



Enhancing Swin Transformer With Semantic Attention For Explainable Prediction

Aneesh Rangnekar, Jue Jiang, Harini Veeraraghavan

July 21, 2024

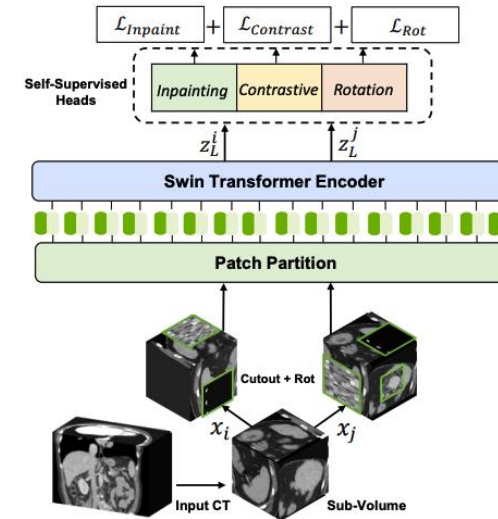
Session Title: Task-Based AI



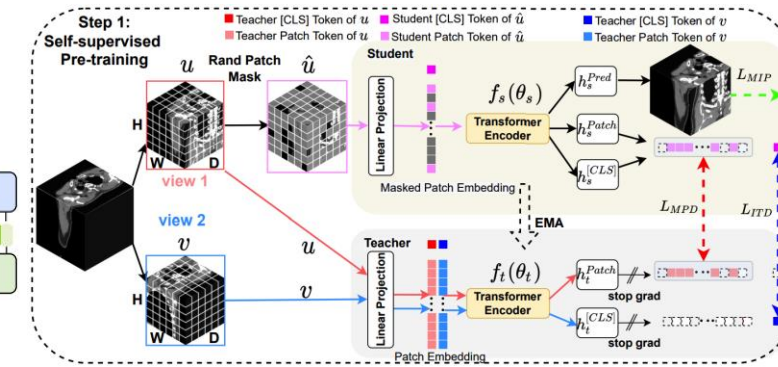
Memorial Sloan Kettering
Cancer Center

Motivation

- State-of-the-art foundation models in for 3D medical imaging
 - Swin UNETR
 - SMIT
 are based on the Swin transformer backbone
- However, the underlying architecture
 - cannot directly visualize the specific areas of interest
 - within the image scans
- We facilitate **visualization** of the regions of interest by modifying the architecture, thereby enabling eXplainable AI

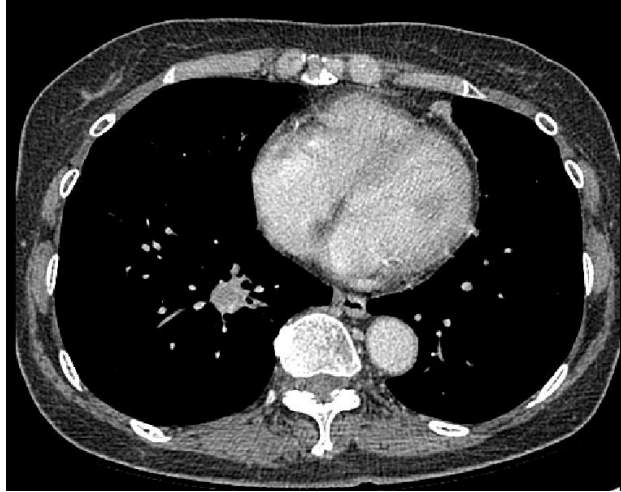


Swin UNETR



SMIT

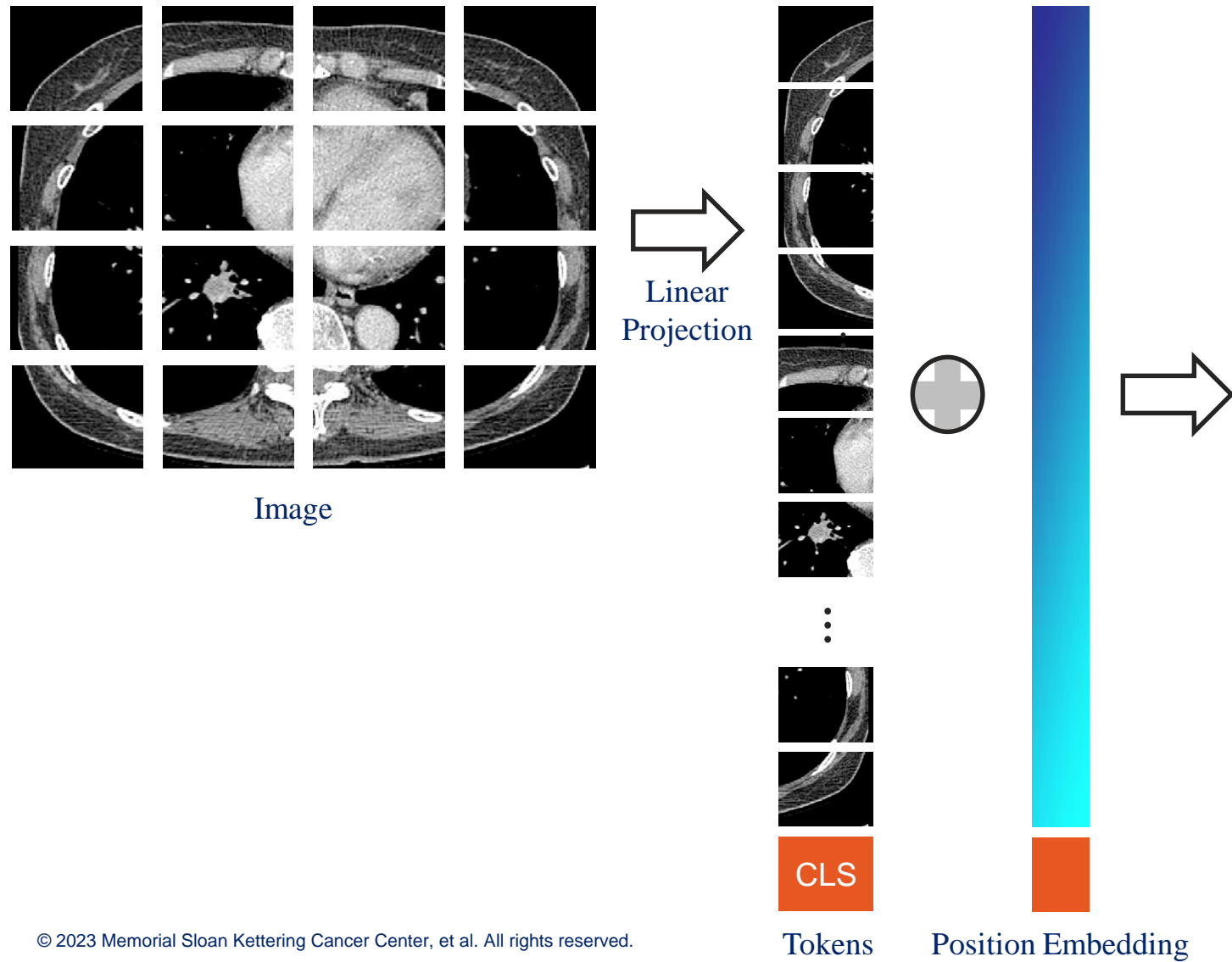
How do [Vision] transformers work?



