

# Test-Time Selection for Robust Skin Lesion Analysis



**Alceu Bissoto<sup>1</sup>, Catarina Barata<sup>2</sup>, Eduardo Valle<sup>3</sup>, Sandra Avila<sup>1</sup>**

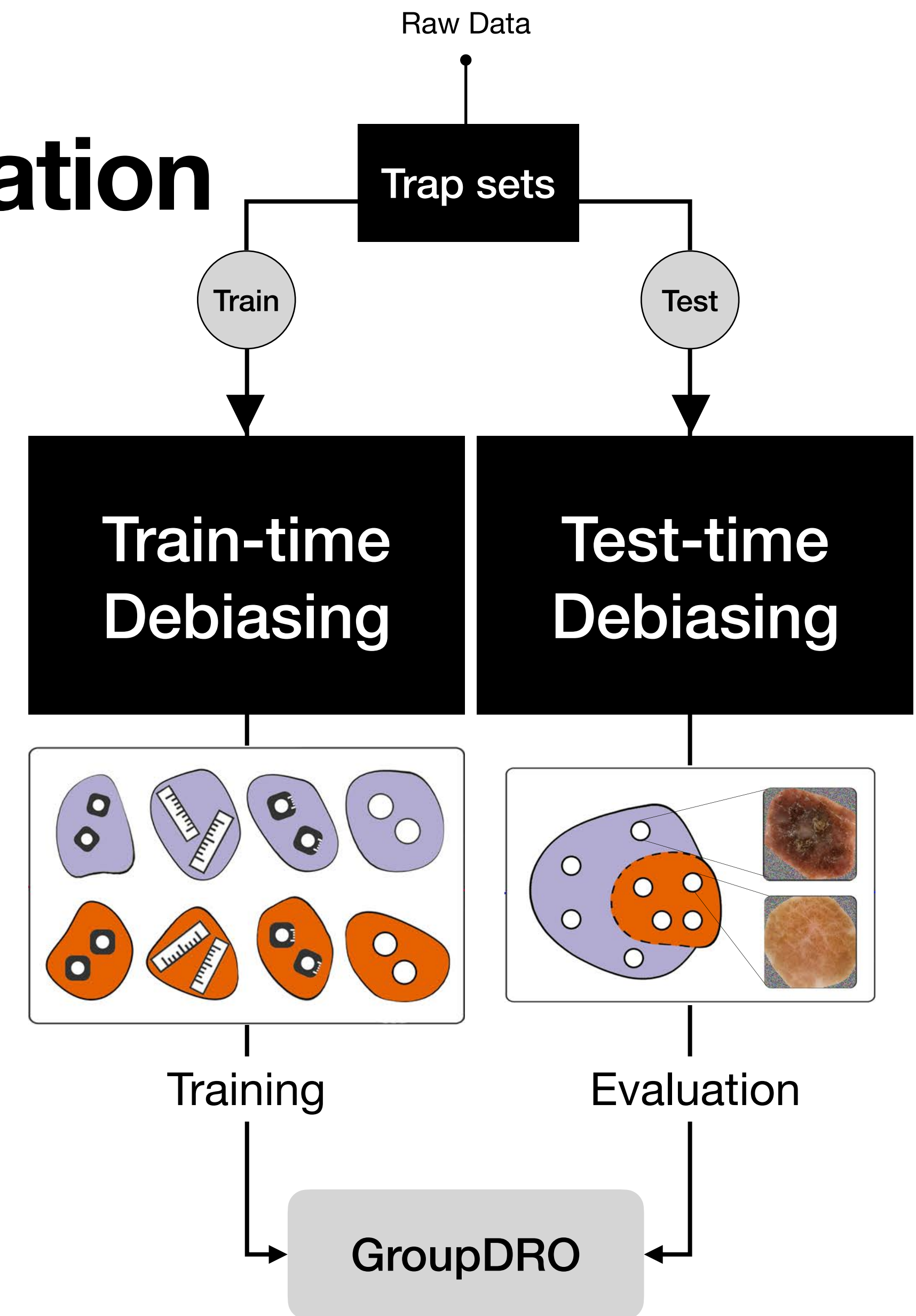
<sup>1</sup>Institute of Computing    <sup>3</sup>School of Electrical and Computing Engineering  
Recod.ai, University of Campinas (UNICAMP), Brazil

<sup>2</sup>Institute for Systems and Robotics, Instituto Superior Técnico, Portugal



# Artifact-based Domain Generalization

ISIC Workshop @ ECCV 2022



## Artifact-based Domain Generalization of Skin Lesion Models

Alceu Bissoto<sup>[0000-0003-2293-6160]1,4</sup>, Catarina Barata<sup>[0000-0002-2852-7723]2</sup>,  
Eduardo Valle<sup>[0000-0001-5396-9868]3,4</sup>, and Sandra Avila<sup>[0000-0001-9068-938X]1,4</sup>

<sup>1</sup> Institute of Computing, University of Campinas, Brazil  
{alceubissoto, sandra}@ic.unicamp.br

<sup>2</sup> Institute for Systems and Robotics, Instituto Superior Técnico, Portugal  
ana.c.fidalgo.barata@tecnico.ulisboa.pt

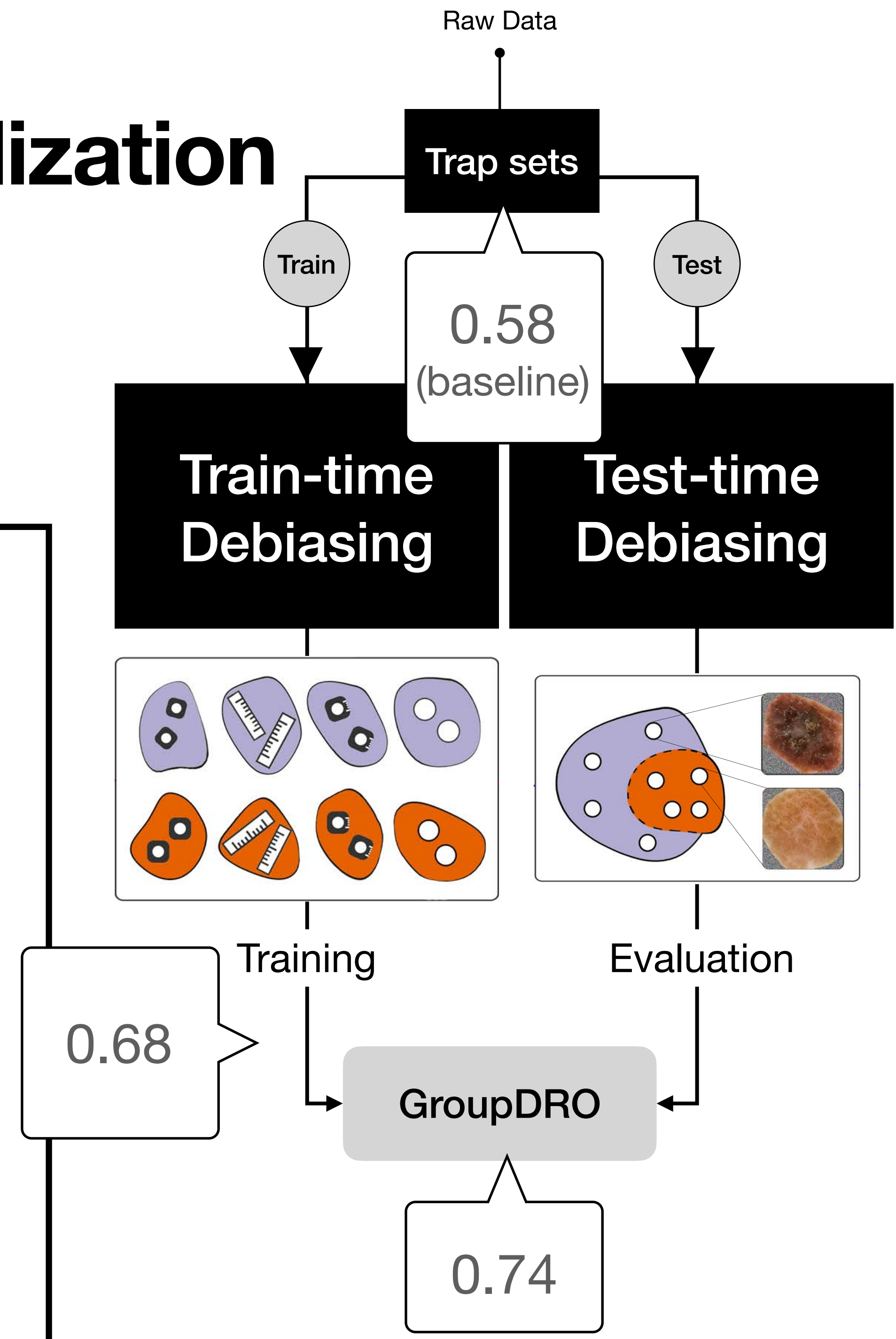
<sup>3</sup> School of Electrical and Computing Engineering, University of Campinas, Brazil  
dovalle@dca.fee.unicamp.br

