

Scaling Self-Supervised and Cross-Modal Pretraining for Volumetric CT Transformers

Supplementary Material

Table 1. Overview of the datasets used for pretraining, summarizing anatomical coverage, the number of CT reconstructions remaining after all exclusions, and whether each dataset is used for self-supervised learning (SSL), vision–language alignment (VLA), or both.

Source dataset	Location	No. Rec.	SSL	VLA
NLST [30]	Chest	132,985	✓	
CT-RATE [8]	Chest	47,149	✓	✓
INSPECT [13]	Chest	23,226	✓	✓
Merlin [5]	Abdomen	15,314	✓	✓
AbdomenAtlas [19]	Abdomen	5,195	✓	
AMOS [16]	Abdomen	2,450	✓	
PANORAMA [3]	Abdomen	2,238	✓	
AbdomenCT-1K [22]	Abdomen	1,062	✓	
Total		229,619		

A. Pretraining Data and Preprocessing

A.1. Datasets

We curated a diverse collection of 3D CT scans from multiple publicly available datasets. Three of these datasets also include accompanying clinical metadata, such as radiology reports and EHR diagnostic codes, and can therefore be used for both SSL and VLA. Imaging data span the thoracic, abdominal, and pelvic regions and comprise a wide range of acquisition settings, including variations in radiation dose and the use of contrast agents. After applying exclusion criteria that are provided below for every dataset, the final set for pretraining comprises 229,619 image series. A summary of the datasets and their characteristics is provided in Tab. 1.

We now provide dataset-level summaries, including the filtering criteria and preprocessing steps used to construct the final pretraining corpus.

- **NLST** (the National Lung Screening Trial) [30] provides lung cancer screening data collected in the United States between 2002 and 2004. The dataset consists of low-dose helical chest CT scans from 26,254 participants with a two-year follow-up, yielding 73,116 studies. Owing to multiple reconstruction settings, the original release contains 203,099 series. For our purposes, we retain only one series per reconstruction kernel; if multiple series are available for the same kernel, we select the one with the largest number of slices. This yields at most two series per CT study, for a total of 132,985 image series. All relevant DICOM series are converted into 3D volumes in NiFTi format.
- **CT-RATE** [8] provides paired 3D chest CT volumes and the corresponding radiology reports of 21,304 patients. It comprises 25,692 non-contrast chest CT studies, expanded to 50,188 series through multiple reconstructions. Each study is accompanied by the radiology report, including both *Findings* and *Impressions* sections, as well as multi-abnormality labels and metadata. The cohort is split into 20,000 patients for training and 1,304 patients for validation. To ensure data consistency, we excluded unintended head CT series by generating segmentation masks with the TotalSegmentator model [34] and removing scans where the proportion of voxels labeled as *brain* or *skull* was an outlier, with outliers defined by Tukey’s rule ($Q3 + 1.5 \times IQR$) of the distribution of relative brain/skull volume. We also respect the original split and use only the original training set for pretraining of our model, resulting in a total of 47,149 image series for training.
- **INSPECT** [13] consists of CT pulmonary angiography (CTPA) scans paired with radiology reports that include the *Impressions* section. It contains imaging data from 19,402 patients with a total of 23,248 studies. To address data issues, we excluded partially uploaded files, resulting in a final set of 23,226 CTPA studies.
- **Merlin** [5] contains abdominal CT scans acquired at the Stanford Hospital Emergency Department between 2012 and 2018. It includes 25,494 studies from 18,317 patients, each paired with a radiology report comprising sections *Findings* and *Impressions*, as well as associated EHR diagnostic codes. The data is split into training, validation, and test sets, with 15,314 studies in the training set. We use the training split of the dataset in its provided form without applying further filtering or modifications to train our model.
- **AbdomenAtlas1.0Mini** [19] is a fully annotated publicly accessible subset of the larger AbdomenAtlas dataset, comprising 5,195 abdominal CT volumes with segmentations at the voxel-level. The annotations cover nine key anatomical structures: spleen, liver, left kidney, right kidney, stomach, gallbladder, pancreas, aorta, and inferior vena cava. The source images are aggregated from multiple existing public datasets, each of which contributes cases with varying imaging protocols, disease states, and anatomical coverage. For pretraining, we use only the raw CT scans without the accompanying segmentation labels.
- **AMOS** [16] is a multi-center dataset designed for multi-organ abdominal segmentation across diverse clinical sce-

narios. It comprises 500 CT and 100 magnetic resonance images (MRI) with voxel-level annotations for 15 abdominal organs, collected from multi-vendor and multi-phase acquisitions spanning a wide range of disease conditions. In addition, AMOS provides 1,900 unlabeled CT and 1,200 unlabeled MRI scans to support semi-supervised and unsupervised learning tasks. After excluding 50 corrupted or incomplete CT files, we retain a total of 2,450 CT scans for pretraining.

- **PANORAMA** [3] is a contrast-enhanced abdominal CT dataset designed to benchmark diagnostic performance for pancreatic ductal adenocarcinoma (PDAC) detection and diagnosis. It includes 2,238 anonymized CT scans acquired at two Dutch medical centers (Radboud University Medical Center and University Medical Center Groningen). The dataset was curated to ensure high-quality imaging and standardized acquisition protocols.
- **AbdomenCT-1K** [22] is a large and diverse abdominal CT dataset comprising 1,062 scans collected by aggregating multiple public single-organ datasets. It includes both contrast-enhanced and non-contrast studies with voxel-level annotations for four major abdominal organs: liver, kidneys, spleen, and pancreas. For our purposes, we use only the raw CT scans for pretraining.

A.2. Image Processing for Self-Supervised Learning

For SSL with the adapted DINOv3, all CT series listed in Tab. 1 are first reoriented to a common anatomical coordinate system (right-left, anterior-posterior, superior-inferior; RAS) to ensure spatial consistency. The voxel intensities in Hounsfield Units (HU) are clipped to the range $[-1000, +1000]$ to remove outliers and normalized to the unit range with 32-bit precision. Each scan is then resampled to a voxel spacing of $0.5 \times 0.5 \times 1.0$ mm using trilinear interpolation and center-cropped to a maximum size of $512 \times 512 \times 384$ voxels, sufficient to cover the FOV of most chest and abdominal scans while minimizing surrounding background. The resulting volumes are converted to 16-bit tensors and stored on disk. During training, these tensors are loaded into memory and a random crop of $256 \times 256 \times 128$ voxels is extracted for each batch element. To further improve I/O efficiency, four random crops are sampled from the same CT scan and treated as separate batch elements.

A.3. Data Processing for Vision-Language Alignment

For VLA with SigLIP, we first select the CT series that contain accompanying text from radiology reports or diagnostic codes (see Tab. 1). Following, all series are re-oriented to the RAS coordinate system, clipped within the range $[-1000, +1000]$ HU and normalized to a unity range, followed by resampling to a voxel spacing of $0.75 \times 0.75 \times$

Table 2. Training hyperparameters during self-supervised learning (SSL) and vision-language alignment (VLA) pretraining stages; LR, LLRD, and WD denote learning rate, linear learning rate decay, and weight decay.

Hyperparameter	SSL	VLA
Epochs	150	100
Steps	83,100	59,000
Batch size	1,536	144
LR begin	4×10^{-3}	1×10^{-4}
LR end	1×10^{-6}	1×10^{-6}
LR schedule	Cosine annealing	Cosine annealing
LR warmup epochs	10	10
LLRD	0.9	0.9
WD begin	0.04	0.01
WD end	0.4	0.01
WD schedule	Inverse cosine annealing	Constant
Optimizer	AdamW	AdamW
	$(\beta_1 = 0.9, \beta_2 = 0.999)$	$(\beta_1 = 0.9, \beta_2 = 0.95)$

1.5 mm using trilinear interpolation. Again, we write the resulting tensors in 16-bit precision to disk and load accordingly during training.