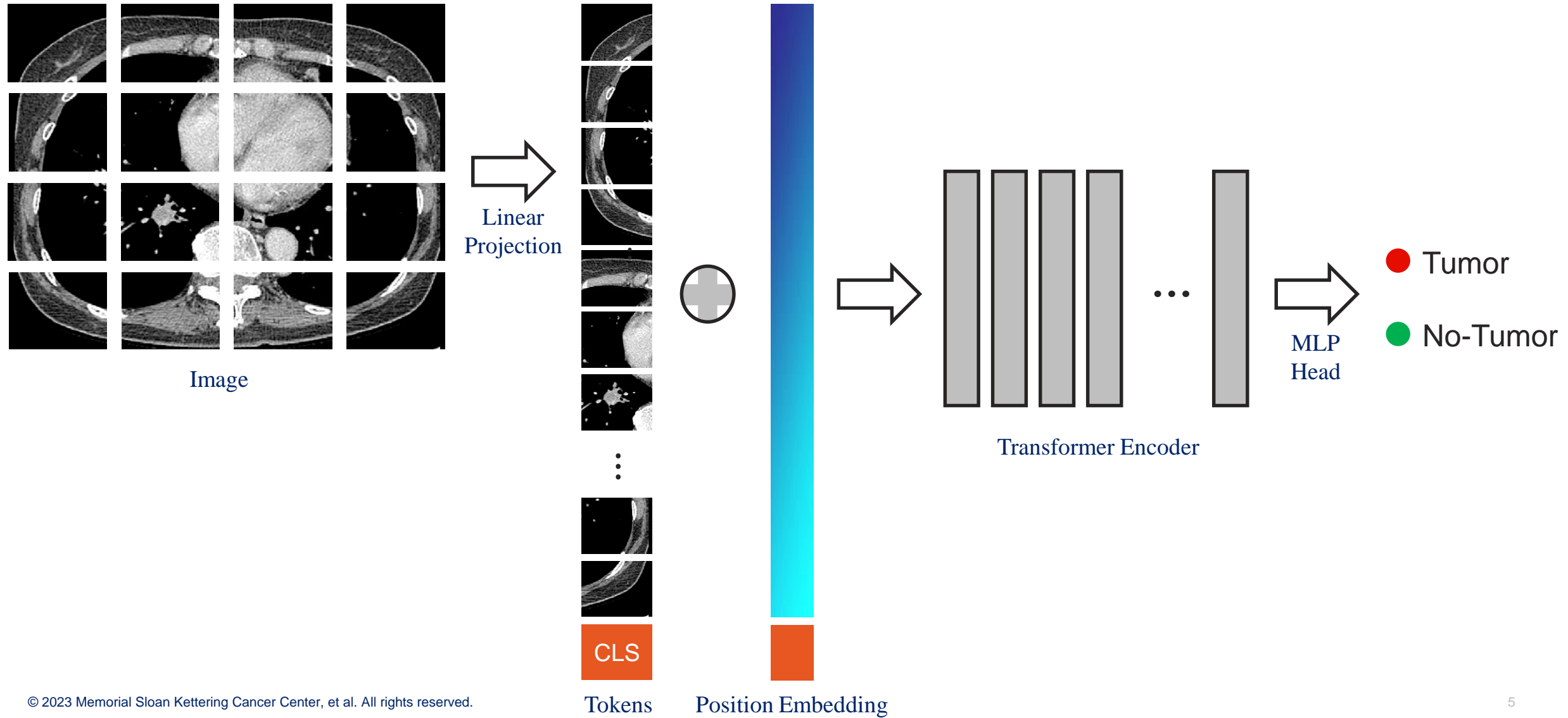
Image*

*Slice from 3D Volume for simplicity

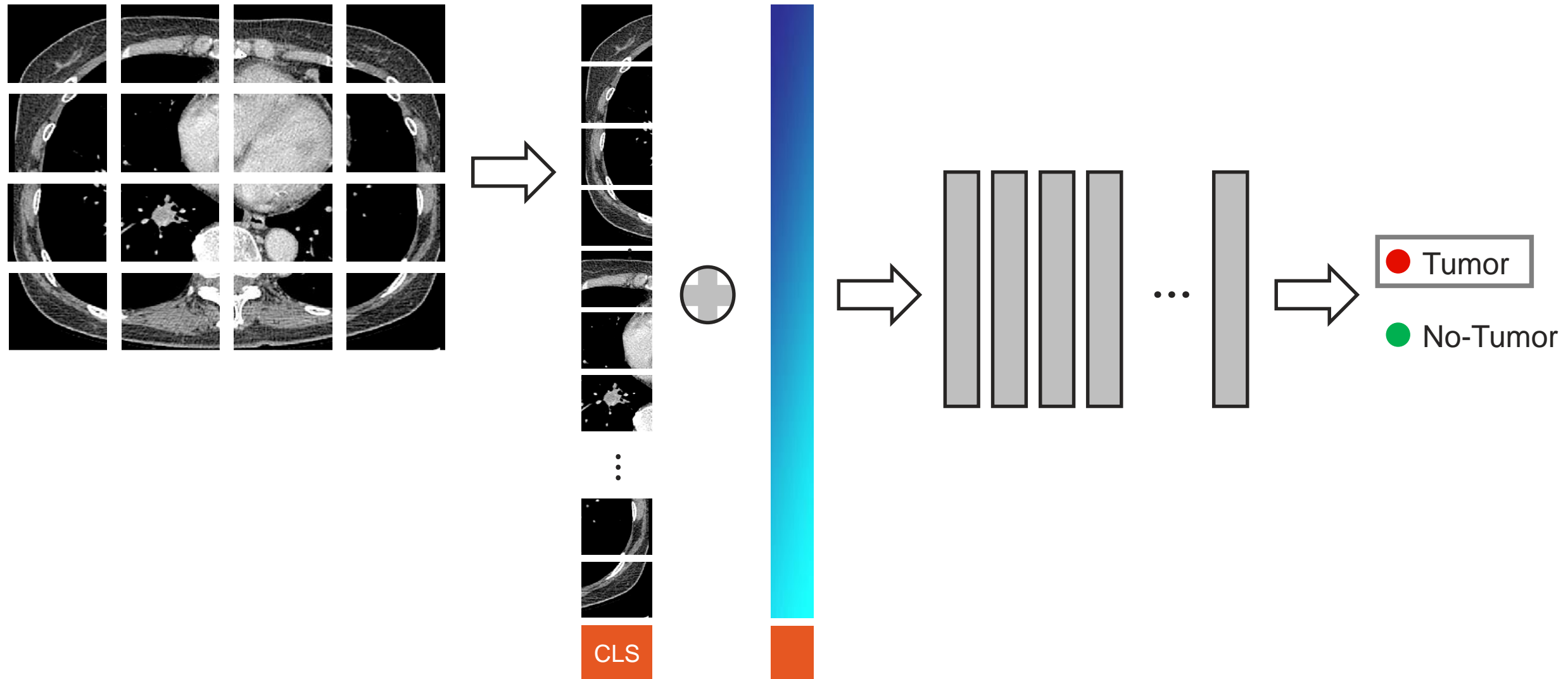
