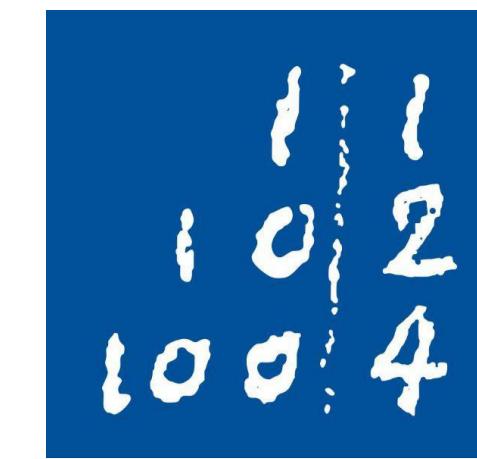


# Explanation-Guided Training for Cross-Domain Few-Shot Classification



Niklas Baack, Jan Malte Töpperwien, Peer Bastian Duensing  
Poster Presentation in context of Interpretable Machine Learning

1

## Summary

- LRP for guided training
- Better performance on cross-domain tasks
- Mostly model agnostic (gradients needed)

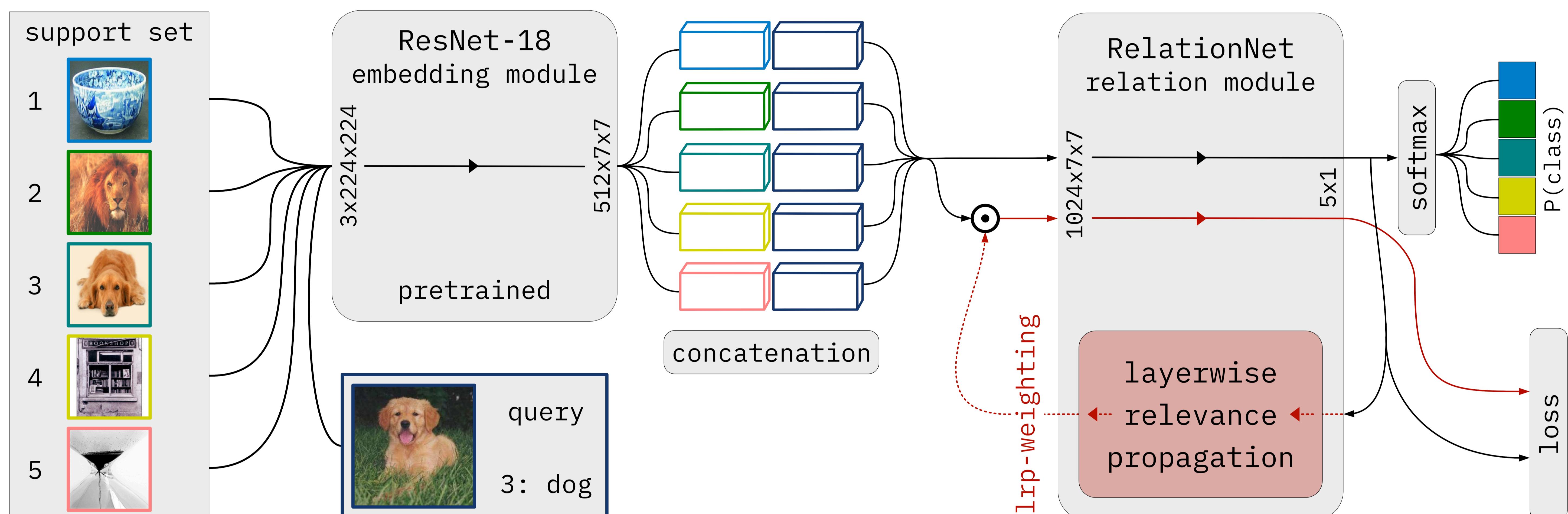
2

## Motivation & Problem Setting

- Humans can identify objects with help of few examples
- Models struggle with cross-domain few-shot classification
- Utilize LRP explanations for performance improvement