To enhance textual descriptions, each paragraph in the radiology reports is expanded to multiple paraphrases. First, the reports are divided into two sections, *Findings* and *Impressions*, whenever these sections are provided. For each section, we prompt a large language model (LLM) to “rephrase clearly and concisely *without changing any medical facts*,” and to “*only return the revised text*”, ensuring clarity and consistency.

This prompt is supported by four examples curated by radiologists (two chest CT, two abdominal CT), demonstrating accurate rewrites that preserve clinical elements such as laterality, measurements, and negations. Using Google’s *Gemini 2 Flash*, we generate two additional paraphrases per section, resulting in three semantically equivalent versions. During training, a version is randomly sampled, following the LaCLIP single positive strategy [7], to provide clinically accurate supervision of vision-language alignment.

In addition to text, some reports include structured EHR diagnosis codes. Since LLMs struggle with raw codes (*e.g.*, J18.9), we replace each with its World Health Organization 2025 short description¹. These are appended as a comma separated list at the end of the report to form a complete text input. The EHR descriptions are not rewritten and are added after LLM rewriting to preserve billing and epidemiological accuracy.

Table 3. Architectural and model-scale specifications of the local (ViT_ℓ) and global (ViT_g) parts of SPECTRE.

Configuration	ViT _ℓ	ViT _g
Trainable parameters	339M	58M
Patch size	16 × 16 × 8	-
Layers	24	4
Embedding dimension	1,080	1,080
Attention heads	12	12
Position embedding	3D RoPE	3D RoPE
LayerScale initialization	0.1	1.0

B. Pretraining & Architectural Details

B.1. Training and Model Configuration

Tab. 2 summarizes the training parameters for both SSL with DINOv3 and VLA with SigLIP. During VLA, the pre-trained ViT_ℓ remains frozen for the first 10 epochs. To improve efficiency, all models are trained with mixed precision (FP16) and optimized using distributed data parallelism. Data loading is performed via GPU Direct Storage, enabling high-throughput I/O directly to device memory and thereby minimizing bottlenecks in large-scale training. Tab. 3 shows the architectural choices of SPECTRE, separated for the local ViT_ℓ and global ViT_g. Note that during VLA, the LayerScale layers of ViT_ℓ are initialized with the weights found during SSL.

B.2. Hardware

Both pretraining phases of the foundation model are conducted on a cluster of three DGX B200 systems (NVIDIA Corp., CA, USA), totaling 24 Blackwell GPUs with 4.32 TB of combined GPU memory. Each system contains 8 B200 GPUs (1.44 TB per DGX), dual Intel Xeon Platinum 8570 processors (112 cores, 224 threads), and 2.16 TB of system memory, yielding a cumulative 6.48 TB across the cluster. The three DGX systems are interconnected via high-speed InfiniBand, enabling efficient distributed training and data exchange.

C. Downstream Experiments

C.1. Cancer Image Biomarker Prediction

This section complements the analyses on the *Cancer Image Biomarker Prediction* experiments and provides additional details on the benchmark tasks, datasets, and evaluation setup used in the downstream cancer imaging experiments.

¹<https://www.who.int/standards/classifications/classification-of-diseases>

C.1.1. Foundation Models

We compare a comprehensive set of eleven publicly available CT foundation models: FMCIB [24], CT-FM [25], CT-CLIP [8], PASTA [18], VISTA3D [11], VOCO [35], SUPREM [20], Merlin [5], MedImageInsight [6], Models-Genesis [38], and our proposed SPECTRE. These models collectively represent the current generation of volumetric CT foundation models, spanning both unimodal (image-only) and multimodal (image-text) pretraining paradigms. More about these models can be found in the *Related Works* of the main paper and in the models’ respective papers.

C.1.2. Evaluation Framework

All models are evaluated within the standardized *TumorImagingBench* reference framework introduced by Pai et al. [26]. This framework ensures that embeddings are extracted under consistent preprocessing conditions, reproducing the *intensity normalization*, *crop sizes*, and *voxel spacings* used during each model’s original pretraining. Such harmonized extraction allows for direct comparison of representation quality across models without task-specific retraining. For SPECTRE, we use a default input size of 128 × 128 × 64 voxels with a voxel spacing of 0.5 × 0.5 × 1.0 mm. Since SPECTRE is trained agnostic to input crop size and spacing, we double the field of view for the *NSCLC-Radiogenomics* [4] and *Colorectal-Liver-Metastases* [29] datasets to ensure that all lesions are fully contained within the input volume. All other datasets use the default configuration.

C.1.3. Tasks & Datasets

The TumorImagingBench spans six public datasets covering diagnostic and prognostic tasks in thoracic, renal, and hepatic oncology. The benchmark includes two task types: (1) lung nodule malignancy classification, and (2) prediction of two-year survival across multiple tumor sites. A brief overview of the datasets used in our experiments is provided below.

- **LUNA16** [28]. A dataset containing 888 CT scans and 1,186 annotated lung nodules. We follow the established subset of 677 nodules enriched for malignancy suspicion. Task: *malignancy classification*.
- **DLCS** (Duke Lung Cancer Screening) [33]. A clinical lung-nodule cohort with 2,487 nodules from 1,613 patients; we adopt the publicly released portion with 1,714 scans and pathology-confirmed malignancy labels. Task: *malignancy classification*.
- **NSCLC-Radiomics** [1]. CT scans from 421 patients with stage I–IIIB non-small cell lung cancer (NSCLC) treated with radiation therapy, including expert Gross Tumor Volume (GTV) segmentations. Task: *two-year survival prediction*.
- **NSCLC-Radiogenomics** [4]. Surgical NSCLC cohort with preoperative CT/PET imaging; we use 133 cases

with curated GTV segmentations. Task: *two-year survival prediction*.

- **C4KC-KiTS** [12]. Renal tumour cohort from partial or radical nephrectomy patients; after filtering for complete segmentations and follow-up, 134 cases remain. Task: *two-year survival prediction*.
- **Colorectal-Liver-Metastases** [29]. Preoperative CT scans from 194 patients undergoing resection of colorectal liver metastases, using the largest lesion per patient. Task: *two-year survival prediction*.

C.1.4. Analysis

Further quantitative analyses on these tasks are provided in Fig. 1, which reports per-model performance with 95% confidence intervals. Notably, for LUNA16, DLCS, and NSCLC-Radiomics, tasks on which our model outperforms all competing approaches, the confidence intervals are narrow, indicating stable performance and low variance across cross-validation folds. In contrast, for NSCLC-Radiogenomics and Colorectal-Liver-metastases, tasks where we do not achieve SOTA performance, all models exhibit large confidence intervals and generally low scores. This is likely due to the smaller dataset sizes, which can introduce quantization noise in the AUC calculation and reflect the inherent difficulty of these tasks.

Additional qualitative evidence is shown in Fig. 2, which visualizes model explanations using saliency maps on non-curated CT samples. Without any task-specific finetuning, the model already attends to pathologic regions associated with tumor presence, indicating that the learned representations encode clinically relevant spatial features. This behavior supports the effectiveness of our pretraining strategy and echoes findings from earlier foundation-model studies demonstrating that large-scale contrastive or multimodal pretraining facilitates robust zero-shot localization and biomarker-related signal emergence [26].

C.2. Semantic Segmentation

This section details the full protocol used to evaluate SPECTRE on volumetric *Semantic Segmentation* and makes explicit the detailed experiments that led to the final SEoMT configuration reported in the main paper and the results across the benchmarks obtained.

C.2.1. Adapting EoMT to 3D Semantic Segmentation

To isolate the segmentation capability of the vision encoder itself, and not any task-specific head, we extend the Encoder-only Mask Transformer (EoMT) paradigm to volumetric (3D) semantic segmentation. In line to the original 2D EoMT, we remove all task-specific decoders and operate entirely on the encoder token space, also in the 3D case. The model starts from our SPECTRE 3D encoder (ViT_ℓ), which produces a sequence of anisotropic 3D tokens (CT patches) using the same tokenizer and 3D RoPE as in pretraining.

