*Supplementary Material for*
# Scalable Whole-Slide Vision-Language Modeling with Learned Token Pruning

**Ali Kerem Bozkurt**[†,⋆]**, Baris Cem Bakay**[†,⋆]**, Ibrahim Kulac**[∗,⋆]**, Cigdem Gunduz Demir**[†,⋆]**,**
**Erkut Erdem**[‡,⋆]**, Aykut Erdem**[†,⋆]

[†] Department of Computer Engineering, Koç University, Istanbul, Turkey
[∗] Department of Pathology, School of Medicine, Koç University, Istanbul, Turkey
[‡] Department of Computer Engineering, Hacettepe University, Ankara, Turkey
[⋆] KUIS AI Center, Koç University, Istanbul, Turkey

This supplementary material provides additional experimental analyzes and supporting details. We begin with zero-shot classification results on the TCGA test set (Sec. A), followed by confusion matrices for TCGA and EBRAINS (Sec. B). We then present ablations on model depth (Sec. C) and positional embeddings (Sec. D), as well as qualitative visualizations of progressive pruning (Sec. E). Finally, we provide the class distributions of the EBRAINS (Sec. F) and TCGA (Sec. G) datasets to support reproducibility and future benchmarking.

## A    Zero-Shot Classification on the Test Set

We evaluate zero-shot classification on the TCGA test set using OncoTree codes, resulting in a total of 747 samples. Table 1 summarizes the performance across multiple metrics, showing that token pruning improves results consistently over the no-pruning baseline.

Table 1: **Zero-shot classification results on TCGA.** Token pruning improves performance across all reported metrics.

| Model | Accuracy | Bal. Acc. | Weighted Precision | Weighted Recall | Weighted F1 | AUROC |
|---|---|---|---|---|---|---|
| SLIM (No Pruning) | 0.641 | 0.498 | 0.692 | 0.641 | 0.605 | 0.981 |
| SLIM | **0.656** | **0.531** | **0.693** | **0.655** | **0.609** | **0.982** |

To generate predictions, we adopt the text prompts introduced in [1], which provide diverse natural language descriptions of each diagnostic class. The full prompt set is listed below:

- 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 slideshows 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.

For each class, embeddings are computed for all prompts and aggregated via average pooling, using the class list in Table 5 as the label space. The resulting confusion matrix is provided in Figure 1, highlighting class-specific strengths and weaknesses.



Figure 1: **Confusion matrix illustrating zero-shot classification performance of our model on the TCGA test set.**

# B Confusion Matrices for Linear Probing

Figures 2 and 3 show confusion matrices for linear probing on the TCGA and EBRAINS datasets, respectively. These visualizations provide a detailed view of class-level performance beyond the aggregate metrics reported in the main text.

