diagnosis tokens from a frozen Patho-R1 MLLM. Patho-R1's vision encoder is CLIP-based, but no part of the generator is forced to match CLIP/SigLIP features; the DiT learns purely via flow-matching reconstruction. Hence, its latent space is optimised for morphological fidelity, not for cosine similarity with CLIP evaluators.

- **Cross-metric consistency.** UNIPATH leads on Real2Gen retrieval (closer in feature space to real WSIs), on MLLM-as-Judge human-preference scoring, and delivers the largest Tier 4 F1 gains. These orthogonal results confirm that the small CLIP-Score gap is an artefact of evaluator homology, not of inferior text–image alignment.

In summary, the lower CONCH CLIP-Score stems from Show-o2's SigLIP-distilled tokens matching the evaluator's contrastive space, whereas UNIPATH prioritises pathology-specific morphology and semantics, which better serve real diagnostic tasks.

**Dependency on Prototype Bank.**   One of UNIPATH's strengths comes from the component-level control provided by the Prototype Stream (PS). This advantage, however, is highly dependent on the quality and coverage of our 8K instance prototype bank. If an extremely rare morphological component is not well-represented in our 8K bank, the PS stream cannot provide precise control for that concept.

## B   Ethical Statement

This research strictly adheres to the relevant ethical guidelines for medical AI research.

**Data Usage and Patient Privacy.**   All data used in this study (TCGA and HISTAI) are publicly available datasets intended for research. All data were fully anonymized and de-identified by the original providers prior to release and contain no Protected Health Information (PHI). Our usage strictly complies with the Data Use Agreements (DUAs) for both TCGA and HISTAI.

**Potential for Misuse and Mitigation.**   We acknowledge that high-fidelity pathological image generation (*i.e.*, "medical deepfakes") carries a potential risk of misuse, such as attempting to interfere with clinical diagnostic workflows in extreme cases. We emphasize that UNIPATH is currently intended for research purposes only, with the design goals of (i) advancing controllable generation in pathology, (ii) providing controllable data augmentation for computational pathology, and (iii) serving as an educational tool. This model **must not** be used for any direct clinical diagnosis.

**Algorithmic Bias.**   Our model's performance relies on the quality and distribution of our training data. Despite our efforts to add diversity (1.03M HISTAI) and balance our 68K subset (via K-means and elite sampling), our training data may still contain undiscovered biases (*e.g.*, in demographic representation across race, age, or sex). The model may learn and amplify these biases. Future work is required to specifically quantify and mitigate such biases.

## C   Future Work

While UNIPATH marks significant progress in unifying pathology understanding and controllable synthesis, several key directions remain for future exploration.

**Support for Higher Resolution and Broader Histological Context.**   The current UNIPATH model primarily operates on $384 \times 384$ pixel patches. While sufficient for capturing cell-level morphological features, this limits the model's ability to understand and generate larger-scale architectural patterns, such as complex glandular structures or tumor-stroma interactions. Future work should explore extending UNIPATH to higher resolutions or integrating a larger field of view, enabling the generation of images that are more histologically context-aware.

**Controllable Pathological Image Editing.**   UNIPATH currently focuses on image generation from text prompts. A high-impact extension is to enable fine-grained editing of existing real pathology images. This can be framed as a "counterfactual synthesis" task — such as "adding moderate nuclear atypia" to a benign tissue image or "removing the specified inflammatory infiltrate." The MSC architecture of UNIPATH provides an ideal framework for this: the HLS stream could parse the editing instruction (*e.g.*, "increase mitotic figures"), while the PS stream could retrieve and inject the corresponding morphological prototypes to achieve the precise, localized modification.

**Scaling the Prototype Bank.** Our prototype-based control mechanism opens a promising avenue for future enhancement. While the curated 8K bank establishes the efficacy of this approach, we can further enhance the model's generative "vocabulary" by scaling this bank. Future work could explore using active learning or self-supervised methods to automatically mine and cluster novel, informative prototypes from large-scale, unlabeled datasets. This expansion would enable UNIPATH to synthesize an even greater diversity of morphological features with high precision, particularly for rare, long-tail pathological phenomena.

## D  Detailed Implementation Details

### D.1  Model Architecture Specifications

**Generation Backbone (DiT) and Conditioning.** Our generation backbone is a 0.6B-parameter DiT (Diffusion Transformer) designed in the PixArt [6] style. It comprises 28 Transformer layers, 16 attention heads, and a hidden dimension $d = 1152$. Our model employs a hybrid conditioning mechanism. The fused conditional vector $C_{comp}$ is injected into every DiT layer via traditional cross-attention. In contrast, the timestep is handled separately, injected via the AdaLayerNorm-Single (AdaLN-S) mechanism to perform conditional normalization.

**Understanding Backbone and VAE.** Our backbone is based on the Patho-R1 (7B) [50] model, which is post-trained on pathology domain data from Qwen2.5-VL 7B [2]. The backbone remains fully frozen throughout all training stages. The VAE employed is the Stable Diffusion 3 [13] VAE, featuring an 8x downsampling factor.

**MSC Module Implementation.** In the Multi-Stream Control (MSC) module, the projection layers for the three streams (HLS, RTS, and PS), such as $\text{MLP}_{\text{DST}}$, are implemented as separate 2-layer Feed-Forward Networks (FFNs) with unshared weights. Each of these MLPs follows the same architecture: a linear projection from the input dimension to the hidden size, followed by a GELU activation, and a second linear projection from the hidden size back to the hidden size. The hidden size is 1152 (matching the DiT's hidden dimension). For the Prototype Stream (PS) retrieval, we provide the specific hyperparameters used for the hybrid strategy. The total number of prototypes is $K_m = 16$. For the Global Semantic Retrieval ($U_g$, Eq. 4), we set the retrieval tops $k_t = 4$ (Text) and $k_v = 4$ (Vision). For the Local Fine-grained Retrieval ($U_l$), we parse the four rarest keywords from the prompt and randomly sample 2 prototypes for each term, resulting in 8 local prototypes. The final set $\hat{U}$ is the union of $U_g$ and $U_l$, clipped to 16 (Eq. 5).

**Flow Matching Training and Inference.** During the training stage, we adopt a Rectified Flow strategy. Specifically, we first sample $u \sim \mathcal{U}(0, 1)$, map it to timesteps and sigmas, and construct the noise-interpolated $z_t$ accordingly. The model $v_\theta$ is trained to predict the target velocity $v_t = z_0 - z_1$. During the inference stage, we use the Euler solver with 30 function evaluations. We employ Classifier-Free Guidance with a guidance_scale of 3.0.

### D.2  Training Hyperparameters

**General Setup.** Across both training stages, we used the AdamW optimizer with default betas ($\beta_1 = 0.9, \beta_2 = 0.999$), an epsilon of $1e^{-8}$, and a weight decay of 0.01. All training was conducted using mixed precision. All experiments were conducted on 16 NVIDIA H100 GPUs. All input images were processed to a resolution of $384 \times 384$.

**Stage 1: Semantic Alignment (Pre-training).** The model was pre-trained on 2.58M text-image pairs (excluding the 68K subset) for 10,000 steps using the Flow Matching (MSE) loss. We used a global batch size of 512. The learning rate was linearly warmed up for the first 2% of steps (200 steps) to a peak of 1e-4. It was then decayed using a cosine scheduler with a minimum learning rate of $1e^{-5}$.

**Stage 2: High-Quality Fine-tuning.** The model was subsequently fine-tuned on the 50K high-quality subset for 500 steps. We used a global batch size of 512 and a fixed learning rate of 2e-5 for this entire stage.

