

Multimodal Whole Slide Foundation Model for Pathology

Tong Ding^{1,2,3,4,*}, Sophia J. Wagner^{1,5,6,*}, Andrew H. Song^{1,2,3,*}, Richard J. Chen^{1,2,3,*}, Ming Y. Lu^{1,2,3,7}, Andrew Zhang^{1,2,3,8,+}, Anurag J. Vaidya^{1,2,3,8,+}, Guillaume Jaume^{1,2,3,+}, Muhammad Shaban^{1,2,3}, Ahrong Kim^{1,9}, Drew F.K. Williamson¹⁰, Bowen Chen^{1,2,3}, Cristina Almagro-Perez^{1,2,3,8}, Paul Doucet^{1,2,3}, Sharifa Sahai^{1,2,3,12}, Chengkuan Chen^{1,2,3}, Daisuke Komura¹³, Akihiro Kawabe¹³, Shumpei Ishikawa^{13,14}, Georg Gerber¹, Tingying Peng^{5,6}, Long Phi Le^{1,8,†}, Faisal Mahmood^{1,2,3,11,†}

¹*Department of Pathology, Mass General Brigham, Harvard Medical School, Boston, MA, USA*

²*Data Science Program, Dana-Farber Cancer Institute, Boston, MA, USA*

³*Cancer Program, Broad Institute of Harvard and MIT, Cambridge, MA, USA*

⁴*John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, USA*

⁵*Helmholtz Munich – German Research Center for Environment and Health, Munich, Germany*

⁶*School of Computation, Information and Technology, Technical University of Munich, Munich, Germany*

⁷*Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA*

⁸*Harvard-MIT Division of Health Sciences and Technology, Massachusetts Institute of Technology, Cambridge, MA, USA*

⁹*Department of Pathology, Pusan National University, Busan, South Korea*

¹⁰*Department of Pathology and Laboratory Medicine, Emory University School of Medicine, Atlanta, GA, USA*

¹¹*Harvard Data Science Initiative, Harvard University, Cambridge, MA, USA*

¹²*Department of Systems Biology, Harvard Medical School, Boston, MA, USA*

¹³*Department of Preventive Medicine, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan*

¹⁴*Division of Pathology, National Cancer Center Exploratory Oncology Research & Clinical Trial Center, Chiba, Japan*

* *Contributed equally (Co-first)*

+ *Contributed equally (Co-third)*

† *Co-senior authors*

Lead Contact:

Faisal Mahmood (FaisalMahmood@bwh.harvard.edu)

Abstract

The field of computational pathology has been transformed with recent advances in foundation models that encode histopathology region-of-interests (ROIs) into versatile and transferable feature representations via self-supervised learning (SSL). However, translating these advancements to address complex clinical challenges at the patient and slide level remains constrained by limited clinical data in disease-specific cohorts, especially for rare clinical conditions. We propose TITAN, a multimodal whole slide foundation model pretrained using 335,645 WSIs via visual self-supervised learning and vision-language alignment with corresponding pathology reports and 423,122 synthetic captions generated from a multimodal generative AI copilot for pathology. Without any finetuning or requiring clinical labels, TITAN can extract general-purpose slide representations and generate pathology reports that generalize to resource-limited clinical scenarios such as rare disease retrieval and cancer prognosis. We evaluate TITAN on diverse clinical tasks and find that TITAN outperforms both ROI and slide foundation models across machine learning settings such as linear probing, few-shot and zero-shot classification, rare cancer retrieval and cross-modal retrieval, and pathology report generation. The model is publicly accessible at <https://github.com/mahmoodlab/TITAN>

Introduction

Foundation models are transforming computational pathology by accelerating the development of AI tools for diagnosis, prognosis, and biomarker prediction from digitized tissue sections¹. Developed using self-supervised learning (SSL) on millions of histology image patches (or regions of interests), these models capture morphological patterns in histology patch embeddings, such as tissue organization and cellular structure²⁻¹⁷. These representations serve as a “foundation” for models that predict clinical endpoints from whole-slide images (WSIs), such as diagnosis or biomarker status¹⁸⁻³⁸. However, translating the capabilities of current patch-based foundation models to address patient- and slide-level clinical challenges still remains complex due to the immense scale of gigapixel WSIs, compounded by the small size of patient cohorts in real-world evidence^{39,40}, posing challenges for disease-specific AI model development⁴¹. As an example, in rare diseases with limited training data⁴²⁻⁴⁴, developing effective predictive models is difficult since the slide encoder—which generates slide-level predictions from patch embeddings—still needs to be trained from scratch^{10,45}. Similarly, given a diagnostically challenging patient slide, retrieving a similar slide via slide search^{5,46-53} or pathology reports through cross-modal report search^{10,54-56} typically requires specialized algorithms to bridge the gap between patch and slide embeddings, introducing hurdles towards clinical adoption.

To overcome these limitations, new types of foundation models have recently been proposed for encoding entire WSIs into slide-level general-purpose feature representations⁵⁷⁻⁶⁶. Instead of training an additional model on top of patch embeddings from scratch^{34,45,67-71}, these whole slide representation models can be pretrained to distill pathology-specific knowledge from large WSI collections, simplifying clinical endpoint prediction. The outstanding challenge then becomes developing whole slide foundation models that faithfully encode the tissue microenvironment based on a set of patch embeddings while also handling arbitrarily large WSIs. Although relatively underexplored, slide-level self-supervision can be performed with vision-only pretraining, either through masked image reconstruction⁵⁸ or intra-slide contrastive learning^{59,60,72}, or in a multimodal fashion involving pathology reports, bulk transcriptomics, or immunohistochemistry^{61-64,66,73}. Furthermore, long-range context modeling can either be neglected, essentially treating a WSI as a bag of independent features^{59,62-64,74}, or explicitly modeled using Transformers^{57,58,60,61}. With efforts to learn general-purpose slide representations intensifying, we believe that adapting successful patch-level recipes to the entire WSI would lead to powerful general-purpose slide representations.

Despite their widespread application potential, previous works on pretraining slide foundation models have several shortcomings. First, these models are predominantly pretrained using vision-only modeling^{57,59,60}, which neglects not only rich supervisory signals found in pathology reports, but also precludes multimodal capabilities such as zero-shot visual-language understanding and cross-modal retrieval – which is a fundamental hallmark in foundation models^{75,76}. Second, whereas current patch foundation models are trained with millions

of histology image patches, slide foundation models are developed with orders of magnitude fewer samples and limited optimization of self-supervised learning recipes, leading to slide representations with restricted generalization capability^{58,62,73,74}. Even with multimodal techniques such as vision-language pretraining that augment the pretraining dataset with pathology reports, current slide foundation models still require end-to-end training or finetuning and lack the capability of learning transferable slide representations for challenging clinical scenarios^{58,73,74}. Finally, the current models are nascent in transforming pathology AI model development due to their limited evaluations in diagnostically relevant settings such as few-shot learning or slide retrieval.

Here, we introduce Transformer-based pathology **Image and Text Alignment Network (TITAN)**, a multimodal whole-slide vision-language model designed for general-purpose slide representation learning in histopathology. Building on the success of knowledge distillation and masked-image modeling^{77,78} for patch encoder pretraining^{21,22}, TITAN introduces a novel paradigm that leverages millions of high-resolution regions-of-interests (ROIs at $8, 192 \times 8, 192$ pixels) for large-scale, resolution-agnostic pretraining and scalable WSI encoding. Trained using 336K WSIs across 20 organ types, vision-only TITAN produces general-purpose slide representations that can readily be applied to slide-level tasks such as cancer subtyping, biomarker prediction, outcome prognosis, and slide retrieval tasks, outperforming supervised baselines and existing multimodal slide foundation models. To augment TITAN with language capabilities, we further finetune by contrasting with 423K synthetic fine-grained ROI-captions generated with PathChat⁷⁹, a multimodal generative AI copilot for pathology, and with 183K pathology reports at slide level. By leveraging free-text morphological descriptions, TITAN gains the ability to generate pathology reports, perform zero-shot classification, and enable cross-modal retrieval between histology slides and clinical reports. Pretraining TITAN on an extensive repository of multimodal pathology data unlocks new levels of performance compared to existing slide foundation models, particularly in low data regimes, language-guided zero-shot classification, and rare cancer retrieval. Additionally, we show the utility of pretraining with synthetic fine-grained morphological descriptions for the first time, hinting at the scaling potential of TITAN pretraining with synthetic data⁸⁰⁻⁸². Through comprehensive evaluation across a large range of clinical tasks, including the first application to rare cancer retrieval across 43 rare cancer types, we demonstrate the efficacy of our vision-language pretraining approach, showcasing the general-purpose capability of our slide foundation model.

Results

Scaling self-supervised learning from histology patches to whole slide images

TITAN is a Vision Transformer (ViT)⁸³ that creates a general-purpose slide representation readily deployable in diverse clinical settings. It is pretrained on an internal dataset (termed **Mass-340K**) consisting of 335,645 WSIs and 182,862 medical reports (**Figure 1A**). To ensure the diversity of the pretraining dataset, which has

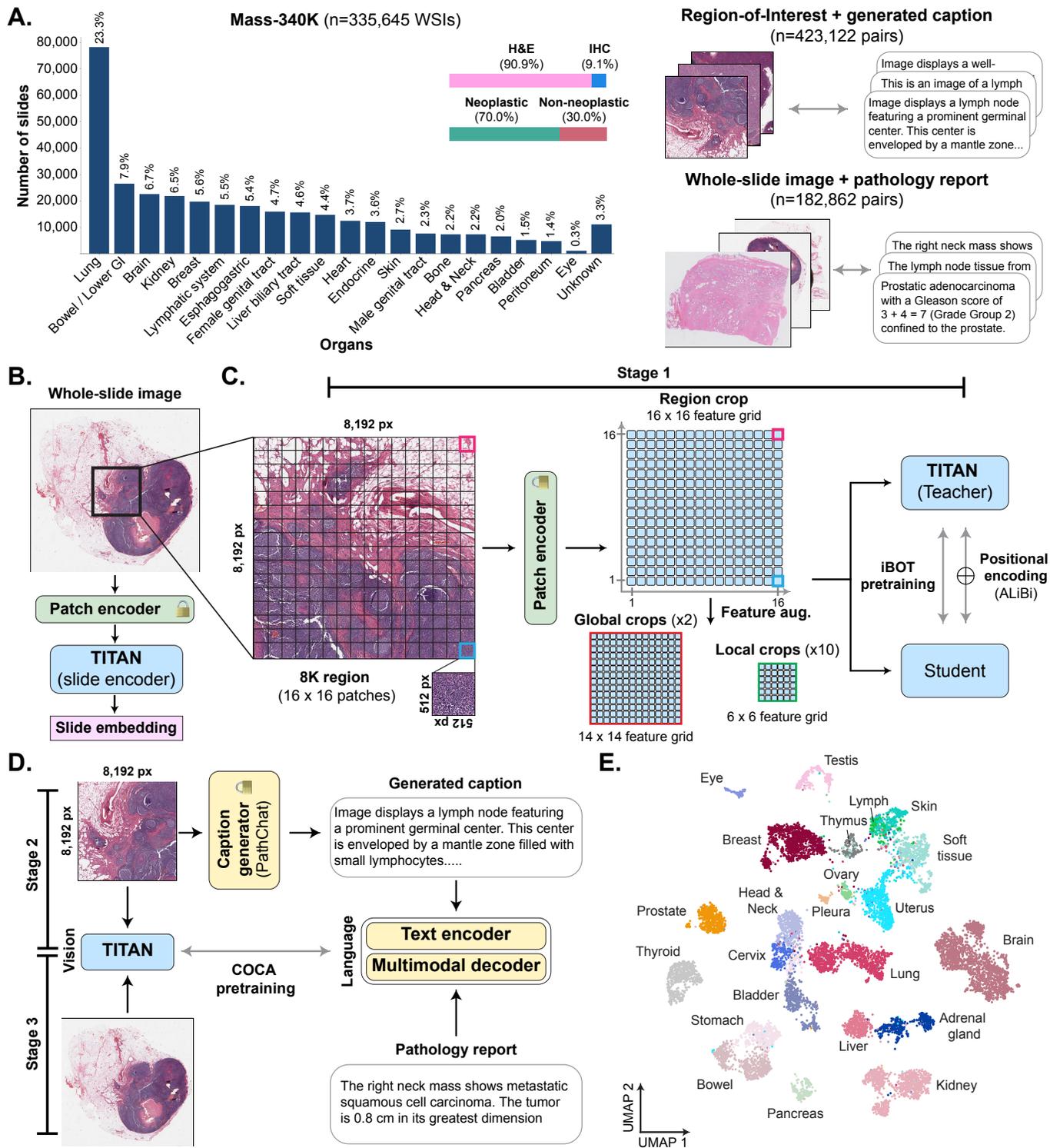


Figure 1: Overview of TITAN. (a) Tissue site distribution of Mass-340K used for TITAN_V pretraining (Stage 1). Mass-340K includes 335,645 WSIs across 20 organs with a mix of hematoxylin-and-eosin-stained (90.9%) and immunohistochemistry-stained tissue sections (9.1%) or a mix of neoplastic (70.0%) and non-neoplastic tissue sections (30.0%). TITAN pretraining (Stages 2 and 3) uses a subset of Mass-340K with paired captions and medical reports. (b–d) Block diagram of TITAN_V pretraining. (b) TITAN uses a Vision Transformer to encode a WSI into a slide embedding. (c) TITAN_V (Stage 1) is pretrained using self-supervised learning with student–teacher knowledge distillation (d) TITAN (Stage 2 and 3) is pretrained using vision-language modeling, first by aligning the slide embedding with synthetic captions (Stage 2) and then with medical reports (Stage 3). (e) UMAP visualization of TCGA slide embeddings obtained with TITAN, color-coded by organ. UMAP: uniform manifold approximation and projection.

proven to be a key factor in successful patch encoders²¹, Mass-340K is distributed across 20 organs, across different stains (Hematoxylin-and-eosin 90.9% and immunohistochemistry 9.1%), and across neoplastic and non-neoplastic tissue (70.0% and 30.0%, respectively). The pretraining strategy consists of three distinct stages to ensure that the resulting slide-level representations capture histomorphological semantics both at the ROI-level ($4\times 4\text{mm}^2$) and at the WSI-level with the help of visual and language supervisory signals: **Stage 1** vision-only unimodal pretraining with Mass-340K on ROI crops (**Figure 1C**), **Stage 2** cross-modal alignment of generated morphological descriptions at ROI-level (423K pairs of $8\text{K}\times 8\text{K}$ ROIs and captions), and **Stage 3** cross-modal alignment at WSI-level (183K pairs of WSIs and clinical reports, **Figure 1D**). For ease of notation, we refer to the model pretrained with vision-only in Stage 1 as TITAN_V and to the full model after all three stages of pretraining as TITAN. A detailed description of the pretraining dataset can be found in **Online Methods** section **Large-scale pretraining datasets**.

The cornerstone of our approach is emulating the patch encoder designed for input patch images at the slide level. Instead of using tokens from a partitioned image patch, the slide encoder takes a sequence of patch features encoded by powerful histology patch encoders^{4,7-14}. Consequently, all of TITAN pretraining stages occur in the embedding space based on pre-extracted patch features, with the patch encoder assuming the role of the “patch embedding layer” in a conventional ViT (**Figure 1B**). To preserve the spatial context of each patch and consequently enable the use of positional encoding in the embedding space, the patch features are spatially arranged in a 2D feature grid replicating the positions of the corresponding patches within the tissue (**Figure 1C**). Following the success of masked image modeling and knowledge distillation in patch encoders²¹, we apply the iBOT⁷⁷ framework for vision-only pretraining of TITAN. The 2D feature grid setup allows us to directly apply student-teacher knowledge distillation approaches which typically require square crop inputs.

While the conceptual transition to slide-level is simple, this presents a new set of model design and pretraining challenges that precludes clinical translation: 1) Handling long and variable input sequences ($> 10^4$ tokens at slide-level vs. 196 to 256 tokens at the patch-level), 2) creating multiple views of one sample for self-supervised learning, and 3) ambiguity over positional encoding schemes that capture local and global context in the tissue microenvironment. First, to tame the computational complexity caused by long input sequence, we construct the input embedding space by dividing each WSI into non-overlapping patches of 512×512 pixels at $20\times$ magnification, followed by extraction of 768-dimensional features for each patch with the extended version of CONCH¹⁰, CONCHv1.5. By increasing the patch size from widely-used 256×256 pixels, we effectively reduce the sequence length by four without impacting the representation quality due to higher resolution patch input, leveraging the robustness of the patch-level foundation models in generalizing to higher resolutions^{9,10,78}. To address the issue of large and irregular-shaped WSIs, we create views of a WSI by sampling a smaller square crop of features (**Figure 1C**). Specifically, at each epoch for a given WSI, a *region crop* of 16×16 features covering a region of $8,192\times 8,192$ pixels is randomly sampled from the WSI

feature grid. From this region crop, two random *global* (14×14) and ten *local* (6×6) crops are sampled for the iBOT training. We augment these feature crops further with vertical and horizontal flipping, followed by posterize feature augmentation⁸⁴. Finally, to ensure that the limited context pretraining translates to slide-level tasks, we use attention with linear bias (ALiBi) for long context extrapolation of TITAN at inference time⁸⁵. Originally proposed for long-context inference in large language models, we extended ALiBi to 2D, where the linear bias is based on the relative Euclidean distance between features in the feature grid, which reflects the actual distances between patches in the tissue. More details of the pretraining dataset and training strategy can be found in **Online Methods** section **Vision-only pretraining dataset** and **Unimodal visual pretraining**, respectively.

To equip our model with language capabilities, we implement two additional multi-resolution pretraining strategies (Stages 2 and 3) using a subset of WSIs in Mass-340K (**Figure 1D**). This is based on the observation that language descriptions exist at multiple morphological scales, from fine-grained descriptions in pathologist annotations or textbooks at the patch- or region-level (Stage 2) to high-level descriptions in pathology reports at the slide-level (Stage 3). For both stages, we use contrastive captioners (CoCa)⁸⁶ as the pretraining strategy that aligns ROI and slide representations with the corresponding captions and reports, while generating accurate descriptions at ROI-level or reports at slide-level, respectively. The slide encoder (weights initialized with TITAN_V), the text encoder, and the multimodal decoder are all finetuned as part of the pretraining. In Stage 2, we pretrain TITAN_V with 423,122 pairs of $8K \times 8K$ ROIs and synthetic captions generated by the vision-language copilot PathChat⁷⁹. In Stage 3, we further pretrain the model with 182,862 pairs of WSIs and corresponding pathology reports, resulting in our final model TITAN. To diversify the captions and reports with data augmentation, we rewrite the text with a locally deployed large language model (LLM)⁸⁷ and select randomly between several versions for vision-language alignment. A detailed description of the vision-language pretraining dataset and training strategy can be found in **Online Methods** sections **Large-scale pretraining datasets** and **Vision-language continual pretraining**, respectively.

TITAN improves region and slide-level diagnostic capabilities

We begin by evaluating TITAN, TITAN_V , and existing slide encoders on a large set of diverse slide-level tasks, including morphological subtyping and molecular classification. Following the standard practice in self-supervised learning^{78,88}, we employ linear probing by fitting a linear model for classification (logistic regression) on frozen slide embeddings. Specifically, we use the linear weights estimated with the ℓ_2 -regularization parameter tuned on a validation set for evaluating the test performance. For tasks with multiple cohorts available, we perform cross-validation on one cohort, e.g., from TCGA^{89,90}, and use the remaining cohorts, e.g., from CPTAC^{91,92} or DHMC^{93,94}, as an external test cohort. As baselines, we evaluate the recent vision-language slide foundation models with model weights available, namely PRISM⁶², GigaPath⁶³,

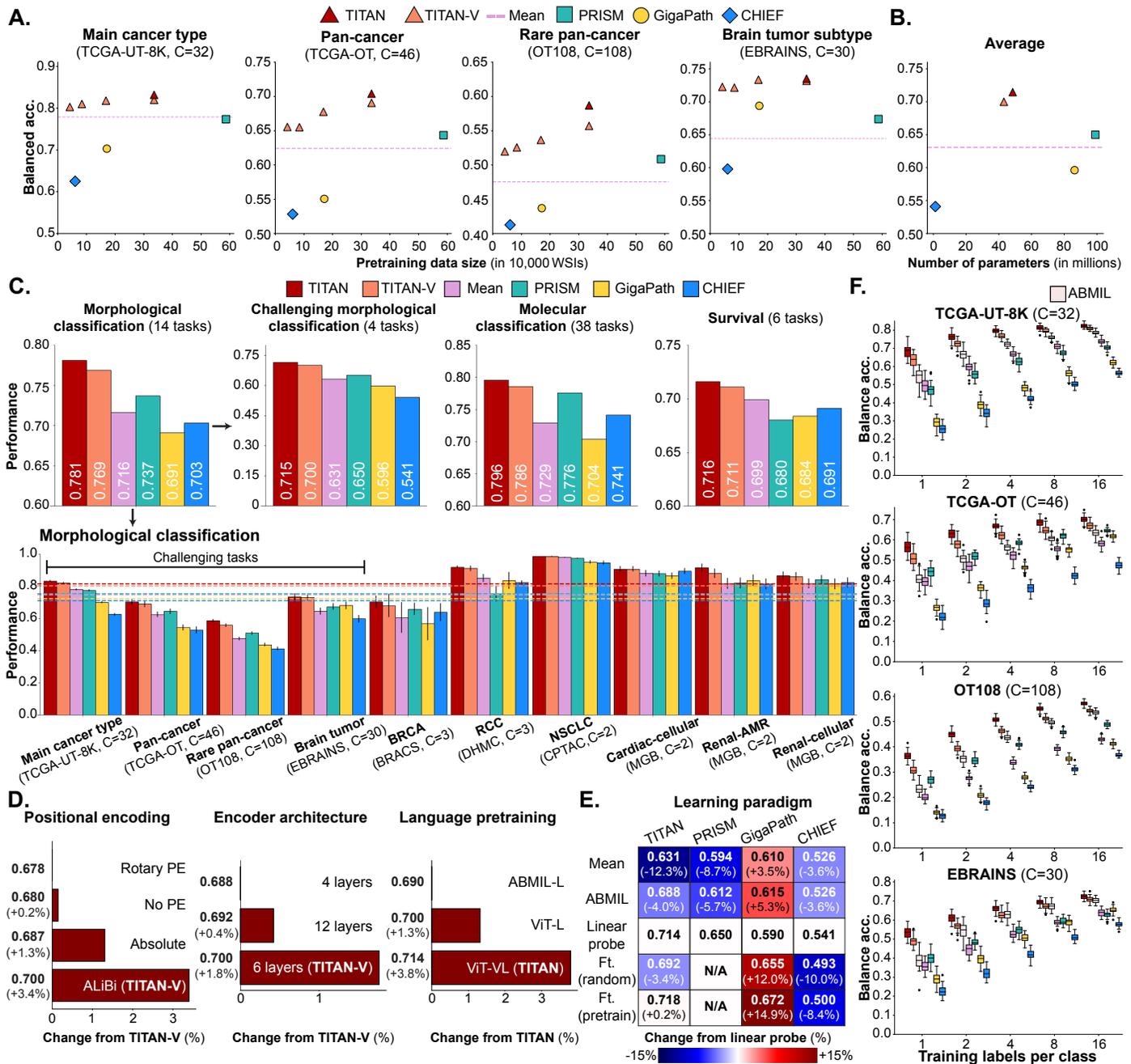


Figure 2: TITAN evaluation. (a) Impact of pretraining data size on TITAN_V and baselines across four challenging subtyping tasks (TCGA-UT-8K, TCGA-OT, OT108 and EBRAINS). TITAN_V is pretrained with 12.5%, 25%, 50%, and 100% of Mass-340K. (b) The average performance of the four tasks against the number of parameters for each baseline. (c) Linear probe evaluation of TITAN and baselines on morphological classification (all and challenging subset), molecular status, and survival prediction tasks. The mean uses the same patch encoder as TITAN (CONCHv1.5). Multi-class tasks are evaluated with balanced accuracy, binary tasks with AUROC, and survival tasks with concordance index. For external cohorts (DHMC, CPTAC), the classifier is trained on the corresponding TCGA cohort. All error bars represent standard deviations based on bootstrapping. (d) Ablation study comparing the impact of positional encoding, number of Transformer layers, and inclusion of vision-pretraining stage. The performance is averaged across the four subtyping tasks. (e) Change in performance of TITAN and baselines averaged across the four subtyping tasks for different learning paradigms. For mean pooling and ABMIL, the respective patch encoder for each framework is used. (f) Linear probe few-shot performance @ K shots, with $K \in \{1, 2, 4, 8, 16\}$, comparing baselines and ABMIL with CONCHv1.5. For each setting, 50 runs were performed. Whiskers extend to data points within $1.5\times$ the interquartile range. C: number of classes. Ft.: finetune. ABMIL: attention-based multiple instance learning.

and CHIEF⁷⁴. Compared to TITAN, these models employ different pretraining strategies for their slide-level encoders (PRISM: WSI-report contrastive pretraining, GigaPath: masked image reconstruction pretraining, CHIEF: supervised contrastive learning of cancerous vs. non-cancerous WSIs), different patch-level encoders pretrained on histology patches trained at different magnifications and patch sizes (PRISM and GigaPath: 256×256 pixels at 20× magnification, CHIEF: 256×256 pixels at 10× magnification), and utilize a varying number of WSIs for pretraining (PRISM: 1.7×, GigaPath: 0.49×, CHIEF: 0.18× the WSIs used for TITAN pretraining). Additionally, we compare against mean pooling with the same CONCHv1.5 patch encoder as TITAN, which has shown to be a simple yet powerful unsupervised slide representation framework^{65,66,95}.

Furthermore, for a comprehensive evaluation of the baselines, we introduce two tumor classification tasks based on the publicly available repository TCGA with two different context lengths: (i) Main cancer type classification on ROIs (*TCGA-Uniform-Tumor-8K* or *TCGA-UT-8K*), a ROI-level cancer subtyping task with 32 classes, where we manually curated 25,495 tumor-containing regions of 8,192×8,192 pixels at 20× magnification ($\sim 4 \times 4$ mm²) across TCGA, covering the same tissue context as the region crops in TITAN_V pretraining (**Extended Data Figure 1, Extended Data Table 1**) and (ii) pan-cancer classification (*TCGA-OncoTree* or *TCGA-OT*), a slide-level OncoTree code⁹⁶ classification task with 46 classes, consisting of 11,186 formalin-fixed paraffin-embedded (FFPE) WSIs from TCGA. TCGA-OT is the largest pan-cancer slide-level classification task that is publicly available (**Extended Data Table 2**). With the exception of CHIEF, the pretraining datasets of TITAN (Mass-340K), PRISM, and GigaPath do not include TCGA and PANDA slides, which allows us to utilize these two datasets as benchmarking tasks without the concern of data leakage⁹⁷. More details on the two tasks can be found in **Online Methods** section **Downstream evaluation datasets**. For the benefit of the community, we plan to release TCGA-UT-8K datasets and TCGA-OT labels.

Prior to the expansive evaluations on diverse tasks, we first assess how the pretraining data scale affects the downstream performance of TITAN_V. We focus on four subtyping tasks – TCGA-UT-8K, TCGA-OT, OT108, and EBRAINS – due to the challenging diagnostic complexity with a large number of diagnostic classes. For the pretraining data scale, we train TITAN_V with 12.5%, 25%, and 50% of Mass-340K, maintaining the same distribution across the organs as the full dataset. We observe that the performance increases on all four tasks as more pretraining data is utilized, where TITAN_V with full Mass-340K exhibits 3.65%, 3.21%, and 1.21% average increase over TITAN_V pretrained with 12.5%, 25%, and 50% of Mass-340K, respectively (**Figure 2A**). Similar to the observation in pretrained patch encoders in pathology⁹ as well as ViTs for natural images^{83,98}, we observe the *scaling law*⁹⁹ for TITAN_V on all four tasks with more pretraining data leads to an increase in performance. Despite the difference in pretraining recipes, we also observe the same general trend for the three other slide encoders, where PRISM outperforms GigaPath and CHIEF by 9.01% and 20.1% on average, having 3.4 times and 9.7 times the number of pretraining WSIs, respectively. Furthermore, we observe that TITAN and TITAN_V, with 48.5 million and 42.1 million parameters respectively, outperform

heavier slide encoders PRISM and GigaPath, with 99.0 million and 86.3 million parameters, demonstrating superior parameter efficiency of our model.

We next evaluate TITAN on an expansive range of tasks comprised of morphological classification (14 tasks), grading (3 tasks), molecular classification (38 tasks), and survival prediction (6 tasks), where the summary of each task can be found in **Extended Data Tables 9 to 13**. On average, we observe that TITAN and TITAN_V outperform other slide encoders (**Figure 2C**). TITAN especially excels at morphological subtyping tasks across the entire spectrum of diagnostic complexities including fine-grade pan-cancer classification (challenging morphological classification tasks in **Extended Data Figure 2C**) and non-cancerous tasks such as allograft rejection, with TITAN and TITAN_V achieving an average of +8.4% and +6.7%, respectively, in performance on multi-class and binary subtyping tasks, averaged over balanced accuracy for multi-class tasks and AUROC for binary tasks over the next-best performing model, PRISM (**Figure 2C**). In particular, TITAN_V (and TITAN) not only outperforms others on TCGA-UT-8K with 8K×8K context that the model was trained on (+ 6% and 7.5% over PRISM), but also on WSI-level tasks that involve the entire tissue context, where TITAN_V benefits from the long-context extrapolation via ALiBi, e.g., TCGA-OT (+ 7% and 9.5% over PRISM), OT108 (+ 10% and 16% over PRISM), and EBRAINS (+ 9% and 9.1% over PRISM). To demonstrate the robustness of TITAN across different evaluation schemes, we run further analyses with prototyping evaluation^{100,101}, where the label of the query slide is determined by the proximity to the mean of the slide embeddings in each diagnostic class, as well as 20 nearest-neighbors evaluation (**Extended Data Tables 14 to 17**). Both TITAN and TITAN_V outperform the next best method PRISM by an even larger margin, + 14% and 9.5% for SimpleShot and + 15% and 9.2% for 20 nearest-neighbor evaluation, averaging balanced accuracy for all morphological subtyping tasks, which manifests that TITAN leads to improved representation quality. On grading tasks, TITAN outperforms the next best models CHIEF on average by + 3.2% and PRISM by + 4% in quadratic-weighted Cohen’s κ , where the high performance of CHIEF can be attributed to including the dataset PANDA in pretraining.

To evaluate the molecular classification performance, we tested the model on tasks from public datasets (BCNB and MUT-HET) and internal-external paired public datasets (TCGA, CPTAC, and EBRAINS), on immunohistochemistry (IHC) tasks, and MGB internal molecular tasks (**Figure 2B, Extended Data Figure 2, Extended Data Tables 18 to 54**). We observe that TITAN consistently performs best with + 0.9% on BCNB and MUT-HET, + 1.7% on TCGA, and +3.7% on internal molecular classification of BRCA and LUAD, in averaged AUROC scores over the next best model PRISM. For IHC quantification, TITAN_V performs best with + 12% in quadratic-weighted Cohen’s κ over the next best model CHIEF since our vision-language alignment does not include IHC reports leading to a slightly lower performance of TITAN of + 7.3 % over CHIEF. In external evaluations, where the linear classifiers trained on TCGA were applied to EBRAINS (IDH) and CPTAC (all other molecular endpoints), TITAN_V outperformed other slide encoders overall (+ 0.9% over

PRISM), whereas TITAN performance slightly worse (- 0.3 %). Since TITAN still outperforms PRISM by + 4.8% in SimpleShot and by + 5.4% in 20-nearest neighbors evaluation (**Extended Data Tables 35 to 43**), this could be attributed to suboptimal generalization of the linear classifier when searching for ℓ_2 -regularization strength over a large range to optimize for linear probing performance. We note that the pretraining dataset of CHIEF includes all WSIs from TCGA, which could partially contribute to its performance in respective tasks.

On survival prediction tasks, we utilize disease-specific survival (DSS)⁹⁰ as the clinical endpoint and the concordance index (c-index) as the evaluation metric and perform 5-fold cross-validation on six TCGA site-preserving stratified cancer cohorts. Specifically, we fit the linear Cox proportional hazards model on the slide embeddings to predict patient-level survival risk. We observe that TITAN and TITAN_v are generally the best-performing baselines, outperforming the next-best performing model CHIEF by +3.62% and +2.90% respectively, even though CHIEF was pretrained on TCGA slides (**Extended Data Table 55**). Interestingly, the mean pooling baseline shows competitive performance. This suggests that the proportion of different morphological phenotypes, which the mean baseline is effectively computing, is an important prognostic factor^{65,95}.

To further understand how the slide embedding space is organized and consequently affects the downstream performance, we investigate UMAP embeddings of WSIs within our largest and most diverse downstream dataset, TCGA-OT, where we color-code by organs instead of OncoTree codes to reduce visual clutter (**Figure 1E, Extended Data Figure 3**). We observe that the embeddings form distinct clusters along organs for TITAN and TITAN_v, with TITAN clusters seemingly better separated than TITAN_v (*e.g.*, breast further separated from bladder, stomach, and lung), confirming superior subtyping performance against other slide encoders. The embedding space for both the mean CONCHv1.5 and PRISM are reasonably separated reflecting good subtyping performances, whereas CHIEF and GigaPath are not able to effectively separate different organ embeddings, consequently leading to poor performance. This demonstrates that the WSIs from diverse organs in our pretraining dataset helps both TITAN and TITAN_v to be able to extract subtle organ-specific morphological cues more effectively.

Algorithmic design considerations for TITAN

To better understand how certain model choices affect the downstream performance, we perform ablation experiments on three design choices of TITAN: the positional encoding, the number of Transformer layers in TITAN_v, and the inclusion of vision pretraining (**Figure 2D**). Similar to previous analyses, we focus on the four subtyping tasks for the ablation experiments.

For the positional encoding, we pretrain TITAN_v with absolute positional encoding, following the ViT design⁸³, rotary positional encoding¹⁰² extended to 2D¹⁰³ (Rotary PE), and without positional encoding (No

PE). Our results show that ALiBi outperforms the other encoding schemes on all four tasks, with an average improvement of + 2.01% over absolute positional encoding, the second-best performing method (**Figure 2D**). This indicates that diagnostic performance can be enhanced with a suitable choice of positional encoding by contextualizing the patch features, with ALiBi helping TITAN_V extrapolate effectively to the entire slide, where the context length is over ten times longer. For the number of layers and consequently the number of parameters for TITAN_V, we observe that 6 layers (43M parameters) on average provide the sweet spot between smaller (4 layers, 29M parameters) and larger models (12 layers, 86M parameters), by outperforming them by 1.72% and 1.31%, respectively (**Extended Data Table 70**).

For the pretraining strategy ablation, we compare the performance between TITAN with the full pretraining and TITAN_L, which only performs vision-language alignment without vision pretraining. We introduce an additional baseline of multiheaded attention-based MIL^{66,67} with the vision-language pretraining (ABMIL-L), to further understand whether a different slide encoder architecture with the same pretraining recipe can be effective. We observe that TITAN outperforms TITAN_L by 2.35% and ABMIL-L by 3.62%. This indicates that the vision pretraining (Stage 1) provides better initialization weights than the random weights for vision-language alignment, leading to improved downstream performance, also observed in patch encoder pretraining¹⁰. This also suggests that the better performance of TITAN over PRISM, which is pretrained only with the vision-language alignment, could be due to the inclusion of the vision pretraining step. Moreover, the worse performance of ABMIL architecture than ViT with the same pretraining recipe justifies the choice of ViT as the architecture for TITAN.

Comparison with different learning paradigms for slide encoding

To further assess the quality of the slide embeddings produced by TITAN_V, we evaluate different learning paradigms by comparing the linear probe performance of each slide encoder against other MIL models comprised of *mean pooling*, i.e., averaging the patch embeddings, *attention-based MIL* (ABMIL)⁶⁷, and task-specific *finetuning* of the slide encoder from random or the pretrained weights. For the mean pooling and ABMIL baselines, we use respective patch encoders for each slide encoder framework. This analysis allows us to gauge whether the pretrained slide encoders have learned meaningful slide representations and consequently outperform the simple yet powerful unsupervised (mean pooling) and supervised (ABMIL) baselines, neither of which involve large-scale pretraining on thousands of WSIs.

We observe several trends with TITAN (**Figure 2E, Extended Data Figure 4, Extended Data Tables 56 to 59**). First, ABMIL outperforms mean pooling, as expected, since ABMIL is supervised and equivalent to weighted averaging of the patch features, which would by default include the simple averaging solution of mean pooling. Next, the linear probe outperforms ABMIL, which demonstrates that TITAN and TITAN_V,

having been pretrained on a large repository of multimodal pathology data much larger than what is provided to ABMIL for each downstream task, can encode additional contextual and semantic morphological details of the slide. This leads to task-agnostic slide embedding of TITAN_V being better equipped for downstream tasks than tasks-specific supervised slide embeddings of ABMIL. Finally, we observe that task-specific finetuning of TITAN leads mostly to performance improvement over linear probe of TITAN and TITAN_V. Furthermore, finetuning the slide encoder from randomly initialized weights yields lower performance (-3.63% on average) than from the TITAN pretrained weights. This suggests that the pretrained weights of TITAN_V can serve as a good initialization set for downstream tasks for typical cohorts with a limited number of patients, in line with previous works^{62,64}. One exception is OT108, which could be attributed to the small number of samples for each class (ranging from 4 to 42), which may lead to overfitting. We observe that similar trends exist for PRISM except for finetuning scenarios, which could not be compared due to PRISM finetuning recipes not being provided.

Interestingly, these trends are not always observed in other slide foundation models. For GigaPath, finetuning the slide encoder significantly improves the performance over linear probe (14.9% on average), but lags behind TITAN and TITAN_V linear probe by 6.25% and 4.17% on average, respectively. That finetuning leads to a significant improvement, combined with low linear probe performance, suggests a lack of generalizability off-the-shelf for Gigapath slide embeddings. This is further supported by the fact that the simple mean pooling, which does not leverage pretraining on large WSI repositories, outperforms the linear probe across all four subtyping tasks (**Extended Data Figure 4**). For CHIEF, the trends are mixed, with ABMIL performing better than the linear probe on half of the tasks. Interestingly, finetuning the slide encoder always yields worse performance than the linear probe (-8.38% on average), suggesting that the pretrained CHIEF weights might be a sub-optimal starting point for task-specific fine-tuning. Nevertheless, finetuning from pretrained weights indeed yields better performance on average than randomly initialized weights for both GigaPath and CHIEF, in line with what is observed for TITAN.

Few-shot learning for low data regime

We also evaluate the data-constrained setting of few-shot learning where only a few examples for each category are provided within the linear probe setting. In the few-shot setting, TITAN_V remains superior across all tasks and number of shots in balanced accuracy when assessed with linear probe (**Figure 2F**). To mitigate sampling bias, we aggregate the results over 50 different runs, with random samples used for training, while fixing the test set. We observe that TITAN is the best-performing model across different tasks, demonstrating the strong generalizability of TITAN. TITAN_V is the second-best performing model, which supports our results that vision-language alignment benefits the downstream task performance. Specifically, TITAN and TITAN_V exhibit especially high performance in one-shot learning, which is on par with other slide encoders

trained on more shots (**Extended Data Tables 72 to 75**). Specifically, TITAN and TITAN_V outperform CHIEF based on 16 shots on TCGA tasks by 22.4% and 13.5% (TCGA-UT-8K) and 18.7% and 6.8% (TCGA-OT) when compared with the median value of 50 runs, respectively, even though CHIEF has been pretrained on TCGA slides.

Interestingly, both TITAN and TITAN_V also outperform ABMIL with the same patch encoder across all settings. While the performance gap with ABMIL shrinks for a higher shot regime as expected, we observe that the gap is indeed wider in the lower shot regime. The largest gap for 1-shot is observed in the OT108 task, where TITAN outperforms ABMIL by 56.7%. These observations underscore the superior data efficiency of a pretrained slide encoder and suggest that TITAN_V can excel in rare cancer settings with a limited number of samples, such as OT108 in our benchmark, where heavily parameterized supervised approaches are inherently restrained. This demonstrates the advantage of TITAN over other slide encoders such as GigaPath and CHIEF, the intended usages of which are in supervised settings with task-specific finetuning, rather than off-the-shelf usage with the frozen slide embeddings. The same trend is observed even when evaluated with prototyping, where K samples (shots) from each class are averaged to construct the prototype (**Extended Data Tables 76 to 79**). More details on the experiments can be found in **Online Methods** section **Few-shot classification**.

Language-aligned TITAN enables cross-modal capabilities

We further assess the language capabilities of TITAN by aligning the slide representations of TITAN_V to language-based morphological descriptions. Specifically, we assess the cross-modal zero-shot classification^{55,56,104} and report-generation capabilities of TITAN and study the effect of Stage 2 pretraining for caption alignment with fine-grained morphological descriptions and Stage 3 pretraining with coarse clinical reports describing the relevant microscopic findings.

