

Interpretability Benchmark for Evaluating Spatial Misalignment of Prototypical Parts Explanations

Mikołaj Sacha^{1,2}, Bartosz Jura³, Dawid Rymarczyk^{1,4}, Łukasz Struski¹, Jacek Tabor¹, Bartosz Zieliński^{1,4}

¹ Faculty of Mathematics and Computer Science, Jagiellonian University

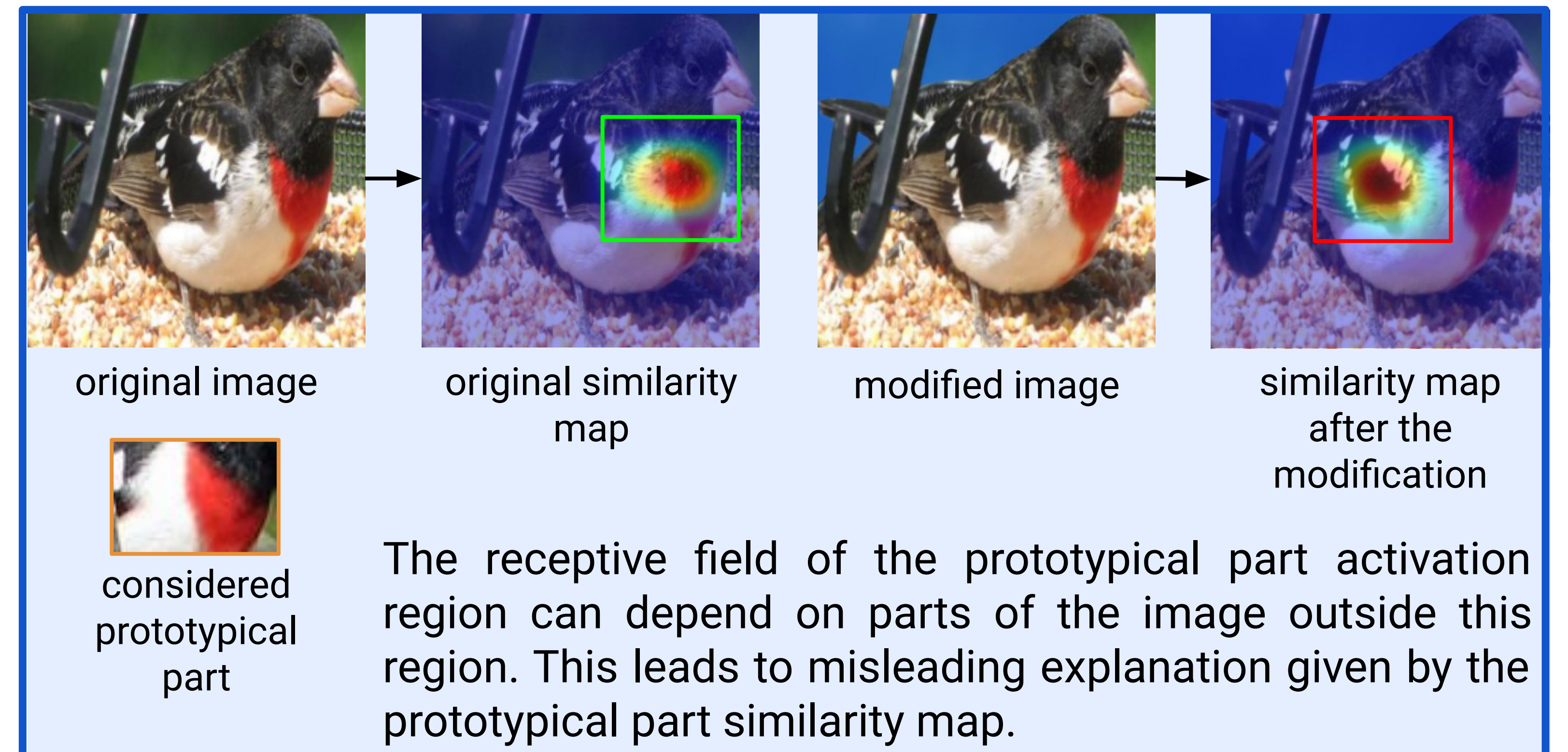
² Doctoral School of Exact and Natural Sciences, Jagiellonian University