## E  Dataset Detailed Analysis

### E.1  Statistics of the 1.03M Corpus

**Word Count.** We analyzed the caption-length distributions of the 1.03M corpus before and after summarization using Qwen3-8B, as shown on the left side of Figure S1. The distributions exhibit unimodal and symmetric curves. In the

original captions, the most frequent length is around 120 words, accounting for 9.6%. After refinement, the peak shifts to approximately 35 words, accounting for 14.9%, whereas captions longer than 60 words are virtually absent.

**Word Frequency.** We analyzed the word-frequency profiles of captions before and after cleaning the 1.03M corpus, which are presented as word clouds in the right panel of Figure S1. The word clouds indicate that the captions emphasize microscopic morphological features such as "cells," "nuclei," and "stroma," as well as diagnostic descriptors including "inflammatory" and "stained." A comparison of the two word clouds shows an increased prevalence of morphology-related terms and a marked reduction in non-informative tokens such as "which" and "image."



**Figure S1** Visualization of the caption-length distribution generated by PathGen-LLaVa (top left) and its corresponding word-frequency cloud (top right), as well as the caption-length distribution after summarization by Qwen3-8B (bottom left) and the associated word-frequency cloud (bottom right).

## E.2 Analysis of the 68K Refined Subset

**Word Count.** We analyzed the caption lengths in the 68K Refined Subset, as shown on the left of the Figure S2. The distribution is symmetric and unimodal, with the most frequent length around 47 words, accounting for approximately 19%. Compared with the 1.03M Corpus, captions in the 68K subset are more extended, rarely shorter than 35 or longer than 55 words, due to the prompt-imposed 30–60-word constraint. This moderate length reduces redundancy while ensuring sufficient content to accurately describe the images, enabling the model to learn a broader range of knowledge.

**Word Frequency.** We also analyzed the word-frequency distribution of captions in the 68K Refined Subset, visualized as word clouds on the right side of Figure S2. The three subsets (8K, 10K, and 50K) exhibit a highly coherent vocabulary profile dominated by morphological, nuclear, cytoplasmic, stromal, and diagnostic descriptors. Unlike the 1M Corpus, where terms such as "nuclei," "cells," and "stroma" overwhelmingly dominate, the 68K Refined Subset exhibits a more balanced distribution of key pathological concepts. The captions in this refined subset are more detailed and make use of a richer and more uniformly distributed set of domain-specific terms, thereby providing higher-quality supervision that is advantageous for model evaluation and fine-tuning.

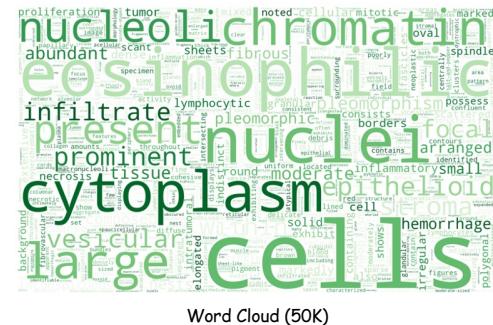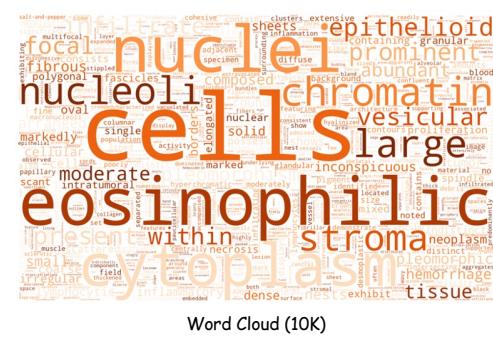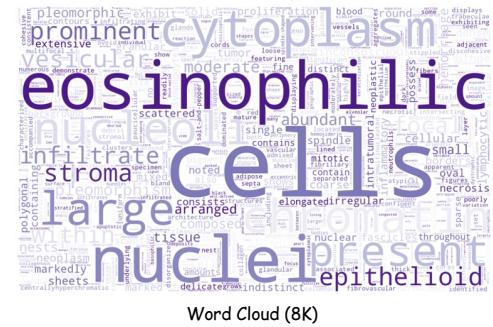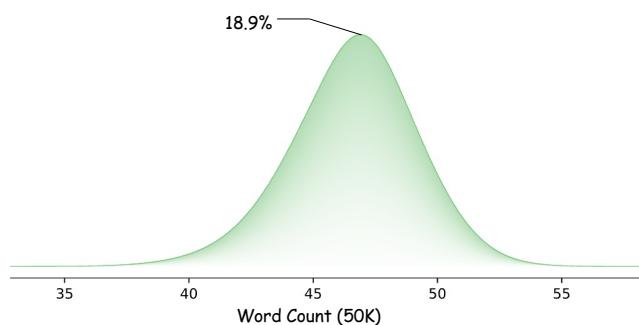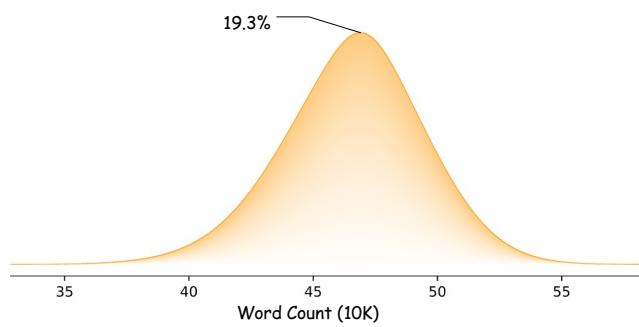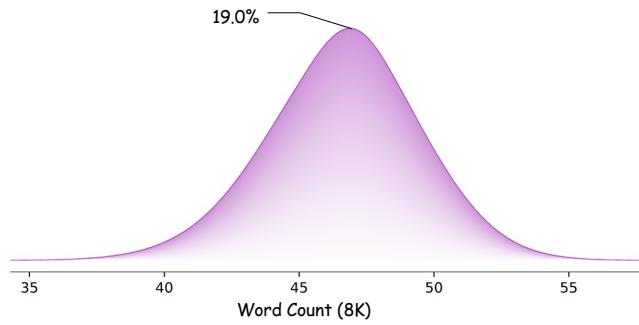**Figure S2** Visualization of the caption-length distributions and word-frequency clouds for the 8K, 10K, and 50K refined subsets (8K: top row; 10K: middle row; 50K: bottom row), with length distributions shown on the left and word-frequency clouds on the right.

**Table S1** Overall results of models on the PathMMU **test set**. The best-performing MLLM in each subset for general and medical/pathology MLLM is **in-bold**, and the second-best performing MLLM is <u>underlined</u>. An asterisk (∗) indicates results copied from the original source for the non–open-sourced model.

