

A Knowledge-enhanced Pathology Vision-language Foundation Model for Cancer Diagnosis

Xiao Zhou¹, Luoyi Sun^{1,2}, Dexuan He³, Wenbin Guan⁴, Ruifen Wang⁴, Lifeng Wang⁴,
Xin Sun⁵, Kun Sun⁶, Ya Zhang^{1,3}, Yanfeng Wang^{1,3}[†] and Weidi Xie^{1,3}[†]

¹Shanghai Artificial Intelligence Laboratory ²Zhejiang University

³School of Artificial Intelligence, Shanghai Jiao Tong University

⁴Department of Pathology, Xin Hua Hospital Affiliated to Shanghai Jiao Tong University School of Medicine ⁵Clinical Research and Innovation Unit, Xinhua Hospital Affiliated to Shanghai Jiao Tong University School of Medicine ⁶Department of Pediatric Cardiology, Xinhua Hospital Affiliated to

Shanghai Jiao Tong University School of Medicine

Deep learning has enabled the development of highly robust foundation models for various pathological tasks across diverse diseases and patient cohorts. Among these models, vision-language pre-training, which leverages large-scale paired data to align pathology image and text embedding spaces, and provides a novel zero-shot paradigm for downstream tasks. However, existing models have been primarily data-driven and lack the incorporation of domain-specific knowledge, which limits their performance in cancer diagnosis, especially for rare tumor subtypes. To address this limitation, we establish a **Knowledge-Enhanced Pathology (KEEP)** foundation model that harnesses disease knowledge to facilitate vision-language pre-training. Specifically, we first construct a disease knowledge graph (KG) that covers 11,454 human diseases with 139,143 disease attributes, including synonyms, definitions, and hypernym relations. We then systematically reorganize the millions of publicly available noisy pathology image-text pairs, into 143K well-structured semantic groups linked through the hierarchical relations of the disease KG. To derive more nuanced image and text representations, we propose a novel knowledge-enhanced vision-language pre-training approach that integrates disease knowledge into the alignment within hierarchical semantic groups instead of unstructured image-text pairs. Validated on 18 diverse benchmarks with more than 14,000 whole slide images (WSIs), KEEP achieves state-of-the-art performance in zero-shot cancer diagnostic tasks. Notably, for cancer detection, KEEP demonstrates an average sensitivity of 89.8% at a specificity of 95.0% across 7 cancer types, significantly outperforming vision-only foundation models and highlighting its promising potential for clinical application. For cancer subtyping, KEEP achieves a median balanced accuracy of 0.456 in subtyping 30 rare brain cancers, indicating strong generalizability for diagnosing rare tumors. All codes and models will be available for reproducing our results.

1 Introduction

Pathology diagnosis remains the golden standard in clinical applications for cancer diagnosis. Over the past decade, advancements in deep learning for computer vision have catalyzed significant progress in computational pathology, resulting in the development of specialized models based on both full supervision [41, 43, 33, 10, 22] or weak supervision [9, 53, 36, 14, 32, 50, 15, 37]. While these approaches show promise, they are generally limited by the high cost and scarcity of annotations, as well as their restricted generalizability across diverse datasets. To address these limitations, self-supervised learning (SSL) strategies [49, 11, 12, 27] have emerged as a promising alternative, enabling to pre-train the model on large collections of unlabeled pathological images, acting as a versatile feature extractor for a series of downstream tasks [17, 13, 47, 51, 57]. However, the vision-only SSL models still require fine-tuning on diverse labeled datasets for specific tasks, limiting their scalability to low-annotation settings, particularly in rare cancer subtype classification tasks.

Recently, the rise of vision-language models [38, 26] has enabled a new paradigm for computational pathology, offering novel avenues in cancer diagnosis. By jointly leveraging visual and textual data, vision-language models introduce free-text descriptions as supervision signals for pathology image representation learning, potentially improving diagnostic accuracy even in data-scarce settings. This approach could enhance generalizability and reduce reliance on extensive labeled datasets, addressing the limitations of vision-only models in distinguishing

* Corresponding author. Email addresses: {wangyanfeng, weidi}@sjtu.edu.cn

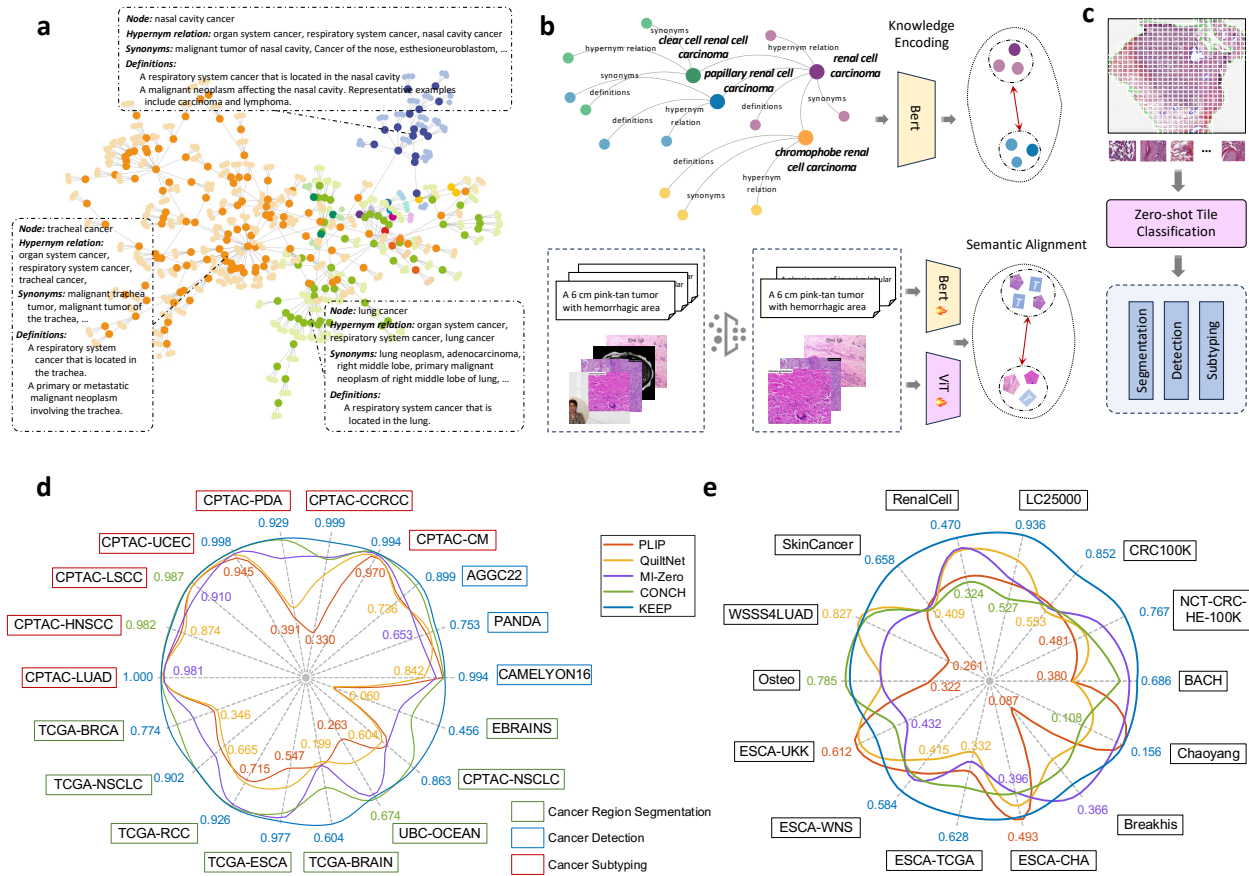


Figure 1 | Overview of KEEP. **a.** Example disease structure in the constructed knowledge graph. Each node represents a disease, consisting of three attributes types: hierarchical relations, synonyms, and definitions, as indicated by the dashed line box. **b.** The knowledge encoding and vision-language alignment stage for the KEEP model. A BERT-based text encoder is initially trained to encode the disease knowledge through metric learning. A knowledge-enhanced vision-language pre-training approach is proposed to align pathology semantic groups with filtered images and augmented captions. **c.** For downstream cancer diagnostic tasks, including cancer region segmentation, cancer detection, and cancer subtyping, whole slide images (WSIs) are divided into tile images for zero-shot classification, with the results of each tile combined to determine the final diagnostic decision. **d.** Performance comparison of cancer diagnosis with the state-of-the-art methods on 18 benchmarks with more than 14,000 WSIs. **e.** Performance comparison of tile-level classification with the state-of-the-art methods on 14 benchmarks. The inner and outer numbers indicate the worst and best results, respectively.

complex cancer subtypes. To create a joint embedding space for vision and language, existing models are trained on pathology image-text pairs gathered from in-house resources (MI-Zero [35], CONCH [34], and PRISM [42]) or public websites, such as Twitter (PLIP [23]) and YouTube videos (QuiltNet [25]), employing straightforward contrastive learning to align images with their corresponding captions.

Despite achieving impressive performance across various downstream tasks, existing pathology vision-language models, including PLIP and QuiltNet, face significant limitations due to the relatively small scale of pathology image-text datasets like OpenPath [23] and Quilt1M [25]. Compared to the expansive datasets used in general computer vision, these pathology-specific resources are orders of magnitude smaller and often sourced from non-professional websites, leading to considerable data noise and limited quality, for example, the annotations accompanying these images tend to be brief, unstructured, and lacking in comprehensive medical knowledge. Such deficiencies hinder the models' ability to accurately recognize and differentiate between various disease manifestations and their corresponding pathological features.

In this paper, we introduce a **Knowledge-Enhanced Pathology** vision-language foundation model, **termed**

as **KEEP**. It leverages public pathology image-text data in conjunction with hierarchical medical knowledge. We begin by constructing a disease knowledge graph (KG) that includes 11,454 human disease entities, their synonyms, definitions, and hypernym relations, derived from authoritative sources like the Disease Ontology (DO) [40] and the Unified Medical Language System (UMLS) [5]. As illustrated in Figure 1a, this extensive KG integrates a comprehensive array of medical information, providing a structured framework to support the model’s learning process. Subsequently, as depicted in Figure 1b, a language model is employed to encode this hierarchical knowledge, which then guides the vision-language representation learning, enhancing the model’s ability to interpret and utilize the complex medical data effectively.

To further enhance the quality of public pathology image-text datasets, we developed a novel framework for meticulously filtering out noisy images and captions. This framework restructures the cleaned data into semantic groups linked by the hierarchical relations in the disease knowledge graph. Such reorganization not only improves the quality of the training data, but also ensures that the model’s pre-training is guided by clinically relevant and semantically rich contexts. The structured pre-training process, shown in Figure 1b, aligns these semantically grouped data, significantly improving the model’s ability to understand and categorize pathology images.

To validate the effectiveness of KEEP, as well as to benchmark other pathology foundation models, we conduct comprehensive evaluations across 18 diverse benchmarks involving over 14,000 whole slide images (WSIs). These evaluations encompass three critical tasks: cancer region segmentation, cancer detection, and cancer subtyping, as depicted in Figure 1c. Additionally, we also assess different models on tile-level tasks, including cross-modal retrieval on four datasets and zero-shot tile image classification across 14 datasets. To our knowledge, this paper presents the most comprehensive evaluation of cancer diagnosis to date. Our quantitative and qualitative experimental results demonstrate that KEEP substantially outperforms state-of-the-art models (Figure 1d,e) by incorporating domain-specific knowledge. Notably, in cancer detection, KEEP’s prediction of tumor area ratios on WSIs achieves an average sensitivity of 89.8% at a specificity of 95% across 7 cancer types, significantly outperforming the CHIEF [50] model. Furthermore, KEEP shows remarkable generalizability in diagnosing rare diseases, achieving a median balanced accuracy of 0.456 in subtyping 30 rare brain cancers—surpassing the CONCH model [34] by 8.5 percentage points. These results underscore KEEP’s superior performance, particularly in tasks that benefit from its integration of domain-specific knowledge.

2 Results

2.1 Overview of KEEP

Zero-shot cancer diagnosis is a key downstream application of pathology vision-language foundation models, that well suits the scenarios for diagnosing rare tumors with very few labeled cases. Current foundation models, typically fed with small gridded tiles from WSIs, integrate embedding features (in vision-only models) or predicted labels (in vision-language models) to derive final diagnostic decisions. While vision-language models offer a more explainable approach by explicitly identifying cancerous tiles, their performance in diagnosing rare diseases is still limited.

In this paper, we introduce KEEP, a novel vision-language foundation model that leverages a disease knowledge graph to enhance both prediction performance and explainability in cancer diagnosis. By aligning semantic groups with a well-defined knowledge structure, KEEP outperforms existing models like PLIP and CONCH, which rely on simple contrastive learning of image-text pairs. This knowledge-driven approach deepens the model’s understanding of various disease characteristics and ensures stronger semantic alignment across diagnostic tasks.

Specifically, we first curate a hierarchical knowledge graph that consists of 11,454 human diseases and corresponding disease attributes, including disease synonyms, definitions, and hypernym relations, as shown in Figure 1a. For instance, *lung squamous cell carcinoma, also known as epidermoid cell carcinoma of the lung* (**synonym**), *is a carcinoma that derives from squamous epithelial cells* (**definition**), *which is a subtype of non-small cell lung cancer and also a subtype of squamous cell carcinoma* (**hypernym relations**). We then train a language model to encode this knowledge graph, with the goal of learning the hierarchical relationships

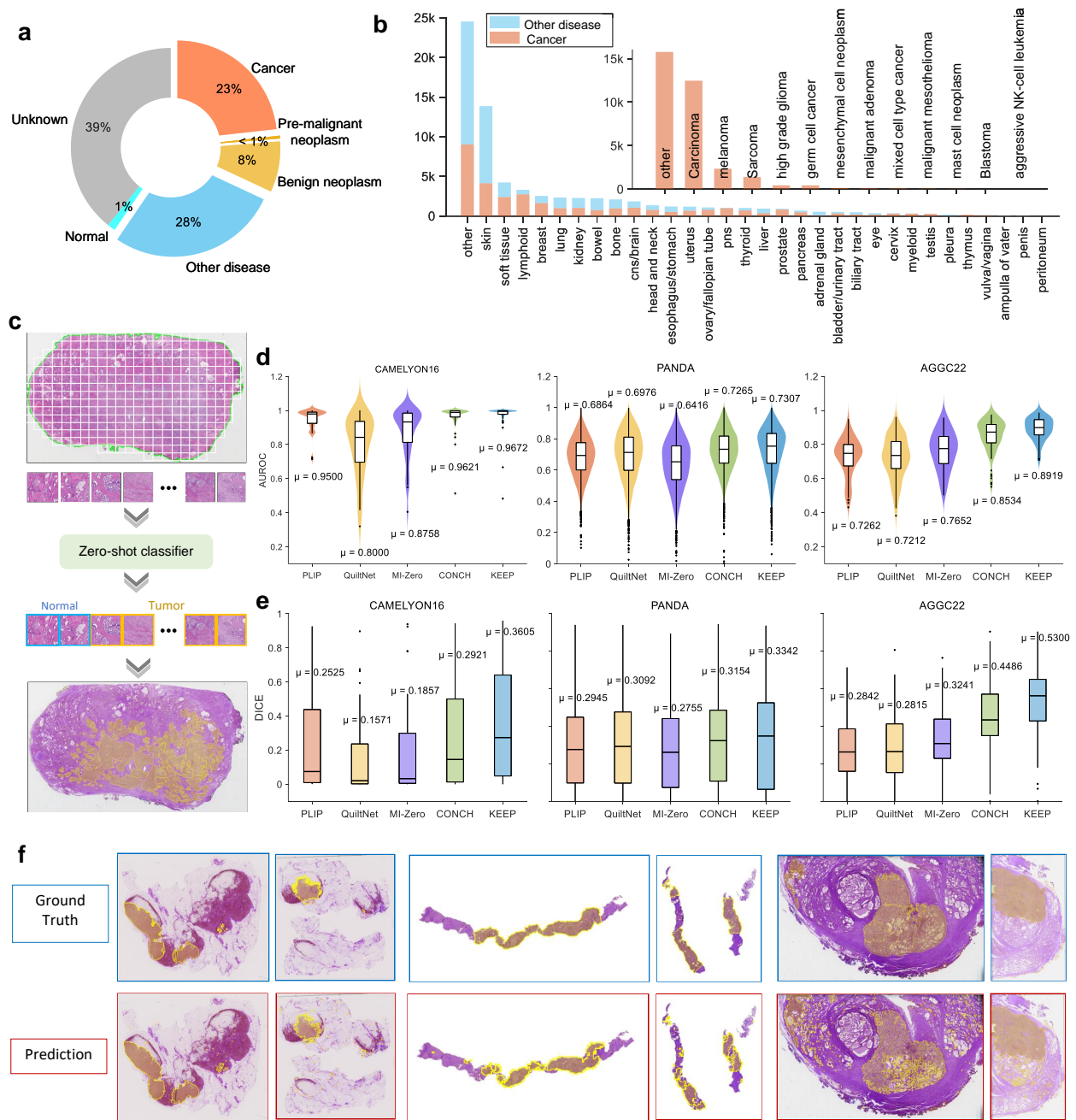


Figure 2 | Statistics of semantic groups and zero-shot cancer region segmentation results. **a.** Statistics of all semantic groups, organized by structuring one million noisy pathology image-text pairs with the guidance of disease KG. More than 60% semantic groups are linked to specific disease nodes. **b.** The anatomy and cell type distribution of the semantic groups with known disease labels, where "other" denotes the anatomy or cell type remains unknown. The anatomical taxonomy is based on OncoTree [31]. The anatomy and tumor types with the largest number of semantic groups are skin and carcinoma, respectively. **c.** The scheme of zero-shot segmentation on WSIs, where individual tiles undergo binary classification and are then combined to delineate the cancerous region. **d,e.** Performance comparisons of AUROC and DICE scores for various models, including PLIP, QuiltNet, MI-Zero, CONCH, and our proposed KEEP, across three WSI datasets: CAMELYON16 (48 WSIs), PANDA (10,494 WSIs), and AGGC22 (128 WSIs). The box plots present the median, first, and third quartiles of results, with μ indicating the average performance. The DICE is calculated using the average threshold corresponding to the optimal cutoff point of ROC curves in each dataset. Our proposed model, KEEP, achieves the best DICE and AUROC performance across all WSI datasets compared to other state-of-the-art models. **f.** Example WSIs from three datasets (the first two for CAMELYON16, the middle two for PANDA, and the last two for AGGC22) showing ground truth and predicted segmentation masks.

between different diseases, as shown in Figure 1b and Figure 6a.

Guided by the curated disease KG, we conduct thorough data cleaning (Figure 6b) and re-organize the noisy image-text data from OpenPath [23] and Quilt1M [25] into 143k semantic groups tied by well-defined hypernym relations. The statistics of semantic groups can be found in Figure 2a,b and Figure S1. Subsequently, the knowledge encoder is utilized to guide vision-language representation learning through a novel semantic-level alignment (Figure 1b and Figure 6c). More details of the KEEP model can be found in Methods. The results of zero-shot cancer region segmentation, cancer detection, cancer subtyping, and tile image profiling are exhibited in the following subsections.

2.2 Results on Zero-shot Cancer Region Segmentation

Segmenting cancerous regions from whole slide images (WSIs) to define the region of interest (ROI) is critical for subsequent morphological profiling in cancer diagnosis. Traditional approaches rely on labor-intensive, task-specific manual annotations to train slide-level segmentation models, a process that is both costly and time-consuming, limiting the scalability of computational pathology. In contrast, vision-language foundation models can perform zero-shot segmentation of cancerous regions by classifying image tiles into binary categories, enabling coarse-grained segmentation of WSIs, as shown in Figure 2c.

Following the common practice [34], we divide the tissue regions of each WSI into small tiles with 75% overlap, and average the classification scores in the overlapping areas to generate the final segmentation map. We compare the performance of KEEP with four other pathology vision-language models—PLIP [23], QuiltNet [25], MI-Zero [35], and CONCH [34]—on three datasets: CAMELYON16 [4] (48 WSIs), PANDA [8] (10,494 WSIs), and AGGC22 [24] (128 WSIs). We adopt the same approach as CONCH [34], ensembling 50 text prompts for each experiment and using the softmax function to normalize the similarities between the tile image and binary text prompts. Segmentation performance is evaluated using the area under the curve (AUC) and dice scores across all WSIs, as shown in Figure 2d,e. KEEP consistently outperforms the other models across all datasets. Notably, it achieves an average DICE score improvement of 0.068 and 0.081 over the state-of-the-art model CONCH on the CAMELYON16 and AGGC22 datasets, respectively.

To validate the enhancement of knowledge in cancer region segmentation, we compare the performance of KEEP with a simple contrastive approach. The experimental results, presented in Supplementary Figure S2 and Supplementary Table S8, demonstrate that KEEP significantly outperforms the simple contrastive method, achieving improvements of 0.13 and 0.10 on the PANDA and AGGC22 datasets, respectively. Example segmentation results are visualized in Figure 2f, with examples from CAMELYON16, PANDA, and AGGC22 in the left, middle, and right two columns, respectively. While KEEP effectively segments large cancerous regions, producing relatively coarse masks, it does exhibit some scattered false positives. This is likely due to the zero-shot binary classification approach, where each patch is processed independently, without incorporating contextual information.

2.3 Results on Zero-shot Cancer Detection

Traditional AI methods for identifying cancerous tissues in whole slide images (WSIs) typically use multiple instance learning (MIL) for weakly supervised classification, which requires frequent retraining across different cancer types or datasets. The pathology foundation model CHIEF [50] improves upon this by combining unsupervised pre-training for extracting tile-level features with weakly supervised training for recognizing WSI patterns, achieving excellent results in various tasks. However, these models often lack interpretability due to their need to amalgamate thousands of tile-level features into a single WSI classification. In contrast, our KEEP model employs a zero-shot, explainable approach using vision-language techniques that enhance both performance and interpretability in cancer detection.

For cancer detection in whole slide images (WSIs), we first perform zero-shot classification to identify cancerous tiles, rather than integrating tile features, we then compute the ratio of cancerous areas to the total tissue area, as shown in Figure 3a. The comparison between cancerous and normal WSIs is visualized in Figure 3b and Supplementary Figure S3a. Our results show that the predicted tumor area ratios for cancerous WSIs differ significantly ($P < 0.001$) from that of normal slides. This predicted ratio is then used as the probability of a WSI being cancerous.