On TCGA, the confusion matrix highlights that most tumor classes are well separated, with errors concentrated among morphologically similar subtypes. On EBRAINS, misclassifications are more frequent, reflecting the smaller dataset size, greater heterogeneity, and cross-domain shift. In particular, several tumor types—such as IDH-mutant versus IDH-wildtype astrocytomas or oligodendrogliomas defined by 1p/19q co-deletion—are nearly indistinguishable from H&E morphology alone without molecular profiling. These patterns are consistent with the known diagnostic challenges in neuropathology, where integrated histological and molecular analysis is required.
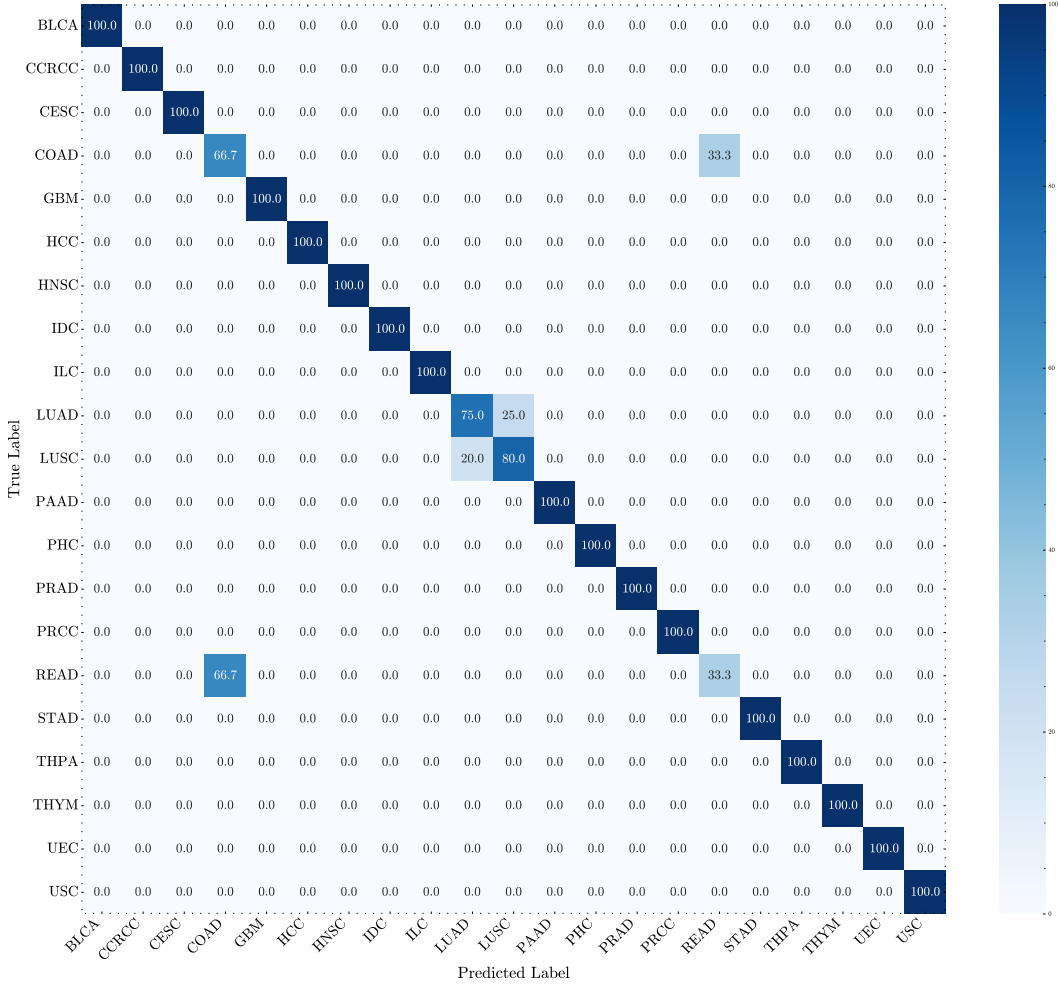


Figure 2: **Confusion matrix illustrating linear probing performance of our model on the TCGA test set.**
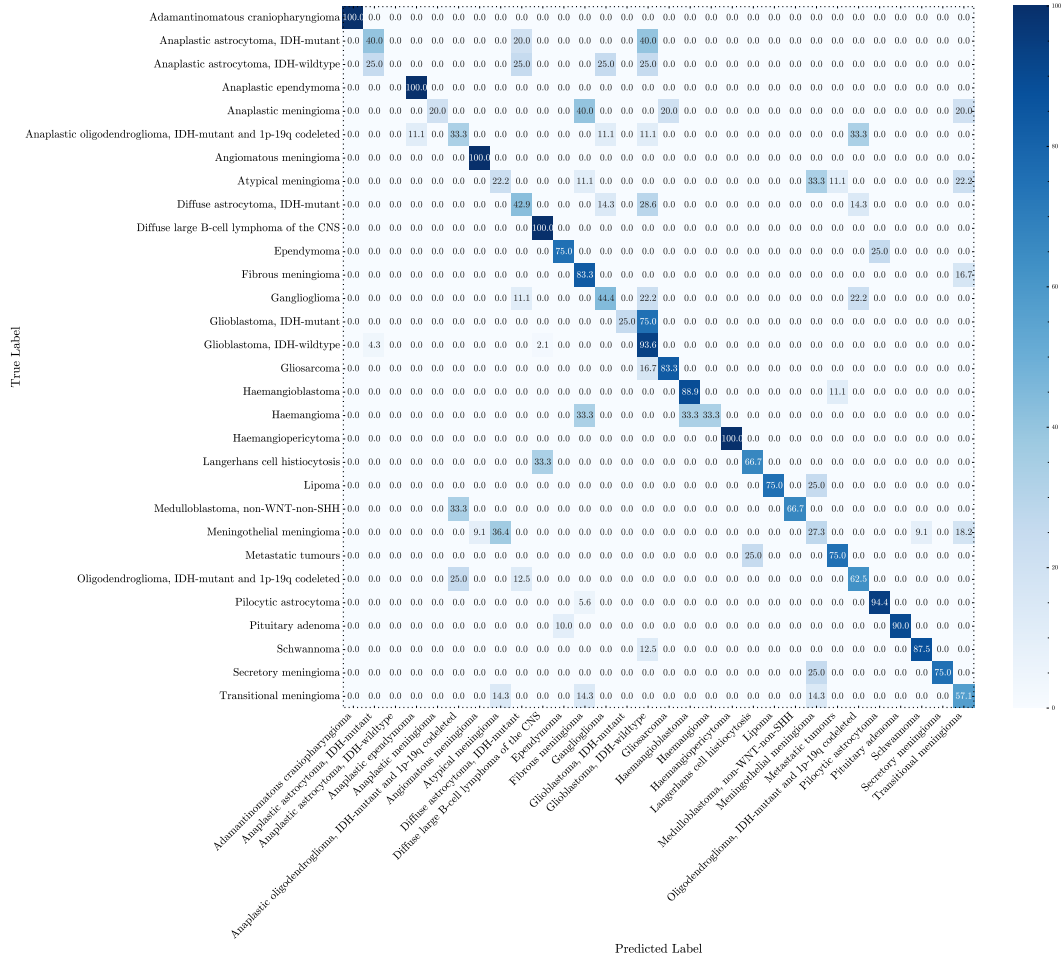
Figure 3: **Confusion matrix illustrating linear probing performance of our model on the EBRAINS dataset.**