How do [Vision] transformers work?



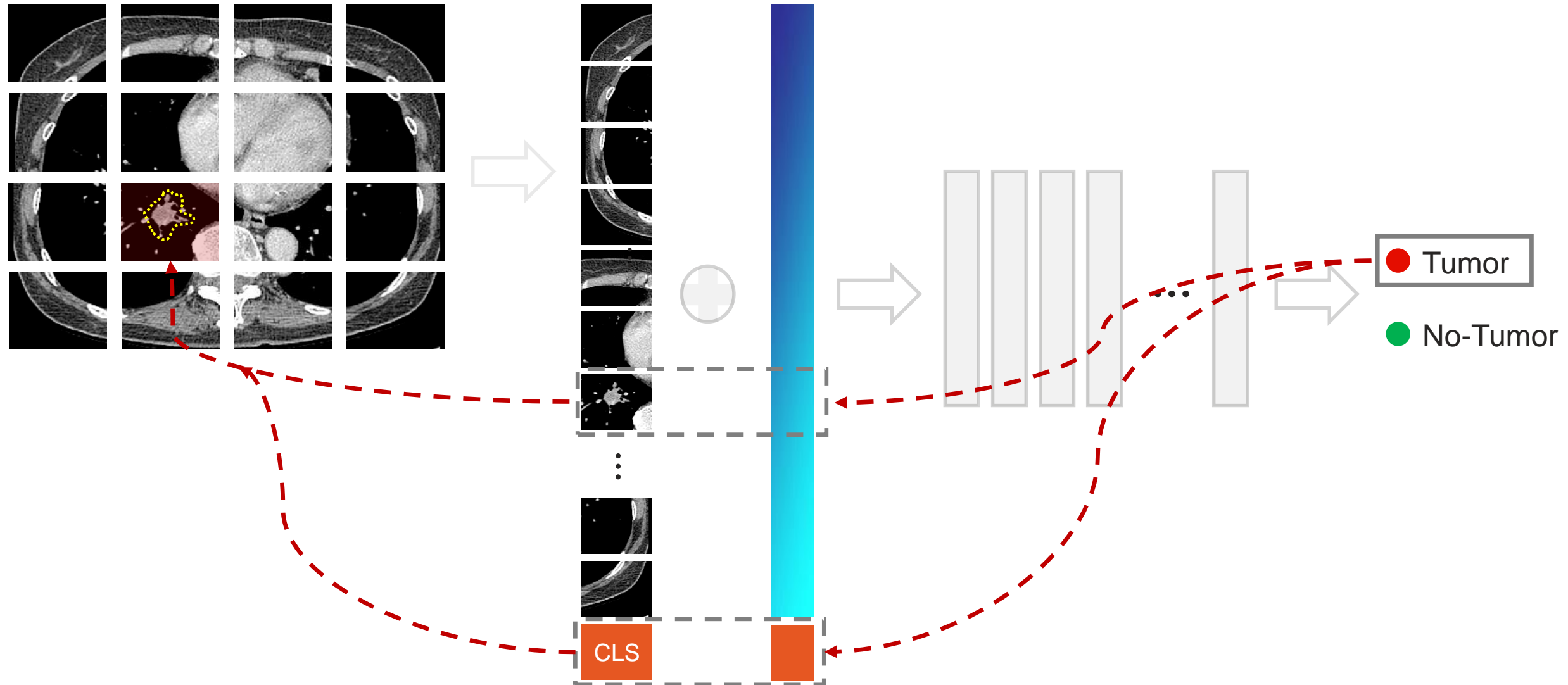
How do [Vision] transformers work?