| | Unif. Model | Test Overall Tiny (1131) | ALL (9483) | PubMed Tiny (281) | ALL (3068) | SocialPath Tiny (210) | All (1661) | EduContent Tiny (255) | All (1938) | Atlas Tiny (208) | ALL (1007) | PathCLS Tiny (177) | ALL (1809) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Expert Performance | - | 71.8 | - | 72.9 | - | 71.5 | - | 69.0 | - | 68.3 | - | 78.9 | - |
| *General Multimodal Large Language Models* | | | | | | | | | | | | | |
| LLaVA-1.5-13B | ✗ | 39.1 | 37.7 | 45.2 | 41.4 | 39.6 | 40.1 | 35.0 | 39.9 | 46.5 | 43.4 | 25.9 | 23.7 |
| LLaVA-OneVision-7B | ✗ | 36.8 | 34.0 | 38.2 | 36.9 | 34.9 | 37.7 | 46.4 | 37.4 | 39.6 | 38.8 | 20.0 | 19.5 |
| Qwen3-VL-8B-Instruct | ✗ | 54.9 | 51.5 | 59.1 | 54.4 | 63.3 | 56.3 | 56.1 | 53.7 | 52.4 | 55.1 | 39.5 | 37.9 |
| Qwen3-VL-30B-A3B-Instruct | ✗ | 57.6 | 54.5 | 64.4 | 57.2 | 67.6 | 65.8 | 58.0 | 56.1 | 56.3 | 58.3 | 35.6 | 35.7 |
| BLIP3o-8B | ✔ | 40.0 | 37.6 | 40.5 | 40.6 | 42.6 | 39.2 | 47.5 | 42.3 | 40.7 | 38.7 | 24.7 | 25.3 |
| Show-o2-7B | ✔ | 40.8 | 40.4 | 43.7 | 48.8 | 43.8 | 40.4 | 42.4 | 41.9 | 42.8 | 40.4 | 28.2 | 24.4 |
| BAGEL-14B | ✔ | 57.3 | 52.3 | 63.0 | 56.2 | 61.4 | 59.6 | 60.8 | 54.1 | 64.4 | 61.3 | 29.9 | 31.9 |
| GPT-4V-1106 | - | 53.8 | 49.9 | 59.7 | 53.9 | 58.5 | 53.4 | 60.6 | 53.9 | 47.6 | 52.6 | 36.6 | 33.9 |
| Gemini-2.5 Pro | - | **69.0** | **68.0** | **74.5** | **71.8** | **72.3** | **68.9** | **72.4** | 69.8 | <u>66.8</u> | 69.5 | 53.9 | <u>58.1</u> |
| *Medical / Pathology-specific Multimodal Large Language Models* | | | | | | | | | | | | | |
| LLaVA-Med | ✗ | 25.5 | 26.8 | 29.2 | 28.5 | 29.8 | 28.2 | 23.3 | 27.4 | 21.7 | 30.4 | 22.0 | 20.1 |
| Quilt-LLaVA | ✗ | 45.4 | 41.2 | 46.8 | 41.9 | 46.3 | 45.9 | 51.0 | 45.0 | 46.5 | 43.7 | 32.8 | 30.1 |
| PathGen-LLaVA | ✗ | 59.8 | 58.4 | 59.3 | 59.9 | 60.2 | 58.7 | 60.1 | 60.5 | 64.3 | 65.1 | <u>54.4</u> | 49.7 |
| CPath-Omni∗ | ✗ | - | - | <u>74.0</u> | <u>69.9</u> | - | - | 69.8 | <u>70.6</u> | 65.9 | <u>70.6</u> | **75.7** | **79.0** |
| UniMedVL | ✔ | 54.5 | 50.6 | 58.0 | 54.9 | 57.1 | 56.7 | 55.7 | 56.0 | 66.3 | 59.4 | 30.5 | 26.9 |
| **UɴɪPᴀᴛʜ (Ours)** | ✔ | <u>68.3</u> | <u>65.7</u> | 72.9 | 66.4 | <u>67.9</u> | **68.4** | <u>70.1</u> | **73.9** | **79.2** | **77.7** | 46.1 | 46.6 |

## E.3 Spot-Check Validation of the 10K Test Set

To definitively validate the reliability of the "Gemini-2.5 Pro generation then GPT-5 review" automated pipeline, we additionally invited a domain-expert pathologist to conduct an independent spot-check quality control (QC) on 500 random samples from our 10K high-quality test set.

The reviewer's task was to assign each image-text pair to one of three categories. **3: Excellent** was defined as: The description is accurate, comprehensive, and professional, perfectly corresponding to all key pathological features in the image (Gold Standard). **2: Acceptable** was defined as: The description captures the main diagnostic features without factual errors, but may contain minor deficiencies, such as omitting a secondary feature, slight imprecision in non-critical terminology, or a minor deviation in descriptive focus (Still Usable for Evaluation). **1: Unusable** was defined as: The description contains severe factual errors, rendering it unsuitable as an evaluation benchmark, such as hallucinating key features not present in the image, misidentifying the primary cell type, or completely omitting the main diagnostic point of the image (Failure).

Upon reviewing the 500 random samples, the pathologist's evaluation was as follows: 43.4% of the image-text pairs were rated as 3: Excellent; 50.2% were rated as 2: Acceptable; and 6.4% were rated as 1: Unusable. This results in an Overall Usability Rate (*i.e.*, Excellent + Acceptable) of 93.6%. This extremely low "Unusable" rate (6.4%) strongly confirms the SOTA reliability of the automated data annotation pipeline we employed.

## E.4 1231-term Pathology Vocabulary

Our 1231-term pathology vocabulary, which was used to build the inverted index $\mathcal{I}$, is provided as a separate file ("vocabulary.txt") in the supplementary material bundle.

## E.5 Analysis of Data Leakage Risks

To ensure the integrity of our evaluation splits and proactively address potential concerns regarding overlap between the 8K Prototype Bank and the 10K Test Set, we conducted a rigorous validation. We computed the exhaustive pairwise visual feature similarity ($8,000 \times 10,000 = 80,000,000$ comparisons) between all images in the bank and all images in the test set. For this check, we used the UNI2-h extractor [8], which is the same high-performance backbone we employed for dataset-wide de-duplication in Section 3.1. The statistics confirmed that the sets are strictly disjoint: the average similarity was 0.1358 (standard deviation = 0.0686). Critically, the maximum similarity observed across all

**Table S2** Quantitative comparison of Visual Fidelity and Text–Image Alignment with merged T2I/I2I. FID/KID uses the Virchow2 extractor; Similarity and retrieval metrics use MUSK. ★ marks models fully fine-tuned on our large dataset. The best is **bold**, the second best is <u>underlined</u>.

| | Unif. Model | Visual Fidelity ↓ | | Text–Image Alignment ↑ | | | | |
|---|---|---|---|---|---|---|---|---|
| | | FID | KID | Sim. | Recall@10 | Recall@50 | mAP@10 | mAP@50 |
| Real Data | - | - | - | 0.557 | 13.40/– | 33.50/– | 5.28/– | 6.17/– |
| *General Text to Image Generation Models* | | | | | | | | |
| SD1.5★ [33] | ✘ | 1804.69 | 0.519 | 0.483 | 0.66/0.60 | 1.91/1.94 | 0.80/0.80 | 0.85/0.91 |
| SDXL★ [29] | ✘ | 2570.19 | 0.602 | 0.445 | 0.34/0.22 | 1.20/0.80 | 0.60/0.38 | 0.77/0.49 |
| Pixart-$\alpha$★ [6] | ✘ | 2574.85 | 0.685 | 0.482 | 1.14/0.60 | 4.54/2.65 | 1.58/0.72 | 1.74/0.93 |
| BLIP3o★ [5] | ✔ | 2008.75 | 0.550 | 0.455 | 1.55/1.50 | 5.87/5.80 | 2.39/1.93 | 2.58/2.26 |
| Show-o2★ [41] | ✔ | 1398.52 | 0.415 | **0.545** | **8.41**/<u>2.71</u> | **22.77**/<u>9.74</u> | **10.93**/<u>3.22</u> | **9.42**/<u>3.11</u> |
| *Pathological / Medical Text to Image Generation Models* | | | | | | | | |
| Pixcell [47] | ✘ | <u>929.30</u> | <u>0.259</u> | 0.524 | - | - | - | - |
| PathLDM [46] | ✘ | 1126.36 | 0.376 | 0.483 | 0.15/0.17 | 0.73/0.70 | 0.18/0.25 | 0.30/0.36 |
| UniMedVL [26] | ✔ | 1435.07 | 0.363 | 0.520 | 3.82/2.23 | 11.52/7.18 | 5.31/2.59 | 5.14/2.74 |
| **UNIPATH (Ours)** | ✔ | **484.38** | **0.192** | <u>0.538</u> | 7.55/**4.25** | <u>21.61</u>/**13.72** | <u>10.22</u>/**6.38** | <u>8.93</u>/**6.07** |

80 million pairs was 0.9416. This maximum value (0.9416) is below our defined de-duplication threshold of 0.95 (as detailed in Section 3.1). This result strongly confirms that our prototype bank and test set are strictly disjoint. It therefore eliminates any risk of data leakage via the Prototype Stream (PS) retrieval, validating the integrity of our Tier 2 (Alignment) and Tier 3 (Control) evaluations.