## C   Effect of Model Depth

We analyze the impact of transformer depth on classification performance, with results summarized in Table 2. Both TCGA and EBRAINS exhibit an optimal range, with six layers providing the best balance of accuracy, balanced accuracy, and F1. This depth appears sufficient to capture global context without introducing unnecessary complexity. Adding more layers does not consistently improve accuracy, and on TCGA performance often plateaus or slightly degrades, indicating diminishing returns once model capacity exceeds dataset scale. On EBRAINS, a smaller and more heterogeneous dataset, depth has a stronger effect: shallow models (4–5 layers) underperform, while deeper configurations beyond six trade off modest drops in accuracy for gains in agreement metrics such as kappa and AUROC. This suggests that while six layers offer a strong efficiency–accuracy trade-off across datasets, additional depth may improve robustness under variability in acquisition conditions.

Table 2: **Effect of model depth on classification performance.** Results across transformer depths on TCGA and EBRAINS. On both datasets, 6 layers provide the best overall balance of accuracy, agreement, and AUROC, though EBRAINS shows larger gains from depth.

| Dataset | Layers | Accuracy | Bal. Acc. | Kappa | Weighted F1 | AUROC |
|---------|--------|----------|-----------|-------|-------------|-------|
| TCGA | 4 | 0.841 | 0.818 | 0.887 | 0.810 | 0.997 |
| | 5 | 0.857 | 0.880 | 0.887 | 0.841 | 0.995 |
| | 6 | **0.921** | **0.931** | 0.888 | **0.919** | 0.996 |
| | 7 | 0.889 | 0.917 | 0.888 | 0.888 | 0.997 |
| | 8 | 0.857 | 0.892 | 0.882 | 0.857 | 0.994 |
| | 10 | 0.905 | 0.926 | **0.925** | 0.901 | **0.997** |
| EBRAINS | 4 | 0.690 | 0.612 | 0.611 | 0.666 | 0.969 |
| | 5 | 0.685 | 0.606 | 0.641 | 0.664 | 0.966 |
| | 6 | **0.711** | **0.654** | 0.650 | **0.693** | 0.970 |
| | 7 | 0.698 | 0.633 | 0.611 | 0.688 | **0.971** |
| | 8 | 0.690 | 0.611 | **0.685** | 0.683 | 0.971 |
| | 10 | 0.681 | 0.623 | 0.641 | 0.666 | 0.967 |

## D   Effect of Positional Embeddings

We compare rotary positional embeddings (RoPE) with standard sine–cosine encodings for slide-level modeling (Table 3). On TCGA, RoPE is consistently stronger across metrics. It improves accuracy by nearly 8 points (0.921 vs. 0.841) and yields stronger agreement metrics and F1. On EBRAINS, RoPE achieves higher accuracy (0.711 vs. 0.681) and balanced accuracy (0.654 vs. 0.599), while maintaining similar AUROC. These results suggest that rotary embeddings are better suited for capturing fine-grained spatial structure in WSIs, particularly under domain shift, where sine–cosine encodings appear less robust.