3

## Approach

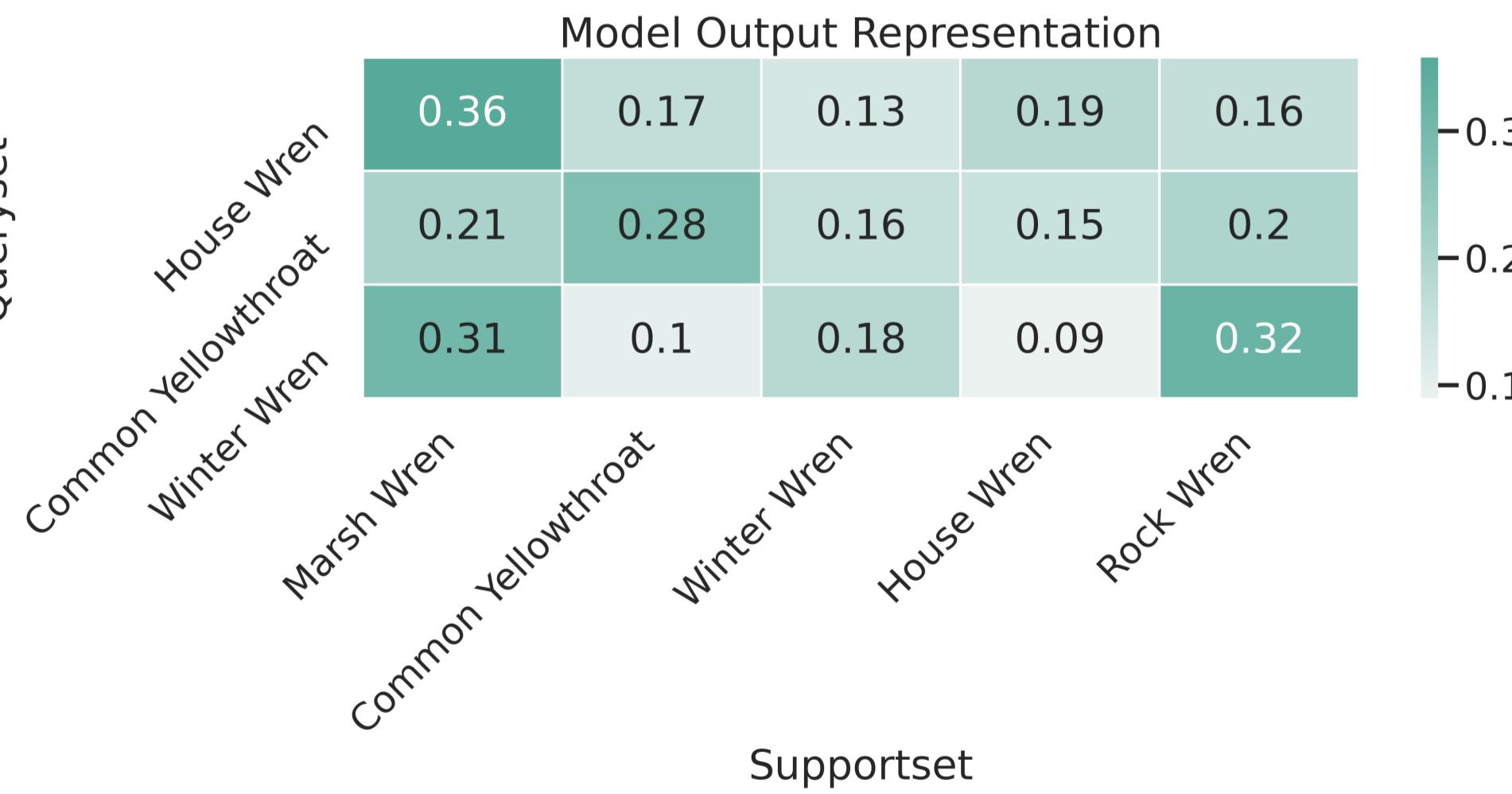


## LRP rules for interpretable backpropagation

- $\alpha$ -rule for convolution layers
- $\varepsilon$ -rule for linear layers
- Variation:  $\gamma$ -rule

$$R_{i \leftarrow j} = R(z_j^{l+1}) \frac{z_i^l w_{ij}}{\sum_i z_i^l w_{ij} + b_j}$$