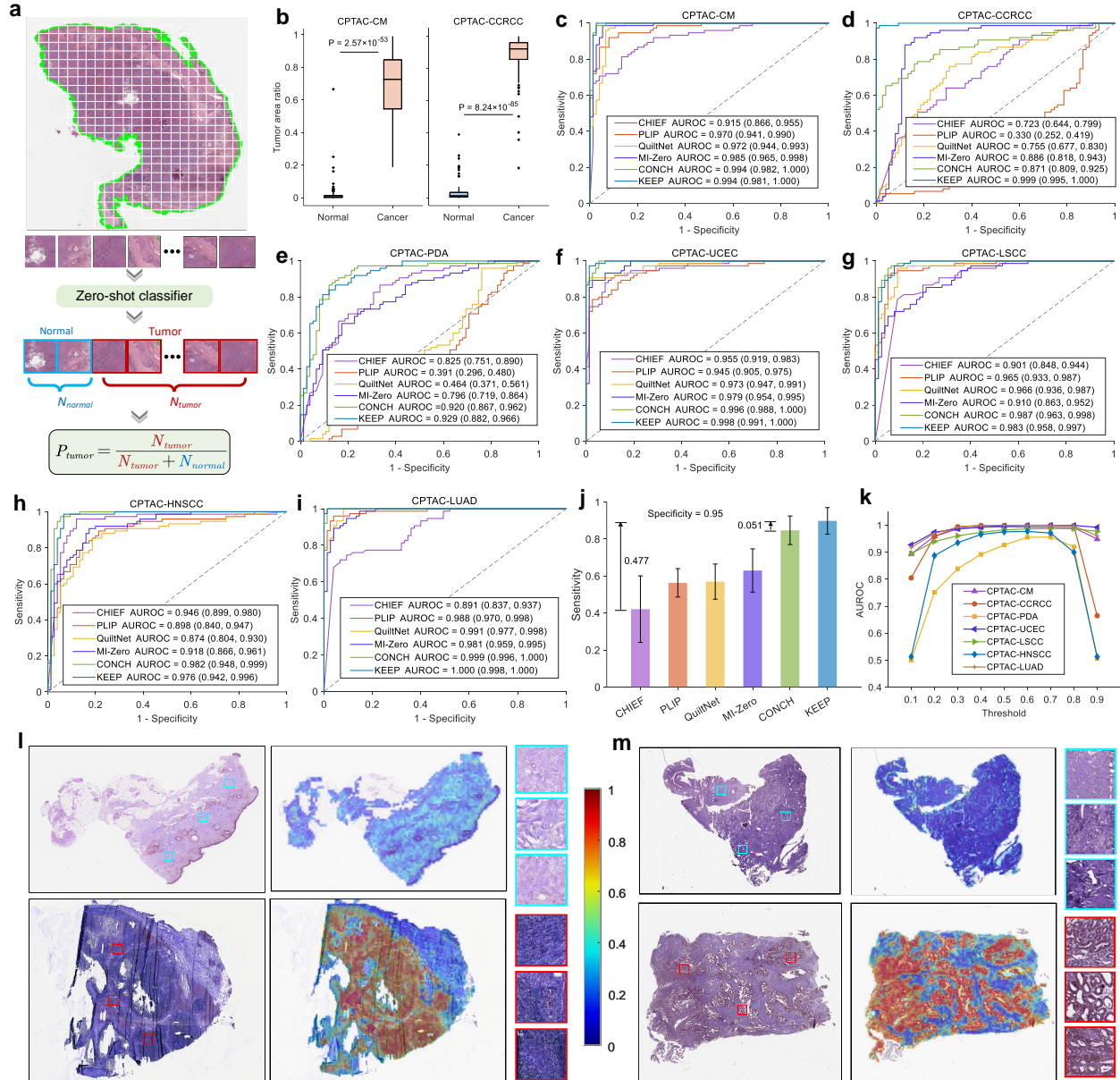


Figure 3 | Zero-shot cancer detection results. **a.** The zero-shot cancer detection scheme on WSIs, where individual tiles undergo binary classification. The probability of a slide being cancerous is determined by the predicted tumor ratio which is calculated by the ratio of tumor tiles to all valid tiles. **b.** The comparison of the predicted tumor ratio between normal and cancer WSIs in CPTAC-CM and CPTAC-CCRCC datasets. Two-sided Welch’s *t* test is used to assess the statistical significance of predicted tumor ratios among different WSIs. **c-i.** Comparison of ROC curves across different models, including CHIEF, PLIP, QuiltNet, MI-Zero, CONCH, and KEEP, evaluated on seven CPTAC datasets across six tissue anatomies: skin, kidney, pancreas, uterine, lung, and head and neck. Each dataset consists of 75 normal WSIs and 75 cancer slides, with each experiment using 1,000 bootstrap iterations. The AUROC for each model is reported as the median along with its 95% confidence intervals (CIs). **j.** Comparison of average sensitivities across all datasets at the specificity of 0.95, the error bar denotes the standard deviation of the performance. **k.** The robustness of our approach towards the threshold of the zero-shot classifier. **l,m.** Example visualizations of cancer detection on CPATC-CM, CPTAC-UCEC datasets. The first and the second rows denote the normal and the cancer WSIs. The heat map is generated by the similarities between the embeddings of tile images and that of "tumor" prompts.

We evaluate this approach on seven datasets from the Clinical Proteomic Tumor Analysis Consortium (CPTAC*), including CPTAC-CM (Cutaneous Melanoma), CPTAC-CCRCC (Clear Cell Renal Cell Carcinoma),

*<https://www.cancerimagingarchive.net/browse-collections/>

CPTAC-PDA (Pancreatic Ductal Adenocarcinoma), CPTAC-UCEC (Uterine Corpus Endometrial Carcinoma), CPTAC-LSCC (Lung Squamous Cell Carcinoma), CPTAC-HNSCC (Head and Neck Squamous Cell Carcinoma), CPTAC-LUAD (Lung Adenocarcinoma). For each dataset, we follow the same approach as CONCH [34], sampling 75 cancer slide images and 75 normal slide images, and ensemble 50 text prompts for each experiment, with results presented in Figure 3c-i.

Both vision-language models, CONCH and KEEP, significantly outperform CHIEF across all datasets, highlighting the utility of the tumor ratio as a promising **biomarker** for distinguishing cancerous from normal WSIs. In particular, KEEP and CONCH achieve a median AUROC of 0.994 on skin cancer (CPTAC-CM), outperforming CHIEF by 8 percentage points. For lung cancer detection (CPTAC-LUAD), KEEP and CONCH achieve near-perfect performance (median AUROC of 1.000), while CHIEF only reaches 0.891. This substantial improvement can be attributed to the fact that, unlike CHIEF, vision-language models like KEEP and CONCH integrate predicted labels rather than embedding tile features, enabling explicit identification of cancerous regions within WSIs.

Compared to other vision-language models, KEEP achieves the best performance on five out of seven datasets. Notably, KEEP attains an average sensitivity of 0.898 at a specificity of 95% across all datasets, as shown in Figure 3j. This is more than twice the sensitivity of CHIEF and 0.051 higher than that of CONCH, underscoring the substantial improvement in cancer detection performance enabled by KEEP. In particular, on the CPTAC-CM, CPTAC-CCRCC, CPTAC-UCEC, and CPTAC-LUAD datasets, KEEP consistently exceeds a sensitivity of 98%, as shown in Supplementary Figure S3b.

We also assess the robustness of KEEP to variations in the threshold of the zero-shot classifier. The threshold is varied from 0.1 to 0.9, and the results are shown in Figure 3k. Notably, KEEP maintains consistent performance across all tasks, with the exception of pancreatic cancer, when the threshold is in the range of [0.4, 0.7], indicating that KEEP is largely robust to threshold variations. In addition, to validate the enhancement of knowledge in cancer detection, we compare the performance of KEEP with a simple contrastive approach. The experimental results, presented in Supplementary Figure S3c and Supplementary Table S9, demonstrate that KEEP achieves comparable or better performance than simple contrastive learning on 6 out of 7 datasets. The similarities between predicted and ground-truth cancerous regions are visualized in heatmaps in Figure 3l,m and Supplementary Figure S4-S10, further demonstrating the consistency of the model’s predictions.

2.4 Results on Zero-shot Cancer Subtyping

Identifying tumor subtypes is essential for accurate cancer diagnosis and personalized treatment. Existing AI models, including multiple instance learning (MIL)-based approaches and the foundation model ProV-GigaPath [51], aggregate image features from individual tiles into a WSI-level representation for multi-class classification. However, these methods require a large number of labeled whole slide images for training, which limits their scalability to new cancer types.

Pathology vision-language models offer a promising zero-shot paradigm, where the predicted labels of individual tile images are aggregated, rather than their features, to determine the final classification. For example, MI-Zero [35] and CONCH [34] predict the subtype probability for each tile and then integrate the top-K predictions to classify the entire WSI. This zero-shot approach is highly adaptable and can be easily extended to new datasets.

To enhance this paradigm in a non-parametric manner, we adopt the same approach as cancer detection that uses the ratio of cancer subtype areas to total tissue areas as the subtype probability, as shown in Figure 4a. We evaluate the performance of different models on both common and rare cancer subtypes. The common cancer subtyping tasks include seven whole slide image (WSI) datasets from The Cancer Genome Atlas (TCGA[†]), CPTAC, and other resources: TCGA-BRCA (invasive breast carcinoma), TCGA-NSCLC (non-small-cell lung carcinoma), TCGA-RCC (renal cell carcinoma), TCGA-ESCA (esophagus carcinoma), TCGA-Brain (brain cancer), UBC-OCEAN [16, 3] (ovarian cancer), and CPTAC-NSCLC (non-small-cell lung carcinoma). Specifically, TCGA-BRCA consists of two subtypes: invasive ductal carcinoma (IDC) and invasive lobular carcinoma (ILC). TCGA-NSCLC contains lung adenocarcinoma (LUAD) and lung squamous cell

[†]<https://portal.gdc.cancer.gov/>

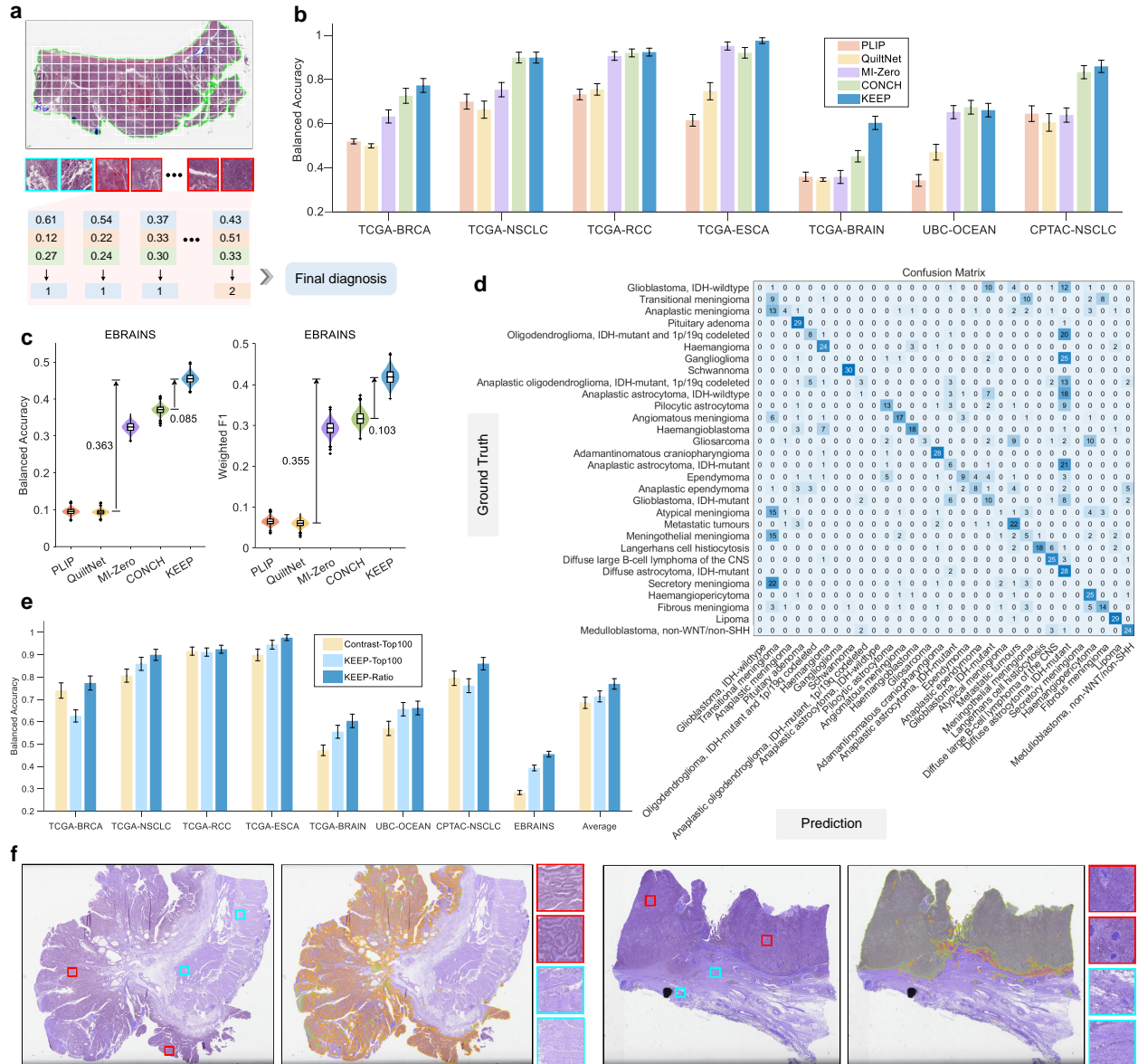


Figure 4 | Zero-shot cancer subtyping results. **a**. The zero-shot cancer subtyping scheme on WSIs, where individual tiles undergo multi-class classification, including a "normal" label and tumor subtype labels. The probability of a slide being classified as type I is determined by the ratio of type I tiles to all valid tiles. **b**. Comparison of balanced accuracy across different models on seven datasets with common cancer subtypes. The TCGA-BRCA, TCGA-NSCLC, TCGA-ESCA, and CPTAC-NSCLC datasets contain two subtypes, while the TCGA-RCC, TCGA-BRAIN, and UBC-OCEAN datasets consist of 3, 3, and 5 subtypes, respectively. Each subtype includes 75 WSIs, except for TCGA-ESCA (65 WSIs) and UBC-OCEAN (35 WSIs), with each experiment using 1,000 bootstrap iterations. **c**. Performance comparison of different models on the rare cancer subtyping dataset, EBRAINS, which consists of 30 rare brain cancer subtypes, each with 30 WSIs. **d**. The confusion matrix of the KEEP model on the rare brain cancer dataset, EBRAINS. **e**. Ablation results. Performance comparison between simple contrastive (Contrastive-Top100), KEEP with knowledge enhancement (KEEP-Top100), and KEEP with tumor-ratio strategy (KEEP-Ratio). Top100 suggests the strategy of top-100 pooling, while Ratio denotes the subtype ratio strategy. **f**. Example WSIs for tumor subtyping. The left and the right WSIs denote esophagus adenocarcinoma and esophagus squamous cell carcinoma, respectively. The orange and the green masks denote the predicted regions of adenocarcinoma and squamous cell carcinoma, respectively. The blue squares denote the tile image from the area with normal predictions.

carcinoma (LUSC). TCGA-RCC is divided into chromophobe renal cell carcinoma (CH RCC), clear-cell renal cell carcinoma (CCRCC), and papillary renal cell carcinoma (PRCC). TCGA-ESCA includes two subtypes:

squamous cell carcinoma and adenocarcinoma. TCGA-BRAIN consists of three subtypes: glioblastoma, astrocytoma, and oligodendroglioma. UBC-OCEAN contains five subtypes: ovarian clear cell carcinoma (CC), ovary endometrioid carcinoma (EC), high-grade ovary serous carcinoma (HGSC), low-grade ovary serous carcinoma (LGSC), ovarian mucinous carcinoma (MC). For the datasets collected from TCGA and CPTAC, we follow CONCH [34] to randomly sample 75 WSIs for each cancer subtype (65 for ESCA due to the number of original whole cases). For UBC-OCEAN, we randomly sample 35 WSIs for each subtype. The details of all the above datasets are listed in Supplementary Table S1.

Following the approach of CONCH, we ensemble 50 text prompts for each task and present the results in Figure 4b and Supplementary Figure S11a. Notably, KEEP outperforms other models on six out of seven datasets. For the brain cancer subtyping task, KEEP achieves an average balanced accuracy of 0.604, which is 0.15 higher than CONCH and 0.25 higher than the other models.

For rare cancer subtyping, we evaluate different models on the EBRAINS [39] dataset, which includes 128 subtypes of rare brain tumors. In accordance with CONCH, we select 30 rare subtypes, each containing more than 30 whole slide images (WSIs), for cancer subtyping evaluation. The results, presented in Figure 4c,d, show that KEEP achieves a median balanced accuracy of 0.456, more than four times higher than PLIP and QuiltNet, and 0.085 higher than CONCH. Figure 4d displays the confusion matrix for KEEP on the EBRAINS dataset, where it performs particularly well on subtypes such as schwannoma, adamantinomatous craniopharyngioma, and lipoma. However, several subtypes are misclassified as transitional meningioma and diffuse astrocytoma, IDH-mutant. This misclassification may be attributed to the limited representation of these brain cancer types in the training pathology image-text dataset.

To validate the enhancement of disease knowledge, we compare the performance of simple contrastive (termed Contrastive-Top100) and that of KEEP (KEEP-Top100), where Top100 suggests the strategy of top-100 tiles’ pooling, developed by MI-Zero and CONCH. The experimental results in Figure 4e, Supplementary Figure S11b and Supplementary Table S10 demonstrate that KEEP-Top100 achieves comparable or better performance than Contrastive-Top100 in 6 out of 8 datasets, with a significant improvement of 0.11 on the rare tumor dataset, EBRAINS, indicating that the disease knowledge can substantially enhance the zero-shot performance of rare tumor subtyping tasks.

We also compare the whole slide image (WSI) subtyping strategy based on top-100 pooling (KEEP-Top100) with our proposed subtype ratio strategy (KEEP-Ratio) across all datasets. The results, presented in Figure 4e and Supplementary Figure S11b, show that our subtype ratio strategy outperforms Top100 pooling on all datasets. Specifically, KEEP-ratio achieves an average balanced accuracy of 0.860 on CPTAC-NSCLC and 0.774 on TCGA-BRCA, surpassing KEEP-top100 by 0.10 and 0.15, respectively. Figure 4f and Supplementary Figure S12-S14 visualize the semantic segmentation of cancer subtypes, combining predicted labels with identified cancerous regions. These visualizations underscore the exceptional interpretability of our approach in cancer subtyping tasks.

2.5 Results on Zero-shot Pathology Tile Image Profiling

In this section, we conduct an evaluation of KEEP on tile-level tasks, including cross-modal retrieval and zero-shot tile image classification.

Cross-modal Retrieval. We evaluate on four pathology image-text datasets: ARCH-PubMed [18], ARCH-Book [18], Pathpair [54], and WebPath, a dataset collected from professional websites. As MI-Zero and CONCH do not release their image-text training data, to avoid unfair comparison from data leakage, we limit our comparison to PLIP and QuiltNet. The experimental results, presented in Figure 5a, show that KEEP outperforms PLIP and QuiltNet by a significant margin across all datasets, in both text-to-image and image-to-text retrieval tasks. This improvement is attributed to the superior initialization of the visual encoder and the knowledge-enhanced text encoder.

Tile Image Classification. For the tile image classification task, we compile 14 datasets covering seven human tissue types: breast (BACH [1] and Breakhis [44]), colon (NCT-CRC-HE-100K [28], CRC100K [29], and Chaoyang [56]), lung (LC25000 [6] and WSSS4LUAD [20]), kidney (RenalCell [7]), bone (Osteo [2]), skin (SkinCancer [30]), and esophagus (ESCA-UKK, ESCA-WNS, ESCA-TCGA, and ESCA-CHA [46]). The number of classes per dataset ranges from 2 to 16. Full dataset details are provided in the Methods section

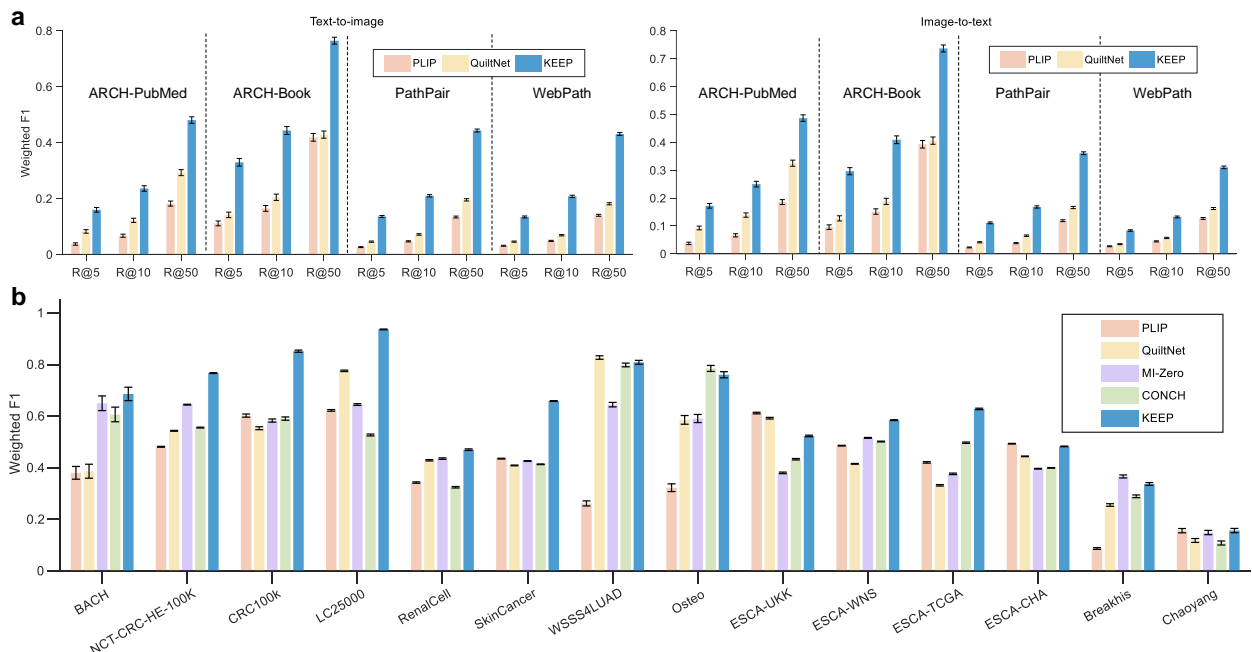


Figure 5 | Zero-shot tile image profiling results. **a.** Performance comparison of different models on the cross-modal retrieval task. R@K ($k = 5, 10, \text{ and } 50$) denotes Recall@K, the ratio of correctly retrieved queries in Top-K retrieved samples. **b.** Performance comparison of different models on the zero-shot tile image classification task. The error bar denotes the standard deviation of the results from 1000 bootstrapping iterations.

and Supplementary Table S1.