Table 3: **Effect of positional embeddings.** RoPE yields clear gains on TCGA and improves accuracy/recall-oriented metrics on EBRAINS.

| Dataset | Positional Embedding | Accuracy | Bal. Acc. | Kappa | Weighted F1 | AUROC |
|---------|---------------------|----------|-----------|-------|-------------|-------|
| TCGA | RoPE | **0.921** | **0.931** | **0.888** | **0.919** | **0.996** |
| | SineCos | 0.841 | 0.858 | 0.805 | 0.839 | 0.994 |
| EBRAINS | RoPE | **0.711** | **0.654** | 0.650 | **0.693** | **0.970** |
| | SineCos | 0.681 | 0.599 | **0.687** | 0.666 | 0.969 |

## E   Visualization of Progressive Token Pruning

Figure 4 illustrates how the Cropr modules gradually reduces sequence length across transformer layers. Beyond showing retained and discarded tokens, the visualizations highlight the dynamics of

|  | Layer 1 | Layer 2 | Layer 3 | Layer 4 | Layer 5 |

Pruning Rate: %27.6   Pruning Rate: %47.5   Pruning Rate : %62.0   Pruning Rate : %72.4   Pruning Rate: %80.0
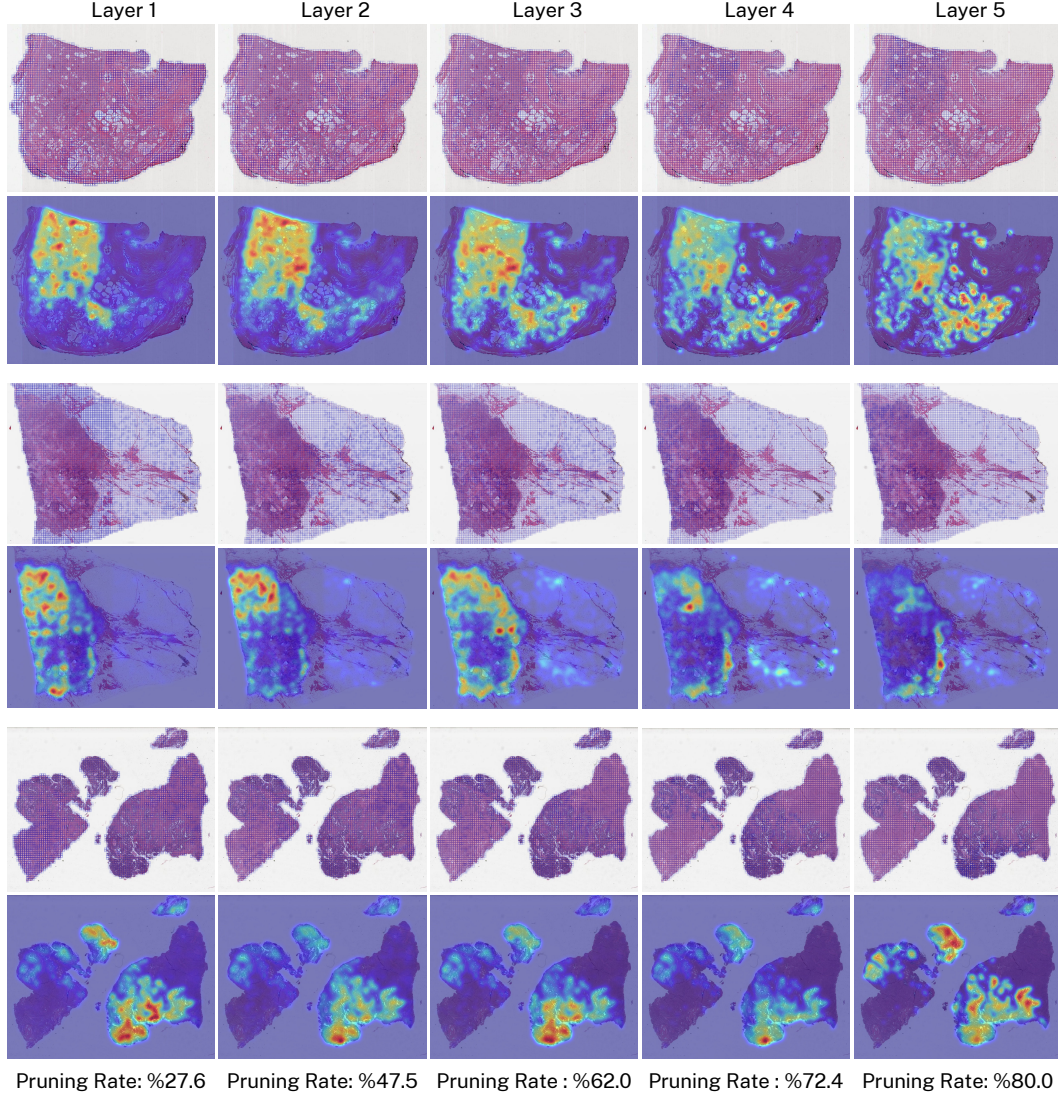
Figure 4: **Visualization of Progressive Token Pruning Across Layers.** Top rows show the retained histology patches (light color), while bottom rows overlay Cropr attention maps, where warmer colors indicate tokens assigned higher salience. As depth increases, uninformative background regions are gradually discarded and the model focuses more tightly on diagnostically relevant tissue areas.

