

Interpretable Models for Segmentation and Classification of Lung Cancer Histological Subtypes

Gal Schneider & Tal Aish, Bar Ilan University

Deep-learning models in digital pathology often operate as black boxes, providing diagnostic labels without explanations grounded in tissue morphology.

We develop an interpretable, concept-based classification pipeline to distinguish non-small-cell lung cancer subtypes—adenocarcinoma (ADC) and squamous-cell carcinoma (SCC)—from haematoxylin-and-eosin (H&E) regions of interest (ROIs).

Utilizing the multi-institutional IGNITE dataset, we establish a leakage-free training environment split strictly at the patient level. Our architecture utilizes a frozen foundation model encoder (Phikon-v2) to extract spatial patch tokens. These tokens are mapped by a non-linear predictor into explicit proportions of high-level, human-understandable tissue concepts (such as stroma, necrosis, and tumor epithelium) derived from segmentation masks. Finally, a logistic regression classifier forms the diagnostic prediction over these intermediate concept proportions, allowing per-concept contribution decomposition and regional explanation maps.

Our concept-based pipeline achieves competitive patient-level Macro-F1 performance, outperforming standard direct image baselines which offer no tissue-grounded context. Remarkably, the pipeline maintains a strong dense concept fidelity (mean Pearson correlation) against expert annotations. A parallel branch derived from a self-configuring nnU-Net reaches higher spatial fidelity but yields a lower classification score, exposing a clear trade-off between segmentation accuracy and diagnostic subtype classification.

This proof-of-concept pipeline demonstrates that integrating local tissue structure preserves diagnostic capacity while surfacing clinically interpretable semantic and regional attribution maps.