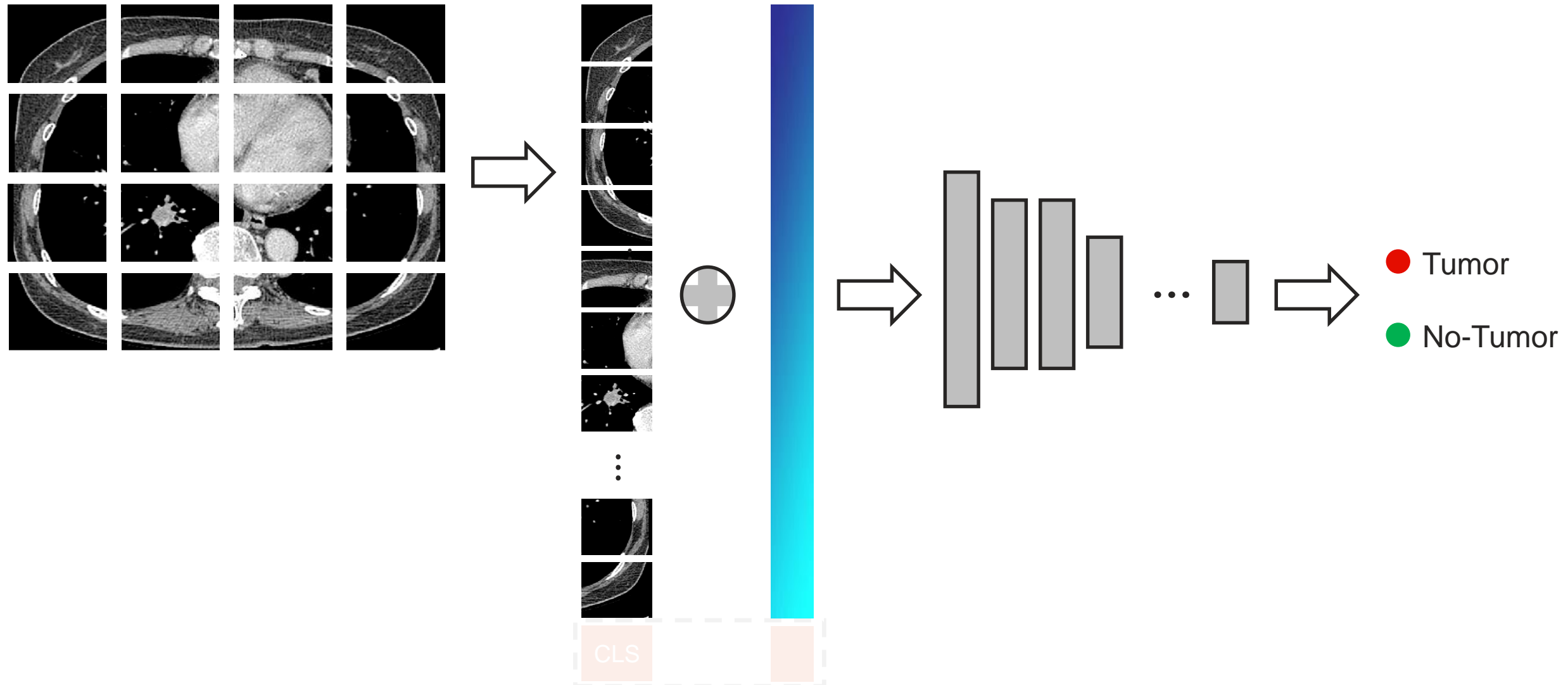
How do [Vision] transformers work?



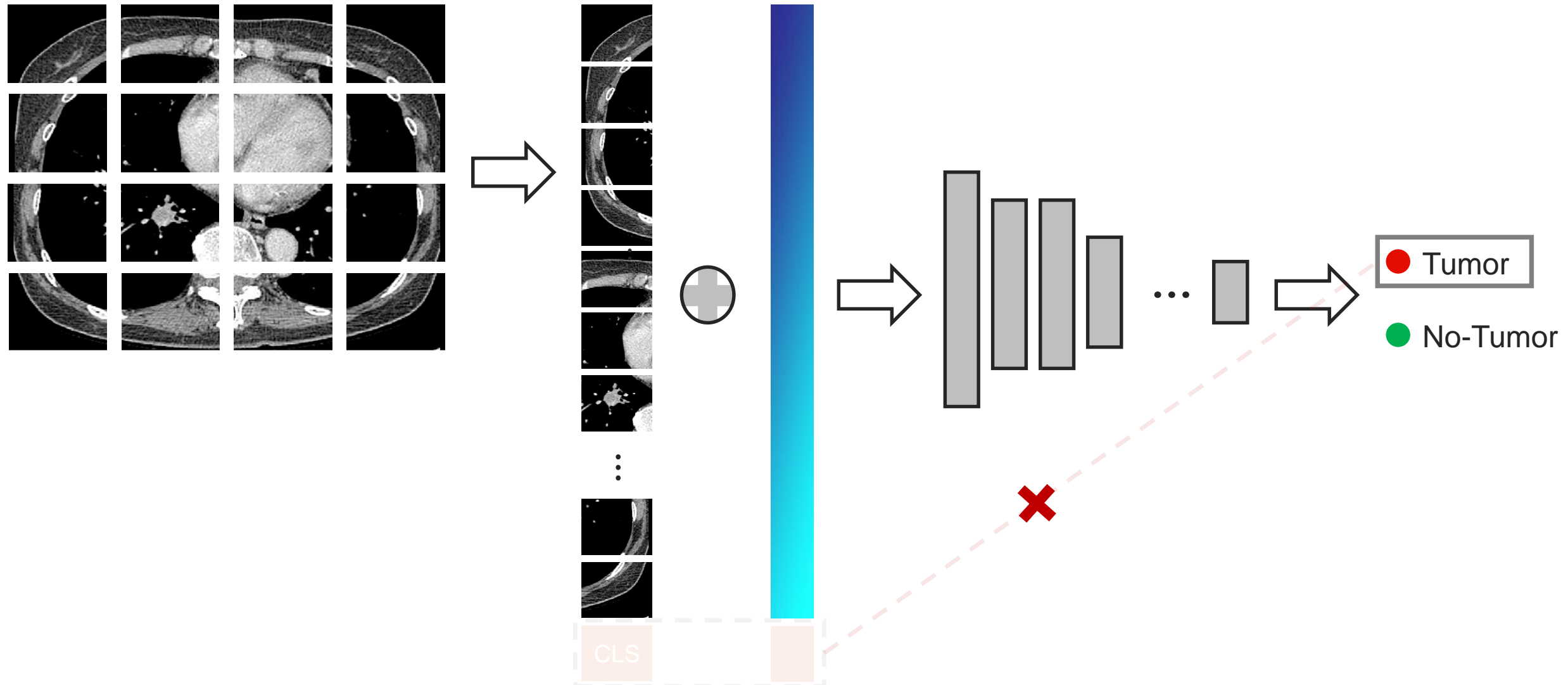
How do [Vision] transformers work?