<sup>4</sup> Recod.ai Lab, University of Campinas, Brazil

**Abstract.** Deep Learning failure cases are abundant, particularly in the medical area. Recent studies in out-of-distribution generalization have advanced considerably on well-controlled synthetic datasets, but they do not represent medical imaging contexts. We propose a pipeline that relies

# Artifact-based Domain Generalization

## ISIC Workshop @ ECCV 2022



### Artifact-based Domain Generalization of Skin Lesion Models

Alceu Bissoto<sup>[0000-0003-2293-6160]1,4</sup>, Catarina Barata<sup>[0000-0002-2852-7723]2</sup>,  
Eduardo Valle<sup>[0000-0001-5396-9868]3,4</sup>, and Sandra Avila<sup>[0000-0001-9068-938X]1,4</sup>

<sup>1</sup> Institute of Computing, University of Campinas, Brazil  
{alceubissoto, sandra}@ic.unicamp.br

<sup>2</sup> Institute for Systems and Robotics, Instituto Superior Técnico, Portugal  
ana.c.fidalgo.barata@tecnico.ulisboa.pt

<sup>3</sup> School of Electrical and Computing Engineering, University of Campinas, Brazil  
dovalle@dca.fee.unicamp.br

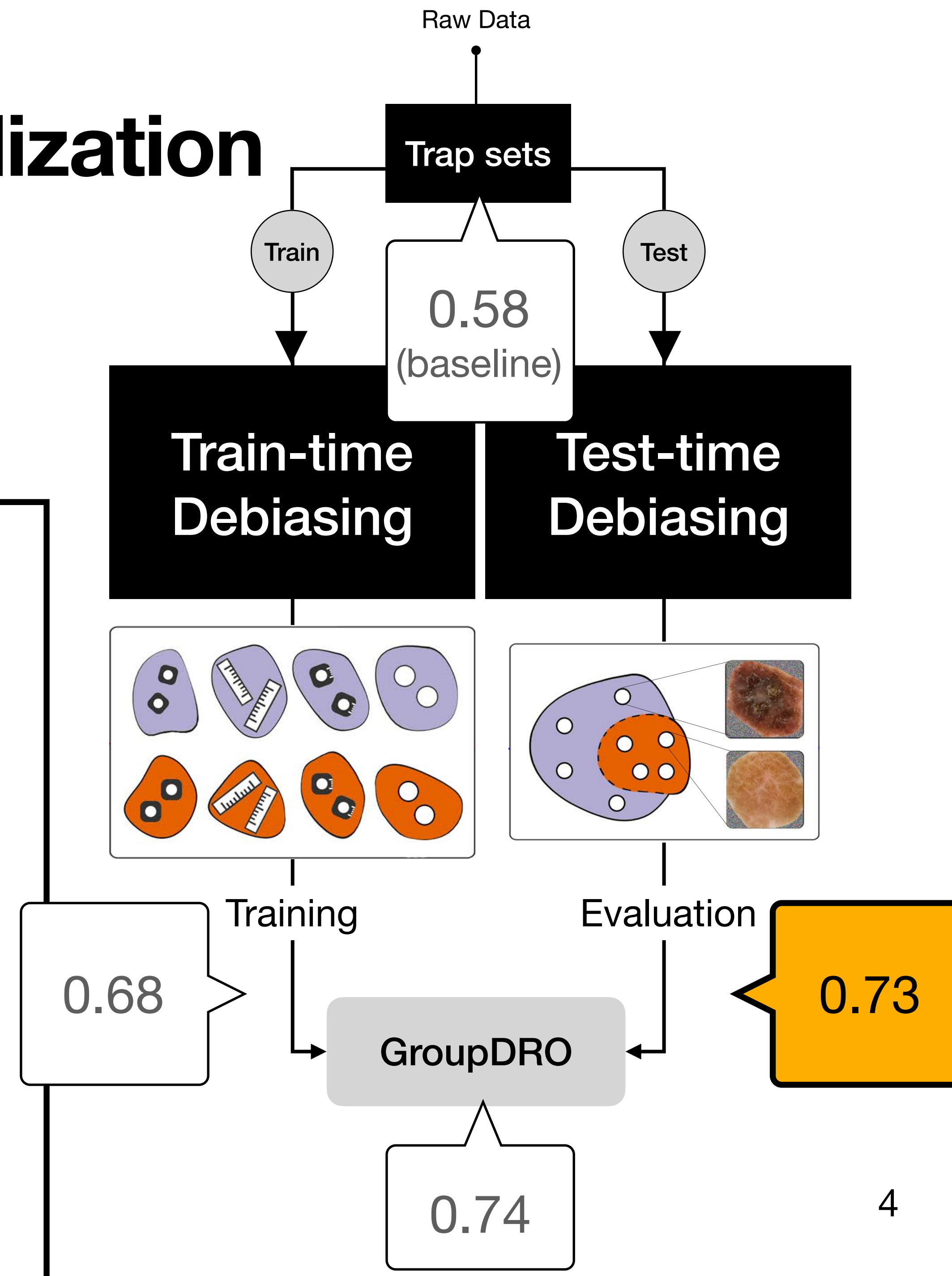
<sup>4</sup> Recod.ai Lab, University of Campinas, Brazil

**Abstract.** Deep Learning failure cases are abundant, particularly in the medical area. Recent studies in out-of-distribution generalization have advanced considerably on well-controlled synthetic datasets, but they do not represent medical imaging contexts. We propose a pipeline that relies



# Artifact-based Domain Generalization

## ISIC Workshop @ ECCV 2022



### Artifact-based Domain Generalization of Skin Lesion Models

Alceu Bissoto<sup>[0000-0003-2293-6160]1,4</sup>, Catarina Barata<sup>[0000-0002-2852-7723]2</sup>,  
Eduardo Valle<sup>[0000-0001-5396-9868]3,4</sup>, and Sandra Avila<sup>[0000-0001-9068-938X]1,4</sup>

<sup>1</sup> Institute of Computing, University of Campinas, Brazil  
{alceubissoto, sandra}@ic.unicamp.br

<sup>2</sup> Institute for Systems and Robotics, Instituto Superior Técnico, Portugal  
ana.c.fidalgo.barata@tecnico.ulisboa.pt

<sup>3</sup> School of Electrical and Computing Engineering, University of Campinas, Brazil  
dovalle@dca.fee.unicamp.br

<sup>4</sup> Recod.ai Lab, University of Campinas, Brazil

**Abstract.** Deep Learning failure cases are abundant, particularly in the medical area. Recent studies in out-of-distribution generalization have advanced considerably on well-controlled synthetic datasets, but they do not represent medical imaging contexts. We propose a pipeline that relies



# Test-time Debiasing Literature

## Alignment with the Clinical Workflow

	Method	#Keypoints	AUC
<b>Baseline</b>	Test-time augmentation	-	58,4
<b>Literature</b>	T3A	-	56,7
<b>Literature</b>	Tent	-	54,1
<b>Literature</b>	NoiseCrop	50.176	72,7

Test-time debiasing literature

- Relies on test batch statistics to update model weights.
- Fail when **only a single image is available.**
- Fail when **test distribution is heterogenous.**