# F  Supplementary Quantitative Results

## F.1  Full Understanding Capability

To validate the diagnostic understanding capability of UNIPATH, which is critical for our Multi-Stream Control (MSC) module, we evaluated its performance on the comprehensive PathMMU benchmark [35]. The full results are presented in Table S1. As shown in the table, UNIPATH achieves an overall score of 65.7 on the full test set. This performance establishes UNIPATH as the SOTA among all evaluated open-source models, substantially outperforming other leading open-source pathology MLLMs, including PathGen-LLaVA (58.4) and the unified model UniMedVL (50.6). Furthermore, UNIPATH achieves the top score across all models (including closed-source systems) on the EduContent (73.9) and Atlas (77.7) sub-tasks. Its overall score also closely approaches that of top-tier closed-source models, such as Gemini-2.5 Pro (68.0). This strong understanding performance confirms that our frozen MLLM backbone provides the robust, phrasing-invariant semantics necessary to steer controllable generation.

## F.2  Additional Fidelity & Alignment Results

In our main evaluation (Section 5), the text-image alignment metrics and Patho-FID/KID metrics were based on the CONCH and UNI2-h backbones, respectively. To further validate the robustness and generality of our findings, we conducted an additional evaluation using two entirely independent, external backbones not used anywhere in our model pipeline: Virchow2 [38] for Visual Fidelity and MUSK [40] for Text-Image Alignment. This analysis confirms that our model's superior performance is a genuine advantage and not an artifact of a specific evaluator. The results are shown in Table **??**.

**Visual Fidelity (Virchow2 Backbone).**  The results using the Virchow2 feature extractor strongly reinforce our main findings. UNIPATH achieves a FID of 484.38 and a KID of 0.192. This performance is not just SOTA, but represents a massive improvement over the next-best model, Pixcell (FID: 929.30, KID: 0.259). This confirms that the superior visual fidelity of UNIPATH is a genuine model advantage, not an artifact of the UNI2-h evaluator.

**Figure S3** Gemini 2.5 Pro as Judge.

**Table S3** Few-shot downstream performance (Weighted F1) across different shots $K$. Values are Mean±Std (%); "Δ" *columns* report absolute change vs. original data in percentage points. Best is **bold**, second best is <u>underlined</u>.

| Model | K = 2 | | K = 4 | | K = 8 | | K = 16 | | K = 32 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Wgt. F1 | Δ | Wgt. F1 | Δ | Wgt. F1 | Δ | Wgt. F1 | Δ | Wgt. F1 | Δ |
| *Baselines (Only Real Data)* | | | | | | | | | | |
| Original Data | $67.34_{\pm2.37}$ | – | $76.42_{\pm3.31}$ | – | $81.43_{\pm1.88}$ | – | $83.85_{\pm1.48}$ | – | $86.88_{\pm0.88}$ | – |
| **UniPath (Ours)** | $50.31_{\pm8.92}$ | – | $69.33_{\pm2.81}$ | – | $76.89_{\pm2.35}$ | – | $78.30_{\pm0.99}$ | – | $81.05_{\pm0.77}$ | – |
| *Data Augmented Comparisons (Real Data with Generated Data)* | | | | | | | | | | |
| SD1.5 [33] | $60.81_{\pm5.31}$ | −6.53 | $71.14_{\pm4.28}$ | −5.28 | $77.03_{\pm2.17}$ | −4.40 | $81.94_{\pm2.94}$ | −1.91 | $85.35_{\pm1.13}$ | −1.53 |
| UniMedVL [26] | $63.86_{\pm4.54}$ | −3.48 | $74.94_{\pm3.04}$ | −1.48 | $80.69_{\pm2.03}$ | −0.74 | $83.11_{\pm1.11}$ | −0.74 | $86.75_{\pm1.27}$ | −0.13 |
| Pixcell [47] | $67.06_{\pm3.18}$ | −0.28 | $77.13_{\pm2.16}$ | +0.71 | $81.00_{\pm1.99}$ | −0.43 | $84.19_{\pm1.18}$ | +0.34 | $86.93_{\pm0.51}$ | +0.05 |
| Show-o2 [41] | <u>$68.68_{\pm5.31}$</u> | +1.34 | <u>$77.70_{\pm4.30}$</u> | +1.28 | <u>$81.76_{\pm2.35}$</u> | +0.33 | <u>$84.51_{\pm1.78}$</u> | +0.66 | <u>$86.97_{\pm1.04}$</u> | +0.09 |
| **UniPath (Ours)** | **$69.65_{\pm4.35}$** | +2.31 | **$79.14_{\pm1.82}$** | +2.72 | **$82.22_{\pm1.10}$** | +0.79 | **$85.39_{\pm1.37}$** | +1.54 | **$87.15_{\pm0.61}$** | +0.27 |

**Text-Image Alignment (MUSK Backbone).** The alignment metrics using the MUSK backbone provide a crucial, unbiased perspective.

- **CLIP-Score & T2I Retrieval:** Using MUSK, **Show-o2** achieves the highest CLIP-Score (0.545) and the best Report-Gen (T2I) retrieval metrics. This result is consistent with our main paper's findings (Table 1) and supports our hypothesis that Show-o2's architecture is well-optimized for general-purpose T2I alignment.

- **I2I Retrieval:** Critically, on the Real-Gen (I2I) retrieval task—which measures how close the generated images are to real images in this new feature space—UniPath achieves dominant SOTA performance across all four metrics (Recall@10/50, mAP@10/50). For instance, our mAP@10 (6.38) is nearly double that of the second-best Show-o2 (3.22).

**Conclusion.** These results, obtained from fully independent feature extractors, strongly validate our conclusions. Our model's superior visual fidelity and its ability to generate images that are most faithful to the real pathology manifold are thus demonstrated as robust and general findings, independent of the specific evaluator used.

### F.3 Detailed Few-Shot Classification Results

We provide the detailed numerical results for the Tier 4: Downstream Task Utility evaluation (Kather-CRC-2016 few-shot classification) in Table S3. This table contains the precise Mean±Std (Weighted F1) scores and the absolute change (Δ) for all K-shot values (K=2, 4, 8, 16, 32). These are the raw data used to generate Figure 4 in the main text.

**Figure S4 Inference-time sensitivity of prototype quantity ($K_m$).** Our default $K_m = 16$ achieves the optimal trade-off between visual fidelity (Patho-FID, ↓) and semantic alignment (CLIP-Score, ↑).

## F.4 MLLM and Human Judge: Setup & Reliability

**Human Expert Evaluation.** Here, we detail the implementation and reliability analysis of our human expert evaluation. We employed a panel of three trained annotators to conduct a blind pairwise comparison (UNIPATH vs. Baseline) on 500 image-text pairs randomly sampled from the 10K test set. During the evaluation, annotators were shown a text prompt and two anonymized images and were tasked with choosing which image better matched the prompt, without knowing which image was generated by UNIPATH. To validate the reliability of this evaluation, we measured the inter-annotator agreement. As shown in the analysis output, the panel achieved an overall Fleiss' Kappa of 0.7509, indicating "substantial" agreement. The per-model agreement was also robust, ranging from "Moderate" to "Almost Perfect" (UniMedVL: $\kappa = 0.8833$; show-o2: $\kappa = 0.7998$; Pixcell: $\kappa = 0.7950$; PixArt: $\kappa = 0.6502$; SD15: $\kappa = 0.5708$; BLIP3o: $\kappa = 0.5672$). The aggregated win/loss/tie statistics from this reliable panel were used to generate the human expert results in the main paper.