In accordance with PLIP, we concatenate a template phrase with the class names (*a histopathology image of {class name}*) in each dataset to construct the text prompts for zero-shot tile classification. The bootstrapping performance of different models is shown in Figure 5b. Notably, KEEP achieves the best performance on 9 out of 14 datasets. Specifically, for the CRC100K, LC25000, and SkinCancer datasets, KEEP shows an improvement of at least 16 percentage points compared to other models. We also conduct ablation experiments of model components on the task of tile image profiling. The experimental results, exhibited in Supplementary Figure S15, demonstrate that knowledge enhancement can improve the performance of both cross-modal retrieval and zero-shot tile image classification tasks.

3 Discussion

In this study, we present KEEP, a novel vision-language foundation model specifically designed to tackle challenges in computational pathology. By incorporating disease-specific knowledge, KEEP achieves state-of-the-art performance in zero-shot cancer diagnosis. In particular, for cancer detection, KEEP significantly outperforms CHIEF, achieving a notable improvement in sensitivity (0.898) at a specificity of 0.95 across multiple cancer types. Similarly, in cancer subtyping, KEEP outperforms CONCH by integrating disease-specific knowledge, which improves the alignment between pathology images and subtype semantics. Notably, in rare cancer subtyping tasks, KEEP achieves a balanced accuracy improvement of 8.5 percentage points over CONCH.

The superior performance of KEEP compared to other models can be primarily attributed to two key factors: the injection of disease knowledge during training and the tumor-ratio-based prediction method employed in downstream tasks. **(i) integration of disease knowledge during training:** disease knowledge is incorporated through three key mechanisms. First, the language model is employed to encode the disease knowledge graph (KG), aligning the representation space of disease names, definitions, synonyms, and hierarchical relationships. This alignment serves as a bridge, implicitly linking pathology images with their corresponding disease types during the vision-language pre-training process. Second, the extensive use of disease synonyms within the

disease KG explicitly highlights critical disease-related information in the textual descriptions of images. This transformation converts weak supervision signals from free-text annotations into stronger disease-level supervision, thereby reinforcing the connection between pathology images and disease entities. Finally, the hierarchical structure of the disease KG organizes image-text pairs into semantic groups with hypernym-hyponym relationships, significantly enhancing alignment accuracy. **(ii)** tumor-ratio-based prediction in downstream tasks: for both cancer detection and subtyping, the tumor-ratio-based prediction method proves to be an intuitive and highly effective biomarker. Unlike methods that aggregate tile-level features, the tumor-ratio-based approach leverages tumor region localization to deliver superior interpretability, a critical requirement for clinical applications. Moreover, compared to approaches that predict tumor subtypes by selecting top-k tumor patches, the tumor-ratio-based method is non-parametric, offering a more straightforward and transparent classification process. This methodology closely mirrors the diagnostic reasoning of human pathologists, combining clinical interpretability with robust performance, ultimately driving substantial improvements across tasks.

While promising, this study also faces several limitations: **(i)** Despite the encouraging zero-shot performance in cancer region segmentation, the DICE score is still hindered by scattered false positives. This limitation stems from the zero-shot nature of vision-language models, where each tile is classified independently, without accounting for the contextual relationships between neighboring regions. Incorporating such contextual information—such as spatial relationships between tiles—could help reduce false positives and improve segmentation accuracy. **(ii)** Although KEEP demonstrates strong performance in rare cancer subtyping, its predictions for certain subtypes remain limited due to the scarcity of these cases in the image-text training data. In these instances, few-shot learning, where at least one whole slide image (WSI) per subtype is available, could enhance model performance by capturing the diversity within rare cancer subtypes. **(iii)** While pathology vision-language models exhibit robust zero-shot ability, their performance is often dependent on prompt engineering, which can limit their adaptability and robustness across diverse datasets. A promising avenue for improvement is prompt learning, wherein the model learns a trainable prompt from a small set of example WSIs, replacing manually designed prompts. This approach would enable the model to adapt more effectively to varied datasets and tasks, enhancing its generalizability and robustness. Another promising avenue for improvement is the multi-modal alignment that introduces genomic [52] or epigenomic [21] information.

In conclusion, our results demonstrate that KEEP offers a powerful tool for cancer diagnosis by injecting domain-specific knowledge into vision-language models. KEEP shows great promise in advancing computational pathology and has the potential to make a significant impact in clinical settings, offering improved accuracy and interpretability in cancer diagnosis.

4 Methods

4.1 Disease Knowledge Graph Construction

We construct a comprehensive disease knowledge graph by integrating data from publicly available databases: the Disease Ontology (DO) [40] and the Unified Medical Language System (UMLS) [5]. The resulting knowledge graph contains 11,454 disease entities and 139,143 associated attributes, including 14,303 definitions, 15,938 hypernym relationships, and 108,902 synonyms.

Existing Medical knowledge Databases. Disease Ontology (DO) was developed to standardize disease nomenclature and classification, DO offers: **(i)** disease classification, that encompasses categories such as infectious diseases, genetic disorders, cancers, and metabolic conditions; **(ii)** disease relationships, that links the diseases as subtypes to their parent categories with a hierarchical structure; **(iii)** database mapping, that extensively maps its terms to other medical vocabularies in Unified Medical Language System (UMLS). UMLS is a comprehensive medical language system crafted by the U.S. National Library of Medicine, that integrates diverse medical terminologies from sources like ICD, SNOMED, and MeSH. It includes: **(i)** metathesaurus, a compilation of medical concepts that encompass diseases, symptoms, treatments, and more; **(ii)** semantic network, a framework that captures semantic relationships among medical concepts, detailing connections between diseases and their symptoms, treatments, and causes.

Constructing the Disease Knowledge Graph. Starting with disease entities and their attributes

from DO, we enriched each entity with additional attributes from UMLS through cross-mapping. We then established hypernym relationships from DO as the edges connecting these entities. In total, the knowledge graph encompasses 11,454 disease entities and 139,143 disease attributes, including 14,303 definitions, 15,938 hypernym relationships, and 108,902 synonyms. Additionally, for each disease entity, we construct a hierarchical disease chain by linking each disease entity to the root through a random hypernym relation path, as shown in Figure 6a. In the disease chain, the name of each disease entity is randomly chosen from its set of synonyms, further enriching the diversity of the disease representations. This process creates a multi-level hierarchy, where each disease entity is connected to its upper-level relations. The resulting hierarchical structure provides a more nuanced understanding of disease relationships, enhancing the contextualization of each entity and enabling better integration of related knowledge for downstream tasks.

4.2 Problem Formulation

Given a set of image-text pairs, denoted by $\mathcal{F} = \{(x_1, c_1), \dots, (x_n, c_n)\}$, conventional vision-language pre-training approaches typically employ simple contrastive learning, that aims to align the embeddings of paired images and texts, while simultaneously separating those of unpaired samples:

$$\text{sim}(\Phi_v(x_i), \Phi_t(c_i)) \gg \text{sim}(\Phi_v(x_i), \Phi_t(c_j)), \quad i \neq j, \quad (1)$$

where Φ_v and Φ_t denote the visual and the text encoder, respectively. Such training procedure suffers from two issues, **(i)** there is no explicit knowledge injection to the training procedure, for example, the intricate relationships between diseases; **(ii)** the noise in existing datasets, for instance, low-quality captions and non-pathology images, further introduces ambiguities and inconsistencies in the alignment process.

Herein, we leverage the constructed disease knowledge graph to enhance the vision-language training procedure, from three key aspects: a strong text encoder via knowledge representation learning, knowledge-guided semantic group construction, and knowledge-guided vision-language representation learning. The target of knowledge-enhanced vision-language training can be formulated as:

$$\min_p \text{sim}(\Phi_v(\tilde{x}_p), \Phi_k(\tilde{c}_i)) \gg \max_q \text{sim}(\Phi_v(\tilde{x}_q), \Phi_k(\tilde{c}_j)), \quad i \neq j, \quad (2)$$

where Φ_k denotes a BERT-based knowledge encoder, pre-trained on our constructed disease knowledge graph (KG). $(\tilde{x}_1, \dots, \tilde{x}_n, \tilde{c}_i)$ represents the i -th semantic group consisting of a single refined caption \tilde{c}_i with a varying number of pathology images $(\tilde{x}_1, \dots, \tilde{x}_n)$. The specific components for knowledge enhancement are detailed in the following subsections.

4.3 Disease Knowledge Encoding

In this section, we describe the procedure for training a language model to construct a knowledge embedding space, in which attributes of the same disease are mapped to similar embeddings. To account for hierarchical relationships between disease entities, we follow the approach outlined in BioCLIP [45], recursively concatenating each disease entity with its hypernyms (parent nodes) to form hierarchical disease chains, as illustrated in Figure 6a. This enables to capture the relationships between diseases at different levels of granularity. We then apply metric learning to align the embeddings of disease attributes, ensuring that entities with similar properties are closer in the embedding space.

Specifically, given a set of disease entities with hypernym relations, we randomly sample parent nodes to construct hierarchical disease chains for each entity. As a result, each disease entity is associated with a varying number of attributes, including synonyms, disease chains, and definitions, which can be denoted by $\mathcal{D} = \{(d_1, \mathbf{a}_1), \dots, (d_n, \mathbf{a}_n)\}$, where d_i denotes the i -th disease entity, and $\mathbf{a}_i = \{a_i^1, \dots, a_i^k\}$ refer to the associated k attributes, both disease and attributes are in the form of natural language. Note that, for different disease entities, k also varies, our goal here is to train a model that satisfies the following condition:

$$\text{sim}(\Phi_k(a_i^p), \Phi_k(a_i^q)) \gg \text{sim}(\Phi_k(a_i^p), \Phi_k(a_j^t)), \quad i \neq j, \quad (3)$$

$$\text{sim}(\Phi_k(a_i^p), \Phi_k(a_i^q)) = \langle \Phi_k(a_i^p), \Phi_k(a_i^q) \rangle \quad (4)$$

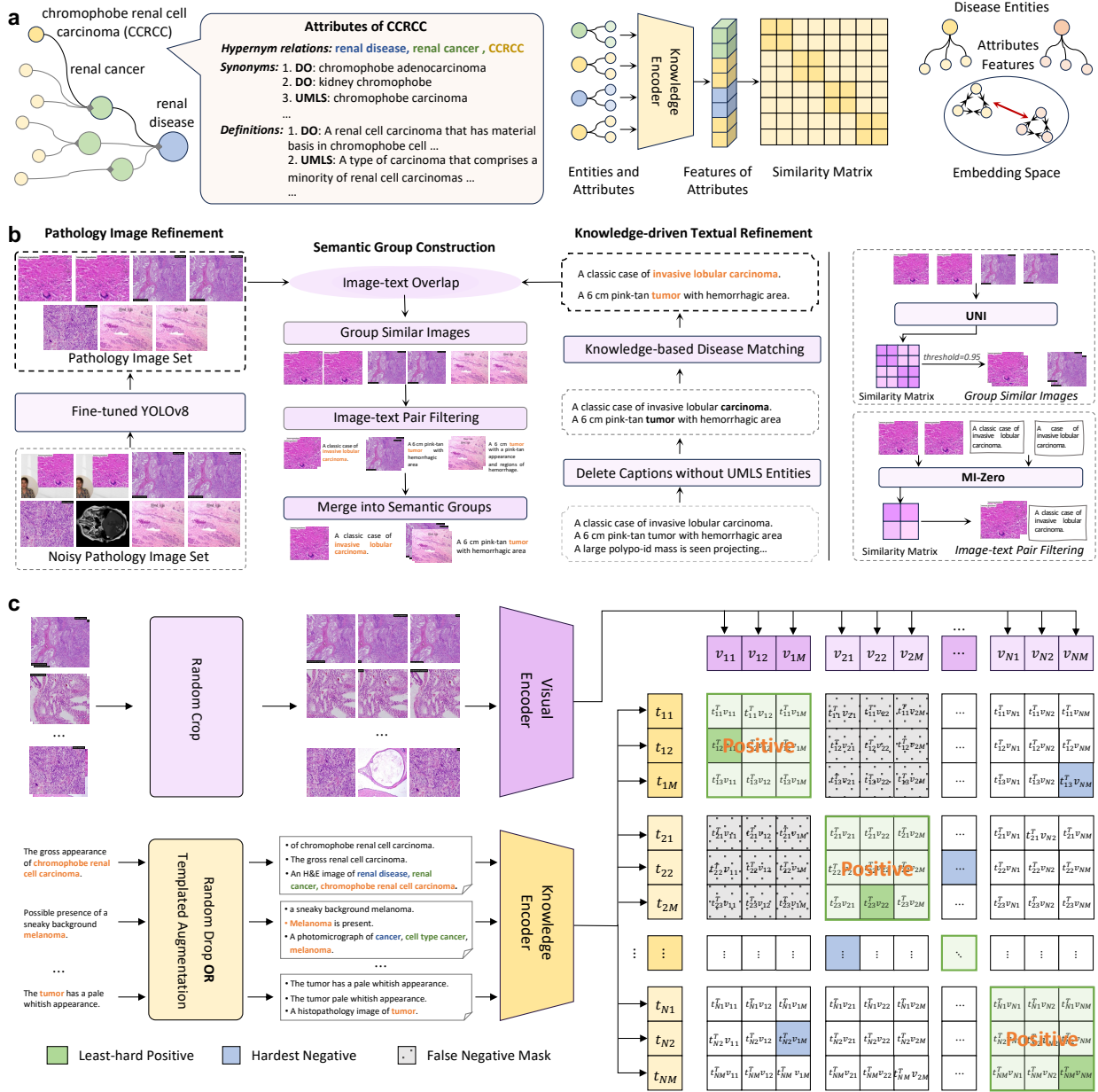


Figure 6 | Architecture of KEEP. **a.** Disease knowledge encoding. We establish a knowledge graph that includes hypernym relations, synonyms, and definitions of diseases, and pre-trained a disease knowledge encoder. Diseases at different levels are represented by different colors. **b.** Knowledge-guided dataset structuring. We fine-tune YOLOv8 to remove noise in the pathology image dataset, extract medical entities from the captions, align the diseases in the captions with the diseases and synonyms in the knowledge graph, and cluster the filtered image and text data into semantic groups. The right side illustrates two specific methods used during the clustering process. **c.** Knowledge-enhanced vision-language pre-training. We perform cropping and random dropping augmentations on the images and texts, and paraphrase captions that contain diseases using templates. During the training process, to mitigate the impact of false negatives, we design strategies for positive mining, hardest negative, and false negative elimination.

where $\Phi_k(\cdot)$ denotes the knowledge encoder, $\langle \cdot \rangle$ refers to the cosine similarity, a_i^p, a_i^q and a_j^t refer to the randomly sampled attributes from the i, j -th disease entity. Intuitively, the knowledge encoder is optimised to pull together the attributes of the same disease, while pushing apart attributes of different diseases.

Metric Loss. At training time, we employ metric learning to construct an embedding space where the intra-class instances are clustered, and inter-class instances are separated. Specifically, given a mini-batch with n randomly selected diseases, each associated with k attributes, we denote the normalized embedding for the p -th attribute of the i -th disease as:

$$\mathbf{z}_p^i = \frac{\Phi_k(a_i^p)}{\|\Phi_k(a_i^p)\|}, \quad (5)$$

where $\Phi_k(a_i^p)$ represents the embedding of the p -th attribute of the i -th disease, and $\|\cdot\|$ denotes the L2-norm. We adopt the recently proposed AdaSP loss [55], which finds out a max-min positive similarity and then shapes a loss with the maximal negative similarity:

$$\mathcal{L}_{\text{metric}} = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp((S_i^- - S_i^+)/\tau)), \quad (6)$$

where τ is a temperature parameter. S_i^+ and S_i^- denote the max-min positive and the maximal negative similarity, which can be computed by the soft version:

$$S_i^+ = \max_p \min_q \langle \mathbf{z}_p^i, \mathbf{z}_q^i \rangle \approx \tau \log \left(\frac{1}{\sum_{p=1}^k \sum_{q=1}^k \exp(-\langle \mathbf{z}_p^i, \mathbf{z}_q^i \rangle / \tau)} \right), \quad (7)$$

$$S_i^- = \max_{j,p,q} \langle \mathbf{z}_p^i, \mathbf{z}_q^j \rangle \approx \tau \log \left(\sum_{p=1}^k \sum_{j=1, j \neq i}^n \sum_{q=1}^k \exp(\langle \mathbf{z}_p^i, \mathbf{z}_q^j \rangle / \tau) \right). \quad (8)$$

4.4 Knowledge-guided Dataset Structuring

In this section, drawing upon the constructed knowledge graph, we propose an automated pipeline for cleaning, and re-organizing the public noisy image-text pairs.

Pathology Image Curation. Quilt-1M [25] and OpenPath [23] are publicly available pathology image-language datasets, curated by sourcing images from online platforms, for example, educational videos on YouTube [‡], or pathology images from Twitter [§]. Due to the diverse data sources, these datasets inevitably contain a high noise ratio, for example, radiological images and pathological images may co-exist in the same slide. We therefore train a detector to crop the pathological part from each of the images. Specifically, we manually annotate 1,000 images and then fine-tune a well-established detection model, YOLOv8 [¶] [48], on this annotated dataset. The refined YOLOv8 model is applied to scan the entire dataset, eliminating samples that fail to detect pathological images, while preserving those with clear pathological detections. Through random sampling and manual verification, we confirm that 99.9% of the reserved samples consist of pure pathological images.

Knowledge-driven Textual Refinement. In open-source datasets, many textual captions are derived from instructional videos, where the spoken language is transcribed directly into texts. As a result, these captions often include substantial amounts of irrelevant information that do not pertain to pathological images. We employ the natural language processing (NLP) tool, SpaCy, to extract named entities from the captions and align them with corresponding entities in the Unified Medical Language System (UMLS) [5]. Sentences that do not contain any UMLS entities are excluded from further processing.

Additionally, as the captions in these datasets are mainly collected from platforms such as Twitter and YouTube, which are typically unstructured and lack explicit disease labels. We perform fuzzy matching between the captions and the synonyms of all disease entities in the pre-constructed disease KG, to identify explicit disease labels for each caption.

[‡]<https://www.youtube.com/>

[§]<https://x.com/>

[¶]<https://docs.ultralytics.com/>

Semantic Group Construction. With the refined images and texts, we first employ the pathology visual encoder UNI [13] to compute image embeddings and group images with high embedding similarity (*e.g.*, consecutive frames in video sequences) using a threshold of 0.95. This step ensures that images of high visual similarity are clustered together into semantic groups. Subsequently, within each group, we employ a pre-trained visual-language model MI-Zero [35] to compute embeddings for both images and texts. We then compute a similarity matrix, and only the caption exhibiting the highest similarity to all images within the group is retained, ensuring the most representative description is selected. Finally, to resolve potential duplicates or redundant captions across groups, we calculate the Intersection over Union (IoU) of their token sets. Groups with an IoU greater than 0.9 are merged, resulting in the formation of the final refined semantic groups.

4.5 Knowledge-enhanced Vision-language Training

In this section, we present our knowledge-enhanced vision-language training framework by introducing semantic-level alignment via metric learning within well-structured semantic groups. The text encoder, pre-trained on the disease KG, ensures that disease attributes are closely aligned in the embedding space.

Semantic Alignment. Given a set of well-structured semantic groups, each with a caption and a varying number of images, we prepare a mini-batch for each iteration in the following steps: (i) we randomly sample N semantic groups, and for each group (denoted by $\{c_i | i \in [1, \dots, N]\}$), we randomly sample M out of G paired images (denoted by $\{x_{ik} | k \in [1, \dots, M]\}$) with replacement, which yields a positive semantic group; (ii) we take random crops for each sampled image, and resize the image from 512×512 pixels to 224×224 pixels; (iii) the caption of each group is augmented M times, with a 50% probability to randomly drop 40% of the words each time. Further, we randomly sample a template from the template set (in Supplementary Table S11) to rephrase the matched captions, *i.e.*, [Template] + disease label/hierarchical disease chain, for instance, *a histopathology image of skin squamous cell carcinoma/ skin disease, skin cancer, skin carcinoma, skin squamous cell carcinoma*. And we randomly sample one of the dropped captions and the paraphrased caption as the semantics for each tile image.

We denote images and augmented captions in i -th semantic group as $(x_i^1, \dots, x_i^n, c_i^1, \dots, c_i^n)$. Correspondingly, the normalized visual embedding of the k -th image in i -th semantic group as $\mathbf{v}_{ik} = \Phi_v(x_i^k)$, where Φ_v suggests the visual encoder. Similarly, the normalized embedding of m -th caption in i -th positive group can be denoted by $\mathbf{t}_{im} = \Phi_k(c_i^m)$, where Φ_k represents the well-trained knowledge encoder. The similarities between positives within each semantic group, correspondingly, can be computed by $\{\mathbf{t}_{im}^T \mathbf{v}_{ik} | m, k \in [1, \dots, M]\}$, which composes a matrix with the size of $M \times M$, marked by green boxes in Figure 6c. The cosine similarities between negatives, *e.g.*, the i -th and the j -th semantic groups, can be computed by $\{\mathbf{t}_{jm}^T \mathbf{v}_{ik} | m, k \in [1, \dots, M], j \neq i\}$, which composes a matrix with the size of $M \times M(N - 1)$. Overall, the similarity scores for a mini-batch shape a $NM \times NM$ matrix with the diagonal blocks suggesting positive semantic groups, shown in Figure 6c.

With this similarity matrix, we aim to exploit a metric loss to increase the similarity scores in positive groups, while decreasing scores between negatives. The metric loss can be formulated by:

$$\mathcal{L}_{\text{metric}} = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp((S_i^- - S_i^+)/\tau)), \quad (9)$$

where S_i^+ and S_i^- represent the positive and the negative similarities for the i -th group, and τ denotes the temperature parameter.