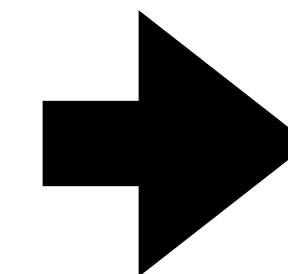
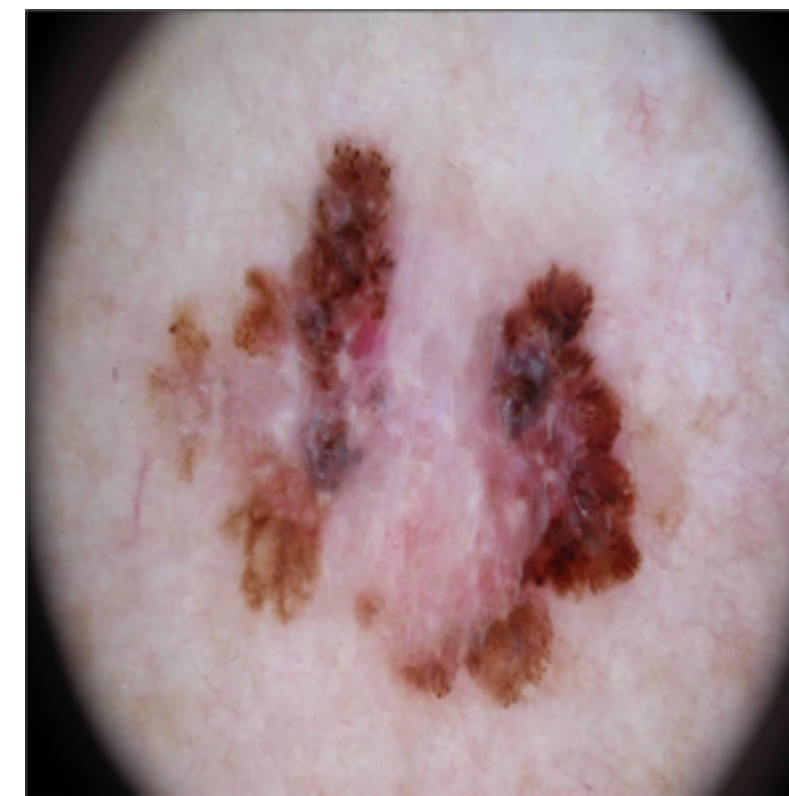
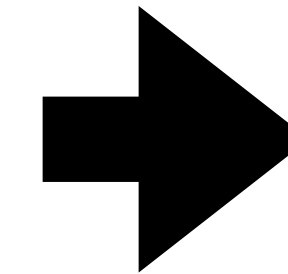
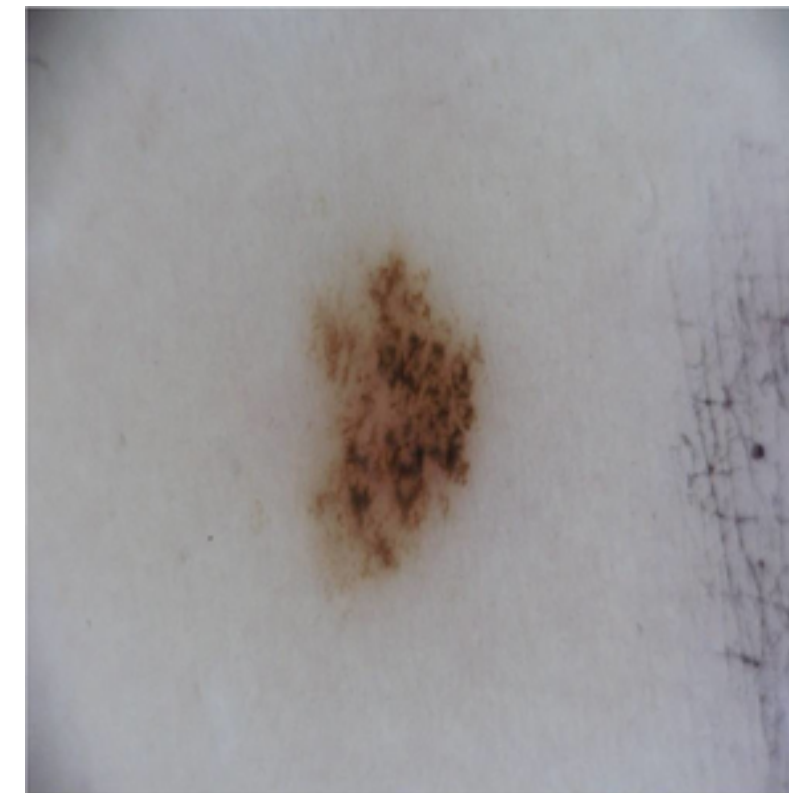


# Test-time Debiasing Literature

## Alignment with the Clinical Workflow

NoiseCrop

- Relies on full segmentation masks, which are hard to annotate.
- Make modifications in the pixel-space, which might introduce unexpected features.





# Test-time Selection (TTS)

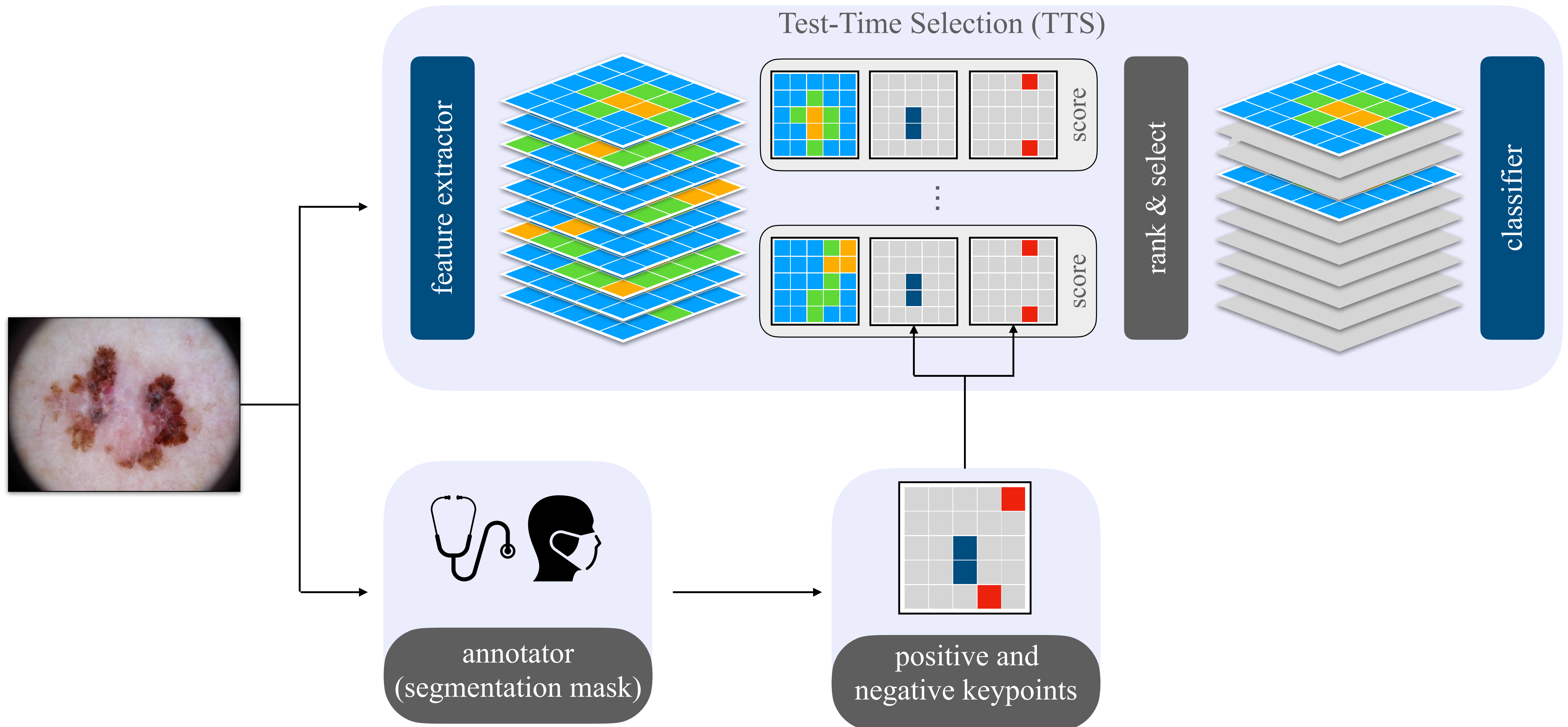
- **Fast to annotate.**
- **Avoid introducing distribution shifts** by intervening on the feature space.
- **Does not rely** on test batch statistics.
- It's **cheap** as there are no model updates.

	Method	#Keypoints	AUC
<b>Baseline</b>	Test-time augmentation	-	58,4
<b>Literature</b>	T3A	-	56,7
<b>Literature</b>	Tent	-	54,1
<b>Literature</b>	NoiseCrop	50.176	72,7
<b>Ours</b>	TTS	40	75,0