the pruning process itself. The top row of each example shows the remaining tissue patches after pruning, while the bottom row overlays token salience scores derived from Cropr's attention weights.

Early layers eliminate broad, homogeneous regions, while later layers focus the model's capacity on localized tissue structures. The progression reveals that pruning operates hierarchically: it first removes obvious background, then progressively sharpens the representation around diagnostically relevant areas.

This behavior is important for two reasons. First, it shows that token pruning is not a blunt reduction mechanism but an adaptive process that reallocates capacity as depth increases. Second, it provides a form of interpretability: the regions that consistently survive pruning align with areas a pathologist would attend to, suggesting that the model's efficiency gains are coupled with meaningful inductive biases.

# F EBRAINS Class Distribution

To support reproducibility, we report the class distribution of the EBRAINS dataset used in our experiments. As described in the main paper, we filtered out categories with fewer than 30 samples, resulting in 2,319 slides spanning 30 diagnostic classes. Table 4 lists the number of samples per class, which may serve as a reference for future benchmarking and comparisons.

Table 4: **EBRAINS class distribution.** Number of samples per class after filtering out categories with fewer than 30 slides. The final dataset contains 2,319 slides across 30 diagnostic classes.

| Class | Sample Count |
|---|---|
| Adamantinomatous craniopharyngioma | 85 |
| Anaplastic astrocytoma, IDH-mutant | 47 |
| Anaplastic astrocytoma, IDH-wildtype | 47 |
| Anaplastic ependymoma | 50 |
| Anaplastic meningioma | 46 |
| Anaplastic oligodendroglioma, IDH-mutant and 1p-19q codeleted | 91 |
| Angiomatous meningioma | 31 |
| Atypical meningioma | 83 |
| Diffuse astrocytoma, IDH-mutant | 70 |
| Diffuse large B-cell lymphoma of the CNS | 59 |
| Ependymoma | 46 |
| Fibrous meningioma | 57 |
| Ganglioglioma | 88 |
| Glioblastoma, IDH-mutant | 34 |
| Glioblastoma, IDH-wildtype | 474 |
| Gliosarcoma | 59 |
| Haemangioblastoma | 88 |
| Haemangioma | 30 |
| Haemangiopericytoma | 34 |
| Langerhans cell histiocytosis | 32 |
| Lipoma | 38 |
| Medulloblastoma, non-WNT-non-SHH | 32 |
| Meningothelial meningioma | 104 |
| Metastatic tumours | 47 |
| Oligodendroglioma, IDH-mutant and 1p-19q codeleted | 85 |
| Pilocytic astrocytoma | 173 |
| Pituitary adenoma | 99 |
| Schwannoma | 81 |
| Secretory meningioma | 41 |
| Transitional meningioma | 68 |

# G TCGA Class Distribution

To support reproducibility, we also provide the distribution of classes in the curated TCGA dataset used in our experiments. After removing slides without valid OncoTree codes, the test set contains 747 slides across 42 diagnostic classes. Table 5 reports the number of samples per class. The prompts are taken from the evaluation setup in [1].