Positive Mining. There are $M \times M$ positive pairs between augmented images and captions for each semantic group. Their positive embedding similarities can be denoted as $\{\mathbf{t}_{im}^T \mathbf{v}_{ik} | i \in [1, \dots, N], m, k \in [1, \dots, M]\}$, where \mathbf{t}_{im} and \mathbf{v}_{ik} suggest the normalized embeddings of the m -th caption and the k -th image in the i -th semantic group. The similarity between the k -th image in group i and its hardest caption can be computed by the smooth approximation of the minimum function:

$$S_{ik}^+ = \min_m(\mathbf{t}_{im}^T \mathbf{v}_{ik}) \approx -\tau \log \left(\sum_{m=1}^M \exp(-\mathbf{t}_{im}^T \mathbf{v}_{ik}/\tau) \right). \quad (10)$$

Intuitively, mining hard positives can accelerate the training procedure and promote the performance of image-text alignment, while the hardest positives, to a large extent, could be false positives due to data noise. As a result, we choose a moderate hard positive by mining the easiest one among the hard positive candidates in Eq. 10:

$$S_i^+ = \max_k (S_{ik}^+) \approx \tau \log \left(\sum_{k=1}^M \exp(S_{ik}^+/\tau) \right). \quad (11)$$

Substituting Eq. 9 into Eq. 10, we have:

$$S_i^+ = \max_k \min_m (\mathbf{t}_{im}^T \mathbf{v}_{ik}) \approx \tau \log \left(\sum_{k=1}^M \frac{1}{\sum_{m=1}^M \exp(-\mathbf{t}_{im}^T \mathbf{v}_{ik}/\tau)} \right). \quad (12)$$

In contrast to mining the hardest positive pair, the least-hard positive is also a moderate choice that can not only reduce the risk of false positives caused by outliers but also prevent the training procedure from degenerating into trivial positives. It is noteworthy that although the mined S_i^+ is the easiest one among hard positives, it is still a hard positive that achieves an elaborated balance between introducing implicit false positives and collapsing to trivial negatives.

False Negative Elimination. To prevent semantic groups with the same disease label or hypernym relations from being misclassified as negatives, we check if one group is reachable from another via their hypernym paths, setting their negative indicator to zero when applicable. In contrast to positive mining, we perform the hardest mining for negative samples since we have eliminated the false negatives above. Correspondingly, the negative similarity for the i -th semantic group can be formulated as:

$$S_i^- = \max_k \max_{j,m} (\mathbf{t}_{jm}^T \mathbf{v}_{ik}) \approx \tau \log \left(\sum_{k=1}^M \mathcal{I}_{ij} \left(\sum_{j=1, j \neq i}^N \sum_{m=1}^M \exp(\mathbf{t}_{jm}^T \mathbf{v}_{ik}/\tau) \right) \right). \quad (13)$$

where \mathbf{t}_{jm} and \mathbf{v}_{ik} suggest the normalized embeddings of the m -th caption in the j -th group and that of the k -th image in the i -th group, respectively. \mathcal{I}_{ij} is a binary indicator that denotes the negative flag between the i -th and the j -th semantic group.

4.6 Model Training Details

Knowledge Encoding. We adopt the architecture of PubMedBERT [19] to conduct knowledge encoding. The embedding dimension is set to 768. The temperature parameter τ is set to 0.04 in Eq. 6. The batch size is set to 256, including 32 disease entities with 8 instances per entity. We train the knowledge encoder for 100 epochs with a maximum learning rate of 3×10^{-5} on 4 A100 GPUs.

Vision-language Pre-training. KEEP consists of a vision encoder based on the backbone of ViT-L-16 and a text encoder based on the architecture of PubMedBERT. We adopt UNI [13] to initialize the image encoder of KEEP and set the size of the input image as 224×224 pixels. Meanwhile, the text encoder of KEEP is initialized by the pre-trained knowledge encoder. The batch size is set to 128, including 32 semantic groups with 4 image-text pairs per group. The temperature in Eq. 9 is set to 0.04 across all experiments. We conduct the vision-language pre-training for 10 epochs with a maximum learning rate of 1×10^{-5} on 1 A100 GPU.

4.7 Zero-shot Evaluation on WSIs

We first evaluate KEEP on whole slide images to segment cancerous cells, detect malignant tumors, and predict subtypes of cancers. For WSI preprocessing, we follow CONCH [34] to divide each whole slide image into 256×256 tiles (224×224 pixels for segmentation) at $20\times$ magnification, which are then fed to KEEP to predict the class label of each tile in a zero-shot manner.

Unsupervised Prompt Screening. We follow CONCH [34] to randomly concatenate one of the 21 templates and cancer synonyms to generate prompt classifiers. As the zero-shot performance can be sensitive to text prompts, we develop an unsupervised prompt screening approach to improve the robustness of performance

on downstream tasks.

Specifically, given a set of prompt classifiers that aim to categorize N tile images into C types, we have C similarities for each tile image by calculating similarity scores between the tile image and the different text prompts in one prompt classifier. Ideally, different categories in one prompt classifier should exhibit consistency in similarity ranges and complementary relationships. For example, in cancer detection tasks, if a tile image has a similarity score of 0.7 with the prompt, *A histopathology image of cancerous tissue*, the corresponding normal prompt (*A histopathology image of normal tissue*) is expected to have a similarity score of 0.3 with the same tile image. Additionally, a larger gap between the similarity scores of different classes indicates better discriminative power. To evaluate and refine the prompt classifiers, we introduce a screening score that ranks each prompt classifier based on its range consistency and discriminability.

$$R_s = \sum_{i=1}^N S_i^* - S_i^{**} - |S_i^* + S_i^{**} - 1| \quad (14)$$

where S_i^* , S_i^{**} denotes the largest and the second-largest of the similarities between the i -th tile image and the category prompts. $S_i^* - S_i^{**}$ measures the discriminability of the prompt classifier, while $|S_i^* + S_i^{**} - 1|$ (smaller is better) suggests the range consistency. In this paper, we use this score to screen the top-50 prompt classifiers to evaluate the zero-shot performance on downstream tasks.

4.8 Evaluation Metrics

Zero-shot Performance on WSI Tasks. We adopt the same metric as MI-Zero [35] and CONCH, namely, balanced accuracy and weighted F1 to measure the cancer subtyping performance on WSIs. We follow CONCH to use the nonparametric bootstrapping with 1000 samples to construct 95% confidence intervals for model performance. For cancer segmentation tasks, we exploit area under curves (AUC) and DICE to evaluate different models. For cancer detection, we use the same metric AUC, sensitivity, and specificity as CHIEF [50] to evaluate the performance of different models.

Zero-shot Performance on Tile-level Tasks. For classification tasks, we adopt the same metric as PLIP [23], namely, weighted F1 (wF1). For cross-modal retrieval tasks, we adopt the metric of Recall@K, suggesting the ratio of correctly retrieved queries in Top-K retrieved samples. We also use the nonparametric bootstrapping with 1000 samples to construct 95% confidence intervals for tile-level tasks.

5 Data Availability

The disease Knowledge, including disease entities and hypernym relations are available in Disease Ontology (DO) (<https://disease-ontology.org/do/>). Disease synonyms and definitions are available in Unified Medical Language System (UMLS) (<http://www.nlm.nih.gov/research/umls/licensedcontent/umlsknowledgesources.html>). The pathology image-text pairs for alignment are available in OpenPath (<https://huggingface.co/vinid/plip>) and Quilt1M (<https://github.com/wisdomikezogwo/quilt1m>). Test datasets for cancer region segmentation are available in CAMELYON16 (<https://camelyon16.grand-challenge.org/>), PANDA (<https://panda.grand-challenge.org/data/>), and AGGC22 (<https://aggc22.grand-challenge.org/>). For cancer detection, test datasets of CPTAC-CM, CPTAC-CCRCC, CPTAC-PDA, CPTAC-UCEC, CPTAC-LSCC, CPTAC-HNSCC, CPTAC-LUAD are available in CPTAC (<https://proteomics.cancer.gov/programs/cptac>). For cancer subtyping, test datasets of TCGA-BRCA, TCGA-NSCLC, TCGA-RCC, TCGA-ESCA, TCGA-BRAIN are available in TCGA (<https://portal.gdc.cancer.gov/>). Other datasets for cancer subtyping are available in CPTAC-NSCLC (<https://proteomics.cancer.gov/programs/cptac>), EBRAINS (<https://data-proxy.ebrains.eu/datasets/>), and UBC-OCEAN (<https://www.kaggle.com/competitions/UBC-OCEAN/>). The sources of all tile datasets are listed in Supplementary Table S2.

6 Code Availability

The source codes for KEEP are available at <https://github.com/MAGIC-AI4Med/KEEP>.

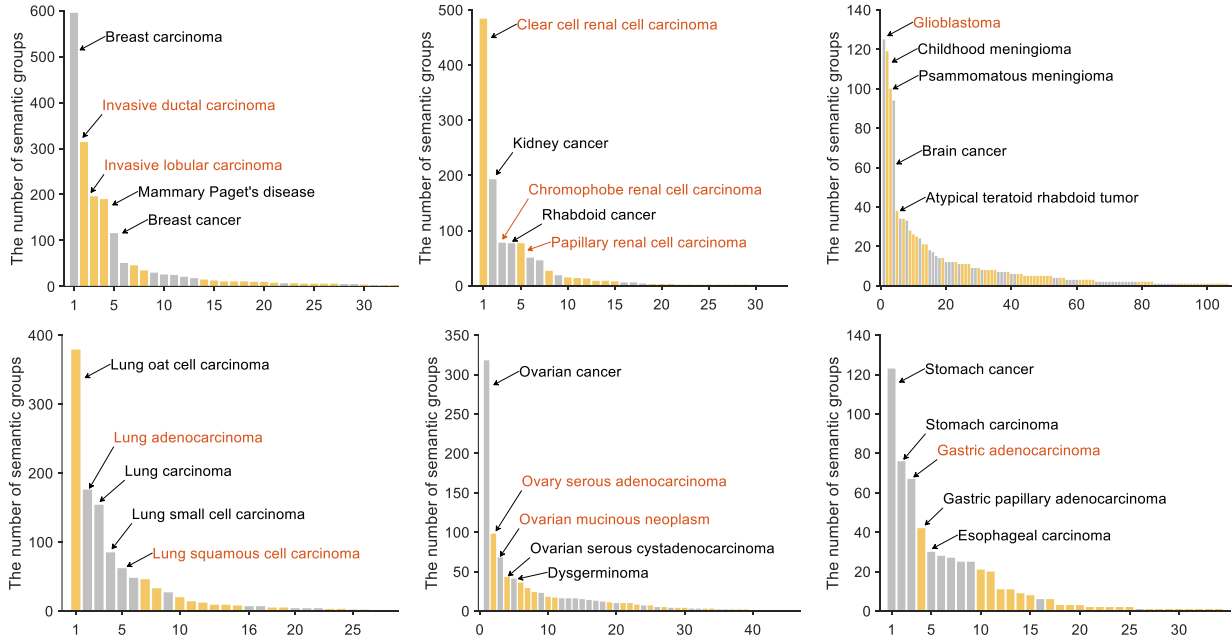
References

- [1] Guilherme Aresta, Teresa Araújo, Scotty Kwok, Sai Saketh Chennamsetty, Mohammed Safwan, Varghese Alex, Bahram Marami, Marcel Prastawa, Monica Chan, Michael Donovan, et al. Bach: Grand challenge on breast cancer histology images. *Medical Image Analysis*, 56:122–139, 2019.
- [2] Harish Babu Arunachalam, Rashika Mishra, Ovidiu Daescu, Kevin Cederberg, Dinesh Rakheja, Anita Sengupta, David Leonard, Rami Hallac, and Patrick Leavey. Viable and necrotic tumor assessment from whole slide images of osteosarcoma using machine-learning and deep-learning models. *PLoS One*, 14(4):e0210706, 2019.
- [3] Maryam Asadi-Aghbolaghi, Hossein Farahani, Allen Zhang, Ardalan Akbari, Sirim Kim, Ashley Chow, Sohier Dane, OCEAN Challenge Consortium, OTTA Consortium, David G Huntsman, et al. Machine learning-driven histotype diagnosis of ovarian carcinoma: Insights from the ocean ai challenge. *medRxiv*, pages 2024–04, 2024.
- [4] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson, Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA*, 318(22):2199–2210, 2017.
- [5] Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic Acids Research*, 32(suppl_1):D267–D270, 2004.
- [6] Andrew A Borkowski, Marilyn M Bui, L Brannon Thomas, Catherine P Wilson, Lauren A DeLand, and Stephen M Mastorides. Lung and colon cancer histopathological image dataset (lc25000). *arXiv preprint arXiv:1912.12142*, 2019.
- [7] Otso Brummer, Petri Pölönen, Satu Mustjoki, and Oscar Brück. Integrative analysis of histological textures and lymphocyte infiltration in renal cell carcinoma using deep learning. *bioRxiv*, pages 2022–08, 2022.
- [8] Wouter Bulten, Kimmo Kartasalo, Po-Hsuan Cameron Chen, Peter Ström, Hans Pinckaers, Kunal Nagpal, Yuannan Cai, David F Steiner, Hester Van Boven, Robert Vink, et al. Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge. *Nature Medicine*, 28(1):154–163, 2022.
- [9] Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Mirafior, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature Medicine*, 25(8):1301–1309, 2019.
- [10] Tsai Hor Chan, Fernando Julio Cendra, Lan Ma, Guosheng Yin, and Lequan Yu. Histopathology whole slide image analysis with heterogeneous graph representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15661–15670, 2023.
- [11] Chengkuan Chen, Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Andrew J Schaumberg, and Faisal Mahmood. Fast and scalable search of whole-slide images via self-supervised deep learning. *Nature Biomedical Engineering*, 6(12):1420–1434, 2022.
- [12] Richard J Chen, Chengkuan Chen, Yicong Li, Tiffany Y Chen, Andrew D Trister, Rahul G Krishnan, and Faisal Mahmood. Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16144–16155, 2022.
- [13] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Bowen Chen, Andrew Zhang, Daniel Shao, Andrew H Song, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*, 2024.
- [14] Yuan-Chih Chen and Chun-Shien Lu. Rankmix: Data augmentation for weakly supervised learning of classifying whole slide images with diverse sizes and imbalanced categories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23936–23945, 2023.
- [15] Omar SM El Nahhas, Marko van Treeck, Georg Wölflin, Michaela Unger, Marta Ligerio, Tim Lenz, Sophia J Wagner, Katherine J Hewitt, Firas Khader, Sebastian Foersch, et al. From whole-slide image to

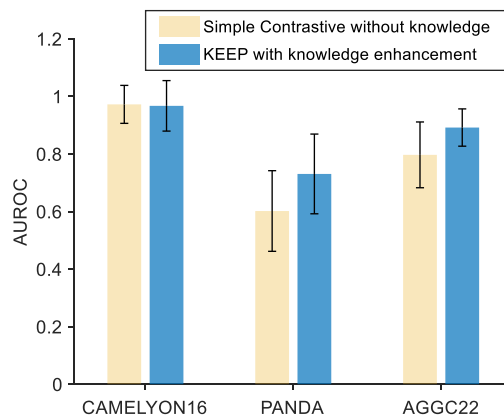
- biomarker prediction: end-to-end weakly supervised deep learning in computational pathology. *Nature Protocols*, pages 1–24, 2024.
- [16] Hossein Farahani, Jeffrey Boschman, David Farnell, Amirali Darbandsari, Allen Zhang, Pouya Ahmadvand, Steven JM Jones, David Huntsman, Martin Köbel, C Blake Gilks, et al. Deep learning-based histotype diagnosis of ovarian carcinoma whole-slide pathology images. *Modern Pathology*, 35(12):1983–1990, 2022.
- [17] Alexandre Filiot, Ridouane Ghermi, Antoine Olivier, Paul Jacob, Lucas Fidon, Alice Mac Kain, Charlie Saillard, and Jean-Baptiste Schiratti. Scaling self-supervised learning for histopathology with masked image modeling. *medRxiv*, pages 2023–07, 2023.
- [18] Jevgenij Gamper and Nasir Rajpoot. Multiple instance captioning: Learning representations from histopathology textbooks and articles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16549–16559, 2021.
- [19] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23, 2021.
- [20] Chu Han, Xipeng Pan, Lixu Yan, Huan Lin, Bingbing Li, Su Yao, Shanshan Lv, Zhenwei Shi, Jinhai Mai, Jiatai Lin, et al. Wss4luad: Grand challenge on weakly-supervised tissue semantic segmentation for lung adenocarcinoma. *arXiv preprint arXiv:2204.06455*, 2022.
- [21] Danh-Tai Hoang, Eldad D Shulman, Rust Turakulov, Zied Abdullaev, Omkar Singh, Emma M Campagnolo, H Lalchungnunga, Eric A Stone, MacLean P Nasrallah, Eytan Ruppim, et al. Prediction of dna methylation-based tumor types from histopathology in central nervous system tumors with deep learning. *Nature Medicine*, pages 1–10, 2024.
- [22] Yanyan Huang, Weiqin Zhao, Shujun Wang, Yu Fu, Yuming Jiang, and Lequan Yu. Conslide: Asynchronous hierarchical interaction transformer with breakup-reorganize rehearsal for continual whole slide image analysis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21349–21360, 2023.
- [23] Zhi Huang, Federico Bianchi, Mert Yuksekogul, Thomas J Montine, and James Zou. A visual–language foundation model for pathology image analysis using medical twitter. *Nature Medicine*, 29(9):2307–2316, 2023.
- [24] Xinmi Huo, Kok Haur Ong, Kah Weng Lau, Laurent Gole, David M Young, Char Loo Tan, Xiaohui Zhu, Chongchong Zhang, Yonghui Zhang, Longjie Li, et al. A comprehensive ai model development framework for consistent gleason grading. *Communications Medicine*, 4(1):84, 2024.
- [25] Wisdom Ikezogwo, Saygin Seyfioglu, Fatemeh Ghezloo, Dylan Geva, Fatwir Sheikh Mohammed, Pavan Kumar Anand, Ranjay Krishna, and Linda Shapiro. Quilt-1m: One million image-text pairs for histopathology. *Advances in Neural Information Processing Systems*, 36, 2024.
- [26] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.
- [27] Mingu Kang, Heon Song, Seonwook Park, Donggeun Yoo, and Sérgio Pereira. Benchmarking self-supervised learning on diverse pathology datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3344–3354, 2023.
- [28] Jakob Nikolas Kather, Niels Halama, and Alexander Marx. 100,000 histological images of human colorectal cancer and healthy tissue. *Zenodo10*, 5281, 2018.
- [29] Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A Valous, Dyke Ferber, et al. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS Medicine*, 16(1):e1002730, 2019.
- [30] Katharina Kriegsmann, Frithjof Lobers, Christiane Zgorzelski, Joerg Kriegsmann, Charlotte Janssen, Rolf Ruedinger Meliss, Thomas Muley, Ulrich Sack, Georg Steinbuss, and Mark Kriegsmann. Deep learning for the detection of anatomical tissue structures and neoplasms of the skin on scanned histopathological tissue sections. *Frontiers in Oncology*, 12:1022967, 2022.

- [31] Ritika Kundra, Hongxin Zhang, Robert Sheridan, Sahussapont Joseph Sirintrapun, Avery Wang, Angelica Ochoa, Manda Wilson, Benjamin Gross, Yichao Sun, Ramyasree Madupuri, et al. Oncotree: a cancer classification system for precision oncology. *JCO Clinical Cancer Informatics*, 5:221–230, 2021.
- [32] Honglin Li, Chenglu Zhu, Yunlong Zhang, Yuxuan Sun, Zhongyi Shui, Wenwei Kuang, Sunyi Zheng, and Lin Yang. Task-specific fine-tuning via variational information bottleneck for weakly-supervised pathology whole slide image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7454–7463, 2023.
- [33] Tiancheng Lin, Zhimiao Yu, Hongyu Hu, Yi Xu, and Chang-Wen Chen. Interventional bag multi-instance learning on whole-slide pathological images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19830–19839, 2023.
- [34] 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.
- [35] Ming Y Lu, Bowen Chen, Andrew Zhang, Drew FK Williamson, Richard J Chen, Tong Ding, Long Phi Le, Yung-Sung Chuang, and Faisal Mahmood. Visual language pretrained multiple instance zero-shot transfer for histopathology images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19764–19775, 2023.
- [36] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering*, 5(6):555–570, 2021.
- [37] Peter Neidlinger, Omar SM El Nahhas, Hannah Sophie Muti, Tim Lenz, Michael Hoffmeister, Hermann Brenner, Marko van Treeck, Rupert Langer, Bastian Dislich, Hans Michael Behrens, et al. Benchmarking foundation models as feature extractors for weakly-supervised computational pathology. *arXiv preprint arXiv:2408.15823*, 2024.
- [38] 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, 2021.
- [39] Thomas Roetzer-Pejrimovsky, Anna-Christina Moser, Baran Atli, Clemens Christian Vogel, Petra A Mercea, Romana Prihoda, Ellen Gelpi, Christine Haberler, Romana Höftberger, Johannes A Hainfellner, et al. The digital brain tumour atlas, an open histopathology resource. *Scientific Data*, 9(1):55, 2022.
- [40] Lynn Marie Schriml, Cesar Arze, Suvarna Nadendla, Yu-Wei Wayne Chang, Mark Mazaitis, Victor Felix, Gang Feng, and Warren Alden Kibbe. Disease ontology: a backbone for disease semantic integration. *Nucleic Acids Research*, 40(D1):D940–D946, 2012.
- [41] Muhammad Shaban, Ruqayya Awan, Muhammad Moazam Fraz, Ayesha Azam, Yee-Wah Tsang, David Snead, and Nasir M Rajpoot. Context-aware convolutional neural network for grading of colorectal cancer histology images. *IEEE Transactions on Medical Imaging*, 39(7):2395–2405, 2020.
- [42] George Shaikovski, Adam Casson, Kristen Severson, Eric Zimmermann, Yi Kan Wang, Jeremy D Kunz, Juan A Retamero, Gerard Oakley, David Klimstra, Christopher Kanan, et al. Prism: A multi-modal generative foundation model for slide-level histopathology. *arXiv preprint arXiv:2405.10254*, 2024.
- [43] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, et al. Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in Neural Information Processing Systems*, 34:2136–2147, 2021.
- [44] Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *IEEE Transactions on Biomedical Engineering*, 63(7):1455–1462, 2015.
- [45] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song, David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, et al. Bioclip: A vision foundation model for the tree of life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19412–19424, 2024.