³ Łukasiewicz Research Network – Poznań

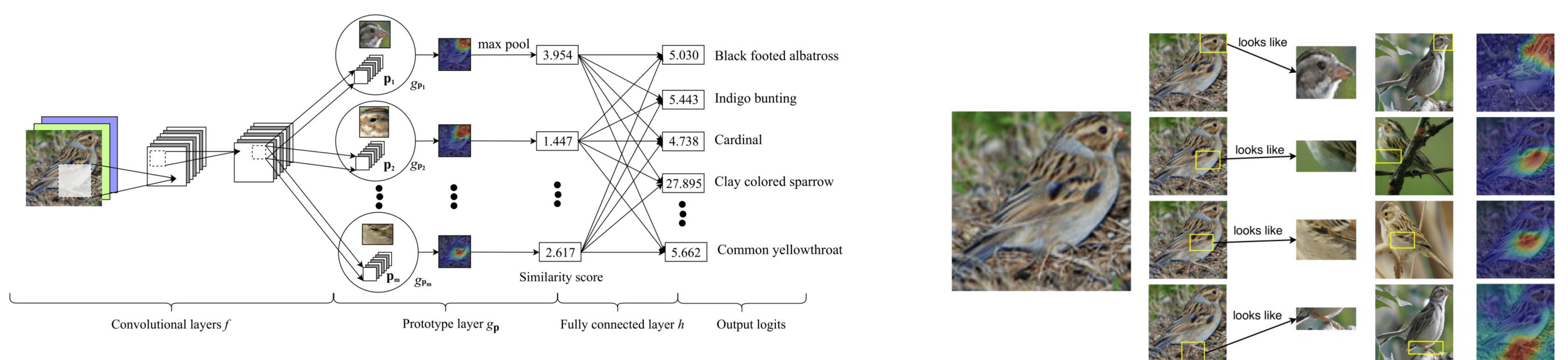
⁴ Ardigen SA

⁵ IDEAS NCBR

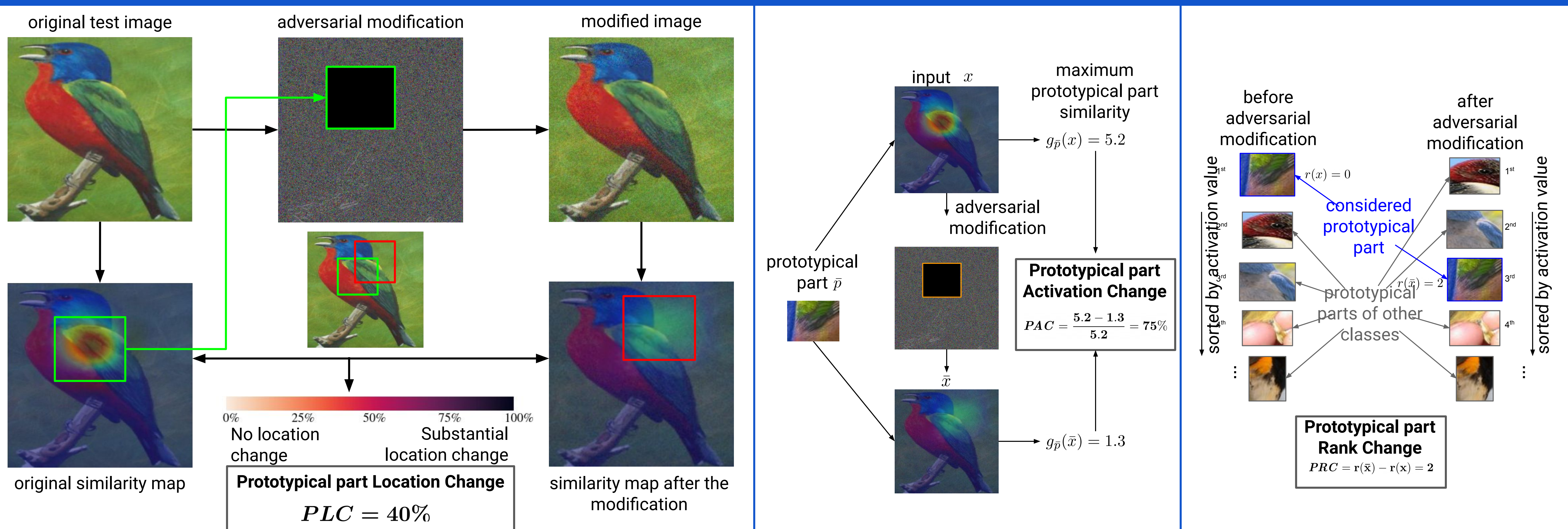
- **Prototypical part-based networks** perform interpretable image classification by comparing regions of images to *prototypical parts*.
- We examine *spatial misalignment* between the explanations provided by prototypical part-based networks and their actual inference mechanism.
- We measure *spatial misalignment* for popular types of prototypical part-based network and **introduce a method for enforcing learning truthful explanations**.