# Methodology

# Methodology



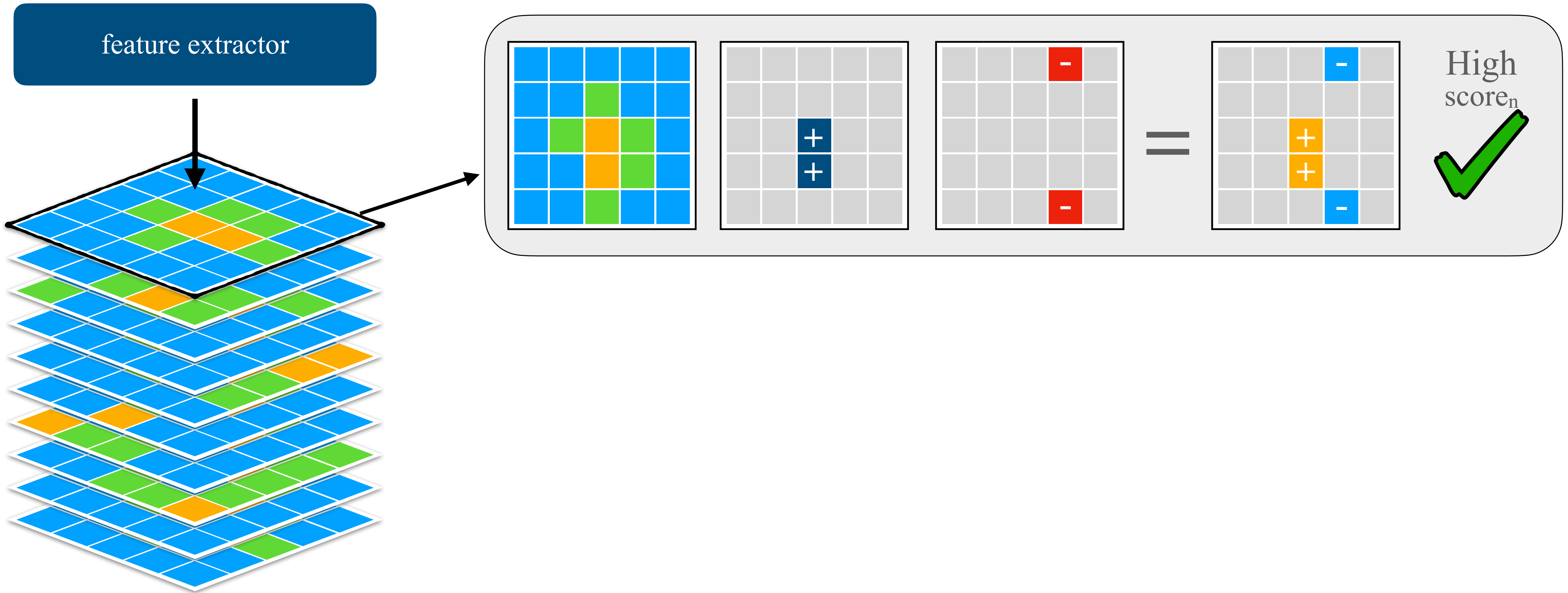


# Methodology

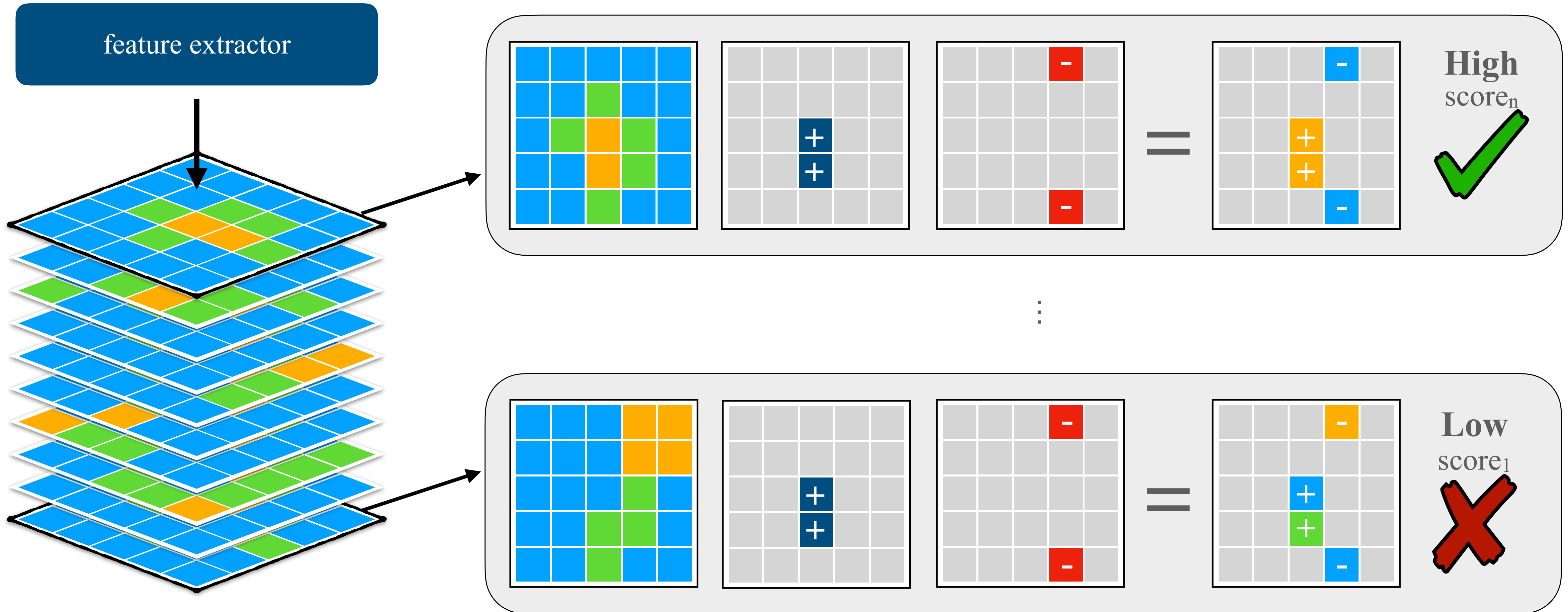
Attention  
Maps

Positive  
Keypoints

Negative  
Keypoints

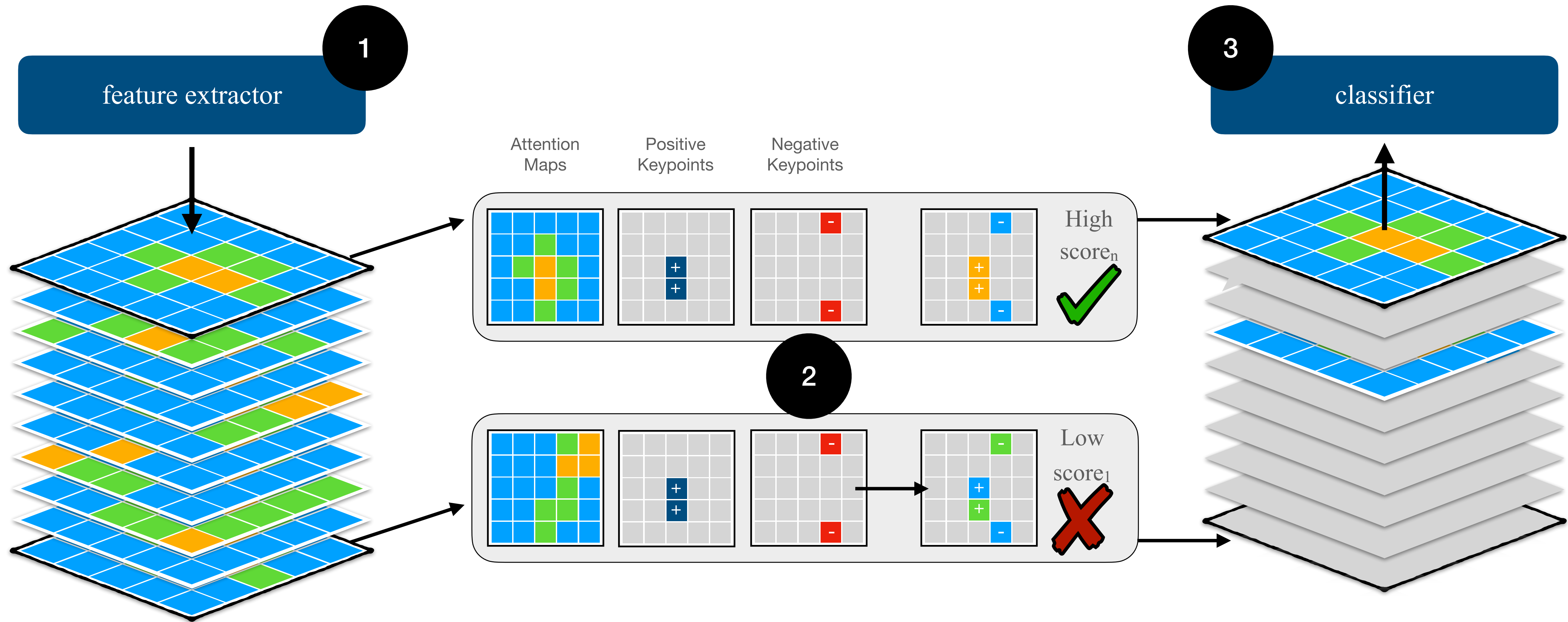


# Methodology





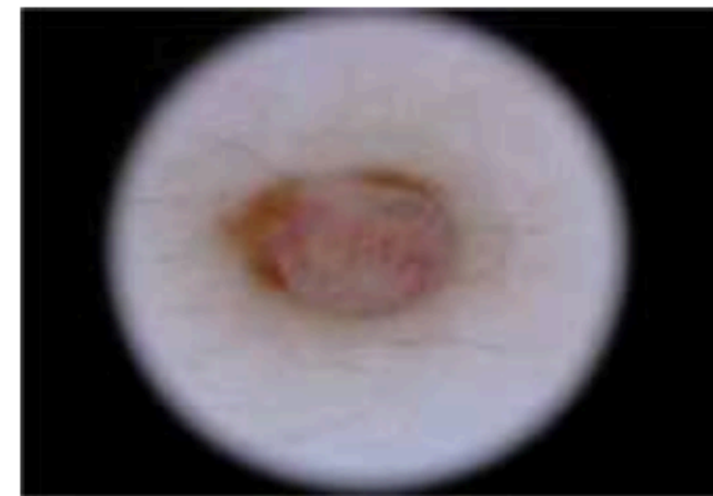
# Methodology



# Evaluation Protocol



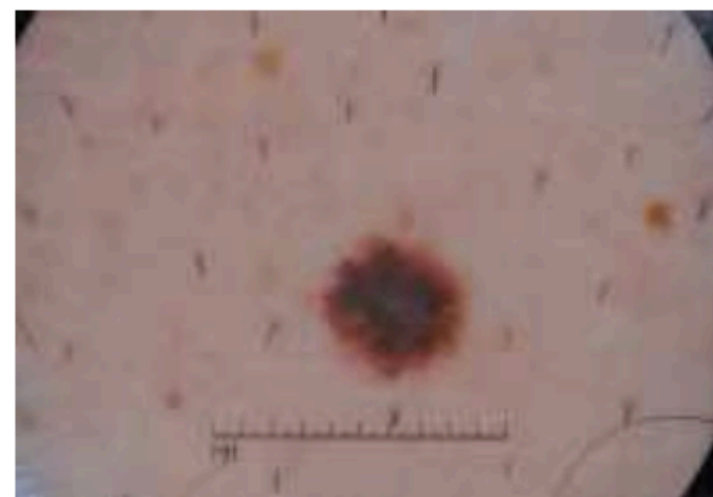
# Artifacts providing spurious correlations