Table 5: **TCGA class distribution.** Number of samples per class in the curated TCGA dataset used in our experiments. The test set contains 747 slides spanning 42 diagnostic classes.

| Class | Sample Count | Prompts |
|---|---|---|
| IDC | 82 | invasive ductal carcinoma<br>breast invasive ductal carcinoma<br>ductal carcinoma, no special type<br>IDC |
| LUSC | 45 | squamous cell carcinoma<br>lung squamous cell carcinoma<br>squamous carcinoma of the lung<br>LUSC |
| CCRCC | 41 | clear cell carcinoma<br>renal clear cell carcinoma<br>clear cell renal cell carcinoma<br>clear cell RCC<br>CCRCC |
| LUAD | 40 | lung adenocarcinoma<br>adenocarcinoma of the lung<br>pulmonary adenocarcinoma<br>peripheral lung adenocarcinoma<br>LUAD |
| HNSC | 40 | head and neck squamous cell carcinoma<br>HNSCC<br>head and neck SCC<br>oropharyngeal squamous cell carcinoma<br>laryngeal squamous cell carcinoma<br>hypopharyngeal squamous cell carcinoma<br>nasopharyngeal squamous cell carcinoma<br>oral cavity squamous cell carcinoma<br>HNSC |
| HCC | 38 | hepatocellular carcinoma<br>liver cancer<br>hepatoma<br>HCC |
| UEC | 35 | endometrioid carcinoma<br>uterine endometrioid carcinoma<br>endometrial endometrioid carcinoma<br>endometrial carcinoma, endometrioid type<br>UEC |
| THPA | 33 | papillary thyroid carcinoma<br>papillary thyroid cancer<br>papillary thyroid neoplasm<br>PTC<br>thyroid papillary carcinoma |
| PRAD | 33 | prostate adenocarcinoma<br>adenocarcinoma of the prostate<br>prostatic adenocarcinoma<br>PRAD |
| COAD | 29 | colon adenocarcinoma<br>adenocarcinoma of the colon<br>colorectal adenocarcinoma<br>COAD |

Continued on next page

8

| Class | Sample Count | Prompts |
|-------|--------------|---------|
| BLCA | 29 | urothelial carcinoma<br>bladder urothelial carcinoma<br>transitional cell carcinoma<br>bladder cancer<br>BLCA |
| CESC | 24 | squamous cell carcinoma of the cervix<br>cervical squamous cell carcinoma<br>cervical SCC<br>CESC |
| PRCC | 23 | papillary renal cell carcinoma<br>renal cell carcinoma, papillary type<br>papillary RCC<br>PRCC |
| READ | 23 | rectal adenocarcinoma<br>adenocarcinoma of the rectum<br>rectal cancer<br>READ |
| PHC | 21 | pheochromocytoma<br>chromaffin cell tumor<br>adrenal pheochromocytoma<br>paraganglioma<br>PHC |
| GBM | 21 | glioblastoma multiforme<br>glioblastoma<br>GBM<br>multiforme glioblastoma<br>grade IV astrocytoma |
| ILC | 20 | invasive lobular carcinoma<br>breast invasive lobular carcinoma<br>invasive lobular carcinoma of the breast<br>lobular carcinoma<br>breast ILC |
| PAAD | 19 | pancreatic adenocarcinoma<br>adenocarcinoma of the pancreas<br>pancreas adenocarcinoma<br>PAAD |
| STAD | 13 | diffuse type stomach adenocarcinoma<br>diffuse gastric adenocarcinoma<br>diffuse-type gastric cancer<br>gastric adenocarcinoma, diffuse type<br>STAD |
| THYM | 11 | thymoma<br>thymic tumor<br>thymic carcinoma<br>thymic epithelial neoplasm<br>THYM |
| USC | 10 | uterine serous carcinoma<br>uterine papillary serous carcinoma<br>serous papillary carcinoma of the uterus<br>serous carcinoma<br>USC |