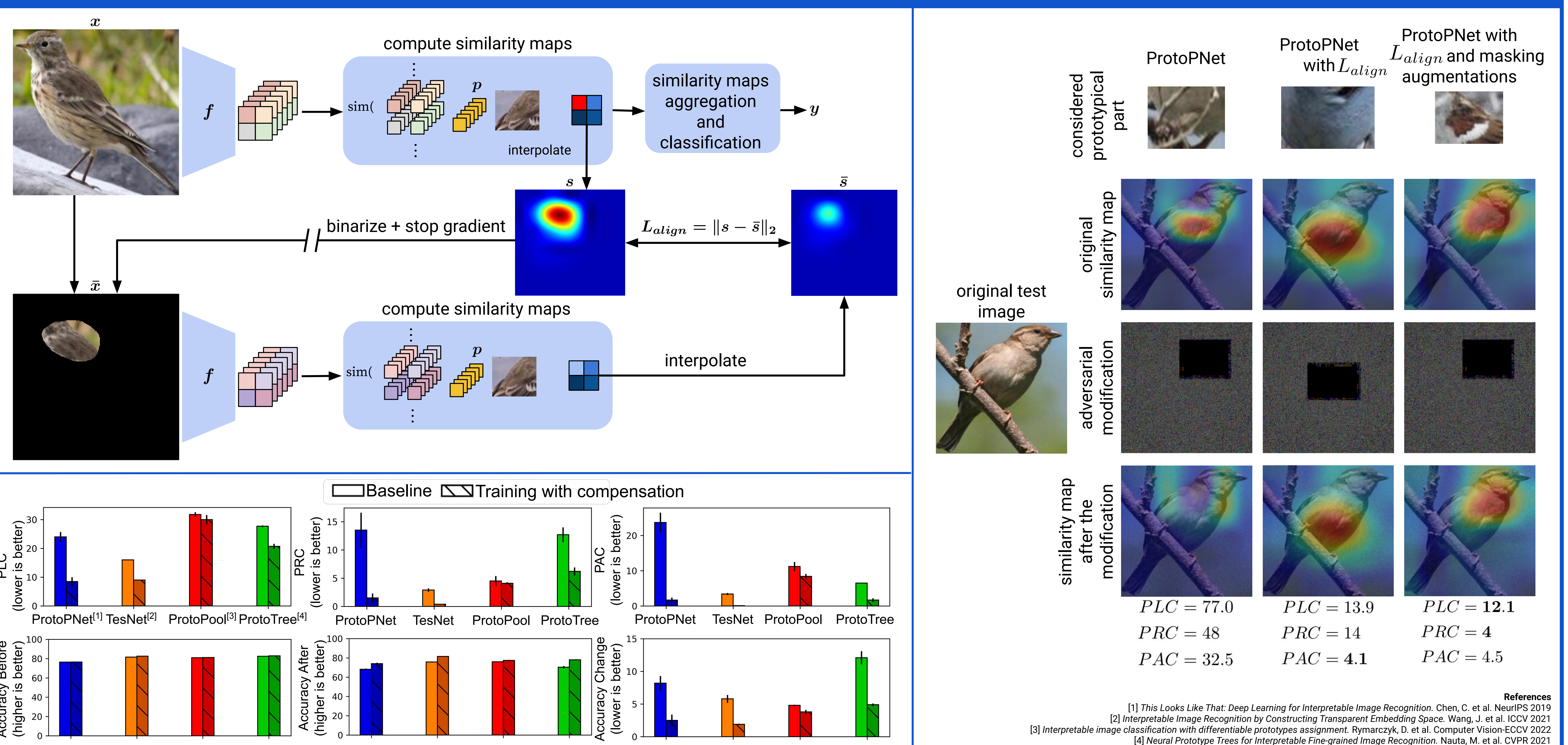
Preliminary: Prototypical Part-Based Networks [1]



Interpretability Benchmark



Training with spatial misalignment compensation



IDEAS
NCBR

group of machine
gmum
learning research

JAGIELLONIAN
UNIVERSITY
IN KRAKÓW