

AN EXPLAINABLE VISION TRANSFORMER FRAMEWORK FOR INTERPRETABLE DISEASE CLASSIFICATION IN MEDICAL IMAGES

PRABAGARANE GANAPADY^{1*}, SUGANYA V.², A. MARTIN³ AND PREMA MAYUDU⁴

Advancements in medical image analysis through deep learning have significantly enhanced disease classification, but the black box nature of conventional neural networks remains a challenge for clinical trust. This work addresses interpretability with Vision Transformers (ViTs) in an explainable manner. The vision transformer-based model uses a self-attention mechanism for long-range spatial dependencies on medical image data. Gradient-weighted class activation method Grad-CAM for overlay saliency maps that pinpoint image regions in which the most important areas for the classification decision. The outcomes substantiate the efficacy of explainable ViT frameworks, ensure safe, trustworthy, and transparent artificial intelligence in healthcare diagnosis.