To evaluate the quality of vision-language alignment, we first perform cross-modal zero-shot experimentation on 13 subtyping tasks of varying difficulties (**Figure 3A**). In cross-modal zero-shot evaluation, the diagnostic labels expressed as text prompts are encoded with the text encoder. The diagnostic prediction of the query slide is then decided by the closest label embedding to the slide embedding encoded with TITAN. The cross-modal zero-shot experiment evaluates how the embedding space with the visual pretraining can be further aligned with the language modality. We compare the zero-shot performance against PRISM, also equipped with cross-modal capabilities. Gigapath is not included as the language-aligned extension has not been publicly released. We observe that TITAN performs best across these tasks, outperforming PRISM by a large margin on multi-class classification tasks (balanced accuracy +56.52%) and binary subtyping tasks (AUROC +13.8%), for both cancer subtyping tasks and non-cancerous tasks (**Figure 3B, Extended Data Tables 80 to 92**). The performance gap between TITAN and PRISM is the widest on the 30-class EBRAINS subtyping task, where

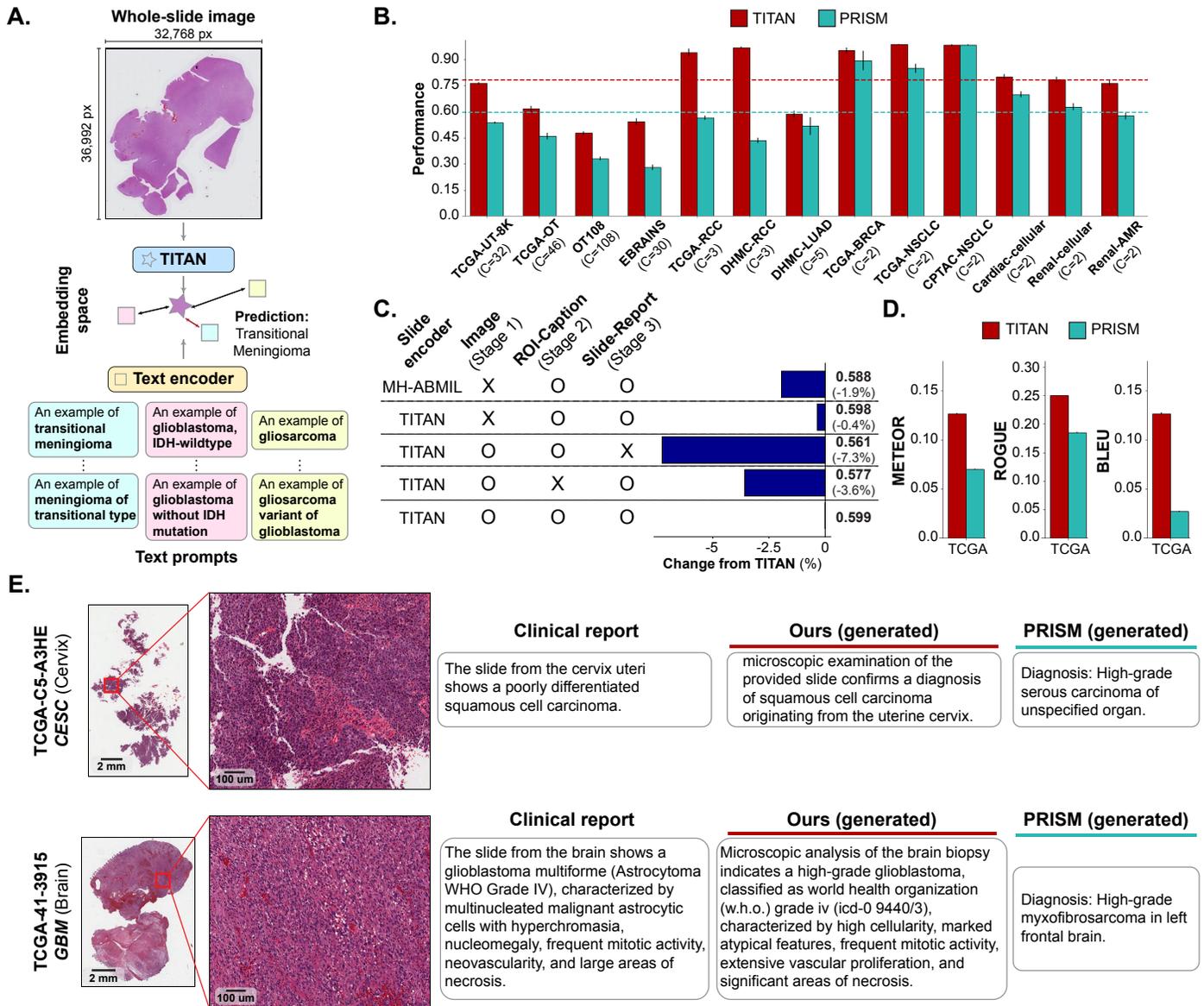


Figure 3: Visual-language evaluation of TITAN. (a) A schematic for zero-shot evaluation. The query slide is classified by identifying the closest text prompt embedding in the slide embedding space. (b) Zero-shot performance of TITAN and PRISM. All multi-class tasks are evaluated with balanced accuracy and binary tasks are evaluated with AUROC. All error bars represent standard deviations based on bootstrapping. (c) Ablation study comparing different pretraining strategies, and assessed with zero-shot performance averaged across TCGA-UT-8K, TCGA-OT, OT108, and EBRAINS. Evaluations are based on the percentage changes of balanced accuracy from the reference zero-shot performance of TITAN. (d) Report generation evaluation on TCGA-Slide-Reports, and evaluated using METEOR, ROGUE, and BLEU. (e) TCGA examples of generated reports of TITAN and PRISM, with the corresponding clinical reports. Additional examples of generated reports are available in **Extended Data Figure 5**. C: number of classes.

the balanced accuracy of TITAN is more than double that of PRISM (balanced accuracy of +121.9%). The text prompts used for zero-shot experiments can be found in **Extended Data Tables 102 to 123**.

To further understand how different design considerations affect the zero-shot performance of TITAN, we ablate over pretraining stages and the slide encoder architecture (**Figure 3C**). In total, we experiment with four different variations of TITAN and present the average performance over four challenging subtyping tasks at slide level, TCGA-UT-8K, TCGA-OT, OT108, and EBRAINS (results for each dataset can be found in **Extended Data Tables 93 to 96**). We observe that TITAN maintains the best overall zero-shot performance. Of the three pretraining stages, Stage 1 vision-pretraining contributes the least (balanced accuracy of -0.4% against TITAN), followed by Stage 2 ROI-caption alignment (-3.6% against TITAN) and Stage 3 slide-report alignment (-7.3% against TITAN). This demonstrates that aligning vision and language at both fine-grained and global levels, thereby combining the insights independently derived at patch-level^{7,10} and slide-level^{58,62}, is necessary, which is lacking in report-only aligned baselines such as PRISM and GigaPath. Finally, the variant using a multi-headed-ABMIL (MH-ABMIL) network as vision backbone with vision-language alignment pretraining lags behind TITAN with and without vision-pretraining by 1.94% and 1.54%, indicating that our ViT architecture using self-attention with ALiBi provides better downstream performance than attention-based networks.

Finally, we assess TITAN’s capabilities of generating pathological reports, utilizing the text decoder trained during CoCa pretraining. To this end, we introduce a report generation task on TCGA, consisting of 10,108 FFPE WSIs with paired slide-level reports parsed from 9,523 patient-level TCGA reports released by a previous study¹⁰⁵. We evaluate the models using three metrics METEOR¹⁰⁶, ROGUE¹⁰⁷, and BLEU¹⁰⁸. We observe that TITAN outperforms PRISM by a large margin, on average by 161% across the three metrics (**Figure 3D**). In addition, TITAN outperforms TITAN without Stage 2 pretraining, agreeing with the previous experiments on the importance of the ROI-level vision-language alignment. Examples of the generated reports for TITAN considered high-quality by the pathologists are shown in **Figure 3E**, often capable of correctly capturing key attributes such as tissue site, diagnosis and tumor grade as well as key representative morphology. Additional examples are illustrated in **Extended Data Figure 5**. More details on the dataset can be found in **Online Methods** section **Downstream evaluation datasets**.

TITAN enables rare cancer retrieval and cross-modal retrieval

Considering cases with similar morphological features and diagnoses is essential for pathologists to make informed decisions, in particular when dealing with complex or rare cases^{5,17,47,48,50,51,53,109,110}. Retrieving similar histology slides or pathology reports facilitates identifying relevant cases from large archival databases, and has become an essential clinical decision support function in digital pathology workflows. This

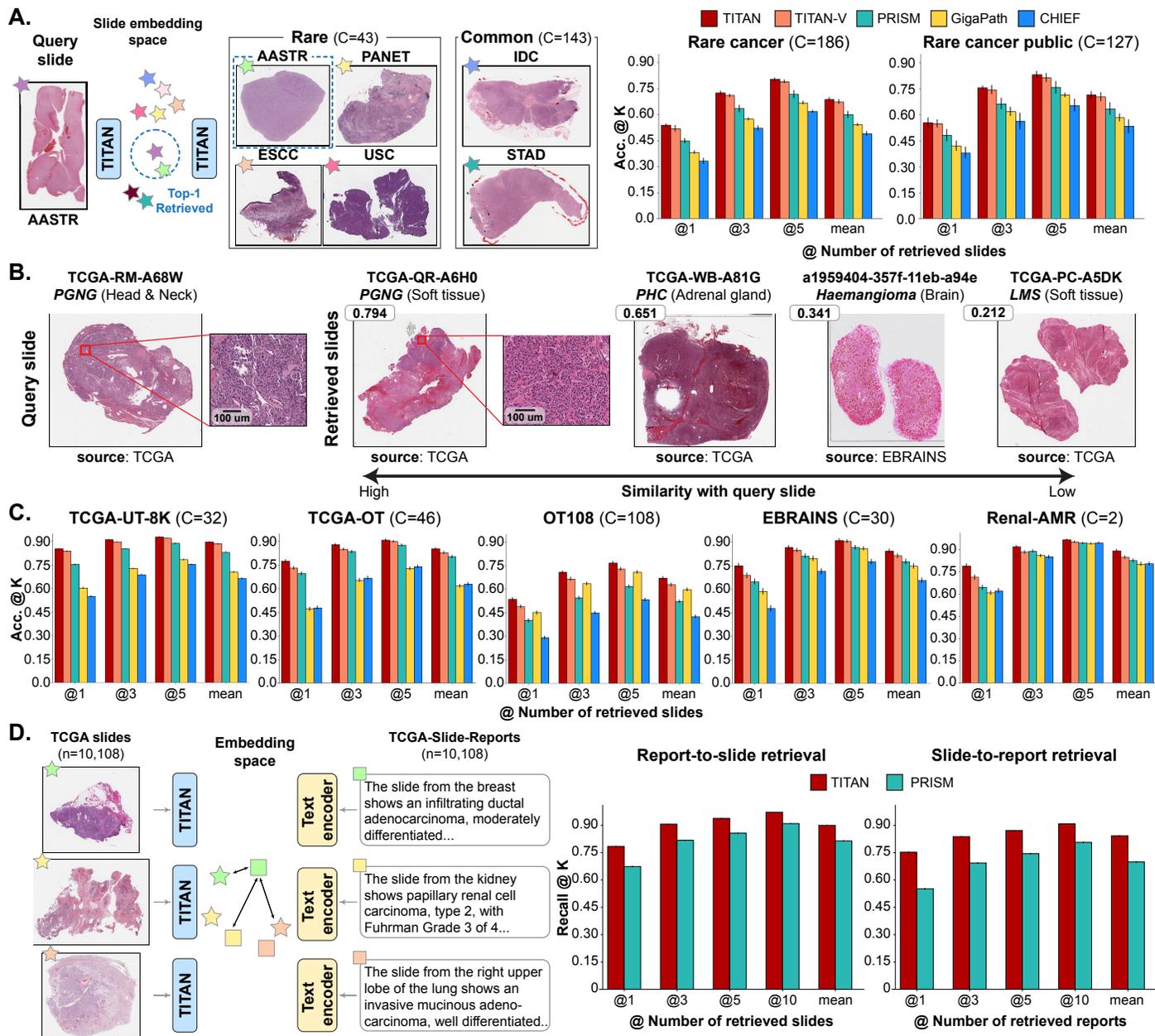


Figure 4: **Retrieval capabilities of TITAN.** (a) Slide retrieval results on rare cancer retrieval tasks assessed with accuracy@ K , with $K = \{1, 3, 5\}$. Rare-Cancer (internal rare cancer cohort) consists of TCGA, EBRAINS, and the MGB internal cohort, with 43 rare and 143 common cancer types for a total of 186 classes. Rare-Cancer-Public (public rare cancer cohort) consists of TCGA and EBRAINS only, with 29 rare and 98 common cancer types for a total of 127 classes. (b) Example of rare cancer retrieval on Rare-Cancer with the query slide and four representative retrieved slides. The number indicates the cosine similarity between the query and the retrieved slide. Additional examples of rare cancer retrieval are available in **Extended Data Figure 6**. (c) Slide retrieval results on five subtyping tasks. Mean represents the average performance across three shots. (d) Report-to-slide and slide-to-report cross-modal retrieval performance assessed with recall @ K , with $K = \{1, 3, 5, 10\}$ on TCGA cohort of 10,108 pairs of WSIs and reports for TITAN and PRISM. Mean represents the average performance across four shots. All error bars represent standard deviations based on bootstrapping. C: number of classes.

is especially beneficial for rare cancers that affect fewer than 15 individuals per 100,000 annually^{42–44}, for which pathologists can identify non-specific malignancies based on WSIs with similar morphologies and their corresponding pathology reports. Slide foundation models readily provide WSI representations for vector database indexing, significantly simplifying the task of histology slide retrieval compared to patch foundation models, which provide more than 10^4 representations per WSI and consequently renders slide-level retrieval non-trivial.

Given a query slide and labeled set of support slides (indexed into a vector database by a slide foundation model), histology slide search is evaluated by assessing the accuracy performance in retrieving similarly labeled slides from the support set. This setting is non-parametric and solely relies on how the slide representations are clustered along different diagnostic labels. Specifically, we test whether the K -closest neighbors of a query slide in the embedding space—determined using cosine similarity with $K = \{1, 3, 5\}$ —include slides sharing the same diagnostic label as the query slide. Performance is assessed using $\text{Accuracy}@K$, which measures whether at least one of the K neighboring slides has the same diagnostic label as the query. We also provide $\text{MVAcc}@K$ which requires the majority vote of top- K neighboring slides is of the same diagnostic label as the query and is therefore more stringent criteria than $\text{Accuracy}@K$. For the rare cancer retrieval task, we create a large database of 186 cancer types with 19,626 WSIs, Rare-Cancer, by combining the *rare cancer set* of 43 cancer types (3,039 WSIs) with the *common cancer set* of 143 more common cancer types (16,587 WSIs) from TCGA, EBRAINS, and MGB internal data (**Figure 4A, Extended Data Table 6**). To assess the performance, we create a query set as the subset of the *rare cancer set*, ensuring all 43 rare cancer types are represented. The support set, from which similar slides are retrieved, is constructed by incorporating the remaining WSIs of the *rare cancer set* into the *common cancer set*, ensuring all 186 cancer types are represented. This design emulates the real-world setting of clinicians interacting with an extensive cancer database encompassing a diverse mix of rare and common cancer types. This procedure is repeated five times, with a different query set each time. We additionally create a public version with 127 cancer types and 14,062 WSIs, Rare-Cancer-Public, using the data from TCGA and EBRAINS resulting in 29 rare cancer types (1,982 WSIs) and a lower diversity in the set of common cancers with 98 types (12,080 WSIs). The same evaluation procedure as for Rare-Cancer is repeated (**Extended Data Table 8**).

We observe that TITAN and TITAN_v outperform other slide encoders with +14.8% and +12.3% in $\text{Accuracy}@K$ and +18.1% and +13.4% in $\text{MVAcc}@K$ to the next best model PRISM (**Extended Data Table 129**) on average. The trends in performance are preserved on the public version of the rare cancer task with slightly higher performance levels as the task is easier with a support set containing fewer cancer types (**Extended Data Table 130**). An example of rare cancer retrieval is demonstrated in **Figure 4B**, where the closest slide to the paraganglioma (PGNG) query is also of PGNG with a high similarity of 0.794 and less similar slides are of different cancer type (Haemangioma from brain, similarity of 0.341). One of the retrieved slides is

Pheochromocytoma (PHC) with a high similarity of 0.651, agreeing with the clinical understanding that both are morphologically tightly connected as rare neuroendocrine tumors¹¹¹. Additional examples of rare cancer retrieval can be found in **Extended Data Figure 6**, where the retrieved slides of high similarity are indeed from the same diagnostic label or organ. Even when further assessed with multi-class cancer subtyping tasks, from relatively simple AMR for renal allograft ($C = 2$) to challenging OT108 ($C = 108$), we observe that both TITAN and TITAN_v outperform other slide encoders (**Figure 4C, Extended Data Tables 131 to 137**).

Encouraged by the unimodal retrieval performance, we further investigate the cross-modal retrieval performance of TITAN, as the slide and report embedding spaces are aligned from the pretraining steps. We perform the cross-modal experiments on TCGA-Slide-Reports, our proposed dataset for report generation with 10,108 slide-report pairs (**Extended Data Table 7**). For the report-to-slide (slide-to-report) retrieval task, we test whether any of the K -closest slides (reports) for the query report (slide) in the embedding space have the same diagnostic label as the query, the metric referred to as Recall@ K with $K = \{1, 3, 5, 10\}$. We observe that TITAN outperforms PRISM on both retrieval tasks across all K retrievals with +10.5% and +20.5% on average for report-to-slide and slide-to-report retrieval tasks, respectively (**Figure 4D, Extended Data Tables 138 and 139**). The largest gap to PRISM is observed for slide-to-report retrieval when only a single report was retrieved, where TITAN outperforms by 36.4%. The strong performance of TITAN even with a single report (0.75) hints at the clinical potential, where for a diagnostically challenging slide clinicians can benefit from sifting through retrieved past medical reports that describe identical diagnoses, and vice versa. More details on the experiments can be found in **Online Methods** section **Slide retrieval** and **Cross-modal retrieval**.

Discussion

We introduce a multimodal whole-slide foundation model for pathology, TITAN, that combines and elevates successful recipes of self-supervised learning (SSL) from the patch level to the slide level. Methodologically, TITAN employs histology knowledge distillation in the feature space (vision-only) and contrastive learning by aligning regions of interest (ROIs) with synthetic captions and whole slide images (WSIs) with reports (vision-language). Pretrained on 336K WSIs, TITAN, a Vision Transformer (ViT) architecture equipped with ALiBi positional encoding for long-context extrapolation, produces powerful general-purpose slide representations for a large variety of downstream tasks even without task-specific finetuning. From cancer subtyping to molecular classification, TITAN consistently outperforms other state-of-the-art slide encoders, such as PRISM⁶², GigaPath⁵⁸, and CHIEF⁷⁴. This superiority is maintained in data-constrained settings such as rare disease classification and histology slide retrieval, which underscores the representation quality of TITAN. Further aligning the vision-pretrained TITAN with 423K ROI-level captions generated by PathChat and 183K pathology reports equips the model with multimodal capabilities such as zero-shot diagnosis, slide-report retrieval, and report generation. We observe that aligning the slide embedding with both the fine-grained (ROI

captions) and coarse-level (pathology reports) morphological descriptions is crucial for handling the multiscale information inherent in tissue slides—an insight made possible for the first time through the use of generated ROI captions. Similar to the unimodal setting, TITAN outperforms PRISM, another language-equipped model, on all cross-modal tasks. To advance the field of slide-representation learning¹¹², we curated and plan to release two challenging multi-class morphology classification tasks beyond patch-level from the publicly available repository TCGA: TCGA-UniformTumor-8K (TCGA-UT-8K) for 32-class tumor-ROI subtyping and TCGA-OncoTree (TCGA-OT) for 46-class WSI-level OncoTree code classification.

Detailed ablation analyses reveal further insights into the properties of TITAN. We observe that Stage 1 unimodal pretraining of TITAN_V captures morphological concepts already with much less data than existing slide encoders, as demonstrated in the few-shot data efficiency experiments. In particular, TITAN_V consistently outperforms its mean pooling and task-specific attention-based pooling baselines that utilize the same patch encoder as TITAN_V, proving that unimodal pretraining effectively captures the context of patch features in contrast to existing unimodal slide encoders. Next, in addition to unlocking language-related capabilities, we observe that the vision-language alignment further enhances the representation quality of our vision-only model. In particular, TITAN improves over TITAN_V on average for slide-level tasks with the strongest improvements in evaluation settings that solely rely on the structure of the slide embedding space without any parameter tuning, such as prototyping or k -nearest neighbor settings. A further sign of the improved representations is that TITAN outperforms all other baselines, including TITAN_V, on slide retrieval and few-shot tasks. While slide embeddings from pretrained TITAN are already promising, especially in the low-data regime, task-specific fine-tuning of the pretrained model can further enhance the downstream performance for tasks with a large enough patient cohort, pointing to the flexibility of TITAN when applied to diverse clinical and data settings.

Providing multimodal slide embedding off-the-shelf presents immediate clinical potential to assist clinicians in their routine diagnostic workflows⁷⁶. Presented with challenging patient tissue slides to diagnose, pathologists and oncologists can hugely benefit from being able to retrieve and analyze diagnostically similar slides or clinical reports, likely leading to a reduction in patient misdiagnosis and interobserver variability. As shown in the extensive set of experiments, TITAN can accurately retrieve similar diagnostic slides and reports for challenging scenarios from a large number of cancer types (> 100), as well as rare cancer types⁴⁴ where the corresponding slides have scarce representation in the database. That all of these could be performed off-the-shelf with pretrained TITAN without a dedicated algorithm for each task underscores both the generalizability of TITAN slide embeddings, as well as how slide-level tasks can become simpler with the advent of pretrained slide encoders.

Despite the encouraging performance of TITAN, our framework has a few shortcomings. First, Mass-

340K contains fewer slides compared to other pretraining datasets used for patch encoders^{12,13,113} and slide encoders such as PRISM⁶². We believe that the already strong performance of TITAN, merged with concurrently ongoing effort to expand Mass-340K, will further allow improved slide-level and cross-modal performance. Next, pretraining on the region crops of $8K \times 8K$ and extrapolating with ALiBi to the entire WSI may still not capture the full contextual information. Larger pretraining contexts such as $16K \times 16K$, and other positional encodings for extrapolation could address this limitation. Finally, preprocessing of clinical reports presents a further challenge for vision-language alignment; Striking a balance between relevant information for contrastive learning while only including morphology-related information is non-trivial and involves a lot of manual tuning despite the automated processing pipelines. Restructuring the reports into distinctive morphology and molecular characteristics could facilitate contrasting with only relevant information.

Promisingly, TITAN can be scaled up in terms of data and architecture to improve performance. WSIs and corresponding medical reports are routinely available and stored in the clinical workflow. The synthetic region-level captions can be generated with the generative AI model in an unlimited manner, providing the model with a wealth of text guidance. Using this additional data, a heavier ViT slide encoder architecture than what is currently being utilized for TITAN can potentially improve the performance, as was already demonstrated at the level of patch encoders^{12,13,113}. In addition, the improved patch representation quality from more powerful patch encoders will likely improve the quality of the downstream slide encoder. We envision TITAN and its future iterations being incorporated into practitioners' everyday toolkits for routine application and comparison with other task-specific supervised frameworks, together reaching new levels of performance in clinically important tasks.

Online Methods

Pretraining dataset

For large-scale visual pretraining, we curated Mass-340K, a diverse dataset consisting of 335,645 WSIs across 20 organs, with 90% hematoxylin and eosin (H&E) stained slides and 10% immunohistochemistry (IHC) slides, sourced from the combination of in-house histology slides and the GTEx consortium¹¹⁴. To explore the effects of data scale at the pretraining stage, we formed three additional partitions of Mass-340K, containing 12.5%, 25%, and 50% of the original dataset. These partitions were sampled to maintain the ratio of different data sources and preserve organ distribution.

Synthetic caption generation using PathChat

For the initial stage of vision-language alignment (Stage 2 of TITAN), we used synthetic captions generated by PathChat, a state-of-the-art multimodal LLM designed for pathology⁷⁹. To go beyond the typically

brief clinical reports focused on the final diagnosis, we prompted PathChat to generate detailed morphological descriptions of ROIs, providing important training data for models to capture complex pathological features. Using PathChat, we generated synthetic captions for 423,122 diverse $8,192 \times 8,192$ ROIs sampled from Mass-340K. Since PathChat cannot process inputs of size $8,192 \times 8,192$ pixels directly, we divide each ROI into 64 $1,024 \times 1,024$ patches. To retain the most representative morphological features, we applied K-means clustering with $K = 16$ to the 64 patches and then randomly sampled one patch from each cluster. The resulting 16 morphologically-representative $1,024 \times 1,024$ patches were subsequently fed to PathChat. To further enhance the diversity of these captions, we utilized Qwen2-7B-Instruct⁸⁷ to rewrite the generated captions, ensuring varied language structures and expressions. Detailed prompts for both PathChat and Qwen2, along with examples of generated and diversified captions, are provided in **Extended Table 124-125**.

Curation of slide-report dataset

For the second stage of vision-language alignment (Stage 3 of TITAN), we curated a dataset of 182,862 slide-report pairs from a combination of in-house clinical reports and pathology notes from the GTEx consortium¹¹⁴. However, clinical reports are often noisy and are typically organized at the patient level, hence contain information on multiple slides from the same patient, complicating the slide-report alignment. To address this, we utilized a locally served Qwen2-7B-Instruct⁸⁷ model to extract slide-specific descriptions and remove sensitive information unrelated to pathological diagnosis, such as gross descriptions, hospital and doctor names, and patient clinical history. Additionally, we applied the same rewriting strategy used for synthetic captions to diversify the report text. Example prompts used for report cleaning and rewriting can be found in **Extended Data Table 126-128**.

Unimodal visual pretraining

Preprocessing

Similar to the previous studies^{9,10,45}, WSIs were preprocessed by tissue segmentation, tiling, and feature extraction using a pretrained patch encoder. We used the CLAM toolbox⁴⁵ for tissue segmentation and tiling. Tissues were segmented by binary thresholding of the saturation channel in HSV color space at a low resolution. Following this, we applied median blurring, morphological closing, and filtering of contours below a minimum area to smooth tissue contours and eliminate artifacts. Non-overlapping 512×512 pixel patches were then extracted from the segmented tissue regions of each WSI at $20\times$ magnification. For feature extraction, we used CONCHv1.5, an extended version of CONCH¹⁰, which was trained with 1.26 million image-caption pairs using the CoCa training objective for 20 epochs. The choice of CONCHv1.5 for feature extraction was due to the fact the model was pretrained on histology regions with diverse stains and tissue types, including FFPE,

frozen tissue, and immunohistochemistry, thereby yielding region features that are robust against diverse tissue processing protocols. Refer to **Extended Data Table 100** for detailed hyperparameters of the patch encoder.

To enhance the effectiveness of the ROI sampling strategy during Stage 1 training of TITAN_V, an additional preprocessing step was performed to group the segmented tissue contours based on their spatial proximity within the slide. This addresses the challenging cases where multiple tissue regions are interspersed with background areas, particularly for biopsy samples where tissue fragments are often widely dispersed and for samples with multiple slices placed on the same slide. Specifically, we grouped tissue contours into clusters based on their coordinates, resulting in tissue groups that contain densely packed tissue regions with minimal background regions between them. Furthermore, tissue groups that contained fewer than 16 patches were filtered out. This grouping operation produced a total of 345,782 tissue groups from Mass-340K.

Pretraining protocol

For training TITAN_V on Mass-340K, we use iBOT, a state-of-the-art self-supervised learning method based on the combination of student-teacher knowledge distillation and masked image modeling⁷⁷. As iBOT is applied in the patch embedding space, instead of the typical use case of the raw image space, we adapt the pretraining recipes as follows.

View generation. During training, we create region crops randomly sampled from the tissue groups, each of which corresponds to a feature grid of size 16×16 , corresponding to a field of view of $8,192 \times 8,192$ at $20 \times$ magnification (**Figure 1B**). The random sampling of region crops, instead of precomputing fixed regions, increases the diversity of the training set and effectively acts as an additional data augmentation, as the model encounters different parts of the same WSI at each training epoch. A region crop contains 256 features, which is equivalent in length to training on images of 256×256 pixels with a token size of 16×16 in the typical natural image setting. From this region crop, two global views (14×14 crops) and ten local views (6×6 crops) are generated by cropping within the region crop without scaling or interpolation and fed to iBOT training.

To achieve realistic augmentations in the embedding space, previous methods have employed offline image augmentations in the pixel space^{34,59} by extracting multiple patch features from different views of a given patch. While effective, this approach limits the number of additional views and becomes computationally infeasible for large training datasets. Additionally, choosing color space augmentations adapted to the use of histopathology that go beyond standard color transformations introduces additional computational overhead. A few recent approaches addressed the difficulty with training generative networks on the feature space to transform the features^{115,116}, but also introduced additional computational cost for training. Instead, we apply frozen feature augmentations, which have been shown to work well for a few-shot classification task in the

feature space of pretrained Vision Transformers⁸⁴.

Positional encoding. Traditional multiple instance learning methods consider the patches to be permutation-invariant within the slide. Despite the promising results, this approach ignores the tissue context which can be essential for capturing the interaction in the tumor micro-environments and can thus affect the model’s performance¹¹⁷. In this context, for TITAN, we employ positional encodings in the patch embedding space to break permutation invariance and encode tissue context. Furthermore, TITAN adopts the strategy of *Train short, test long* to ease the computational burden, which also requires positional information via positional encodings. Trained at the region crops (ROIs) of $8,192 \times 8,192$ pixels (*Train short*), we directly apply TITAN on the whole slide during inference (*Test long*). We used Attention with Linear Biases (ALiBi), a method originally proposed for 1D sequence in large language models (LLMs)⁸⁵. Absolute positional encoding, another popular alternative that works well for images at training sizes, was shown to have weak extrapolation abilities⁸⁵. Unlike other positional encodings applied to the input features, ALiBi adds a bias to the query-key dot product during the computation of attention scores. ALiBi effectively penalizes the attention score for tokens which are further apart from each other. Formally, let $q_i \in \mathbb{R}^d$ and $k_j \in \mathbb{R}^d$ represent the i -th query and j -th key, respectively. The attention score, which is typically computed as $\text{softmax}(q_i k_j^T)$, is modified with 1D ALiBi as $\text{softmax}(q_i k_j^T - m|i - j|)$, where m is a predefined slope specific to each attention head. Since the feature grids and the resulting views are of 2D grid structure, we extend ALiBi to 2D by incorporating the Euclidean distance between patches i and j . The 2D ALiBi can be written as

$$\text{softmax}\left(q_i k_j^T - m\sqrt{(i_x - j_x)^2 + (i_y - j_y)^2}\right), \quad (1)$$

where i_x, i_y and j_x, j_y are the 2D grid coordinates of patches i and j . The x and y coordinates are defined as the 2D patch coordinates (at magnification $20\times$) divided by the patch size of 512.

Network architecture and training details. For the slide encoder, we use a Vision Transformer (ViT)⁸³ with 6 Transformer layers, 12 attention heads of dimension 64 resulting in an embedding dimension of 768, and a hidden dimension of 3,072. This smaller architecture, compared to typical ViTs used in patch encoders, is chosen based on previous studies⁵⁷, which suggest that a compact network suffices for slide representation learning on the embedding space, especially given the limited data scale of WSIs compared to histology patch datasets at the scale of billions. The patch embedding layer is replaced by an MLP to process the feature inputs. We train the model for 270 epochs (equivalent to 91,260 iterations), distributed across four NVIDIA A100 80GB graphics processing units (GPUs) with a local batch size of 256 per GPU. For all training hyperparameters, refer to **Extended Data Table 97**.

Vision-language continual pretraining

To enhance the unimodal capabilities of TITAN_V, we further explored the multimodal vision-language alignment of TITAN_V with clinical text. Training a multimodal foundation model, however, faces several limitations related to data and compute. First, paired slide-report data are scarce compared to the scale of millions of image-caption pairs for patches. Additionally, real-world clinical reports typically only contain brief diagnostic information, unlike the detailed morphological descriptions in educational captions for histology regions of interest (ROI) images. Finally, contrastive learning-based cross-modal training typically requires a large batch size, which is computationally infeasible for WSIs.

To address these issues, we propose a two-stage continual pretraining approach (referred to **Stage 2** and **Stage 3** for TITAN) that progressively aligns the model with increasing context. We first align synthetic captions for $8,192 \times 8,192$ ROIs, followed by real clinical reports for WSIs. With emphasis on detailed morphological descriptions, the first vision-language alignment stage allows the model to learn fine-grained pathological concepts using a large batch size. In the next stage, we further augment the model’s understanding of diagnostic terminology and reasoning, targeted to enhance its zero-shot understanding in downstream tasks. The second stage also serves as a “high-resolution fine-tuning” phase, adapting the model from the local contexts of ROIs to the full-scale global context of WSIs. Altogether, these two stages are designed to gradually build the model’s ability to comprehend and generate meaningful vision-language representations for WSIs.

Network architecture and training details

Following the success of previous studies¹⁰, we use CoCa⁸⁶, a state-of-the-art visual-language foundation model pretraining method, for both stages of vision-language alignment. The model consists of an image encoder, a text encoder, and a multimodal text decoder. Using our unimodal TITAN_V as the image backbone, we add two attentional pooler components on top. The first attentional pooler uses a single query (contrastive query) to pool a single global representation of the feature grids and enable cross-modal contrastive learning with text embeddings. This global WSI representation can then be used for zero-shot or unsupervised evaluation of TITAN on downstream tasks. The second attentional pooler uses $n = 128$ queries (reconstruction queries) to generate a set of 128 image tokens designed for interacting with the multimodal text decoder for caption generation. We use the pretrained text encoders and multimodal decoders of CONCHv1.5¹⁰, each consisting of 12 Transformer layers with an embedding dimension of 768 and a hidden dimension of 3,072.

For both stages, we used 8 NVIDIA A100 80GB GPUs. During Stage 2 vision-caption pretraining, we used a local batch size of 196 per GPU, with gradient accumulation of 2 resulting in an effective batch size of 3,136. For Stage 3 vision-report pretraining, we randomly crop the WSIs to 64×64 feature grids to allow for larger batch sizes while maintaining a large field-of-views, corresponding to $32,768 \times 32,768$ pixels, which already covers most slides in our pretraining dataset. We used a local batch size of 16 per GPU, with a gradient

accumulation of 2 to achieve an effective batch size of 256. To avoid deteriorating the quality of the pretrained vision encoder, we used a smaller learning rate and weight decay, as well as a slow warm-up strategy for the vision backbone, following previous work¹¹⁸. For all hyperparameters, refer to **Extended Data Table 98-99**.

Evaluation setting

Baselines

We compare TITAN_V against 1) *unsupervised* baselines with four other slide encoders, Prov-GigaPath (referred to as GigaPath in the manuscript)⁵⁸, PRISM⁶², CHIEF⁷⁴, and the mean pooling baselines with features from the respective patch encoders, 2) *supervised* baselines, and 3) our vision-language model TITAN against zero-shot baseline PRISM.

Unsupervised baselines. GigaPath uses LongNet architecture as the slide encoder, a ViT⁸³ in the “base configuration”, replacing the vanilla dense attention with dilated attention. It was trained on 171,189 in-house WSIs from Providence via masked autoencoder¹¹⁹. As patch encoder, GigaPath uses ViT-G/14 pretrained with DINOv2⁷⁸ on the same in-house dataset. While GigaPath further performed continual vision-language pretraining, we only assess the unimodal model, as the multimodal model is not publicly available. For performance analysis, we use the output of the Transformer layer 11 as slide representation, which yields the best results on downstream tasks and also agrees with the provided finetuning recipe. PRISM⁶² employs the Perceiver architecture¹²⁰ as the slide encoder with CoCa-based vision and language alignment⁸⁶ on 195,344 specimen-report pairs, where each specimen contains one or more WSIs with a total of 587,196 WSIs. As for the patch encoder, PRISM uses Virchow¹¹, a ViT-H/14 pretrained with DINOv2⁷⁸ on an in-house dataset. CHIEF⁷⁴ applies attention-based feature aggregation, trained via slide-level contrastive learning and anatomic site information. The patch encoder is based on CTransPath⁴, a self-supervised SwinTransformer¹²¹ trained on 15 million patches. In addition to the pretrained slide encoders, we evaluate mean pooling as baseline, where the patch features are averaged within each slide, as it serves as a strong unsupervised baseline despite its simplicity⁶⁴⁻⁶⁶. While we mainly compare with mean pooling based on CONCHv1.5 patch features, we also provide results for mean pooling with the corresponding patch encoders of each slide encoder for a subset of analyses.

Supervised baselines. We compare TITAN against attention-based MIL (ABMIL)^{45,67} and finetuning of the pretrained slide encoders. For ABMIL, the model was trained with a batch size of 1 using the AdamW optimizer with weight decay 10^{-5} and a Cosine annealing learning rate scheduler with peak learning rate 10^{-4} over 20 epochs. The patch encoders were selected accordingly for each analysis. For GigaPath finetuning, we used the publicly available code, which uses a batch size of 1, AdamW optimizer with weight decay 0.05,

and Cosine annealing learning rate scheduler with warm-up and base learning rate $2 \cdot 10^{-3}$ over 5 epochs. For CHIEF finetuning, we also used the publicly available finetuning code. For tasks with a validation set, the best model is chosen based on the validation loss.

Cross-modal baselines. For cross-modal zero-shot retrieval and clinical report generation, we compare TITAN against PRISM⁶².

Linear and K-nearest neighbors probe evaluation

To evaluate the transfer capabilities and representation quality of slide encoders, we adopt recent work in representation learning with self-supervised frameworks and perform linear (logistic regression) and k -nearest neighbor (k -NN) probing. For linear probing, we minimize cross-entropy loss using the scikit-learn L-BFGS solver with ℓ_2 -regularization, selecting ℓ_2 from 45 logarithmically spaced values between 10^{-6} and 10^5 based on the validation loss. The maximum number of L-BFGS iterations is set to 500. For datasets without validation set, e.g., in small datasets or in few-shot experiments, we choose default values of $\ell_2 = 1$ with 1,000 iterations. We additionally evaluated with k -NN probing, a non-parametrized measure to quantify the representation quality of fixed embeddings. We apply it in two settings: First, we follow SimpleShot to create a prototypical class representation by averaging all slide embeddings per diagnostic class¹⁰⁰. Second, we use the scikit-learn implementation of k -NN with $k = 20$ following stability observations from self-supervised learning literature^{78,122}. In both settings, Euclidean distance is used as the distance metric based on the centered and normalized slide embeddings.

Slide retrieval

To further evaluate the representation quality of different slide encoders, we perform content-based slide retrieval using slide-level classification datasets, where we retrieve slides with the same class label as a given query slide. Specifically, we extract slide features for all WSIs. The training and validation sets are combined to serve as the database of candidate slides (keys), and we treat each slide in the test set as a query slide. Prior to retrieval, we preprocess both keys and queries to center the slide embeddings by subtracting their Euclidean centroid, followed by ℓ_2 normalization. The similarity between the query and each candidate in the database is computed using the ℓ_2 distance metric, where a smaller distance indicates a higher similarity. The retrieved slides are then sorted based on their similarities to the query. The class labels are used to evaluate the retrieval performance using $\text{Acc}@K$ for $K \in \{1, 3, 5\}$, which measures whether at least one of the top K retrieved slides shared the same class label as the query, and $\text{MVAcc}@5$, which considers the majority class label among the top 5 retrieved slides. Detailed descriptions of these metrics are provided in the section **Evaluation Metrics**.

Cross-modal retrieval

Leveraging the vision-language aligned embedding space, we also evaluate cross-modal retrieval performance on TCGA-Slide-Reports. Specifically, we assess both slide-to-report and report-to-slide retrieval tasks. All slides and reports are embedded into a shared space using the vision and text encoders, respectively, followed by ℓ_2 normalization. Retrieval is performed by calculating pairwise cosine similarity between the slide and report embeddings. Our class-based approach mirrors the uni-modal slide retrieval, where retrieval is successful if the retrieved slide or report belongs to the same diagnostic class as the query. Performance is quantified using Recall@K for $K \in \{1, 3, 5, 10\}$ for the class-based approach, which measures the proportion of queries for which the correct result appears among the top-K retrieved items. Additionally, we report the mean recall, computed as the average of the Recall@K values across the four K levels. Further details on these metrics are provided in Section **Evaluation Metrics**.

Few-shot slide classification

We evaluate few-shot classification by varying the number of shots k in $\{1, 2, 4, 8, 16, 32\}$. For each k , we select k shots per class or all samples per class if the class has less than k samples. We follow previous studies that used the SimpleShot¹⁰⁰ framework for evaluation of the few-shot learning performance of self-supervised models⁹. SimpleShot computes a prototypical representation per class by averaging all samples within that class. The distances to the class prototypes are then computed on the test set. All embeddings are centered and normalized based on the few-shot samples. To make the evaluation better comparable to supervised baselines such as ABMIL, we additionally assess few-shot classification with linear probing. As no validation set is available in few-shot experiments, we use the default scikit-learn recipe with regularization strength $\ell_2 = 1$ and up to 1,000 iterations of the L-BFGS solver.

Survival analysis

For survival analysis, we employed the linear Cox proportional hazards model on the disease-specific survival (DSS) clinical endpoint. We note that this is different from typical MIL survival prediction with negative log likelihood^{65, 123}, as we deal with a single embedding for the slide (as opposed to a bag of patch embeddings) and patients can be batched (as opposed to the single patient per batch due to memory usage). To reduce the impact of batch effects, we performed a five-fold site-preserved stratification¹²⁴. Due to the small cohort size for reliable survival prediction modeling, we used four folds for training and the remaining fold for evaluation, without employing the validation fold. A hyperparameter α was searched over 25 logarithmically spaced values between 10^1 and 10^5 , with the ℓ_2 coefficient defined as $C = \alpha$. For each combination of encoder and cancer type, we chose C that yielded the best average test metric across the five folds. For fitting and testing the Cox model, we used the scikit-surv package.