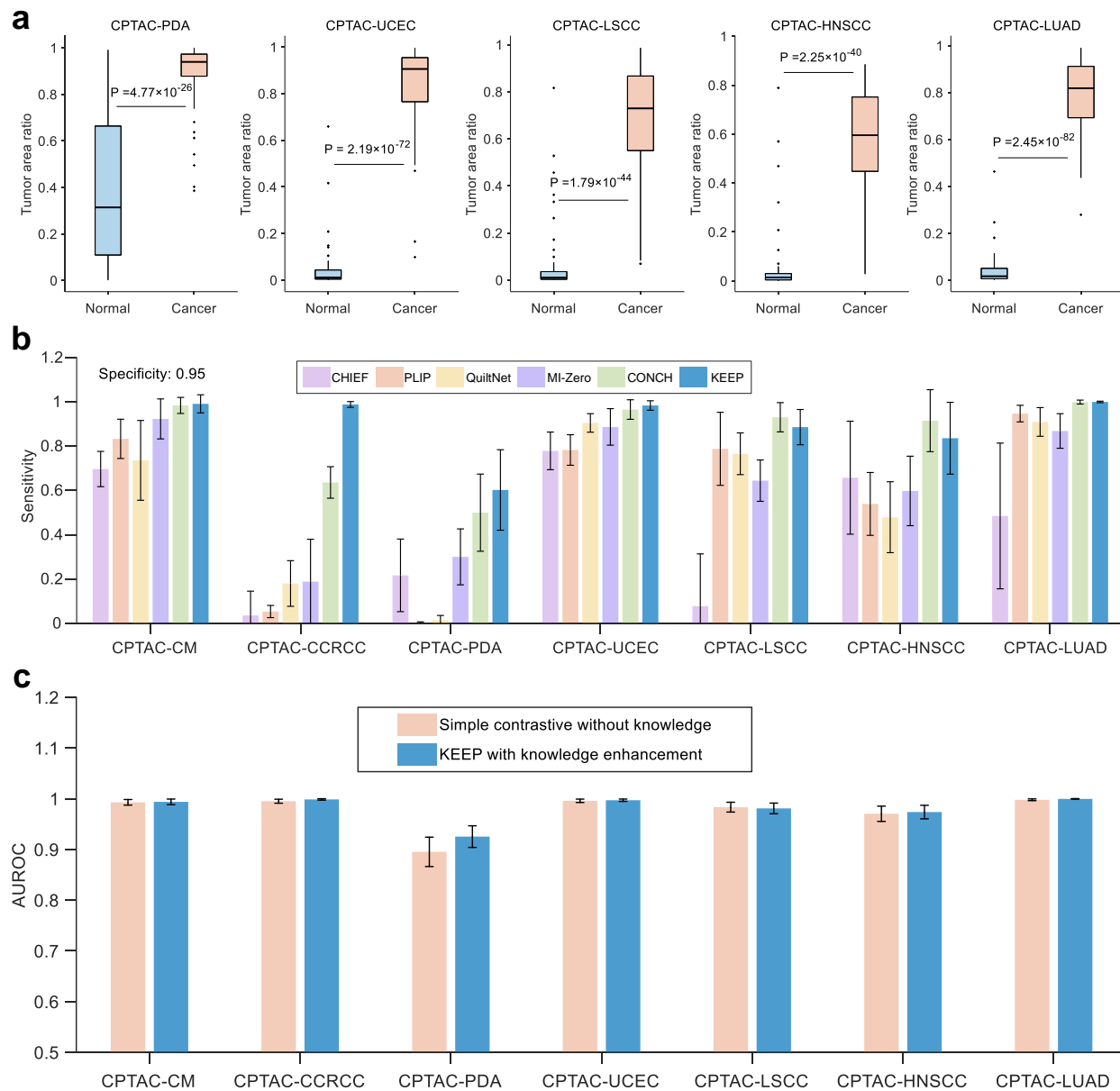
- [46] Yuri Tolkach, Lisa Marie Wolgast, Alexander Damanakis, Alexey Pryalukhin, Simon Schallenberg, Wolfgang Hulla, Marie-Lisa Eich, Wolfgang Schroeder, Anirban Mukhopadhyay, Moritz Fuchs, et al. Artificial intelligence for tumour tissue detection and histological regression grading in oesophageal adenocarcinomas: a retrospective algorithm development and validation study. *The Lancet Digital Health*, 5(5):e265–e275, 2023.
- [47] 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*, pages 1–12, 2024.
- [48] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7464–7475, 2023.
- [49] Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical Image Analysis*, 81:102559, 2022.
- [50] Xiyue Wang, Junhan Zhao, Eliana Marostica, Wei Yuan, Jietian Jin, Jiayu Zhang, Ruijiang Li, Hongping Tang, Kanran Wang, Yu Li, et al. A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature*, pages 1–9, 2024.
- [51] 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*, pages 1–8, 2024.
- [52] Yingxue Xu, Yihui Wang, Fengtao Zhou, Jiabo Ma, Shu Yang, Huangjing Lin, Xin Wang, Jiguang Wang, Li Liang, Anjia Han, et al. A multimodal knowledge-enhanced whole-slide pathology foundation model. *arXiv preprint arXiv:2407.15362*, 2024.
- [53] Xiao Zhou, Zhen Cheng, Miao Gu, and Fei Chang. Lirnet: Local integral regression network for both strongly and weakly supervised nuclei detection. In *2020 IEEE International Conference on Bioinformatics and Biomedicine*, pages 945–951. IEEE, 2020.
- [54] Xiao Zhou, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yanfeng Wang. Knowledge-enhanced visual-language pretraining for computational pathology. *arXiv preprint arXiv:2404.09942*, 2024.
- [55] Xiao Zhou, Yujie Zhong, Zhen Cheng, Fan Liang, and Lin Ma. Adaptive sparse pairwise loss for object re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19691–19701, 2023.
- [56] Chuang Zhu, Wenkai Chen, Ting Peng, Ying Wang, and Mulan Jin. Hard sample aware noise robust learning for histopathology image classification. *IEEE Transactions on Medical Imaging*, 41(4):881–894, 2021.
- [57] 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.



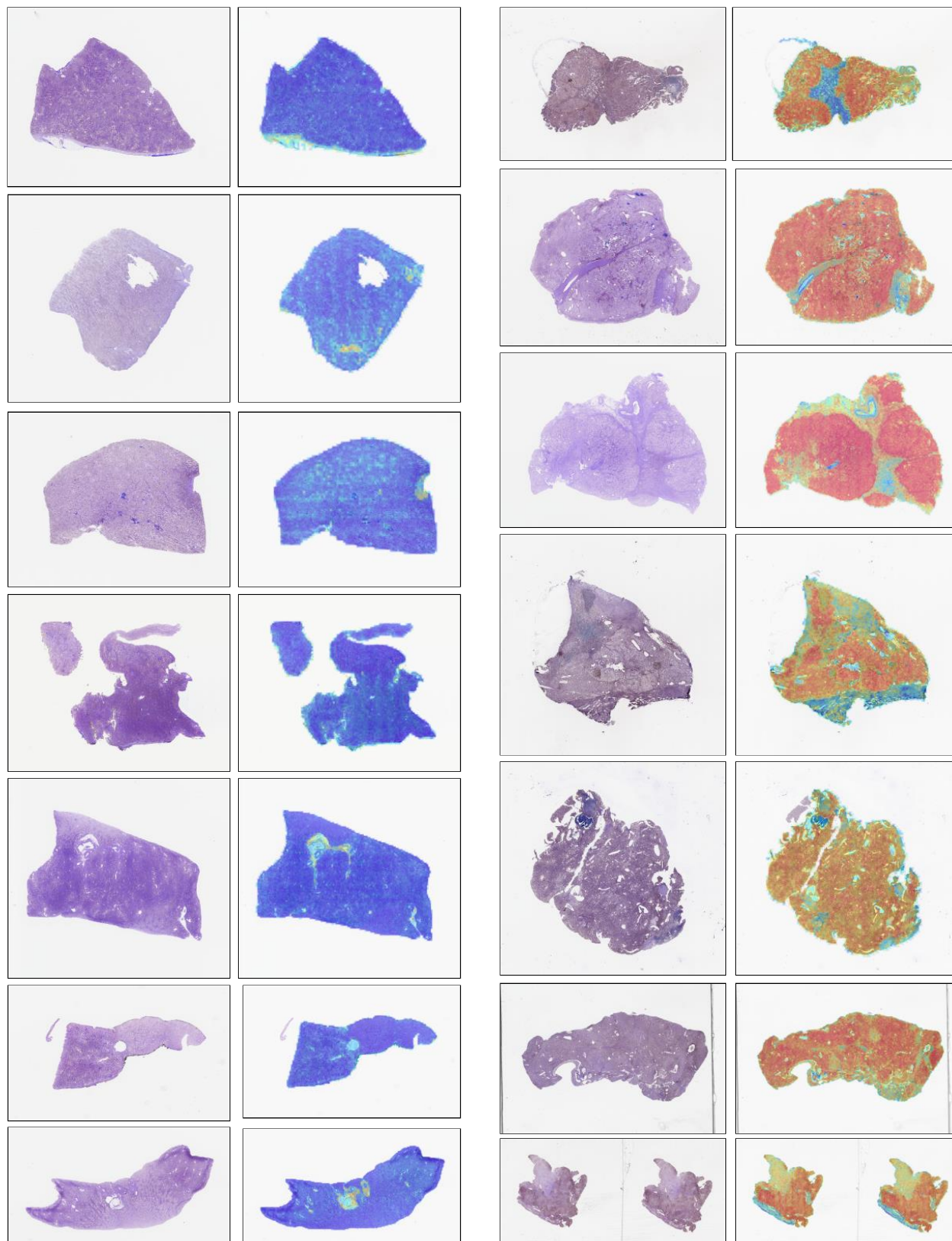
Supplementary Figure S1 | Statistics of semantic groups. Distributions of cancer types in example anatomies, including breast, kidney, brain/cns, lung, ovarian, and esophagus/stomach. Top-5 cancer subtypes are listed by the arrowed text. The gray and yellow bars represent the non-leaf and leaf nodes in the constructed knowledge base. The red text suggests the subtypes overlapped with the cancer subtyping tasks in this paper.



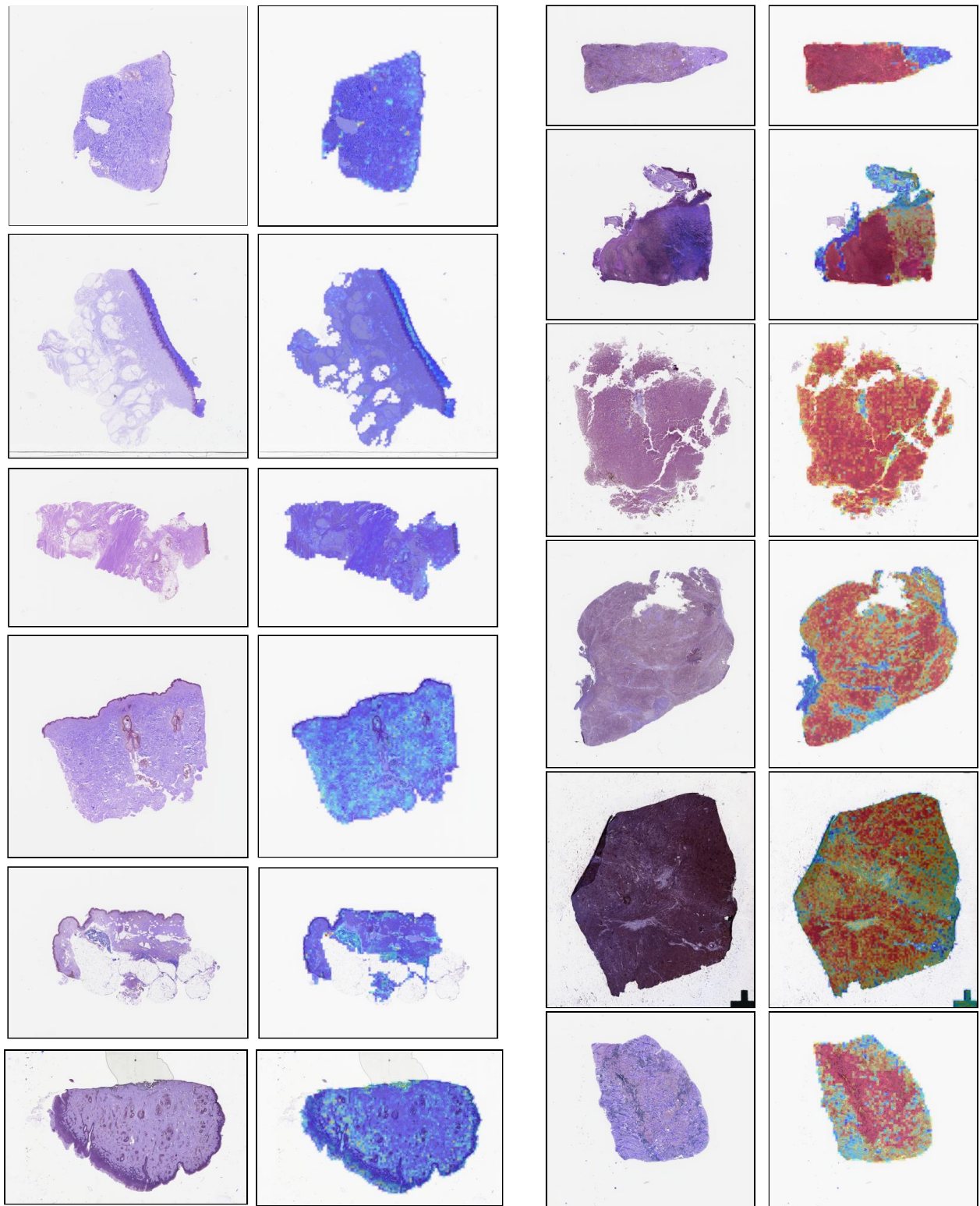
Supplementary Figure S2 | Ablation results. Performance comparison between simple contrastive without knowledge and KEEP with knowledge enhancement on cancer region segmentation. KEEP significantly outperforms simple contrastive on PANDA and AGGC22 datasets.



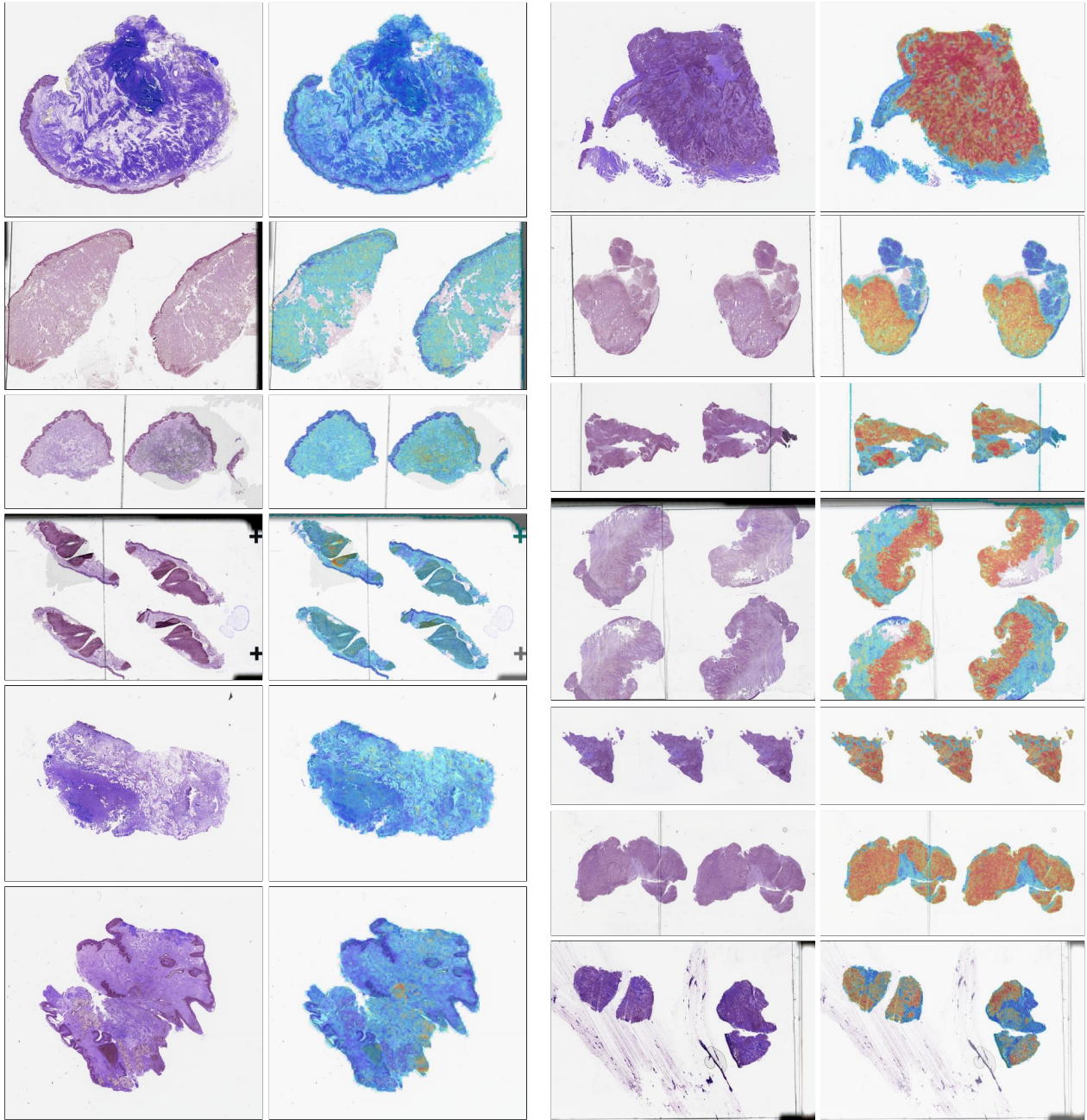
Supplementary Figure S3 | Additional results of cancer detection. **a.** The comparison of the predicted tumor ratio between normal and cancer WSIs in the datasets of CPTAC-PDA, CPTAC-UCEC, CPTAC-LSCC, CPTAC-HNSCC, and CPTAC-LUAD. Each dataset consists of 75 normal WSIs and 75 cancer WSIs. **b.** Comparison of cancer detection sensitivities on each dataset at the specificity of 0.95, the error bar denotes the standard deviation of the performance with 1,000 bootstrap iterations. **c.** Performance comparison between simple contrastive without knowledge and KEEP with knowledge enhancement on cancer detection. KEEP outperforms simple contrastive on 6 out of 7 datasets.



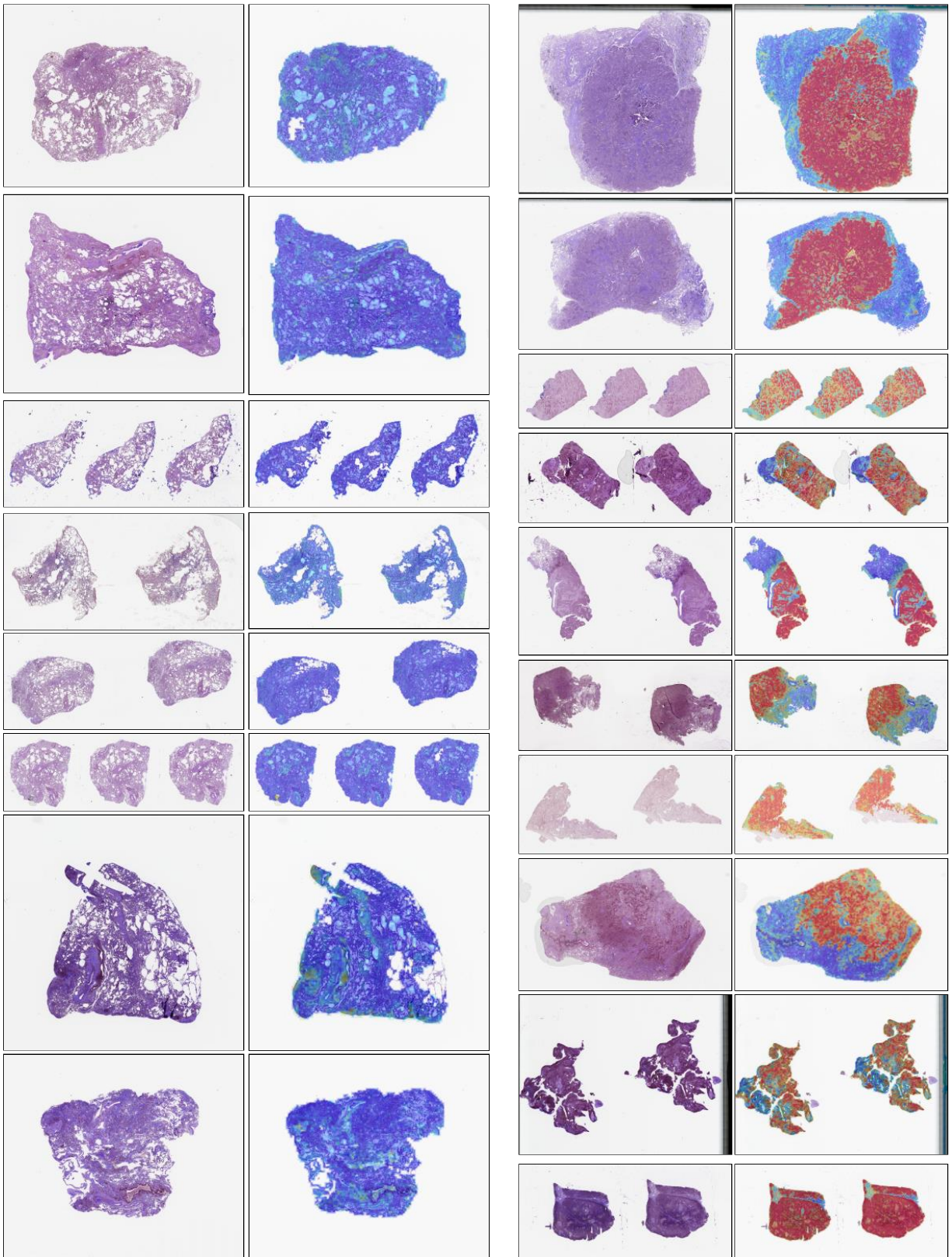
Supplementary Figure S4 | Additional visualization of cancer detection. Example results on CPATC-CCRCC, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