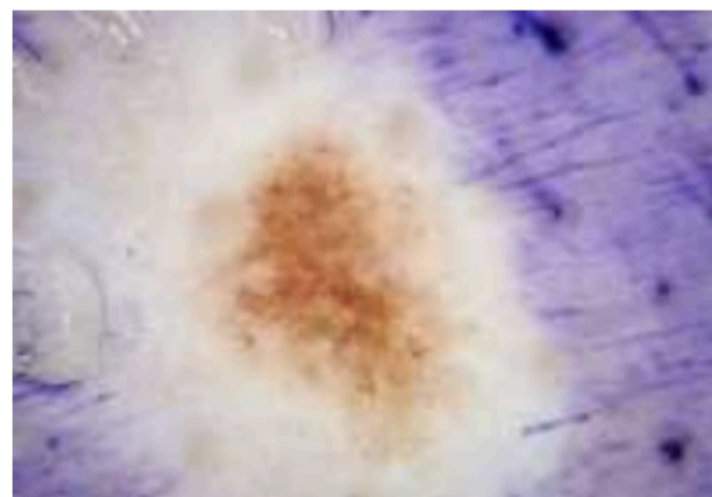
Dark Corners



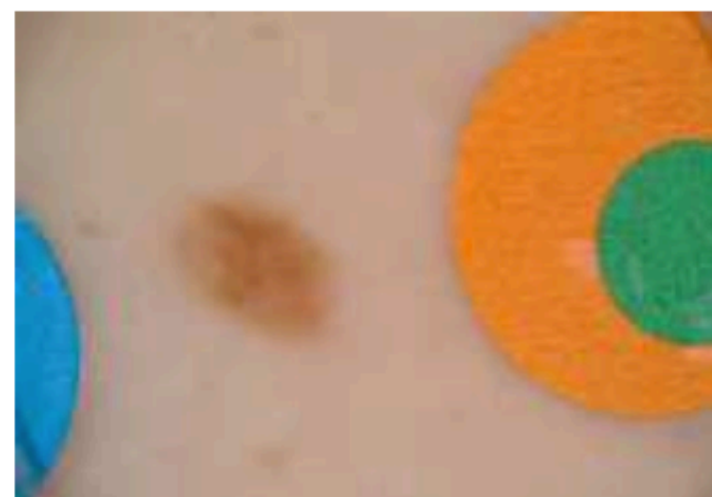
Hair



Ruler



Ink markings

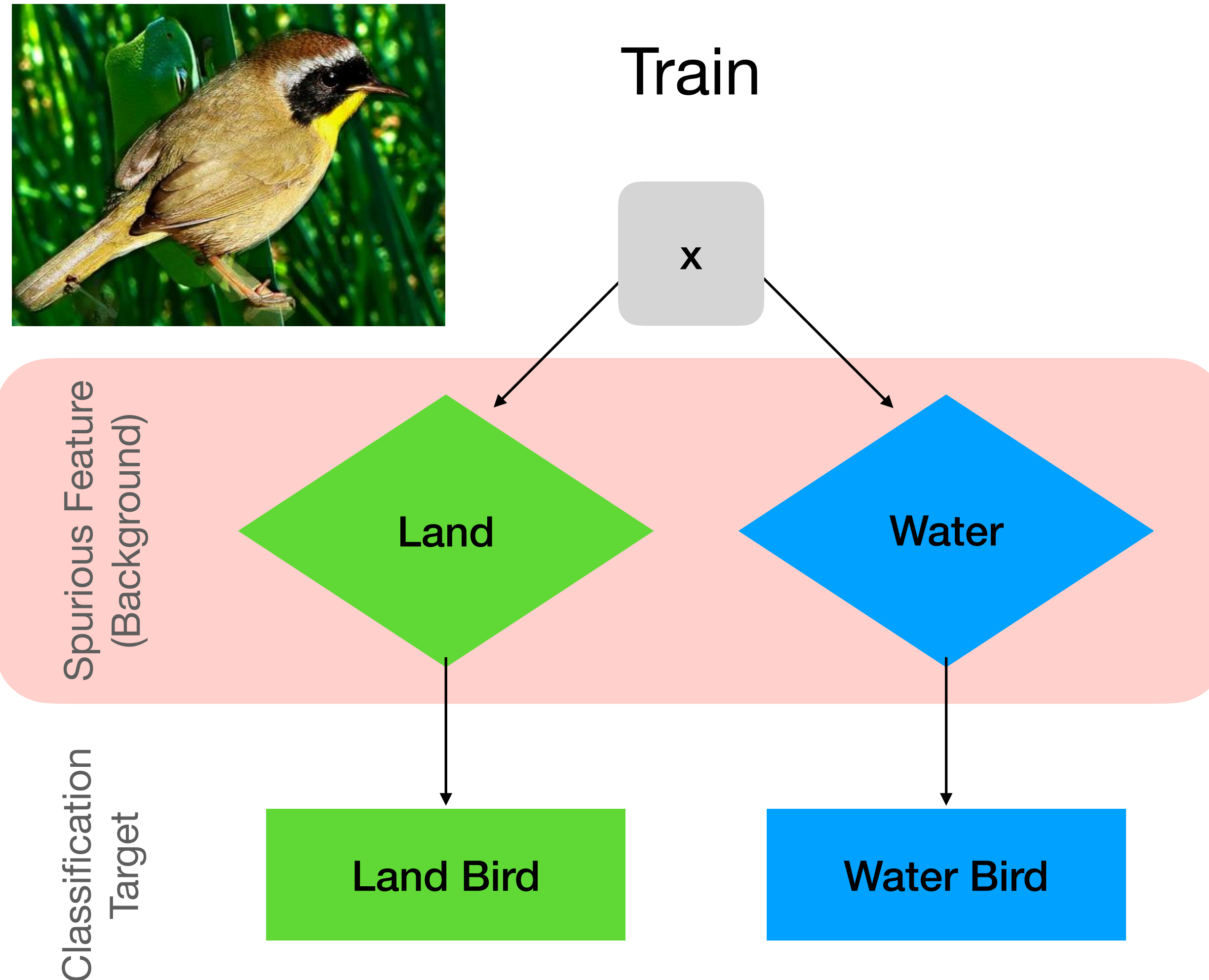


Patches

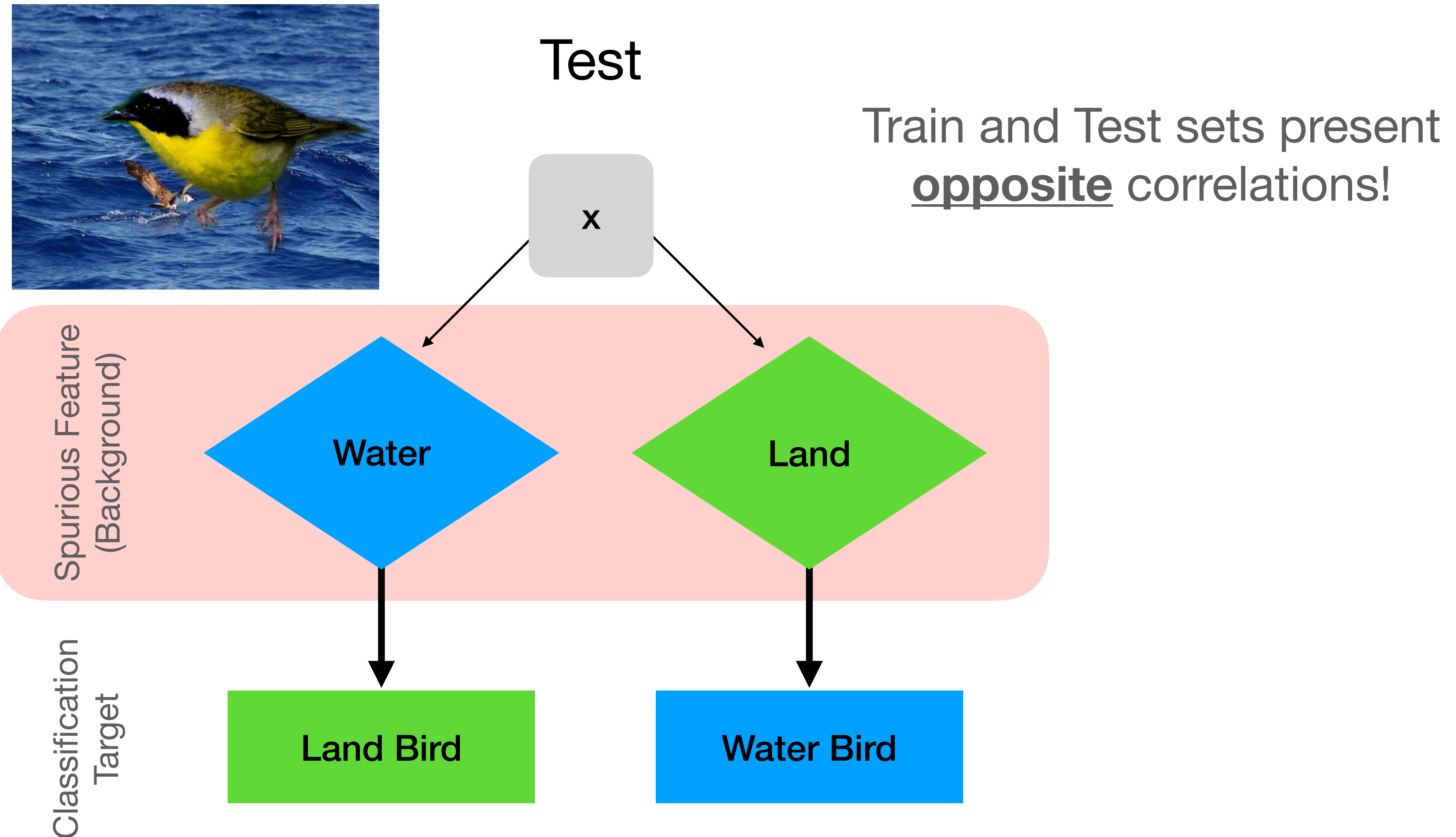
factor	set	dark corner	hair	ruler	ink	patches
0	train	0.119	-0.104	0.142	0.023	-0.138
	test	0.135	-0.112	0.162	0.030	-0.149