Zero-shot slide classification

For zero-shot slide classification, we adopted the method described in CLIP¹⁰⁴ to use the similarities between a given slide and the text prompts of each class as its prediction logits. Specifically, for a class $c \in \{1, 2, \dots, C\}$, we first created the text prompts for each class, followed by extracting their ℓ_2 -normalized text embeddings \mathbf{v}_c using the text encoder. Since the model could be sensitive to the specific choice of text prompts, we created an ensemble of prompts for each class. The complete set of prompt ensembles are provided in **Extended Data Table 101**. For each WSI, we similarly computed a ℓ_2 -normalized embedding \mathbf{u}_i using the slide encoder. We then calculated the cosine similarity between the slide embedding and each class text embedding. The predicted class for a slide was the one with the highest cosine similarity score:

$$\hat{y}_i = \operatorname{argmax}_c \mathbf{u}_i^T \mathbf{v}_c \quad (2)$$

Report generation

Slide captioning provides concise and interpretable summaries of visual findings in pathology, potentially enhancing clinical workflows. The generative objective of CoCa enabled the model’s capabilities of generating pathological reports, which we explored on 10,108 slide-report pairs from TCGA. We performed zero-shot captioning using TITAN and compared the quality of the generated report against PRISM⁶². Specifically, we use a beam search decoding strategy with 5 beams and 1 beam group, where the model explores five potential sequences at each step and retains only the most likely sequence within a single group to maximize quality while minimizing redundancy.

Evaluation metrics

We report balanced accuracy and weighted F1-score for all classification tasks with more than two classes. For ordinal multiclass classification tasks, we report balanced accuracy and quadratic weighted Cohen’s κ . For binary classification tasks, we report balanced accuracy and area-under-the-receiver-operator-curve (AUROC). For survival tasks, we report the concordance index (c-index), which measures the agreement between the model’s predicted risks and the actual survival times. For slide retrieval tasks, we report $\text{Acc}@K$ for $K \in 1, 3, 5$, which measures if at least one slide among the top K retrieved slides has the same class label as the query. We also report $\text{MVAcc}@5$, which is a more strict metric that considers whether the majority vote of the top 5 retrieved slides is in the same class as the query. For cross-modal retrieval tasks, we report $\text{Recall}@K$ for $K \in 1, 3, 5, 10$, which measures the proportion of queries for which the correct result appears in the top- K retrieved items. We also report mean recall, which is calculated as the average of the four $\text{Recall}@K$ values. For report generation, we compare the generated reports with the ground truth pathologi-

cal reports using METEOR, ROUGE, and BLEU. METEOR¹⁰⁶ is a metric that evaluates text quality through unigram matching by considering both precision and recall while also accounting for synonyms, stemming, and word order between the candidate and reference texts. ROUGE¹⁰⁷ compares the overlap of n-grams, word sequences, and word pairs between the generated and reference texts, focusing on recall. We use, ROUGE-1, which specifically measures the overlap of unigrams. BLEU¹⁰⁸ measures the quality of generated text based on unigram overlap, focusing on precision. We use BLEU-1, which evaluates the extent of word-level matches between the generated and reference texts.

Statistical analysis

For the datasets with five-fold splits, where we employ 5-fold cross-validation, we report the mean performance and the standard deviations across all folds. For the datasets with a single split, we use non-parametric bootstrapping with 1,000 samples to calculate the mean and standard deviation.

Downstream evaluation datasets

For the evaluation of TITAN on a diverse set of downstream tasks (**Extended Data Tables 9, 10 and 12**, we re-arrange the pre-extracted CONCHv1.5 features from 512×512 patches to feature grids cropped around the tissue regions of the WSIs. Additionally, background masks are created to mask out features corresponding to background patches. Each WSI is then one single input image to TITAN. For downstream tasks with patient-level annotations, we create the patient embeddings by averaging all slide embeddings of TITAN corresponding to one patient. In the following, we detail all datasets used in our downstream evaluations including splits and targets. We first describe the five datasets that we introduce in our study, TCGA-UniformTumor-8K, TCGA-OncoTree, TCGA-Slide-Reports, Rare-Cancer, and Rare-Cancer-Public, followed by existing datasets in alphabetical order.

TCGA-UniformTumor-8K (TCGA-UT-8K) is a pan-cancer subtyping dataset at region-level consisting of 25,495 $8,192 \times 8,192$ pixel cancer regions of 9,662 H&E FFPE diagnostic histopathology WSIs from TCGA. The tumor regions were manually annotated by two expert pathologist, with slide exclusion due to poor staining, poor focus, lacking cancerous regions and incorrect cancer types. Approximately three representative tumor regions per WSI were annotated with pixel-level contours. For each contour, we center cropped a $8,192 \times 8,192$ image region in order to contain both the dense tumor and surrounding tissue context. We split the regions into train-val-test split (13,853:3,434:8,208 slides) preserving the source site. Refer to **Extended Data Table 1** for a detailed overview of all classes contained in this dataset.

TCGA-OncoTree (TCGA-OT) is a pan-cancer subtyping dataset of 11,186 H&E FFPE diagnostic histopathology WSIs from TCGA⁸⁹. All WSIs are classified into 46 classes according to the OncoTree classification

system such that every class is represented by at least 50 samples. We select all diagnostic H&E FFPE WSIs from TCGA with primary tumors. Concretely, we exclude frozen tissue slides, slides without magnification information, metastatic or recurrent tumor slides, slides without tumor tissue, and IHC slides. For training and evaluation, we split the dataset into training-validation-test folds of 8,226:1,612:1,348 samples while preserving the source site, i.e., all slides from one source site are in one split. Refer to **Extended Data Table 2** for a detailed overview of all classes contained in this dataset.

TCGA-Slide-Reports is a pan-cancer slide-report dataset of H&E FFPE diagnostic histopathology WSIs from TCGA⁸⁹. The dataset consists of 10,108 WSIs with paired pathological reports at slide-level. The dataset is built on the TCGA-Reports dataset, which consists of 9,523 patient-level reports released by a previous study¹⁰⁵. The dataset TCGA-Reports was created using 11,108 pathology report PDFs, corresponding to 11,010 patients, available on the TCGA data portal. The raw reports were preprocessed by removing 82 patients with multiple reports, 399 patients with non-primary tumors, 72 patients with no survival data, 381 “Missing Pathology” reports, and 212 “TCGA Pathologic Diagnosis Discrepancy Form” reports, resulting in 9,850 reports. Optical character recognition (OCR) was then performed for text extraction from the PDFs, followed by the removal of “Consolidated Diagnostic Pathology Form” reports, “Synoptic Translated” forms, within-report TCGA metadata insertions, and clinically irrelevant reports, resulting in 9,523 patient-level reports. While these reports are clean and clinically relevant, they often contain descriptions of multiple tissue blocks per patient. This lack of one-to-one mapping between slides and reports poses a challenge for slide-level report generation and cross-modal retrieval, which require distinct slide-to-report alignment. Since block IDs are unavailable in TCGA metadata, we used the slide-level diagnoses to map diagnoses in each tissue block description. Specifically, if a block’s diagnosis matched the slide-level diagnosis, we designated it as corresponding to the slide. This process was automated with GPT4o-mini, producing a final set of 10,108 slide-report pairs. These paired slides are all H&E FFPE WSIs from primary tumors adhering to the same exclusion criteria as mentioned for TCGA-OT. We excluded all frozen tissue slides, slides without magnification information, metastatic or recurrent tumor slides, slides without tumor tissue, and IHC slides. Refer to **Extended Data Table 7** for a detailed overview of the diagnosis distribution of this dataset.

Rare-Cancer-Public is a pan-cancer dataset of H&E FFPE diagnostic WSIs from TCGA⁸⁹. The dataset consists of 1,982 WSIs, with 1,548 WSIs from TCGA and 434 WSIs from EBRAINS, representing 28 rare cancer types. According to the National Institute of Health, rare cancers are defined as those occurring in fewer than 15 individuals per 100,000 annually⁴³. The OncoTree codes of WSIs from TCGA and EBRAINS were manually curated for this criterion by two expert pathologists (A.K., D.F.K.W.). EBRAINS provides more granular diagnostic classifications than the OncoTree codes, enabling the dataset to include finer distinctions for rare brain tumors. We split the dataset into five folds on the patient level. The dataset was divided into five patient-level folds. To assess the retrieval performance for rare cancers within a clinically representative dataset, we

use one fold of the rare cancer dataset as the query set and the remaining folds combined with the common cancer types as a support set. In total, the support and query datasets contain 14,062 slides, including 11,646 WSIs from TCGA and 2,416 from EBRAINS.

Rare-Cancer is an in-house extension of the public dataset Rare-Cancer-Public with MGB internal cases. This dataset comprises 43 rare cancer types and 3,039 H&E FFPE diagnostic histopathology WSIs, where 1,056 additional cases were added from Brigham and Women’s Hospital (BWH). The entire dataset including common cancer types contains 19,626 WSIs with 5,564 WSIs from BWH from 186 OncoTree codes.

BCNB consist of 1,058 H&E FFPE WSIs of early breast cancer core-needle biopsies¹²⁵. All cases are annotated with ER (WT: 227, MUT: 831), PR (WT: 268, MUT: 790), HER2 (WT: 781, MUT: 277) expressions. We split the dataset label-stratified by a ratio of 60:20:20 (676:170:212 slides).

BRACS consists of 547 H&E FFPE WSIs of benign (including normal), atypical, and malignant breast tumors from 189 patients¹²⁶. The cases are annotated in coarse and fine-grained subtypes of three classes (benign tumors: 265, atypical tumors: 89, malignant tumors: 193) and six classes (atypical ductal hyperplasia: 48, ductal carcinoma in situ: 61, flat epithelial atypia: 41, invasive carcinoma: 132, normal: 44, pathological benign: 147, usual ductal hyperplasia: 74). We split the dataset label-stratified on patient level into five splits by a ratio of 60:20:20 (approx. 302:94:151 slides).

Cardiac allograft rejection consists of 5,021 H&E FFPE WSIs of 1,688 patient biopsies collected from BWH²⁴. Each biopsy is labeled for the presence of cardiac rejection characterized by acute cellular rejection (no rejection: 866 patients, rejection: 822 patients) . We split the dataset label-stratified on patient level into train, val, and test splits by ratio 70:10:20 (3547:484:990 slides).

DHMC LUAD consists of 143 H&E FFPE WSIs of lung adenocarcinoma from the Department of Pathology and Laboratory Medicine at Dartmouth-Hitchcock Medical Center (DHMC)¹²⁷. All WSIs are labeled into five classes of the predominant patterns of lung adenocarcinoma (acinar: 59, lepidic: 19, micropapillary: 9, papillary: 5, solid: 51). Given the limited size of the dataset, we use it exclusively for evaluation in a zero-shot setting, where we use the entire dataset as test set.

DHMC RCC consists of 563 H&E FFPE WSIs of renal cell carcinoma (RCC) from DHMC¹²⁸. All slides are labeled into the four predominant patterns of RCC including one benign class (renal oncocytoma, chromophobe RCC, clear cell RCC, papillary RCC). We use the three RCC subtypes as external test set for the three class subtyping task TCGA RCC.

EBRAINS consists of 2,319 H&E FFPE diagnostic histopathology WSIs from the EBRAINS Digital Tumor

Atlas sourced from the University of Vienna¹²⁹. Due to small sample size we exclude two classes and predict a fine-grained 30 class brain tumor subtyping task. All brain tumors in these tasks are designated as rare cancers by the RARECARE project and the NCI-SEER program. For training and evaluation, we approximately label-stratified the dataset into a train-validation-test fold with a 50:25:25 ratio (1,151:595:573 slides). Additionally, we use 873 samples with annotations for isocitrate dehydrogenase (IDH) mutation as external test set for IDH mutation prediction on the TCGA-GBMLGG cohort.

IMP-CRC consists of 5,333 H&E FFPE colorectal biopsy and polypectomy WSIs retrieved from the data archive of IMP Diagnostics laboratory, Portugal¹³⁰⁻¹³². All cases are classified within one of three categories: Non-neoplastic (847 slides), low-grade lesions (2847 slides), i.e., conventional adenomas with low-grade dysplasia, and high-grade lesions (1639 slides), i.e., conventional adenomas with high-grade dysplasia, intramucosal carcinomas, and invasive adenocarcinomas. We split the dataset label-stratified by a ratio of 60:20:20 into train:val:test set (3546:887:900 slides).

MGB-BRCA consists of 1,264 H&E FFPE WSIs of biopsies and resections invasive breast cancers (BRCA) from BWH^{66,133}. Each case is annotated with three IHC status prediction tasks: estrogen receptor (ER) status prediction (0: 261, 1: 613), progesterone receptor (PR) status prediction (0: 37, 1: 504), and human epidermal growth factor receptor 2 (HER2) status prediction (0: 665, 1: 151), where ER, PR, and HER2 status were manually extracted from pathology reports.

MGB-LUAD consists of 1,939 H&E FFPE WSIs of lung adenocarcinoma from BWH^{66,133}. The WSIs are annotated by five molecular tasks with ground truth from IHC: protein 40 (P40) status prediction (0: 113, 1: 72), protein 63 (P63) status prediction (0: 72, 1: 81), Napsin A status prediction (0: 60, 1: 66), caudal type homeobox 2 (CDX2) status prediction (0: 55, 1: 24), and cytokeratin 5 and 6 (CK-5&6) status prediction (0: 29, 1: 29).

MGH-BRCA consists of 1,071 IHC FFPE WSIs of invasive breast breast carcinoma from Mass General Hospital (MGH)⁶⁶. The cases contain annotations for IHC quantification in six expression levels of ER abundance (1: 168, 2: 169, 3: 219, 4: 170, 5: 175, 6: 169) and PR abundance (1: 2603, 2: 2397, 3: 1209, 4: 1118, 5: 1124, 6: 1101).

MUT-HET consists of 1,291 H&E FFPE WSIs of clear-cell RCC (ccRCC), each representing a single patient treated at the Mayo Clinic^{134,135}. All cases are labeled with the following mutations, determined from matched IHC slides: BAP1 mutation (WT: 1130, MUT: 162), PBRM1 mutation (WT: 622, WT: 670), and SETD2 mutation (WT: 943, MUT: 349). We split the dataset into five splits with train:val:test ratio of 60:20:20 (774:258:259 slides) in each split.

OT-108 is an in-house pan-cancer subtyping dataset consisting of 5,564 H&E FFPE diagnostic WSIs from BWH classified into 108 classes according to the OncoTree classification⁹⁶. We split the dataset into train-val-test (3,164:780:1,620 slides). The test set is balanced across the classes and contains 15 slides per class.

PANDA consists of 10,616 H&E FFPE diagnostic histopathology WSIs of core needle biopsies of prostate cancer sourced from the Radboud University Medical Center and the Karolinska Institute. Each slide is assigned a score recommended by the International Society of Urological Pathology (ISUP) that defines prostate cancer grade (6-class grading task). For quality control, we follow prior work¹³⁶ in excluding slides which were erroneously annotated or had noisy labels resulting in overall 9,555 slides (grade 0: 2,603, grade 1: 2,399, grade 2: 1,209, grade 3: 1,118, grade 4: 1,124, grade 5: 1,102). For training and evaluation, we label-stratified PANDA into 80:10:10 train-validation-test folds (7,645:954:953 slides).

Renal allograft rejection consists of 4,847 H&E FFPE WSIs of renal allograft biopsies from 1,118 patients collected at BWH between 2013 and 2022. Each case has associated labels for antibody-mediated rejection (AMR) status (AMR: 286 patients, no AMR: 832 patients), cellular-mediated rejection (cellular rejection: 341, no cellular rejection: 777), and interstitial fibrosis and tubular atrophy (IFTA) status (advanced IFTA: 162 patients, mild IFTA: 706 patients, moderate IFTA: 250 patients). We split the dataset label-stratified into train:val:test set (3002:376:824 slide).

TCGA BRCA consists of 1,049 invasive breast carcinoma (BRCA) H&E FFPE diagnostic histopathology WSIs from TCGA. The WSIs are classified into two classes: invasive ductal carcinoma (IDC) and invasive lobular carcinoma (ILC).

TCGA NSCLC consists of 1,043 H&E FFPE diagnostic histopathology WSIs from TCGA of 946 patients with non-small cell lung cancer (NSCLC). The WSIs are classified into two classes: lung adenocarcinoma (LUAD, 531 slides) and lung squamous cell carcinoma (LUSC, 512 slides). We split the dataset into 5-fold cross validation, stratified by labels with ratio 60:20:20 (e.g., 659:191:193 for fold 0). **CPTAC-NSCLC** serves as external dataset with 1,091 H&E FFPE diagnostic histopathology WSIs from CPTAC of 422 patients with NSCLC.

TCGA LUAD consists of 524 H&E FFPE diagnostic histopathology WSIs from TCGA of 462 patients with lung adenocarcinoma (LUAD). We predict the mutations in the genes *EGFR* (wildtype (WT): 404 patients, mutated (MUT): 58 patients), *KRAS* (WT: 317, MUT: 145), *STK11* (WT: 391, MUT: 71), and *TP53* (WT: 222, MUT: 240). We split the dataset into 5-fold cross validation, stratified by labels with ratio 60:20:20 (e.g., 659:191:193 for fold 0). **CPTAC-LUAD** serves as external dataset with 324 H&E FFPE diagnostic histopathology WSIs from CPTAC of 108 patients with LUAD.

TCGA CRC consists of 549 H&E FFPE diagnostic histopathology WSIs from TCGA of 543 patients with colorectal cancer (CRC). We predict the presence of microsatellite instability (MSI: 61 patients, microsatellite stable (MSS): 353 patients), mutations in the genes *BRAF* (WT: 429 patients, MUT: 58 patients) and *KRAS* (WT: 286 patients, MUT: 201 patients), and tumor staging (T1: 16 slides, T2: 97 slides, T3: 372 slides, T4: 64 slides). **CPTAC-COAD** with 107 H&E FFPE diagnostic histopathology WSIs from CPTAC of 103 patients with colon adenocarcinoma serves as external validation dataset for all tasks (MSI: 24 patients, MSS: 79 patients, *BRAF* WT: 16 patients, *BRAF* MUT: 87 patients, *KRAS* WT: 36 patients, *KRAS* MUT: 58 patients, T2: 17 slides, T2: 77 slides, T4: 13 slides).

TCGA GBMLGG consists of 1,123 H&E FFPE diagnostic histopathology WSIs from TCGA of 558 patients with gliomas, more specifically glioblastomas multiforme and lower-grad gliomas (GBMLGG). The WSIs are classified into two classes: Isocitrate Dehydrogenase (IDH) mutation (425 slides) and no IDH mutation (698 slides). **EBRAINS** serves as an external cohort for this task (IDH MUT: 333 slides, IDH WT 540 slides).

Computing Software and Hardware

We used Python (version 3.9.16) for all experiments and analyses in the study, which can be replicated using open-source libraries as outlined below. We used PyTorch (version 2.0.1, CUDA 11.8) for deep learning model training and inference. To train TITAN_v and TITAN, we modified the public implementation of iBOT (github.com/bytedance/ibot) and CoCa (github.com/mlfoundations/open_clip). We used four and eight \times 80GB NVIDIA A100 GPUs configured for multi-GPU training using distributed data-parallel (DDP) for TITAN_v and TITAN training, respectively. All downstream experiments were conducted on single 24GB NVIDIA 3090 GPUs. All WSI processing was supported by OpenSlide (version 4.3.1), openslide-python (version 1.2.0), and CLAM (github.com/mahmoodlab/CLAM). We used Scikit-learn (version 1.2.2) for its implementation of K-Nearest Neighbors, and the logistic regression implementation and SimpleShot implementation provided by the LGSSL codebase (github.com/mbanani/lgssl). For survival tasks, we used scikit-survival (Version 0.23.1). Implementations of other slide encoders benchmarked in the study are found at the following links: GigaPath (github.com/prov-gigapath/prov-gigapath), PRISM ([huggingface.coco/paige-ai/Prism](https://huggingface.co/paige-ai/Prism)), and CHIEF (github.com/hms-dbmi/CHIEF). For training weakly-supervised ABMIL models, we adapted the training scaffold code from the CLAM codebase (github.com/mahmoodlab/CLAM). Matplotlib (version 3.8.4) and Seaborn (version 0.13.2) were used to create plots and figures. Usage of other miscellaneous Python libraries is listed in the **Reporting Summary**.

Data availability

GTEx data used in pretraining can be accessed through the GTEx portal (<https://www.gtexportal.org/home/>).

For benchmarks, TCGA and CPTAC data can be accessed through the NIH genomic data commons (<https://portal.gdc.cancer.gov>) and proteomics data commons (<https://proteomic.datacommons.cancer.gov>) respectively. Coordinates and labels of TCGA-UniformTumor-8K dataset is made publicly available in the TITAN GitHub repository. All other publicly-available datasets benchmarked in this work can be accessed in their respective data portals: EBRAINS (<https://doi.org/10.25493/WQ48-ZGX>), DHMC-RCC (<https://bmirids.github.io/KidneyCancer>), DHMC-LUAD (<https://bmirids.github.io/LungCancer/>), BRACS (<https://bracs.icar.cnr.it>), PANDA (<https://panda.grand-challenge.org>), IMP (<https://rdm.inesctec.pt/dataset/nis-2023-008>), BCNB (<https://bupt-ai-cz.github.io/BCNB/>), MUT-HET-RCC (<https://aacrjournals.org/cancerres/article/82/15/2792/707325/Intratumoral-Resolution-of-Driver-Gene-Mutation>). Links for all public datasets are also presented in **Extended Data Table 13**.

Code availability

Code and model weights for loading both TITAN and TITAN_v can be accessed for academic research purposes at <https://github.com/mahmoodlab/TITAN>.

Ethics Statement The retrospective analysis of internal pathology images and associated reports used in this study received approval from Mass General Brigham institutional review board. Prior to the computational analysis and model development, all internal digital data, including whole slide images (WSIs), pathology reports, and electronic medical records, were anonymized. Since the study did not involve direct patient participation or recruitment, informed consent was waived for the analysis of archival pathology slides.

Author Contributions

T.D., S.J.W., A.H.S, R.J.C., F.M. conceived the study and designed the experiments. L.P.L., T.D., R.J.C., B.C. curated the Mass-340K whole-slide images and corresponding pathology reports. R.J.C., S.J.W., A.H.S., A.J.V., G.J., C.A.P., P.D. scanned the whole-slide images. T.D., S.J.W., R.J.C., A.H.S developed the stage 1 vision-only TITAN model. T.D., R.J.C, M.Y.L., S.J.W., A.H.S. developed the stage 2 and stage 3 vision-language TITAN models. T.D., S.J.W., and M.Y.L. developed the codebase for zero-shot vision-language slide understanding. T.D., S.J.W., A.H.S., A.Z., A.J.V., G.J. implemented the benchmarking codebase for pretrained slide models. A.K., D.F.K.W. evaluated the synthetic captions and generated reports and helped with the study design for slide retrieval. A.K., D.F.K.W., C.C. curated the rare disease retrieval dataset. R.J.C., D.K., A.K., S.I. curated and annotated the TCGA-Uniform-8K dataset. S.S. curated the renal allograft rejection dataset. T.D., S.J.W., A.H.S, R.J.C., F.M. prepared the manuscript. All authors contributed to the writing. L.P.L., F.M. supervised the research.

Acknowledgements

This work was funded in part by the Brigham and Women's Hospital (BWH) President's Fund, Mass General Hospital (MGH) Pathology and by the National Institute of Health (NIH) National Institute of General Medical Sciences (NIGMS) through R35GM138216. S.J.W. was supported by the Helmholtz Association under the joint research school "Munich School for Data Science - MUDS" and the Add-on Fellowship of the Joachim Herz Foundation. This work was supported by a fellowship of the German Academic Exchange Service (DAAD). M.Y.L. was supported by the Tau Beta Pi Fellowship and the Siebel Foundation. This work was additionally supported by the AMED Practical Research for Innovative Cancer Control under grant number JP 24ck0106873 to SI. The content is solely the responsibility of the authors and does not reflect the official views of the NIH, NIGMS, NCI, DoD.

References

1. Song, A. H. *et al.* Artificial intelligence for digital and computational pathology. *Nature Reviews Bio-engineering* **1**, 930–949 (2023).
2. Riasatian, A. *et al.* Fine-tuning and training of densenet for histopathology image representation using tcga diagnostic slides. *Medical image analysis* **70**, 102032 (2021).
3. Ciga, O., Xu, T. & Martel, A. L. Self supervised contrastive learning for digital histopathology. *Machine Learning with Applications* **7**, 100198 (2022).
4. Wang, X. *et al.* Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical image analysis* **81**, 102559 (2022).
5. Wang, X. *et al.* Retccl: Clustering-guided contrastive learning for whole-slide image retrieval. *Medical image analysis* **83**, 102645 (2023).
6. Kang, M., Song, H., Park, S., Yoo, D. & Pereira, S. Benchmarking self-supervised learning on diverse pathology datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3344–3354 (2023).
7. Huang, Z., Bianchi, F., Yuksekgonul, M., Montine, T. J. & Zou, J. A visual–language foundation model for pathology image analysis using medical twitter. *Nature medicine* **29**, 2307–2316 (2023).
8. Azizi, S. *et al.* Robust and data-efficient generalization of self-supervised machine learning for diagnostic imaging. *Nature Biomedical Engineering* **7**, 756–779 (2023).
9. Chen, R. J. *et al.* Towards a general-purpose foundation model for computational pathology. *Nature Medicine* **30**, 850–862 (2024).
10. Lu, M. Y. *et al.* A visual-language foundation model for computational pathology. *Nature Medicine* **30**, 863–874 (2024).
11. Vorontsov, E. *et al.* A foundation model for clinical-grade computational pathology and rare cancers detection. *Nature Medicine* 1–12 (2024).
12. Zimmermann, E. *et al.* Virchow 2: Scaling self-supervised mixed magnification models in pathology. *arXiv preprint arXiv:2408.00738* (2024).
13. Saillard, C. *et al.* H-optimus-0 (2024). URL <https://github.com/bioptimus/releases/tree/main/models/h-optimus/v0>.
14. Filiot, A. *et al.* Scaling self-supervised learning for histopathology with masked image modeling. *medRxiv* 2023–07 (2023).
15. Filiot, A., Jacob, P., Kain, A. M. & Saillard, C. Phikon-v2, a large and public feature extractor for biomarker prediction (2024). 2409.09173.
16. Juyal, D. *et al.* Pluto: Pathology-universal transformer. *arXiv preprint arXiv:2405.07905* (2024).
17. Dippel, J. *et al.* Rudolfov: A foundation model by pathologists for pathologists. *arXiv preprint arXiv:2401.04079* (2024).

18. Bera, K., Schalper, K. A., Rimm, D. L., Velcheti, V. & Madabhushi, A. Artificial intelligence in digital pathology—new tools for diagnosis and precision oncology. *Nature reviews Clinical oncology* **16**, 703–715 (2019).
19. Echle, A. *et al.* Deep learning in cancer pathology: a new generation of clinical biomarkers. *British journal of cancer* **124**, 686–696 (2021).
20. Shmatko, A., Ghaffari Laleh, N., Gerstung, M. & Kather, J. N. Artificial intelligence in histopathology: enhancing cancer research and clinical oncology. *Nature cancer* **3**, 1026–1038 (2022).
21. Campanella, G. *et al.* A clinical benchmark of public self-supervised pathology foundation models. *arXiv preprint arXiv:2407.06508* (2024).
22. Neidlinger, P. *et al.* Benchmarking foundation models as feature extractors for weakly-supervised computational pathology. *arXiv preprint arXiv:2408.15823* (2024).
23. Yu, K.-H. *et al.* Predicting non-small cell lung cancer prognosis by fully automated microscopic pathology image features. *Nature communications* **7**, 12474 (2016).
24. Lipkova, J. *et al.* Deep learning-enabled assessment of cardiac allograft rejection from endomyocardial biopsies. *Nature medicine* **28**, 575–582 (2022).
25. Coudray, N. *et al.* Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nature medicine* **24**, 1559–1567 (2018).
26. Campanella, G. *et al.* Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine* **25**, 1301–1309 (2019).
27. Kather, J. N. *et al.* Deep learning can predict microsatellite instability directly from histology in gastrointestinal cancer. *Nature medicine* **25**, 1054–1056 (2019).
28. Lu, M. Y. *et al.* AI-based pathology predicts origins for cancers of unknown primary. *Nature* **594**, 106–110 (2021).
29. Fu, Y. *et al.* Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis. *Nature cancer* **1**, 800–810 (2020).
30. Bulten, W. *et al.* Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge. *Nature medicine* **28**, 154–163 (2022).
31. Zheng, Y. *et al.* A graph-transformer for whole slide image classification. *IEEE transactions on medical imaging* **41**, 3003–3015 (2022).
32. Skrede, O.-J. *et al.* Deep learning for prediction of colorectal cancer outcome: a discovery and validation study. *The Lancet* **395**, 350–360 (2020).
33. Bejnordi, B. E. *et al.* Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama* **318**, 2199–2210 (2017).
34. Wagner, S. J. *et al.* Transformer-based biomarker prediction from colorectal cancer histology: A large-scale multicentric study. *Cancer Cell* **41**, 1650–1661 (2023).
35. Foersch, S. *et al.* Multistain deep learning for prediction of prognosis and therapy response in colorectal cancer. *Nature medicine* **29**, 430–439 (2023).

36. Courtiol, P. *et al.* Deep learning-based classification of mesothelioma improves prediction of patient outcome. *Nature medicine* **25**, 1519–1525 (2019).
37. Lee, Y. *et al.* Derivation of prognostic contextual histopathological features from whole-slide images of tumours via graph deep learning. *Nature Biomedical Engineering* 1–15 (2022).
38. Niehues, J. M. *et al.* Generalizable biomarker prediction from cancer pathology slides with self-supervised deep learning: A retrospective multi-centric study. *Cell reports Medicine* **4** (2023).
39. Yu, K.-H., Beam, A. L. & Kohane, I. S. Artificial intelligence in healthcare. *Nature biomedical engineering* **2**, 719–731 (2018).
40. Krishnan, R., Rajpurkar, P. & Topol, E. J. Self-supervised learning in medicine and healthcare. *Nature Biomedical Engineering* **6**, 1346–1352 (2022).
41. Van der Laak, J., Litjens, G. & Ciompi, F. Deep learning in histopathology: the path to the clinic. *Nature medicine* **27**, 775–784 (2021).
42. Gatta, G. *et al.* Burden and centralised treatment in europe of rare tumours: results of rarecarenet—a population-based study. *The Lancet Oncology* **18**, 1022–1039 (2017).
43. NCI Dictionary of Cancer Terms. Rare cancer. *National Cancer Institute* Accessed: 2024-11-20.
44. Surveillance, Epidemiology, an End Results Program. Rare cancer classification. *National Cancer Institute* Accessed: 2024-11-21.
45. Lu, M. Y. *et al.* Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering* **5**, 555–570 (2021).
46. Lew, M. S., Sebe, N., Djeraba, C. & Jain, R. Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* **2**, 1–19 (2006).
47. Cruz-Roa, A., Caicedo, J. C. & González, F. A. Visual pattern mining in histology image collections using bag of features. *Artificial intelligence in medicine* **52**, 91–106 (2011).
48. Caicedo, J. C., Cruz, A. & Gonzalez, F. A. Histopathology image classification using bag of features and kernel functions. In *Artificial Intelligence in Medicine: 12th Conference on Artificial Intelligence in Medicine, AIME 2009, Verona, Italy, July 18-22, 2009. Proceedings 12*, 126–135 (Springer, 2009).
49. Sridhar, A., Doyle, S. & Madabhushi, A. Content-based image retrieval of digitized histopathology in boosted spectrally embedded spaces. *Journal of pathology informatics* **6**, 41 (2015).
50. Kalra, S. *et al.* Yottixel—an image search engine for large archives of histopathology whole slide images. *Medical Image Analysis* **65**, 101757 (2020).
51. Chen, C. *et al.* Fast and scalable search of whole-slide images via self-supervised deep learning. *Nature Biomedical Engineering* **6**, 1420–1434 (2022).
52. Zheng, Y. *et al.* Histopathological whole slide image analysis using context-based cbir. *IEEE transactions on medical imaging* **37**, 1641–1652 (2018).
53. Shang, H. H. *et al.* Histopathology slide indexing and search—are we there yet? *NEJM AI* **1**, AIcs2300019 (2024).

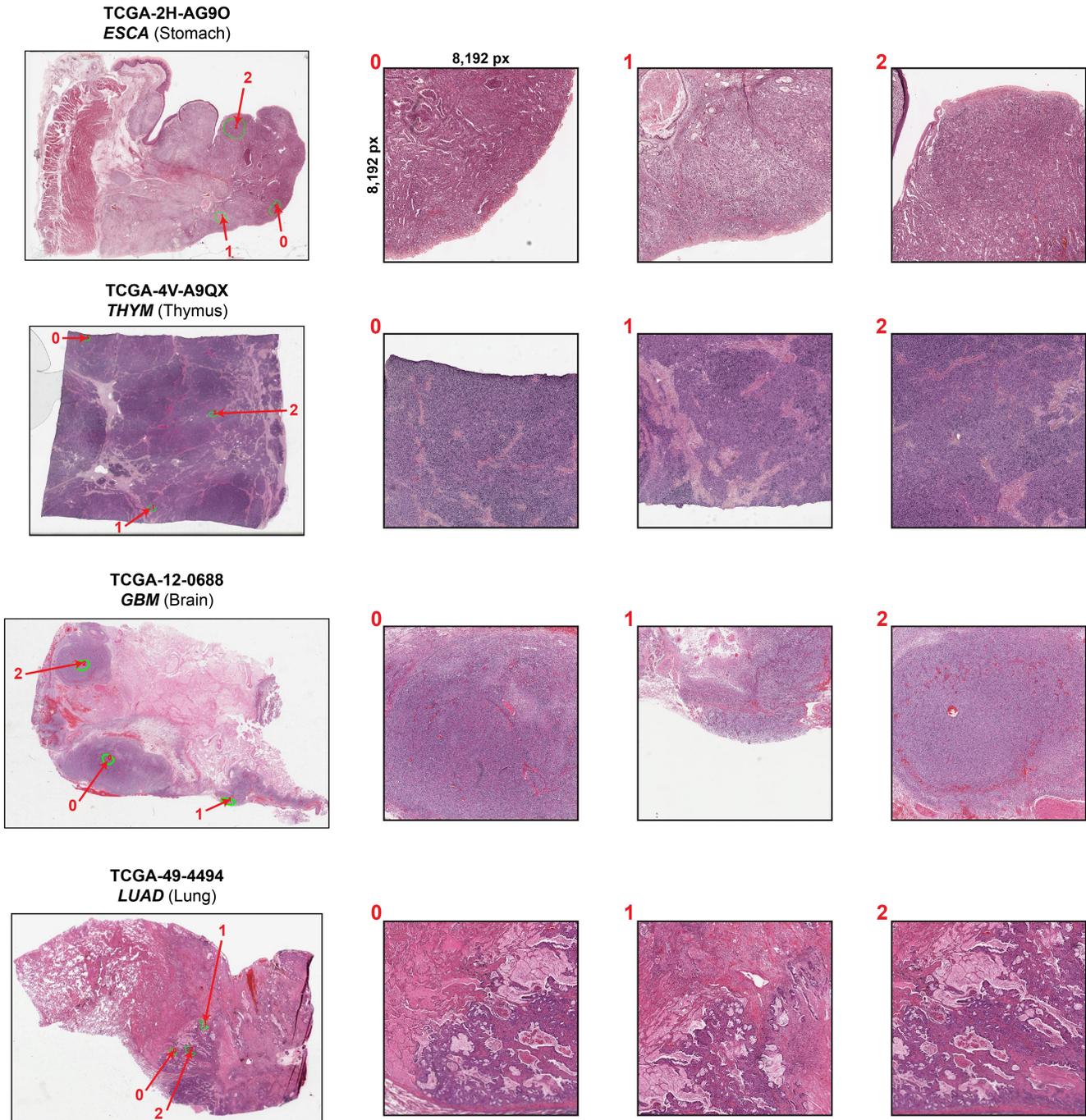
54. Zhang, Z., Xie, Y., Xing, F., McGough, M. & Yang, L. Mdnnet: A semantically and visually interpretable medical image diagnosis network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 6428–6436 (2017).
55. Lu, M. Y. *et al.* Visual Language Pretrained Multiple Instance Zero-Shot Transfer for Histopathology Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 19764–19775 (2023).
56. Ikezogwo, W. *et al.* Quilt-1m: One million image-text pairs for histopathology. *Advances in neural information processing systems* **36** (2024).
57. Chen, R. J. *et al.* Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16144–16155 (2022).
58. Xu, H. *et al.* A whole-slide foundation model for digital pathology from real-world data. *Nature* 1–8 (2024).
59. Lazard, T., Lerousseau, M., Decencière, E. & Walter, T. Giga-ssl: Self-supervised learning for gigapixel images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4304–4313 (2023).
60. Hou, X. *et al.* A self-supervised framework for learning whole slide representations (2024). 2402.06188.
61. Tran, M. *et al.* Generating clinical-grade pathology reports from gigapixel whole slide images with HistoGPT. *medRxiv* 2024–03 (2024).
62. Shaikovski, G. *et al.* PRISM: A Multi-Modal Generative Foundation Model for Slide-Level Histopathology. *arXiv preprint arXiv:2405.10254* (2024).
63. Xu, Y. *et al.* A multimodal knowledge-enhanced whole-slide pathology foundation model. *arXiv preprint arXiv:2407.15362* (2024).
64. Jaume, G. *et al.* Transcriptomics-guided slide representation learning in computational pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9632–9644 (2024).
65. Song, A. H. *et al.* Morphological prototyping for unsupervised slide representation learning in computational pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11566–11578 (2024).
66. Jaume, G. *et al.* Multistain pretraining for slide representation learning in pathology. In *European Conference on Computer Vision*, 19–37 (Springer, 2024).
67. Ilse, M., Tomczak, J. & Welling, M. Attention-based deep multiple instance learning. In *International conference on machine learning*, 2127–2136 (PMLR, 2018).
68. Li, B., Li, Y. & Eliceiri, K. W. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 14318–14328 (2021).

69. Chen, R. J. *et al.* Whole slide images are 2d point clouds: Context-aware survival prediction using patch-based graph convolutional networks. In *Medical Image Computing and Computer Assisted Intervention–MICCAI*, 339–349 (2021).
70. Shao, Z. *et al.* Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in neural information processing systems* **34**, 2136–2147 (2021).
71. Xiang, J. & Zhang, J. Exploring Low-Rank Property in Multiple Instance Learning for Whole Slide Image Classification. In *The Eleventh International Conference on Learning Representations* (2023).
72. Kondepudi, A. *et al.* Foundation models for fast, label-free detection of glioma infiltration. *Nature* 1–7 (2024).
73. Ahmed, F. *et al.* PathAlign: A vision-language model for whole slide images in histopathology (2024). 2406.19578.
74. Wang, X. *et al.* A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature* (2024).
75. Bommasani, R. *et al.* On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).
76. Moor, M. *et al.* Foundation models for generalist medical artificial intelligence. *Nature* **616**, 259–265 (2023).
77. Zhou, J. *et al.* ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832* (2021).
78. Oquab, M. *et al.* Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023).
79. Lu, M. Y. *et al.* A Multimodal Generative AI Copilot for Human Pathology. *Nature* 1–3 (2024).
80. Chen, R. J., Lu, M. Y., Chen, T. Y., Williamson, D. F. & Mahmood, F. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering* **5**, 493–497 (2021).
81. Kokosi, T. & Harron, K. Synthetic data in medical research. *BMJ medicine* **1** (2022).
82. Carrillo-Perez, F. *et al.* Generation of synthetic whole-slide image tiles of tumours from RNA-sequencing data via cascaded diffusion models. *Nature Biomedical Engineering* 1–13 (2024).
83. Dosovitskiy, A. *et al.* An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
84. Bär, A., Houlsby, N., Dehghani, M. & Kumar, M. Frozen feature augmentation for few-shot image classification. *arXiv preprint arXiv:2403.10519* (2024).
85. Press, O., Smith, N. & Lewis, M. Train short, test long: Attention with linear biases enables input length extrapolation. In *International Conference on Learning Representations* (2022).
86. Yu, J. *et al.* Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research* (2022).
87. Yang, A. *et al.* Qwen2 technical report. *arXiv preprint arXiv:2407.10671* (2024).