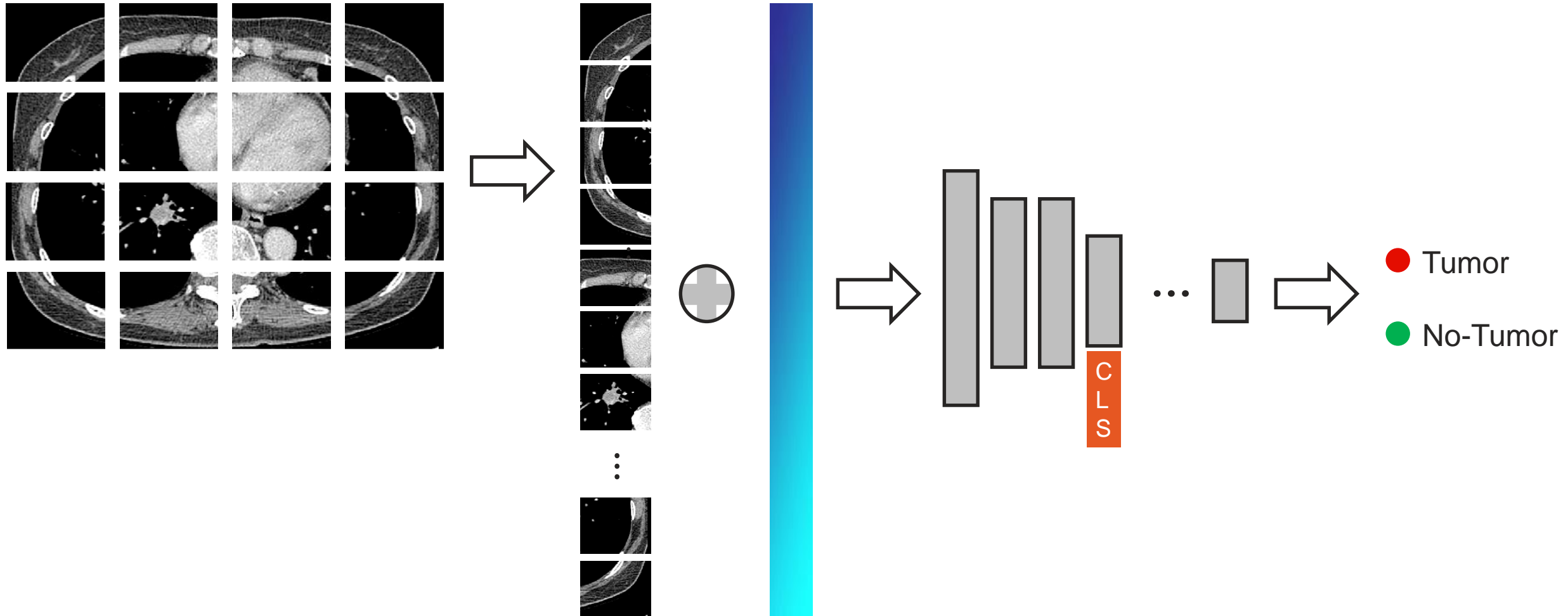
How do [Swin] transformers work?



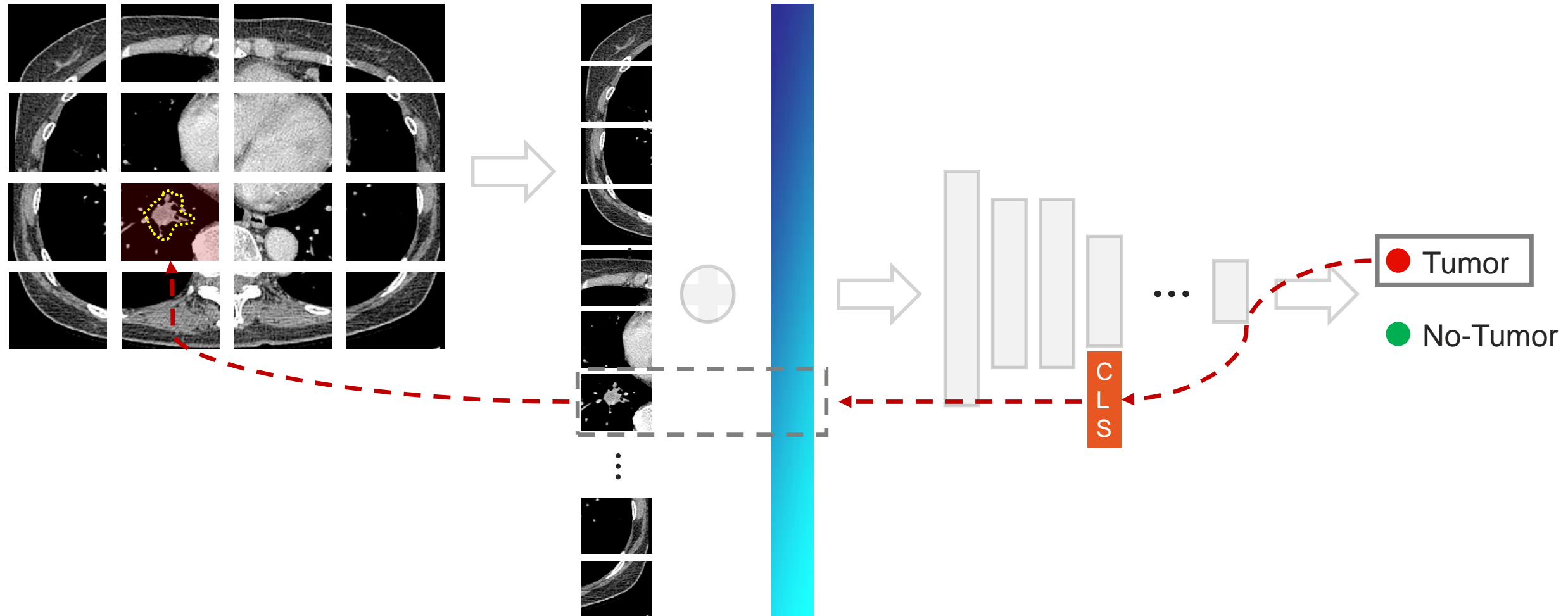
How do [Swin] transformers work?