- Mild correlations.
- Gaining robustness to artifacts will hardly impact any metric.
- **We need to control/amplify the correlations.**

# Spurious Features vs. Generalization

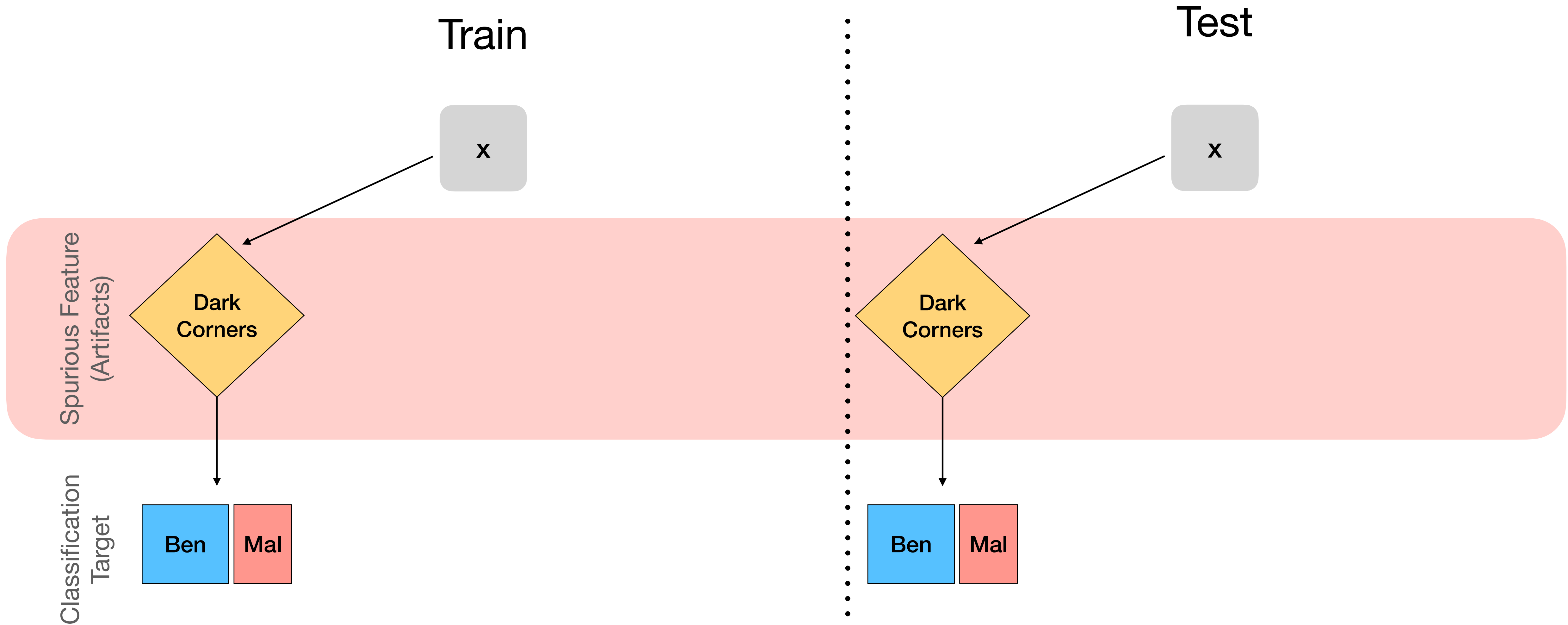


# Spurious Features vs. Generalization

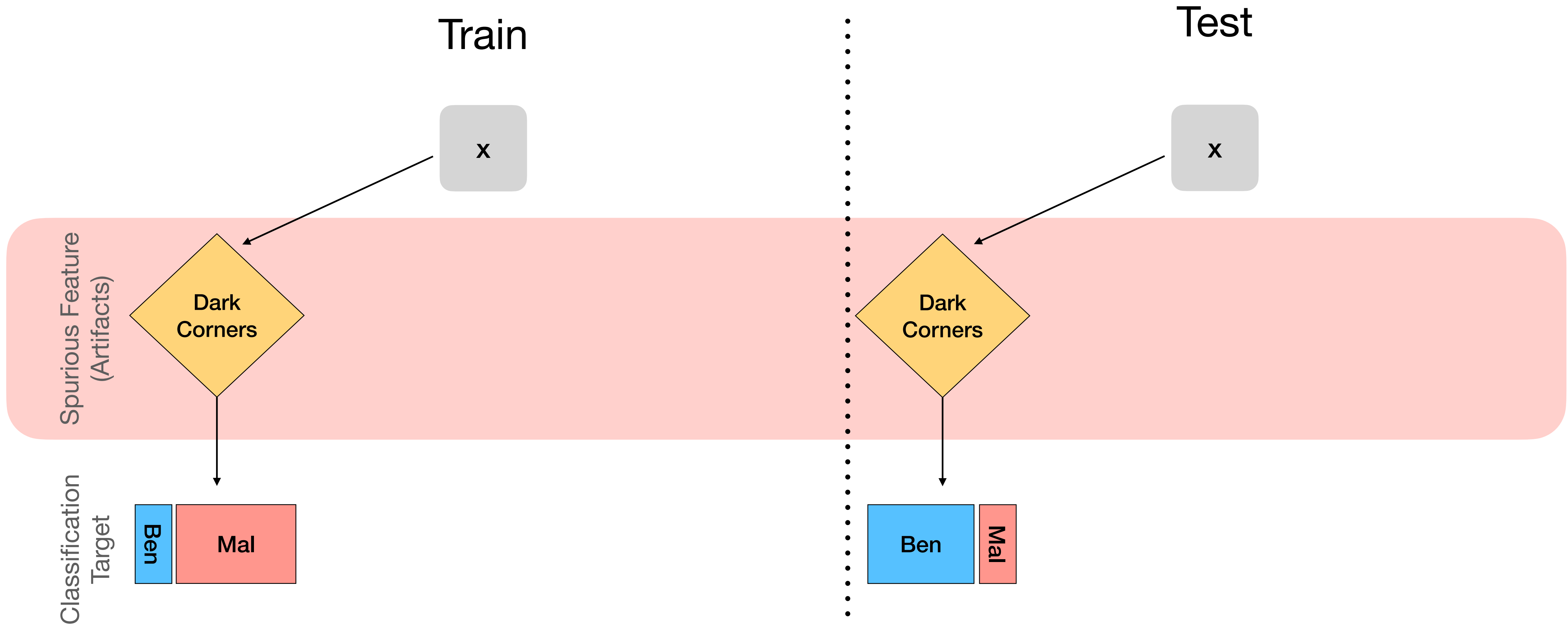




