

# A Multicenter Benchmark of Multiple Instance Learning Models for Lymphoma Subtyping from HE-stained Whole Slide Images

## Supplementary Material

Rao Muhammad Umer<sup>1</sup> Daniel Sens<sup>1</sup> Jonathan Noll<sup>1</sup> Sohom Dey<sup>1</sup> Christian Matek<sup>1,2,7</sup> Lukas Wolfseher<sup>8</sup> Rainer Spang<sup>8</sup> Ralf Huss<sup>9</sup>  
Johannes Raffler<sup>9</sup> Sarah Reinke<sup>10</sup> Ario Sadafi<sup>1,4,5</sup> Wolfram Klapper<sup>10</sup> Katja Steiger<sup>6</sup> Kristina Schwamborn<sup>6</sup> Carsten Marr<sup>1,2,3,4</sup>

<sup>1</sup> Institute of AI for Health, Helmholtz Munich, Neuherberg, Germany

<sup>2</sup> Department of Medicine III, Ludwig-Maximilian-University Hospital, Munich, Germany

<sup>3</sup> German Cancer Consortium (DKTK), Partner Site Munich, Germany

<sup>4</sup> Munich Center for Machine Learning (MCML), Munich, Germany

<sup>5</sup> Computer Aided Medical Procedures, Technical University of Munich, Munich, Germany

<sup>6</sup> Institute of Pathology, Technical University of Munich, School of Medicine and Health, Munich, Germany

<sup>7</sup> Institute of Pathology, Erlangen, Germany

<sup>8</sup> University of Regensburg, Regensburg, Germany

<sup>9</sup> Institute for Digital Medicine, University Hospital, Augsburg, Germany

<sup>10</sup> Department of Pathology, Hematopathology Section and Lymph Node Registry,

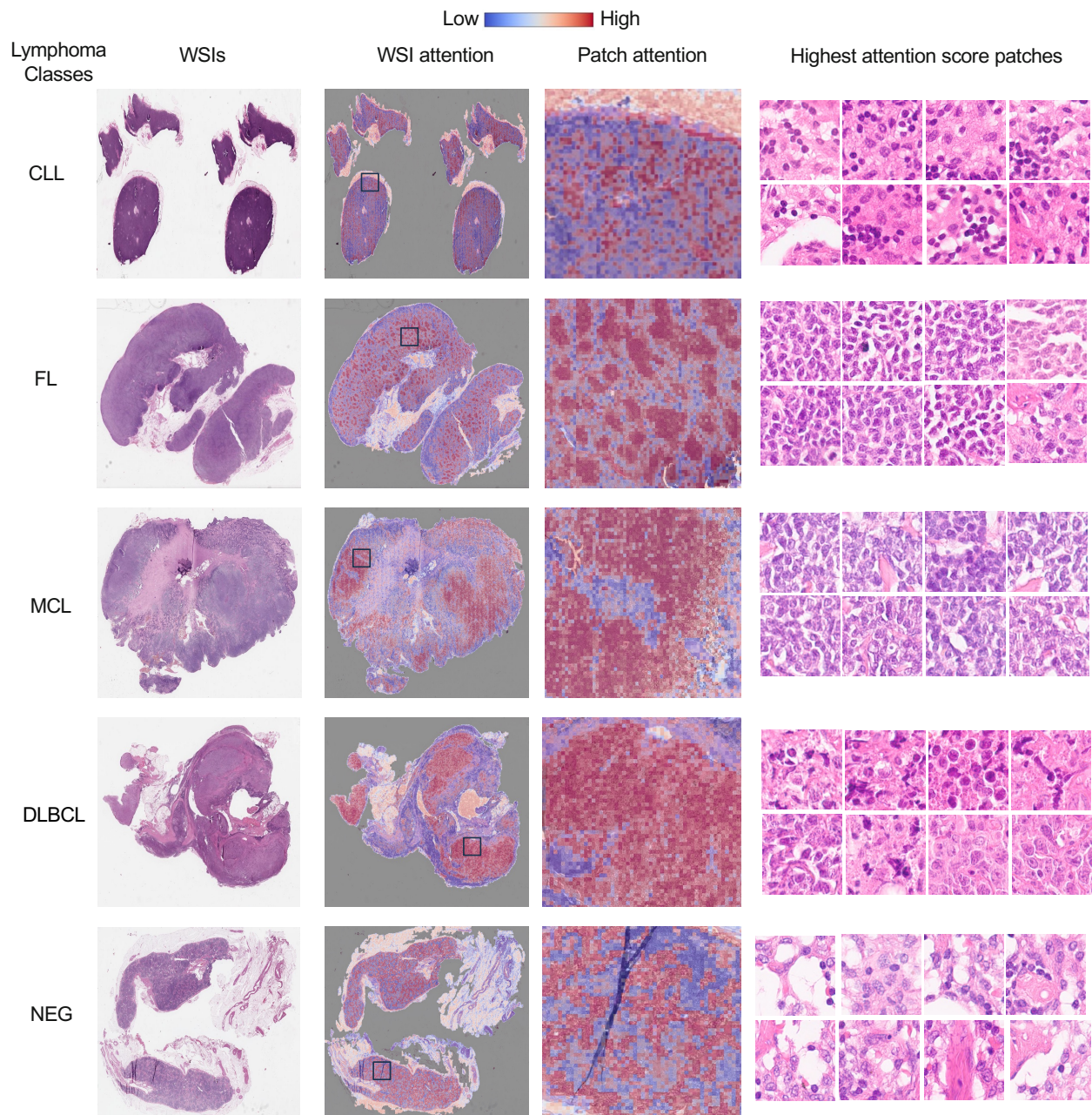
University Hospital Schleswig-Holstein, Kiel, Germany