Table 5 – continued from previous page

| Class | Sample Count | Prompts |
|-------|--------------|---------|
| LMS | 9 | leiomyosarcoma<br>smooth muscle sarcoma<br>leiomyosarcoma of the soft tissue<br>LMS |
| ODG | 9 | oligodendroglioma<br>oligodendroglial tumor<br>oligodendroglioma grade II<br>oligodendroglioma grade III<br>ODG |
| CHRCC | 9 | chromophobe renal cell carcinoma<br>chromophobe RCC<br>chromophobe carcinoma<br>renal cell carcinoma, chromophobe type<br>CHRCC |
| UCS | 9 | uterine carcinosarcoma<br>uterine malignant mixed Müllerian tumor<br>malignant mixed Müllerian tumor<br>carcinosarcoma of the uterus<br>UCS |
| MACR | 9 | mucinous adenocarcinoma<br>mucinous adenocarcinoma of the colon<br>mucinous adenocarcinoma of the rectum<br>mucinous colorectal adenocarcinoma<br>MACR |
| THFO | 9 | follicular thyroid carcinoma<br>follicular thyroid cancer<br>follicular carcinoma of the thyroid<br>thyroid follicular cancer<br>THFO |
| AASTR | 7 | anaplastic astrocytoma<br>astrocytoma, anaplastic<br>grade III astrocytoma<br>AASTR |
| UM | 6 | uveal melanoma<br>uvea melanoma<br>ocular melanoma<br>choroidal melanoma<br>ciliary body melanoma<br>iris melanoma<br>UM |
| OAST | 6 | oligoastrocytoma<br>oligoastrocytic tumor<br>oligodendroglial-astrocytic tumor<br>OAST |
| TSTAD | 6 | tubular stomach adenocarcinoma<br>tubular adenocarcinoma of the stomach<br>stomach tubular adenocarcinoma<br>gastric tubular adenocarcinoma<br>TSTAD |
| DSTAD | 5 | diffuse type adenocarcinoma of the stomach<br>diffuse stomach adenocarcinoma<br>diffuse gastric adenocarcinoma<br>stomach adenocarcinoma, diffuse type<br>DSTAD |

| Class | Sample Count | Prompts |
|---|---|---|
| ASTR | 5 | astrocytoma<br>diffuse astrocytoma<br>fibrillary astrocytoma<br>anaplastic astrocytoma<br>grade II astrocytoma<br>grade III astrocytoma<br>ASTR |
| PLEMESO | 5 | pleural mesothelioma<br>epithelioid mesothelioma<br>epithelioid pleural mesothelioma<br>pleural mesothelioma, epithelioid type<br>PLEMESO |
| MEL | 4 | acral melanoma<br>melanoma<br>acral lentiginous melanoma<br>MEL |
| ESCA | 4 | esophageal adenocarcinoma<br>adenocarcinoma of the esophagus<br>esophageal cancer, adenocarcinoma type<br>ESCA |
| ESCC | 4 | esophageal squamous cell carcinoma<br>esophageal SCC<br>squamous cell carcinoma of the esophagus<br>SCC of the esophagus<br>ESCC |
| SKCM | 3 | cutaneous melanoma<br>skin melanoma<br>melanoma of the skin<br>malignant melanoma<br>SKCM |
| HGSOC | 3 | high-grade serous ovarian cancer<br>high-grade serous carcinoma<br>HGSOC<br>serous papillary carcinoma<br>serous ovarian carcinoma |
| AOAST | 2 | anaplastic oligoastrocytoma<br>oligoastrocytoma<br>anaplastic mixed glioma<br>AOAST |
| DDLS | 2 | dedifferentiated liposarcoma<br>dedifferentiated liposarcoma variant<br>liposarcoma, dedifferentiated<br>DDLS |
| ACC | 1 | adrenocortical carcinoma<br>adrenal cortical carcinoma<br>adrenal cortex carcinoma<br>ACC |

# References

[1] Tong Ding, Sophia J. Wagner, Andrew H. Song, Richard J. Chen, Ming Y. Lu, Andrew Zhang, Anurag J. Vaidya, Guillaume Jaume, Muhammad Shaban, Ahrong Kim, Drew F. K. Williamson, Bowen Chen, Cristina Almagro-Perez, Paul Doucet, Sharifa Sahai, Chengkuan Chen, Daisuke Komura, Akihiro Kawabe, Shumpei Ishikawa, Georg Gerber, Tingying Peng, Long Phi Le, and Faisal Mahmood. Multimodal whole slide foundation model for pathology, 2024.