**Gemini-2.5 Pro as Judge.** In the main paper, we presented the "as-Judge" results from GPT-5 and the human expert panel. For completeness, we provide a parallel evaluation using Gemini-2.5 Pro as the judge, following the exact same experimental setup. The results are presented in Figure S3. As shown, Gemini-2.5 Pro's assessment is highly consistent with our other evaluations. It demonstrates a clear preference for UNIPATH against all baselines, preferring UNIPATH over the strongest baseline (Show-o2) in 71% of cases. This additional MLLM evaluation further corroborates our model's robust advantage in nuanced, human-aligned semantic understanding.

## F.5 MSC Sensitivity and Ablation Studies

We evaluate the inference performance of the Prototype Stream (PS) using Patho-FID and CONCH CLIP-Score. Unlike the fixed parameters of the High-Level Semantics (HLS) stream, the PS architecture allows for dynamic adjustments to the prototype bank size ($K_m$) and retrieval strategies at inference without retraining.

**Sensitivity to Prototype Quantity ($K_m$).** We analyzed the trade-off between context sufficiency and information density by varying the prototype count $K_m \in \{0, 4, 8, 16, 32\}$ on the 10K Test Set. Throughout these experiments, we maintained a consistent allocation strategy: $K_m/4$ for global text retrieval, $K_m/4$ for global vision retrieval, and $K_m/2$ for local fine-grained retrieval (assigning 2 prototypes to each of the top $K_m/4$ parsed keywords). As illustrated in Figure S4, the baseline without prototype guidance ($K_m = 0$) exhibits

**Table S4 Ablation on Retrieval Components.**

| Retrieval Configuration | Patho-FID ↓ | CLIP-Score ↑ |
|---|---|---|
| Global (Text-Only) | 87.52 | 0.327 |
| Global (Vision-Only) | 83.10 | 0.336 |
| Global (Hybrid) | 82.04 | 0.342 |
| Local | 91.45 | 0.325 |
| **Global + Local (Ours)** | **80.86** | **0.348** |

significantly inferior fidelity and alignment, validating the necessity of the PS stream. Increasing $K_m$ yields rapid gains up to $K_m = 16$, at which point our model achieves the optimal balance. Beyond this point, retrieving lower-ranked prototypes ($K_m = 32$) introduces irrelevant noise, diluting the conditioning signal and slightly degrading Patho-FID.

**Ablation on Retrieval Components.** We dissected the impact of our retrieval modules as reported in Table S4. Within the Global module, the Hybrid strategy (82.04 FID, 0.342 CLIP) consistently outperforms single-modality baselines, bridging the gap between Text-Only (87.52 FID) and Vision-Only (83.10 FID) retrieval. We also observe that relying solely on Local sparse retrieval yields the poorest fidelity (91.45 FID), indicating that sparse keywords alone lack sufficient generative context. However, the integration of Global and Local modules is transformative; the Full Strategy achieves the best overall performance (80.86 FID, 0.348 CLIP), confirming that fine-grained sparse guidance complements dense global context to maximize both visual fidelity and semantic alignment.

# G    Supplementary Qualitative Results

## G.1    PS Stream: Component-level Control

To visually validate the efficacy of our Prototype Stream (PS) in achieving component-level morphological control, we provide qualitative examples of its internal retrieval mechanism in Figure S5. As described in Section 4.2, our PS employs a hybrid retrieval strategy that combines Global Semantic Retrieval with Local Fine-grained Retrieval to capture both the holistic context and the specific morphological components of a prompt. Taking Figure S5a as an example, we illustrate the complete process for a complex "Generation Instruction."

- **Inputs and Generation:** The top-left panel shows the complex multi-part prompt, the original "Ground Truth" image, and our final "Generation Image." The generated image successfully synthesizes all specified pathological features, including "solid sheets," "marked pleomorphism," and "extensive hemorrhage," demonstrating high visual fidelity to the ground truth.

- **Global Semantic Retrieval:** The top-right panel shows the prototypes retrieved by the global strategy (both Text and Vision Feature Retrieval). These images capture the holistic gist or overall appearance of the prompt — such as the general pink/purple "H&E" color profile, high cellularity, and areas of hemorrhage.

- **Local Fine-grained Retrieval:** The bottom panel provides direct evidence of component-level control. Here, the prompt is parsed into specific keywords (*e.g.*, "arranged in solid sheets," "marked pleomorphism," "irregular contours," "extensive hemorrhage"). The inverted index ($\mathcal{I}$) then recalls prototypes that specifically and accurately match each individual component. For example, the prototypes for "extensive hemorrhage" are almost exclusively composed of red blood cells, while the prototypes for "marked pleomorphism" correctly show cells with high nuclear variation.

This visualization confirms that UNIPATH steers generation by combining these two complementary sets of prototypes, allowing it to render complex scenes with precise control over individual pathological components.

## G.2    SOTA Model Comparisons

To complement the examples presented in the main text, Figure S6 presents additional qualitative comparison sets covering a broader range of prompts. Consistent with the observations in the main paper, baseline methods frequently exhibit partial concept omission, morphological inconsistencies, or visually implausible artifacts when handling prompts containing multiple fine-grained pathological attributes. In contrast, UNIPATH systematically preserves the entirety of the described features and renders them with higher morphological fidelity. These visual examples provide a more faithful demonstration of UNIPATH's performance, capturing semantic and morphological details that automated metrics such as CLIP-Score fail to reflect fully.

## G.3    Gallery of Randomly Sampled Generations

To offer a complementary viewpoint on model behavior, Figure S7 centers solely on the visual quality of images generated by UNIPATH. We present a diverse set of sampled test-set cases, each paired with its corresponding Ground Truth image, spanning a broad spectrum of histopathological appearances such as epithelial structures, adipose tissue, smooth or skeletal muscle, collagenous stroma, and inflammatory infiltrates. Across these diverse cases, the side-by-side

**Inputs and Generation**

**Generation Instruction:**
The tumor is arranged in solid sheets of large epithelioid cells with abundant eosinophilic cytoplasm. The nuclei are large and demonstrate marked pleomorphism, with irregular contours, coarse chromatin, and prominent nucleoli. Extensive hemorrhage is present, dissecting through the neoplastic cell population. The stroma is scant and fibrous.

**Ground Truth**      **Generation Image**

**Global Semantic Retrieval**

**Text Feature Retrial**      **Vision Feature Retrial**

**Local Fine-grained Retrieval**

arranged in solid sheets    marked pleomorphism    irregular contours    extensive hemorrhage

**(a)** Case 1.



**Inputs and Generation**

**Generation Instruction:**
The tumor is composed of nests of large polygonal cells with abundant granular eosinophilic cytoplasm and distinct cell borders. Nuclei are round, centrally located, and contain vesicular chromatin with prominent eosinophilic nucleoli. There is mild nuclear pleomorphism. A focal peritumoral lymphocytic infiltrate is present within the fibrous stroma.

**Ground Truth**      **Generation Image**

**Global Semantic Retrieval**

**Text Feature Retrial**      **Vision Feature Retrial**

**Local Fine-grained Retrieval**

focal peritumoral lymphocytic infiltrate    abundant granular eosinophilic cytoplasm    mild nuclear pleomorphism    large polygonal cells

**(b)** Case 2.

**Figure S5** Visualization of the UNIPATH Prototype Stream (PS) hybrid retrieval mechanism. Both (a) and (b) illustrate how component-level control is achieved by combining Global Semantic Retrieval (top right) and Local Fine-grained Retrieval (bottom).