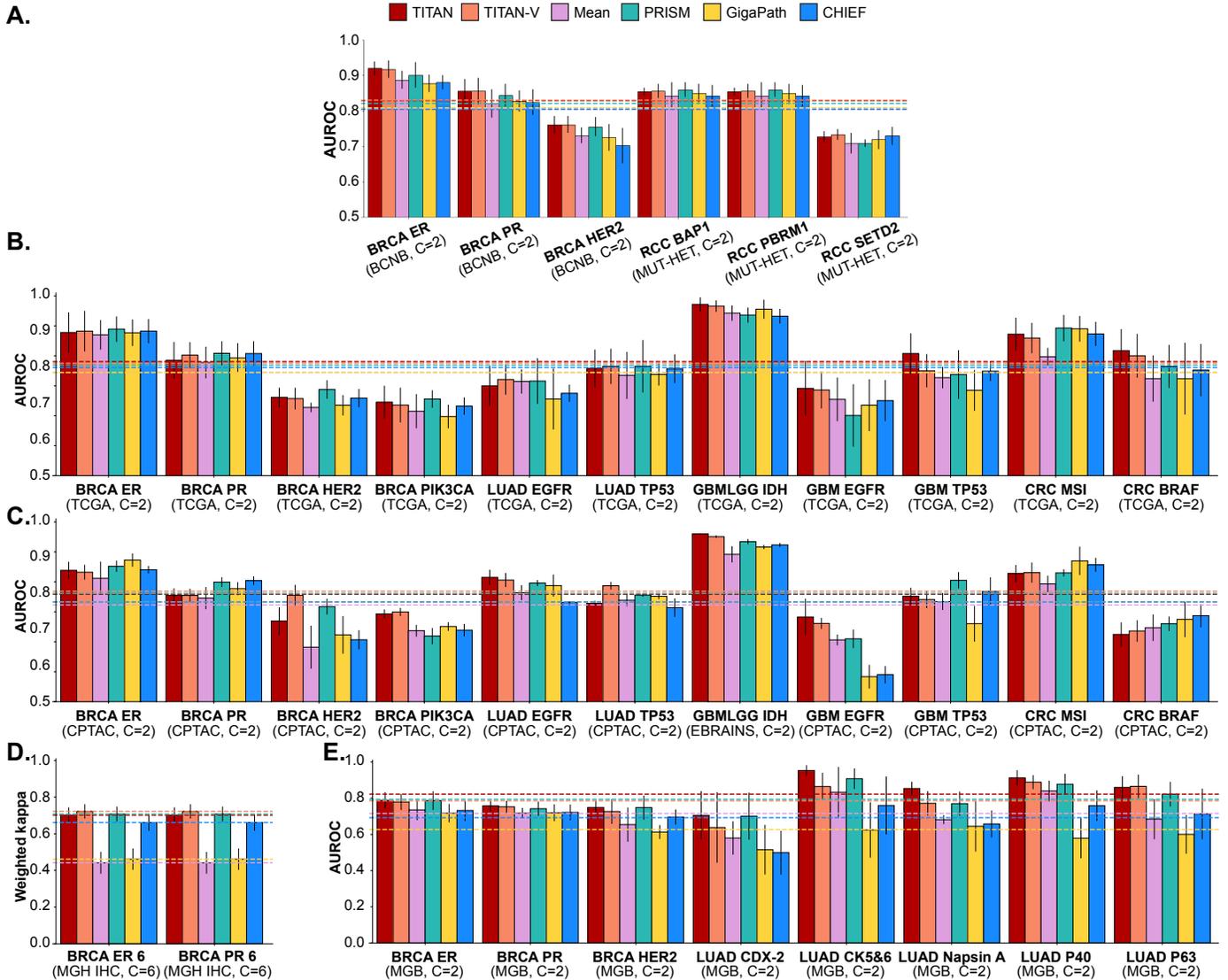
88. Tian, Y., Fan, L., Isola, P., Chang, H. & Krishnan, D. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. *Advances in Neural Information Processing Systems* **36** (2024).
89. Weinstein, J. N. *et al.* The cancer genome atlas pan-cancer analysis project. *Nature genetics* **45**, 1113–1120 (2013).
90. Liu, J. *et al.* An integrated TCGA pan-cancer clinical data resource to drive high-quality survival outcome analytics. *Cell* **173**, 400–416 (2018).
91. Edwards, N. J. *et al.* The CPTAC Data Portal: A Resource for Cancer Proteomics Research. *Journal of Proteome Research* **14**, 2707–2713 (2015).
92. Thangudu, R. R. *et al.* Abstract LB-242: Proteomic Data Commons: A resource for proteogenomic analysis. *Cancer Research* **80** (2020).
93. Wei, J. W. *et al.* Pathologist-level classification of histologic patterns on resected lung adenocarcinoma slides with deep neural networks. *Scientific reports* **9**, 1–8 (2019).
94. Zhu, M. *et al.* Development and evaluation of a deep neural network for histologic classification of renal cell carcinoma on biopsy and surgical resection slides. *Scientific reports* **11**, 1–9 (2021).
95. Song, A. H. *et al.* Multimodal prototyping for cancer survival prediction. In *Forty-first International Conference on Machine Learning*.
96. Kundra, R. *et al.* Oncotree: a cancer classification system for precision oncology. *JCO clinical cancer informatics* **5**, 221–230 (2021).
97. Kapoor, S. & Narayanan, A. Leakage and the reproducibility crisis in machine-learning-based science. *Patterns* **4** (2023).
98. Zhai, X., Kolesnikov, A., Houlsby, N. & Beyer, L. Scaling vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12104–12113 (2022).
99. Kaplan, J. *et al.* Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).
100. Wang, Y., Chao, W.-L., Weinberger, K. Q. & Van Der Maaten, L. Simpleshot: Revisiting nearest-neighbor classification for few-shot learning. *arXiv preprint arXiv:1911.04623* (2019).
101. Snell, J., Swersky, K. & Zemel, R. Prototypical networks for few-shot learning. *Advances in neural information processing systems* **30** (2017).
102. Su, J. *et al.* Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* **568**, 127063 (2024).
103. Heo, B., Park, S., Han, D. & Yun, S. Rotary position embedding for vision transformer. *arXiv preprint arXiv:2403.13298* (2024).
104. Radford, A. *et al.* Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763 (PMLR, 2021).
105. Kefeli, J. & Tatonetti, N. TCGA-Reports: A machine-readable pathology report resource for benchmarking text-based AI models. *Patterns* **5** (2024).

106. Banerjee, S. & Lavie, A. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 65–72 (2005).
107. Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 74–81 (2004).
108. Papineni, K., Roukos, S., Ward, T. & Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318 (2002).
109. Komura, D. & Ishikawa, S. Machine learning methods for histopathological image analysis. *Computational and structural biotechnology journal* **16**, 34–42 (2018).
110. Hegde, N. *et al.* Similar image search for histopathology: SMILY. *NPJ digital medicine* **2**, 56 (2019).
111. Neumann, H. P., Young Jr, W. F. & Eng, C. Pheochromocytoma and paraganglioma. *New England journal of medicine* **381**, 552–565 (2019).
112. Wagner, S. J. *et al.* Make deep learning algorithms in computational pathology more reproducible and reusable. *Nature Medicine* **28**, 1744–1746 (2022).
113. Nechaev, D., Pchelnikov, A. & Ivanova, E. Hibou: A Family of Foundational Vision Transformers for Pathology (2024). 2406.05074.
114. Consortium, G. *et al.* The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
115. Zaffar, I., Jaume, G., Rajpoot, N. & Mahmood, F. Embedding space augmentation for weakly supervised learning in whole-slide images. In *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, 1–4 (IEEE, 2023).
116. Shao, Z., Dai, L., Wang, Y., Wang, H. & Zhang, Y. Augdiff: Diffusion based feature augmentation for multiple instance learning in whole slide image. *arXiv preprint arXiv:2303.06371* (2023).
117. Jaume, G., Song, A. H. & Mahmood, F. Integrating context for superior cancer prognosis. *Nature Biomedical Engineering* **6**, 1323–1325 (2022).
118. Beyer, L. *et al.* PaliGemma: A versatile 3B VLM for transfer. *arXiv preprint arXiv:2407.07726* (2024).
119. He, K. *et al.* Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16000–16009 (2022).
120. Jaegle, A. *et al.* Perceiver: General perception with iterative attention. In *International conference on machine learning*, 4651–4664 (PMLR, 2021).
121. Liu, Z. *et al.* Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, 10012–10022 (2021).
122. Caron, M. *et al.* Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9650–9660 (2021).
123. Zadeh, S. G. & Schmid, M. Bias in cross-entropy-based training of deep survival networks. *IEEE transactions on pattern analysis and machine intelligence* **43**, 3126–3137 (2020).

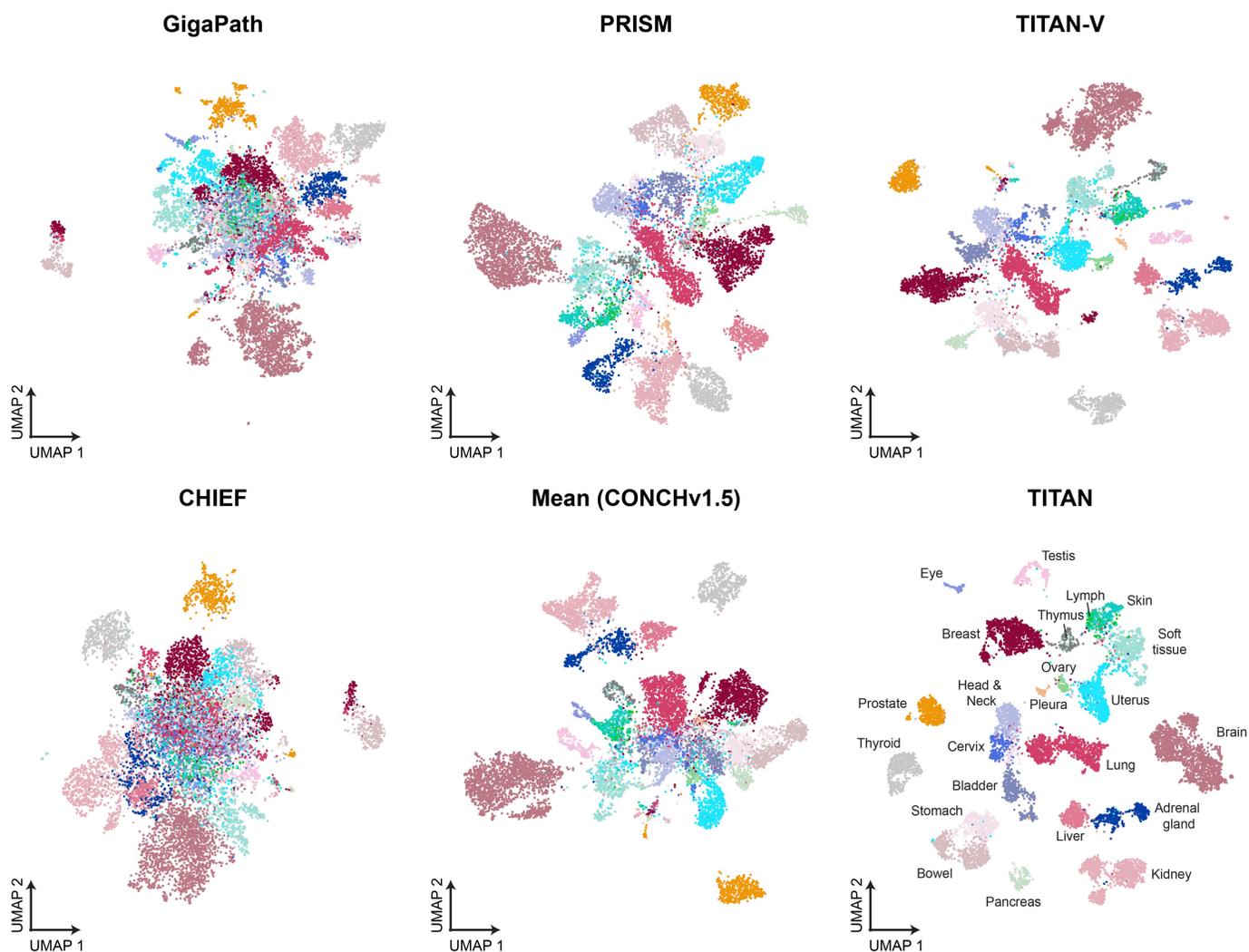
124. Howard, F. M. *et al.* The impact of site-specific digital histology signatures on deep learning model accuracy and bias. *Nature communications* **12**, 4423 (2021).
125. Xu, F. *et al.* Predicting axillary lymph node metastasis in early breast cancer using deep learning on primary tumor biopsy slides. *Frontiers in Oncology* 4133 (2021).
126. Brancati, N. *et al.* Bracs: A dataset for breast carcinoma subtyping in h&e histology images. *Database* **2022**, baac093 (2022).
127. Wei, J. W. *et al.* Pathologist-level classification of histologic patterns on resected lung adenocarcinoma slides with deep neural networks. *Scientific reports* **9**, 1–8 (2019).
128. Zhu, M. *et al.* Development and evaluation of a deep neural network for histologic classification of renal cell carcinoma on biopsy and surgical resection slides. *Scientific reports* **11**, 1–9 (2021).
129. Roetzer-Pejrimovsky, T. *et al.* The digital brain tumour atlas, an open histopathology resource. *Scientific Data* **9**, 55 (2022).
130. Oliveira, S. P. *et al.* CAD systems for colorectal cancer from WSI are still not ready for clinical acceptance. *Scientific Reports* **11**, 1–15 (2021).
131. Neto, P. C. *et al.* iMIL4PATH: A semi-supervised interpretable approach for colorectal whole-slide images. *Cancers* **14**, 2489 (2022).
132. Neto, P. C. *et al.* An interpretable machine learning system for colorectal cancer diagnosis from pathology slides. *npj Precision Oncology* **8**, 56 (2024).
133. Vaidya, A. *et al.* Demographic bias in misdiagnosis by computational pathology models. *Nature Medicine* **30**, 1174–1190 (2024).
134. Joseph, R. W. *et al.* Clear cell renal cell carcinoma subtypes identified by BAP1 and PBRM1 expression. *The Journal of urology* **195**, 180–187 (2016).
135. Acosta, P. H. *et al.* Intratumoral resolution of driver gene mutation heterogeneity in renal cancer using deep learning. *Cancer research* **82**, 2792–2806 (2022).
136. Pati, P. *et al.* Weakly supervised joint whole-slide segmentation and classification in prostate cancer. *Medical Image Analysis* **89**, 102915 (2023).



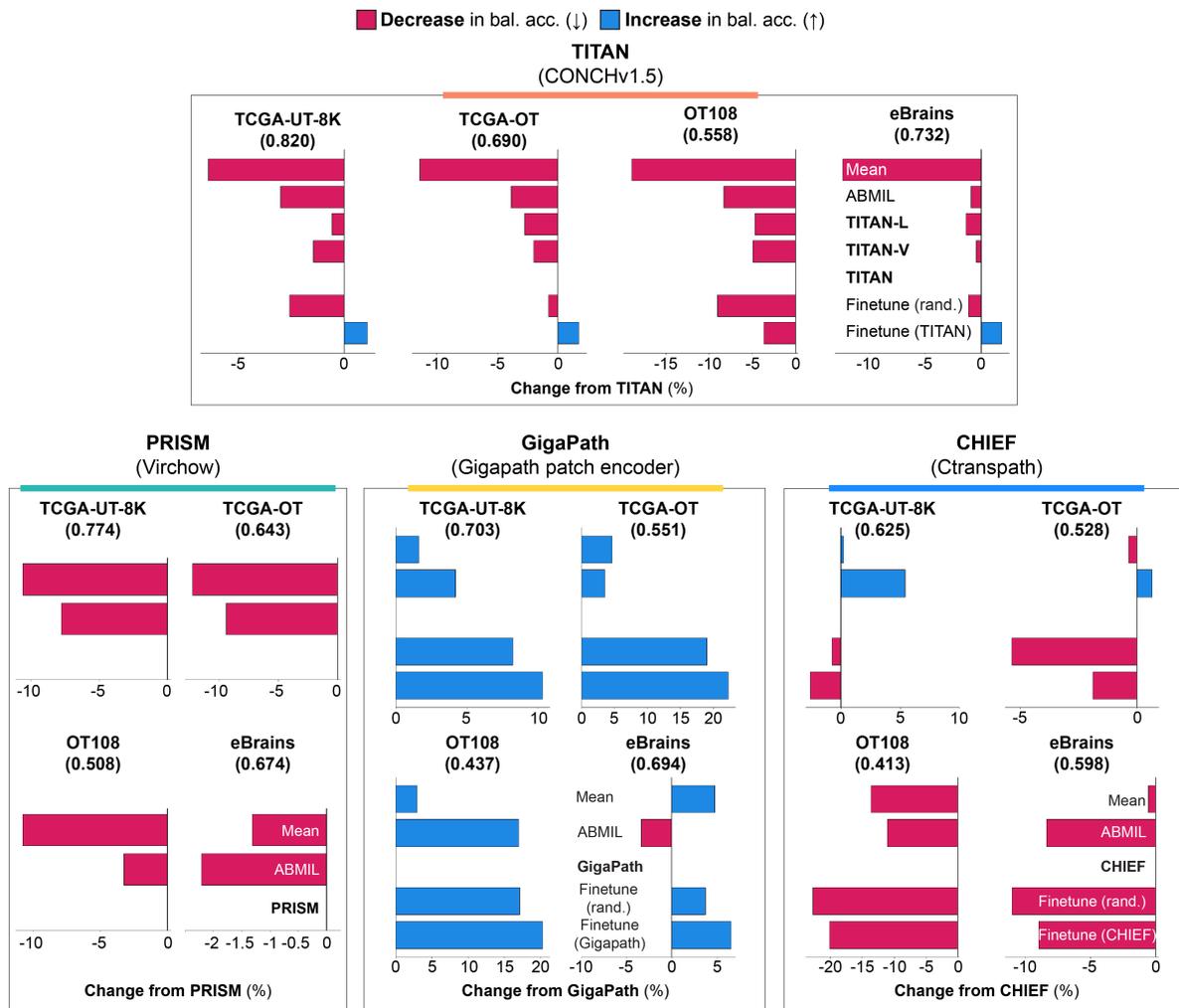
Extended Data Figure 1: Examples of TCGA-UT-8K dataset. Examples of TCGA-UT-8K, which are ROIs of $8,192 \times 8,192$ pixel selected by the pathologists. The green contours illustrate the cancer region annotations, with the red number indicating the ROI index within a given TCGA slide.



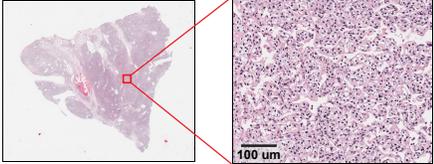
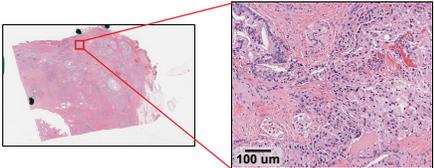
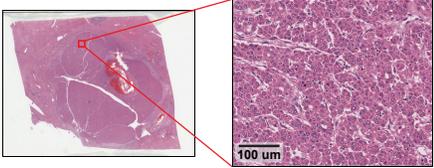
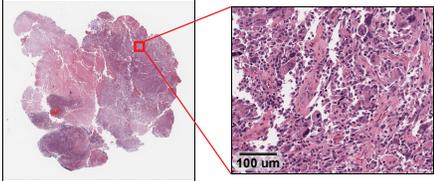
Extended Data Figure 2: Linear probe results for molecular classification tasks. (a) Linear models are fitted and evaluated on binary molecular status predictions for BCNB and MUT-HET. (b) Linear models are fitted and evaluated on five fold-splits on TCGA, (c) the same models are evaluated on the corresponding external datasets from CPTAC and EBRAINS. (d) 6-level ER and PR prediction from immunohistochemistry (IHC) slides from Mass General Hospital (MGH). (e) molecular classification tasks for BRCA and LUAD from Mass General Brigham (MGB).



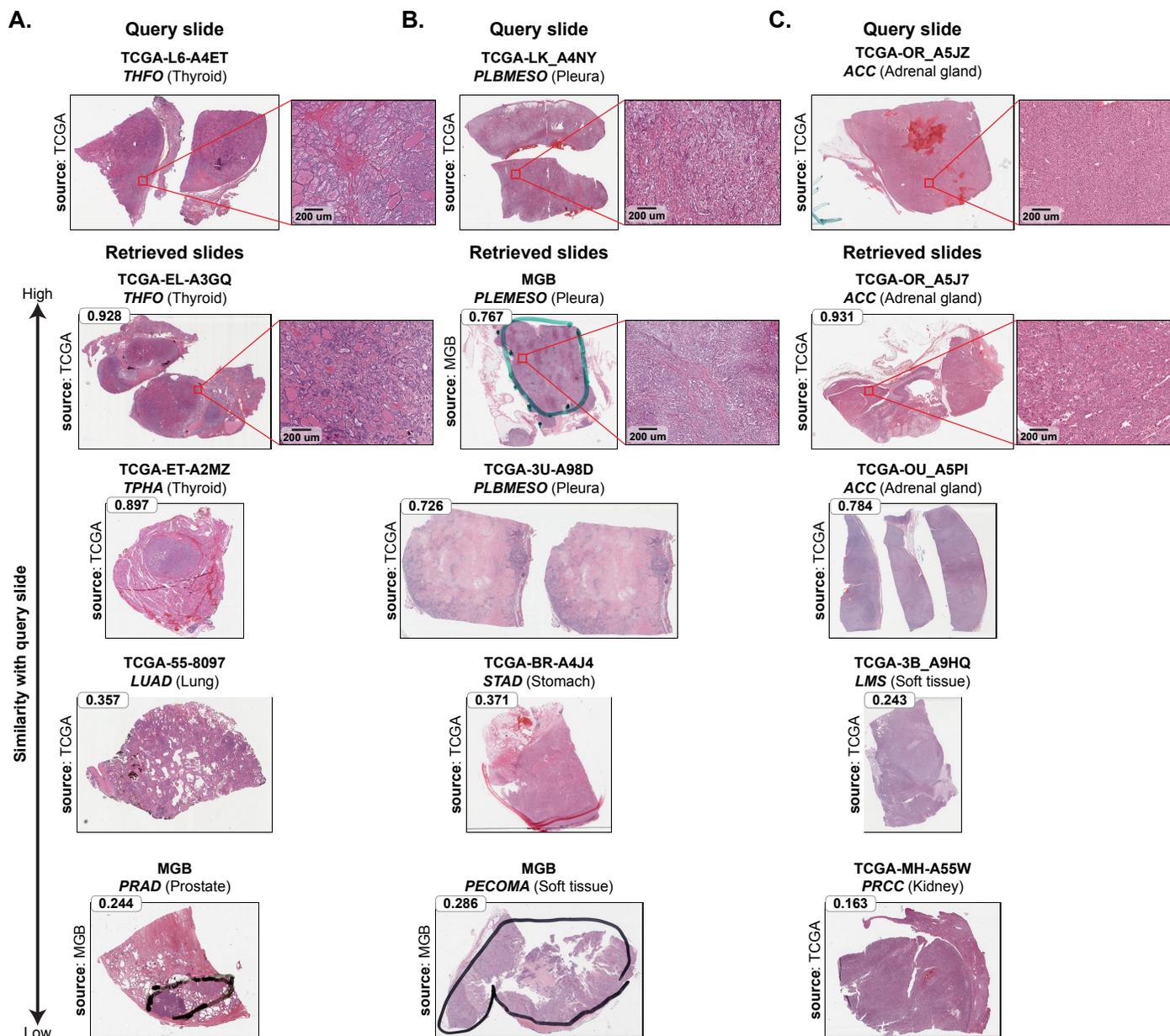
Extended Data Figure 3: UMAP of slide embedding space for TCGA-OT. UMAP visualization of slide embeddings in TCGA-OT cohort for all slide encoder baselines including TITAN and TITAN_V, color-coded by different organs for visual decluttering.



Extended Data Figure 4: Ablation experiments on different learning paradigms. Change in balanced accuracy performance for several learning paradigms on four subtyping tasks with respect to the linear probe. The baselines include mean pooling, ABMIL, linear probe, and finetuned from pretrained or randomly initialized weights. The number under each task name indicates the linear probe performance. TITAN_L represents the variation of TITAN without vision-pretraining. For mean pooling and ABMIL, we use the respective patch encoder for each framework, as specified under each slide encoder name. Finetuning results are not provided for PRISM, as the finetuning recipes were not available.

<p>TCGA-AK-3450 CCRCC (Kidney)</p> 	<p>Clinical report</p> <p>The slide from the kidney shows a renal cell carcinoma, clear cell type, Fuhrman nuclear grade II/IV, confined to renal parenchyma with no angiolymphatic invasion. Surgical resection margins are free of carcinoma.</p>	<p>Ours (generated)</p> <p>clear cell renal cell renal cell carcinoma, fuhrman nuclear grade ii/iv, confined within the renal parenchyma without renal capsule invasion.</p>	<p>PRISM (generated)</p> <p>Examined tissue: Benign adrenal gland.</p>
<p>TCGA-V1-A8WW PRAD (Prostate)</p> 	<p>Clinical report</p> <p>The slide from the prostate gland shows prostatic adenocarcinoma (Gleason score 4+5=9) with bilateral seminal vesicle invasion, extensive extraprostatic extension, and presence of perineural and lymphovascular invasion.</p>	<p>Ours (generated)</p> <p>prostatic adenocarcinoma (gleason score 4+5=9 (grade group 5) involving both sides of the seminal vesicle and extraprostatic soft tissue on the right side, with lymphovascular and perineural invasion, extraprostatic extension, seminal vesicle invasion, and extraprostatic extension on the left side. pik3r1-pik3r1-pik3-tmprss2 rearrangement detected.</p>	<p>PRISM (generated)</p> <p>Diagnosis: Prostatic adenocarcinoma with perineural invasion; organ: prostate.</p>
<p>TCGA-2Y-A9H4 HCC (Liver)</p> 	<p>Clinical report</p> <p>The slide from the liver shows a moderately differentiated hepatocellular carcinoma, 3.3 x 2.5 x 1.5 cm in size, confined to the hepatic parenchyma, with no evidence of vascular or perineural invasion.</p>	<p>Ours (generated)</p> <p>moderately differentiated hepatocellular carcinoma (g2) identified in the liver, confined to the hepatic parenchyma without evidence of vascular or perineural invasion. molecular testing reveals cttnb1 and cttnb1 mutations.</p>	<p>PRISM (generated)</p> <p>Diagnosis: Well-differentiated hepatocellular carcinoma in liver tissue.</p>
<p>TCGA-19-A60I GBM (Brain)</p> 	<p>Clinical report</p> <p>The slide from the brain shows glioblastoma, WHO grade IV, characterized by marked pleomorphism, numerous atypical mitoses, vascular proliferation, necrosis, and infiltration into the brain matter.</p>	<p>Ours (generated)</p> <p>glioblastoma, who grade iv, identified in brain tissue with dense cellularity, severe atypia, numerous mitoses, vascular proliferation, and necrosis. mgtm promoter is methylated and tp53 mutations are present.</p>	<p>PRISM (generated)</p> <p>Diagnosis: Metastatic high-grade sarcoma in examined tissue.</p>

Extended Data Figure 5: Examples of generated reports. TCGA examples of generated reports of TITAN and PRISM, with the corresponding clinical reports.



Extended Data Figure 6: Rare cancer retrieval with TITAN. (a-c) Examples of slide retrieval on Rare-Cancer. The number for each retrieved slide represents the cosine similarity between the query and the retrieved slide. The retrieved slides with high similarity are either of the same diagnostic label or from the same organ as the query slide.

Tissue site	Class	OncoTree code	#samples (train:val:test)
Adrenal gland	Adrenocortical carcinoma	ACC	493 (371:74:48)
Biliary tract	Cholangiocarcinoma	CHOL	90 (57:3:30)
Bladder	Bladder urothelial carcinoma	BLCA	943 (535:103:305)
Bowel	Colon adenocarcinoma	COAD	798 (623:120:55)
Brain	Glioblastoma multiforme	GBM	2283 (1223:342:718)
	Lower grade glioma	–	2098 (1113:278:707)
Breast	Invasive carcinoma	BRCA	2196 (1086:261:849)
Cervix	Squamous cell carcinoma and endocervical adenocarcinoma	CESC, ECAD	591 (340:78:173)
Esophagus	Esophageal carcinoma	–	294 (117:27:150)
Eye	Uveal melanoma	UM	147 (65:16:66)
Head and neck	Head and neck squamous cell carcinoma	HNSC	1159 (555:158:446)
Kidney	Renal clear cell carcinoma	CCRCC	798 (406:100:292)
	Papillary renal cell carcinoma	PRCC	443 (249:43:151)
	Chromophobe Renal Cell Carcinoma	CHRCC	180 (103:33:44)
Liver	Hepatocellular carcinoma	HCC	841 (565:128:148)
Lung	Lung adenocarcinoma	LUAD	1217 (676:144:397)
	Lung squamous cell carcinoma	LUSC	1213 (606:195:412)
Lymph	Diffuse large B-cell lymphoma	–	72 (37:12:23)
–	Mesothelioma	–	163 (75:31:57)
Ovary	Serous Cystadenocarcinoma	–	220 (95:24:101)
Pancreas	Pancreatic adenocarcinoma	PAAD	341 (190:42:109)
–	Pheochromocytoma and paraganglioma	PHC, PGNG	128 (81:17:30)
Prostate	Prostate adenocarcinoma	PRAD	815 (411:117:287)
Rectum	Rectum adenocarcinoma	READ	150 (93:29:28)
Soft tissue	Sarcoma	–	1270 (850:186:234)
Skin	Cutaneous Melanoma	SKCM	931 (313:94:524)
Stomach	Stomach adenocarcinoma	STAD	2306 (1482:335:489)
Testis	Testicular germ cell tumor	–	551 (375:79:97)
Thymus	Thymoma	THYM	328 (138:34:156)
Thyroid	Thyroid carcinoma	–	1063 (528:166:369)
Uterus	Uterine corpus endometrial carcinoma	UCEC	1172 (424:144:604)
	Uterine carcinosarcoma	UCS	201 (71:21:109)

Extended Data Table 1: Overview of the dataset TCGA-UniformTumor-8K with 32 classes grouped by tissue site and sorted by largest class. This dataset is exclusively curated for ROIs (8,192×8,192 pixels) subtyping task. Not every class has a one-to-one mapping to a tissue site or a single OncoTree code.

Tissue site	OncoTree code	#samples (train:val:test)
Adrenal gland	ACC	227 (158:50:19)
	PHC	163 (118:27:18)
Bladder	BLCA	457 (376:49:32)
Bowel	COAD	375 (278:49:48)
	READ	156 (113:25:18)
	MACR	63 (48:11:4)
Brain	GBM	858 (798:50:10)
	OAST	217 (165:45:7)
	ODG	203 (152:46:5)
	AASTR	164 (113:39:12)
	AOAST	155 (112:41:2)
Breast	ILC	211 (154:36:21)
	IDC	838 (743:49:46)
Cervix	CESC	229 (134:40:55)
Eye	UM	79 (34:10:35)
Head and neck	HNSC	472 (338:47:87)
Kidney	CCRCC	519 (455:49:15)
	PRCC	297 (197:46:54)
	CHRCC	109 (65:19:25)
Liver	HCC	362 (303:46:13)
Lung	LUAD	531 (305:48:178)
	LUSC	512 (364:48:100)
Melanoma	MEL	393 (347:31:15)
Ovary	HGSOC	107 (75:23:9)
Pancreas	PAAD	194 (135:32:27)
Pleura	PLEMESO	62 (43:17:2)
Prostate	PRAD	449 (382:44:23)
Skin	SKCM	75 (36:29:10)
Soft tissue	MFH	165 (119:40:6)
	LMS	155 (108:37:10)
	MFS	141 (91:41:9)
	DDL5	87 (56:29:2)
	THYM	34 (18:6:10)
Stomach	STAD	193 (114:46:33)
	ESCC	92 (60:12:20)
	TSTAD	87 (45:14:28)
	DSTAD	74 (57:14:3)
	ESCA	66 (44:10:12)
Testis	NSGCT	128 (78:43:7)
	SEM	91 (55:32:4)
Thymus	THYM	146 (68:35:43)
Thyroid	THPA	408 (253:45:110)
	THFO	109 (41:11:57)
Uterus	UEC	422 (308:44:70)
	USC	120 (64:40:16)
	UCS	87 (42:34:11)

Extended Data Table 2: Overview of the dataset TCGA-OT with 46 OncoTree codes grouped by tissue site and sorted by largest classes. This dataset is used for slide-level evaluations. Some cancer types can occur at multiple tissue sites and are listed in the tissue sites with the most samples, e.g., Leiomyosarcoma (LMS) contains samples from the uterus, stomach, bone, ovary, and head and neck. Additionally, all melanomas (MEL) are listed as a separate site with samples from skin (241), lymph (110), soft tissue (21), bowel (8), spleen (2), adrenal gland, brain, head and neck, thorax, and vulva (each 1). For every OncoTree code, we list the total number of samples and the number of samples contained in train, val, and test folds.

Tissue Site	Oncotree code	#samples
Adrenal Gland	ACC	221
	PHC	155
	MNET	1
Biliary tract	IHCH	30
	PHCH	4
	EHCH	2
	CHOL	1
Bladder	BLCA	424
Bone	LMS	3
	DLBCLNOS	1
Bowel	COAD	347
	READ	149
	MACR	53
	DLBCLNOS	5
	COADREAD	3
	DDLS	3
Brain	GBM	521
	OAST	203
	ODG	184
	AASTR	153
	AOAST	148
	ASTR	96
Breast	IDC	788
	ILC	203
	MDLC	27
	IMMC	16
	MBC	12
	BRCA	6
	PD	3
	MPT	2
	SPC	2
	BRCNOS	1
	BCC	1
	ACBC	1
DLBCLNOS	1	
Cervix	CESC	219
	ECAD	21
	CEMU	13
	CEAS	4
	CEEN	3
Eye	UM	64
	PHC	1
Head and neck	HNSC	464
	PGNG	2
	LMS	1
	DLBCLNOS	1
Heart	PGNG	1
Kidney	CCRCC	504
	PRCC	284
	CHRCC	108
	SYNS	5
	MFH	2
Liver	HCC	326
	HCCIHCH	7
	FLC	3
Lung	LUAD	500
	LUSC	476
	PLEMESO	1
Lymph	DLBCLNOS	25
	PGNG	1

Extended Data Table 3: Overview of the dataset TCGA-Slide-Reports with 89 OncoTree codes grouped by 30 tissue sites.

Nervous system	DLBCLNOS	2
	MPNST	1
Ovary	HGSOC	63
	LMS	1
Pancreas	PAAD	182
	PANET	8
	UCP	1
Peritoneum	DDL	1
	DLBCLNOS	1
Pleura	PLEMESO	56
	PLBMESO	18
	PLMESO	3
	PLSMESO	1
Peripheral nervous system	MPNST	5
	PGNG	1
Prostate	PRAD	406
Skin	SKCM	73
	MEL	29
	DESM	3
	SKLMM	1
	ACRM	1
Soft tissue	MFH	150
	MFS	136
	LMS	109
	DDL	80
	THYM	27
	PGNG	23
	SYNS	15
	MPNST	15
	DES	2
	DLBCLNOS	2
	PHC	1
SARCNO	1	
Stomach	STAD	162
	TSTAD	70
	ESCA	66
	ESCC	52
	DSTAD	50
	MSTAD	19
	SSRCC	13
	PSTAD	7
	LMS	5
	MBN	1
	DDL	1
	DLBCLNOS	1
Testis	SEM	80
	NSGCT	28
	EMBCA	11
	DLBCLNOS	1
	DDL	1
Thorax, not otherwise specified	PGNG	2
Thymus	THYM	140
	THPA	385
	THFO	107
	DLBCLNOS	2
	THPD	1
Uterus	UEC	419
	USC	120
	UCS	83
	LMS	34
	UCEC	23

Extended Data Table 4: Overview of the dataset TCGA-Slide-Reports. Continued.

Tissue site	OncoTree code	# samples
Adrenal gland	ACC	248
	PAAD	1
Biliary tract	IHCH	32
	GBC	27
	CHOL	5
Bladder	UTUC	28
	LUCA	1
Bone	THFO	2
	CHOL	2
	PAAD	1
	ESCC	1
Bowel	ANSC	41
	GBC	4
	USC	1
	EOV	1
	PANET	1
	ACC	1
	LNET	1
Brain	ODG	256
	AASTR	214
	ASTR	167
	PAST	33
	AODG	30
	LUNE	4
	PAAD	2
	LNET	1
	ACC	1
	WT	1
Breast	MBC	13
Cervix	ECAD	24
	CEMU	17
Head and neck	PGNG	2
	THAP	1
Heart	PGNG	1
Kidney	CHRCC	128
	WT	49
	UTUC	26
	SYNS	5

Extended Data Table 5: Overview of the dataset Rare-Cancers with 43 OncoTree codes grouped by 28 tissue sites.

Liver	PANET	23
	CHOL	14
	GBC	11
	PAAD	10
	ANSC	6
	LUNE	3
	ESCC	3
	LUCA	2
	CCOV	2
	THME	2
	ACC	1
	UTUC	1
	LNET	1
CHRCC	1	
Lung	LUCA	55
	LUNE	11
	LNET	6
	ACC	4
	ESCC	4
	USC	2
	CHOL	2
	PAAD	2
	ANSC	2
	WT	1
	THAP	1
	UTUC	1
Lymph	THME	9
	LUNE	6
	USC	5
	ANSC	4
	CHOL	3
	GBC	2
	PANET	2
	THFO	2
	CCOV	2
	EOV	2
	PGNG	1
	UTUC	1
	LNET	1
	THAP	1
	LUCA	1
PAAD	1	
CHRCC	1	
Nervous system, not otherwise specified	PGNG	1
	MPNST	1
Ovary	EOV	51
	CCOV	51
	USC	1
	ACC	1
	PANET	1
Pancreas	PAAD	224
	PANET	42
	ESCC	1

Extended Data Table 6: Overview of the dataset Rare-Cancers. Continued.

Peritoneum	USC	3
	GBC	3
	CHOL	3
	EOV	2
	CCOV	2
	PANET	1
	PAAD	1
Pleura	PLBMESO	19
	LUNE	1
	THAP	1
Peripheral neural sytem	MPNST	5
	PGNG	1
Skin	THME	1
	ANSC	1
	GBC	1
Soft tissue	PGNG	23
	MPNST	21
	SYNS	15
	GBC	6
	CHOL	5
	EOV	4
	USC	4
	WT	2
	PAAD	2
	LUNE	1
	PANET	1
	LNET	1
	THME	1
	CHRCC	1
Stomach	ESCC	125
	PAAD	1
Testis	EMBCA	35
Thorax, not otherwise specified	PGNG	2
Thyroid	THFO	131
	THAP	23
	THME	21
Uterus	USC	152
	CCOV	1
Vulva	ANSC	1
N/A	GBC	5
	PAAD	3
	USC	3
	ANSC	2
	CCOV	1
	ESCC	1
	PAST	1
	THAP	1
	UTUC	1

Extended Data Table 7: Overview of the dataset Rare-Cancers. Continued.

Tissue Site	Oncotree code	#samples
Adrenal gland	ACC	227
Biliary tract	IHCH	32
Brain	ODG	203
	AASTR	164
	ASTR	104
Breast	MBC	13
Cervix	ECAD	24
	CEMU	17
Head and neck	PGNG	2
Heart	PGNG	1
Kidney	CHRCC	109
	SYNS	5
Lymph	PGNG	1
Nervous system, not otherwise specified	PGNG	1
	MPNST	1
Pancreas	PAAD	194
	PANET	8
Pleura	PLBMESO	19
Peripheral neural system	MPNST	5
	PGNG	1
Soft tissue	PGNG	23
	MPNST	21
	SYNS	15
Stomach	ESCC	92
Testis	EMBCA	35
Thorax, not otherwise specified	PGNG	2
Thyroid	THFO	109
Uterus	USC	120

Extended Data Table 8: Overview of the dataset Rare-Cancers-Public with 29 OncoTree codes grouped by 19 tissue sites.

Dataset	Tissue site	Target	Label-level	# Classes	# Patients	# WSIs	Metric	External test
OT108	Pan-cancer	OncoTree	WSI	108	5,510	5,564	Bal. acc.	
TCGA-OT	Pan-cancer	OncoTree	WSI	46	9,149	11,186	Bal. acc.	
TCGA-UT-8K	Pan-cancer	Subtype	WSI	32	7,784	24,392	Bal. acc.	✓
EBRAINS	Brain	Diagnosis	WSI	30	2,147	2,319	Bal. acc.	
BRACS	Breast	Subtype	WSI	7	189	547	Bal. acc.	
BRACS	Breast	Subtype (coarse)	WSI	3	189	547	Bal. acc.	
TCGA-BRCA	Breast	Subtype	WSI	2	984	1,049	AUROC	
CRANE	Heart	Cellular	Patient	2	1,688	5,021	AUROC	
TCGA-RCC	Kidney	OncoTree	WSI	3	895	922	Bal. acc.	
CPTAC-DHMC	Kidney	OncoTree	WSI	3	673	872	Bal. acc.	✓
Renal allograft rejection	Kidney	AMR	Patient	2	1,118	4,847	AUROC	
Renal allograft rejection	Kidney	Cellular rejection	Patient	2	1,118	4,847	AUROC	
TCGA-NSCLC	Lung	Subtype	WSI	2	946	1,043	AUROC	
CPTAC-NSCLC	Lung	Subtype	WSI	2	422	1,091	AUROC	✓

Extended Data Table 9: Overview of the morphological tasks sorted by organ.

Dataset	Tissue site	Target	Label-level	# Classes	# Patients	# WSIs	Metric
IMP	Colorectal	Dysplasia grading	WSI	3	5,333	5,333	Cohen's κ
Renal allograft rejection	Kidney	IFTA status	Patient	3	1,118	4,847	Cohen's κ
PANDA	Prostate	Gleason grading	WSI	6	9,555	9,555	Cohen's κ

Extended Data Table 10: Overview of the grading tasks sorted by organ. The metric Cohen's κ is quadratic weighted.

Dataset	Tissue Site	Target	Label-level	# Patients	# WSIs
TCGA-BRCA	Breast	Survival	Patient	1,022	1,092
TCGA-BLCA	Bladder	Survival	Patient	360	424
TCGA-CRC	Colorectal	Survival	Patient	545	554
TCGA-KIRC	Kidney	Survival	Patient	502	508
TCGA-UCEC	Uterus	Survival	Patient	504	564
TCGA-NSCLC	Lung	Survival	Patient	844	920

Extended Data Table 11: Overview of the survival prediction tasks for measuring disease-specific survival (DSS) with c-index sorted by organ. All tasks are evaluated in five-fold site-preserving cross-validation splits.

