

Improving Explainability of Disentangled Representations using Multipath-Attribution Mappings

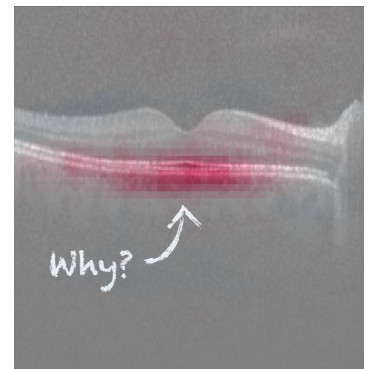
Novel framework combining disentangled representations with multipath-attribution, yielding enhanced interpretability and generalisation on medical datasets.

Authors: Lukas Klein^{1,2,3}, João Carvalho³, Mennatallah El-Assady³, Paolo Penna³, Joachim Buhmann³, Paul Jäger^{1,2}
Affiliations: ¹Helmholtz Imaging, ²DKFZ, ³ETH Zürich



Problem Statement

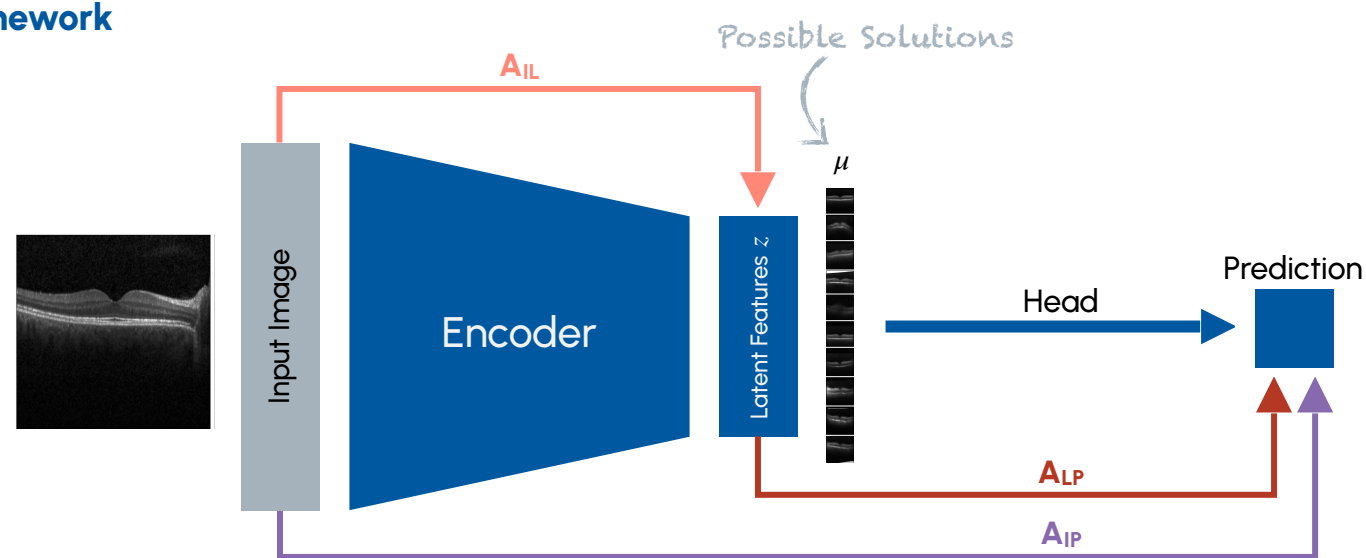
- Feature importance through **attribution maps** fall short when explaining **why a visual feature is used**, limiting causal statements and model interpretability.
- The attribution map on the right does not explain why the white retinal layer is important to the model's prediction of a healthy state. Could it be due to e.g. its shape, brightness, or thickness?