C.2.2. Query Design for Semantic Segmentation (3D)

After an initial set of encoder blocks, we append a fixed set of learnable query tokens to this sequence. Because semantic segmentation does not require instance enumeration, we set the number of learnable query tokens equal to the number of semantic classes in the dataset (e.g., 3 for liver, kidney, tumor). The remaining 3D encoder blocks then run joint self-attention over both volume tokens and class queries. This allows the queries to attend to spatial tokens and, symmetrically, lets spatial tokens condition on the queries, so no extra transformer decoder is required. At the output, we obtain (1) per-class embeddings from the queries and (2) a dense 3D feature grid from the encoder tokens. We project the 1/4-resolution feature grid to per-voxel class logits and trilinearly upsample it back to the original CT resolution to compute Dice and Cross-Entropy losses. Because the number of classes in medical CT is small and fixed, every query is forced to explain a coherent anatomical or lesion region, which stabilizes training and removes the need for Hungarian matching or instance-slot allocation.

C.2.3. Integration into nnU-Net

To position this as a fair encoder-only test, we integrate SPECTRE directly into nnU-Net as a drop-in encoder replacement. Apart from replacing the encoder, no architecture-specific components (multi-scale FPN-style features, convolutions for scale mixing, mask transformer decoders, etc.) are introduced. This ensures the comparison measures “representation quality of the encoder” – not engineering around it. We overwrite some of the suggested training plans with new SPECTRE plans. The images are resampled to $0.75 \times 0.75 \times 1.5\text{mm}$ and intensities are rescaled to 0-1 using the 0.5% and 99.5% datasets intensity profiles. Additionally we employ the optimizers and learning rate schedulers as suggested in [17], with an AdamW optimizer with a learning rate of 1×10^{-5} , weight decay 3×10^{-5} and gradient clipping of 1.0. Models are trained for 150 epochs with 250 steps per epoch and a batch size of 2, following the noSLL [32] finetuning pipeline. The nnU-Net with SPECTRE integration is publicly available at <https://github.com/cvivierv/nnUNet>.

C.2.4. Datasets

To avoid unstable conclusions caused by noisy, small, or historically under-annotated radiology datasets, we follow the recommended large-scale segmentation benchmarks by Isensee et al. [15]. However, to avoid data contamination, we drop the AMOS dataset as we used it for pretraining. Since our approach focuses purely on CT imaging, we also drop the datasets that contain MRI and add an additional CT dataset. Specifically we consider the datasets as provided in Tab. 4.

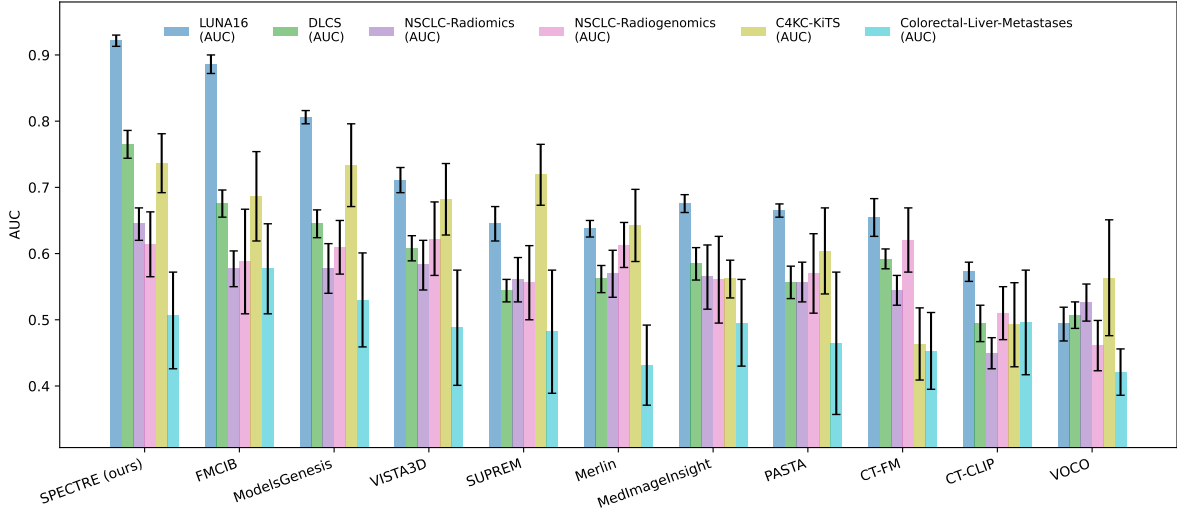


Figure 1. Quantitative comparison of 11 CT foundation models across six biomarker classification benchmarks using frozen-embedding kNN classifiers. Bars represent mean performance for each task, with error bars indicating 95% confidence intervals across cross-validation folds.

Table 4. Segmentation benchmark datasets. TS = TotalSegmentator.

Dataset	# Volumes	Classes	Description
KiTS23	489	3	kidney / tumor / cyst
LiTS	131	2	liver + tumor
WORD	120	16	16 abdominal organs
TS v1-Full	1204	104	anatomical structures
TS v201-Full	1228	117	anatomical structures
TS v201-Merlin	1228	20	anatomical structures

C.2.5. Evaluation Protocol & Results

We adopt the evaluation protocol employed in Wald et al. [31]. All experiments are conducted within the nnU-Net framework [14], with the training set randomly divided into 80%/20% train/validation splits across 5 folds. After training, the model with the best pseudo Dice is used and validation on the validation set is automatically performed. We directly record the outcome of that result and thus the average of 5-fold cross-validation. For KiTS23, we tuned SPECTRE on *fold-0* during development and exclude that fold from the final reported cross-validation to avoid optimism. All other folds and datasets use exactly the same hyperparameters to make the cross-dataset comparison meaningful.

We compare SPECTRE against various 3D domain-specific segmentation architectures. Specifically, we consider nnU-Net [14] and the updated state-of-the-art ResNet-based nnU-Net ResEnc Large [15] for comparison with convolutional-based models. Recently many transformer-based models have been developed for segmentation in 3D data. We include CoTr [36], nnFormer [37], Swin-

UNETRv2 [10], UNETR [9], WaveFormer [2] and the recent Primus [31] model for comparison. The results of these models on KiTS23 and LiTS and WORD are obtained from Wald et al. [31]. In their experiments, the models were tuned on fold-0 of KiTS23 and LiTS, and thus, the average of the other four folds are used as the baselines. All results are reported in average Dice across all classes.

During model development, we evaluated the impact of optimization strategies and input crop sizes on downstream segmentation performance after finetuning. The corresponding ablation results are reported in Tab. 5. Within the nnU-Net framework, Stochastic Gradient Descent (SGD) constitutes the default optimizer; however, our experiments show that AdamW [21] consistently outperforms SGD. Additionally, applying Deep Supervision, the default implementation in nnU-Net which computes weighed segmentation losses on intermediate lower resolution layers of the model, further improves training stability and final segmentation performance. We additionally observed that, at smaller crop sizes, models initialized with VLA (SigLIP) weights outperform counterparts initialized with just SSL (DINO). Increasing the crop size improves overall performance across initializations, while the performance gap between SSL- and VLA-initialized models narrows, rendering them comparable. Based on these observations, all subsequent experiments employ AdamW with a learning rate of 1×10^{-5} , Deep Supervision enabled, and an input crop size of $320 \times 320 \times 128$.