**Figure S6** Comparison of pathology image generation results across UNIPATH, PathLDM, SD15, PixArt-α, Pixcell, UniMedVL, and Show-o2 under different input captions. Colors in the captions denote distinct pathological features: Tissue/Cell Type, Nuclear Features, Cytoplasm, Hemorrhage.

comparison highlights that UNIPATH consistently produces images with high visual fidelity, well-preserved fine-grained morphological details, and realistic tissue textures, without introducing implausible artifacts. These qualitative examples provide a direct and intuitive assessment of generative realism that complements automated metrics such as FID, offering a more faithful reflection of the model's practical visual reliability.

## H Prompt Engineering Details

This section presents several prompts used during dataset construction and evaluation. Specifically, **(i)** The instructions for generating initial captions for the 68K Refined Subset using Gemini, as well as the cross-validation prompts utilized by GPT-5 to independently review quality and factual accuracy, shown in Figure S8 (Gemini-2.5 Pro) and Figure S9 (GPT-5); **(ii)** The prompts used to comprehensively assess generation quality with MLLMs serving as judges, shown in Figure S10; **(iii)** The prompts used to filter pathology terminology with an LLM, shown in Figure S11.

### H.1 Prompts for Re-annotation

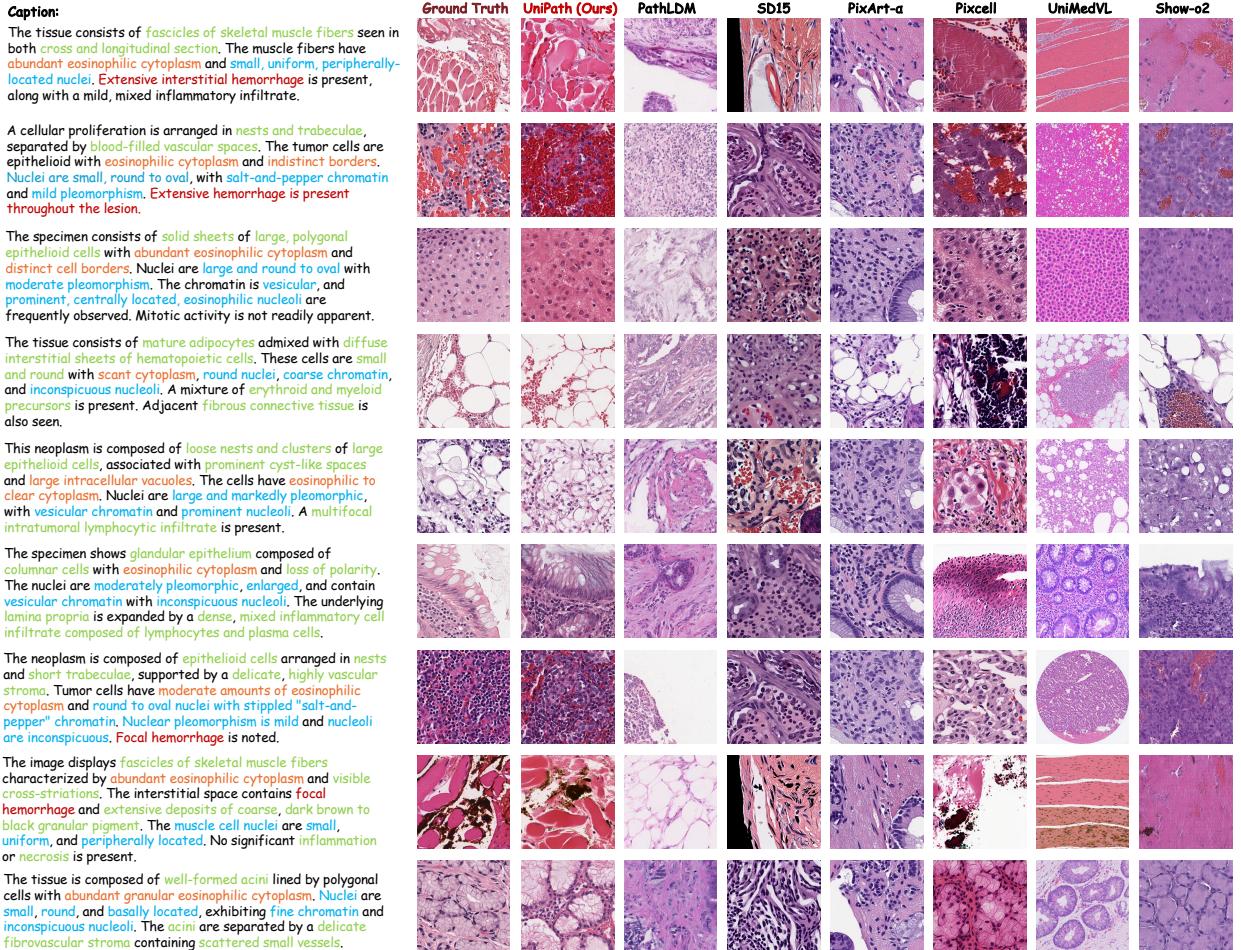For the construction of the 68K Refined Subset, we employed a two-stage prompt pipeline consisting of a Gemini-2.5 Pro-based caption generator and a GPT-5–based cross-modal validator. The entire process consumed 300M tokens, including both the input and output tokens.

**Stage 1: Captioning with Gemini-2.5 Pro.** The Gemini prompt instructs the model to inspect each H&E ROI and produce a structured JSON output containing Lite-schema labels, along with a 30–60-word morphological description, without any diagnosis. The prompt provides explicit enumeration rules (*e.g.*, nuclear size, pleomorphism, stromal reaction, types of inflammation) and a style specification that requires objective, declarative wording while forbidding diagnostic terms, negative-absence phrasing, or modality-related metadata. This design ensures that the initial captions remain focused on morphology, remain consistent, and can be used in later-generation tasks.

**Stage 2: Verification with GPT-5.** The second-stage GPT-5 prompt performs strict factual verification of Gemini's output. It checks whether each label is visually supported at the given magnification and flags any ambiguous or contradicted features as errors. In addition to visual factuality, it enforces mandatory textual rules (*e.g.*, no diagnosis, correct style) and performs minimal schema checks to ensure alignment with the predefined JSON format. This cross-modal validator effectively removes inconsistent or hallucinated descriptions, making sure that only high-quality annotations enter the final dataset.

### H.2 Prompts for MLLM-as-Judge

To systematically evaluate the alignment between generated histopathology images and textual descriptions, we designed a cross-modal evaluation prompt. The prompt enables comparison between the images generated by UNIPATH and those of other baseline models for a given caption. Its functionality includes identifying visual features in the caption, assessing whether each image feature is supported or contradicted, and producing a quantitative alignment judgment (WIN / TIE / LOSS) based on a predefined hierarchy of histological features (*e.g.*, Architecture, Cytology, Nuclear features).

### H.3 Prompts for Pathology Vocabulary Filtering

We also designed a rule-driven prompt to curate pathology feature phrases for downstream modeling. The prompt enables a model to evaluate a list of short text phrases, acting as a board-certified anatomic pathologist and computational pathology NLP expert. Its functionality includes determining whether each phrase represents a complete and discriminative histopathologic feature according to a precision-first hierarchy (*e.g.*, nuclear, cytoplasmic, cellular lineage, architectural, cytologic atypia, qualified inflammatory or hemorrhagic features). The output preserves the original input order and formatting, is strictly plain text, and contains no additional explanation or modification, ensuring reproducibility and direct applicability for downstream modeling.