How do we enhance the [Swin] transformer?

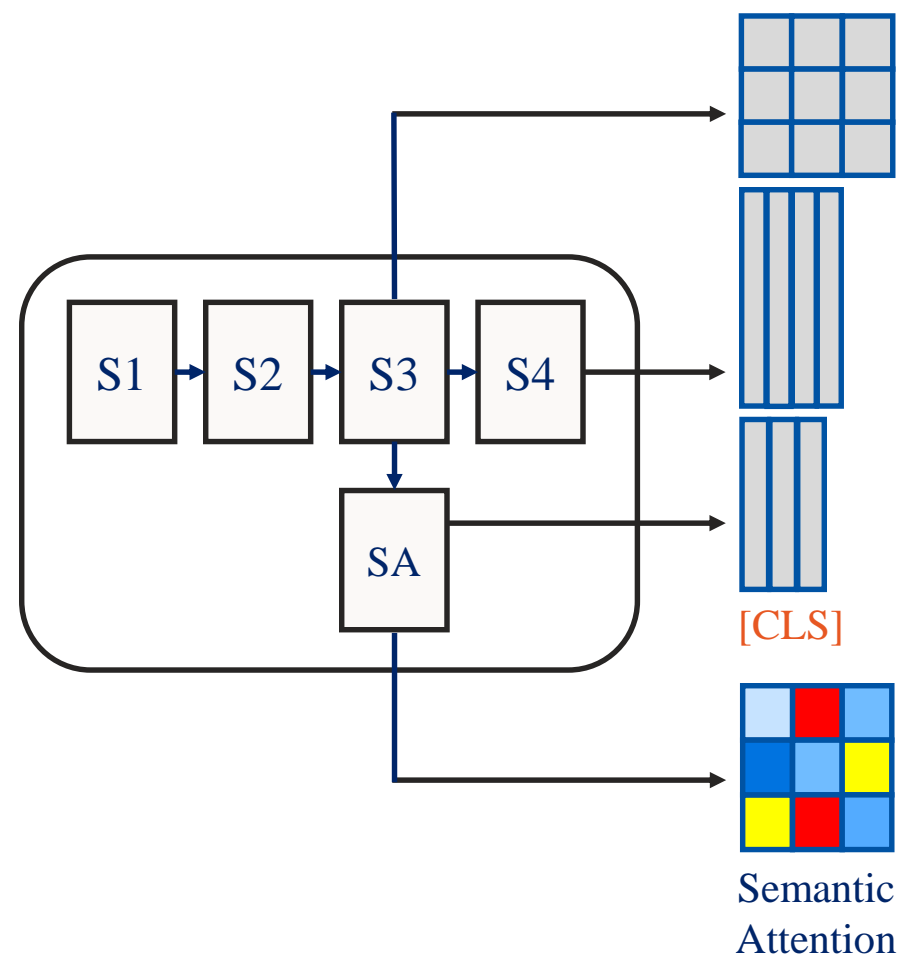


How do we enhance the [Swin] transformer?



Foundation Model: SMART

[S]elf-distilled [M]asked [A]ttention guided masked image modeling with noise [R]egularized [T]eacher



*built on SMIT



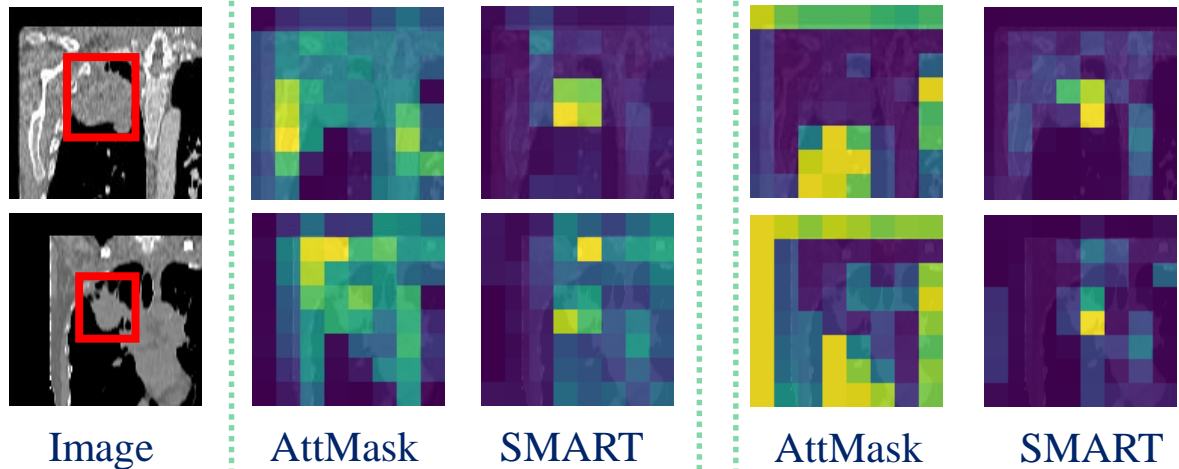
Results

We compare SMART against two other approaches:

- AttMask
- SMIT

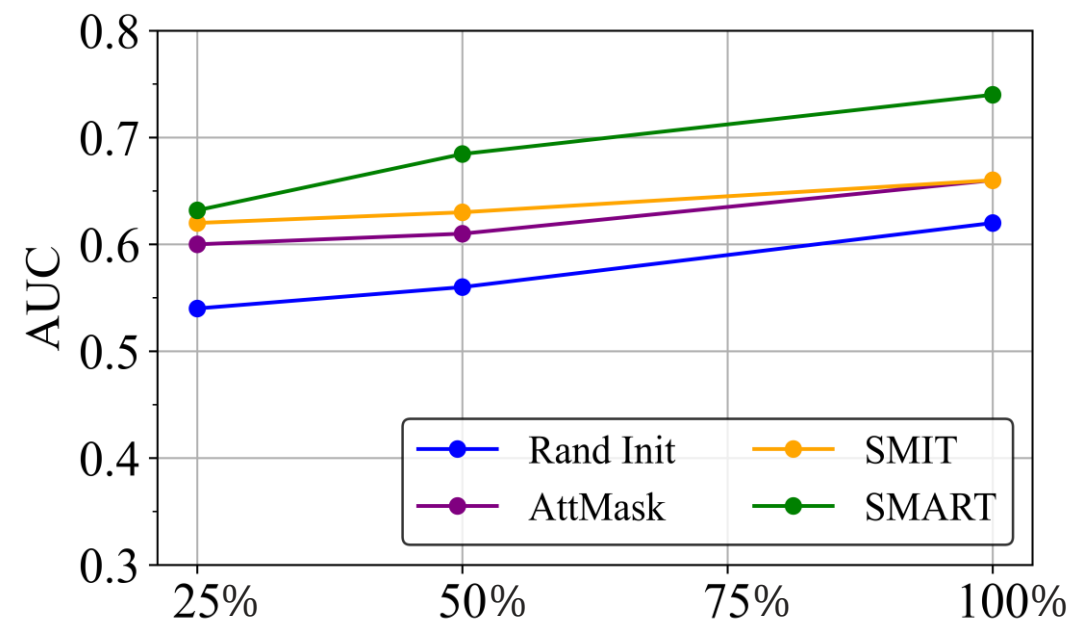
Pretext Task	Linear Probing			Fine-tuning		
	AP ₅₀	AR ₅₀	AUC	AP ₅₀	AR ₅₀	AUC
AttMask	44.9	54.0	0.570	56.0	68.5	0.660
SMIT	46.4	58.1	0.620	56.3	67.0	0.660
SMART	54.5	68.4	0.660	57.4	71.7	0.740

Durable Clinical Benefit Prediction with 100% data utilization



Prediction Attention Visualization post Fine-Tuning

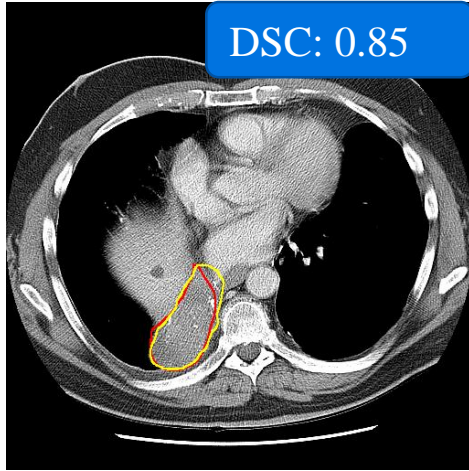
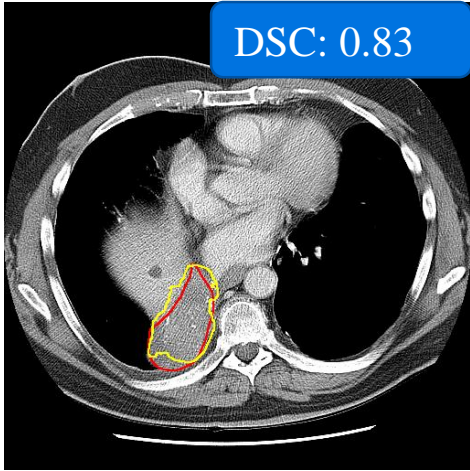
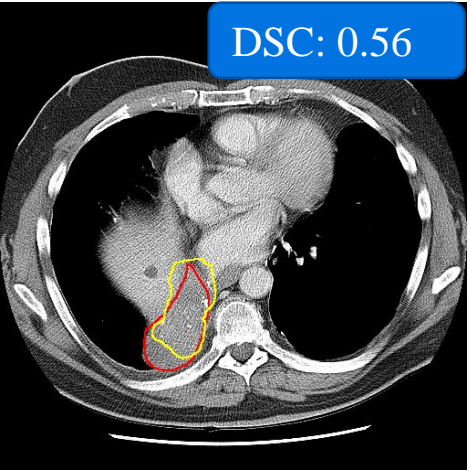
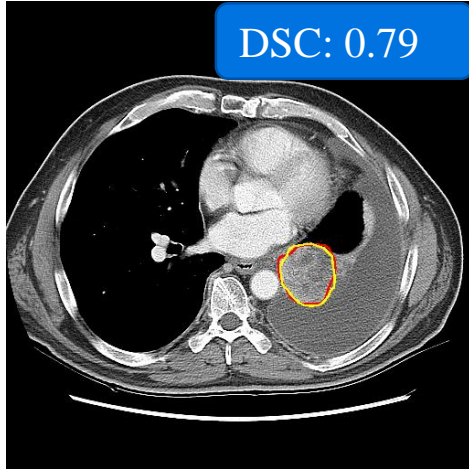
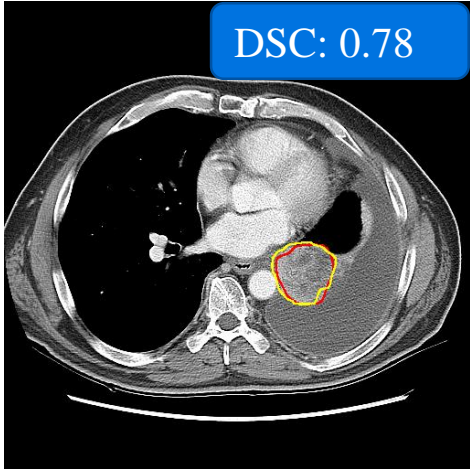
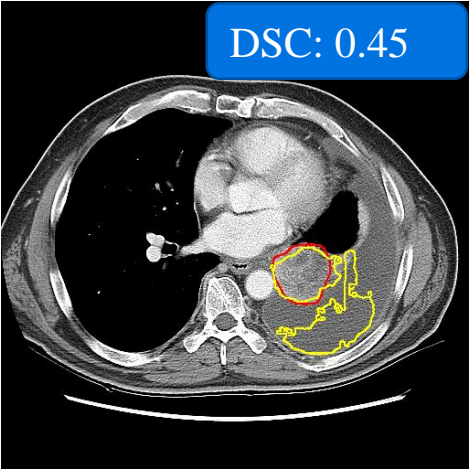
Results