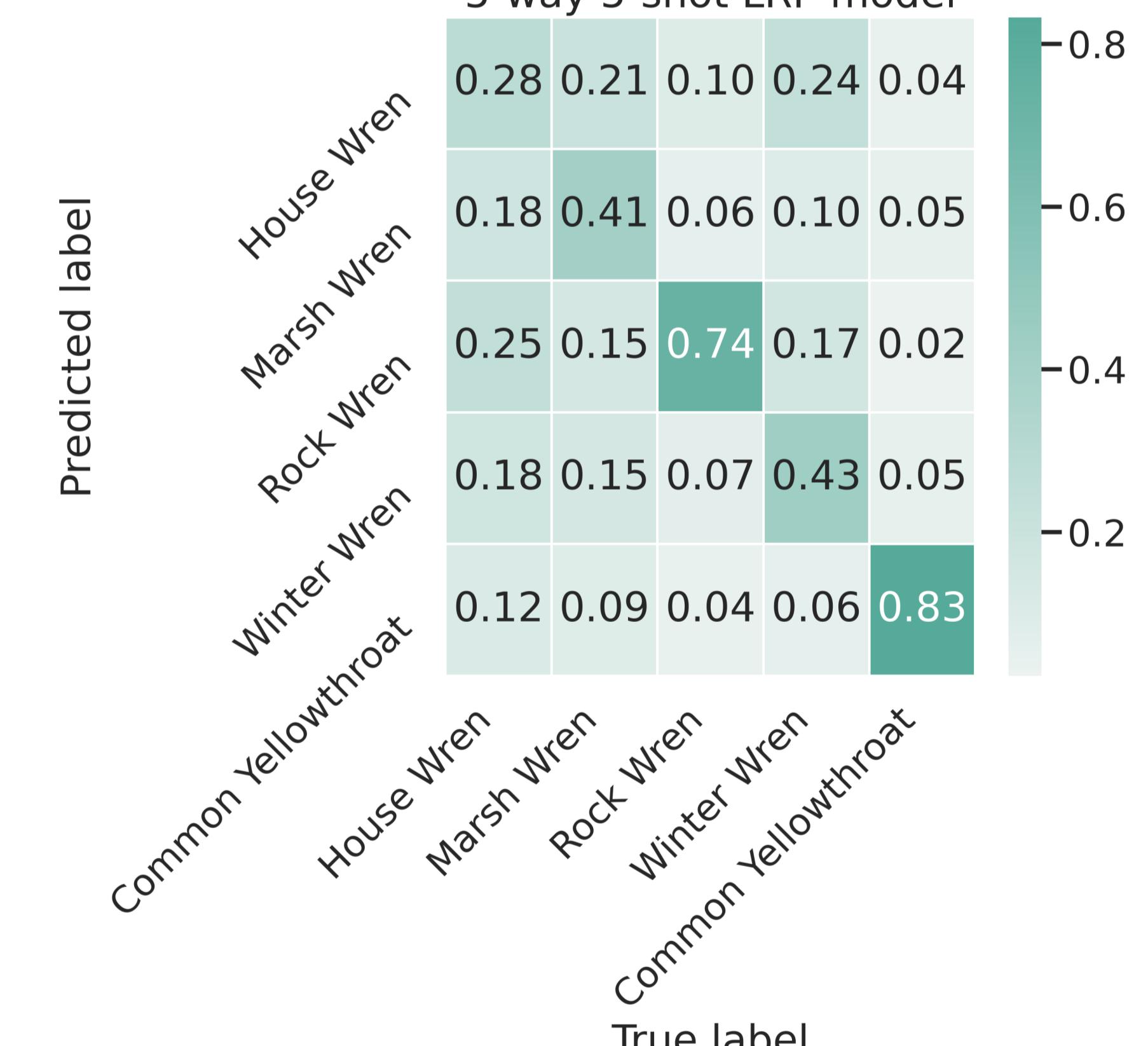
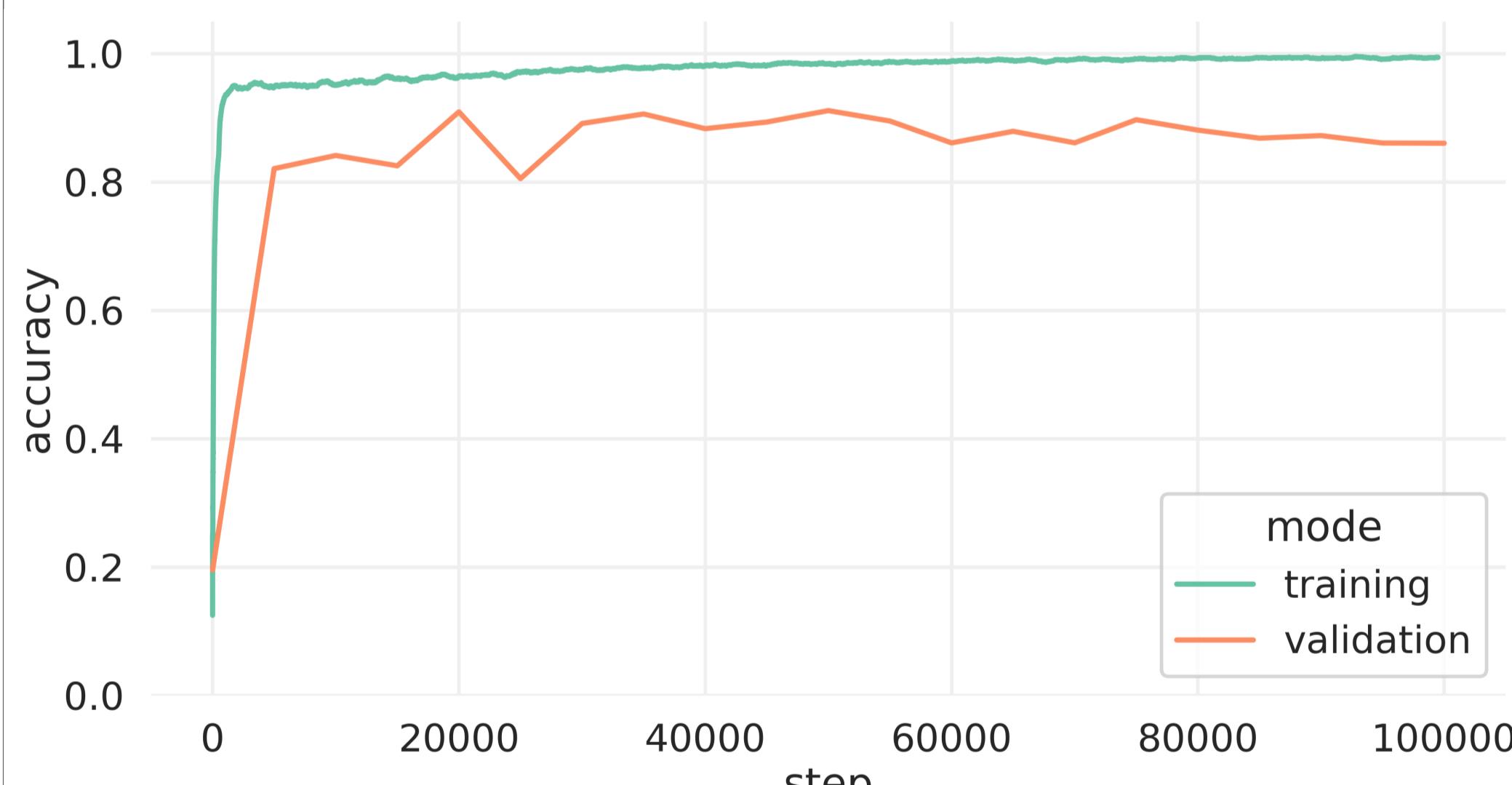
upstream relevance      gradient  
relevance                      layer normalization