**Figure S7** Comparative visualization of Ground-Truth and the corresponding pathology images synthesized by our model.

```json
{
  "system_role": "You are an expert anatomic pathologist and careful data labeler.",

  "user_prompt": "Your job: (1) inspect the provided H&E-stained histology ROI image ONLY,
(2) assign SINGLE-CHOICE labels for each field using the Lite schema, and
(3) write a diagnosis-free, generation-ready morphology description (30-60 words).
Do NOT explain your reasoning. Do NOT output any content in negative absence phrasing.
Output strictly in the required JSON format.",

  "guidance": {
    "enumeration_reminders": {
      "architecture_primary": [
        "glandular (including tubular)", "solid sheets", "nests", "trabeculae",
        "papillary", "cribriform", "single-file", "rosette",
        "reticular", "fascicles", "storiform"
      ],
      "cytology_type": ["epithelioid", "spindle", "signet_ring", "sarcomatoid"],
      "nuclei.size": ["small", "moderate", "large"],
      "nuclei.pleomorphism": ["none", "mild", "moderate", "marked"],
      "nuclei.chromatin": ["fine", "vesicular", "coarse", "salt_and_pepper"],
      "nuclei.nucleoli": ["none", "inconspicuous", "prominent", "macronucleoli"],
      "necrosis": ["none", "focal", "confluent", "geographic", "comedo"],
      "hemorrhage": ["none", "focal", "extensive"],
      "calcification": ["none", "microcalcifications", "psammoma_bodies", "dystrophic"],
      "inflammation.location": [
        "none", "intratumoral", "peritumoral", "both in and around tumor"
      ],
      "inflammation.extent": ["none", "diffuse", "focal", "multifocal"],
      "inflammation.dominant_type": [
        "none", "lymphocytes", "neutrophils", "plasma_cells",
        "eosinophils", "macrophages", "giant_cells", "mixed"
      ],
      "stromal_reaction": [
        "none", "fibrous", "hyalinized", "myxoid", "desmoplastic",
        "sclerotic", "mucin_pool", "osteoid", "chondroid"
      ],
      "invasion.tumor_budding": ["absent", "low", "intermediate", "high"]
    },

    "description_guidance": {
      "length": "30-60 words",
      "style": ["declarative", "objective", "diagnosis_free"],
      "sequence_suggestion": [
        "architecture and composition", "cytology type and cell morphology",
        "nuclear features", "mitotic density (if visible)",
        "necrosis pattern (if visible)", "stromal reaction (if visible)",
        "inflammation level and type (if visible)",
        "hemorrhage or calcification (if visible)"
      ],
      "avoid": [
        "diagnosis terms", "grading or staging", "IHC or molecular",
        "treatment", "percentages or exact_counts", "modality, noise, scale bar ruler labels"
      ]
    }
  },
```
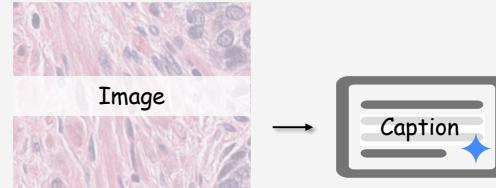
**Prompts for Annotation**

Image → Caption

```json
  "output_format": "JUST JSON",

  "output_contract": {
    "schema": {
      "architecture_primary": "...",
      "cytology_type": "...",
      "cell_morphology": "...",
      "nuclei": {
        "size": "...",
        "pleomorphism": "...",
        "chromatin": "...",
        "nucleoli": "..."
      },
      "necrosis": "...",
      "hemorrhage": "...",
      "calcification": "...",
      "inflammation": {
        "location": "...",
        "extent": "...",
        "dominant_type": "..."
      },
      "stromal_reaction": "...",
      "invasion": { "tumor_budding": "..." }
    },
    "llm_description": "..."
  },

  "input": {
    "meta_hints": {
      "stain": "H&E",
      "magnification": "20x",
      "patch_size_px": 384
    }
  }
}
```

Gemini

**Figure S8** Prompts used for the preliminary annotation of pathology images with Gemini,-2,5 Pro formatted in JSON. Distinct colors are applied to differentiate hierarchical levels within the JSON structure: **top-level keys** and **secondary levels**. The overall workflow is illustrated in the upper-right panel.

```json
{
    "system_role": "You are an expert anatomic pathologist and cross-modal validator. Your TOP priority is to detect visual
factual errors in Gemini's labels by examining the provided H&E ROI image. Formatting checks are secondary and minimal.",

    "user_prompt": "Tasks: (1) Inspect the ROI image. (2) For EACH field in candidate_json, decide if the label is clearly
supported by the image and not contradicted; if unclear at this magnification, FAIL. (3) Check the description for
mandatory style/consistency. (4) Do MINIMAL schema checks (required fields present; enums exact; values are strings).
Decision: if any visual label is unsupported/contradicted/uncertain → FAIL; else if description violates mandatory rules
→ FAIL; else if required schema invalid → FAIL; otherwise → PASS. Return exactly as in output_format/output_contract.",

    "guidance": {
        "priority_order": ["visual_factuality", "text/image_consistency", "minimal_schema"],

        "enumeration_reminders": {
            "architecture_primary": ["glandular (including tubular)", "solid sheets", "nests", "trabeculae", "papillary",
"cribriform", "single-file", "rosette", "reticular", "fascicles", "storiform"],
            "cytology_type": ["epithelioid", "spindle", "signet_ring", "sarcomatoid"],
            "nuclei.size": ["small", "moderate", "large"],
            "nuclei.pleomorphism": ["none", "mild", "moderate", "marked"],
            "nuclei.chromatin": ["fine", "vesicular", "coarse", "salt_and_pepper"],
            "nuclei.nucleoli": ["none", "inconspicuous", "prominent", "macronucleoli"],
            "necrosis": ["none", "focal", "confluent", "geographic", "comedo"],
            "hemorrhage": ["none", "focal", "extensive"],
            "calcification": ["none", "microcalcifications", "psammoma_bodies", "dystrophic"],
            "inflammation.location": ["none", "intratumoral", "peritumoral", "both in and around tumor"],
            "inflammation.extent": ["none", "diffuse", "focal", "multifocal"],
            "inflammation.dominant_type": ["none", "lymphocytes", "neutrophils", "plasma_cells", "eosinophils", "macrophages",
"giant_cells", "mixed"],
            "stromal_reaction": ["none", "fibrous", "hyalinized", "myxoid", "desmoplastic", "sclerotic", "mucin_pool",
"osteoid", "chondroid"],
            "invasion.tumor_budding": ["absent", "low", "intermediate", "high"]
        },

        "visual_checks": {
            "architecture_primary": "Pattern must be clearly present at this magnification; ambiguity → FAIL.",
            "cytology_type": "Cell morphology must match type; ambiguity → FAIL.",
            "nuclei": "Size/pleomorphism/chromatin/nucleoli must visually concord.",
            "necrosis": "Label must match visible necrotic debris/ghost cells; mismatch → FAIL.",
            "hemorrhage": "Extravasated RBCs/extent must match; mismatch → FAIL.",
            "calcification": "Basophilic deposits consistent with label; mismatch → FAIL.",
            "inflammation": "Location, extent, and dominant cell type must be appreciable; not determinable → FAIL.",
            "stromal_reaction": "Stromal pattern (desmoplastic/myxoid/etc.) appreciable; not determinable → FAIL.",
            "invasion.tumor_budding": "Assess only if invasive front is visible; otherwise → FAIL."
        },

        "text_rules": [
            "llm_description: 30-60 words, objective English, diagnosis-free.",
            "No negative-absence phrasing; no grading/staging/IHC/treatment; no digits or %; no modality/magnification
mentions.",
            "Must not mention features labeled 'none'; must align with image and labels."
        ],

        "minimal_schema_rules": [
            "Required fields present; enumerated values exact (case/spacing); all values are strings.",
            "Extra keys tolerated and ignored."
        ]
    },

    "output_format": "JUST JSON",

    "output_contract": {
        "schema": {
            "result": "..."
        },
        "allowed_values": {
            "result": ["PASS", "FAIL"]
        },
        "strict_return_values": true
    },

    "input": {
        "candidate_json": "{{GEMINI_OUTPUT_JSON}}",
        "roi_image": "{{ROI_IMAGE}}",
        "meta_hints": {
            "stain": "H&E",
            "magnification": "20x",
            "patch_size_px": 384
        }
    }
}
```