Fig. 4 shows qualitative masks for multiple datasets; they illustrate the same pattern as the numbers: large organs are clean and contiguous with the ground truth labels, while

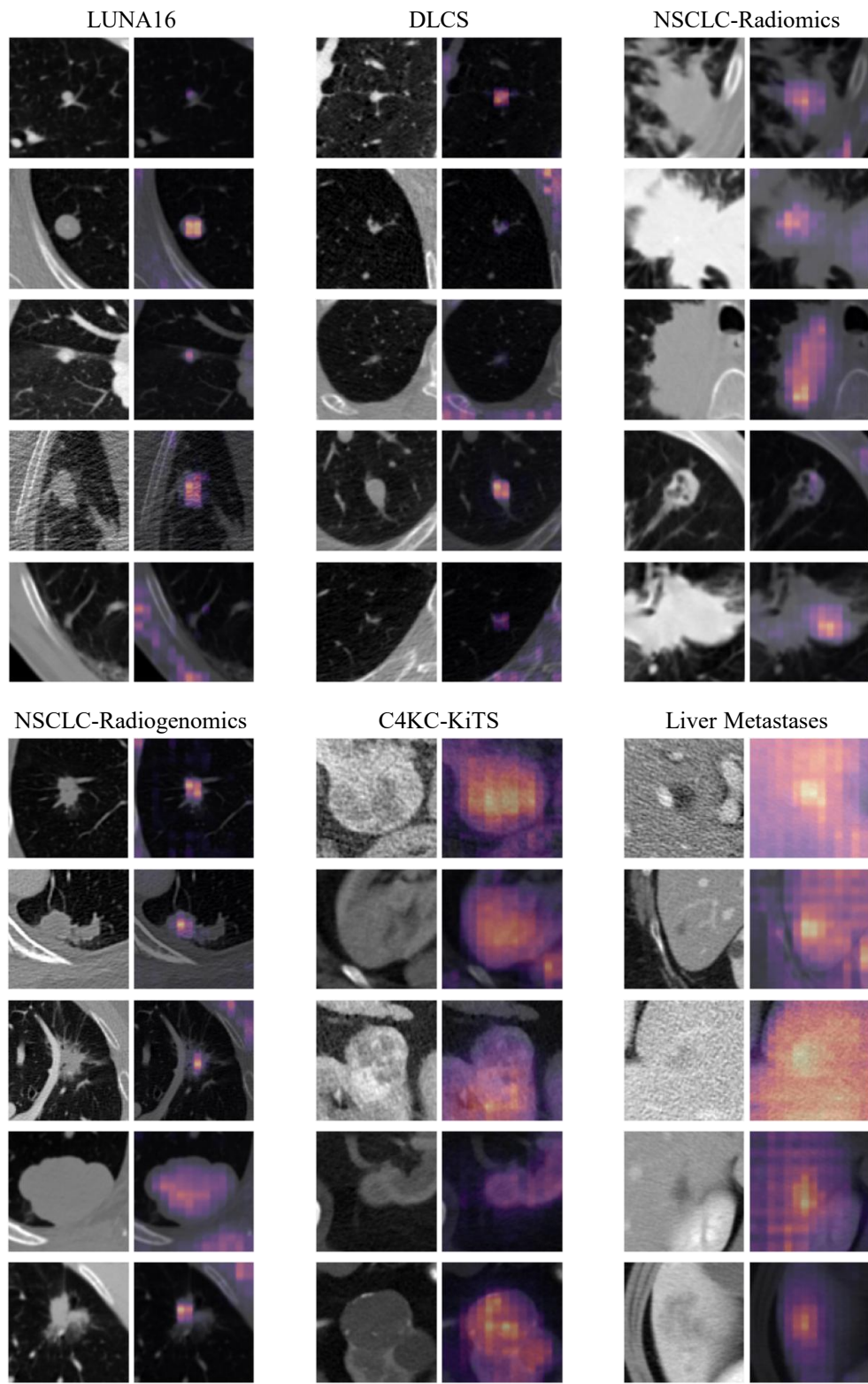


Figure 2. Non-curated saliency maps of SPECTRE on six tumor image biomarker datasets, obtained by occlusion sensitivity.

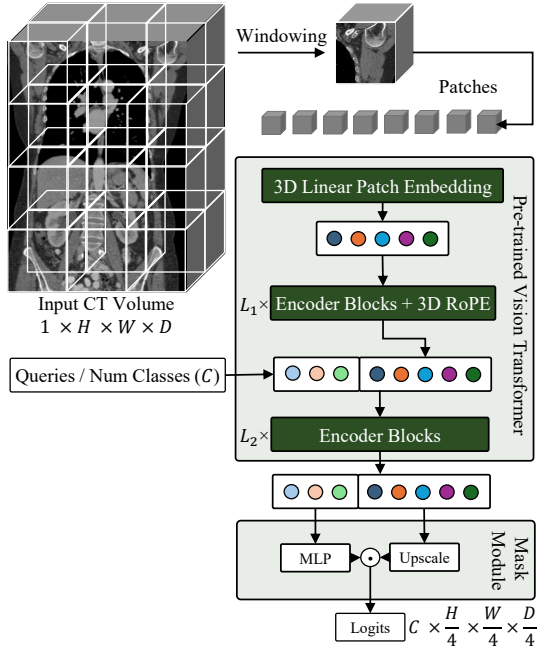


Figure 3. SEoMT architecture, derived from the EoMT. A learnable query for each class C is initialized and concatenated to the patch tokens. The new set of tokens are jointly processed by the last L_2 blocks and used to predict logits corresponding to the semantic masks.

Table 5. Ablation experiments with the Kidney tumor segmentation on KiTS23 [12] dataset. Results in Dice on fold-0. DNF = “did not finish”.

	SSL	VLA	Dice \uparrow
<i>Optimizer @ 320 × 320 × 128 input crop</i>			
SGD with LR 1×10^{-4}	✓	✓	DNF
SGD with LR 1×10^{-5}	✓	✓	0.763
AdamW with LR 1×10^{-5}	✓	✓	0.868
+ Deep Supervision	✓	✓	0.871
<i>Crop Size</i>			
128 × 128 × 64	✓		0.854
128 × 128 × 64	✓	✓	0.862
320 × 320 × 128	✓		0.871
320 × 320 × 128	✓	✓	0.871

small tumors are present but slightly smoothed—consistent with predicting at 1/4 resolution and upsampling. We chose not to add an extra refinement head to keep the experiment honest.

C.2.6. Additional Comparison on TotalSegmentator Benchmarks

To further strengthen the segmentation study beyond KiTS23, LiTS, and WORD, we additionally evaluate

Table 6. Segmentation results (Dice) over the last 3 datasets. TS is TotalSegmentator.

	Method	Dice (%) \uparrow		
		TSv1-Full	TSv2-Full	TSv2-Merlin
Conv.	SuPreM (U-Net)	-	86.95	-
	CT-FM (Res.U-Net)	-	89.81	90.17
	Merlin (Res.U-Net)	-	-	86.20
	nnU-Net	85.22	-	-
Trans.	SAM-Med3D (1 click)	84.68	-	-
	SAM-Med3D (10 clicks)	87.59	-	-
	SPECTRE (SEoMT)	87.34	88.85	87.29

SPECTRE on the three TotalSegmentator-based benchmarks reported in Tab. 6: *TSv1-Full*, *TSv2-Full*, and the more distribution-aligned *TSv2-Merlin* subset. These experiments are included to assess whether the encoder-only SEoMT formulation remains competitive on broader anatomical segmentation tasks and to compare against recent CT foundation models that were explicitly designed for segmentation. In particular, we compare against SuPreM, CT-FM, and Merlin, and we additionally include SAM-Med3D as a strong interactive transformer baseline where available.

The results in Tab. 6 show that SPECTRE with SEoMT remains consistently competitive across all three benchmarks. On *TSv1-Full*, SPECTRE achieves a Dice score of 87.34%, improving over the conventional nnU-Net baseline (85.22%) and also exceeding the one-click SAM-Med3D result (84.68%). The ten-click SAM-Med3D result is slightly higher (87.59%), but this comes at the cost of interactive prompting, whereas SPECTRE operates fully automatically in a feed-forward manner. On *TSv2-Full*, SPECTRE reaches 88.85%, outperforming SuPreM (86.95%) while trailing CT-FM (89.81%). On the *TSv2-Merlin* subset, SPECTRE obtains 87.29%, again remaining competitive and improving over Merlin, though CT-FM achieves the strongest score (90.17%).