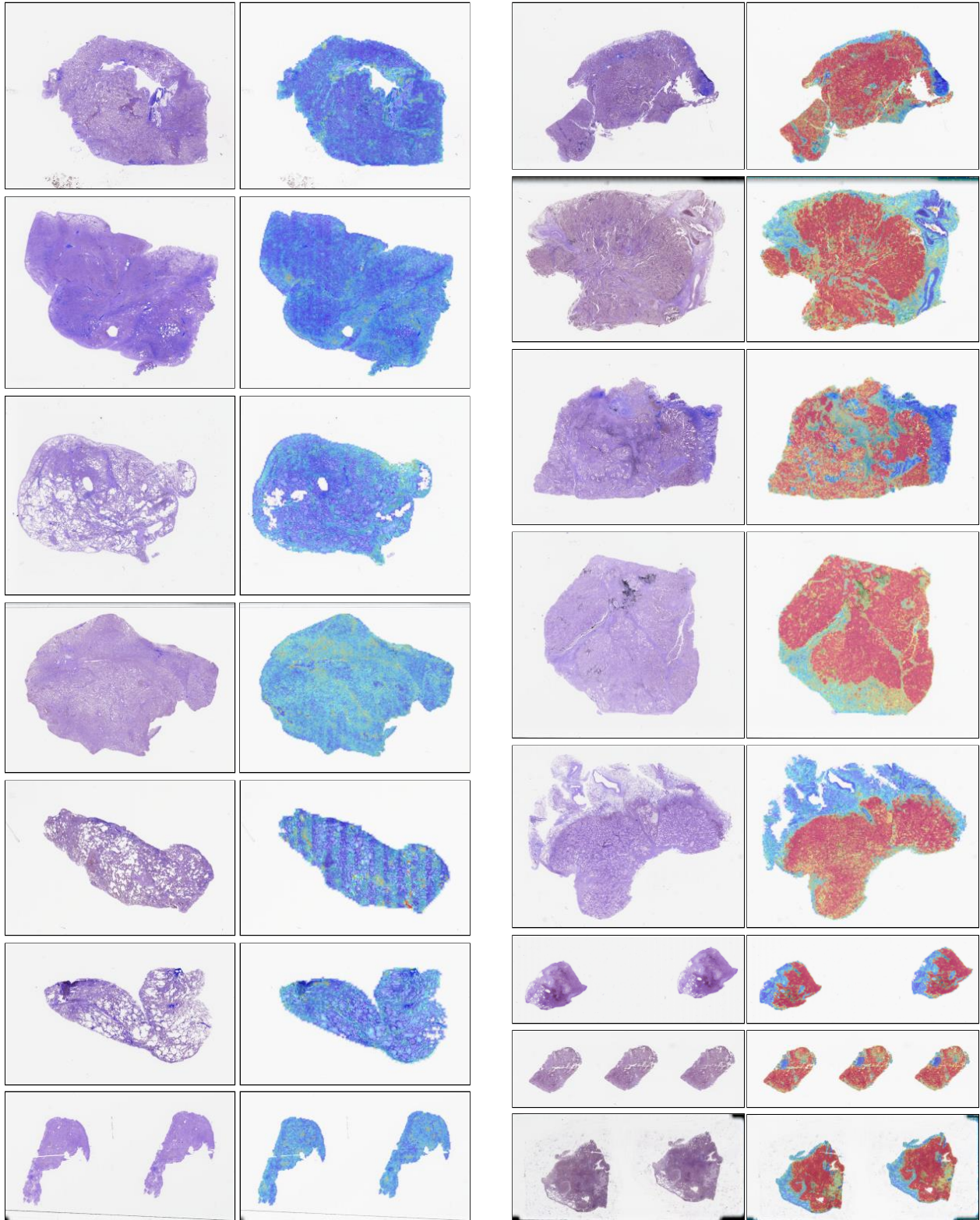
Supplementary Figure S5 | Additional visualization of cancer detection. Example results on CPATC-CM, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



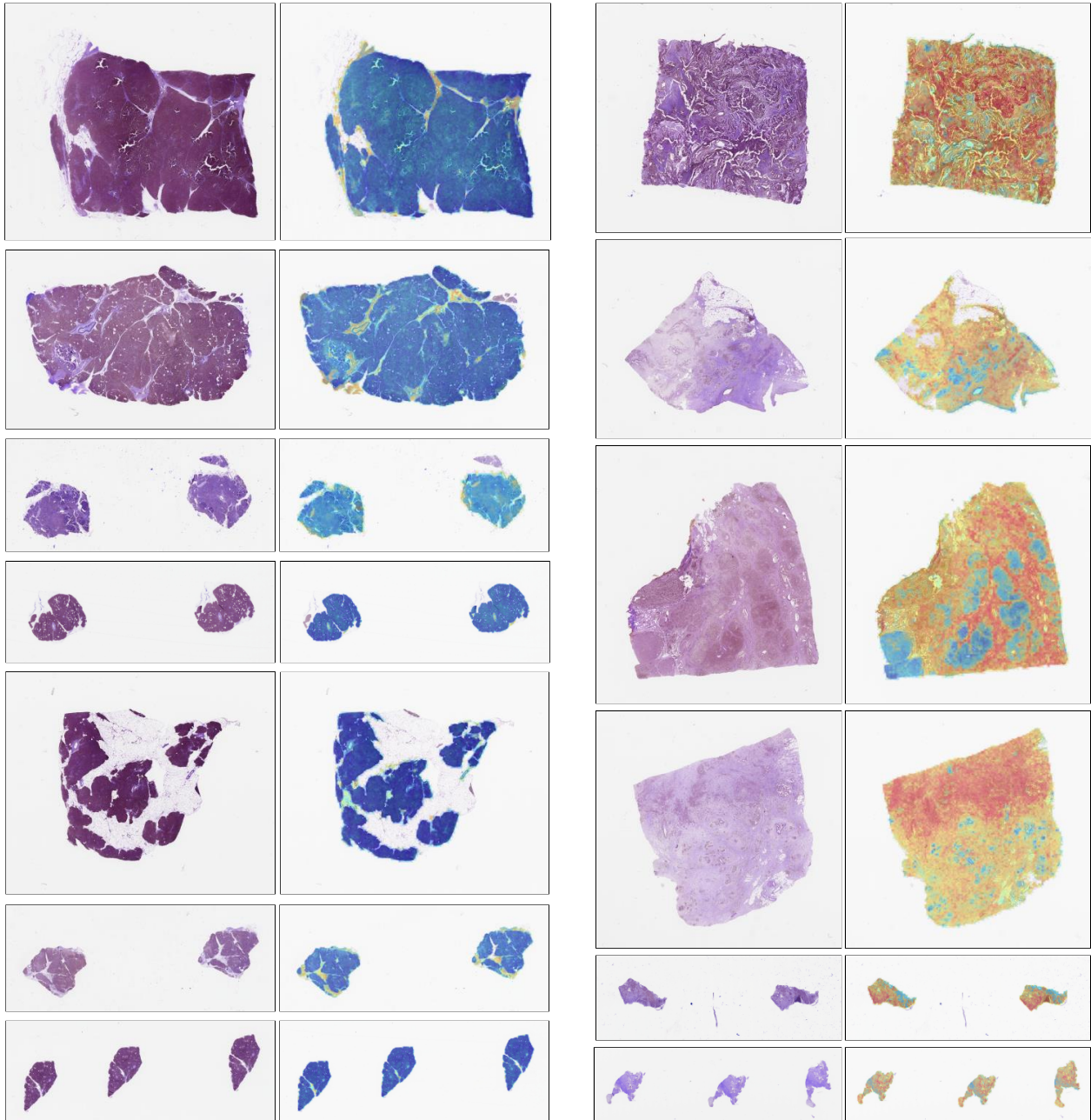
Supplementary Figure S6 | Additional visualization of cancer detection. Example results on CPATC-HNSCC, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



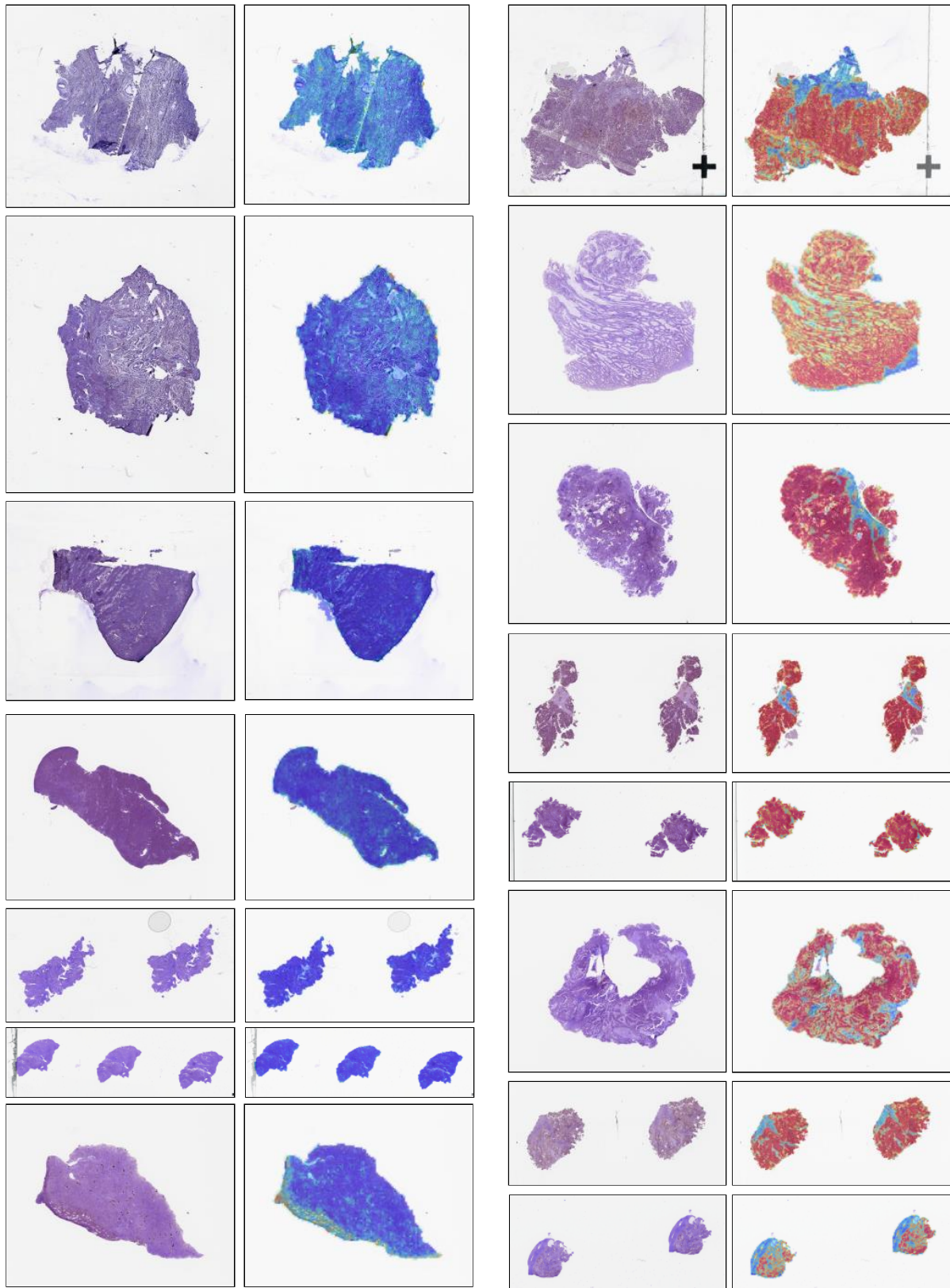
Supplementary Figure S7 | Additional visualization of cancer detection. Example results on CPATC-LSCC, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



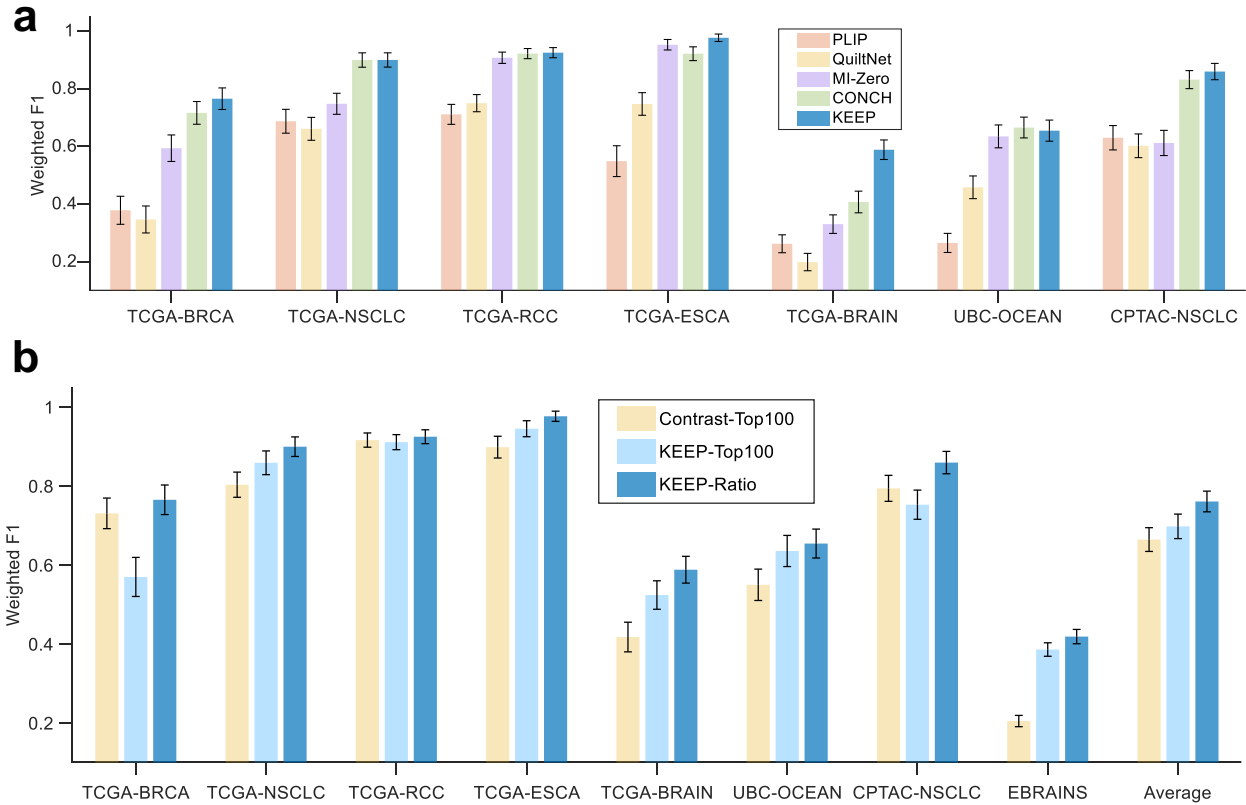
Supplementary Figure S8 | Additional visualization of cancer detection. Example results on CPATC-LUAD, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



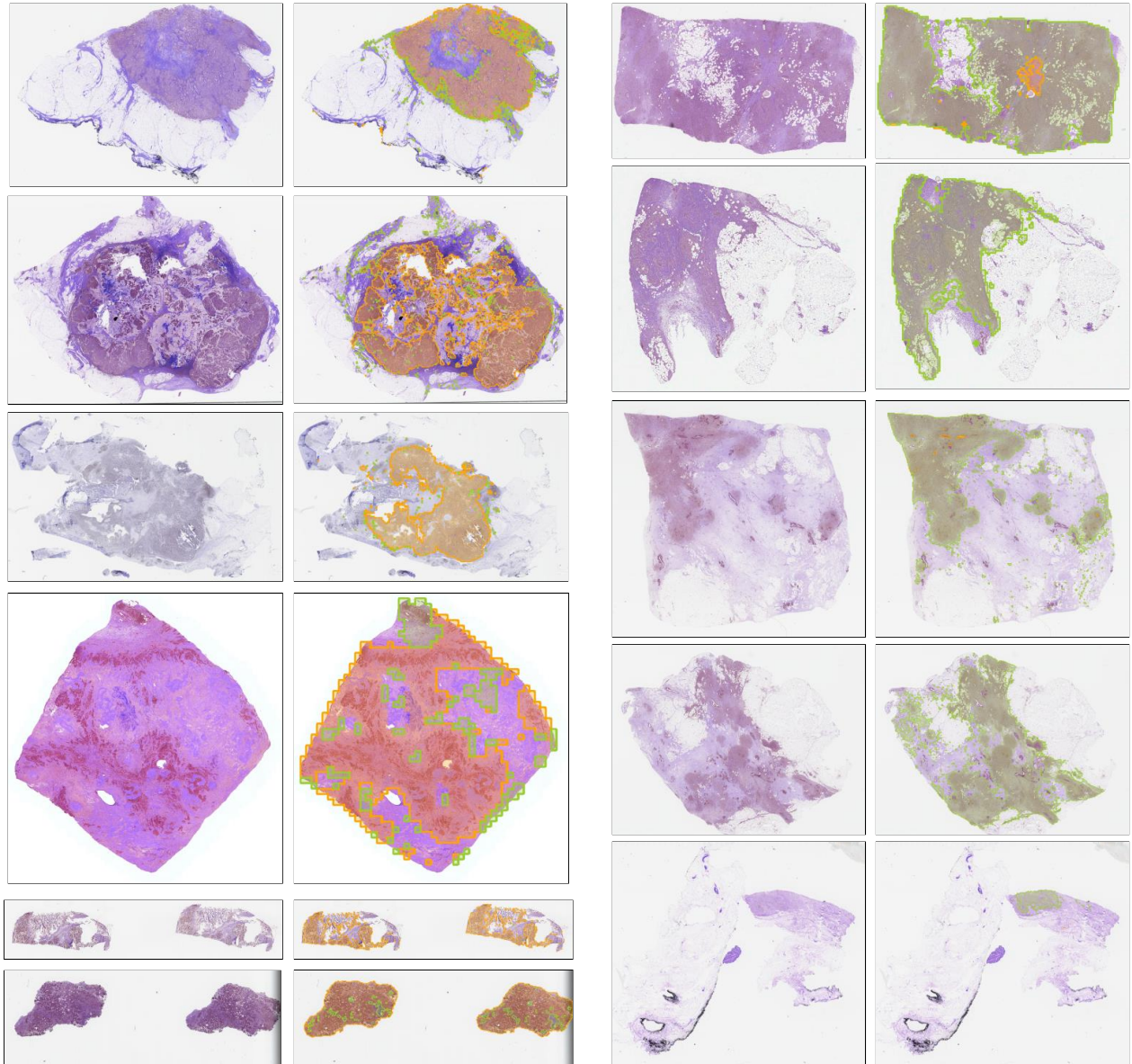
Supplementary Figure S9 | Additional visualization of cancer detection. Example results on CPATC-PDA, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* [Template] + *cancerous tissue*.



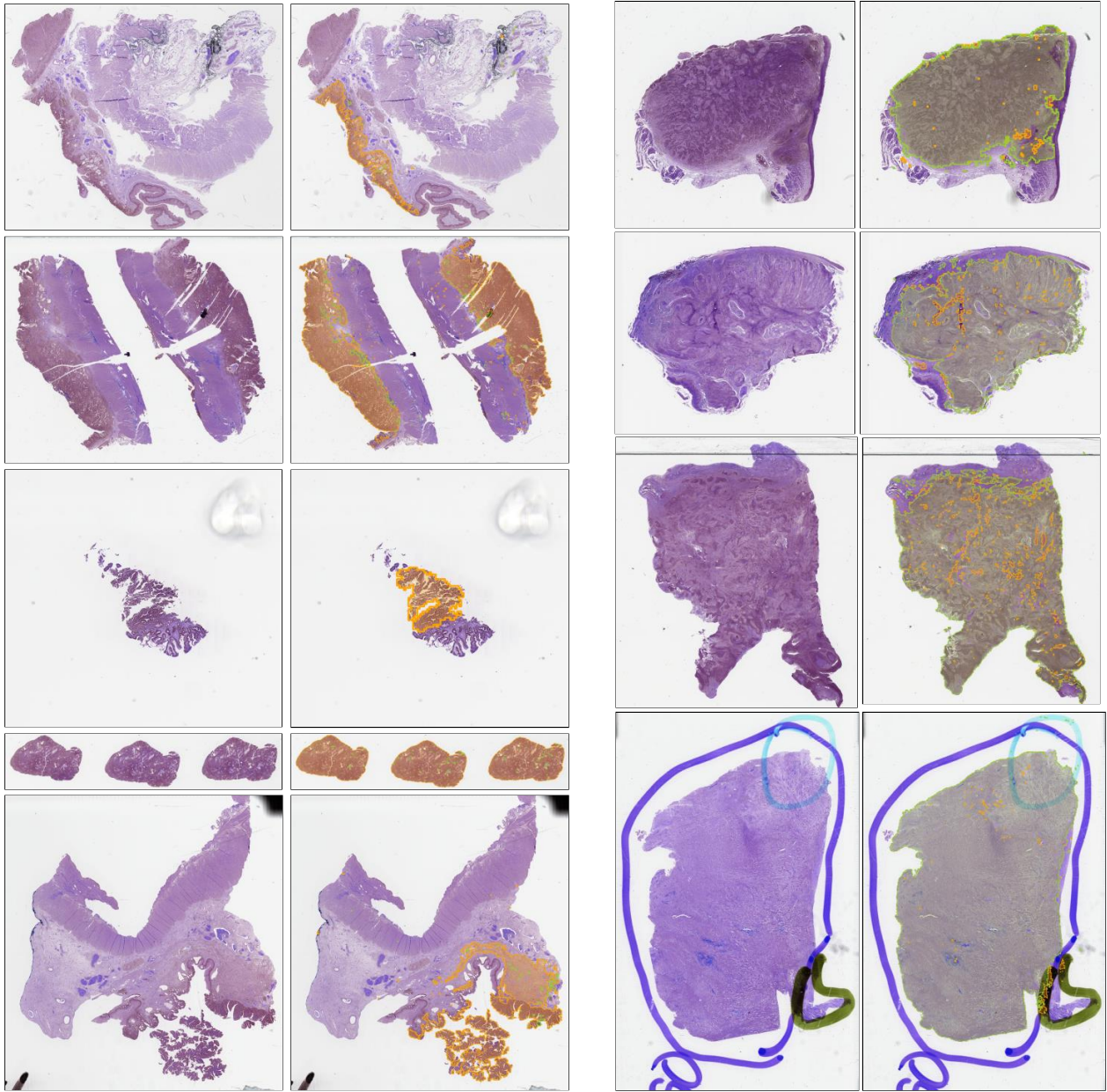
Supplementary Figure S10 | Additional visualization of cancer detection. Example results on CPATC-UCEC, the left and the right two columns represent the normal and the cancer WSIs, respectively. The heatmap is generated by the similarities between the embeddings of tile images and that of "tumor" prompts, *i.e.* {Template} + *cancerous tissue*.