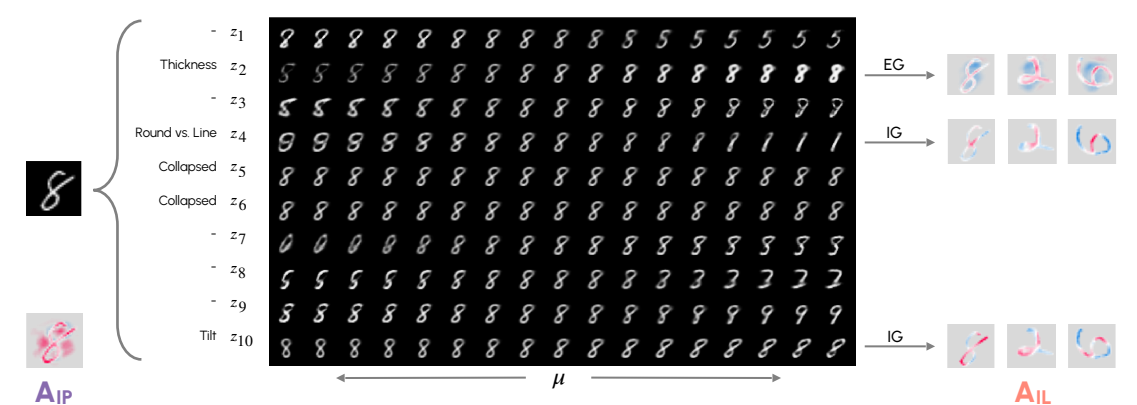
Our Approach

- Framework** based on capturing semantically meaningful features in disentangled latent features.
- Verifying (**PoC 1**) and enhancing (**PoC 2**) interpretations through multi-path attribution maps.
- Allowing for qualitative shortcut detection w/o OOD test set and explaining why a model fails or generalises.

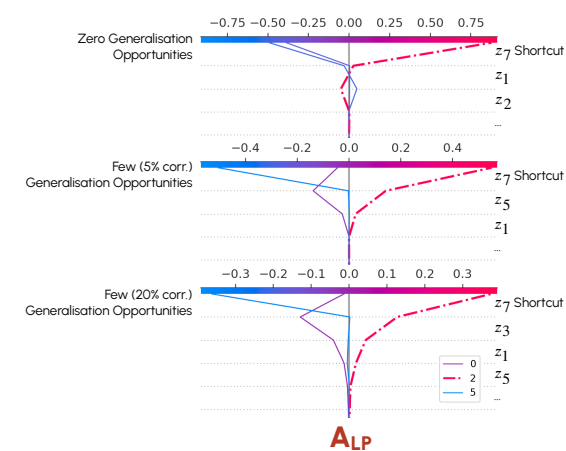
Framework



PoC 1: Verifying the interpretation of disentangled Latent Features by AIL.



PoC 2: Using enhanced interpretability by ALP for shortcut detection and generalisation evaluation.



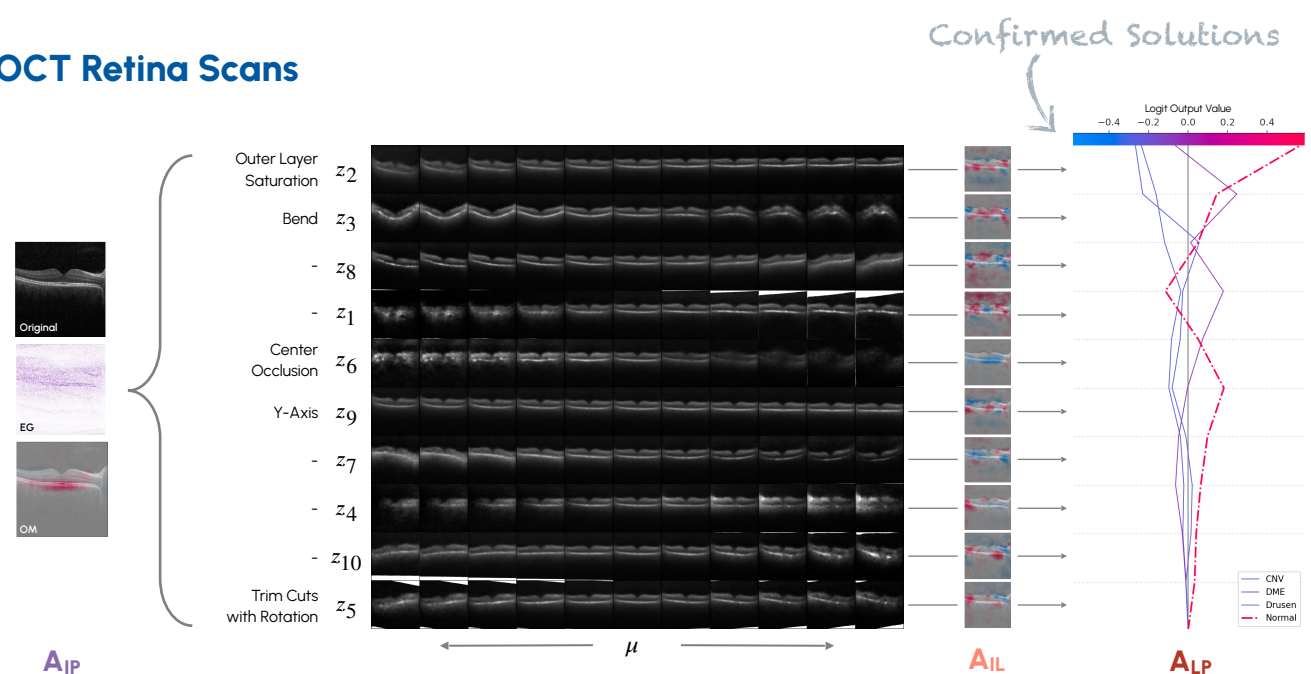
Attribution-paths:

AIP: Input Image into Prediction. **ALP**: Disentangled Latent Features into Prediction. **AIL**: Input Image into dis. Latent Features.