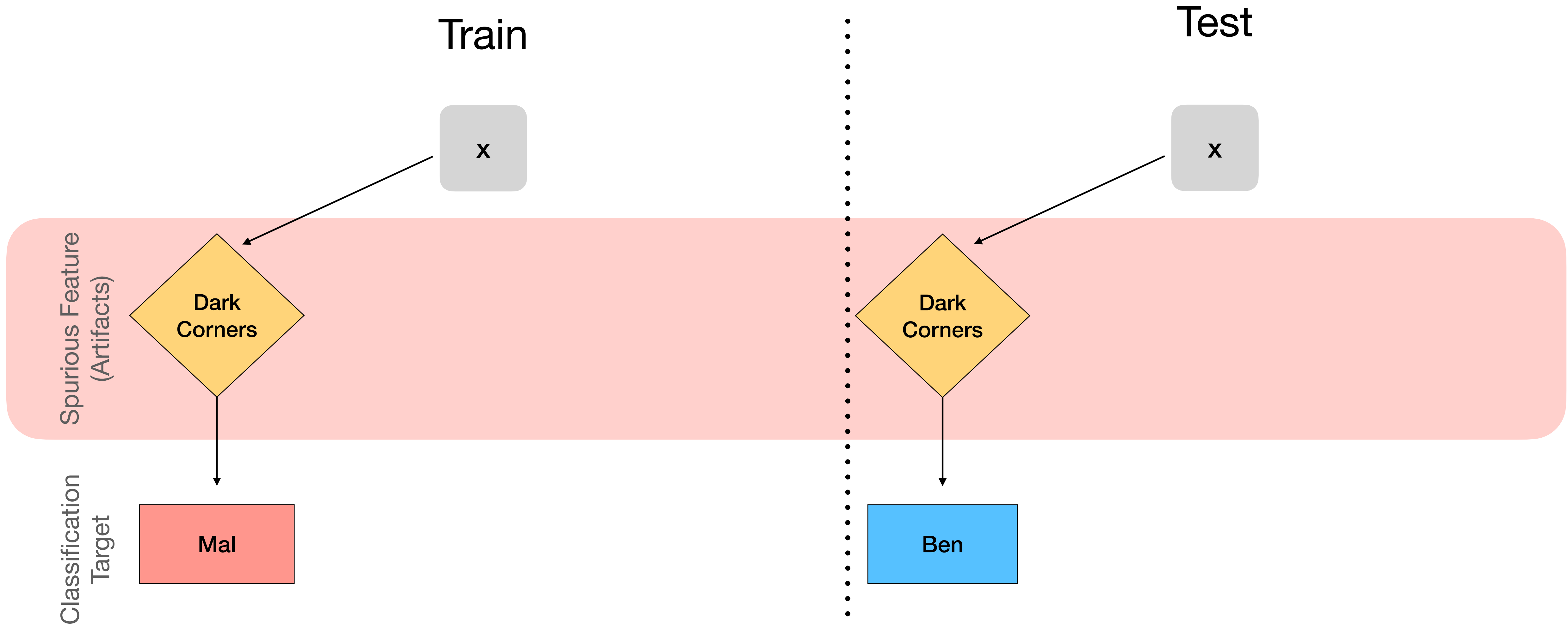
# Trap Sets - Controllable known biases



# Trap Sets - Controllable known biases

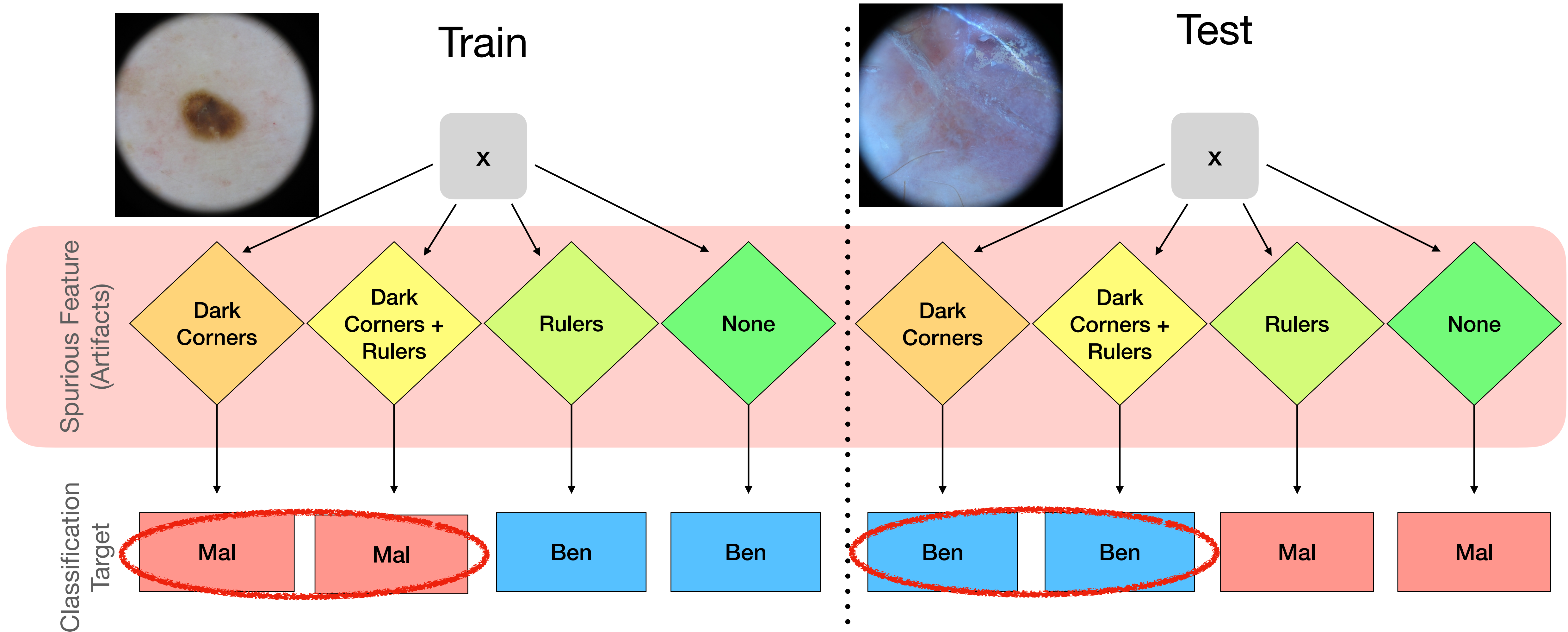


# Trap Sets - Controllable known biases



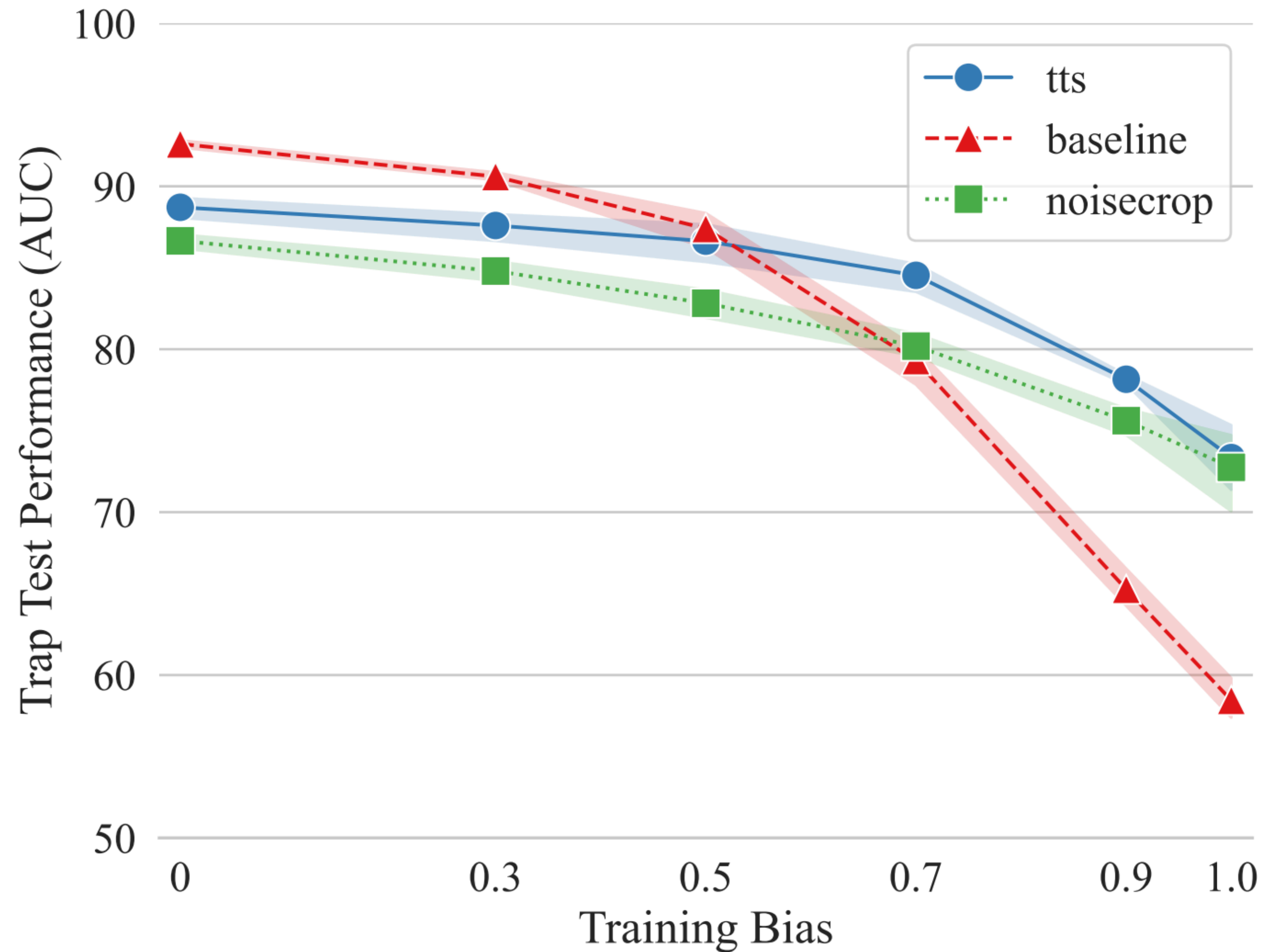