4

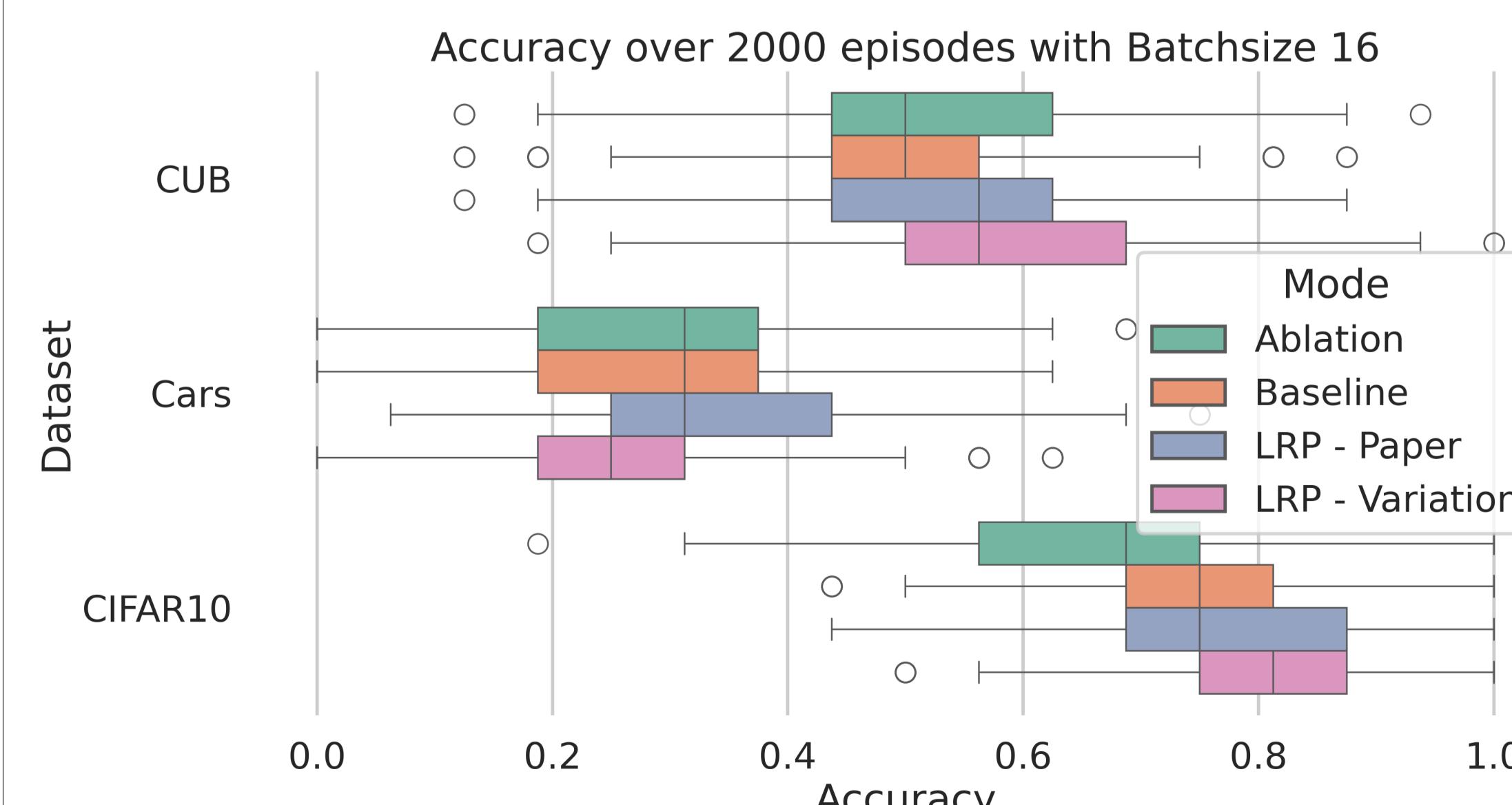
## Key Insights

### Training accuracy (5-way 5-shot)

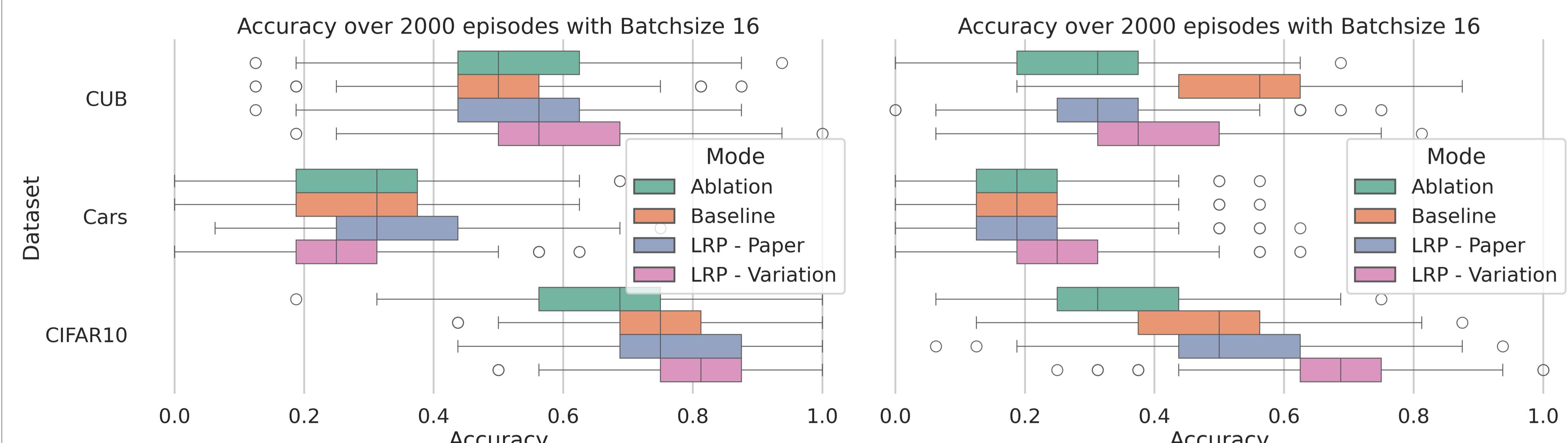


### Cross-Domain Classification

#### 5-way 5-shot



#### 5-way 1-shot



5

## Future Works

- Apply on other architectures
- Use on harder domains
- Use with more classes

- Advanced generalization by noise suppression (LRP rules)
- Better embeddings boost performance significantly