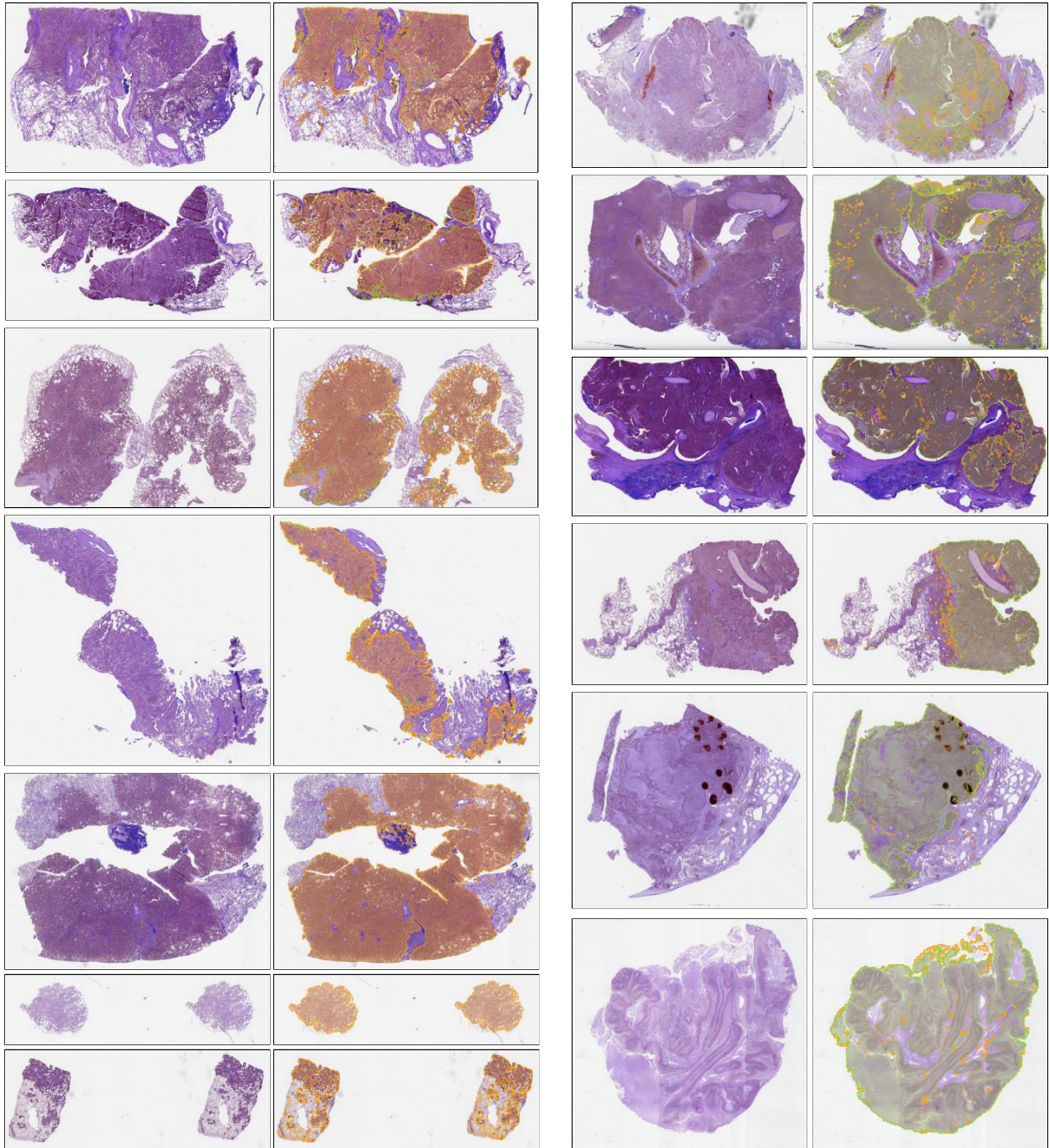
Supplementary Figure S11 | Additional results of cancer subtyping. a. Comparison of weighted F1 across different models on seven datasets with common cancer subtypes. The TCGA-BRCA, TCGA-NSCLC, TCGA-ESCA, and CPTAC-NSCLC datasets contain two subtypes, while the TCGA-RCC, TCGA-BRAIN, and UBC-OCEAN datasets consist of 3, 3, and 5 subtypes, respectively. Each subtype includes 75 WSIs, except for TCGA-ESCA (65 WSIs) and UBC-OCEAN (35 WSIs), with each experiment using 1,000 bootstrap iterations. **b.** Ablation results. Performance comparison between simple contrastive (Contrastive-Top100) and KEEP with knowledge enhancement (KEEP-Top100). Top100 suggests the strategy of top-100 pooling, while Ratio denotes the subtype ratio strategy.



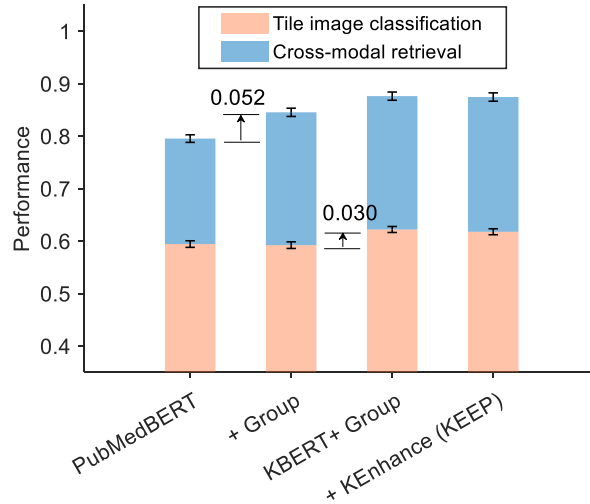
Supplementary Figure S12 | Additional visualization of cancer subtyping. Example results on TCGA-BRCA, the left and the right two columns represent the cancer WSIs from breast invasive ductal carcinoma (IDC) and breast invasive lobular carcinoma (ILC), respectively. The orange and the green masks denote the predicted regions of IDC and ILC, respectively.



Supplementary Figure S13 | Additional visualization of cancer subtyping. Example results on TCGA-ESCA, the left and the right two columns represent the cancer WSIs from esophagus adenocarcinoma and esophagus squamous cell carcinoma, respectively. The orange and green masks denote the predicted regions of esophagus adenocarcinoma and esophagus squamous cell carcinoma, respectively.



Supplementary Figure S14 | Additional visualization of cancer subtyping. Example results on TCGA-NSCLC, the left and the right two columns represent the cancer WSIs from lung adenocarcinoma and lung squamous cell carcinoma, respectively. The orange and the green masks denote the predicted regions of lung adenocarcinoma and lung squamous cell carcinoma, respectively.



Supplementary Figure S15 | Ablation results. Comparison of performance on tile-level tasks. "PubMedBERT" indicates that the text encoder is initialized using PubMedBERT, with simple contrastive learning applied to align images and their paired captions. "+ Group" signifies that semantic group alignment is employed instead of direct image-text pair alignment. "KBERT + Group" incorporates a knowledge encoder alongside semantic group alignment. "+ KEnhance (KEEP)" further builds on "KBERT + Group" by introducing knowledge-enhanced caption augmentation and strategies for eliminating false negatives.

Supplementary Table S1 | Statistics of test datasets used in this paper.

Task	Dataset	Anatomy	#images	#types
Cancer Region Segmentation	CAMELYON16	Breast	48 (Test)	-
	PANDA	Prostate	10,494	-
	AGGC22	Prostate	128 (Test)	-
Cancer Detection	CPTAC-CM	Skin	150	2
	CPTAC-CCRCC	Skin	150	2
	CPTAC-PDA	Pancreas	150	2
	CPTAC-UCEC	Uterine	150	2
	CPTAC-LSCC	Lung	150	2
	CPTAC-HNSCC	Head and neck	150	2
	CPTAC-LUAD	Lung	150	2
Cancer Subtyping	TCGA-BRCA	Breast	150	2
	TCGA-NSCLC	Lung	150	2
	TCGA-RCC	Kidney	225	3
	TCGA-ESCA	Esophagus	130	2
	TCGA-BRAIN	Brain	225	3
	UBC-OCEAN	Kidney	175	5
	CPTAC-NSCLC	Lung	150	2
	EBRAINS	Brain	900	30
	Cross-modal Retrieval	Arch-PubMed	Multiple	1,923
Arch-book		Multiple	1,306	-
PathPair		Multiple	9,358	-
WebPath		Multiple	12,192	-
Tile Classification	BACH	Breast	400	4
	NCT-CRC-HE-100K	Colorectal	100K	9
	CRC100K	Colon	7,180	9
	LC25000	Lung, Colon	25K	5
	RenalCell	Renal	36,687	5
	SkinCancer	Skin	129,369	16
	WSSS4LUAD	Lung	4,693	3
	Osteo	Bone	1,144	3
	ESCA-UKK	Esophagus	34,704	11
	ESCA-WNS	Esophagus	121,642	11
	ESCA-TCGA	Esophagus	32,796	11
	ESCA-CHA	Esophagus	178,187	11
	Breakhis	Breast	7,909	8
	Chaoyang	Colon	2,139	4

Supplementary Table S2 | URLs of test datasets used in this paper.

Dataset	URL
CAMELYON16	https://CAMELYON16.grand-challenge.org/Data/
PANDA	https://www.kaggle.com/c/prostate-cancer-grade-assessment/
AGGC22	https://aggc22.grand-challenge.org/Download/
CPTAC-CM	https://www.cancerimagingarchive.net/collection/cptac-cm/
CPTAC-CCRCC	https://www.cancerimagingarchive.net/collection/cptac-ccrcc/
CPTAC-PDA	https://www.cancerimagingarchive.net/collection/cptac-pda/
CPTAC-UCEC	https://www.cancerimagingarchive.net/collection/cptac-ucec/
CPTAC-LSCC	https://www.cancerimagingarchive.net/collection/cptac-lsc/
CPTAC-HNSCC	https://www.cancerimagingarchive.net/collection/cptac-hnsc/
CPTAC-LUAD	https://www.cancerimagingarchive.net/collection/cptac-luad/
TCGA-BRCA	https://portal.gdc.cancer.gov/projects/TCGA-BRCA
TCGA-NSCLC	https://portal.gdc.cancer.gov/projects/TCGA-LUAD , https://portal.gdc.cancer.gov/projects/TCGA-LUSC
TCGA-RCC	https://portal.gdc.cancer.gov/projects/TCGA-KIRP , https://portal.gdc.cancer.gov/projects/TCGA-KIRC , https://portal.gdc.cancer.gov/projects/TCGA-KICH
TCGA-ESCA	https://portal.gdc.cancer.gov/projects/TCGA-ESCA
TCGA-BRAIN	https://portal.gdc.cancer.gov/projects/TCGA-LGG , https://portal.gdc.cancer.gov/projects/TCGA-GBM
UBC-OCEAN	https://www.kaggle.com/competitions/UBC-OCEAN/
CPTAC-NSCLC	https://www.cancerimagingarchive.net/collection/cptac-lsc/ , https://www.cancerimagingarchive.net/collection/cptac-luad/
EBRAINS	https://data-proxy.ebrains.eu/datasets/8fc108ab-e2b4-406f-8999-60269dc1f994
Arch-PubMed	https://warwick.ac.uk/fac/cross_fac/tia/data/arch
Arch-book	https://warwick.ac.uk/fac/cross_fac/tia/data/arch
PathPair	https://github.com/MAGIC-AI4Med/KEP
WebPath	-
BACH	https://iciar2018-challenge.grand-challenge.org/Dataset/
NCT-CRC-HE-100K	https://zenodo.org/records/1214456
CRC100K	https://zenodo.org/records/1214456
LC25000	https://academictorrents.com/details/7a638ed187a6180fd6e464b3666a6ea0499af4af
RenalCell	https://zenodo.org/records/6528599
SkinCancer	https://doi.org/10.11588/data/7QCR8S
WSSS4LUAD	https://wsss4luad.grand-challenge.org/
Osteo	https://doi.org/10.7937/tcia.2019.bvhhjhdas
ESCA-UKK	https://zenodo.org/records/7548828
ESCA-WNS	https://zenodo.org/records/7548828
ESCA-TCGA	https://zenodo.org/records/7548828
ESCA-CHA	https://zenodo.org/records/7548828
Breakhis	http://web.inf.ufpr.br/vri/breast-cancer-database
Chaoyang	https://bupt-ai-cz.github.io/HSA-NRL/

Supplementary Table S3 | Zero-shot performance of the cancer region segmentation task. Bold suggests the best performance.

Metric	Models	PLIP	QuiltNet	MI-Zero (Pub)	CONCH	KEEP(Ours)
AUROC	CAMELYON16 [4]	0.950±0.064	0.800±0.174	0.876±0.145	0.962±0.078	0.967±0.088
	PANDA [8]	0.685±0.124	0.698±0.149	0.642±0.152	0.727±0.125	0.731±0.139
	AGGC22 [24]	0.726±0.113	0.721±0.123	0.765±0.112	0.853±0.095	0.892±0.065
DICE	CAMELYON16 [4]	0.253±0.312	0.157±0.233	0.186±0.262	0.292±0.333	0.361±0.332
	PANDA [8]	0.295±0.234	0.309±0.243	0.276±0.218	0.315±0.232	0.334±0.261
	AGGC22 [24]	0.284±0.158	0.282±0.166	0.324±0.147	0.449±0.177	0.530±0.179

Supplementary Table S4 | Zero-shot performance (AUROC) of the cancer detection task. The performance is represented by the median and its 95% CIs. Bold suggests the best performance.

Models	CHIEF	PLIP	QuiltNet	MI-Zero (Pub)	CONCH	KEEP(Ours)
CPTAC-CM	0.915 (0.866, 0.955)	0.970 (0.941, 0.990)	0.972 (0.944, 0.993)	0.985 (0.965, 0.998)	0.994 (0.982, 1.000)	0.994 (0.981, 1.000)
CPTAC-CCRCC	0.723 (0.644, 0.799)	0.330 (0.252, 0.419)	0.755 (0.677, 0.830)	0.886 (0.818, 0.943)	0.871 (0.809, 0.925)	0.999 (0.995, 1.000)
CPTAC-PDA	0.825 (0.751, 0.890)	0.391 (0.296, 0.480)	0.464 (0.371, 0.561)	0.796 (0.719, 0.864)	0.920 (0.867, 0.962)	0.929 (0.882, 0.966)
CPTAC-UCEC	0.955 (0.919, 0.983)	0.945 (0.905, 0.975)	0.973 (0.947, 0.991)	0.979 (0.954, 0.995)	0.996 (0.988, 1.000)	0.998 (0.991, 1.000)
CPTAC-LSCC	0.901 (0.848, 0.944)	0.965 (0.933, 0.987)	0.966 (0.936, 0.987)	0.910 (0.863, 0.952)	0.987 (0.963, 0.998)	0.983 (0.958, 0.997)
CPTAC-HNSCC	0.946 (0.899, 0.980)	0.898 (0.840, 0.947)	0.874 (0.804, 0.930)	0.918 (0.866, 0.961)	0.982 (0.948, 0.999)	0.976 (0.942, 0.996)
CPTAC-LUAD	0.891 (0.837, 0.937)	0.988 (0.970, 0.998)	0.991 (0.977, 0.998)	0.981 (0.959, 0.995)	0.999 (0.996, 1.000)	1.000 (0.998, 1.000)

Supplementary Table S5 | Zero-shot performance (weighted F1 and balanced accuracy) of the cancer subtyping task. The performance is represented by the median and its 95% CIs. Bold suggests the best performance.

Metric	Models	PLIP	QuiltNet	MI-Zero (Pub)	CONCH	KEEP(Ours)
WF1	TCGA-BRCA	0.376 (0.287, 0.475)	0.346 (0.256, 0.445)	0.593 (0.498, 0.682)	0.717 (0.644, 0.790)	0.765 (0.695, 0.831)
	TCGA-NSCLC	0.687 (0.603, 0.763)	0.665 (0.585, 0.743)	0.746 (0.669, 0.817)	0.900 (0.853, 0.947)	0.900 (0.853, 0.947)
	TCGA-RCC	0.715 (0.650, 0.777)	0.750 (0.691, 0.806)	0.907 (0.868, 0.942)	0.921 (0.886, 0.951)	0.925 (0.890, 0.960)
	TCGA-ESCA	0.547 (0.443, 0.653)	0.746 (0.670, 0.817)	0.954 (0.916, 0.985)	0.923 (0.876, 0.969)	0.977 (0.946, 1.000)
	TCGA-BRAIN	0.261 (0.203, 0.331)	0.199 (0.146, 0.260)	0.332 (0.269, 0.394)	0.409 (0.339, 0.478)	0.588 (0.520, 0.652)
	UBC-OCEAN	0.263 (0.202, 0.332)	0.455 (0.378, 0.530)	0.635 (0.550, 0.707)	0.666 (0.591, 0.732)	0.655 (0.580, 0.727)
	CPTAC-NSCLC	0.631 (0.548, 0.709)	0.604 (0.516, 0.682)	0.615 (0.527, 0.694)	0.837 (0.761, 0.893)	0.861 (0.803, 0.913)
	EBRAINS	0.065 (0.050, 0.082)	0.060 (0.044, 0.078)	0.294 (0.264, 0.326)	0.317 (0.285, 0.349)	0.420 (0.386, 0.454)
BACC	TCGA-BRCA	0.519 (0.500, 0.544)	0.500 (0.480, 0.519)	0.633 (0.576, 0.689)	0.727 (0.666, 0.787)	0.774 (0.712, 0.829)
	TCGA-NSCLC	0.699 (0.631, 0.766)	0.667 (0.592, 0.740)	0.753 (0.692, 0.816)	0.901 (0.854, 0.947)	0.902 (0.852, 0.946)
	TCGA-RCC	0.735 (0.690, 0.780)	0.755 (0.705, 0.803)	0.908 (0.869, 0.943)	0.921 (0.885, 0.951)	0.926 (0.890, 0.959)
	TCGA-ESCA	0.614 (0.565, 0.670)	0.746 (0.670, 0.820)	0.954 (0.919, 0.985)	0.923 (0.874, 0.968)	0.977 (0.946, 1.000)
	TCGA-BRAIN	0.361 (0.321, 0.405)	0.346 (0.333, 0.363)	0.361 (0.304, 0.419)	0.453 (0.407, 0.503)	0.604 (0.547, 0.661)
	UBC-OCEAN	0.343 (0.289, 0.391)	0.469 (0.396, 0.539)	0.652 (0.590, 0.709)	0.674 (0.610, 0.737)	0.661 (0.602, 0.721)
	CPTAC-NSCLC	0.647 (0.580, 0.715)	0.607 (0.524, 0.683)	0.643 (0.576, 0.697)	0.836 (0.770, 0.889)	0.863 (0.806, 0.914)
	EBRAINS	0.096 (0.080, 0.110)	0.093 (0.081, 0.107)	0.325 (0.300, 0.349)	0.371 (0.351, 0.390)	0.456 (0.432, 0.479)

Supplementary Table S6 | Zero-shot performance of cross-modal retrieval. The mean and the standard deviation represent the performance. Bold suggests the best performance.

Dataset	Model	Text-to-image			Image-to-text		
		Recall@5	Recall@10	Recall@50	Recall@5	Recall@10	Recall@50
ARCH-PubMed	PLIP	0.037 ± 0.004	0.067 ± 0.006	0.182 ± 0.009	0.037 ± 0.004	0.066 ± 0.006	0.186 ± 0.009
	QuiltNet	0.082 ± 0.006	0.122 ± 0.008	0.292 ± 0.011	0.093 ± 0.007	0.139 ± 0.008	0.325 ± 0.011
	KEEP	0.160 ± 0.008	0.236 ± 0.010	0.480 ± 0.012	0.172 ± 0.008	0.250 ± 0.010	0.487 ± 0.012
ARCH-book	PLIP	0.111 ± 0.009	0.164 ± 0.010	0.419 ± 0.014	0.096 ± 0.008	0.152 ± 0.010	0.393 ± 0.014
	QuiltNet	0.141 ± 0.010	0.204 ± 0.011	0.428 ± 0.013	0.127 ± 0.009	0.188 ± 0.011	0.406 ± 0.013
	KEEP	0.329 ± 0.014	0.443 ± 0.014	0.764 ± 0.012	0.296 ± 0.013	0.409 ± 0.014	0.737 ± 0.013
PathPair	PLIP	0.026 ± 0.002	0.047 ± 0.002	0.134 ± 0.003	0.023 ± 0.002	0.038 ± 0.002	0.119 ± 0.003
	QuiltNet	0.046 ± 0.002	0.071 ± 0.003	0.195 ± 0.004	0.041 ± 0.002	0.065 ± 0.003	0.166 ± 0.004
	KEEP	0.136 ± 0.004	0.210 ± 0.004	0.443 ± 0.005	0.111 ± 0.003	0.168 ± 0.004	0.361 ± 0.005
WebPath	PLIP	0.031 ± 0.002	0.048 ± 0.002	0.140 ± 0.004	0.027 ± 0.002	0.045 ± 0.002	0.127 ± 0.003
	QuiltNet	0.046 ± 0.002	0.069 ± 0.002	0.181 ± 0.004	0.035 ± 0.002	0.057 ± 0.002	0.162 ± 0.004
	KEEP	0.134 ± 0.003	0.208 ± 0.004	0.431 ± 0.005	0.083 ± 0.003	0.132 ± 0.003	0.310 ± 0.005

Supplementary Table S7 | Zero-shot performance (weighted F1) of the tile classification task. Bold suggests the best performance.