carsten.marr@helmholtz-munich.de

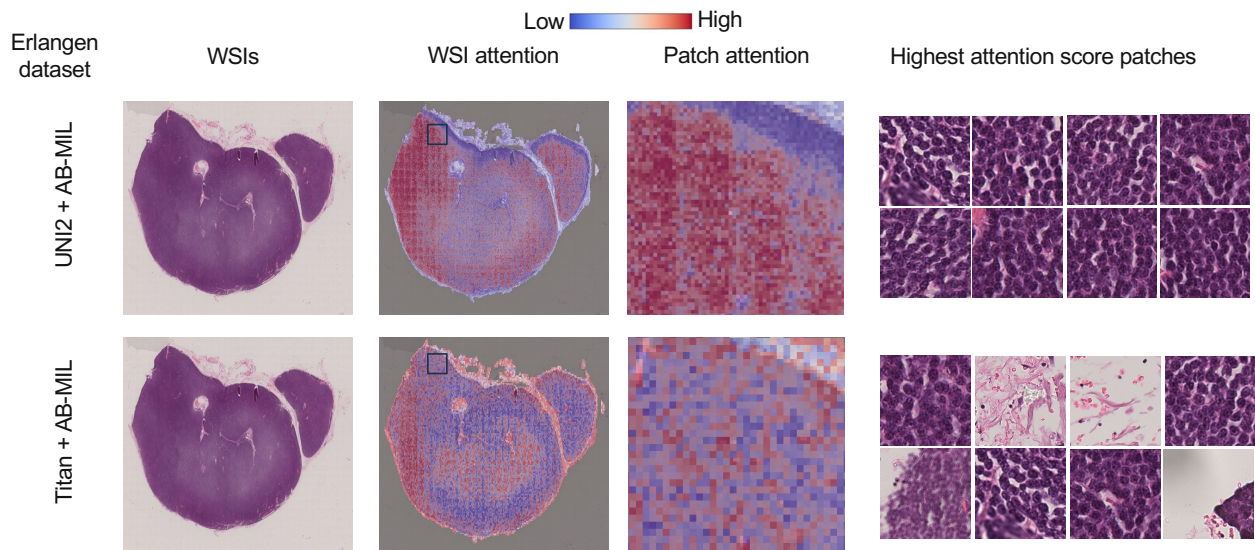
### 1. Whole slide image attention visualization

Attention scores can be visualized as heatmaps to highlight diagnostically informative regions, i.e., areas assigned high attention, while de-emphasizing regions of low relevance, such as normal tissue or background artifacts. To interpret the spatial distribution of model attention across a WSI, we convert the attention scores corresponding to the predicted class into percentiles and map these normalized values back to their original coordinates on the slide. Fine-grained heatmaps are produced by extracting overlapping