Dataset	Tissue Site	Target	Label-level	# Classes	# Patients	# WSIs	Metric	External test
TCGA-GBM	Brain	<i>EGFR</i>	Patient	2	228	560	AUROC	
CPTAC-GBM	Brain	<i>EGFR</i>	Patient	2	99	243	AUROC	✓
TCGA-GBM	Brain	<i>TP53</i>	Patient	2	228	560	AUROC	
CPTAC-GBM	Brain	<i>TP53</i>	Patient	2	99	243	AUROC	
TCGA-GBMLGG	Brain	IDH	WSI	2	558	1,123	AUROC	
EBRAINS	Brain	IDH	WSI	2	795	873	AUROC	✓
BCNB	Breast	ER	Patient	2	1,058	1,058	AUROC	
BCNB	Breast	PR	Patient	2	1,058	1,058	AUROC	
BCNB	Breast	HER2	Patient	2	1,058	1,058	AUROC	
MGB-BRCA	Breast	ER	Patient	2	874	874	AUROC	
MGB-BRCA	Breast	PR	Patient	2	874	874	AUROC	
MGB-BRCA	Breast	HER2	Patient	2	816	816	AUROC	
MGH-BRCA (IHC)	Breast	ER level	Patient	6	962	962	Cohen's κ	
MGH-BRCA (IHC)	Breast	PR level	Patient	6	1,070	1,070	Cohen's κ	
TCGA-BRCA	Breast	ER	Patient	2	937	996	AUROC	
CPTAC-BRCA	Breast	ER	Patient	2	102	111	AUROC	✓
TCGA-BRCA	Breast	PR	Patient	2	934	993	AUROC	
CPTAC-BRCA	Breast	PR	Patient	2	97	106	AUROC	✓
TCGA-BRCA	Breast	HER2	Patient	2	647	693	AUROC	
CPTAC-BRCA	Breast	HER2	Patient	2	103	112	AUROC	✓
TCGA-BRCA	Breast	<i>PIK3CA</i>	Patient	2	970	1,034	AUROC	
CPTAC-BRCA	Breast	<i>PIK3CA</i>	Patient	2	103	112	AUROC	✓
TCGA-CRC	Colorectal	MSI	Patient	2	414	419	AUROC	
CPTAC-COAD	Colorectal	MSI	Patient	2	103	107	AUROC	✓
TCGA-CRC	Colorectal	<i>BRAF</i>	Patient	2	487	492	AUROC	
CPTAC-COAD	Colorectal	<i>BRAF</i>	Patient	2	103	107	AUROC	✓
MUT-HET-RCC	Kidney	<i>BAP1</i>	Patient	2	1,292	1,292	AUROC	
MUT-HET-RCC	Kidney	<i>PBRM1</i>	Patient	2	1,292	1,292	AUROC	
MUT-HET-RCC	Kidney	<i>SETD2</i>	Patient	2	1,292	1,292	AUROC	
MGB-LUAD	Lung	CDX-2	Patient	2	79	79	AUROC	
MGB-LUAD	Lung	CK-5&6	Patient	2	58	58	AUROC	
MGB-LUAD	Lung	Napsin A	Patient	2	126	126	AUROC	
MGB-LUAD	Lung	P40	Patient	2	185	185	AUROC	
MGB-LUAD	Lung	P63	Patient	2	153	153	AUROC	
TCGA-NSCLC	Lung	<i>EGFR</i>	Patient	2	462	524	AUROC	
CPTAC-LUAD	Lung	<i>EGFR</i>	Patient	2	108	324	AUROC	✓
TCGA-NSCLC	Lung	<i>TP53</i>	Patient	2	462	524	AUROC	
CPTAC-LUAD	Lung	<i>TP53</i>	Patient	2	108	324	AUROC	✓

Extended Data Table 12: Overview of the molecular tasks sorted by organ. The metric Cohen's κ is quadratic weighted.

Dataset	Link
TCGA	https://portal.gdc.cancer.gov/
CPTAC	https://proteomic.datacommons.cancer.gov/pdc
EBRAINS	https://doi.org/10.25493/WQ48-ZGX
DHMC-RCC	https://bmirids.github.io/KidneyCancer
DHMC-LUAD	https://bmirids.github.io/LungCancer/
BRACS	https://www.bracs.icar.cnr.it/
PANDA	https://panda.grand-challenge.org
IMP	https://rdm.inesctec.pt/dataset/nis-2023-008
BCNB	https://bupt-ai-cz.github.io/BCNB/
MUT-HET-RCC	https://aacrjournals.org/cancerres/article/82/15/2792/707325/Intratumoral-Resolution-of-Driver-Gene-Mutation

Extended Data Table 13: Summary of publicly available datasets.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.711±0.0069	0.791±0.0046	0.606±0.0073	0.674±0.0050	0.563±0.0071	0.639±0.0052
GigaPath	0.700±0.0069	0.782±0.0046	0.591±0.0074	0.668±0.0051	0.564±0.0070	0.644±0.0051
Mean pool (Virchow)	0.691±0.0072	0.771±0.0047	0.524±0.0079	0.579±0.0055	0.523±0.0070	0.615±0.0054
PRISM	0.774±0.0062	0.824±0.0042	0.708±0.0063	0.737±0.0049	0.702±0.0068	0.777±0.0048
Mean pool (CHIEF)	0.627±0.0077	0.729±0.0050	0.477±0.0076	0.560±0.0054	0.495±0.0072	0.581±0.0056
CHIEF	0.625±0.0079	0.723±0.0051	0.503±0.0072	0.567±0.0053	0.506±0.0071	0.594±0.0054
Mean pool (CONCH)	0.779±0.0064	0.833±0.0043	0.757±0.0062	0.785±0.0045	0.743±0.0067	0.796±0.0045
TITAN _v	<u>0.820±0.0055</u>	<u>0.870±0.0037</u>	<u>0.812±0.0050</u>	<u>0.842±0.0038</u>	<u>0.818±0.0053</u>	<u>0.864±0.0039</u>
TITAN	0.832±0.0056	0.881±0.0036	0.843±0.0056	0.865±0.0037	0.828±0.0052	0.875±0.0037

Extended Data Table 14: Linear probing results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.568±0.0167	0.662±0.0132	0.479±0.0146	0.562±0.0137	0.419±0.0138	0.521±0.0142
GigaPath	0.543±0.0178	0.659±0.0134	0.477±0.0149	0.538±0.0138	0.402±0.0132	0.516±0.0145
Mean pool (Virchow)	0.564±0.0165	0.656±0.0135	0.379±0.0141	0.454±0.0150	0.371±0.0109	0.479±0.0148
PRISM	0.643±0.0181	0.732±0.0129	0.617±0.0201	0.703±0.0130	0.606±0.0166	0.710±0.0133
Mean pool (CHIEF)	0.526±0.0188	0.643±0.0141	0.416±0.0195	0.493±0.0136	0.381±0.0143	0.496±0.0146
CHIEF	0.528±0.0205	0.640±0.0142	0.442±0.0207	0.478±0.0138	0.403±0.0146	0.504±0.0145
Mean pool (CONCH)	0.624±0.0169	0.700±0.0131	0.615±0.0171	0.680±0.0130	0.573±0.0153	0.666±0.0132
TITAN _v	<u>0.690±0.0196</u>	<u>0.758±0.0123</u>	<u>0.706±0.0183</u>	<u>0.775±0.0116</u>	<u>0.639±0.0149</u>	<u>0.744±0.0124</u>
TITAN	0.704±0.0192	0.764±0.0116	0.744±0.0152	0.802±0.0108	0.685±0.0183	0.774±0.0119

Extended Data Table 15: Linear probing results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.450±0.0111	0.436±0.0127	0.365±0.0107	0.353±0.0123	0.318±0.0098	0.293±0.0114
GigaPath	0.437±0.0110	0.423±0.0124	0.354±0.0107	0.344±0.0125	0.308±0.0097	0.283±0.0116
Mean pool (Virchow)	0.454±0.0111	0.439±0.0127	0.315±0.0102	0.304±0.0121	0.283±0.0096	0.255±0.0117
PRISM	0.508±0.0110	0.481±0.0133	0.483±0.0119	0.469±0.0131	0.434±0.0102	0.384±0.0129
Mean pool (CHIEF)	0.356±0.0112	0.340±0.0125	0.226±0.0097	0.206±0.0114	0.210±0.0095	0.185±0.0103
CHIEF	0.413±0.0109	0.394±0.0130	0.295±0.0108	0.282±0.0120	0.262±0.0093	0.231±0.0108
Mean pool (CONCH)	0.476±0.0112	0.459±0.0129	0.411±0.0107	0.404±0.0127	0.368±0.0103	0.341±0.0122
TITAN _v	<u>0.558±0.0104</u>	<u>0.536±0.0127</u>	<u>0.524±0.0106</u>	<u>0.517±0.0129</u>	<u>0.464±0.0105</u>	<u>0.430±0.0128</u>
TITAN	0.587±0.0103	0.563±0.0130	0.580±0.0103	0.567±0.0129	0.527±0.0101	0.489±0.0132

Extended Data Table 16: Linear probing results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.712±0.0215	0.766±0.0193	0.626±0.0238	0.627±0.0197	0.543±0.0221	0.633±0.0221
GigaPath	0.680±0.0217	0.746±0.0197	0.626±0.0235	0.610±0.0201	0.530±0.0213	0.621±0.0219
Mean pool (Virchow)	0.665±0.0220	0.735±0.0189	0.559±0.0236	0.553±0.0207	0.481±0.0200	0.573±0.0220
PRISM	0.674±0.0200	0.732±0.0191	0.625±0.0223	0.583±0.0204	0.609±0.0198	0.660±0.0201
Mean pool (CHIEF)	0.594±0.0240	0.663±0.0205	0.487±0.0242	0.463±0.0214	0.368±0.0193	0.443±0.0217
CHIEF	0.598±0.0237	0.670±0.0206	0.478±0.0239	0.453±0.0213	0.415±0.0198	0.498±0.0219
Mean pool (CONCH)	0.644±0.0215	0.715±0.0200	0.550±0.0221	0.480±0.0217	0.551±0.0221	0.631±0.0211
TITAN _v	<u>0.732±0.0208</u>	<u>0.785±0.0175</u>	<u>0.742±0.0192</u>	<u>0.727±0.0182</u>	<u>0.656±0.0187</u>	<u>0.709±0.0199</u>
TITAN	0.735±0.0204	0.786±0.0182	0.754±0.0203	0.748±0.0182	0.695±0.0196	0.746±0.0200

Extended Data Table 17: Linear probing results for tumor type ($C = 30$) prediction on EBRAINS. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.938±0.0204	0.780±0.0358	0.792±0.0354	0.799±0.0428	0.639±0.0460	0.812±0.0346
GigaPath	0.935±0.0243	0.780±0.0463	0.776±0.0542	0.783±0.0493	0.607±0.0247	0.794±0.0208
Mean pool (Virchow)	0.904±0.0178	0.749±0.0431	0.718±0.0583	0.725±0.0804	0.540±0.0200	0.744±0.0255
PRISM	0.944±0.0241	0.850±0.0398	0.877±0.0412	0.906±0.0197	<u>0.879±0.0400</u>	0.919±0.0195
Mean pool (CHIEF)	0.930±0.0137	0.786±0.0593	0.733±0.0702	0.741±0.0863	0.621±0.0375	0.801±0.0297
CHIEF	0.935±0.0219	0.843±0.0534	0.822±0.0261	0.836±0.0269	0.775±0.0708	0.879±0.0378
Mean pool (CONCH)	0.927±0.0229	0.776±0.0626	0.803±0.0443	0.810±0.0723	0.767±0.0600	0.878±0.0334
TITAN _v	0.948±0.0146	<u>0.848±0.0376</u>	<u>0.885±0.0290</u>	<u>0.905±0.0265</u>	0.875±0.0428	<u>0.921±0.0180</u>
TITAN	<u>0.945±0.0174</u>	<u>0.848±0.0293</u>	0.892±0.0271	0.903±0.0192	0.902±0.0268	0.926±0.0162

Extended Data Table 18: Linear probing results for tumor subtype ($C = 2$) prediction on TCGA-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.948±0.0129	0.883±0.0278	0.795±0.0636	0.795±0.0656	0.783±0.0109	0.771±0.0125
GigaPath	TCGA	0.934±0.0142	0.869±0.0251	0.759±0.0605	0.758±0.0608	0.744±0.0075	0.727±0.0101
Mean pool (Virchow)	TCGA	0.944±0.0153	0.866±0.0339	0.684±0.0620	0.682±0.0643	0.668±0.0593	0.646±0.0831
PRISM	TCGA	0.977±0.0058	<u>0.934±0.0195</u>	0.914±0.0189	0.915±0.0185	<u>0.927±0.0018</u>	0.928±0.0019
Mean pool (CHIEF)	TCGA	0.947±0.0121	0.881±0.0245	0.765±0.0582	0.764±0.0595	0.803±0.0070	0.793±0.0099
CHIEF	TCGA	0.963±0.0087	0.897±0.0197	0.846±0.0276	0.845±0.0289	0.862±0.0081	0.859±0.0091
Mean pool (CONCH)	TCGA	0.965±0.0147	0.904±0.0207	0.831±0.0370	0.831±0.0381	0.904±0.0035	0.903±0.0034
TITAN _v	TCGA	0.979±0.0059	0.932±0.0190	<u>0.929±0.0147</u>	<u>0.930±0.0148</u>	0.933±0.0032	<u>0.934±0.0030</u>
TITAN	TCGA	0.982±0.0058	0.941±0.0230	0.944±0.0198	0.944±0.0193	0.933±0.0031	0.935±0.0029
Mean pool (GigaPath)	CPTAC	0.954±0.0105	0.889±0.0166	0.755±0.0367	0.732±0.0432	0.844±0.0086	0.843±0.0087
GigaPath	CPTAC	0.950±0.0105	0.882±0.0165	0.722±0.0310	0.692±0.0396	0.581±0.0097	0.470±0.0188
Mean pool (Virchow)	CPTAC	0.953±0.0151	0.864±0.0208	0.814±0.0300	0.803±0.0364	0.672±0.0118	0.653±0.0159
PRISM	CPTAC	0.973±0.0036	0.932±0.0124	<u>0.941±0.0017</u>	<u>0.941±0.0019</u>	0.922±0.0039	0.921±0.0040
Mean pool (CHIEF)	CPTAC	0.945±0.0130	0.867±0.0186	0.801±0.0226	0.787±0.0266	0.792±0.0149	0.787±0.0159
CHIEF	CPTAC	0.945±0.0116	0.867±0.0175	0.871±0.0115	0.867±0.0125	0.844±0.0092	0.840±0.0107
Mean pool (CONCH)	CPTAC	0.978±0.0032	0.934±0.0069	0.876±0.0088	0.872±0.0100	0.847±0.0163	0.851±0.0163
TITAN _v	CPTAC	0.986±0.0008	<u>0.939±0.0046</u>	0.940±0.0019	0.940±0.0022	0.939±0.0025	0.939±0.0024
TITAN	CPTAC	<u>0.985±0.0010</u>	0.950±0.0031	0.953±0.0025	0.954±0.0022	<u>0.929±0.0031</u>	<u>0.932±0.0028</u>

Extended Data Table 19: Linear probing results for tumor subtype ($C = 2$) prediction on TCGA-NSCLC and external dataset CPTAC-NSCLC. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.841±0.0470	0.896±0.0334	0.856±0.0561	0.843±0.0425	0.776±0.0262	0.899±0.0304
GigaPath	TCGA	0.823±0.0343	0.892±0.0290	0.859±0.0520	0.863±0.0366	0.745±0.0473	0.905±0.0137
Mean pool (Virchow)	TCGA	0.819±0.0879	0.886±0.0543	0.759±0.1060	0.739±0.1575	0.570±0.0407	0.841±0.0898
PRISM	TCGA	0.866±0.0360	0.915±0.0272	0.859±0.0509	0.855±0.0415	0.713±0.0168	0.917±0.0032
Mean pool (CHIEF)	TCGA	0.900±0.0473	0.921±0.0366	0.897±0.0279	0.881±0.0308	0.832±0.0071	0.910±0.0163
CHIEF	TCGA	0.913±0.0268	0.928±0.0280	0.910±0.0346	0.891±0.0298	0.841±0.0123	0.890±0.0140
Mean pool (CONCH)	TCGA	0.871±0.0498	0.910±0.0437	0.900±0.0306	<u>0.911</u> ±0.0350	0.839±0.0108	0.945±0.0017
TITAN _v	TCGA	<u>0.920</u> ±0.0331	<u>0.931</u> ±0.0354	<u>0.954</u> ±0.0148	0.939 ±0.0239	<u>0.945</u> ±0.0131	0.966 ±0.0028
TITAN	TCGA	0.942 ±0.0222	0.941 ±0.0333	0.955 ±0.0202	0.939 ±0.0341	0.971 ±0.0025	<u>0.965</u> ±0.0009
Mean pool (GigaPath)	CPTAC-DHMC	0.854±0.0327	0.919±0.0288	0.840±0.0137	0.835±0.0205	0.841±0.0228	0.950±0.0047
GigaPath	CPTAC-DHMC	0.833±0.0521	0.933±0.0187	0.850±0.0107	0.882±0.0175	0.615±0.0368	0.911±0.0079
Mean pool (Virchow)	CPTAC-DHMC	0.701±0.0716	0.799±0.0827	0.488±0.0135	0.584±0.0332	0.666±0.0884	0.694±0.0091
PRISM	CPTAC-DHMC	0.753±0.0508	0.934±0.0054	0.724±0.0071	0.902±0.0024	0.547±0.0381	0.897±0.0072
Mean pool (CHIEF)	CPTAC-DHMC	0.840±0.0289	0.938±0.0108	0.868±0.0069	0.883±0.0167	0.655±0.0173	0.908±0.0063
CHIEF	CPTAC-DHMC	0.822±0.0143	0.925±0.0115	<u>0.869</u> ±0.0059	0.842±0.0123	0.727±0.0091	0.926±0.0047
Mean pool (CONCH)	CPTAC-DHMC	0.850±0.0238	0.944±0.0089	0.836±0.0033	0.936±0.0025	0.668±0.0079	0.897±0.0026
TITAN _v	CPTAC-DHMC	<u>0.910</u> ±0.0169	<u>0.959</u> ±0.0063	0.971 ±0.0011	<u>0.959</u> ±0.0029	0.922 ±0.0067	0.971 ±0.0021
TITAN	CPTAC-DHMC	0.918 ±0.0073	0.967 ±0.0040	0.971 ±0.0026	0.962 ±0.0012	<u>0.890</u> ±0.0089	<u>0.959</u> ±0.0013

Extended Data Table 20: Linear probing results for OncoTree code ($C = 3$) prediction on TCGA-RCC and external dataset CPTAC-CCRCC and DHMC-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.321±0.0482	0.433±0.0998	0.308±0.0424	0.357±0.0863	0.239±0.0416	0.352±0.0967
GigaPath	0.305±0.0465	0.429±0.1034	0.287±0.0355	0.332±0.0837	0.282±0.0813	0.367±0.1239
Mean pool (Virchow)	0.307±0.0345	0.415±0.1096	0.275±0.0430	0.342±0.0766	0.262±0.0613	0.375±0.0937
PRISM	0.411 ±0.0569	0.521 ±0.0584	<u>0.474</u> ±0.0594	0.535 ±0.0760	0.408 ±0.0524	0.522 ±0.0662
Mean pool (CHIEF)	0.322±0.0617	0.443±0.1139	0.281±0.0424	0.309±0.0906	0.282±0.0426	0.352±0.0864
CHIEF	<u>0.405</u> ±0.0553	<u>0.513</u> ±0.0837	0.418±0.0529	0.430±0.1012	<u>0.368</u> ±0.0257	0.449±0.0785
Mean pool (CONCH)	0.332±0.0509	0.450±0.0869	0.336±0.0432	0.396±0.0696	0.313±0.0619	0.412±0.0980
TITAN _v	0.404±0.0409	0.521 ±0.0836	0.424±0.0554	0.478±0.1046	0.349±0.0375	0.466±0.0884
TITAN	0.400±0.0259	0.511±0.0686	0.483 ±0.0819	<u>0.491</u> ±0.1154	0.367±0.0387	<u>0.479</u> ±0.0797

Extended Data Table 21: Linear probing results for subtype ($C = 7$) prediction on BRACS. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.560±0.0906	0.662±0.1228	0.565±0.0755	0.557±0.0847	0.509±0.0904	0.589±0.1080
GigaPath	0.568±0.1028	0.661±0.1221	0.562±0.0620	0.548±0.0766	0.515±0.0774	0.599±0.1054
Mean pool (Virchow)	0.571±0.0792	0.674±0.1144	0.581±0.0863	0.580±0.0946	0.514±0.0633	0.618±0.1045
PRISM	0.658±0.0393	0.759±0.0771	<u>0.726</u> ±0.0456	0.761 ±0.0517	<u>0.672</u> ±0.0319	<u>0.766</u> ±0.0608
Mean pool (CHIEF)	0.563±0.0811	0.666±0.0884	0.567±0.0566	0.545±0.0834	0.492±0.0530	0.586±0.0793
CHIEF	0.638±0.0547	0.739±0.0695	0.675±0.0985	0.696±0.0844	0.625±0.0738	0.721±0.0897
Mean pool (CONCH)	0.606±0.0958	0.702±0.1145	0.645±0.0455	0.643±0.0592	0.516±0.0471	0.624±0.0799
TITAN _v	<u>0.679</u> ±0.0712	<u>0.780</u> ±0.0664	0.739 ±0.0964	0.738±0.1010	0.630±0.0650	0.737±0.0869
TITAN	0.702 ±0.0395	0.803 ±0.0410	0.739 ±0.0543	<u>0.755</u> ±0.0713	0.682 ±0.0120	0.790 ±0.0545

Extended Data Table 22: Linear probing results for coarse subtype ($C = 3$) prediction on BRACS. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.882±0.0187	0.811±0.0216	0.692±0.0242	0.690±0.0260	0.513±0.0194	0.354±0.0519
GigaPath	0.865±0.0204	0.801±0.0230	<u>0.721±0.0229</u>	<u>0.717±0.0264</u>	0.528±0.0150	0.489±0.0240
Mean pool (Virchow)	0.872±0.0189	0.805±0.0214	0.706±0.0245	0.706±0.0254	0.499±0.0119	0.308±0.0510
PRISM	0.877±0.0182	0.801±0.0213	0.692±0.0256	0.692±0.0256	<u>0.657±0.0190</u>	0.660±0.0198
Mean pool (CHIEF)	0.887±0.0186	0.817±0.0205	0.654±0.0258	0.654±0.0261	0.562±0.0460	0.565±0.0480
CHIEF	<u>0.893±0.0178</u>	0.822±0.0208	0.667±0.0253	0.667±0.0260	0.504±0.0334	0.465±0.0533
Mean pool (CONCH)	0.880±0.0196	0.826±0.0212	0.716±0.0235	0.714±0.0253	0.658±0.0419	<u>0.658±0.0413</u>
TITAN _V	0.906±0.0166	0.836±0.0204	0.713±0.0229	0.710±0.0251	0.554±0.0190	0.410±0.0544
TITAN	0.906±0.0162	<u>0.830±0.0207</u>	0.766±0.0222	0.765±0.0237	0.614±0.0343	0.557±0.0500

Extended Data Table 23: Linear probing results for cellular rejection ($C = 2$) prediction on CRANE. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.851±0.0318	<u>0.764±0.0354</u>	0.591±0.0389	0.598±0.0340	0.546±0.0368	0.636±0.0379
GigaPath	0.836±0.0354	0.758±0.0360	0.566±0.0403	0.571±0.0356	0.518±0.0348	0.619±0.0398
Mean pool (Virchow)	0.855±0.0302	0.751±0.0360	0.600±0.0384	0.585±0.0354	0.524±0.0340	0.622±0.0375
PRISM	0.820±0.0350	0.687±0.0351	0.633±0.0383	0.666±0.0327	0.569±0.0363	0.659±0.0358
Mean pool (CHIEF)	0.834±0.0322	0.727±0.0365	0.627±0.0389	0.638±0.0336	0.575±0.0381	0.661±0.0383
CHIEF	0.813±0.0336	0.649±0.0355	0.606±0.0390	0.645±0.0334	0.590±0.0360	0.678±0.0358
Mean pool (CONCH)	0.817±0.0343	0.594±0.0354	0.644±0.0389	0.668±0.0330	0.541±0.0351	0.642±0.0388
TITAN _V	<u>0.880±0.0293</u>	0.744±0.0356	<u>0.681±0.0370</u>	<u>0.688±0.0319</u>	<u>0.721±0.0371</u>	<u>0.770±0.0303</u>
TITAN	0.915±0.0238	0.807±0.0331	0.709±0.0357	0.704±0.0319	0.802±0.0335	0.844±0.0263

Extended Data Table 24: Linear probing results for antibody-mediated rejection ($C = 2$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.802±0.0331	0.667±0.0344	0.541±0.0384	0.538±0.0373	0.572±0.0361	0.615±0.0359
GigaPath	0.814±0.0331	0.715±0.0364	0.557±0.0389	0.559±0.0371	0.627±0.0356	0.667±0.0347
Mean pool (Virchow)	0.802±0.0323	0.669±0.0347	0.556±0.0386	0.548±0.0371	0.531±0.0376	0.568±0.0366
PRISM	0.840±0.0288	<u>0.736±0.0338</u>	0.692±0.0355	0.713±0.0323	0.670±0.0372	0.687±0.0344
Mean pool (CHIEF)	0.772±0.0360	0.669±0.0346	0.567±0.0387	0.571±0.0368	0.598±0.0371	0.634±0.0355
CHIEF	0.823±0.0307	0.681±0.0347	0.633±0.0384	0.656±0.0356	0.598±0.0374	0.634±0.0356
Mean pool (CONCH)	0.817±0.0310	0.686±0.0348	<u>0.721±0.0345</u>	<u>0.736±0.0316</u>	0.682±0.0358	<u>0.710±0.0334</u>
TITAN _V	<u>0.860±0.0270</u>	0.716±0.0348	0.719±0.0339	0.718±0.0322	<u>0.691±0.0370</u>	0.708±0.0341
TITAN	0.864±0.0264	0.743±0.0340	0.753±0.0335	0.753±0.0316	0.741±0.0346	0.764±0.0316

Extended Data Table 25: Linear probing results for cellular rejection ($C = 2$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.
Mean pool (GigaPath)	0.678±0.0476	0.630±0.0406	0.317±0.0604	0.477±0.0391	0.423±0.0736	0.516±0.0383
GigaPath	0.666±0.0497	0.632±0.0398	0.396±0.0584	0.518±0.0406	0.333±0.0725	0.449±0.0369
Mean pool (Virchow)	0.640±0.0542	0.597±0.0405	0.309±0.0596	0.485±0.0396	0.210±0.0785	0.416±0.0364
PRISM	0.603±0.0528	0.557±0.0396	0.521±0.0608	0.554±0.0413	0.549 ±0.0689	<u>0.583</u> ±0.0406
Mean pool (CHIEF)	0.673±0.0495	0.613±0.0395	0.419±0.0598	0.529±0.0366	0.383±0.0714	0.463±0.0366
CHIEF	0.665±0.0527	0.644±0.0385	0.569±0.0571	0.638 ±0.0355	0.506±0.0688	0.556±0.0377
Mean pool (CONCH)	<u>0.711</u> ±0.0482	0.712 ±0.0376	<u>0.600</u> ±0.0545	0.600±0.0406	0.411±0.0756	0.536±0.0414
TITAN _v	0.685±0.0519	0.681±0.0394	0.531±0.0641	0.613±0.0416	0.401±0.0702	0.532±0.0389
TITAN	0.723 ±0.0470	<u>0.688</u> ±0.0384	0.627 ±0.0506	<u>0.636</u> ±0.0392	<u>0.544</u> ±0.0645	0.590 ±0.0394

Extended Data Table 26: Linear probing results for Interstitial fibrosis and tubular atrophy (IFTA) level ($C = 3$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.
Mean pool (GigaPath)	0.911±0.0110	0.923±0.0085	0.805±0.0188	0.848±0.0123	0.841±0.0162	0.859±0.0129
GigaPath	0.909±0.0112	0.919±0.0092	0.790±0.0199	0.834±0.0129	0.852±0.0151	0.863±0.0125
Mean pool (Virchow)	0.910±0.0118	0.917±0.0100	0.772±0.0214	0.831±0.0133	0.838±0.0169	0.854±0.0136
PRISM	0.946 ±0.0076	<u>0.945</u> ±0.0079	0.912 ±0.0115	0.929 ±0.0081	0.944 ±0.0080	0.943 ±0.0079
Mean pool (CHIEF)	0.886±0.0128	<u>0.894</u> ±0.0111	0.642±0.0272	0.767±0.0150	0.780±0.0216	0.834±0.0140
CHIEF	0.909±0.0129	0.925±0.0094	0.856±0.0160	<u>0.886</u> ±0.0107	0.895±0.0139	<u>0.914</u> ±0.0098
Mean pool (CONCH)	0.899±0.0126	0.907±0.0104	0.716±0.0227	0.790±0.0141	0.831±0.0165	0.859±0.0122
TITAN _v	0.935±0.0097	0.939±0.0086	0.804±0.0160	0.820±0.0138	0.892±0.0120	0.898±0.0104
TITAN	<u>0.944</u> ±0.0091	0.946 ±0.0084	<u>0.857</u> ±0.0139	0.859±0.0124	<u>0.906</u> ±0.0111	0.908±0.0100

Extended Data Table 27: Linear probing results for dysplasia grading ($C = 3$) prediction on IMP-CRC biopsies. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.
Mean pool (GigaPath)	0.897±0.0112	0.727 ±0.0149	0.477±0.0289	0.373±0.0143	0.716±0.0221	0.526±0.0169
GigaPath	0.869±0.0136	0.676±0.0163	0.461±0.0287	0.354±0.0139	0.718±0.0228	0.532±0.0169
Mean pool (Virchow)	0.896±0.0102	<u>0.693</u> ±0.0160	0.468±0.0262	0.326±0.0131	0.694±0.0225	0.508±0.0168
PRISM	0.918 ±0.0078	<u>0.663</u> ±0.0157	<u>0.815</u> ±0.0171	0.573 ±0.0161	0.894 ±0.0101	0.652 ±0.0155
Mean pool (CHIEF)	0.843±0.0145	0.625±0.0165	<u>0.506</u> ±0.0269	0.370±0.0149	0.722±0.0222	0.504±0.0169
CHIEF	<u>0.911</u> ±0.0098	0.684±0.0158	0.712±0.0209	0.486±0.0155	0.864±0.0142	<u>0.639</u> ±0.0162
Mean pool (CONCH)	0.847±0.0130	0.602±0.0160	0.703±0.0192	0.442±0.0160	0.744±0.0209	<u>0.538</u> ±0.0164
TITAN _v	0.888±0.0102	0.648±0.0158	0.704±0.0199	0.472±0.0159	0.800±0.0172	0.571±0.0166
TITAN	0.899±0.0097	0.657±0.0162	0.831 ±0.0128	<u>0.533</u> ±0.0156	<u>0.867</u> ±0.0118	0.622±0.0158

Extended Data Table 28: Linear probing results for Gleason grading ($C = 6$) prediction on PANDA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.860 \pm 0.0161	0.643 \pm 0.0277	0.738 \pm 0.0392	0.745 \pm 0.0253	0.532 \pm 0.0060	0.833 \pm 0.0029
GigaPath	0.850 \pm 0.0261	0.564 \pm 0.0243	0.760 \pm 0.0454	0.769 \pm 0.0262	0.527 \pm 0.0146	0.831 \pm 0.0077
Mean pool (Virchow)	<u>0.856</u> \pm 0.0310	<u>0.617</u> \pm 0.0405	0.746 \pm 0.0494	0.746 \pm 0.0297	0.520 \pm 0.0193	0.827 \pm 0.0122
PRISM	0.860 \pm 0.0204	0.560 \pm 0.0428	<u>0.778</u> \pm 0.0615	0.779 \pm 0.0305	0.518 \pm 0.0065	0.826 \pm 0.0038
Mean pool (CHIEF)	0.840 \pm 0.0337	0.552 \pm 0.0295	0.754 \pm 0.0383	0.746 \pm 0.0195	0.515 \pm 0.0170	0.823 \pm 0.0103
CHIEF	0.842 \pm 0.0302	0.577 \pm 0.0484	0.764 \pm 0.0399	0.745 \pm 0.0212	0.522 \pm 0.0096	0.828 \pm 0.0065
Mean pool (CONCH)	0.842 \pm 0.0386	0.557 \pm 0.0227	0.775 \pm 0.0484	0.796 \pm 0.0230	0.519 \pm 0.0130	0.826 \pm 0.0084
TITAN _V	<u>0.856</u> \pm 0.0194	0.598 \pm 0.0225	0.783 \pm 0.0502	0.796 \pm 0.0214	0.560 \pm 0.0204	0.845 \pm 0.0097
TITAN	0.854 \pm 0.0108	0.580 \pm 0.0107	0.777 \pm 0.0351	<u>0.791</u> \pm 0.0249	<u>0.553</u> \pm 0.0176	<u>0.842</u> \pm 0.0086

Extended Data Table 29: Linear probing results for BAP1 mutation ($C = 2$) prediction on MUT-HET-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.800 \pm 0.0174	0.728 \pm 0.0199	0.613 \pm 0.0123	0.610 \pm 0.0132	0.649 \pm 0.0252	0.650 \pm 0.0254
GigaPath	0.788 \pm 0.0164	0.722 \pm 0.0227	0.623 \pm 0.0221	0.620 \pm 0.0240	0.654 \pm 0.0208	0.654 \pm 0.0207
Mean pool (Virchow)	<u>0.799</u> \pm 0.0132	<u>0.726</u> \pm 0.0150	0.627 \pm 0.0144	0.624 \pm 0.0147	0.644 \pm 0.0138	0.645 \pm 0.0134
PRISM	0.784 \pm 0.0138	0.711 \pm 0.0196	0.635 \pm 0.0273	0.633 \pm 0.0281	0.680 \pm 0.0110	<u>0.681</u> \pm 0.0109
Mean pool (CHIEF)	0.750 \pm 0.0178	0.674 \pm 0.0157	0.642 \pm 0.0128	0.638 \pm 0.0146	0.641 \pm 0.0182	0.640 \pm 0.0174
CHIEF	0.789 \pm 0.0277	0.714 \pm 0.0252	0.646 \pm 0.0120	0.645 \pm 0.0119	0.669 \pm 0.0174	0.668 \pm 0.0169
Mean pool (CONCH)	0.771 \pm 0.0187	0.710 \pm 0.0142	0.625 \pm 0.0222	0.624 \pm 0.0222	0.635 \pm 0.0167	0.635 \pm 0.0164
TITAN _V	0.790 \pm 0.0081	0.705 \pm 0.0106	<u>0.672</u> \pm 0.0098	<u>0.669</u> \pm 0.0104	<u>0.681</u> \pm 0.0209	<u>0.681</u> \pm 0.0212
TITAN	0.775 \pm 0.0205	0.713 \pm 0.0192	0.674 \pm 0.0114	0.674 \pm 0.0110	0.685 \pm 0.0237	0.685 \pm 0.0232

Extended Data Table 30: Linear probing results for PBRM1 mutation ($C = 2$) prediction on MUT-HET-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.714 \pm 0.0193	0.559 \pm 0.0158	0.626 \pm 0.0339	0.636 \pm 0.0319	0.512 \pm 0.0261	0.639 \pm 0.0267
GigaPath	0.719 \pm 0.0268	0.574 \pm 0.0266	<u>0.651</u> \pm 0.0340	0.655 \pm 0.0314	0.526 \pm 0.0109	0.652 \pm 0.0130
Mean pool (Virchow)	0.715 \pm 0.0199	0.559 \pm 0.0297	0.628 \pm 0.0282	0.638 \pm 0.0252	0.517 \pm 0.0058	0.644 \pm 0.0055
PRISM	0.709 \pm 0.0103	0.563 \pm 0.0106	0.633 \pm 0.0252	0.651 \pm 0.0238	0.552 \pm 0.0149	0.677 \pm 0.0142
Mean pool (CHIEF)	0.711 \pm 0.0120	0.560 \pm 0.0120	0.645 \pm 0.0232	0.656 \pm 0.0277	0.511 \pm 0.0070	0.637 \pm 0.0079
CHIEF	<u>0.730</u> \pm 0.0249	0.565 \pm 0.0155	0.644 \pm 0.0174	0.673 \pm 0.0039	0.541 \pm 0.0253	0.667 \pm 0.0253
Mean pool (CONCH)	0.709 \pm 0.0290	0.568 \pm 0.0142	0.647 \pm 0.0296	0.669 \pm 0.0201	0.559 \pm 0.0360	0.684 \pm 0.0360
TITAN _V	0.734 \pm 0.0155	0.592 \pm 0.0218	<u>0.651</u> \pm 0.0195	0.686 \pm 0.0140	<u>0.561</u> \pm 0.0145	<u>0.686</u> \pm 0.0128
TITAN	0.727 \pm 0.0154	<u>0.586</u> \pm 0.0158	0.661 \pm 0.0147	<u>0.685</u> \pm 0.0069	0.572 \pm 0.0190	0.695 \pm 0.0173

Extended Data Table 31: Linear probing results for SETD2 mutation ($C = 2$) prediction on MUT-HET-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.893±0.0285	0.711±0.0441	0.732±0.0451	0.767±0.0376	0.595±0.0255	0.771±0.0193
GigaPath	0.878±0.0251	0.684±0.0331	0.724±0.0494	0.764±0.0377	0.560±0.0219	0.744±0.0177
Mean pool (Virchow)	0.879±0.0150	0.691±0.0433	0.666±0.0477	0.701±0.0442	0.534±0.0184	0.723±0.0163
PRISM	0.901±0.0365	0.748±0.0535	0.779±0.0456	0.792±0.0348	0.727±0.0168	0.843±0.0132
Mean pool (CHIEF)	0.852±0.0331	0.660±0.0462	0.686±0.0461	0.719±0.0417	0.558±0.0206	0.743±0.0165
CHIEF	0.880±0.0200	0.727±0.0361	0.749±0.0202	0.766±0.0206	0.654±0.0313	0.801±0.0207
Mean pool (CONCH)	0.887±0.0249	0.732±0.0367	0.751±0.0195	0.768±0.0149	0.678±0.0320	0.806±0.0157
TITAN _V	<u>0.917±0.0248</u>	0.769±0.0394	<u>0.800±0.0230</u>	0.813±0.0239	<u>0.764±0.0359</u>	0.848±0.0190
TITAN	0.920±0.0191	<u>0.764±0.0280</u>	0.805±0.0207	<u>0.802±0.0263</u>	0.775±0.0451	<u>0.844±0.0313</u>

Extended Data Table 32: Linear probing results for ER ($C = 2$) prediction on BCNC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.841±0.0263	0.716±0.0364	0.704±0.0533	0.742±0.0431	0.595±0.0212	0.731±0.0174
GigaPath	0.827±0.0298	0.705±0.0325	0.698±0.0540	0.736±0.0453	0.568±0.0139	0.708±0.0121
Mean pool (Virchow)	0.820±0.0196	0.682±0.0232	0.638±0.0486	0.674±0.0474	0.546±0.0115	0.690±0.0110
PRISM	0.844±0.0322	0.734±0.0279	0.730±0.0345	0.749±0.0257	0.714±0.0325	<u>0.807±0.0221</u>
Mean pool (CHIEF)	0.808±0.0441	0.697±0.0392	0.657±0.0535	0.694±0.0475	0.563±0.0106	0.705±0.0107
CHIEF	0.824±0.0362	0.713±0.0196	0.717±0.0299	0.738±0.0235	0.654±0.0292	0.770±0.0206
Mean pool (CONCH)	0.821±0.0400	0.699±0.0467	0.703±0.0398	0.727±0.0279	0.659±0.0451	0.764±0.0358
TITAN _V	<u>0.856±0.0378</u>	0.741±0.0332	<u>0.765±0.0474</u>	0.783±0.0355	0.752±0.0346	0.822±0.0267
TITAN	0.857±0.0336	<u>0.739±0.0423</u>	0.770±0.0362	<u>0.775±0.0282</u>	<u>0.746±0.0435</u>	<u>0.807±0.0337</u>

Extended Data Table 33: Linear probing results for PR ($C = 2$) prediction on BCNC. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.745±0.0390	0.593±0.0322	0.631±0.0402	0.662±0.0198	0.538±0.0212	0.672±0.0209
GigaPath	0.725±0.0375	0.569±0.0355	0.620±0.0393	0.656±0.0215	0.528±0.0086	0.662±0.0103
Mean pool (Virchow)	0.722±0.0233	0.557±0.0170	0.589±0.0282	0.624±0.0304	0.514±0.0213	0.646±0.0230
PRISM	0.755±0.0268	0.615±0.0373	0.694±0.0326	0.727±0.0331	0.567±0.0094	0.699±0.0070
Mean pool (CHIEF)	0.695±0.0427	0.559±0.0260	0.590±0.0247	0.621±0.0198	0.518±0.0131	0.652±0.0149
CHIEF	0.702±0.0496	0.573±0.0264	0.637±0.0347	0.676±0.0302	0.533±0.0267	0.667±0.0246
Mean pool (CONCH)	0.731±0.0225	0.584±0.0335	0.632±0.0332	0.666±0.0172	0.558±0.0099	0.691±0.0083
TITAN _V	0.762±0.0242	<u>0.630±0.0141</u>	<u>0.699±0.0154</u>	<u>0.723±0.0121</u>	<u>0.594±0.0167</u>	<u>0.720±0.0110</u>
TITAN	<u>0.761±0.0244</u>	0.639±0.0420	0.703±0.0277	0.718±0.0202	0.612±0.0199	0.731±0.0160