Models	PLIP	QuiltNet	MI-Zero (Pub)	CONCH	KEEP(Ours)
BACH [1]	0.380 ± 0.025	0.386 ± 0.027	0.650 ± 0.029	0.606 ± 0.028	0.686 ± 0.026
NCT-CRC-HE-100K [28]	0.481 ± 0.002	0.543 ± 0.002	0.644 ± 0.002	0.556 ± 0.002	0.767 ± 0.002
CRC100K [29]	0.602 ± 0.006	0.553 ± 0.006	0.583 ± 0.006	0.590 ± 0.006	0.852 ± 0.004
LC25000 [6]	0.622 ± 0.003	0.775 ± 0.003	0.645 ± 0.003	0.527 ± 0.004	0.936 ± 0.002
RenalCell [7]	0.342 ± 0.003	0.429 ± 0.003	0.435 ± 0.003	0.324 ± 0.003	0.470 ± 0.003
SkinCancer [30]	0.435 ± 0.002	0.409 ± 0.002	0.426 ± 0.002	0.413 ± 0.002	0.658 ± 0.002
WSSS4LUAD [20]	0.261 ± 0.010	0.827 ± 0.007	0.644 ± 0.009	0.798 ± 0.007	0.809 ± 0.008
Osteo [2]	0.322 ± 0.015	0.585 ± 0.017	0.591 ± 0.016	0.785 ± 0.012	0.760 ± 0.012
ESCA-UKK [46]	0.612 ± 0.003	0.592 ± 0.003	0.380 ± 0.003	0.433 ± 0.003	0.523 ± 0.003
ESCA-WNS [46]	0.485 ± 0.002	0.415 ± 0.002	0.516 ± 0.002	0.501 ± 0.002	0.584 ± 0.002
ESCA-TCGA [46]	0.421 ± 0.003	0.332 ± 0.003	0.376 ± 0.003	0.497 ± 0.003	0.628 ± 0.003
ESCA-CHA [46]	0.493 ± 0.001	0.444 ± 0.001	0.396 ± 0.001	0.399 ± 0.001	0.483 ± 0.001
Breakhis [44]	0.087 ± 0.003	0.255 ± 0.005	0.366 ± 0.006	0.289 ± 0.005	0.337 ± 0.006
Chaoyang [56]	0.156 ± 0.008	0.118 ± 0.008	0.149 ± 0.008	0.108 ± 0.008	0.156 ± 0.008

Supplementary Table S8 | Ablation results of cancer region segmentation. Bold suggests the best performance.

Metric	Dataset	Simple Contrastive	KEEP
AUROC	CAMELYON16	0.972 ± 0.066	0.967 ± 0.088
	PANDA	0.602 ± 0.140	0.731 ± 0.139
	AGGC22	0.797 ± 0.114	0.891 ± 0.065
DICE	CAMELYON16	0.354 ± 0.348	0.361 ± 0.332
	PANDA	0.223 ± 0.205	0.334 ± 0.261
	AGGC22	0.347 ± 0.170	0.530 ± 0.179

Supplementary Table S9 | Ablation results of cancer detection. Bold suggests the best performance.

Metric	Dataset	Simple Contrastive	KEEP
AUROC	CPTAC-CM	0.993 ± 0.006	0.994 ± 0.006
	CPTAC-CCRCC	0.995 ± 0.004	0.999 ± 0.001
	CPTAC-PDA	0.895 ± 0.029	0.925 ± 0.021
	CPTAC-UCEC	0.996 ± 0.003	0.997 ± 0.002
	CPTAC-LSCC	0.983 ± 0.010	0.981 ± 0.010
	CPTAC-HNSCC	0.970 ± 0.015	0.974 ± 0.013
	CPTAC-LUAD	0.998 ± 0.002	0.999 ± 0.001
Sensitivity	CPTAC-CM	0.978 ± 0.059	0.988 ± 0.053
	CPTAC-CCRCC	0.981 ± 0.019	0.988 ± 0.013
	CPTAC-PDA	0.345 ± 0.159	0.591 ± 0.181
	CPTAC-UCEC	0.980 ± 0.034	0.983 ± 0.023
	CPTAC-LSCC	0.908 ± 0.099	0.882 ± 0.094
	CPTAC-HNSCC	0.861 ± 0.143	0.842 ± 0.161
	CPTAC-LUAD	0.992 ± 0.019	0.999 ± 0.004

Supplementary Table S10 | Ablation results of cancer subtyping. Performance comparison between simple contrastive (Contrastive-Top100) and KEEP with knowledge enhancement (KEEP-Top100). Top100 suggests the strategy of top-100 pooling, while Ratio denotes the subtype ratio strategy. Bold suggests the best performance.

Metric	Dataset	Contrast-Top100	KEEP-Top100	KEEP-Ratio
Balanced accuracy	TCGA-BRCA	0.740 ± 0.034	0.626 ± 0.027	0.773 ± 0.031
	TCGA-NSCLC	0.806 ± 0.029	0.859 ± 0.029	0.900 ± 0.025
	TCGA-RCC	0.916 ± 0.018	0.912 ± 0.018	0.925 ± 0.018
	TCGA-ESCA	0.899 ± 0.026	0.945 ± 0.020	0.977 ± 0.013
	TCGA-BRAIN	0.472 ± 0.024	0.555 ± 0.029	0.604 ± 0.030
	UBC-OCEAN	0.570 ± 0.032	0.656 ± 0.030	0.662 ± 0.031
	CPTAC-NSCLC	0.795 ± 0.032	0.761 ± 0.031	0.860 ± 0.028
	EBRAINS	0.283 ± 0.009	0.393 ± 0.013	0.455 ± 0.013
Weighted F1	TCGA-BRCA	0.731 ± 0.039	0.570 ± 0.050	0.765 ± 0.037
	TCGA-NSCLC	0.803 ± 0.032	0.859 ± 0.030	0.900 ± 0.025
	TCGA-RCC	0.916 ± 0.018	0.911 ± 0.019	0.925 ± 0.018
	TCGA-ESCA	0.898 ± 0.028	0.945 ± 0.020	0.977 ± 0.013
	TCGA-BRAIN	0.418 ± 0.038	0.524 ± 0.036	0.588 ± 0.034
	UBC-OCEAN	0.550 ± 0.040	0.636 ± 0.039	0.654 ± 0.037
	CPTAC-NSCLC	0.794 ± 0.033	0.753 ± 0.037	0.859 ± 0.028
	EBRAINS	0.205 ± 0.014	0.386 ± 0.017	0.419 ± 0.018

Supplementary Table S11 | Prompt templates used in this paper, consistent with CONCH [34], CLASSNAME is replaced by the names/synonyms of classes.

CLASSNAME.
 a photomicrograph showing CLASSNAME.
 a photomicrograph of CLASSNAME.
 an image of CLASSNAME.
 an image showing CLASSNAME.
 an example of CLASSNAME.
 CLASSNAME is shown.
 this is CLASSNAME.
 there is CLASSNAME.
 a histopathological image showing CLASSNAME.
 a histopathological image of CLASSNAME.
 a histopathological photograph of CLASSNAME.
 a histopathological photograph showing CLASSNAME.
 shows CLASSNAME.
 presence of CLASSNAME.
 CLASSNAME is present.
 an H&E stained image of CLASSNAME.
 an H&E stained image showing CLASSNAME.
 an H&E image showing CLASSNAME.
 an H&E image of CLASSNAME.
 CLASSNAME, H&E stain.
 CLASSNAME, H&E.

Supplementary Table S12 | Class names of WSI datasets in cancer region segmentation tasks.

Dataset	Class	Names/Synonyms
CAMELYON16	Tumor	'tumor tissue', 'tumor epithelial tissue', 'cancerous tissue', 'breast tumor tissue', 'breast tumor epithelial tissue', 'breast cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal breast tissue', 'breast non-cancerous tissue', 'benign breast tissue', 'benign tissue'
PANDA	Tumor	'tumor tissue', 'tumor epithelial tissue', 'cancerous tissue', 'prostate tumor tissue', 'prostate tumor epithelial tissue', 'prostate cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal prostate tissue', 'prostate non-cancerous tissue', 'benign prostate tissue', 'benign tissue'
AGGC22	Tumor	'tumor tissue', 'tumor epithelial tissue', 'cancerous tissue', 'prostate tumor tissue', 'prostate tumor epithelial tissue', 'prostate cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal prostate tissue', 'prostate non-cancerous tissue', 'benign prostate tissue', 'benign tissue'

Supplementary Table S13 | Class names of WSI datasets in cancer detection tasks.

Dataset	Class	Names/Synonyms
CPTAC-CM	Tumor	'cutaneous melanoma', 'skin cancer, melanoma', 'malignant cutaneous melanoma', 'malignant melanoma of skin', 'melanoma, cutaneous', 'malignant melanoma', 'skin melanoma tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal skin tissue', 'normal cutaneous tissue', 'cutaneous non-cancerous tissue', 'benign cutaneous tissue', 'benign skin tissue', 'benign tissue'
CPTAC-CCRCC	Tumor	'chromophobe renal cell carcinoma', 'renal cell carcinoma, chromophobe type', 'renal cell carcinoma of the chromophobe type', 'chromophobe cell carcinoma of the kidney', 'renal cell carcinoma, chromophobe cell', 'chromophobe carcinoma of the kidney', 'chromophobe renal cell cancer', 'kidney tumor tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal renal tissue', 'normal kidney tissue', 'renal non-cancerous tissue', 'benign renal tissue', 'benign kidney tissue', 'benign tissue'
CPTAC-PDA	Tumor	'pancreatic ductal adenocarcinoma', 'ductal adenocarcinoma of the pancreas', 'pancreas cancer, duct cell adenocarcinoma', 'pancreatic cancer, duct cell adenocarcinoma', 'pancreatic cancer', 'pancreatic tumor tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal pancreatic tissue', 'pancreatic non-cancerous tissue', 'benign pancreatic tissue', 'benign tissue'
CPTAC-UCEC	Tumor	'uterine corpus endometrial carcinoma', 'endometrial carcinoma', 'uterine endometrial cancer, carcinoma', 'endometrial cancer', 'carcinoma of the endometrium', 'endometrial cancerous tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal endometrial tissue', 'normal uterine tissue', 'endometrial non-cancerous tissue', 'benign endometrial tissue', 'benign uterine tissue', 'benign tissue'
CPTAC-LSCC	Tumor	'lung squamous cell carcinoma', 'squamous cell carcinoma of the lung', 'squamous cell lung cancer', 'epidermoid carcinoma of lung', 'epidermoid cell lung carcinoma', 'squamous cell lung carcinoma', 'lung cancerous tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal lung tissue', 'lung non-cancerous tissue', 'benign lung tissue', 'benign tissue'
CPTAC-HNSCC	Tumor	'head and neck cancer', 'head and neck tumor', 'head and neck squamous cell carcinoma', 'squamous cell carcinoma of head and neck', 'neoplasm of the head and neck', 'head and neck cancerous tissue', 'squamous cell carcinoma', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal head and neck tissue', 'head and neck non-cancerous tissue', 'benign head and neck tissue', 'benign tissue'
CPTAC-LUAD	Tumor	'lung adenocarcinoma', 'adenocarcinoma of the lung', 'nonsmall cell lung adenocarcinoma', 'non-oat cell adenocarcinoma of the lung', 'pulmonary adenocarcinoma', 'lung cancer, adenocarcinoma', 'lung cancerous tissue', 'tumor tissue', 'cancerous tissue'
	Normal	'normal tissue', 'non-cancerous tissue', 'normal lung tissue', 'lung non-cancerous tissue', 'benign lung tissue', 'benign tissue'

Supplementary Table S14 | Class names of WSI datasets in cancer subtyping tasks.

Dataset	Class	Names/Synonyms
TCGA-BRCA	IDC	'breast invasive ductal carcinoma', 'invasive ductal carcinoma of the breast'
	ILC	'breast invasive lobular carcinoma', 'invasive lobular carcinoma of the breast'
	Normal	'normal breast tissue', 'breast normal tissue', 'breast non-cancerous tissue'
TCGA-NSCLC	LUAD	'lung adenocarcinoma', 'adenocarcinoma of the lung'
	LUSC	'lung squamous cell carcinoma', 'squamous cell carcinoma of the lung'
	Normal	'normal lung tissue', 'lung normal tissue', 'lung non-cancerous tissue'
TCGA-RCC	CCRCC	'clear cell renal cell carcinoma', 'renal cell carcinoma of the clear cell type', 'clear cell RCC'
	PRCC	'papillary renal cell carcinoma', 'renal cell carcinoma of the papillary type', 'papillary RCC'
	CHRCC	'chromophobe renal cell carcinoma', 'renal cell carcinoma of the chromophobe type', 'chromophobe RCC'
	Normal	'normal renal tissue', 'renal normal tissue', 'renal non-cancerous tissue'
TCGA-ESCA	Adenocarcinoma, NOS	'esophagus adenocarcinoma', 'adenocarcinoma of the esophagus'
	Squamous cell carcinoma, NOS	'esophagus squamous cell carcinoma', 'squamous cell carcinoma of the esophagus'
	Normal	'normal esophagus tissue', 'esophagus normal tissue', 'esophagus non-cancerous tissue'
TCGA-BRAIN	Glioblastoma	'brain glioblastoma', 'glioblastoma of the brain'
	Astrocytoma, NOS	'brain astrocytoma', 'astrocytoma of the brain'
	Oligodendroglioma, NOS	'brain oligodendroglioma', 'oligodendroglioma of the brain'
	Normal	'normal brain tissue', 'brain normal tissue', 'brain non-cancerous tissue'
UBC-OCEAN	CC	'ovarian clear cell carcinoma', 'clear cell carcinoma of the ovary'
	EC	'ovary endometrioid carcinoma', 'endometrioid carcinoma of the ovary'
	HGSC	'high-grade ovary serous carcinoma', 'high-grade serous carcinoma of the ovary'
	LGSC	'low-grade ovary serous carcinoma', 'low-grade serous carcinoma of the ovary'
	MC	'ovarian mucinous carcinoma', 'mucinous carcinoma of the ovary'
	Normal	'normal ovarian tissue', 'ovary normal tissue', 'ovary non-cancerous tissue'
CPTAC-NSCLC	LUAD	'lung adenocarcinoma', 'adenocarcinoma of the lung'
	LUSC	'lung squamous cell carcinoma', 'squamous cell carcinoma of the lung'
	Normal	'normal lung tissue', 'lung normal tissue', 'lung non-cancerous tissue'

Supplementary Table S15 | Class names of EBRAINS dataset, consistent with CONCH [34].

Dataset	Subtype	Names/Synonyms
EBRAINS	Glioblastoma, IDH-wildtype	'glioblastoma, IDH-wildtype', 'glioblastoma without IDH mutation', 'glioblastoma with retained IDH', 'glioblastoma, IDH retained'
	Transitional meningioma	'transitional meningioma', 'meningioma, transitional type', 'meningioma of transitional type', 'meningioma, transitional'
	Anaplastic meningioma	'anaplastic meningioma', 'meningioma, anaplastic type', 'meningioma of anaplastic type', 'meningioma, anaplastic'
	Pituitary adenoma	'pituitary adenoma', 'adenoma of the pituitary gland', 'pituitary gland adenoma', 'pituitary neuroendocrine tumor', 'neuroendocrine tumor of the pituitary', 'neuroendocrine tumor of the pituitary gland'
	Oligodendroglioma, IDH-mutant and 1p/19q codeleted	'oligodendroglioma, IDH-mutant and 1p/19q codeleted', 'oligodendroglioma', 'oligodendroglioma with IDH mutation and 1p/19q codeletion'
	Haemangioma	'hemangioma', 'haemangioma of the CNS', 'hemangioma of the CNS', 'haemangioma of the central nervous system', 'hemangioma of the central nervous system'
	Ganglioglioma	'gangliocytoma', 'glioneuronal tumor', 'circumscribed glioneuronal tumor'
	Schwannoma	'schwannoma', 'Antoni A', 'Antoni B', 'neurilemoma'
	Anaplastic oligodendroglioma, IDH-mutant, 1p/19q codeleted	'anaplastic oligodendroglioma, IDH-mutant and 1p/19q codeleted', 'anaplastic oligodendroglioma', 'anaplastic oligodendroglioma with IDH mutation and 1p/19q codeletion'
	Anaplastic astrocytoma, IDH-wildtype	'anaplastic astrocytoma, IDH-wildtype', 'anaplastic astrocytoma without IDH mutation', 'anaplastic astrocytoma, IDH retained', 'anaplastic astrocytoma with retained IDH'
	Pilocytic astrocytoma	'pilocytic astrocytoma', 'juvenile pilocytic astrocytoma', 'spongioblastoma', 'pilomyxoid astrocytoma'
	Angiomatous meningioma	'angiomatous meningioma', 'meningioma, angiomatous type', 'meningioma of angiomatous type', 'meningioma, angiomatous'
	Haemangioblastoma	'haemangioblastoma', 'capillary hemangioblastoma', 'lindau tumor', 'angioblastoma'
	Gliosarcoma	'gliosarcoma', 'gliosarcoma variant of glioblastoma'
	Adamantinomatous craniopharyngioma	'adamantinomatous craniopharyngioma', 'craniopharyngioma'
	Anaplastic astrocytoma, IDH-mutant	'anaplastic astrocytoma, IDH-mutant', 'anaplastic astrocytoma with IDH mutation', 'anaplastic astrocytoma with mutant IDH', 'anaplastic astrocytoma with mutated IDH'
	Ependymoma	'ependymoma', 'subependymoma', 'myxopapillary ependymoma'
	Anaplastic ependymoma	'anaplastic ependymoma', 'ependymoma, anaplastic', 'ependymoma, anaplastic type'
	Glioblastoma, IDH-mutant	'glioblastoma, IDH-mutant', 'glioblastoma with IDH mutation', 'glioblastoma with mutant IDH', 'glioblastoma with mutated IDH'
	Atypical meningioma	'atypical meningioma', 'meningioma, atypical type', 'meningioma of atypical type', 'meningioma, atypical'
	Metastatic tumours	'metastatic tumors', 'metastases to the brain', 'metastatic tumors to the brain', 'brain metastases', 'brain metastatic tumors'
	Meningothelial meningioma	'meningothelial meningioma', 'meningioma, meningothelial type', 'meningioma of meningothelial type', 'meningioma, meningothelial'
	Langerhans cell histiocytosis	'langerhans cell histiocytosis', 'histiocytosis X', 'eosinophilic granuloma', 'Hand-Schüller-Christian disease', 'Hashimoto-Pritzker disease', 'Letterer-Siwe disease'
	Diffuse large B-cell lymphoma of the CNS	'diffuse large B-cell lymphoma of the CNS', 'DLBCL', 'DLBCL of the CNS', 'DLBCL of the central nervous system'
	Diffuse astrocytoma, IDH-mutant	'diffuse astrocytoma, IDH-mutant', 'diffuse astrocytoma with IDH mutation', 'diffuse astrocytoma with mutant IDH', 'diffuse astrocytoma with mutated IDH'
	Secretory meningioma	'secretory meningioma', 'meningioma, secretory type', 'meningioma of secretory type', 'meningioma, secretory'
	Haemangiopericytoma	'haemangiopericytoma', 'solitary fibrous tumor', 'hemangiopericytoma', 'angioblastic meningioma'
	Fibrous meningioma	'fibrous meningioma', 'meningioma, fibrous type', 'meningioma of fibrous type', 'meningioma, fibrous'
	Lipoma	'lipoma', 'CNS lipoma', 'lipoma of the CNS', 'lipoma of the central nervous system'
	Medulloblastoma, non-WNT/non-SHH	'medulloblastoma, non-WNT/non-SHH', 'medulloblastoma', 'medulloblastoma group 3', 'medulloblastoma group 4'
Normal	'normal brain tissue', 'brain normal tissue', 'brain non-cancerous tissue'	

Supplementary Table S16 | Class names of tile image datasets.

Dataset	Class Names/Synonyms
BACH	'Benign': 'breast non-malignant benign tissue', 'InSitu': 'breast malignant in-situ carcinoma', 'Invasive': 'breast malignant invasive carcinoma', 'Normal': 'normal breast tissue'
NCT-CRC-HE-100K, CRC100K	'ADI': 'adipose', 'BACK': 'background', 'DEB': 'debris', 'LYM': 'lymphocytes', 'MUC': 'mucus', 'MUS': 'smooth muscle', 'NORM': 'normal colon mucosa', 'STR': 'cancer-associated stroma', 'TUM': 'colorectal adenocarcinoma epithelium'
LC25000	'lung_aca': 'lung adenocarcinoma', 'lung_n': 'benign lung', 'lung_scc': 'lung squamous cell carcinoma', 'colon_aca': 'colon adenocarcinoma', 'colon_n': 'benign colon'
RenalCell	'blood': 'red blood cells', 'cancer': 'renal cancer', 'normal': 'non-tumor', 'other': 'torn adipose necrotic tissue', 'stroma': 'muscle fibrous stroma blood vessels'
RenalCell	'blood': 'red blood cells', 'cancer': 'renal cancer', 'normal': 'non-tumor', 'other': 'torn adipose necrotic tissue', 'stroma': 'muscle fibrous stroma blood vessels'
SkinCancer	'necrosis': 'necrosis', 'skeletal': 'skeletal muscle', 'sweatglands': 'eccrine sweat glands', 'vessel': 'vessels', 'elastosis': 'elastosis', 'chondraltissue': 'chondral tissue', 'hairfollicle': 'hair follicle', 'epidermis': 'epidermis', 'nerves': 'nerves', 'subcutis': 'subcutis', 'dermis': 'dermis', 'sebaceousglands': 'sebaceous', 'sqcc': 'squamous-cell carcinoma', 'melanoma': 'melanoma in-situ', 'bcc': 'basal-cell carcinoma', 'naevus': 'naevus'
WSSS4LUAD	'normal': 'non-tumor', 'stroma': 'tumor-associated stroma', 'tumor': 'tumor tissue'
Osteo	'Normal non-tumor': 'normal non-tumor', 'Necrotic': 'necrosis', 'Tumor': 'viable tumor'
ESCA-UKK, ESCA-WNS, ESCA-TCGA, ESCA-CHA	'TUMOR': 'vital tumor tissue', 'REGR_TU': 'regression areas', 'SH_OES': 'oesophageal mucosa', 'SH_MAG': 'gastric mucosa', 'LAM_PROP': 'lamina propria mucosae', 'SUBMUC': 'submucosa', 'SUB_GL': 'submucosal glands', 'MUSC_MUC': 'lamina muscularis mucosae', 'MUSC_PROP': 'muscularis propria', 'ADVENT': 'adventitial tissue', 'ULCUS': 'areas of ulceration'
Breakhis	'adenosis': 'adenosis', 'fibroadenoma': 'fibroadenoma', 'phyllodes_tumor': 'phyllodes tumor', 'tubular_adenoma': 'tubular adenoma', 'ductal_carcinoma': 'ductal carcinoma', 'lobular_carcinoma': 'lobular carcinoma', 'mucinous_carcinoma': 'mucinous carcinoma', 'papillary_carcinoma': 'papillary carcinoma'
Chaoyang	'normal': 'normal', 'serrated': 'serrated', 'adenocarcinoma': 'adenocarcinoma', 'adenoma': 'adenoma'