Overall, these results support two conclusions. First, the proposed encoder-only adaptation is not limited to kidney or lesion segmentation, but transfers well to large-scale multi-structure CT benchmarks. Second, although SEoMT is not intended as an aggressively optimized decoder for state-of-the-art segmentation, it provides strong performance with minimal decoder-specific bias and therefore offers a more direct probe of encoder feature quality. We therefore position these experiments primarily as evidence that the learned SPECTRE representation is broadly useful for segmentation, rather than as a claim that SEoMT is the final or optimal decoder for 3D CT. Future work should investigate stronger task-specific 3D decoders built on top of the same pretrained encoder.

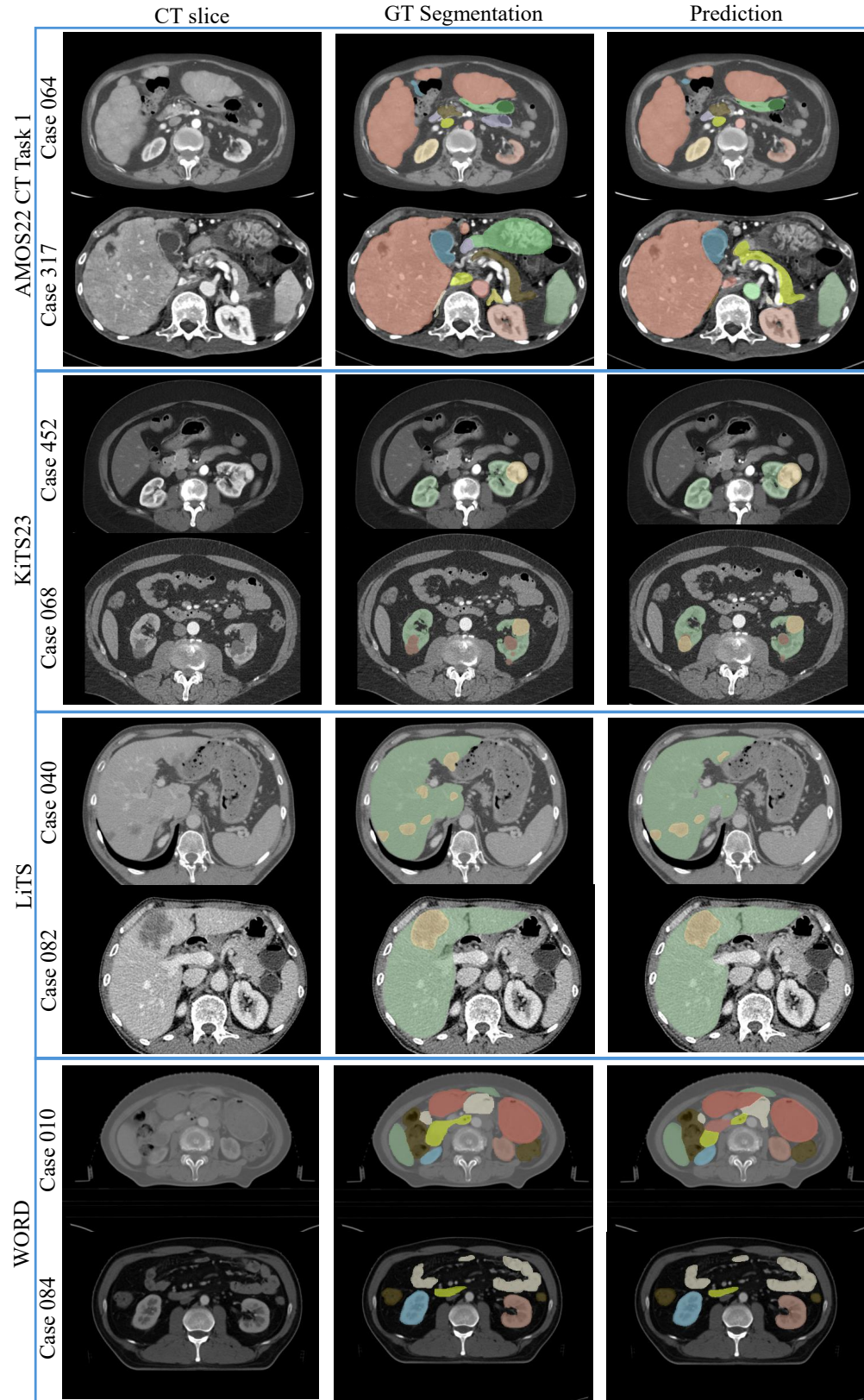


Figure 4. Curated example semantic segmentation predictions of SPECTRE on the different datasets employed in this work. Predictions with good performance and with the worst performance are depicted. Window settings optimized for organs of interest.

Table 7. Segmentation results (DSC, NSD) using different decoders. All fold 0 retrained.

Method	KiTS23		KiTS23 (<10 mm masses)		
	Dice (%) \uparrow	NSD \uparrow	Dice (%) \uparrow	NSD \uparrow	Detect Rate (%)
37.5k Training Steps (150 epochs \times 250 iterations)					
SPECTRE (Linear)	84.11	0.929	3.27×10^{-5}	0.01	2.0
SPECTRE (SEoMT)	86.70	0.943	3.76×10^{-4}	0.02	1.53
SPECTRE (UNETR)	87.53	0.945	1.34×10^{-3}	0.07	11.12
75k Training Steps (300 epochs \times 250 iterations)					
SPECTRE (Linear)	85.42	0.934	4.72×10^{-5}	0.03	1.81
SPECTRE (SEoMT)	87.13	0.948	4.11×10^{-4}	0.03	5.36
SPECTRE (UNETR)	87.82	0.948	1.71×10^{-3}	0.08	15.12
250k Training Steps (1000 epochs \times 250 iterations)					
nnU-Net ResEnc L	88.26	0.954	2.18×10^{-3}	0.11	22.79

C.2.7. Decoder Variants and Small-Structure Analysis

A potential concern with the encoder-only SEoMT design is that its simplicity may understate the true segmentation potential of the pretrained SPECTRE features. To study this explicitly, we compare three decoder choices on KiTS23 in Tab. 7: (1) a *Linear* decoder, which projects the encoder features directly to class logits; (2) the proposed *SEoMT* decoder, which appends class queries and performs joint self-attention in the final transformer blocks; and (3) a stronger *UNETR*-style decoder that introduces a more conventional task-specific decoding pathway. All models are retrained on fold 0 under matched training budgets, and we additionally report performance on the particularly challenging subset of lesions smaller than 10 mm.

The full KiTS23 results show a clear ranking across decoder complexity. Under the 37.5k-step setting (150 epochs), Linear reaches 84.11% Dice and 0.929 NSD, SEoMT improves this to 86.70% Dice and 0.943 NSD, and UNETR further increases performance to 87.53% Dice and 0.945 NSD. The same ordering remains after extending training to 75k steps (300 epochs), where Linear obtains 85.42%, SEoMT 87.13%, and UNETR 87.82% Dice. This consistent progression indicates that the SPECTRE encoder exposes useful dense features and that stronger decoders can indeed extract additional segmentation performance from them.

The analysis on tiny masses is even more informative. For lesions smaller than 10 mm, all models perform substantially worse than on the full benchmark, confirming that this regime is intrinsically difficult. Nevertheless, the same decoder trend persists: Linear yields near-zero Dice and NSD, SEoMT improves modestly, and UNETR provides the best small-structure sensitivity, increasing the detection rate from 1.53% with SEoMT to 11.12% at 37.5k steps and to 15.12% at 75k steps. For reference, the much longer-trained nnU-Net ResEnc L baseline, optimized for 250k steps, reaches a detection rate of 22.79%. These results indicate that the limitation on very small structures is not caused solely by the pretrained representation, but also by the decoding strategy and the training budget.