Extended Data Table 34: Linear probing results for HER2 ($C = 2$) prediction on BCNC. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.960±0.0255	0.886±0.0426	0.795±0.0258	0.813±0.0252	0.737±0.0105	0.747±0.0151
GigaPath	TCGA	0.956±0.0318	0.887±0.0432	0.783±0.0391	0.795±0.0440	0.723±0.0132	0.726±0.0219
Mean pool (Virchow)	TCGA	0.925±0.0499	0.851±0.0522	0.722±0.0626	0.725±0.0730	0.671±0.0119	0.692±0.0071
PRISM	TCGA	0.936±0.0253	0.851±0.0435	0.774±0.0731	0.776±0.0749	0.851±0.0040	0.854±0.0057
Mean pool (CHIEF)	TCGA	0.922±0.0197	0.819±0.0361	0.665±0.0224	0.690±0.0491	0.720±0.0137	0.735±0.0157
CHIEF	TCGA	0.932±0.0245	0.848±0.0332	0.704±0.0254	0.725±0.0477	0.765±0.0152	0.769±0.0217
Mean pool (CONCH)	TCGA	0.942±0.0256	0.865±0.0356	0.785±0.0462	0.788±0.0531	0.757±0.0070	0.783±0.0046
TITAN _v	TCGA	<u>0.966±0.0192</u>	<u>0.896±0.0394</u>	<u>0.851±0.0528</u>	<u>0.864±0.0563</u>	<u>0.867±0.0046</u>	<u>0.871±0.0065</u>
TITAN	TCGA	0.972±0.0237	0.919±0.0321	0.873±0.0521	0.880±0.0583	0.896±0.0047	0.899±0.0021
Mean pool (GigaPath)	EBRAINS	0.924±0.0074	0.838±0.0223	0.737±0.0447	0.753±0.0409	0.642±0.0288	0.574±0.0589
GigaPath	EBRAINS	0.916±0.0069	0.833±0.0144	0.744±0.0367	0.742±0.0364	0.659±0.0282	0.688±0.0269
Mean pool (Virchow)	EBRAINS	0.873±0.0196	0.788±0.0135	0.733±0.0111	0.729±0.0090	0.528±0.0079	0.298±0.0223
PRISM	EBRAINS	0.934±0.0085	<u>0.874±0.0100</u>	0.814±0.0105	0.794±0.0143	0.821±0.0232	0.799±0.0339
Mean pool (CHIEF)	EBRAINS	0.908±0.0041	0.805±0.0310	0.729±0.0140	0.731±0.0090	0.680±0.0301	0.625±0.0487
CHIEF	EBRAINS	0.923±0.0068	0.836±0.0144	0.753±0.0103	0.748±0.0040	0.754±0.0259	0.724±0.0392
Mean pool (CONCH)	EBRAINS	0.892±0.0270	0.788±0.0306	0.761±0.0051	0.780±0.0039	0.760±0.0435	0.727±0.0727
TITAN _v	EBRAINS	<u>0.951±0.0044</u>	0.848±0.0557	<u>0.874±0.0034</u>	<u>0.875±0.0019</u>	<u>0.877±0.0066</u>	<u>0.875±0.0095</u>
TITAN	EBRAINS	0.960±0.0026	0.908±0.0077	0.883±0.0023	0.881±0.0053	0.899±0.0044	0.893±0.0056

Extended Data Table 35: Linear probing results for Isocitrate dehydrogenase (IDH) mutation ($C = 2$) prediction on TCGA-GBMLGG and external dataset EBRAINS. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.871±0.0499	0.734±0.0761	0.751±0.0543	0.758±0.0809	0.503±0.0069	0.602±0.0090
GigaPath	TCGA	0.877±0.0443	0.723±0.0615	0.744±0.0540	0.760±0.0855	0.525±0.0248	0.628±0.0307
Mean pool (Virchow)	TCGA	0.868±0.0543	0.707±0.0467	0.641±0.0658	0.648±0.1091	0.500±0.0000	0.597±0.0000
PRISM	TCGA	0.889±0.0425	<u>0.773±0.0385</u>	0.800±0.0416	0.808±0.0383	0.650±0.0361	0.745±0.0258
Mean pool (CHIEF)	TCGA	0.859±0.0467	0.703±0.0668	0.686±0.0301	0.708±0.0521	0.599±0.0169	0.702±0.0157
CHIEF	TCGA	<u>0.882±0.0403</u>	0.759±0.0731	0.725±0.0366	0.738±0.0571	<u>0.692±0.0158</u>	0.778±0.0100
Mean pool (CONCH)	TCGA	0.870±0.0494	0.756±0.0417	0.765±0.0616	0.782±0.0674	0.632±0.0390	0.717±0.0099
TITAN _v	TCGA	<u>0.882±0.0678</u>	0.782±0.0486	0.815±0.0763	0.828±0.0761	0.657±0.0830	0.727±0.0390
TITAN	TCGA	0.878±0.0675	0.772±0.0656	<u>0.814±0.0605</u>	<u>0.813±0.0588</u>	0.759±0.0206	0.718±0.0446
Mean pool (GigaPath)	CPTAC	0.881±0.0266	0.618±0.0450	0.792±0.0140	0.813±0.0379	0.507±0.0084	0.606±0.0110
GigaPath	CPTAC	<u>0.873±0.0222</u>	0.713±0.0352	<u>0.783±0.0230</u>	<u>0.775±0.0515</u>	0.555±0.0233	0.663±0.0266
Mean pool (Virchow)	CPTAC	0.860±0.0170	<u>0.750±0.0423</u>	0.617±0.0521	0.660±0.0329	0.508±0.0152	0.607±0.0204
PRISM	CPTAC	0.852±0.0187	0.708±0.0554	0.743±0.0149	0.708±0.0397	0.568±0.0128	0.679±0.0128
Mean pool (CHIEF)	CPTAC	0.808±0.0128	0.664±0.0457	0.645±0.0184	0.528±0.0365	0.555±0.0207	0.665±0.0224
CHIEF	CPTAC	0.840±0.0125	0.686±0.0597	0.746±0.0157	0.738±0.0205	0.529±0.0112	0.636±0.0130
Mean pool (CONCH)	CPTAC	0.812±0.0554	0.765±0.0424	0.603±0.0061	0.386±0.0132	0.510±0.0084	0.611±0.0110
TITAN _v	CPTAC	0.832±0.0249	0.685±0.0548	0.670±0.0182	0.541±0.0310	<u>0.591±0.0266</u>	<u>0.699±0.0252</u>
TITAN	CPTAC	0.838±0.0284	0.735±0.0315	0.690±0.0109	0.571±0.0241	0.776±0.0167	0.757±0.0082

Extended Data Table 36: Linear probing results for ER ($C = 2$) prediction on TCGA-BRCA and external dataset CPTAC-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.791±0.0396	0.688±0.0509	0.685±0.0420	0.686±0.0530	0.536±0.0268	0.553±0.0406
GigaPath	TCGA	0.793±0.0495	0.705±0.0499	0.686±0.0356	0.693±0.0460	0.563±0.0314	0.594±0.0454
Mean pool (Virchow)	TCGA	0.779±0.0394	0.692±0.0364	0.598±0.0665	0.587±0.1042	0.510±0.0137	0.515±0.0235
PRISM	TCGA	0.808 ±0.0403	0.725±0.0341	<u>0.739</u> ±0.0487	<u>0.749</u> ±0.0465	0.643±0.0216	0.655±0.0220
Mean pool (CHIEF)	TCGA	0.762±0.0513	0.668±0.0329	0.644±0.0401	0.653±0.0478	<u>0.667</u> ±0.0456	<u>0.707</u> ±0.0396
CHIEF	TCGA	<u>0.807</u> ±0.0413	0.728 ±0.0288	0.673±0.0277	0.685±0.0315	0.702 ±0.0401	0.741 ±0.0293
Mean pool (CONCH)	TCGA	0.778±0.0528	0.692±0.0457	0.687±0.0546	0.699±0.0853	0.615±0.0280	0.657±0.0237
TITAN _v	TCGA	0.802±0.0436	<u>0.727</u> ±0.0368	0.735±0.0631	0.752 ±0.0824	0.630±0.0334	0.646±0.0251
TITAN	TCGA	0.785±0.0609	0.703±0.0335	0.745 ±0.0498	0.747±0.0643	0.659±0.0186	0.608±0.0220
Mean pool (GigaPath)	CPTAC	0.784±0.0119	0.561±0.0453	0.738 ±0.0296	0.771 ±0.0163	0.541±0.0269	0.562±0.0386
GigaPath	CPTAC	0.777±0.0217	0.667±0.0378	<u>0.733</u> ±0.0284	<u>0.733</u> ±0.0526	0.568±0.0361	0.601±0.0468
Mean pool (Virchow)	CPTAC	0.773±0.0098	<u>0.702</u> ±0.0363	0.584±0.0423	0.573±0.0685	0.527±0.0340	0.541±0.0513
PRISM	CPTAC	<u>0.799</u> ±0.0162	0.713 ±0.0164	0.627±0.0208	0.528±0.0374	0.635±0.0464	<u>0.683</u> ±0.0453
Mean pool (CHIEF)	CPTAC	<u>0.799</u> ±0.0161	0.651±0.0397	0.592±0.0288	0.443±0.0585	0.542±0.0145	0.568±0.0240
CHIEF	CPTAC	0.804 ±0.0145	0.686±0.0400	0.679±0.0087	0.653±0.0143	0.549±0.0270	0.579±0.0329
Mean pool (CONCH)	CPTAC	0.746±0.0371	0.695±0.0632	0.544±0.0079	0.329±0.0177	0.534±0.0173	0.560±0.0195
TITAN _v	CPTAC	0.754±0.0222	0.692±0.0442	0.541±0.0132	0.387±0.0302	<u>0.666</u> ±0.0382	0.712 ±0.0343
TITAN	CPTAC	0.755±0.0230	0.677±0.0235	0.597±0.0174	0.476±0.0302	0.679 ±0.0309	0.655±0.0341

Extended Data Table 37: Linear probing results for PR ($C = 2$) prediction on TCGA-BRCA and external dataset CPTAC-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.643±0.0608	0.509±0.0176	0.581±0.0620	0.652±0.0428	0.500±0.0000	0.815±0.0000
GigaPath	TCGA	0.635±0.0339	0.497±0.0039	0.557±0.0346	0.629±0.0509	0.492±0.0057	0.808±0.0050
Mean pool (Virchow)	TCGA	0.522±0.0475	0.500±0.0000	0.484±0.0348	0.594±0.0931	0.500±0.0000	0.815±0.0000
PRISM	TCGA	0.687 ±0.0309	0.547 ±0.0362	0.604 ±0.0230	0.652±0.0418	0.556 ±0.0373	0.830 ±0.0119
Mean pool (CHIEF)	TCGA	0.649±0.0382	0.510±0.0280	0.565±0.0424	0.633±0.0716	0.500±0.0000	0.815±0.0000
CHIEF	TCGA	0.659±0.0305	0.514±0.0371	0.548±0.0346	0.636±0.0965	0.500±0.0000	0.815±0.0000
Mean pool (CONCH)	TCGA	0.628±0.0160	0.502±0.0074	0.582±0.0764	0.605±0.0432	0.498±0.0027	0.813±0.0024
TITAN _v	TCGA	0.657±0.0363	0.513±0.0389	<u>0.599</u> ±0.0272	0.660 ±0.0557	0.502±0.0266	0.810±0.0105
TITAN	TCGA	<u>0.661</u> ±0.0324	<u>0.529</u> ±0.0318	0.598±0.0570	<u>0.657</u> ±0.0635	<u>0.504</u> ±0.0145	<u>0.816</u> ±0.0078
Mean pool (GigaPath)	CPTAC	0.625±0.0185	0.534±0.0438	0.535±0.0247	0.180±0.1111	0.491±0.0134	0.793±0.0214
GigaPath	CPTAC	0.623±0.0629	<u>0.554</u> ±0.0640	0.516±0.0349	0.210±0.1196	0.500±0.0000	0.815±0.0000
Mean pool (Virchow)	CPTAC	0.512±0.0292	0.495±0.0094	0.501±0.0275	0.249±0.1759	0.498±0.0176	0.806±0.0105
PRISM	CPTAC	<u>0.717</u> ±0.0264	0.555 ±0.0821	0.649 ±0.0243	<u>0.449</u> ±0.0737	0.521 ±0.0112	0.822 ±0.0056
Mean pool (CHIEF)	CPTAC	0.505±0.0450	0.500±0.0030	0.445±0.0312	0.406±0.0976	0.500±0.0000	0.815±0.0000
CHIEF	CPTAC	0.606±0.0319	0.523±0.0548	0.496±0.0317	0.544 ±0.1194	0.500±0.0000	0.815±0.0000
Mean pool (CONCH)	CPTAC	0.582±0.0718	0.508±0.0095	0.480±0.0271	0.298±0.1414	<u>0.511</u> ±0.0357	0.816±0.0216
TITAN _v	CPTAC	0.755 ±0.0348	0.534±0.0775	0.576±0.0549	0.299±0.1709	<u>0.511</u> ±0.0226	<u>0.820</u> ±0.0142
TITAN	CPTAC	0.669±0.0463	0.544±0.1170	<u>0.617</u> ±0.0282	0.397±0.0564	0.499±0.0022	0.814±0.0019

Extended Data Table 38: Linear probing results for HER2 ($C = 2$) prediction on TCGA-BRCA and external dataset CPTAC-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		5-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.597±0.0319	0.507±0.0110	0.546±0.0183	0.577±0.0178	0.477±0.0149	0.525±0.0154
GigaPath	TCGA	0.597±0.0392	0.504±0.0129	0.558±0.0294	0.583±0.0267	0.497±0.0092	0.547±0.0178
Mean pool (Virchow)	TCGA	0.597±0.0290	0.508±0.0112	0.557±0.0390	0.590±0.0224	0.493±0.0429	0.546±0.0465
PRISM	TCGA	0.656 ±0.0292	0.543 ±0.0299	0.603±0.0186	0.614±0.0252	0.525±0.0132	<u>0.576</u> ±0.0137
Mean pool (CHIEF)	TCGA	0.587±0.0400	0.502±0.0081	0.576±0.0223	0.589±0.0318	0.509±0.0234	<u>0.562</u> ±0.0224
CHIEF	TCGA	0.632±0.0289	<u>0.522</u> ±0.0187	0.606±0.0166	<u>0.615</u> ±0.0156	0.492±0.0316	0.546±0.0340
Mean pool (CONCH)	TCGA	0.615±0.0569	0.503±0.0124	0.593±0.0355	0.587±0.0382	0.515±0.0285	0.556±0.0145
TITAN _v	TCGA	0.636±0.0582	0.505±0.0134	<u>0.609</u> ±0.0472	0.605±0.0250	<u>0.537</u> ±0.0159	0.573±0.0278
TITAN	TCGA	<u>0.645</u> ±0.0545	0.513±0.0140	0.619 ±0.0395	0.620 ±0.0193	0.560 ±0.0207	0.605 ±0.0199
Mean pool (GigaPath)	CPTAC	0.639±0.0222	0.500±0.0000	<u>0.564</u> ±0.0296	0.613±0.0320	0.497±0.0116	0.551±0.0114
GigaPath	CPTAC	0.651±0.0145	0.499±0.0030	0.573 ±0.0297	0.620 ±0.0324	<u>0.556</u> ±0.0096	0.603 ±0.0111
Mean pool (Virchow)	CPTAC	0.652±0.0162	0.509 ±0.0227	0.554±0.0336	0.602±0.0384	0.486±0.0061	0.512±0.0056
PRISM	CPTAC	0.619±0.0271	0.503±0.0056	0.506±0.0088	0.529±0.0133	0.532±0.0078	0.579±0.0089
Mean pool (CHIEF)	CPTAC	0.620±0.0162	0.503±0.0051	0.517±0.0186	0.559±0.0268	0.456±0.0220	0.507±0.0221
CHIEF	CPTAC	0.639±0.0206	0.503±0.0050	0.541±0.0125	0.590±0.0154	0.562 ±0.0325	<u>0.602</u> ±0.0280
Mean pool (CONCH)	CPTAC	0.637±0.0197	0.499±0.0030	0.531±0.0110	0.570±0.0169	0.508±0.0168	0.519±0.0217
TITAN _v	CPTAC	0.699 ±0.0136	<u>0.507</u> ±0.0133	0.547±0.0218	0.592±0.0290	0.477±0.0363	0.449±0.0309
TITAN	CPTAC	<u>0.693</u> ±0.0150	0.506±0.0111	<u>0.564</u> ±0.0223	<u>0.619</u> ±0.0229	0.521±0.0215	0.575±0.0223

Extended Data Table 39: Linear probing results for PIK3CA ($C = 2$) prediction on TCGA-BRCA and external dataset CPTAC-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		5-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.638±0.0945	0.509 ±0.0175	0.565±0.0733	0.708±0.0464	0.530±0.0135	0.551±0.0213
GigaPath	TCGA	0.656±0.1019	0.499±0.0025	0.557±0.0908	0.698±0.0534	0.520±0.0183	0.539±0.0305
Mean pool (Virchow)	TCGA	0.628±0.0470	0.498±0.0049	0.544±0.0747	0.658±0.0779	0.524±0.0188	0.539±0.0286
PRISM	TCGA	<u>0.716</u> ±0.0764	0.500±0.0000	0.608±0.0738	0.703±0.0243	0.593 ±0.0124	0.636 ±0.0147
Mean pool (CHIEF)	TCGA	0.641±0.0577	0.496±0.0074	0.584±0.0741	0.696±0.0447	0.516±0.0207	0.534±0.0325
CHIEF	TCGA	0.675±0.0286	0.496±0.0049	0.630±0.0727	0.731 ±0.0374	0.544±0.0365	0.576±0.0442
Mean pool (CONCH)	TCGA	0.714±0.0404	0.499±0.0025	0.644±0.0996	0.706±0.0271	0.552±0.0261	0.580±0.0356
TITAN _v	TCGA	0.721 ±0.0510	0.500±0.0000	0.672 ±0.0215	<u>0.725</u> ±0.0231	0.572±0.0176	0.607±0.0236
TITAN	TCGA	0.700±0.0670	<u>0.508</u> ±0.0167	<u>0.661</u> ±0.0772	0.717±0.0262	<u>0.589</u> ±0.0212	<u>0.628</u> ±0.0283
Mean pool (GigaPath)	CPTAC	0.780±0.0365	0.526 ±0.0281	<u>0.726</u> ±0.0329	0.759 ±0.0322	0.614 ±0.0195	0.661 ±0.0210
GigaPath	CPTAC	0.787±0.0368	<u>0.521</u> ±0.0264	0.718±0.0274	<u>0.753</u> ±0.0233	0.553±0.0180	0.582±0.0250
Mean pool (Virchow)	CPTAC	0.697±0.0509	0.506±0.0164	0.593±0.0545	0.541±0.1452	0.517±0.0404	0.532±0.0615
PRISM	CPTAC	0.796±0.0101	0.500±0.0000	0.665±0.0116	0.680±0.0143	<u>0.608</u> ±0.0593	<u>0.636</u> ±0.0601
Mean pool (CHIEF)	CPTAC	0.730±0.0209	0.513±0.0256	0.679±0.0035	0.705±0.0076	0.533±0.0281	0.566±0.0355
CHIEF	CPTAC	0.731±0.0041	0.515±0.0150	0.653±0.0215	0.668±0.0180	0.547±0.0173	0.583±0.0176
Mean pool (CONCH)	CPTAC	0.764±0.0249	0.513±0.0140	0.697±0.0081	0.736±0.0087	0.526±0.0081	0.539±0.0124
TITAN _v	CPTAC	<u>0.806</u> ±0.0234	0.503±0.0051	0.711±0.0092	0.742±0.0092	0.587±0.0274	0.625±0.0339
TITAN	CPTAC	0.815 ±0.0252	0.505±0.0103	0.727 ±0.0168	0.748±0.0129	0.516±0.0207	0.525±0.0304

Extended Data Table 40: Linear probing results for epidermal growth factor receptor (EGFR) ($C = 2$) prediction on TCGA-NSCLC and external dataset CPTAC-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	<u>0.763</u> ±0.0522	0.693±0.0390	0.652±0.0200	0.649±0.0203	0.653±0.0227	0.657±0.0249
GigaPath	TCGA	0.738±0.0372	0.670±0.0371	0.639±0.0240	0.634±0.0271	0.610±0.0305	0.613±0.0366
Mean pool (Virchow)	TCGA	0.702±0.0687	0.642±0.0535	0.564±0.0176	0.553±0.0207	0.603±0.0224	0.593±0.0406
PRISM	TCGA	0.765 ±0.0875	0.706 ±0.0708	0.677±0.0706	0.678±0.0704	0.658±0.0239	<u>0.666</u> ±0.0229
Mean pool (CHIEF)	TCGA	0.722±0.0584	0.660±0.0692	0.641±0.0439	0.638±0.0486	0.616±0.0363	<u>0.628</u> ±0.0347
CHIEF	TCGA	0.757±0.0481	0.690±0.0407	0.666±0.0631	0.665±0.0683	0.601±0.0195	0.616±0.0182
Mean pool (CONCH)	TCGA	0.734±0.0784	0.663±0.0508	<u>0.682</u> ±0.0549	<u>0.680</u> ±0.0553	0.625±0.0178	0.640±0.0188
TITAN _v	TCGA	0.765 ±0.0589	<u>0.702</u> ±0.0487	0.683 ±0.0515	0.682 ±0.0501	0.710 ±0.0239	0.719 ±0.0233
TITAN	TCGA	0.758±0.0628	<u>0.699</u> ±0.0535	0.671±0.0491	0.671±0.0493	<u>0.662</u> ±0.0094	0.664±0.0093
Mean pool (GigaPath)	CPTAC	0.751±0.0174	0.670±0.0301	<u>0.686</u> ±0.0109	<u>0.690</u> ±0.0070	0.607±0.0152	0.606±0.0170
GigaPath	CPTAC	0.752±0.0093	0.677±0.0248	0.684±0.0159	<u>0.690</u> ±0.0217	0.634±0.0176	0.655±0.0181
Mean pool (Virchow)	CPTAC	0.728±0.0176	0.629±0.0429	0.556±0.0617	0.407±0.1610	0.572±0.0422	0.530±0.0883
PRISM	CPTAC	<u>0.756</u> ±0.0056	<u>0.693</u> ±0.0222	0.668±0.0091	0.671±0.0105	0.692 ±0.0156	0.712 ±0.0150
Mean pool (CHIEF)	CPTAC	0.675±0.0163	0.625±0.0237	0.618±0.0170	0.620±0.0167	0.558±0.0106	0.578±0.0104
CHIEF	CPTAC	0.714±0.0307	0.647±0.0329	0.615±0.0195	0.610±0.0207	0.565±0.0134	0.578±0.0163
Mean pool (CONCH)	CPTAC	0.738±0.0204	0.669±0.0285	0.714 ±0.0122	0.729 ±0.0118	0.606±0.0105	0.622±0.0134
TITAN _v	CPTAC	0.787 ±0.0126	0.713 ±0.0148	0.679±0.0017	0.686±0.0004	<u>0.670</u> ±0.0208	<u>0.689</u> ±0.0204
TITAN	CPTAC	0.728±0.0088	0.689±0.0137	0.669±0.0048	0.673±0.0045	0.667±0.0289	0.687±0.0268

Extended Data Table 41: Linear probing results for *TP53* mutation ($C = 2$) prediction on TCGA-NSCLC and external dataset CPTAC-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		5-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.906 ±0.0328	0.712 ±0.0843	<u>0.780</u> ±0.0388	<u>0.827</u> ±0.0681	0.542±0.0236	0.708±0.0214
GigaPath	TCGA	0.890±0.0428	0.687±0.0861	0.750±0.0370	0.800±0.0603	0.585±0.0344	0.735±0.0251
Mean pool (Virchow)	TCGA	0.795±0.0480	0.557±0.0476	0.602±0.0944	0.617±0.1341	0.550±0.0145	0.703±0.0183
PRISM	TCGA	<u>0.893</u> ±0.0440	0.695±0.0459	0.694±0.0400	0.703±0.0895	0.601 ±0.0398	0.758 ±0.0358
Mean pool (CHIEF)	TCGA	0.829±0.0586	0.598±0.0562	0.614±0.0248	0.670±0.1197	<u>0.597</u> ±0.0294	0.706±0.0372
CHIEF	TCGA	0.873±0.0403	0.655±0.0850	0.740±0.0765	0.748±0.1431	0.581±0.0276	0.686±0.0202
Mean pool (CONCH)	TCGA	0.797±0.0296	0.566±0.0624	0.701±0.0260	0.751±0.0252	0.535±0.0186	0.703±0.0177
TITAN _v	TCGA	0.859±0.0507	0.712 ±0.0439	0.728±0.0520	0.780±0.0612	0.574±0.0282	<u>0.737</u> ±0.0255
TITAN	TCGA	0.872±0.0551	<u>0.700</u> ±0.0730	0.800 ±0.0581	0.831 ±0.0480	0.549±0.0114	0.713±0.0093
Mean pool (GigaPath)	CPTAC	0.890 ±0.0274	0.557±0.0360	0.818 ±0.0277	0.828 ±0.0320	<u>0.561</u> ±0.0197	<u>0.725</u> ±0.0177
GigaPath	CPTAC	<u>0.870</u> ±0.0474	0.553±0.0344	0.732±0.0430	0.693±0.0885	0.544±0.0380	0.703±0.0308
Mean pool (Virchow)	CPTAC	0.771±0.0677	0.681 ±0.0648	0.530±0.0033	0.189±0.0122	0.508±0.0224	0.585±0.0525
PRISM	CPTAC	0.830±0.0122	0.615±0.0314	0.592±0.0220	0.450±0.0860	0.535±0.0138	0.703±0.0131
Mean pool (CHIEF)	CPTAC	0.802±0.0423	0.646±0.0516	0.556±0.0336	0.290±0.0996	0.542±0.0161	0.703±0.0152
CHIEF	CPTAC	0.857±0.0234	0.609±0.0631	0.591±0.0520	0.389±0.1169	0.595 ±0.0261	0.749 ±0.0198
Mean pool (CONCH)	CPTAC	0.793±0.0275	<u>0.676</u> ±0.0621	0.636±0.0227	0.540±0.0836	0.550±0.0165	0.704±0.0288
TITAN _v	CPTAC	0.830±0.0337	0.614±0.0472	0.676±0.0369	0.678±0.0871	0.543±0.0159	0.710±0.0143
TITAN	CPTAC	0.828±0.0287	0.653±0.0302	<u>0.771</u> ±0.0105	<u>0.745</u> ±0.0238	0.545±0.0254	0.707±0.0237

Extended Data Table 42: Linear probing results for microsatellite instability (MSI) high ($C = 2$) prediction on TCGA-CRC and external dataset CPTAC-COAD. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		5-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	0.761±0.0916	0.515±0.0188	0.657±0.0818	0.566±0.2283	0.473±0.0183	0.746±0.0156
GigaPath	TCGA	0.724±0.1200	0.506±0.0164	0.596±0.0812	0.530±0.2415	0.513±0.0069	0.773±0.0074
Mean pool (Virchow)	TCGA	0.630±0.1246	0.500±0.0000	0.529±0.0550	0.422±0.2608	0.559±0.0328	0.797±0.0262
PRISM	TCGA	0.764±0.0719	<u>0.532</u> ±0.0434	0.631±0.0604	0.634±0.1285	0.566±0.0284	0.803±0.0210
Mean pool (CHIEF)	TCGA	0.698±0.0901	0.521±0.0354	0.542±0.0633	0.387±0.3078	0.614 ±0.0727	0.780±0.0438
CHIEF	TCGA	0.752±0.0874	0.523±0.0498	0.609±0.0445	0.557±0.1888	0.560±0.0342	0.745±0.0285
Mean pool (CONCH)	TCGA	0.724±0.0771	0.513±0.0180	0.649±0.0843	0.611±0.1636	0.565±0.0088	0.816±0.0035
TITAN _v	TCGA	<u>0.800</u> ±0.0731	0.516±0.0439	<u>0.658</u> ±0.0323	<u>0.694</u> ±0.0701	<u>0.595</u> ±0.0377	0.826 ±0.0216
TITAN	TCGA	0.817 ±0.0717	0.542 ±0.0343	0.770 ±0.0413	0.771 ±0.0502	0.581±0.0187	<u>0.819</u> ±0.0106
Mean pool (GigaPath)	CPTAC	0.697 ±0.0451	0.547±0.0535	0.538±0.0494	0.349±0.1680	0.521±0.0156	0.782±0.0068
GigaPath	CPTAC	0.675±0.0593	0.552±0.0667	0.482±0.0216	0.255±0.1239	0.493±0.0127	0.754±0.0214
Mean pool (Virchow)	CPTAC	0.604±0.0351	0.510±0.0204	0.543±0.0144	0.255±0.1432	0.512±0.0496	0.752±0.0489
PRISM	CPTAC	0.660±0.0243	<u>0.571</u> ±0.0275	<u>0.597</u> ±0.0554	0.461±0.1346	0.490±0.0183	0.761±0.0135
Mean pool (CHIEF)	CPTAC	0.620±0.0485	0.563±0.0584	0.501±0.0108	0.065±0.0139	<u>0.583</u> ±0.0608	0.797±0.0123
CHIEF	CPTAC	<u>0.687</u> ±0.0384	0.549±0.0502	0.507±0.0043	0.065±0.0144	0.598 ±0.0414	0.820 ±0.0235
Mean pool (CONCH)	CPTAC	0.647±0.0442	0.602 ±0.0598	0.512±0.0133	0.372±0.1000	0.504±0.0218	0.775±0.0167
TITAN _v	CPTAC	0.636±0.0361	0.566±0.0390	0.569±0.0372	0.586 ±0.1379	0.497±0.0028	0.771±0.0024
TITAN	CPTAC	0.624±0.0402	0.555±0.0254	0.622 ±0.0232	<u>0.556</u> ±0.0429	0.570±0.0296	<u>0.805</u> ±0.0206

Extended Data Table 43: Linear probing results for *BRAF* mutation ($C = 2$) prediction on TCGA-CRC and external dataset CPTAC-COAD. The best result is marked in bold and the second best is underlined.

Encoder	Cohort	Logistic regression		SimpleShot		20-nearest neighbors	
		AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	TCGA	<u>0.605</u> ±0.0604	<u>0.567</u> ±0.0377	0.510±0.0481	0.383±0.1048	0.497±0.0065	0.480±0.0187
GigaPath	TCGA	0.569±0.0359	0.552±0.0425	0.521±0.0457	0.401±0.1042	0.518±0.0450	0.504±0.0446
Mean pool (Virchow)	TCGA	0.509±0.1003	0.508±0.0166	0.506±0.0403	0.392±0.1362	0.484±0.0120	0.472±0.0265
PRISM	TCGA	0.555±0.0445	0.535±0.0320	0.529±0.0426	0.461±0.0905	0.509±0.0294	<u>0.533</u> ±0.0309
Mean pool (CHIEF)	TCGA	0.575±0.0731	0.510±0.0149	0.528±0.0490	0.411±0.1003	0.495±0.0205	0.501±0.0324
CHIEF	TCGA	0.603±0.0701	0.534±0.0468	0.529±0.0521	0.432±0.0868	0.536 ±0.0368	0.553 ±0.0507
Mean pool (CONCH)	TCGA	0.573±0.0267	0.561±0.0222	<u>0.531</u> ±0.0293	<u>0.478</u> ±0.1036	0.512±0.0299	0.501±0.0490
TITAN _v	TCGA	0.601±0.0315	0.537±0.0220	0.535 ±0.0473	0.501 ±0.0848	0.520±0.0391	0.503±0.0643
TITAN	TCGA	0.630 ±0.0710	0.569 ±0.0560	0.530±0.0440	0.464±0.0984	<u>0.522</u> ±0.0357	0.553 ±0.0366
Mean pool (GigaPath)	CPTAC	0.561±0.0375	0.507±0.0233	0.539±0.0256	0.464±0.0595	0.469±0.0543	0.460±0.0337
GigaPath	CPTAC	0.590±0.0295	0.520±0.0259	0.531±0.0219	0.447±0.1018	0.507±0.0133	0.494±0.0259
Mean pool (Virchow)	CPTAC	0.517±0.0795	0.498±0.0034	0.490±0.0073	0.485±0.0244	0.481±0.0137	0.468±0.0127
PRISM	CPTAC	0.616±0.0508	0.532 ±0.0338	0.528±0.0282	0.544±0.0225	0.503±0.0202	0.511±0.0210
Mean pool (CHIEF)	CPTAC	0.646±0.0376	0.501±0.0021	0.593±0.0154	0.565±0.0603	0.497±0.0160	0.511±0.0184
CHIEF	CPTAC	<u>0.659</u> ±0.0182	0.502±0.0034	0.593±0.0359	0.561±0.0662	0.508±0.0150	0.503±0.0323
Mean pool (CONCH)	CPTAC	0.641±0.0681	0.525±0.0208	0.664 ±0.0554	<u>0.618</u> ±0.0957	0.532 ±0.0052	0.535±0.0171
TITAN _v	CPTAC	0.675 ±0.0350	0.532 ±0.0176	<u>0.626</u> ±0.0147	0.654 ±0.0107	<u>0.526</u> ±0.0206	<u>0.542</u> ±0.0283
TITAN	CPTAC	0.625±0.0192	<u>0.526</u> ±0.0252	0.558±0.0126	0.460±0.0599	0.513±0.0377	0.547 ±0.0372

Extended Data Table 44: Linear probing results for *KRAS* mutation ($C = 2$) prediction on TCGA-CRC and external dataset CPTAC-COAD. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.
Mean pool (GigaPath)	0.555±0.0558	0.364±0.0331	0.373±0.0655	0.330±0.0338	0.283±0.0683	0.269±0.0303
GigaPath	0.462±0.0583	0.311±0.0310	0.354±0.0666	0.314±0.0329	0.270±0.0686	0.290±0.0309
Mean pool (Virchow)	0.432±0.0637	0.294±0.0297	0.124±0.0704	0.194±0.0267	0.081±0.0704	0.233±0.0306
PRISM	<u>0.708±0.0410</u>	0.389±0.0329	0.671±0.0396	0.385±0.0338	0.653±0.0440	0.377±0.0303
Mean pool (CHIEF)	0.533±0.0523	0.316±0.0319	0.339±0.0636	0.297±0.0289	0.321±0.0646	0.285±0.0299
CHIEF	0.662±0.0449	0.379±0.0331	0.462±0.0644	0.343±0.0314	0.492±0.0596	<u>0.352±0.0306</u>
Mean pool (CONCH)	0.442±0.0594	0.344±0.0320	0.383±0.0652	0.285±0.0317	0.388±0.0641	0.256±0.0291
TITAN _V	0.723±0.0390	0.425±0.0345	0.585±0.0540	<u>0.365±0.0325</u>	<u>0.618±0.0492</u>	0.342±0.0306
TITAN	0.702±0.0423	<u>0.394±0.0328</u>	<u>0.631±0.0430</u>	0.329±0.0330	0.613±0.0482	0.377±0.0347

Extended Data Table 45: Linear probing results for ER expression level ($C = 6$) prediction on MGH-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.	Cohen's κ	Bal. acc.
Mean pool (GigaPath)	0.538±0.0555	0.307±0.0307	0.458±0.0540	0.259±0.0291	0.414±0.0637	0.244±0.0310
GigaPath	0.560±0.0514	0.306±0.0298	0.451±0.0576	0.247±0.0293	0.310±0.0693	0.215±0.0281
Mean pool (Virchow)	0.443±0.0622	0.247±0.0277	0.137±0.0694	0.221±0.0279	0.172±0.0683	0.165±0.0242
PRISM	0.594±0.0502	0.345±0.0327	0.536±0.0527	<u>0.359±0.0320</u>	0.556±0.0505	<u>0.350±0.0311</u>
Mean pool (CHIEF)	0.536±0.0575	0.350±0.0315	0.393±0.0575	0.263±0.0277	0.408±0.0598	0.264±0.0309
CHIEF	0.674±0.0478	0.417±0.0341	0.543±0.0534	0.331±0.0300	0.565±0.0551	0.349±0.0305
Mean pool (CONCH)	0.582±0.0512	0.330±0.0320	0.414±0.0552	0.234±0.0295	0.462±0.0548	0.256±0.0287
TITAN _V	0.767±0.0375	0.476±0.0329	0.694±0.0408	0.381±0.0318	0.632±0.0471	0.357±0.0314
TITAN	<u>0.731±0.0361</u>	<u>0.418±0.0318</u>	<u>0.544±0.0512</u>	0.336±0.0313	<u>0.578±0.0465</u>	0.311±0.0305

Extended Data Table 46: Linear probing results for PR expression level ($C = 6$) prediction on MGH-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.732±0.0548	0.595±0.0309	0.657±0.0402	0.669±0.0295	0.567±0.0281	0.663±0.0270
GigaPath	0.717±0.0510	0.580±0.0224	0.642±0.0337	0.659±0.0247	0.555±0.0176	0.652±0.0186
Mean pool (Virchow)	0.730±0.0509	0.588±0.0203	0.573±0.0447	0.605±0.0347	0.558±0.0194	0.653±0.0180
PRISM	<u>0.787±0.0517</u>	0.632±0.0486	0.736±0.0264	0.744±0.0184	0.688±0.0341	0.758±0.0283
Mean pool (CHIEF)	0.685±0.0463	0.562±0.0351	0.585±0.0288	0.608±0.0270	0.557±0.0084	0.653±0.0056
CHIEF	0.732±0.0576	0.612±0.0220	0.635±0.0197	0.654±0.0215	0.595±0.0342	0.685±0.0317
Mean pool (CONCH)	0.735±0.0573	0.577±0.0200	0.647±0.0255	0.664±0.0195	0.612±0.0166	0.703±0.0156
TITAN _V	0.779±0.0442	0.651±0.0269	0.712±0.0164	0.735±0.0117	0.650±0.0294	0.732±0.0264
TITAN	0.789±0.0445	<u>0.648±0.0397</u>	<u>0.729±0.0182</u>	<u>0.738±0.0158</u>	<u>0.684±0.0352</u>	<u>0.752±0.0323</u>

Extended Data Table 47: Linear probing results for ER ($C = 2$) prediction on MGB-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.733±0.0444	0.668±0.0205	0.637±0.0334	0.636±0.0325	0.615±0.0165	0.629±0.0160
GigaPath	0.720±0.0464	0.669±0.0263	0.634±0.0291	0.637±0.0259	0.616±0.0271	0.630±0.0272
Mean pool (Virchow)	0.697±0.0385	0.645±0.0283	0.573±0.0161	0.580±0.0161	0.563±0.0411	0.576±0.0402
PRISM	0.742±0.0398	0.680±0.0342	0.693 ±0.0380	0.694 ±0.0336	<u>0.675</u> ±0.0306	<u>0.686</u> ±0.0304
Mean pool (CHIEF)	0.715±0.0409	0.639±0.0325	0.585±0.0265	0.585±0.0274	0.580±0.0311	0.594±0.0311
CHIEF	0.722±0.0434	0.662±0.0283	0.610±0.0185	0.615±0.0184	0.617±0.0254	0.632±0.0261
Mean pool (CONCH)	0.717±0.0287	0.649±0.0250	0.646±0.0338	0.650±0.0317	0.631±0.0298	0.645±0.0290
TITAN _v	<u>0.752</u> ±0.0329	<u>0.685</u> ±0.0277	0.672±0.0316	0.680±0.0301	0.637±0.0368	0.651±0.0365
TITAN	0.758 ±0.0335	0.686 ±0.0275	<u>0.682</u> ±0.0193	<u>0.685</u> ±0.0160	0.681 ±0.0179	0.693 ±0.0175

Extended Data Table 48: Linear probing results for PR ($C = 2$) prediction on MGB-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		5-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.630±0.0366	0.506±0.0052	0.550±0.0249	0.626±0.0107	0.516±0.0117	0.734±0.0081
GigaPath	0.613±0.0385	0.513±0.0143	0.557±0.0403	0.642±0.0179	0.523±0.0198	0.738±0.0125
Mean pool (Virchow)	0.575±0.0789	0.507±0.0080	0.512±0.0587	0.620±0.0272	0.526±0.0180	0.738±0.0093
PRISM	<u>0.748</u> ±0.0669	0.536 ±0.0334	<u>0.653</u> ±0.0471	0.682 ±0.0240	0.594 ±0.0217	0.781 ±0.0150
Mean pool (CHIEF)	0.568±0.0478	0.500±0.0000	0.527±0.0480	0.593±0.0164	0.516±0.0242	0.731±0.0147
CHIEF	0.696±0.0418	0.513±0.0194	0.613±0.0469	0.663±0.0198	0.529±0.0338	0.740±0.0191
Mean pool (CONCH)	0.654±0.0947	0.516±0.0204	0.533±0.0534	0.605±0.0301	0.530±0.0205	0.741±0.0134
TITAN _v	0.727±0.0681	<u>0.517</u> ±0.0102	0.627±0.0492	<u>0.673</u> ±0.0237	<u>0.554</u> ±0.0527	<u>0.758</u> ±0.0306
TITAN	0.749 ±0.0416	0.502±0.0076	0.662 ±0.0564	0.667±0.0400	<u>0.554</u> ±0.0351	0.753±0.0155