Durable Clinical Benefit Prediction with limited-data learning protocol

Results

Tumor Segmentation	
AttMask	0.69 ± 0.21
SMIT	0.76 ± 0.13
SMART	0.77 ± 0.11



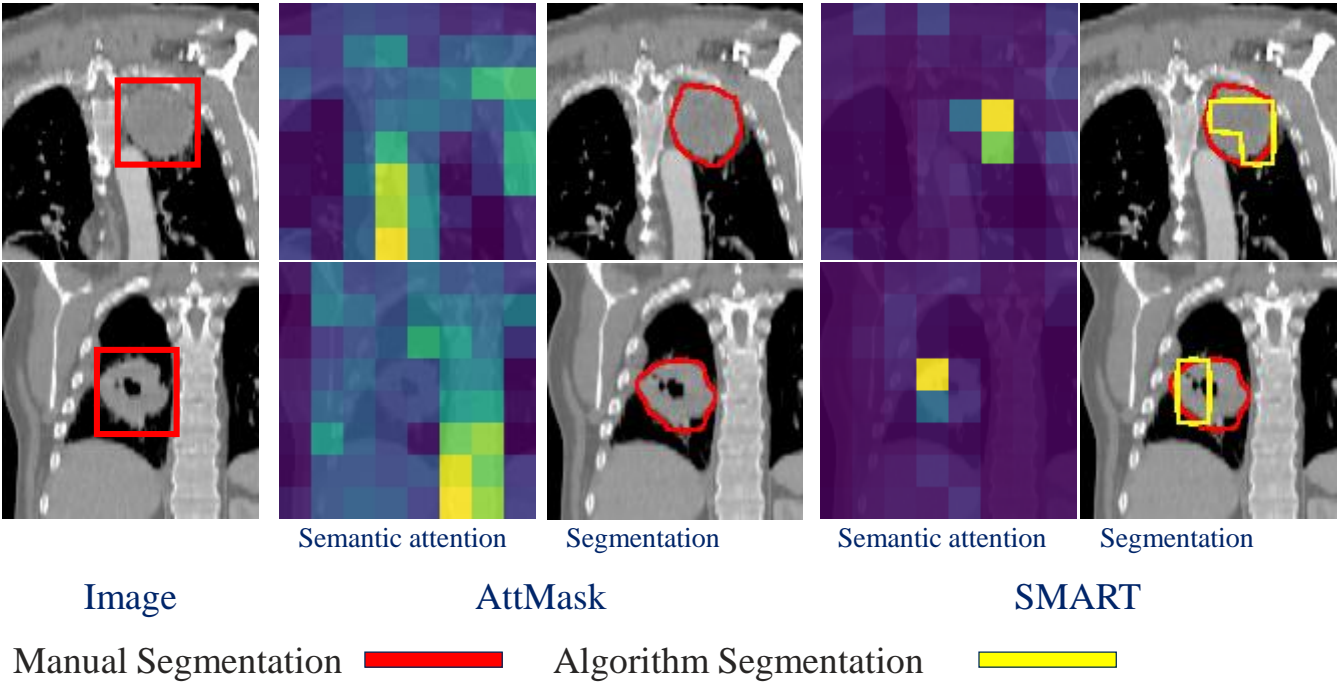
AttMask

SMIT

SMART

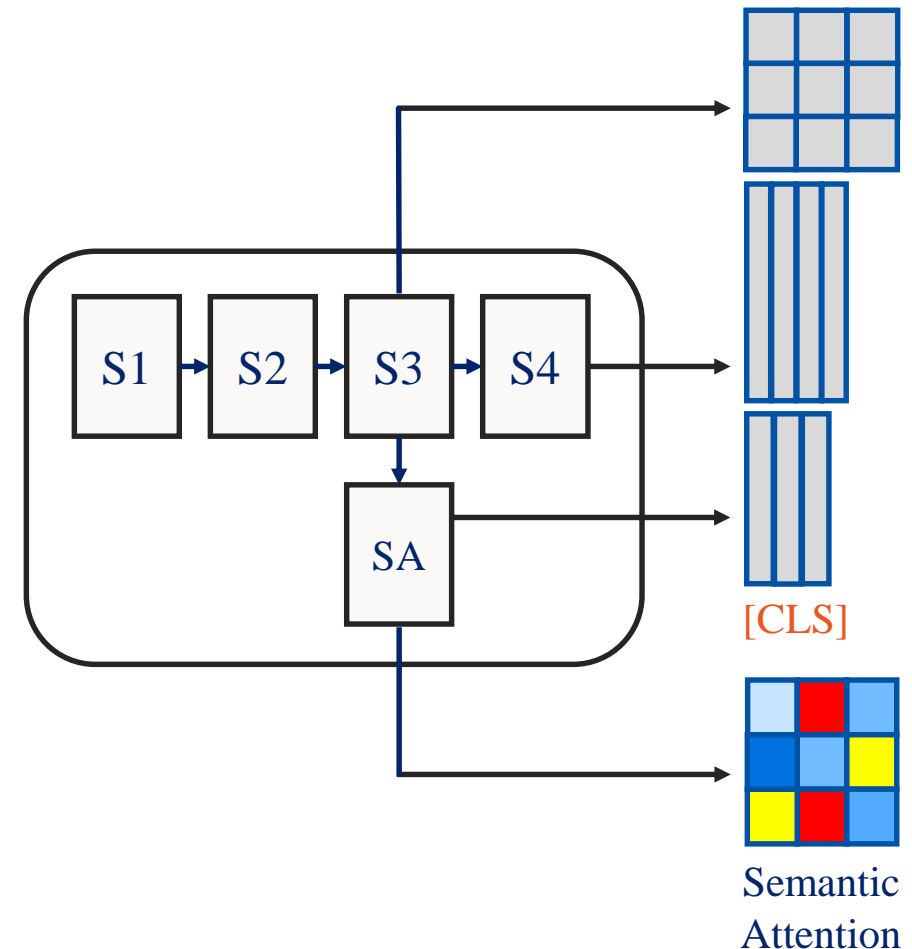
Results

Additional Results for zero-shot localization on held-out TCIA dataset



Summary

- We introduced a semantic attention block that enables CLS token
 - Helps pretraining the network with informative token masking
 - Helps post fine-tuning XAI
- Further boosts fine-tuning performance under limited data
- Also applied to
 - LIDC dataset for tumor malignancy classification and segmentation
 - TCIA dataset for zero-shot tumor localization
 - Lung Radiomics and Radiogenomics dataset for tumor segmentation
 - on arXiv: 2310.01209



Thank you

Questions?