Importantly, SEoMT was designed to evaluate encoder feature quality with minimal decoder bias rather than to maximize segmentation performance at all costs. In our implementation, trilinear interpolation is applied to the 1/4-resolution *feature maps*, not to already discretized masks, which preserves more fine-grained spatial information despite the lightweight decoding path. Even so, Tab. 7 shows that a stronger decoder such as UNETR is beneficial, especially for tiny lesions. We therefore view SEoMT as a clean and informative encoder-centric evaluation protocol, while the decoder ablation confirms that future work on SPECTRE should investigate more expressive 3D decoders when absolute downstream segmentation performance is the primary objective.

C.3. Zero-Shot Text-to-Image Retrieval

We finally report additional details on the zero-shot text-to-image retrieval experiments conducted in parallel to the downstream evaluations. Our goal is to align the protocol as closely as possible with prior work; all retrieval metrics and data splits follow the procedures described in CT-RATE [8] and MERLIN [5].

C.3.1. Retrieval on CT-RATE validation cohort

For CT-RATE, we evaluate retrieval performance using both the *Impressions* and *Findings* sections of each radiology report. Following the original setup, each report is treated as a single textual query, and we compute $\text{Recall}@ \{5, 10, 50, 100\}$ on the full validation set of $N = 1,564$ studies. Retrieval is based on cosine similarity in the shared image-text embedding space, and a query is counted as correct if the paired CT scan appears among the top- K nearest neighbors.

Fig. 5 visualizes the joint embedding distribution of CT-RATE after UMAP projection, showing extensive cross-modal overlap and a smooth trajectory correlated with the total number of abnormalities described in the reports. This overlap indicates that the model learns a coherent shared latent space in which radiology images and their associated reports are embedded consistently, suggesting that the representations are largely modality-agnostic and capture clinically meaningful semantics rather than modality-specific artifacts. The continuous progression along the manifold with an increasing abnormality count further supports the notion that the embedding space encodes a graded representation of pathological severity or complexity.

We repeat the same experiment using MedSigLIP², which forms the visual encoder of Google’s MedGemma model [27]. Retrieval performance is low, with $\text{Recall}@ \{5, 10, 50, 100\} = \{0.3, 0.7, 4.8, 8.2\}\%$, only slightly above random chance. This limited performance can likely be attributed to the model’s restriction to 128 input tokens, which truncates longer radiological reports and prevents the

²<https://github.com/Google-Health/medsiglip>

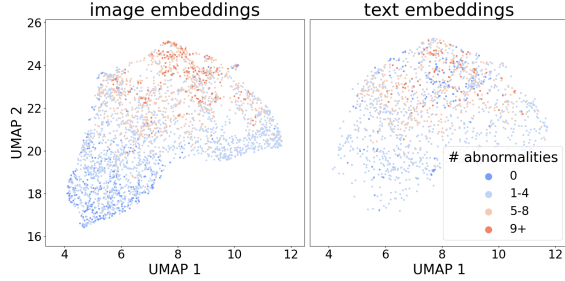


Figure 5. UMAP [23] visualization of image and text embeddings from the CT-RATE validation set [8]. Each point represents a sample categorized by the number of abnormalities noted in the corresponding radiology report.

model from accessing much of the available descriptive information.

All retrieval experiments are conducted using a fixed voxel spacing of $0.5 \times 0.5 \times 1.0$ mm. To assess the robustness of our model to variations in scan resolution and anisotropy, we also perform the CT-RATE experiment using each scan’s native spacing. We observe minimal performance change (-0.6% in Recall@5), demonstrating that our model is largely insensitive to differences in voxel resolution and anisotropy, which is important for real-world clinical applicability.

C.3.2. Retrieval on MERLIN test cohort

For MERLIN, we mirror the evaluation strategy presented in the original paper. Retrieval is conducted separately for the *Impressions* and the *Findings* sections, and performance is quantified using Recall@{1, 8}. Rather than the full test set, MERLIN evaluates retrieval over sampling pools of fixed sizes $N \in \{32, 64, 128\}$, each representing a different difficulty level. Cosine similarity is again used to rank image–text pairs, and correctness is assessed based on whether the paired CT volume is returned within the top- K matches.

Medical reports are inherently noisy due to variability in clinicians’ writing styles, abbreviations, and selective reporting. To assess the robustness of our text-to-image retrieval model under such realistic noise conditions, we simulate report corruption in two complementary ways. First, we perform *random token dropout*, which models inconsistencies in clinical phrasing. For instance, a report might mention “tumor” rather than the more specific “lung tumor,” reflecting incomplete or abbreviated descriptions. Second, we apply *random span dropout*, where contiguous spans of 10–50 tokens are removed throughout the report to simulate missing observations or unrecorded findings. The results of these experiments are shown in Fig. 6. As anticipated, model performance degrades more significantly under span dropout than token dropout, reflecting the greater impact of missing semantic content. Interestingly, perfor-

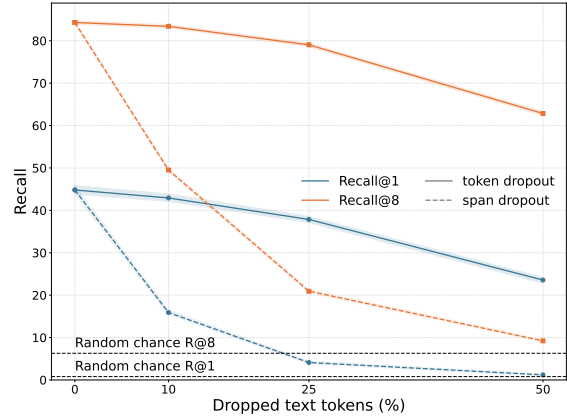


Figure 6. Impact of text dropout on retrieval performance.

mance remains relatively stable under token dropout: even when 25% of all tokens are removed, the decline in Recall@{1, 8} both remain below 10%. This demonstrates the robustness of the Qwen3 Embedding model with LoRA adapters in capturing medical language semantics, maintaining meaningful retrieval even when reports are partially incomplete. These findings highlight the model’s potential for real-world clinical applications, where reports are often imperfect or partially specified.

We further analyze the model performance with respect to report length by splitting the dataset into long reports (more than 500 tokens) and short reports (fewer than 500 tokens). We observe a notable difference in retrieval performance, with Recall@1=48.7% for long reports compared to Recall@1=34.7% for short reports. This suggests that the model effectively leverages the richer, more detailed information present in longer reports, allowing for more precise alignment with corresponding images. In contrast, shorter reports provide less context and fewer descriptive cues, which limits the model’s ability to establish strong associations. We note, however, that shorter reports often correspond to healthy subjects, where findings are minimal and reports tend to be more uniform, which could also contribute to the observed performance gap.

C.4. Hardware

All downstream and ablation experiments are performed on a single H100 GPU (NVIDIA Corp., CA, USA) containing 96 GB of GPU memory, hosted in a system equipped with an AMD 4th Gen EPYC processor (18 cores, 36 threads) and 180 GB of system memory.