Extended Data Table 49: Linear probing results for HER2 ($C = 2$) prediction on MGB-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		5-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.536±0.0813	0.500±0.0000	0.514±0.1075	0.563±0.0831	0.585±0.0477	0.656±0.0594
GigaPath	0.515±0.1371	0.500±0.0000	0.512±0.0640	0.568±0.0377	0.536±0.0679	0.619±0.0604
Mean pool (Virchow)	0.572±0.1610	0.500±0.0000	0.441±0.1094	0.468±0.1234	0.532±0.1186	0.612±0.1092
PRISM	<u>0.701</u> ±0.1304	0.607 ±0.1044	0.710 ±0.0930	0.748 ±0.0725	<u>0.688</u> ±0.0560	<u>0.714</u> ±0.0322
Mean pool (CHIEF)	0.453±0.1495	0.491±0.0182	0.472±0.0990	0.480±0.1053	0.494±0.0940	0.588±0.0919
CHIEF	0.499±0.1212	0.473±0.0545	0.492±0.0943	0.560±0.0763	0.565±0.0706	0.644±0.0598
Mean pool (CONCH)	0.580±0.0921	0.500±0.0000	0.579±0.1528	0.608±0.1384	0.608±0.0821	0.679±0.0711
TITAN _v	0.638±0.1947	0.525±0.0500	<u>0.672</u> ±0.1171	<u>0.712</u> ±0.0937	0.672±0.1141	0.708±0.0966
TITAN	0.705 ±0.1338	<u>0.567</u> ±0.1233	0.663±0.1054	0.702±0.0910	0.755 ±0.0949	0.774 ±0.0870

Extended Data Table 50: Linear probing results for CDX-2 expression ($C = 2$) prediction on MGB-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.656±0.1128	0.527±0.1181	0.607±0.1957	0.595±0.1991	0.643±0.1580	0.636±0.1555
GigaPath	0.624±0.1530	0.510±0.0873	0.590±0.0742	0.584±0.0754	0.527±0.1618	0.491±0.1760
Mean pool (Virchow)	0.738±0.1436	0.657±0.1138	0.623±0.1158	0.622±0.1173	0.627±0.1587	0.623±0.1569
PRISM	<u>0.908±0.0570</u>	0.833±0.0913	0.917±0.0527	0.914±0.0543	<u>0.850±0.0816</u>	<u>0.838±0.0935</u>
Mean pool (CHIEF)	0.619±0.1362	0.530±0.1481	0.577±0.2099	0.560±0.2227	0.510±0.1555	0.491±0.1659
CHIEF	0.760±0.1603	0.747±0.1376	0.637±0.1267	0.630±0.1264	0.597±0.1540	0.591±0.1527
Mean pool (CONCH)	0.833±0.1394	0.710±0.1659	0.680±0.2031	0.674±0.2068	0.607±0.0910	0.593±0.0796
TITAN _V	0.864±0.0761	0.760±0.0904	0.847±0.0636	0.844±0.0625	0.830±0.0748	0.828±0.0737
TITAN	0.953±0.0295	<u>0.813±0.0968</u>	<u>0.897±0.0323</u>	<u>0.896±0.0317</u>	0.897±0.0323	0.896±0.0317

Extended Data Table 51: Linear probing results for Cytokeratin 5 & 6 expression ($C = 2$) prediction on MGB-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.662±0.1120	0.600±0.0825	0.604±0.0567	0.602±0.0540	0.496±0.0417	0.495±0.0456
GigaPath	0.645±0.1395	0.625±0.1053	0.562±0.0926	0.557±0.0965	0.515±0.0967	0.513±0.0970
Mean pool (Virchow)	0.611±0.0824	0.562±0.0447	0.550±0.0683	0.548±0.0682	0.529±0.0855	0.530±0.0851
PRISM	0.770±0.0670	0.657±0.0578	0.635±0.0644	0.628±0.0657	<u>0.638±0.1001</u>	0.627±0.0958
Mean pool (CHIEF)	0.582±0.1049	0.551±0.0498	0.524±0.0484	0.522±0.0468	0.500±0.0570	0.497±0.0579
CHIEF	0.658±0.0749	0.603±0.0671	0.569±0.1033	0.564±0.1053	0.585±0.0721	0.578±0.0761
Mean pool (CONCH)	0.682±0.0249	0.543±0.0384	0.590±0.0642	0.582±0.0626	0.538±0.1001	0.535±0.1000
TITAN _V	<u>0.772±0.0673</u>	<u>0.687±0.1010</u>	<u>0.684±0.1227</u>	<u>0.680±0.1201</u>	0.633±0.0736	<u>0.631±0.0730</u>
TITAN	0.854±0.0373	0.761±0.0716	0.692±0.1049	0.690±0.1040	0.700±0.0706	0.697±0.0688

Extended Data Table 52: Linear probing results for Napsin A expression ($C = 2$) prediction on MGB-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.646±0.0920	0.591±0.0629	0.472±0.0325	0.477±0.0315	0.503±0.0643	0.519±0.0681
GigaPath	0.579±0.1112	0.549±0.0512	0.459±0.0894	0.470±0.0885	0.473±0.0245	0.471±0.0279
Mean pool (Virchow)	0.673±0.1061	0.597±0.0408	0.444±0.0216	0.449±0.0127	0.507±0.0316	0.515±0.0238
PRISM	0.877±0.0576	<u>0.793±0.0861</u>	0.784±0.0791	0.796±0.0677	<u>0.805±0.0852</u>	<u>0.826±0.0758</u>
Mean pool (CHIEF)	0.503±0.0518	0.522±0.0271	0.393±0.0437	0.397±0.0355	0.470±0.0679	0.479±0.0580
CHIEF	0.759±0.0841	0.680±0.0627	0.609±0.0796	0.613±0.0661	0.567±0.0399	0.582±0.0439
Mean pool (CONCH)	0.840±0.0571	0.754±0.0591	0.567±0.0584	0.578±0.0502	0.566±0.0357	0.585±0.0379
TITAN _V	<u>0.889±0.0387</u>	0.771±0.0602	<u>0.831±0.0498</u>	<u>0.832±0.0461</u>	0.738±0.1012	0.765±0.0949
TITAN	0.913±0.0422	0.803±0.0558	0.872±0.0427	0.876±0.0367	0.870±0.0512	0.884±0.0449

Extended Data Table 53: Linear probing results for P40 expression ($C = 2$) prediction on MGB-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	AUROC	Bal. acc.	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
Mean pool (GigaPath)	0.606±0.1164	0.543±0.0379	0.562±0.1023	0.558±0.0991	0.582±0.0772	0.581±0.0800
GigaPath	0.600±0.1071	0.557±0.0759	0.552±0.1093	0.552±0.1129	0.518±0.0994	0.518±0.1027
Mean pool (Virchow)	0.611±0.1399	0.536±0.1112	0.605±0.1078	0.602±0.1099	0.586±0.0891	0.586±0.0885
PRISM	0.822±0.0681	0.735±0.0846	0.744±0.0822	0.735±0.0915	<u>0.754±0.0974</u>	<u>0.740±0.1031</u>
Mean pool (CHIEF)	0.517±0.0872	0.512±0.0244	0.498±0.0532	0.489±0.0409	0.522±0.0637	0.519±0.0629
CHIEF	0.713±0.1397	0.583±0.0918	0.672±0.1047	0.672±0.1053	0.618±0.0973	0.619±0.0962
Mean pool (CONCH)	0.684±0.1107	0.600±0.0889	0.682±0.1037	0.678±0.1090	0.630±0.0703	0.626±0.0707
TITAN _v	0.866±0.0649	0.782±0.0743	<u>0.767±0.0671</u>	<u>0.763±0.0720</u>	0.713±0.0961	0.713±0.0959
TITAN	<u>0.860±0.0622</u>	<u>0.770±0.1035</u>	0.784±0.0786	0.777±0.0826	0.771±0.0733	0.763±0.0776

Extended Data Table 54: Linear probing results for P63 expression ($C = 2$) prediction on MGB-LUAD. The best result is marked in bold and the second best is underlined.

Extended Data Table 55: Survival prediction Results for TITAN and other slide encoders for measuring disease-specific survival (DSS) with c-index on six TCGA cancer cohorts. Except for CHIEF, all slide encoders have not used TCGA as the pretraining dataset. The best and second-best performances are denoted by **bold** and underlined, respectively.

	CHIEF	GigaPath	PRISM	Mean	TITAN _v	TITAN
BLCA	0.599 ± 0.019	0.589 ± 0.054	<u>0.656</u> ± 0.023	0.630 ± 0.058	0.657 ± 0.055	<u>0.656</u> ± 0.040
BRCA	<u>0.737</u> ± 0.039	0.687 ± 0.083	0.685 ± 0.076	0.711 ± 0.048	0.713 ± 0.044	0.757 ± 0.015
CRC	0.680 ± 0.079	0.628 ± 0.085	0.622 ± 0.125	0.720 ± 0.089	<u>0.710</u> ± 0.110	0.705 ± 0.097
KIRC	0.736 ± 0.061	0.751 ± 0.072	<u>0.769</u> ± 0.061	0.774 ± 0.070	0.774 ± 0.060	0.768 ± 0.064
NSCLC	<u>0.633</u> ± 0.037	0.670 ± 0.031	0.600 ± 0.074	0.596 ± 0.049	0.624 ± 0.082	0.623 ± 0.041
UCEC	0.758 ± 0.109	0.779 ± 0.105	0.747 ± 0.026	0.761 ± 0.133	0.789 ± 0.091	<u>0.787</u> ± 0.067
Avg. (↑)	0.691	0.684	0.680	0.699	<u>0.711</u>	0.716

Encoder	Bal. acc.	Weighted F1	AUROC
ABMIL (CONCHv1.5) (finetuned)	0.807±0.0058	0.854±0.0039	<u>0.991</u> ±0.0004
GigaPath	0.700±0.0069	0.782±0.0046	0.979±0.0010
GigaPath (finetuned)	0.772±0.0061	0.841±0.0039	0.983±0.0008
PRISM	0.774±0.0062	0.824±0.0042	0.989±0.0005
PRISM (finetuned)	0.746±0.0065	0.811±0.0045	0.987±0.0007
CHIEF	0.625±0.0079	0.723±0.0051	0.973±0.0011
CHIEF (finetuned)	0.609±0.0078	0.708±0.0051	0.960±0.0017
TITAN _v	0.820±0.0055	0.870±0.0037	0.994 ±0.0003
TITAN _v (finetuned)	<u>0.836</u> ±0.0054	0.874±0.0037	0.984±0.0010
TITAN	0.832±0.0056	<u>0.881</u> ±0.0036	0.994 ±0.0003
TITAN (finetuned)	0.841 ±0.0055	0.882 ±0.0036	0.989±0.0008

Extended Data Table 56: Finetuning results for tumor subtype ($C = 2$) prediction on TCGA-UT-8K for all baselines. Finetuned results are trained on the specific task, while non-finnetuned results are logistic regression fits of the frozen slide embeddings. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
ABMIL (CONCHv1.5) (finetuned)	0.677±0.0204	0.735±0.0126	0.982±0.0026
GigaPath	0.543±0.0178	0.659±0.0134	0.976±0.0021
GigaPath (finetuned)	0.664±0.0175	0.736±0.0126	0.979±0.0022
PRISM	0.643±0.0181	0.732±0.0129	0.986±0.0014
PRISM (finetuned)	0.548±0.0172	0.644±0.0134	0.980±0.0015
CHIEF	0.528±0.0205	0.640±0.0142	0.968±0.0028
CHIEF (finetuned)	0.518±0.0200	0.630±0.0143	0.944±0.0043
TITAN _v	0.690±0.0196	0.758±0.0123	<u>0.989</u> ±0.0014
TITAN _v (finetuned)	0.695±0.0170	<u>0.767</u> ±0.0117	0.974±0.0028
TITAN	<u>0.704</u> ±0.0192	0.764±0.0116	0.990 ±0.0012
TITAN (finetuned)	0.716 ±0.0182	0.774 ±0.0113	0.985±0.0016

Extended Data Table 57: Finetuning results for OncoTree code ($C = 2$) prediction on TCGA-OT for all baselines. Finetuned results are trained on the specific task, while non-finnetuned results are logistic regression fits of the frozen slide embeddings. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
ABMIL (CONCHv1.5) (finetuned)	0.538±0.0109	0.521±0.0134	0.977±0.0016
GigaPath	0.437±0.0110	0.423±0.0124	0.949±0.0030
GigaPath (finetuned)	0.525±0.0109	0.508±0.0134	0.972±0.0019
PRISM	0.508±0.0110	0.481±0.0133	0.976±0.0018
PRISM (finetuned)	0.492±0.0110	0.478±0.0128	0.967±0.0022
CHIEF	0.413±0.0109	0.394±0.0130	0.944±0.0029
CHIEF (finetuned)	0.330±0.0108	0.319±0.0121	0.916±0.0038
TITAN _v	0.558±0.0104	0.536±0.0127	<u>0.979</u> ±0.0017
TITAN _v (finetuned)	<u>0.575</u> ±0.0103	0.554±0.0129	0.966±0.0023
TITAN	0.587 ±0.0103	0.563 ±0.0130	0.983 ±0.0014
TITAN (finetuned)	0.565±0.0105	0.541±0.0130	<u>0.979</u> ±0.0017

Extended Data Table 58: Finetuning results for OncoTree code ($C = 2$) prediction on OT108 for all baselines. Finetuned results are trained on the specific task, while non-finnetuned results are logistic regression fits of the frozen slide embeddings. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
ABMIL (CONCHv1.5) (finetuned)	0.728±0.0203	0.778±0.0184	0.982±0.0024
GigaPath	0.680±0.0217	0.746±0.0197	0.978±0.0027
GigaPath (finetuned)	0.725±0.0209	0.780±0.0185	0.983±0.0028
PRISM	0.674±0.0200	0.732±0.0191	0.978±0.0026
PRISM (finetuned)	0.652±0.0224	0.707±0.0196	0.967±0.0039
CHIEF	0.598±0.0237	0.670±0.0206	0.969±0.0032
CHIEF (finetuned)	0.544±0.0247	0.622±0.0208	0.941±0.0076
TITAN _v	0.732±0.0208	0.785±0.0175	0.985 ±0.0021
TITAN _v (finetuned)	<u>0.746</u> ±0.0200	0.794 ±0.0171	0.985 ±0.0022
TITAN	0.735±0.0204	0.786±0.0182	0.983±0.0023
TITAN (finetuned)	0.748 ±0.0192	<u>0.787</u> ±0.0176	<u>0.984</u> ±0.0030

Extended Data Table 59: Finetuning results for tumor type ($C = 2$) prediction on EBRAINS for all baselines. Finetuned results are trained on the specific task, while non-finnetuned results are logistic regression fits of the frozen slide embeddings. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (abs)	0.816±0.0059	<u>0.870</u> ±0.0037	0.813±0.0053	0.858±0.0038	0.811±0.0050	0.839±0.0038
TITAN _v (none)	<u>0.819</u> ±0.0054	0.867±0.0038	<u>0.815</u> ±0.0054	<u>0.859</u> ±0.0039	0.814 ±0.0048	0.851 ±0.0037
TITAN _v (rope)	0.818±0.0059	0.873 ±0.0037	0.812±0.0053	<u>0.859</u> ±0.0040	0.804±0.0056	0.836±0.0040
TITAN _v	0.820 ±0.0055	<u>0.870</u> ±0.0037	0.818 ±0.0053	0.864 ±0.0039	<u>0.812</u> ±0.0050	<u>0.842</u> ±0.0038

Extended Data Table 60: Ablation study of positional encodings. Linear probing results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (abs)	<u>0.673</u> ±0.0157	<u>0.753</u> ±0.0126	0.631±0.0184	0.721±0.0130	<u>0.706</u> ±0.0190	0.777 ±0.0116
TITAN _v (none)	0.664±0.0188	0.749±0.0127	<u>0.633</u> ±0.0175	<u>0.726</u> ±0.0129	0.708 ±0.0187	0.777 ±0.0112
TITAN _v (rope)	0.672±0.0170	0.735±0.0127	0.627±0.0164	0.719±0.0129	0.679±0.0207	0.765±0.0115
TITAN _v	0.690 ±0.0196	0.758 ±0.0123	0.639 ±0.0149	0.744 ±0.0124	<u>0.706</u> ±0.0183	<u>0.775</u> ±0.0116

Extended Data Table 61: Ablation study of positional encodings. Linear probing results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (abs)	0.532±0.0106	0.508±0.0133	0.446±0.0102	0.411±0.0130	<u>0.506</u> ±0.0102	<u>0.497</u> ±0.0128
TITAN _v (none)	<u>0.536</u> ±0.0103	<u>0.512</u> ±0.0130	0.450±0.0104	0.414±0.0130	0.488±0.0108	0.483±0.0130
TITAN _v (rope)	0.526±0.0103	0.505±0.0132	0.446±0.0102	<u>0.414</u> ±0.0129	0.495±0.0109	<u>0.497</u> ±0.0129
TITAN _v	0.558 ±0.0104	0.536 ±0.0127	0.464 ±0.0105	0.430 ±0.0128	0.524 ±0.0106	0.517 ±0.0129

Extended Data Table 62: Ablation study of positional encodings. Linear probing results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (abs)	<u>0.728</u> ±0.0203	<u>0.781</u> ±0.0178	0.658 ±0.0190	0.713 ±0.0209	<u>0.722</u> ±0.0212	0.735 ±0.0185
TITAN _v (none)	0.693±0.0197	0.743±0.0195	0.648±0.0195	0.702±0.0209	0.714±0.0204	<u>0.728</u> ±0.0189
TITAN _v (rope)	0.702±0.0198	0.763±0.0187	0.591±0.0212	0.666±0.0212	0.675±0.0205	0.669±0.0198
TITAN _v	0.732 ±0.0208	0.785 ±0.0175	<u>0.656</u> ±0.0187	<u>0.709</u> ±0.0199	0.742 ±0.0192	0.727±0.0182

Extended Data Table 63: Ablation study of positional encodings. Linear probing results for tumor type ($C = 30$) prediction on EBRAINS. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (12.5)	0.803±0.0061	0.862±0.0039	0.784±0.0061	0.829±0.0043	0.791±0.0058	0.826±0.0041
TITAN _v (25)	0.810±0.0059	0.860±0.0040	0.791±0.0063	0.841±0.0042	0.809±0.0057	0.837±0.0041
TITAN _v (50)	<u>0.818±0.0057</u>	<u>0.865±0.0039</u>	<u>0.802±0.0058</u>	<u>0.850±0.0040</u>	0.816±0.0050	<u>0.839±0.0041</u>
TITAN _v	0.820±0.0055	0.870±0.0037	0.818±0.0053	0.864±0.0039	<u>0.812±0.0050</u>	0.842±0.0038

Extended Data Table 64: Ablation study of dataset size. Linear probing results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (12.5)	0.656±0.0184	0.741±0.0123	0.610±0.0173	0.700±0.0132	<u>0.684±0.0187</u>	0.759±0.0118
TITAN _v (25)	0.655±0.0194	0.726±0.0126	0.626±0.0174	0.706±0.0130	0.674±0.0197	0.755±0.0120
TITAN _v (50)	<u>0.678±0.0179</u>	<u>0.743±0.0119</u>	0.644±0.0178	<u>0.735±0.0125</u>	0.706±0.0198	0.783±0.0113
TITAN _v	0.690±0.0196	0.758±0.0123	<u>0.639±0.0149</u>	0.744±0.0124	0.706±0.0183	<u>0.775±0.0116</u>

Extended Data Table 65: Ablation study of dataset size. Linear probing results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (12.5)	0.520±0.0111	0.502±0.0134	0.448±0.0098	0.415±0.0123	0.511±0.0112	0.511±0.0131
TITAN _v (25)	0.526±0.0109	0.507±0.0130	0.452±0.0102	0.421±0.0129	0.515±0.0110	0.514±0.0127
TITAN _v (50)	<u>0.537±0.0102</u>	<u>0.515±0.0127</u>	0.468±0.0101	0.435±0.0123	0.537±0.0109	0.532±0.0127
TITAN _v	0.558±0.0104	0.536±0.0127	<u>0.464±0.0105</u>	<u>0.430±0.0128</u>	<u>0.524±0.0106</u>	<u>0.517±0.0129</u>

Extended Data Table 66: Ablation study of dataset size. Linear probing results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (12.5)	0.723±0.0209	0.773±0.0177	<u>0.630±0.0216</u>	<u>0.690±0.0207</u>	0.708±0.0211	0.714±0.0186
TITAN _v (25)	0.721±0.0204	0.774±0.0191	0.620±0.0203	0.689±0.0202	<u>0.725±0.0208</u>	<u>0.724±0.0184</u>
TITAN _v (50)	0.733±0.0206	<u>0.784±0.0180</u>	0.618±0.0207	0.684±0.0207	0.719±0.0202	0.727±0.0184
TITAN _v	<u>0.732±0.0208</u>	0.785±0.0175	0.656±0.0187	0.709±0.0199	0.742±0.0192	0.727±0.0182

Extended Data Table 67: Ablation study of dataset size. Linear probing results for tumor type ($C = 30$) prediction on EBRAINS. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v (tbase)	0.818±0.0060	<u>0.873±0.0037</u>	0.809±0.0058	0.857±0.0040	0.813±0.0052	0.843±0.0039
TITAN _v	<u>0.820±0.0055</u>	0.870±0.0037	<u>0.818±0.0053</u>	<u>0.864±0.0039</u>	<u>0.812±0.0050</u>	<u>0.842±0.0038</u>
TITAN _v (base)	0.827±0.0052	0.876±0.0037	0.821±0.0052	0.868±0.0037	0.803±0.0053	0.841±0.0039

Extended Data Table 68: Ablation study of model sizes. Linear probing results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _V (tbase)	<u>0.675</u> ±0.0190	<u>0.747</u> ±0.0117	0.643 ±0.0155	<u>0.732</u> ±0.0127	<u>0.706</u> ±0.0189	<u>0.771</u> ±0.0114
TITAN _V	0.690 ±0.0196	0.758 ±0.0123	<u>0.639</u> ±0.0149	0.744 ±0.0124	<u>0.706</u> ±0.0183	<u>0.775</u> ±0.0116
TITAN _V (base)	0.670±0.0178	<u>0.750</u> ±0.0119	<u>0.621</u> ±0.0165	0.723±0.0126	0.709 ±0.0186	0.779 ±0.0117

Extended Data Table 69: Ablation study of model sizes. Linear probing results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _V (tbase)	<u>0.551</u> ±0.0107	<u>0.529</u> ±0.0130	0.467 ±0.0102	<u>0.428</u> ±0.0131	0.528 ±0.0108	0.523 ±0.0127
TITAN _V	0.558 ±0.0104	0.536 ±0.0127	<u>0.464</u> ±0.0105	0.430 ±0.0128	<u>0.524</u> ±0.0106	<u>0.517</u> ±0.0129
TITAN _V (base)	0.539±0.0104	0.521±0.0126	<u>0.449</u> ±0.0100	0.416±0.0128	<u>0.484</u> ±0.0105	<u>0.474</u> ±0.0131

Extended Data Table 70: Ablation study of model sizes. Linear probing results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Logistic regression		SimpleShot		20-nearest neighbors	
	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _V (tbase)	<u>0.709</u> ±0.0211	0.761±0.0187	<u>0.636</u> ±0.0202	<u>0.696</u> ±0.0201	0.693±0.0219	<u>0.701</u> ±0.0192
TITAN _V	0.732 ±0.0208	0.785 ±0.0175	0.656 ±0.0187	0.709 ±0.0199	0.742 ±0.0192	0.727 ±0.0182
TITAN _V (base)	0.732 ±0.0199	<u>0.781</u> ±0.0177	<u>0.621</u> ±0.0203	0.682±0.0209	<u>0.698</u> ±0.0201	0.691±0.0192

Extended Data Table 71: Ablation study of model sizes. Linear probing results for tumor type ($C = 30$) prediction on EBRAINS. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.495±0.0367	0.597±0.0289	0.667±0.0201	0.711±0.0156	0.738±0.0132	0.755±0.0105
ABMIL (CONCH)	0.539±0.0459	0.662±0.0275	0.724±0.0133	0.759±0.0124	0.787±0.0111	0.805±0.0109
GigaPath	0.282±0.0286	0.381±0.0282	0.476±0.0193	0.557±0.0204	0.620±0.0138	0.659±0.0112
PRISM	0.469±0.0380	0.560±0.0284	0.627±0.0253	0.674±0.0181	0.704±0.0151	0.725±0.0121
CHIEF	0.255±0.0292	0.339±0.0304	0.424±0.0198	0.503±0.0153	0.564±0.0135	0.609±0.0113
TITAN _v	<u>0.634±0.0389</u>	<u>0.723±0.0236</u>	<u>0.767±0.0149</u>	<u>0.795±0.0118</u>	<u>0.810±0.0084</u>	<u>0.819±0.0094</u>
TITAN	0.683±0.0357	0.761±0.0223	0.796±0.0134	0.814±0.0131	0.824±0.0106	0.830±0.0096

Extended Data Table 72: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for tumor subtype ($C = 2$) prediction on TCGA-UT-8K. Results are given in balanced accuracy of linear probing evaluation or ABMIL training, respectively. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.396±0.0293	0.466±0.0302	0.521±0.0233	0.555±0.0217	0.583±0.0193	0.600±0.0173
ABMIL (CONCH)	0.407±0.0333	0.505±0.0257	0.563±0.0158	0.607±0.0170	0.634±0.0192	0.659±0.0138
GigaPath	0.210±0.0268	0.289±0.0294	0.373±0.0192	0.438±0.0193	0.502±0.0189	0.537±0.0171
PRISM	0.441±0.0282	0.519±0.0196	0.584±0.0159	0.621±0.0148	0.645±0.0110	0.662±0.0098
CHIEF	0.219±0.0248	0.287±0.0301	0.362±0.0256	0.423±0.0207	0.476±0.0177	0.510±0.0152
TITAN _v	<u>0.508±0.0292</u>	<u>0.580±0.0257</u>	<u>0.622±0.0198</u>	<u>0.647±0.0210</u>	<u>0.670±0.0164</u>	<u>0.681±0.0191</u>
TITAN	0.566±0.0333	0.629±0.0239	0.667±0.0196	0.687±0.0207	0.701±0.0191	0.702±0.0164

Extended Data Table 73: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for OncoTree code ($C = 2$) prediction on TCGA-OT. Results are given in balanced accuracy of linear probing evaluation or ABMIL training, respectively. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.201±0.0144	0.276±0.0123	0.338±0.0110	0.392±0.0094	0.430±0.0087	0.462±0.0065
ABMIL (CONCH)	0.236±0.0186	0.352±0.0161	0.437±0.0104	0.496±0.0086	0.535±0.0098	0.553±0.0065
GigaPath	0.136±0.0110	0.202±0.0121	0.276±0.0102	0.353±0.0106	0.416±0.0089	0.452±0.0054
PRISM	0.271±0.0161	0.346±0.0137	0.409±0.0096	0.457±0.0099	0.489±0.0086	0.501±0.0058
CHIEF	0.126±0.0111	0.179±0.0108	0.242±0.0089	0.311±0.0109	0.368±0.0072	0.405±0.0055
TITAN _v	<u>0.307±0.0160</u>	<u>0.394±0.0132</u>	<u>0.462±0.0097</u>	<u>0.512±0.0075</u>	<u>0.543±0.0079</u>	<u>0.554±0.0063</u>
TITAN	0.364±0.0164	0.449±0.0135	0.508±0.0098	0.550±0.0099	0.572±0.0081	0.578±0.0061

Extended Data Table 74: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for OncoTree code ($C = 2$) prediction on OT108. Results are given in balanced accuracy of linear probing evaluation or ABMIL training, respectively. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.354±0.0267	0.445±0.0310	0.526±0.0212	0.589±0.0147	0.638±0.0176	0.668±0.0127
ABMIL (CONCH)	0.388±0.0417	0.544±0.0372	<u>0.629±0.0230</u>	0.671±0.0195	0.702±0.0138	0.729±0.0118
GigaPath	0.292±0.0284	0.398±0.0250	0.507±0.0206	0.589±0.0190	0.658±0.0175	0.709±0.0122
PRISM	0.399±0.0328	0.484±0.0192	0.548±0.0225	0.595±0.0193	0.628±0.0173	0.659±0.0156
CHIEF	0.226±0.0280	0.319±0.0278	0.418±0.0237	0.509±0.0171	0.578±0.0158	0.637±0.0122
TITAN _v	<u>0.489±0.0261</u>	<u>0.568±0.0208</u>	0.627±0.0195	<u>0.672±0.0184</u>	<u>0.708±0.0154</u>	<u>0.730±0.0114</u>
TITAN	0.535±0.0303	0.609±0.0255	0.660±0.0206	0.695±0.0149	0.724±0.0146	0.733±0.0127

Extended Data Table 75: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for tumor type ($C = 2$) prediction on EBRAINS. Results are given in balanced accuracy of linear probing evaluation or ABMIL training, respectively. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.505±0.0387	0.597±0.0284	0.666±0.0173	0.711±0.0186	0.735±0.0116	0.748±0.0126
GigaPath	0.292±0.0297	0.366±0.0289	0.441±0.0205	0.505±0.0187	0.543±0.0145	0.562±0.0152
PRISM	0.478±0.0403	0.569±0.0276	0.641±0.0191	0.677±0.0137	0.694±0.0101	0.702±0.0091
CHIEF	0.262±0.0295	0.323±0.0286	0.391±0.0220	0.440±0.0162	0.468±0.0148	0.482±0.0121
TITAN _v	<u>0.639±0.0385</u>	<u>0.727±0.0243</u>	<u>0.776±0.0155</u>	0.798±0.0122	<u>0.809±0.0094</u>	<u>0.816±0.0100</u>
TITAN	0.684±0.0366	0.767±0.0232	0.808±0.0133	0.824±0.0113	0.832±0.0109	0.838±0.0097

Extended Data Table 76: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for tumor subtype ($C = 2$) prediction on TCGA-UT-8K. Results are given in balanced accuracy of SimpleShot evaluation. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.397±0.0312	0.466±0.0303	0.525±0.0223	0.565±0.0189	0.590±0.0152	0.603±0.0105
GigaPath	0.217±0.0285	0.276±0.0296	0.343±0.0206	0.391±0.0168	0.432±0.0203	0.455±0.0154
PRISM	0.455±0.0288	0.528±0.0224	0.589±0.0148	0.621±0.0112	0.637±0.0084	0.649±0.0077
CHIEF	0.223±0.0246	0.275±0.0290	0.329±0.0245	0.369±0.0168	0.401±0.0197	0.414±0.0118
TITAN _v	<u>0.517±0.0296</u>	<u>0.587±0.0233</u>	<u>0.633±0.0185</u>	<u>0.664±0.0164</u>	<u>0.685±0.0135</u>	<u>0.696±0.0106</u>
TITAN	0.568±0.0316	0.639±0.0215	0.685±0.0201	0.713±0.0132	0.730±0.0125	0.745±0.0074

Extended Data Table 77: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for OncoTree code ($C = 2$) prediction on TCGA-OT. Results are given in balanced accuracy of SimpleShot evaluation. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.203±0.0158	0.261±0.0139	0.315±0.0124	0.361±0.0091	0.392±0.0080	0.406±0.0043
GigaPath	0.139±0.0123	0.189±0.0108	0.240±0.0107	0.293±0.0100	0.330±0.0081	0.348±0.0054
PRISM	0.278±0.0165	0.347±0.0121	0.412±0.0110	0.455±0.0089	0.478±0.0062	0.486±0.0042
CHIEF	0.124±0.0102	0.162±0.0093	0.208±0.0097	0.249±0.0086	0.279±0.0062	0.293±0.0052
TITAN _v	<u>0.299±0.0148</u>	<u>0.380±0.0136</u>	<u>0.445±0.0104</u>	<u>0.491±0.0085</u>	<u>0.516±0.0065</u>	<u>0.525±0.0054</u>
TITAN	0.370±0.0165	0.452±0.0133	0.514±0.0112	0.552±0.0082	0.573±0.0058	0.581±0.0035

Extended Data Table 78: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for OncoTree code ($C = 2$) prediction on OT108. Results are given in balanced accuracy of SimpleShot evaluation. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	$k=1$	$k=2$	$k=4$	$k=8$	$k=16$	$k=32$
Mean pool (CONCH)	0.358±0.0248	0.417±0.0261	0.480±0.0233	0.519±0.0229	0.540±0.0151	0.548±0.0110
GigaPath	0.304±0.0276	0.388±0.0220	0.483±0.0203	0.544±0.0172	0.583±0.0165	0.610±0.0104
PRISM	0.425±0.0326	0.503±0.0229	0.549±0.0187	0.589±0.0187	0.609±0.0146	0.628±0.0099
CHIEF	0.221±0.0248	0.287±0.0268	0.355±0.0243	0.410±0.0221	0.445±0.0149	0.473±0.0131
TITAN _v	<u>0.502±0.0266</u>	<u>0.572±0.0242</u>	<u>0.633±0.0183</u>	<u>0.669±0.0163</u>	<u>0.703±0.0132</u>	<u>0.723±0.0080</u>
TITAN	0.549±0.0318	0.615±0.0193	0.669±0.0179	0.709±0.0139	0.731±0.0112	0.749±0.0082

Extended Data Table 79: Few shot experiments with $k \in \{1, 2, 4, 8, 16, 32\}$ shots per class for tumor type ($C = 2$) prediction on EBRAINS. Results are given in balanced accuracy of SimpleShot evaluation. If a class contains less than k samples, all samples are used. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.536±0.0061	0.599±0.0055	0.966±0.0010
TITAN	0.761 ±0.0063	0.798 ±0.0043	0.989 ±0.0006

Extended Data Table 80: Zero-shot results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.460±0.0185	0.538±0.0143	0.968±0.0023
TITAN	0.616 ±0.0178	0.713 ±0.0127	0.985 ±0.0014

Extended Data Table 81: Zero-shot results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.331±0.0100	0.296±0.0116	0.959±0.0021
TITAN	0.477 ±0.0108	0.452 ±0.0133	0.980 ±0.0015

Extended Data Table 82: Zero-shot results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.279±0.0161	0.263±0.0189	0.899±0.0054
TITAN	0.543 ±0.0201	0.365 ±0.0210	0.968 ±0.0025

Extended Data Table 83: Zero-shot results for tumor type ($C = 30$) prediction on eBrains. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.517±0.0507	0.390±0.0460	0.909 ±0.0167
TITAN	0.585 ±0.0199	<u>0.489</u> ±0.0466	<u>0.859</u> ±0.0238

Extended Data Table 84: Zero-shot results for histological pattern ($C = 5$) prediction on DHMC-LUAD. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.565±0.0124	0.724±0.0234	0.886±0.0225
TITAN	0.940 ±0.0220	0.917 ±0.0368	0.985 ±0.0154

Extended Data Table 85: Zero-shot results for OncoTree code ($C = 3$) prediction on TCGA-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.435±0.0155	<u>0.862</u> ±0.0136	0.839±0.0173
TITAN	0.966 ±0.0064	0.940 ±0.0074	0.993 ±0.0018

Extended Data Table 86: Zero-shot results for OncoTree code ($C = 3$) prediction on CPTAC-DHMC-RCC. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	<u>0.894</u> ±0.0135	<u>0.895</u> ±0.0213	0.893±0.0552
TITAN	0.895 ±0.0254	0.903 ±0.0216	0.952 ±0.0167

Extended Data Table 87: Zero-shot results for tumor subtype ($C = 2$) prediction on TCGA-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.918 ±0.0192	0.918 ±0.0194	0.848±0.0270
TITAN	<u>0.907</u> ±0.0196	<u>0.906</u> ±0.0196	0.985 ±0.0039

Extended Data Table 88: Zero-shot results for tumor subtype ($C = 2$) prediction on TCGA-NSCLC. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.938 ±0.0069	0.936 ±0.0072	0.982 ±0.0037
TITAN	<u>0.919</u> ±0.0078	<u>0.916</u> ±0.0083	<u>0.981</u> ±0.0044

Extended Data Table 89: Zero-shot results for tumor subtype ($C = 2$) prediction on CPTAC-NSCLC. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	<u>0.584</u> ±0.0179	<u>0.573</u> ±0.0163	<u>0.576</u> ±0.0210
TITAN	0.693 ±0.0148	0.633 ±0.0166	0.762 ±0.0159

Extended Data Table 90: Zero-shot results for antibody-mediated rejection ($C = 2$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	<u>0.562</u> ±0.0144	<u>0.610</u> ±0.0194	<u>0.627</u> ±0.0193
TITAN	0.737 ±0.0159	0.732 ±0.0153	0.784 ±0.0154

Extended Data Table 91: Zero-shot results for cellular rejection ($C = 2$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Bal. acc.	Weighted F1	AUROC
PRISM	0.612 ±0.0127	0.569 ±0.0183	<u>0.699</u> ±0.0165
TITAN	<u>0.603</u> ±0.0103	<u>0.528</u> ±0.0193	0.801 ±0.0136

Extended Data Table 92: Zero-shot results for cellular rejection ($C = 2$) prediction on CRANE. The best result is marked in bold and the second best is underlined.

Encoder	S1	S2	S3	Logistic regression		SimpleShot		20-nearest neighbors	
				Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v	✓			0.820±0.0055	0.870±0.0037	0.818±0.0053	0.864±0.0039	0.812±0.0050	0.842±0.0038
TITAN	✓	✓		<u>0.827</u> ±0.0056	<u>0.875</u> ±0.0038	0.817±0.0057	0.866±0.0039	0.822±0.0060	0.839±0.0041
TITAN	✓		✓	0.832 ±0.0055	0.881 ±0.0036	0.821±0.0057	0.864±0.0039	0.829±0.0053	<u>0.857</u> ±0.0039
TITAN _L		✓	✓	<u>0.827</u> ±0.0053	0.873±0.0038	<u>0.822</u> ±0.0053	<u>0.867</u> ±0.0038	<u>0.834</u> ±0.0057	0.848±0.0040
MH-ABMIL _L		✓	✓	0.815±0.0054	0.862±0.0038	0.801±0.0056	0.850±0.0040	0.823±0.0059	0.833±0.0040
TITAN	✓	✓	✓	0.832 ±0.0056	0.881 ±0.0036	0.828 ±0.0052	0.875 ±0.0037	0.843 ±0.0056	0.865 ±0.0037

Extended Data Table 93: Ablation study of vision-language design. Linear probing results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	S1	S2	S3	Logistic regression		SimpleShot		20-nearest neighbors	
				Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v	✓			0.690±0.0196	0.758±0.0123	0.639±0.0149	0.744±0.0124	0.706±0.0183	0.775±0.0116
TITAN	✓	✓		0.651±0.0163	0.733±0.0127	0.646±0.0162	0.716±0.0130	0.674±0.0190	0.723±0.0123
TITAN	✓		✓	<u>0.693</u> ±0.0161	0.752±0.0119	0.666±0.0177	0.756±0.0120	0.755 ±0.0159	<u>0.801</u> ±0.0107
TITAN _L		✓	✓	0.685±0.0156	<u>0.761</u> ±0.0120	0.699 ±0.0159	0.775 ±0.0121	0.731±0.0168	0.799±0.0111
MH-ABMIL _L		✓	✓	0.663±0.0160	0.760±0.0123	0.636±0.0176	0.740±0.0125	0.719±0.0173	0.773±0.0116
TITAN	✓	✓	✓	0.704 ±0.0192	0.764 ±0.0116	<u>0.685</u> ±0.0183	<u>0.774</u> ±0.0119	<u>0.744</u> ±0.0152	0.802 ±0.0108

Extended Data Table 94: Ablation study of vision-language design. Linear probing results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	S1	S2	S3	Logistic regression		SimpleShot		20-nearest neighbors	
				Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v	✓			0.558±0.0104	0.536±0.0127	0.464±0.0105	0.430±0.0128	0.524±0.0106	0.517±0.0129
TITAN	✓	✓		0.539±0.0108	0.519±0.0131	0.452±0.0101	0.413±0.0130	0.508±0.0109	0.501±0.0130
TITAN	✓		✓	<u>0.564</u> ±0.0103	<u>0.540</u> ±0.0128	0.498±0.0103	0.458±0.0134	0.553±0.0108	0.542±0.0130
TITAN _L		✓	✓	0.559±0.0107	0.538±0.0130	<u>0.515</u> ±0.0101	<u>0.475</u> ±0.0135	<u>0.569</u> ±0.0108	<u>0.557</u> ±0.0135
MH-ABMIL _L		✓	✓	0.557±0.0109	0.535±0.0130	0.491±0.0102	0.452±0.0131	0.550±0.0109	0.540±0.0136
TITAN	✓	✓	✓	0.587 ±0.0103	0.563 ±0.0130	0.527 ±0.0101	0.489 ±0.0132	0.580 ±0.0103	0.567 ±0.0129

Extended Data Table 95: Ablation study of vision-language design. Linear probing results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	S1	S2	S3	Logistic regression		SimpleShot		20-nearest neighbors	
				Bal. acc.	Weighted F1	Bal. acc.	Weighted F1	Bal. acc.	Weighted F1
TITAN _v	✓			0.732±0.0208	0.785±0.0175	0.656±0.0187	0.709±0.0199	<u>0.742</u> ±0.0192	<u>0.727</u> ±0.0182
TITAN	✓	✓		0.707±0.0212	0.766±0.0187	0.596±0.0211	0.664±0.0213	0.654±0.0216	0.634±0.0203
TITAN	✓		✓	<u>0.734</u> ±0.0212	0.791 ±0.0169	<u>0.685</u> ±0.0202	<u>0.739</u> ±0.0197	0.725±0.0208	0.721±0.0191
TITAN _L		✓	✓	0.725±0.0206	0.771±0.0178	0.682±0.0196	0.733±0.0198	0.712±0.0214	0.704±0.0192
MH-ABMIL _L		✓	✓	0.724±0.0204	0.779±0.0189	0.642±0.0204	0.706±0.0210	0.711±0.0217	0.707±0.0189
TITAN	✓	✓	✓	0.735 ±0.0204	<u>0.786</u> ±0.0182	0.695 ±0.0196	0.746 ±0.0200	0.754 ±0.0203	0.748 ±0.0182

Extended Data Table 96: Ablation study of vision-language design. Linear probing results for tumor type ($C = 30$) prediction on EBRAINS. The best result is marked in bold and the second best is underlined.