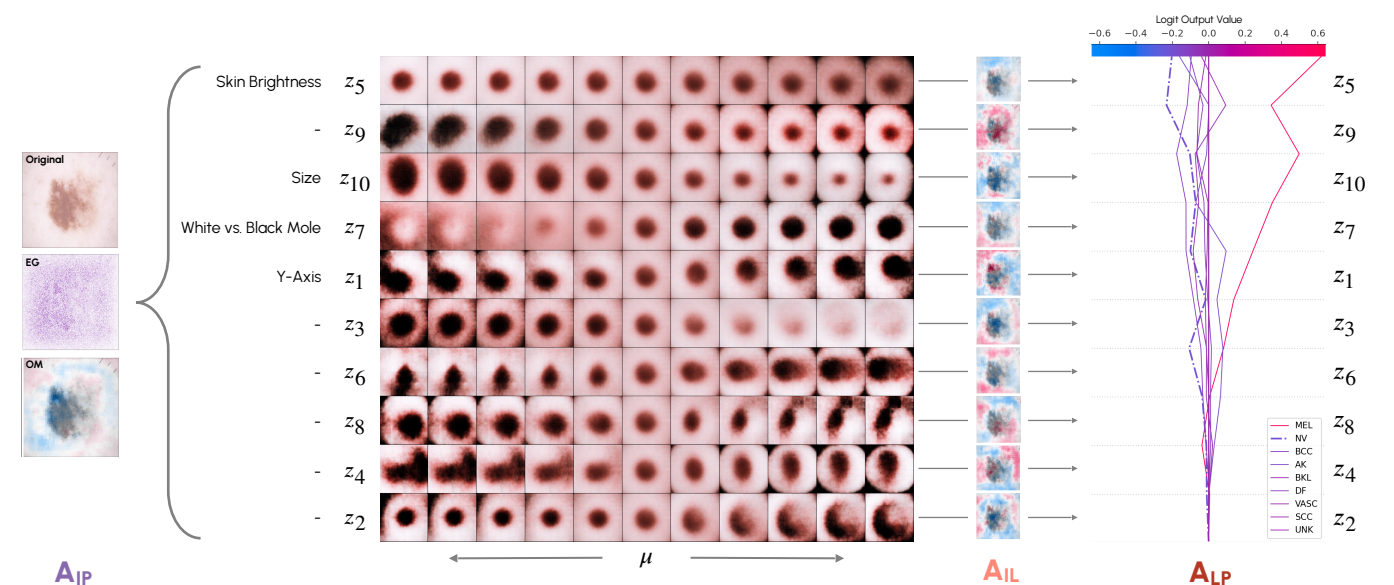
Experiments

- OCT Retina Scans**: (1) disentangles a possible shortcut (2) enhances interpretability by explaining why the white retinal layer is used for prediction (3) corrects a wrong interpretation based on the classical attribution map, indicating that the trim cuts are used as shortcuts.
- ISIC Skin Lesions**: (1) enables interpretability where classical attribution maps are uninformative (2) disentangles possible shortcuts (3) explains why the model fails.

OCT Retina Scans



ISIC Skin Lesions

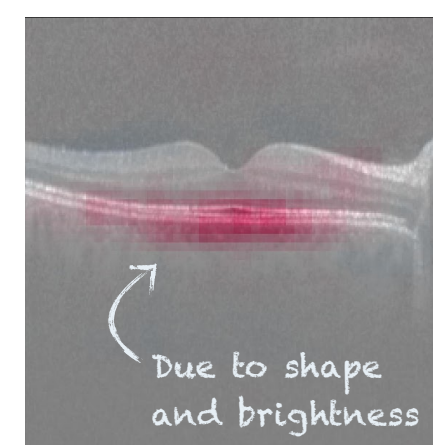


Results

The framework:

- catalyses more informative causality statements than classical saliency-maps
- facilitates qualitative detection of shortcut learning, and
- enables verification of model generalisation,

all combined and in an interactive setting.



<https://bit.ly/3w34pHF>



Paper

<https://bit.ly/3Fcz1M>



Code