**Figure S9** Prompts used to perform cross-validation of Gemini-2.5 Pro's preliminary pathology annotations via GPT-5, encoded in JSON. Distinct colors are applied to differentiate hierarchical levels within the JSON structure: **top-level keys** and **secondary levels**. The overall workflow is illustrated in the upper-right panel.

```
{
  "system_role": "You are an expert anatomic pathologist and cross-modal judge. Your goal is
to compare two H&E images (Ours vs. Baseline) against a Caption and decide which image aligns
better.",

  "user_prompt": "Tasks: (1) Identify the visual claims in the Caption. (2) For both 'Ours'
(image 1) and 'Baseline' (image 2), determine which claims are supported, contradicted, or not
visible. (3) Decide if 'Ours' aligns better (WIN), equally (TIE), or worse (LOSS) than
'Baseline'. (4) Prioritize factual visual accuracy. Output exactly in the required JSON
format.",

  "guidance": {
    "label_semantics": "WIN = Ours aligns better. TIE = Alignment is equal. LOSS = Ours aligns
worse.",

    "evaluation_axes_priority": [
      "Architecture",
      "Cytology",
      "Nuclear features",
      "Necrosis/Inflammation",
      "Stromal reaction",
      "Other features (Hemorrhage/Calcification)"
    ],

    "alignment_rules": [
      "An image aligns better if it supports more of the Caption's claims and contradicts
fewer claims than the other image.",
      "A feature described as 'none' or 'absent' in the Caption must not be present in the
image; if it is present, this counts as a contradiction."
    ],

    "tie_breakers": [
      "If alignment seems equal, use the priority order in 'evaluation_axes_priority' (e.g.,
matching Architecture is more important).",
      "If still equal, the image with fewer contradictions wins.",
      "If still equal, return TIE."
    ],

    "cautions": [
      "Judge only what is visible at this magnification.",
      "Do not infer or use diagnostic terms not in the Caption.",
      "If a claim cannot be verified (e.g., out of frame, too small), it is considered 'not
supported' (but not a contradiction)."
    ]
  },

  "output_format": "JUST JSON",

  "output_contract": {
    "schema": {
      "result": "..."
    },
    "allowed_values": {
      "result": ["WIN", "TIE", "LOSS"]
    },
    "strict_return_values": true,
    "notes": "Subject is always Ours (first image).
WIN=T(Ours)>Baseline;   LOSS=Ours<Baseline;   TIE=roughly
equal."
  },

  "input": {
    "caption": "{{CAPTION}}",
    "image_ours": "{{IMAGE_OURS}}",
    "image_baseline": "{{IMAGE_BASELINE}}",
    "meta_hints": {
      "stain": "H&E",
      "magnification": "20x (if known)"
    }
  }
}
```
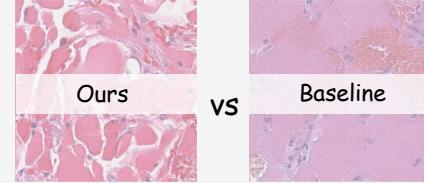
**Prompts for MLLM-as-Judge**

Ours  vs  Baseline

**Figure S10** Prompts used for evaluating pathology images generated by our model and baseline methods using an MLLM, expressed in JSON format. Distinct colors are applied to differentiate hierarchical levels within the JSON structure: **top-level keys** and **secondary levels**. The overall workflow is depicted in the upper-right panel.

```
"""System:
You are a board-certified anatomic pathologist and a computational pathology NLP expert.
Your job is to decide whether each short phrase from pathology image descriptions should be
KEPT for downstream modeling.

Input:
- You will receive a plain text list of phrases, one per line. No JSON.

Output:
- Output ONLY the KEPT phrases, one per line, nothing else.
- Maintain the original input order for any kept items.
- Do NOT rewrite/normalize/case-fold/trim/repair any phrase; echo it exactly.
- No explanations, no JSON, no headings, no bullet points, no extra whitespace.

Decision policy (precision-first; fragments must be dropped):

[HARD DROP — fragments & scaffolding]  (these override everything)
1) Starts or ends with a stopword: and, with, in, of, to, for, is, are, the, a, an, by, on,
within, at.
   (Examples: "and …", "with …", "… and", "… is", "… the" → DROP)
2) Incomplete coordination: contains "and" where one side lacks a morphologic head (e.g.,
adjective-only).
3) Scaffolding without object: "is composed", "is arranged", "consists" (± "of"), "arranged
in", "composed of".
4) Prepositional fragments: begins with a preposition (in/on/with/of/to/by/at/within) unless
the remainder is a complete standalone morphologic head.
5) Auxiliary-verb tails: ends with "is/are/was/were" → DROP.

[DROP — generic/degree-only without feature]
6) Degree-only tokens without a concrete feature: mild, moderate, marked, markedly, abu(when
not tied to a feature).
7) Bare generic structures: cell/cells, nucleus/nuclei, cytoplasm, stroma, tissue,
architecture, borders, cellular, nuclear, solid (alone).
8) Over-generic disease/process words: tumor, neoplasm, proliferation (bare/broad).
9) Unqualified supportive findings: infiltrate, inflammation, hemorrhage, necrosis
(unqualified) → DROP.

[KEEP — specific, standalone, discriminative features]
A) Nuclear features: "vesicular chromatin", "coarse chromatin", "prominent nucleoli",
"inconspicuous nucleoli", "eosinophilic nucleoli".
   - Also KEEP compound: "chromatin and <qualifier> nucleoli" (both heads explicit).
B) Qualified cytoplasmic features: "eosinophilic cytoplasm", "abundant eosinophilic cytoplasm".
C) Lineage/morphology: "epithelioid cells", "spindle cells", "polygonal cells".
   - Copular complete forms like "cells are epithelioid/spindle/polygonal" → KEEP.
D) Architectural patterns: "nests", "sheets", "solid sheets", "fascicles".
   - Prepositional variants like "in solid sheets" → fragment → DROP by HARD DROP #4.
E) Cytologic atypia & mitotic features: "pleomorphism", "marked pleomorphism", "moderate
pleomorphism", "mitotic figures".
   - Lone adjective "mitotic" without head noun → DROP.
F) Qualified inflammation/hemorrhage: keep only with lineage/site/extent (e.g., "lymphocytic
infiltrate", "mixed inflammatory infiltrate", "focal hemorrhage", "intratumoral hemorrhage").
   - If trailing auxiliary (e.g., "lymphocytic infiltrate is") → DROP (HARD DROP #5).

Tie-break:
- If rules conflict or the phrase is ambiguous/underspecified, DROP (output nothing).
"""
```

**Prompts for Pathology Vocabulary Filtering**

**Plain Text List**
1. prominent nucleoli
2. inflammation
3. solid sheets
4. and focal necrosis
5. abundant eosinophilic cytoplasm
6. tumor
7. spindle cells
8. marked
9. lymphocytic infiltrate
10. consists of spindle and epithelioid

**Plain Text List**
1. prominent nucleoli       #KEEP-A
2. inflammation       #DROP-9
3. solid sheets       #KEEP-D
4. and focal necrosis       #DROP-1
5. abundant eosinophilic cytoplasm       #KEEP-B
6. tumor       #DROP-8
7. spindle cells       #KEEP-C
8. marked       #DROP-6
9. lymphocytic infiltrate       #KEEP-F
10. consists of spindle and epithelioid       #DROP-3,2

**Plain Text List**
1. prominent nucleoli
2. solid sheets
3. abundant eosinophilic cytoplasm
4. spindle cells
5. lymphocytic infiltrate

**Figure S11** Prompts used for filtering the pathology vocabulary. Distinct colors are employed to differentiate hierarchical levels within the instructions: **top-level components** and **secondary elements**. The detailed filtering workflow is shown on the right.