Hyperparameter	Value
Layers	6
Heads	12
Head activation	GELU
Embedding dimension	768
Drop path rate	0.1
Global crop size	14
Global crop number	2
Local crop size	6
Local crop number	10
Partial prediction shape	Block
Partial prediction ratio	0.3
Partial prediction variance	0.2
Gradient clipping max norm	3.0
Normalize last layer	False
Shared head	True
AdamW β	(0.9, 0.999)
Batch size	1024
Freeze last layer epochs	3
Warmup epochs	30
Warmup teacher temperature epochs	90
Max epochs	300
Learning rate schedule	Cosine
Learning rate (start)	0
Learning rate (post warmup)	5e-4
Learning rate (final)	1e-6
Teacher temperature (start)	0.04
Teacher temperature (final)	0.07
Teacher momentum (start)	0.996
Teacher momentum (final)	1.000
Weight decay (start)	0.04
Weight decay (end)	0.4
Automatic mixed precision	bf16

Extended Data Table 97: Hyperparameters used in pretraining the unimodal vision model. $4 \times 80\text{GB}$ NVIDIA A100 GPUs were used for training. Batch size refers to the total batch size across GPUs.

Hyperparameter	Value
Automatic mixed precision	fp16
Batch size	1568
Gradient accumulation	2
Vision weight decay	1e-6
Other weight decay	5e-5
Vision peak learning rate	5e-6
Other peak learning rate	5e-5
Vision warmup steps	600
Other warmup steps	200
AdamW β	(0.9, 0.999)
Temperature	Learned
Learning rate schedule	Cosine
Epochs	10

Extended Data Table 98: Hyperparameters used in visual-caption pretraining. 8×80 GB NVIDIA A100 GPUs were used for training. Effective batch size used for optimization is batch size \times gradient accumulation steps. The maximum sequence length for captions is set to 128.

Hyperparameter	Value
Automatic mixed precision	fp16
Batch size	128
Gradient accumulation	2
Vision weight decay	1e-5
Other weight decay	5e-5
Vision peak learning rate	5e-5
Other peak learning rate	5e-5
Vision warmup steps	1800
Other warmup steps	200
AdamW β	(0.9, 0.999)
Temperature	Learned
Learning rate schedule	Cosine
Epochs	5

Extended Data Table 99: Hyperparameters used in visual-report pretraining. 8×80 GB NVIDIA A100 GPUs were used for training. Effective batch size used for optimization is batch size \times gradient accumulation steps. The maximum sequence length for captions is set to 128.

Hyperparameter	Value
Image size	448 × 448
Automatic mixed precision	FP16
Batch size	256
Gradient accumulation steps	3
Learning rate scheduler	Cosine
Warmup steps	250
Peak learning rate	1e-4
AdamW β	(0.9, 0.999)
AdamW ϵ	1e-8
Weight decay	0.2
Softmax temperature	Learned
Epochs	20

Extended Data Table 100: Hyperparameters used in CONCH v1.5 training. $8 \times 80\text{GB}$ NVIDIA GPUs were used for training CONCH v1.5. We initialized vision backbone with UNI⁹, while keeping the text tower configuration of CONCH v1 unchanged. Effective batch size used for optimization is batch size \times gradient accumulation steps. The maximum sequence length for captions is set to 128.

CLASSNAME.
 an image of CLASSNAME.
 the image shows CLASSNAME.
 the image displays CLASSNAME.
 the image exhibits CLASSNAME.
 an example of CLASSNAME.
 CLASSNAME is shown.
 this is CLASSNAME.
 I observe CLASSNAME.
 the pathology image shows CLASSNAME.
 a pathology image shows CLASSNAME.
 the pathology slide shows CLASSNAME.
 shows CLASSNAME.
 contains CLASSNAME.
 presence of CLASSNAME.
 CLASSNAME is present.
 CLASSNAME is observed.
 the pathology image reveals CLASSNAME.
 a microscopic image of showing CLASSNAME.
 histology shows CLASSNAME.
 CLASSNAME can be seen.
 the tissue shows CLASSNAME.
 CLASSNAME is identified.

Extended Data Table 101: Prompt templates used for zero-shot slide classification. The name of the class replaces CLASSNAME. See **Tables 102-123** for class prompts of each task.

Task	Class	Class names
TCGA NSCLC	LUAD	adenocarcinoma lung adenocarcinoma adenocarcinoma of the lung LUAD
	LUSC	squamous cell carcinoma lung squamous cell carcinoma squamous cell carcinoma of the lung LUSC
TCGA BRCA	IDC	invasive ductal carcinoma breast invasive ductal carcinoma invasive ductal carcinoma of the breast invasive carcinoma of the breast, ductal pattern breast IDC
	ILC	invasive lobular carcinoma breast invasive lobular carcinoma invasive lobular carcinoma of the breast invasive carcinoma of the breast, lobular pattern breast ILC

Extended Data Table 102: Class prompts for TCGA NSCLC and TCGA BRCA subtyping..

Task	Class	Class names
OT108	ANSC	anal squamous cell carcinoma squamous cell carcinoma of the anal canal anal SCC squamous carcinoma of the anal region anal squamous carcinoma
	CCOV	clear cell ovarian cancer clear cell carcinoma of the ovary clear cell adenocarcinoma of the ovary CCOV
	CHOL	cholangiocarcinoma bile duct carcinoma biliary tract cancer cholangiocellular carcinoma CHOL
	CSCC	cutaneous squamous cell carcinoma skin squamous cell carcinoma squamous cell carcinoma of the skin CSCC
	DLBCLNOS	diffuse large B-cell lymphoma DLBCL diffuse large B-cell lymphoma, not otherwise specified DLBCLNOS
	EOV	endometrioid ovarian cancer endometrioid carcinoma of the ovary endometrioid tumor EOV
	FL	follicular lymphoma follicular lymphoid neoplasm follicular non-Hodgkin lymphoma FL
	GBC	gallbladder cancer gallbladder adenocarcinoma gallbladder carcinoma GBC
	GEJ	adenocarcinoma of the gastroesophageal junction gastroesophageal junction adenocarcinoma GEJ adenocarcinoma gastroesophageal junction cancer adenocarcinoma at the gastroesophageal junction
GINET	gastrointestinal neuroendocrine tumors GI neuroendocrine tumors gastrointestinal carcinoid tumors gastrointestinal neuroendocrine neoplasms GINET	

Extended Data Table 103: Class prompts for OT-108 and TCGA-OncoTree subtyping.

Task	Class	Class names
OT108	GIST	gastrointestinal stromal tumor gastric stromal tumor intestinal stromal tumor gastrointestinal mesenchymal tumor GIST
	HEMA	hemangioma vascular tumor hemangioma of soft tissue hemangioma of skin cavernous hemangioma capillary hemangioma HEMA
	HGGNOS	high-grade glioma high-grade glioma, not otherwise specified high-grade glioma NOS HGGNOS
	LGGNOS	low-grade glioma low-grade glioma not otherwise specified low-grade glioma, NOS LGG NOS
	LGSOC	low-grade serous ovarian cancer low-grade serous carcinoma serous papillary ovarian carcinoma, low grade LGSOC
	MCC	Merkel cell carcinoma neuroendocrine carcinoma of the skin cutaneous neuroendocrine carcinoma MCC
	MDLC	mixed ductal and lobular carcinoma breast mixed ductal and lobular carcinoma mixed carcinoma of the breast MDLC
	MNG	meningioma meningeal tumor benign meningioma malignant meningioma MNG
	NBL	neuroblastoma neuroblastic tumor neuroblastic neoplasm NBL
NSCLCPD	poorly differentiated non-small cell lung cancer poorly differentiated NSCLC non-small cell lung cancer (poorly differentiated) NSCLCPD	

Extended Data Table 104: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108	PANET	pancreatic neuroendocrine tumor pancreatic endocrine tumor pancreatic neuroendocrine neoplasm neuroendocrine tumor of the pancreas PANET
	PTAD	pituitary adenoma adenoma of the pituitary gland pituitary gland adenoma PTAD
	SCLC	small cell lung cancer small cell carcinoma of the lung SCLC
	UTUC	urothelial carcinoma upper tract urothelial carcinoma urothelial carcinoma of the upper urinary tract UTUC
	BLAD	bladder adenocarcinoma adenocarcinoma of the bladder bladder cancer, adenocarcinomatous type BLAD
	MBL	medulloblastoma primitive neuroectodermal tumor (PNET) medulloblastoma variant MBL
	SBC	duodenal adenocarcinoma adenocarcinoma of the duodenum duodenum cancer SBC
	OS	osteosarcoma osteogenic sarcoma primary osteosarcoma OS
	ULMS	leiomyosarcoma uterine leiomyosarcoma leiomyosarcoma of the uterus ULMS
AMPCA	ampullary carcinoma carcinoma of the ampulla of Vater ampullary adenocarcinoma AMPCA	

Extended Data Table 105: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108	MGCT	mixed germ cell tumor mixed germ cell neoplasm germ cell tumor, mixed type MGCT
	WT	Wilms' tumor nephroblastoma renal embryonal tumor WT
	OCS	ovarian carcinosarcoma malignant mixed mesodermal tumor mixed mesodermal tumor of the ovary ovarian mixed tumor OCS
	AODG	anaplastic oligodendroglioma oligodendroglioma anaplastic oligodendroglioma WHO grade III AODG
	ES	Ewing sarcoma Ewing's sarcoma primitive neuroectodermal tumor (PNET) ES
	ALUCA	atypical lung carcinoid atypical carcinoid tumor lung carcinoid tumor (atypical) atypical neuroendocrine tumor of the lung ALUCA
	LUCA	lung carcinoid carcinoid tumor of the lung pulmonary carcinoid LUCA
	EPM	ependymoma ependymal tumor central nervous system ependymoma EPM
	LUAS	adenosquamous carcinoma lung adenosquamous carcinoma adenosquamous cell carcinoma of the lung LUAS
UCCC	clear cell carcinoma uterine clear cell carcinoma clear cell carcinoma of the uterus UCCC	

Extended Data Table 106: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108	MXOV	mixed ovarian carcinoma mixed epithelial ovarian carcinoma ovarian mixed carcinoma MXOV
	PLBMESO	pleural mesothelioma biphasic mesothelioma biphasic pleural mesothelioma pleural biphasic mesothelioma PLBMESO
	DES	desmoid tumor aggressive fibromatosis desmoid fibromatosis extra-abdominal desmoid tumor desmoid-type fibromatosis DES
	SYNS	synovial sarcoma synovial sarcoma with biphasic morphology monophasic synovial sarcoma synovial sarcoma, high grade SYNS
	SFT	solitary fibrous tumor hemangiopericytoma solitary fibrous tumor of the pleura hemangiopericytoma of the soft tissue SFT
	MAAP	mucinous adenocarcinoma of the appendix appendiceal mucinous adenocarcinoma mucinous carcinoma of the appendix MAAP
	ECAD	endocervical adenocarcinoma adenocarcinoma of the cervix cervical adenocarcinoma ECAD
	MYCF	mycosis fungoides cutaneous T-cell lymphoma mycosis fungoides/Sezary syndrome MF
	CHS	chondrosarcoma primary chondrosarcoma secondary chondrosarcoma dedifferentiated chondrosarcoma chondrosarcoma of bone CHS
ANGS	angiosarcoma vascular sarcoma hemangiosarcoma angiosarcoma of soft tissue ANGS	

Extended Data Table 107: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108	PAST	pilocytic astrocytoma juvenile pilocytic astrocytoma cystic astrocytoma PAST
	MOV	mucinous ovarian cancer mucinous cystadenocarcinoma mucinous neoplasm of the ovary mucinous ovarian carcinoma MOV
	ACYC	adenoid cystic carcinoma cylindroma adenoid cystic adenocarcinoma ACC
	SCHW	schwannoma neurilemmoma peripheral nerve sheath tumor schwannoma tumor SCHW
	UMEC	mixed endometrial carcinoma uterine mixed endometrial carcinoma mixed carcinoma of the uterine corpus endometrial carcinoma, mixed subtype UMEC
	MPNST	malignant peripheral nerve sheath tumor malignant schwannoma peripheral nerve sheath tumor MPNST
	SBOV	serous borderline ovarian tumor serous borderline tumor borderline serous ovarian neoplasm SBOV
	LUNE	large cell neuroendocrine carcinoma large cell neuroendocrine tumor neuroendocrine carcinoma, large cell type large cell carcinoma with neuroendocrine features LCNEC
	CLLSLL	chronic lymphocytic leukemia small lymphocytic lymphoma chronic lymphocytic leukemia/small lymphocytic lymphoma CLL / SLL
GRCT	granulosa cell tumor granulosa cell neoplasm granulosa tumor granulosa cell tumor of the ovary GRCT	

Extended Data Table 108: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108	THME	medullary thyroid cancer medullary thyroid carcinoma thyroid medullary carcinoma medullary carcinoma of the thyroid MTC
	VSC	squamous cell carcinoma of the vulva squamous cell carcinoma of the vagina vulvar squamous cell carcinoma vaginal squamous cell carcinoma VSC
	THAP	anaplastic thyroid carcinoma anaplastic thyroid cancer undifferentiated thyroid carcinoma THAP
	PGNG	paraganglioma extra-adrenal paraganglioma head and neck paraganglioma carotid body tumor glomus tumor PGNG
	SSRCC	signet ring cell carcinoma signet ring cell carcinoma of the stomach gastric signet ring cell carcinoma signet ring cell gastric cancer SSRCC
	PECOMA	perivascular epithelioid cell tumor PEComa epithelioid smooth muscle tumor angiomyolipoma clear cell sugar tumor
	WDLS	well-differentiated liposarcoma differentiated liposarcoma liposarcoma, well-differentiated WDLS
	ATM	atypical meningioma meningioma WHO grade II atypical meningeal tumor atypical tumor of the meninges ATM
	ROCY	oncocytoma renal oncocytoma oncocytic tumor of the kidney kidney oncocytoma ROCY
	ERMS	embryonal rhabdomyosarcoma rhabdomyosarcoma, embryonal type embryonal rhabdomyoblastoma ERMS

Extended Data Table 109: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108 & TCGA-OT	ACC	adrenocortical carcinoma adrenal cortical carcinoma adrenal cortex carcinoma ACC
	ASTR	astrocytoma diffuse astrocytoma fibrillary astrocytoma anaplastic astrocytoma grade II astrocytoma grade III astrocytoma ASTR
	BLCA	urothelial carcinoma bladder urothelial carcinoma transitional cell carcinoma bladder cancer BLCA
	CCRCC	clear cell carcinoma renal clear cell carcinoma clear cell renal cell carcinoma clear cell RCC CCRCC
	CHRCC	chromophobe renal cell carcinoma chromophobe RCC chromophobe carcinoma renal cell carcinoma, chromophobe type CHRCC
	COAD	colon adenocarcinoma adenocarcinoma of the colon colorectal adenocarcinoma COAD
	DDL5	dedifferentiated liposarcoma dedifferentiated liposarcoma variant liposarcoma, dedifferentiated DDL5
	ESCA	esophageal adenocarcinoma adenocarcinoma of the esophagus esophageal cancer, adenocarcinoma type ESCA
	ESCC	esophageal squamous cell carcinoma esophageal SCC squamous cell carcinoma of the esophagus SCC of the esophagus ESCC
	GBM	glioblastoma multiforme glioblastoma GBM multiforme glioblastoma grade IV astrocytoma

Extended Data Table 110: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108 & TCGA-OT	HCC	hepatocellular carcinoma liver cancer hepatoma HCC
	HGSOC	high-grade serous ovarian cancer high-grade serous carcinoma HGSOC serous papillary carcinoma serous ovarian carcinoma
	HNSC	head and neck squamous cell carcinoma HNSCC head and neck SCC oropharyngeal squamous cell carcinoma laryngeal squamous cell carcinoma hypopharyngeal squamous cell carcinoma nasopharyngeal squamous cell carcinoma oral cavity squamous cell carcinoma HNSC
	IDC	invasive ductal carcinoma breast invasive ductal carcinoma ductal carcinoma, no special type IDC
	ILC	invasive lobular carcinoma breast invasive lobular carcinoma invasive lobular carcinoma of the breast lobular carcinoma breast ILC
	LMS	leiomyosarcoma smooth muscle sarcoma leiomyosarcoma of the soft tissue LMS
	LUAD	lung adenocarcinoma adenocarcinoma of the lung pulmonary adenocarcinoma peripheral lung adenocarcinoma LUAD
	LUSC	squamous cell carcinoma lung squamous cell carcinoma squamous carcinoma of the lung LUSC
	MEL	acral melanoma melanoma acral lentiginous melanoma MEL
	MFH	undifferentiated pleomorphic sarcoma malignant fibrous histiocytoma high-grade spindle cell sarcoma MFH

Extended Data Table 111: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108 & TCGA-OT	PAAD	pancreatic adenocarcinoma adenocarcinoma of the pancreas pancreas adenocarcinoma PAAD
	PLEMESO	pleural mesothelioma epithelioid mesothelioma epithelioid pleural mesothelioma pleural mesothelioma, epithelioid type PLEMESO
	PRAD	prostate adenocarcinoma adenocarcinoma of the prostate prostatic adenocarcinoma PRAD
	PRCC	papillary renal cell carcinoma renal cell carcinoma, papillary type papillary RCC PRCC
	READ	rectal adenocarcinoma adenocarcinoma of the rectum rectal cancer READ
	STAD	diffuse type stomach adenocarcinoma diffuse gastric adenocarcinoma diffuse-type gastric cancer gastric adenocarcinoma, diffuse type STAD
	THPA	papillary thyroid carcinoma papillary thyroid cancer papillary thyroid neoplasm thyroid papillary carcinoma PTC
	UCS	uterine carcinosarcoma uterine malignant mixed Müllerian tumor malignant mixed Müllerian tumor carcinosarcoma of the uterus UCS
	UEC	endometrioid carcinoma uterine endometrioid carcinoma endometrial endometrioid carcinoma endometrial carcinoma, endometrioid type UEC
	USC	uterine serous carcinoma uterine papillary serous carcinoma serous papillary carcinoma of the uterus serous carcinoma USC

Extended Data Table 112: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
OT108 & TCGA-OT	THFO	follicular thyroid carcinoma follicular thyroid cancer follicular carcinoma of the thyroid thyroid follicular cancer THFO
	AASTR	anaplastic astrocytoma astrocytoma, anaplastic grade III astrocytoma AASTR
	CESC	squamous cell carcinoma of the cervix cervical squamous cell carcinoma cervical SCC CESC
	ODG	oligodendroglioma oligodendroglial tumor oligodendroglioma grade II oligodendroglioma grade III ODG
	MACR	mucinous adenocarcinoma mucinous adenocarcinoma of the colon mucinous adenocarcinoma of the rectum mucinous colorectal adenocarcinoma MACR
	SEM	seminoma testicular seminoma pure seminoma classic seminoma seminomatous germ cell tumor SEM
	MFS	myxofibrosarcoma myxoid fibrosarcoma myxofibrosarcomatosis MFS
TCGA-OT	THYM	thymoma thymic tumor thymic carcinoma thymic epithelial neoplasm THYM
	NSGCT	non-seminomatous germ cell tumor germ cell tumor, non-seminomatous mixed germ cell tumor embryonal carcinoma yolk sac tumor teratoma NSGCT
	OAST	oligoastrocytoma oligoastrocytic tumor oligodendroglial-astrocytic tumor OAST

Extended Data Table 113: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Task	Class	Class names
TCGA-OT	PHC	pheochromocytoma chromaffin cell tumor adrenal pheochromocytoma paranglioma PHC
	DSTAD	diffuse type adenocarcinoma of the stomach diffuse stomach adenocarcinoma diffuse gastric adenocarcinoma stomach adenocarcinoma, diffuse type DSTAD
	AOAST	anaplastic oligoastrocytoma oligoastrocytoma anaplastic mixed glioma AOAST
	UM	uveal melanoma uvea melanoma ocular melanoma choroidal melanoma ciliary body melanoma iris melanoma UM
	SKCM	cutaneous melanoma skin melanoma melanoma of the skin malignant melanoma SKCM
	TSTAD	tubular stomach adenocarcinoma tubular adenocarcinoma of the stomach stomach tubular adenocarcinoma gastric tubular adenocarcinoma TSTAD

Extended Data Table 114: Class prompts for OT-108 and TCGA-OncoTree subtyping. Continued.

Class	Class names
Adrenocortical carcinoma	adrenocortical carcinoma adrenal cortex cancer adrenal cortical carcinoma ACC
Bladder Urothelial Carcinoma	bladder urothelial carcinoma bladder cancer urothelial carcinoma of the bladder BUC
Brain Lower Grade Glioma	brain lower grade glioma low grade glioma brain glioma with low grade BLGG
Breast invasive carcinoma	breast invasive carcinoma breast cancer invasive carcinoma of the breast BIC
Cervical squamous cell carcinoma and endocervical adenocarcinoma	cervical squamous cell carcinoma cervical cancer endocervical adenocarcinoma cervical carcinoma
Cholangiocarcinoma	cholangiocarcinoma bile duct cancer cholangiocellular carcinoma CCA
Colon adenocarcinoma	colon adenocarcinoma colon cancer adenocarcinoma in the colon COAD
Esophageal carcinoma	esophageal carcinoma esophageal cancer esophageal squamous cell carcinoma ESCA
Glioblastoma multiforme	glioblastoma multiforme GBM brain tumor glioblastoma
Head and Neck squamous cell carcinoma	head and neck squamous cell carcinoma head and neck cancer HNSCC

Extended Data Table 115: Class prompts for TCGA Uniform Tumor subtyping.

Class	Class names
Kidney Chromophobe	kidney chromophobe chromophobe renal cell carcinoma KICH
Kidney renal clear cell carcinoma	kidney renal clear cell carcinoma clear cell renal cell carcinoma KIRC CCRCC
Kidney renal papillary cell carcinoma	kidney renal papillary cell carcinoma papillary renal cell carcinoma KIRP PRCC
Liver hepatocellular carcinoma	liver hepatocellular carcinoma hepatocellular carcinoma hepatocellular carcinoma in the liver LIHC HCC
Lung adenocarcinoma	lung adenocarcinoma adenocarcinoma of the lung LUAD
Lung squamous cell carcinoma	lung squamous cell carcinoma squamous cell carcinoma of the lung LUSC
Lymphoid Neoplasm Diffuse Large B-cell Lymphoma	lymphoid neoplasm diffuse large B-cell lymphoma diffuse large B-cell lymphoma DLBCL DLBC
Mesothelioma	mesothelioma cancer in the mesothelium MESO
Ovarian serous cystadenocarcinoma	ovarian serous cystadenocarcinoma serous cystadenocarcinoma OV
Pancreatic adenocarcinoma	pancreatic adenocarcinoma adenocarcinoma in the pancreas pancreatic cancer PAAD PDAC

Extended Data Table 116: Class prompts for TCGA Uniform Tumor subtyping. Continued.

Class	Class names
Pheochromocytoma and Paraganglioma	pheochromocytoma and paraganglioma pheochromocytoma paraganglioma PCPG
Prostate adenocarcinoma	prostate adenocarcinoma prostate cancer adenocarcinoma in the prostate PRAD
Rectum adenocarcinoma	rectum adenocarcinoma rectal cancer rectal adenocarcinoma READ
Sarcoma	sarcoma soft tissue sarcoma bone sarcoma SARC
Skin Cutaneous Melanoma	skin cutaneous melanoma cutaneous melanoma skin cancer SKCM
Stomach adenocarcinoma	stomach adenocarcinoma gastric adenocarcinoma gastric cancer STAD
Testicular Germ Cell Tumors	testicular germ cell tumors testicular cancer germ cell tumors TGCT
Thymoma	thymoma thymic carcinoma thymus cancer THYM
Thyroid carcinoma	thyroid carcinoma thyroid cancer carcinoma of the thyroid gland THCA
Uterine Carcinosarcoma	uterine carcinosarcoma carcinosarcoma UCS

Extended Data Table 117: Class prompts for TCGA Uniform Tumor subtyping. Continued.

Class	Class names
Uterine Corpus Endometrial Carcinoma	uterine corpus endometrial carcinoma endometrial carcinoma UCEC
Uveal Melanoma	uveal melanoma eye melanoma ocular melanoma UVM

Extended Data Table 118: Class prompts for TCGA UniformTumor subtyping. Continued.

Class	Class names
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

Extended Data Table 119: Class prompts for EBRAINS subtyping.

Class	Class names
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

Extended Data Table 120: Class prompts for EBRAINS subtyping. Continued.

Class	Class names
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

Extended Data Table 121: Class prompts for EBRAINS subtyping. Continued.

Task	Class	Class names
Antibody-mediated Rejection (AMR)	AMR	antibody-mediated rejection AMR in kidney transplant acute antibody-mediated rejection chronic antibody-mediated rejection graft rejection caused by donor-specific antibodies complement-mediated kidney allograft rejection kidney transplant rejection with microvascular inflammation C4d-positive antibody-mediated rejection
	no AMR	no antibody-mediated rejection absence of AMR in kidney transplant no evidence of donor-specific antibody-mediated rejection absence of C4d deposition in kidney allograft no signs of complement activation or C4d deposition
T-Cell Mediated Rejection (TCMR)	TCMR	T cell-mediated rejection TCMR in kidney transplant acute T cell-mediated rejection chronic T cell-mediated rejection kidney transplant rejection caused by T cells cellular-mediated rejection inflammatory rejection involving T cells kidney allograft rejection with tubulitis kidney transplant rejection with interstitial inflammation
	no TCMR	no T cell-mediated rejection absence of TCMR in kidney transplant kidney transplant without T cell-mediated rejection no evidence of cellular-mediated rejection absence of tubulitis or interstitial inflammation absence of T cell involvement in kidney transplant rejection

Extended Data Table 122: Class prompts for renal allograft rejection..

Task	Class	Class names
Cellular-mediated Rejection (CMR)	CMR	antibody-mediated rejection CMR in heart transplant acute cellular-mediated rejection cellular rejection detected in endomyocardial biopsy heart transplant rejection involving T cells lymphocytic infiltration in the myocardium interstitial inflammation in the heart tissue rejection driven by cellular immune response
	no CMR	no cellular-mediated rejection absence of CMR in heart transplant no evidence of cellular rejection in endomyocardial biopsy heart tissue with no signs of cellular-mediated rejection absence of interstitial inflammation or lymphocytic infiltration normal endomyocardial biopsy findings without rejection no cellular immune response in the myocardium

Extended Data Table 123: Class prompts for cellular-mediated heart allograft rejection..

Example pathchat prompt	<p>Provided are H&E-stained images from a {tissue site} tissue region. The first image shows the microview of the region for you to have an overview of what this region looks like, and the following images are the representative features in this region at higher resolution for you to check the morphologies more clearly. Your task is to write a short morphological caption that describes the unique microscopic findings (cellular and tissue characteristics) in this region as a whole.</p> <ol style="list-style-type: none"> 1. Do not describe each image individually. 2. Do not list your findings one by one. You must summarize all observed microscopic findings in a single paragraph and answer holistically. 3. Pretend you are a busy pathologist writing microscopic findings for this tissue region. Make your description as short as possible and only include the most important findings. 4. Format your reply like a pathology report. Title is not needed.
Example synthetic caption	<p>Microscopic examination of the bladder tissue region reveals a complex histological landscape. The region is characterized by a mixture of squamous and glandular components, indicative of adenosquamous carcinoma. The squamous component is evident by the presence of keratinizing cells, while the glandular component is identified by the formation of gland-like structures. Additionally, there are areas of necrosis and inflammation, suggesting an aggressive tumor behavior. The stroma surrounding the tumor is fibrotic, and there is a notable presence of desmoplasia, which is a reactive fibrous tissue response to the tumor. The tumor cells exhibit pleomorphism, with some cells showing prominent nucleoli, indicative of high-grade malignancy.</p>

Extended Data Table 124: PathChat prompts for synthetic caption generation for 8192 × 8192 pixels ROIs at 20× magnification.

System prompt	<p>You are assisting me summarizing detailed pathology descriptions into concise, clear, and holistic reports suitable for a busy pathologist. The summaries should focus on the key diagnostic features and findings, avoiding unnecessary details while maintaining accuracy.</p> <p>Instructions:</p> <ol style="list-style-type: none"> 1. Read the detailed description and identify the key diagnostic features, such as growth patterns, cellular characteristics, presence or absence of malignant features, and any significant findings. 2. Write a summary that captures the essential information in a clear and concise manner. Aim to keep the summary brief and to the point, using incomplete sentences if needed for brevity. It should provide a holistic view of the pathology findings. 3. You can choose the formats from a) concise summarization in one paragraph; b) summarization using keywords and incomplete sentences similar to a busy and lazy pathologist; and c) formatted similar to a pathologic report. 4. Ignore "grade -" in the description.
Example rewritten caption	<p>Microscopic examination of the bladder tissue shows a complex mix of squamous and glandular elements, characteristic of adenosquamous carcinoma. The presence of keratinized cells identifies the squamous component, while gland-like structures signify the glandular aspect. Necrosis and inflammation indicate an invasive tumor. The surrounding stroma is fibrotic, with marked desmoplasia. Tumor cells display pleomorphism, including cells with large nuclei, suggesting high malignancy grade.</p>

Extended Data Table 125: Prompts for synthetic caption rewriting using Qwen-2-7b instruct.

System prompt	<p>You are assisting me extracting slide description from patient-level pathology reports. The slides of the patient are named by block ids such "A". Each block represents a tissue taken from the patient.</p> <p>You'll be given the entire pathology report of a patient which contains descriptions of multiple blocks, and the specific block we are interested in. Your task to extract the block description from the pathology report, summarizing the morphological details and diagnoses.</p> <p>Instructions:</p> <ol style="list-style-type: none"> 1. Locate the block description in the report. 2. Record Metadata: Identify and record the organ and diagnosis for each slide.
Example user prompt	<p>block: A report: A. SPECIMEN LABELED "SUBGALEAL COLLECTION": Fibrous tissue with florid acute and chronic inflammation and Gram-positive bacilli (highlighted by Gram, MSS, and Warthin-Starry Stains; FITE stain is negative for microorganisms). B. SPECIMEN LABELED "RIGHT FRONTAL BONE": Bone with no significant pathologic change. C. SPECIMEN LABELED "RIGHT FRONTAL LESION": Meninges and brain parenchyma with necrosis, florid acute and chronic inflammation and Gram-positive microorganisms (highlighted by Gram, MSS, and Warthin-Starry Stains; FITE stain is negative for microorganisms). NOTE: Concurrent cultures are positive for Cutibacterium acnes. Correlation with final culture results and other laboratory and clinical findings is recommended.</p>
Example LLM response	<p>Fibrous tissue with florid acute and chronic inflammation.</p>

Extended Data Table 126: Prompts for parsing slide description from patient-level reports.

System prompt	<p>You are assisting me cleaning noisy pathologic reports for H&E stained slides. You should remove the following information from the report if they are present: 1. sensitive patient information including name, age, and address; 2. hospital name; 3. Results of immunohistochemical stain / Immuno stain / other stainings and molecular test; 4. patient clinical history; 5. date. Rewrite the report if it's incomplete or has grammatical mistakes.</p> <p>If the report does not have any sensitive information, just reply the original report. Do not reply "The report does not contain any sensitive information".</p> <p>Along with the original noisy report, you'll be provided diagnosis and tissue site for your information. Reply the cleaned report only. Do not include any explanation.</p>
Example user prompt	<p>ANUS BIOPSY: Diagnosis of Extramammary Paget disease supported by immunohistochemical stains showing positivity for CK7, CEA, EMA and negativity for SOX-10 and CK20. Mucicarmine stain highlights intracytoplasmic mucin in the tumor cells. Slides were reviewed by Dr. {Doctor Name} from the section of dermatopathology.</p>
Example LLM response	<p>The biopsy from the anus shows features consistent with Extramammary Paget disease.</p>

Extended Data Table 127: Prompts for removing sensitive information and information that is not directly visible from the slides.

System prompt	<p>You are assisting me enhancing the diversity of pathology reports while maintaining a formal and professional tone. Reply the diversified report only.</p>
Example user prompt	<p>The biopsy from the kidney shows poorly differentiated carcinoma, consistent with renal cell carcinoma, subtype not determined.</p>
Example rewritten report	<p>Kidney biopsy reveals a poorly differentiated carcinoma, suggestive of renal cell carcinoma, with the subtype yet to be identified.</p>

Extended Data Table 128: Prompts for rewriting parsed slide-level reports.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
GigaPath	0.384±0.0102	0.574±0.0086	0.486±0.0050	0.669±0.0108	0.465±0.0065
PRISM	0.449±0.0135	0.634±0.0216	0.540±0.0174	0.717±0.0238	0.519±0.0201
CHIEF	0.333±0.0166	0.523±0.0139	0.446±0.0135	0.618±0.0090	0.430±0.0138
TITAN _v	<u>0.520±0.0184</u>	<u>0.711±0.0097</u>	<u>0.607±0.0143</u>	<u>0.791±0.0094</u>	<u>0.594±0.0090</u>
TITAN	0.539±0.0100	0.725±0.0125	0.631±0.0221	0.804±0.0099	0.620±0.0156

Extended Data Table 129: Slide retrieval results for rare cancer retrieval ($C = 43$) prediction on Rare-Cancer. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
GigaPath	0.420±0.0285	0.618±0.0258	0.507±0.0273	0.713±0.0141	0.481±0.0313
PRISM	0.482±0.0353	0.663±0.0334	0.547±0.0375	0.757±0.0372	0.530±0.0327
CHIEF	0.380±0.0336	0.563±0.0455	0.470±0.0197	0.654±0.0357	0.446±0.0202
TITAN _v	<u>0.549±0.0236</u>	<u>0.743±0.0269</u>	<u>0.614±0.0301</u>	<u>0.814±0.0273</u>	<u>0.601±0.0370</u>
TITAN	0.555±0.0289	0.754±0.0126	0.647±0.0244	0.832±0.0192	0.636±0.0187

Extended Data Table 130: Slide retrieval results for rare cancer retrieval ($C = 29$) prediction on Rare-Cancer-Public. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.776±0.0047	0.866±0.0038	0.807±0.0044	0.893±0.0035	0.804±0.0044
GigaPath	0.605±0.0053	0.728±0.0049	0.645±0.0052	0.784±0.0047	0.646±0.0054
PRISM	0.754±0.0046	0.854±0.0037	0.788±0.0044	0.888±0.0033	0.787±0.0044
CHIEF	0.552±0.0054	0.690±0.0050	0.609±0.0053	0.755±0.0046	0.613±0.0054
TITAN _v	<u>0.838±0.0041</u>	<u>0.898±0.0034</u>	<u>0.857±0.0039</u>	<u>0.920±0.0030</u>	<u>0.864±0.0038</u>
TITAN	0.854±0.0039	0.912±0.0032	0.875±0.0036	0.931±0.0028	0.876±0.0036

Extended Data Table 131: Slide retrieval results for tumor subtype ($C = 32$) prediction on TCGA-UT-8K. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.649±0.0133	0.803±0.0108	0.712±0.0125	0.864±0.0095	0.711±0.0125
GigaPath	0.483±0.0135	0.666±0.0123	0.572±0.0129	0.732±0.0118	0.565±0.0136
PRISM	0.697±0.0126	0.836±0.0098	0.755±0.0116	0.878±0.0086	0.749±0.0116
CHIEF	0.479±0.0136	0.669±0.0124	0.602±0.0132	0.742±0.0118	0.573±0.0134
TITAN _v	<u>0.732±0.0123</u>	<u>0.852±0.0096</u>	<u>0.773±0.0115</u>	<u>0.902±0.0080</u>	<u>0.766±0.0118</u>
TITAN	0.775±0.0112	0.880±0.0086	0.807±0.0103	0.910±0.0078	0.808±0.0106

Extended Data Table 132: Slide retrieval results for OncoTree code ($C = 46$) prediction on TCGA-OT. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.402±0.0123	0.546±0.0126	0.494±0.0128	0.619±0.0123	0.473±0.0127
GigaPath	0.293±0.0115	0.450±0.0124	0.414±0.0121	0.535±0.0122	0.429±0.0119
PRISM	0.452±0.0125	0.636±0.0126	0.547±0.0127	0.708±0.0119	0.519±0.0128
CHIEF	0.264±0.0112	0.442±0.0131	0.400±0.0129	0.519±0.0127	0.386±0.0126
TITAN _v	<u>0.490±0.0124</u>	<u>0.666±0.0122</u>	<u>0.587±0.0125</u>	<u>0.730±0.0113</u>	<u>0.556±0.0130</u>
TITAN	0.537±0.0129	0.707±0.0116	0.621±0.0126	0.768±0.0108	0.606±0.0126

Extended Data Table 133: Slide retrieval results for OncoTree code ($C = 108$) prediction on OT108. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.633±0.0203	0.813±0.0163	0.719±0.0190	0.873±0.0140	0.710±0.0193
GigaPath	0.593±0.0206	0.806±0.0163	0.733±0.0180	0.860±0.0142	0.705±0.0189
PRISM	0.648±0.0196	0.811±0.0160	0.751±0.0177	0.867±0.0143	0.720±0.0188
CHIEF	0.479±0.0207	0.713±0.0181	0.631±0.0193	0.776±0.0169	0.625±0.0197
TITAN _v	<u>0.687±0.0191</u>	<u>0.848±0.0150</u>	<u>0.776±0.0180</u>	<u>0.904±0.0123</u>	<u>0.763±0.0181</u>
TITAN	0.750±0.0185	0.865±0.0147	0.809±0.0168	0.911±0.0123	0.811±0.0167

Extended Data Table 134: Slide retrieval results for tumor type ($C = 30$) prediction on eBrains. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.840±0.0437	0.950±0.0207	0.870±0.0291	0.972±0.0106	0.879±0.0213
GigaPath	0.776±0.0177	0.933±0.0184	0.819±0.0252	0.961±0.0192	0.827±0.0075
PRISM	0.859±0.0322	0.939±0.0091	0.897±0.0195	0.963±0.0116	0.907±0.0227
CHIEF	0.835±0.0212	0.954±0.0053	0.869±0.0286	0.976±0.0063	0.880±0.0158
TITAN _v	<u>0.878±0.0133</u>	<u>0.957±0.0040</u>	<u>0.909±0.0041</u>	0.969±0.0038	<u>0.910±0.0185</u>
TITAN	0.888±0.0267	0.960±0.0181	0.921±0.0249	<u>0.973±0.0169</u>	0.921±0.0278

Extended Data Table 135: Slide retrieval results for tumor subtype ($C = 2$) prediction on TCGA-BRCA. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.821±0.0222	0.945±0.0041	0.844±0.0216	0.977±0.0079	0.862±0.0188
GigaPath	0.626±0.0778	0.878±0.0487	0.680±0.0524	0.942±0.0300	0.682±0.0585
PRISM	0.882±0.0143	<u>0.967±0.0112</u>	0.898±0.0141	<u>0.980±0.0097</u>	0.913±0.0138
CHIEF	0.790±0.0184	0.944±0.0097	0.824±0.0236	0.972±0.0061	0.830±0.0189
TITAN _v	<u>0.890±0.0126</u>	<u>0.967±0.0091</u>	<u>0.912±0.0158</u>	0.982±0.0053	<u>0.919±0.0154</u>
TITAN	0.920±0.0088	0.968±0.0104	0.930±0.0107	0.979±0.0089	0.935±0.0104

Extended Data Table 136: Slide retrieval results for tumor subtype ($C = 2$) prediction on TCGA-NSCLC. The best result is marked in bold and the second best is underlined.

Encoder	Top-1 acc.	Top-3 acc.	MV@3 acc.	Top-5 acc.	MV@5 acc.
Mean pool (CONCHv1.5)	0.672±0.0158	0.884±0.0113	0.689±0.0155	0.945±0.0079	0.699±0.0158
GigaPath	0.608±0.0170	0.857±0.0116	0.630±0.0162	0.938±0.0082	0.636±0.0162
PRISM	0.644±0.0171	<u>0.887</u> ±0.0107	0.666±0.0167	0.942±0.0080	0.691±0.0157
CHIEF	0.621±0.0167	0.848±0.0119	0.646±0.0162	0.943±0.0077	0.632±0.0164
TITAN _v	<u>0.712</u> ±0.0160	0.881±0.0113	<u>0.723</u> ±0.0157	<u>0.949</u> ±0.0076	<u>0.756</u> ±0.0153
TITAN	0.789 ±0.0147	0.919 ±0.0094	0.785 ±0.0146	0.963 ±0.0066	0.791 ±0.0144

Extended Data Table 137: Slide retrieval results for antibody-mediated rejection ($C = 2$) prediction on Renal allograft rejection. The best result is marked in bold and the second best is underlined.

Encoder	Recall@1	Recall@3	Recall@5	Recall@10	Mean recall
PRISM	<u>0.674</u> ±0.0048	<u>0.818</u> ±0.0039	<u>0.857</u> ±0.0035	<u>0.910</u> ±0.0028	<u>0.815</u> ±0.0032
TITAN	0.784 ±0.0042	0.906 ±0.0029	0.938 ±0.0024	0.972 ±0.0017	0.900 ±0.0023

Extended Data Table 138: Cross-modal retrieval results (report-to-slide) on TCGA-Slide-Reports. The best result is marked in bold and the second best is underlined.

Encoder	Recall@1	Recall@3	Recall@5	Recall@10	Mean recall
PRISM	<u>0.552</u> ±0.0051	<u>0.693</u> ±0.0046	<u>0.744</u> ±0.0043	<u>0.807</u> ±0.0039	<u>0.699</u> ±0.0040
TITAN	0.752 ±0.0042	0.838 ±0.0037	0.871 ±0.0033	0.908 ±0.0028	0.842 ±0.0031

Extended Data Table 139: Cross-modal retrieval results (slide-to-report) on TCGA-Slide-Reports. The best result is marked in bold and the second best is underlined.