References

- [1] Hugo J. W. L. Aerts, Emmanuel Rios Velazquez, Ralph T. H. Leijenaar, Chintan Parmar, Patrick Grossmann, Sara Carvalho, Johan Bussink, René Monshouwer, Benjamin Haibe-Kains, Derek Rietveld, Frank Hoebbers, Michelle M. Rietbergen, C. René Leemans, Andre Dekker, John Quackenbush, Robert J. Gillies, and Philippe Lambin. Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach. *Nature Communications*, 5(1): 4006, 2014. 3
- [2] Md Mahfuz Al Hasan, Mahdi Zaman, Abdul Jawad, Alberto Santamaria-Pang, Ho Hin Lee, Ivan Tarapov, Kyle B. See, Md Shah Imran, Antika Roy, Yaser Pourmohammadi Fallah, Navid Asadizanjani, and Reza Forghani. WaveFormer: A 3D Transformer with Wavelet-Driven Feature Representation for Efficient Medical Image Segmentation. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2025*, pages 684–694. 2026. 5
- [3] Natália Alves, Megan Schuurmans, Dawid Rutkowski, Derya Yakar, Ingfrid Haldorsen, Marjolein Liedenbaum, Anders Molven, Pierpaolo Vendittelli, Geert Litjens, John Hermans, and Henkjan Huisman. The PANORAMA Study Protocol: Pancreatic Cancer Diagnosis - Radiologists Meet AI. Technical report, Zenodo, 2024. 1, 2
- [4] Shaimaa Bakr, Olivier Gevaert, Sebastian Echegaray, Kelsey Ayers, Mu Zhou, Majid Shafiq, Hong Zheng, Jalen Anthony Benson, Weiruo Zhang, Ann N. C. Leung, Michael Kadoch, Chuong D. Hoang, Joseph Shrager, Andrew Quon, Daniel L. Rubin, Sylvia K. Plevritis, and Sandy Napel. A radiogenomic dataset of non-small cell lung cancer. *Scientific Data*, 5:180202, 2018. 3
- [5] Louis Blankemeier, Joseph Paul Cohen, Ashwin Kumar, Dave Van Veen, Syed Jamal Safdar Gardezi, Magdalini Paschali, Zhihong Chen, Jean-Benoit Delbrouck, Eduardo Reis, Cesar Truys, Christian Bluethgen, Malte Engmann Kjeldskov Jensen, Sophie Ostmeier, Maya Varma, Jeya Maria Jose Valanarasu, Zhongnan Fang, Zepeng Huo, Zaid Nabulsi, Diego Ardila, Wei-Hung Weng, Edson Amaro Junior, Neera Ahuja, Jason Fries, Nigam H. Shah, Andrew Johnston, Robert D. Boutin, Andrew Wentland, Curtis P. Langlotz, Jason Hom, Sergios Gatidis, and Akshay S. Chaudhari. Merlin: A Vision Language Foundation Model for 3D Computed Tomography, 2024. 1, 3, 9
- [6] Noel C. F. Codella, Ying Jin, Shrey Jain, Yu Gu, Ho Hin Lee, Asma Ben Abacha, Alberto Santamaria-Pang, Will Guyman, Naiteek Sangani, Sheng Zhang, Hoifung Poon, Stephanie Hyland, Shruthi Bannur, Javier Alvarez-Valle, Xue Li, John Garrett, Alan McMillan, Gaurav Rajguru, Madhu Maddi, Nilesh Vijayrania, Rehaan Bhimai, Nick Mecklenburg, Rupal Jain, Daniel Holstein, Naveen Gaur, Vijay Aski, Jenq-Neng Hwang, Thomas Lin, Ivan Tarapov, Matthew Lungren, and Mu Wei. MedImageInsight: An Open-Source Embedding Model for General Domain Medical Imaging, 2024. 3
- [7] Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi, and Yonglong Tian. Improving CLIP Training with Language Rewrites. In *Advances in Neural Information Processing Systems*, 2023. 2
- [8] Ibrahim Ethem Hamamci, Sezgin Er, Furkan Almas, Ayse Gulnihhan Simsek, Sevval Nil Esirgun, Irem Dogan, Muhammed Furkan Dasdelen, Bastian Wittmann, Enis Simsar, Mehmet Simsar, Emine Benu Erdemir, Abdullah Alanbay, Anjany Sekuboyina, Berkan Lafci, Mehmet K. Ozdemir, and Bjoern Menze. A foundation model utilizing chest CT volumes and radiology reports for supervised-level zero-shot detection of abnormalities, 2024. 1, 3, 9, 10
- [9] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R. Roth, and Daguang Xu. UNETR: Transformers for 3D Medical Image Segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 574–584, 2022. 5
- [10] Yufan He, Vishwesh Nath, Dong Yang, Yucheng Tang, Andriy Myronenko, and Daguang Xu. SwinUNETR-V2: Stronger Swin Transformers with Stagewise Convolutions for 3D Medical Image Segmentation. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2023*, pages 416–426. 2023. 5
- [11] Yufan He, Pengfei Guo, Yucheng Tang, Andriy Myronenko, Vishwesh Nath, Ziyue Xu, Dong Yang, Can Zhao, Benjamin Simon, Mason Belue, Stephanie Harmon, Baris Turkbey, Daguang Xu, and Wenqi Li. VISTA3D: A Unified Segmentation Foundation Model For 3D Medical Imaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20863–20873, 2025. 3
- [12] Nicholas Heller, Fabian Isensee, Klaus H. Maier-Hein, Xi-aoshuai Hou, Chunmei Xie, Fengyi Li, Yang Nan, Guangrui Mu, Zhiyong Lin, Miofei Han, Guang Yao, Yaozong Gao, Yao Zhang, Yixin Wang, Feng Hou, Jiawei Yang, Guangwei Xiong, Jiang Tian, Cheng Zhong, Jun Ma, Jack Rickman, Joshua Dean, Bethany Stai, Resha Tejpal, Makinna Oestreich, Paul Blake, Heather Kaluzniak, Shaneabbas Raza, Joel Rosenberg, Keenan Moore, Edward Walczak, Zachary Rengel, Zach Edgerton, Ranveer Vasdev, Matthew Peterson, Sean McSweeney, Sarah Peterson, Arveen Kalapara, Niranjana Sathianathan, Nikolaos Papanikolopoulos, and Christopher Weight. The state of the art in kidney and kidney tumor segmentation in contrast-enhanced CT imaging: Results of the KiTS19 challenge. *Medical Image Analysis*, 67:101821, 2021. 4, 7
- [13] Shih-Cheng Huang, Zepeng Huo, Ethan Steinberg, Chia-Chun Chiang, Curtis Langlotz, Matthew Lungren, Serena Yeung, Nigam Shah, and Jason Fries. INSPECT: A Multimodal Dataset for Patient Outcome Prediction of Pulmonary Embolisms. *Advances in Neural Information Processing Systems*, 36:17742–17772, 2023. 1
- [14] Fabian Isensee, Paul F. Jaeger, Simon A. A. Kohl, Jens Petersen, and Klaus H. Maier-Hein. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, 18(2):203–211, 2021. 5
- [15] Fabian Isensee, Tassilo Wald, Constantin Ulrich, Michael Baumgartner, Saikat Roy, Klaus Maier-Hein, and Paul F. Jäger. nnU-Net Revisited: A Call for Rigorous Validation in 3D Medical Image Segmentation. In *Medical Image Com-*