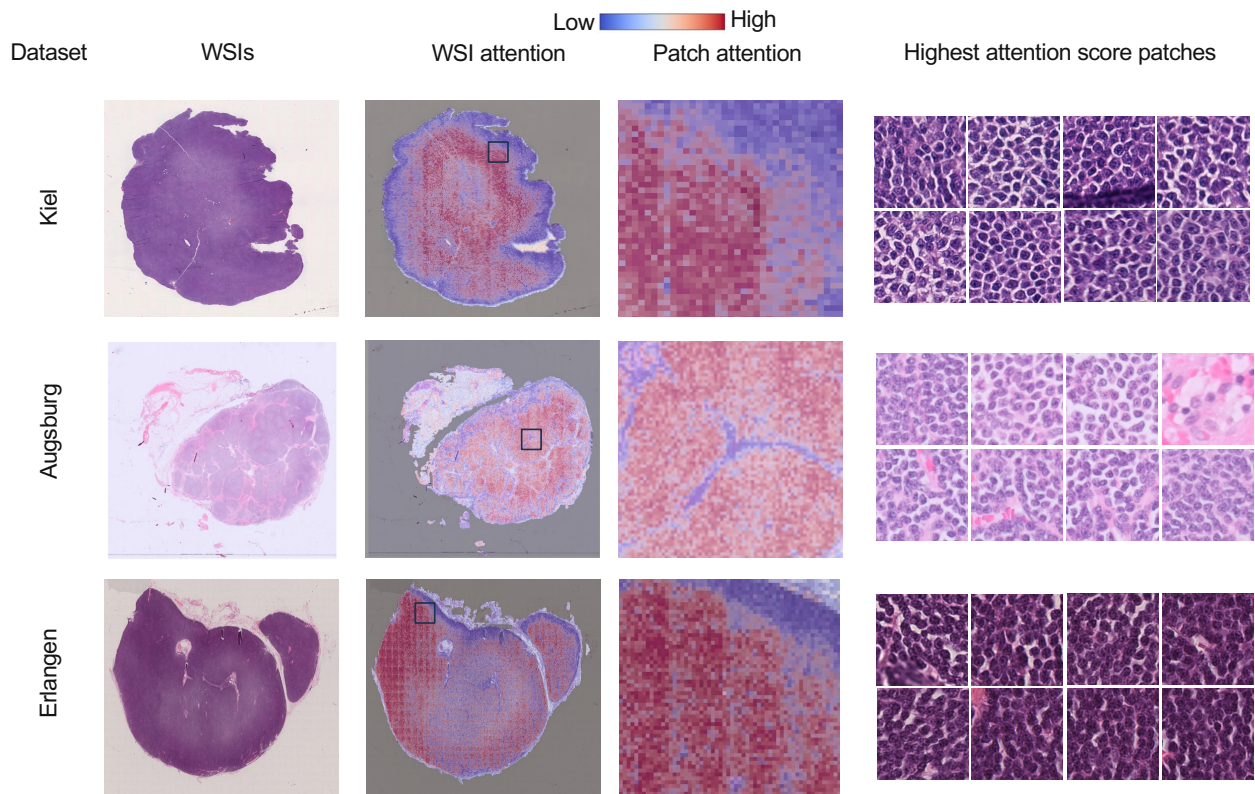
patches and averaging the attention values within the overlapping areas. In our workflow, we employ Trident [24] to generate WSI-level attention visualizations (Supplementary Figures 1, 2, and 3) and to display the patches with the highest attention scores. In Supplementary Figure 2, UNI2 features generalize more robustly than Titan, and avoid artifact patches. We note that a high-attention score in CLL is given to non-lymphoma areas and in DLBCL to poorly preserved areas (Supplementary Figure 1). This behavior might be avoidable when training the algorithm with more cases.



Supplementary Figure 1. Attention visualization on IID Munich testset.



Supplementary Figure 2. Attention visualization of a CLL lymphoma WSI from the OOD Erlangen cohort. UNI2 avoids artifact patches with the highest attention scores as compared to Titan.



Supplementary Figure 3. Attention visualization of CLL lymphoma subtype on OOD lymphoma testsets.