- puting and Computer Assisted Intervention – MICCAI 2024, pages 488–498. 2024. 4, 5
- [16] Yuanfeng Ji, Haotian Bai, Chongjian Ge, Jie Yang, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhanng, Wanling Ma, Xiang Wan, and Ping Luo. AMOS: A Large-Scale Abdominal Multi-Organ Benchmark for Versatile Medical Image Segmentation. In *Advances in Neural Information Processing Systems*, 2022. 1
- [17] Tommie Kerssies, Niccolò Cavagnero, Alexander Hermans, Narges Norouzi, Giuseppe Averta, Bastian Leibe, Gijs Dubbelman, and Daan de Geus. Your ViT is Secretly an Image Segmentation Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25303–25313, 2025. 4
- [18] Wenhui Lei, Hanyu Chen, Zitian Zhang, Luyang Luo, Qiong Xiao, Yannian Gu, Peng Gao, Yankai Jiang, Ci Wang, Guangtao Wu, Tongjia Xu, Yingjie Zhang, Pranav Rajpurkar, Xiaofan Zhang, Shaoting Zhang, and Zhenning Wang. A Synthetic Data-Driven Radiology Foundation Model for Pan-tumor Clinical Diagnosis, 2025. 3
- [19] Wenxuan Li, Chongyu Qu, Xiaoxi Chen, Pedro R. A. S. Bassi, Yijia Shi, Yuxiang Lai, Qian Yu, Huimin Xue, Yixiong Chen, Xiaorui Lin, Yutong Tang, Yining Cao, Haoqi Han, Zheyuan Zhang, Jiawei Liu, Tiezheng Zhang, Yujie Ma, Jincheng Wang, Guang Zhang, Alan Yuille, and Zongwei Zhou. AbdomenAtlas: A large-scale, detailed-annotated, & multi-center dataset for efficient transfer learning and open algorithmic benchmarking. *Medical Image Analysis*, 97:103285, 2024. 1
- [20] Wenxuan Li, Alan Yuille, and Zongwei Zhou. How Well Do Supervised 3D Models Transfer to Medical Imaging Tasks?, 2025. 3
- [21] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *Proceedings of the International Conference on Learning Representations*, 2018. 5
- [22] Jun Ma, Yao Zhang, Song Gu, Cheng Zhu, Cheng Ge, Yichi Zhang, Xingle An, Congcong Wang, Qiyuan Wang, Xin Liu, Shucheng Cao, Qi Zhang, Shangqing Liu, Yunpeng Wang, Yuhui Li, Jian He, and Xiaoping Yang. AbdomenCT-1K: Is Abdominal Organ Segmentation a Solved Problem? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6695–6714, 2022. 1, 2
- [23] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. UMAP: Uniform Manifold Approximation and Projection. *Journal of Open Source Software*, 3(29):861, 2018. 10
- [24] Suraj Pai, Dennis Bontempi, Ibrahim Hadzic, Vasco Prudente, Mateo Sokač, Tafadzwa L. Chaunzwa, Simon Bernatz, Ahmed Hosny, Raymond H. Mak, Nicolai J. Birkbak, and Hugo J. W. L. Aerts. Foundation model for cancer imaging biomarkers. *Nature Machine Intelligence*, pages 1–14, 2024. 3
- [25] Suraj Pai, Ibrahim Hadzic, Dennis Bontempi, Keno Bresslem, Benjamin H. Kann, Andriy Fedorov, Raymond H. Mak, and Hugo J. W. L. Aerts. Vision Foundation Models for Computed Tomography, 2025. 3
- [26] Suraj Pai, Ibrahim Hadzic, Andrey Fedorov, Raymond H. Mak, and Hugo JWL Aerts. Foundation model embeddings for quantitative tumor imaging biomarkers, 2025. 3, 4
- [27] Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau, Justin Chen, Fereshteh Mahvar, Liron Yatziv, Tiffany Chen, Bram Sterling, Stefanie Anna Baby, Susanna Maria Baby, Jeremy Lai, Samuel Schmidgall, Lu Yang, Kejia Chen, Per Bjornsson, Shashir Reddy, Ryan Brush, Kenneth Philbrick, Mercy Asiedu, Ines Mezerreg, Howard Hu, Howard Yang, Richa Tiwari, Sunny Jansen, Preeti Singh, Yun Liu, Shekoofeh Azizi, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Riviere, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Elena Buchatskaya, Jean-Baptiste Alayrac, Dmitry Lepikhin, Vlad Feinberg, Sebastian Borgeaud, Alek Andreev, Cassidy Hardin, Robert Dadashi, Léonard Hussenot, Armand Joulin, Olivier Bachem, Yossi Matias, Katherine Chou, Avinatan Hassidim, Kavi Goel, Clement Farabet, Joelle Barral, Tris Warkentin, Jonathon Shlens, David Fleet, Victor Cotruta, Omar Sanseviero, Gus Martins, Phoebe Kirk, Anand Rao, Shravya Shetty, David F. Steiner, Can Kirmizibayrak, Rory Pilgrim, Daniel Golden, and Lin Yang. MedGemma Technical Report, 2025. 9
- [28] Arnaud Arindra Adiyoso Setio, Alberto Traverso, Thomas de Bel, Moira S. N. Berens, Cas van den Bogaard, Piergiorgio Cerello, Hao Chen, Qi Dou, Maria Evelina Fantacci, Bram Geurts, Robbert van der Gugten, Pheng Ann Heng, Bart Jansen, Michael M. J. de Kaste, Valentin Kotov, Jack Yu-Hung Lin, Jeroen T. M. C. Manders, Alexander Sónora-Mengana, Juan Carlos García-Naranjo, Evgenia Papavasileiou, Mathias Prokop, Marco Saletta, Cornelia M Schaefer-Prokop, Ernst T. Scholten, Luuk Scholten, Miranda M. Snoeren, Ernesto Lopez Torres, Jef Vandemeulebroucke, Nicole Walasek, Guido C. A. Zuidhof, Bram van Ginneken, and Colin Jacobs. Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: The LUNA16 challenge. *Medical Image Analysis*, 42:1–13, 2017. 3
- [29] Amber L. Simpson, Jacob Peoples, John M. Creasy, Gabor Fichtinger, Natalie Gangai, Krishna N. Keshavamurthy, Andras Lasso, Jinru Shia, Michael I. D’Angelica, and Richard K. G. Do. Preoperative CT and survival data for patients undergoing resection of colorectal liver metastases. *Scientific Data*, 11(1):172, 2024. 3, 4
- [30] The NLST Research Team. Reduced Lung-Cancer Mortality with Low-Dose Computed Tomographic Screening. *New England Journal of Medicine*, 365(5):395–409, 2011. 1
- [31] Tassilo Wald, Saikat Roy, Fabian Isensee, Constantin Ulrich, Sebastian Ziegler, Dasha Trofimova, Raphael Stock, Michael Baumgartner, Gregor Köhler, and Klaus H. Maier-Hein. Primus: Enforcing Attention Usage for 3D Medical Image Segmentation. *CoRR*, 2025. 5
- [32] Tassilo Wald, Constantin Ulrich, Jonathan Suprijadi, Sebastian Ziegler, Michal Nohel, Robin Peretzke, Gregor Kohler,

- and Klaus Maier-Hein. An OpenMind for 3D Medical Vision Self-supervised Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23839–23879, 2025. 4
- [33] Avivah J. Wang, Fakrul Islam Tushar, Michael R. Harowicz, Betty C. Tong, Kyle J. Lafata, Tina D. Tailor, and Joseph Y. Lo. The Duke Lung Cancer Screening (DLCS) Dataset: A Reference Dataset of Annotated Low-Dose Screening Thoracic CT. *Radiology: Artificial Intelligence*, 7(4):e240248, 2025. 3
- [34] Jakob Wasserthal, Hanns-Christian Breit, Manfred T. Meyer, Maurice Pradella, Daniel Hinck, Alexander W. Sauter, Tobias Heye, Daniel T. Boll, Joshy Cyriac, Shan Yang, Michael Bach, and Martin Segeroth. TotalSegmentator: Robust Segmentation of 104 Anatomic Structures in CT Images. *Radiology: Artificial Intelligence*, 5(5):e230024, 2023. 1
- [35] Linshan Wu, Jiaxin Zhuang, and Hao Chen. VoCo: A Simple-yet-Effective Volume Contrastive Learning Framework for 3D Medical Image Analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22873–22882, 2024. 3
- [36] Yutong Xie, Jianpeng Zhang, Chunhua Shen, and Yong Xia. CoTr: Efficiently Bridging CNN and Transformer for 3D Medical Image Segmentation. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*, pages 171–180. 2021. 5
- [37] Hong-Yu Zhou, Jiansen Guo, Yinghao Zhang, Xiaoguang Han, Lequan Yu, Liansheng Wang, and Yizhou Yu. nnFormer: Volumetric Medical Image Segmentation via a 3D Transformer. *IEEE Transactions on Image Processing*, 32: 4036–4045, 2023. 5
- [38] Zongwei Zhou, Vatsal Sodha, Jiaxuan Pang, Michael B. Gotway, and Jianming Liang. Models Genesis. *Medical Image Analysis*, 67:101840, 2021. 3