# Trap Sets - Controllable known biases



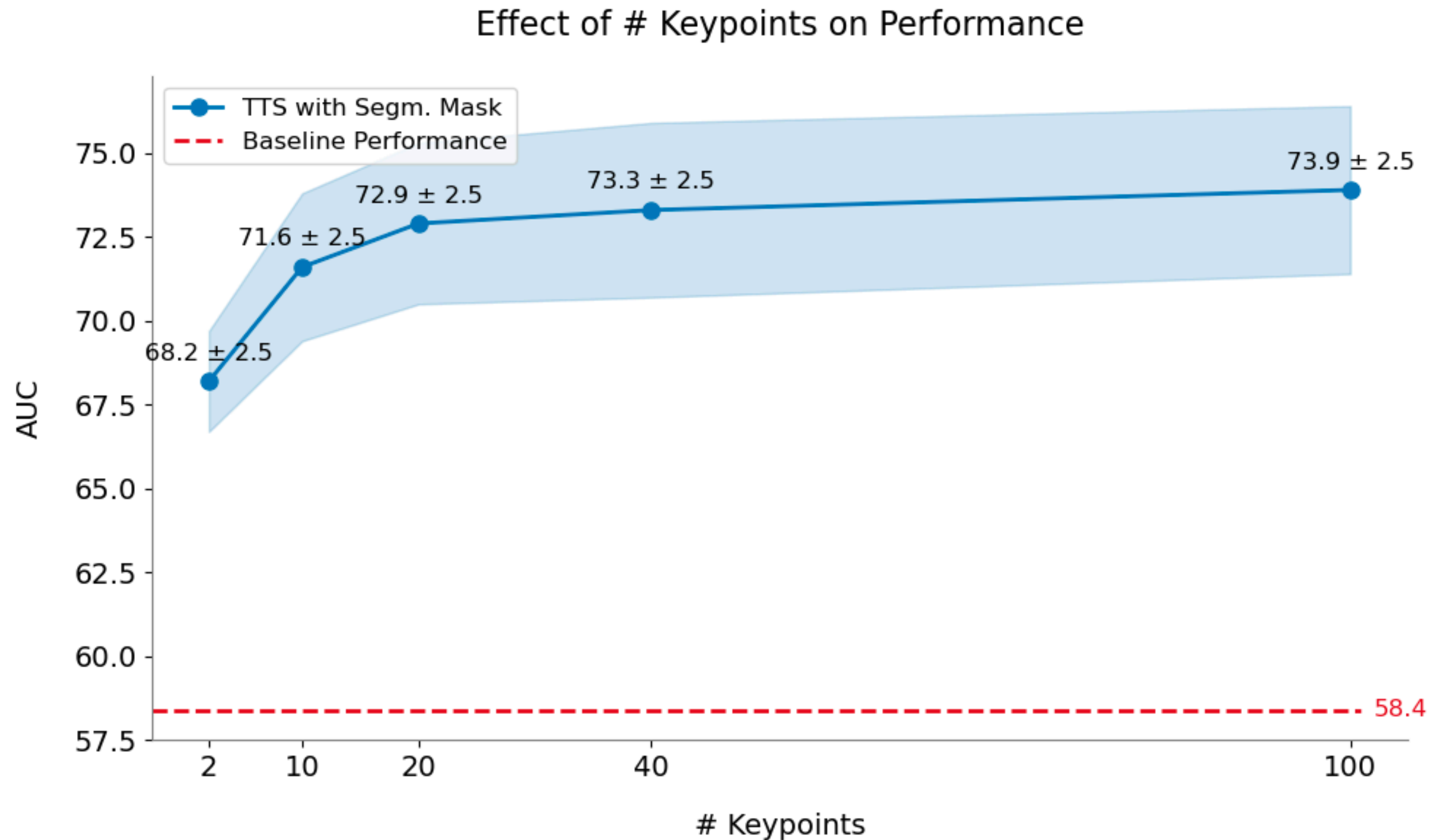
# Results

# Effective throughout different training biases

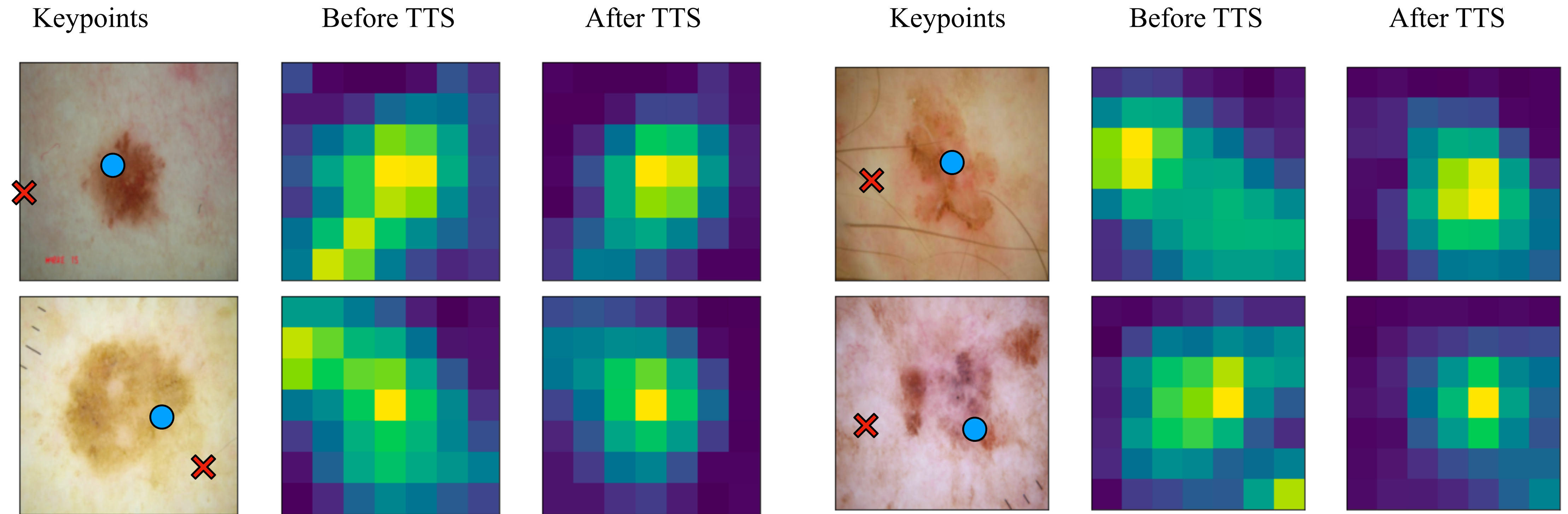




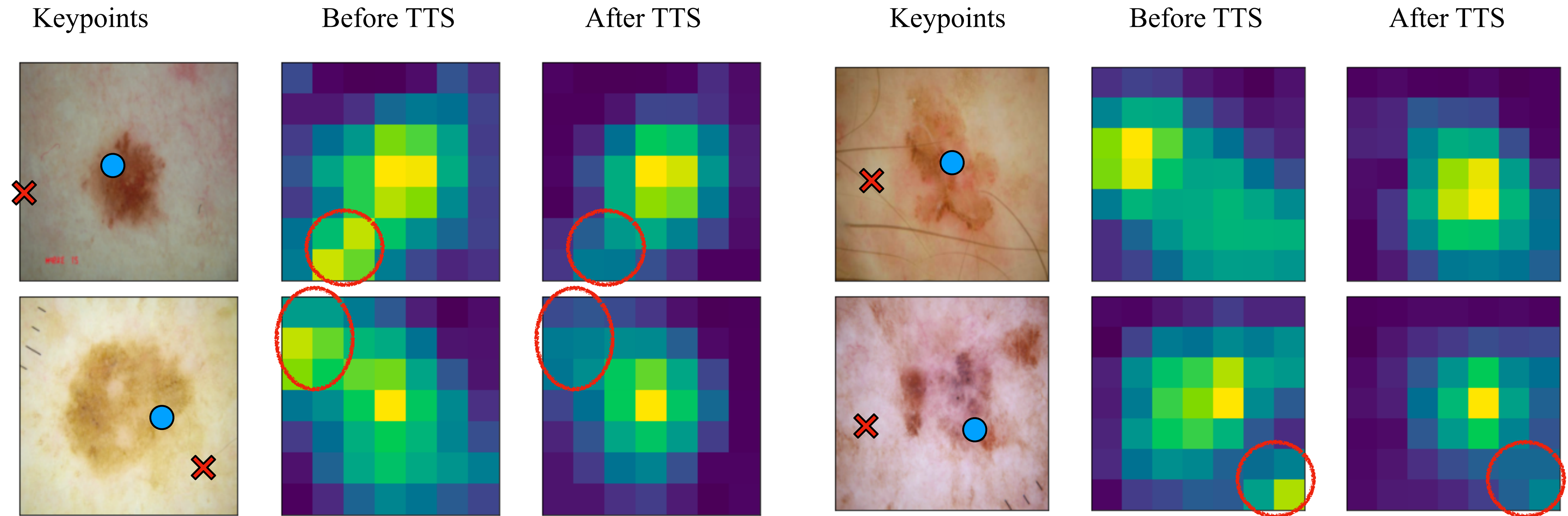
# Effective even with a single pair of keypoints



# Visualization



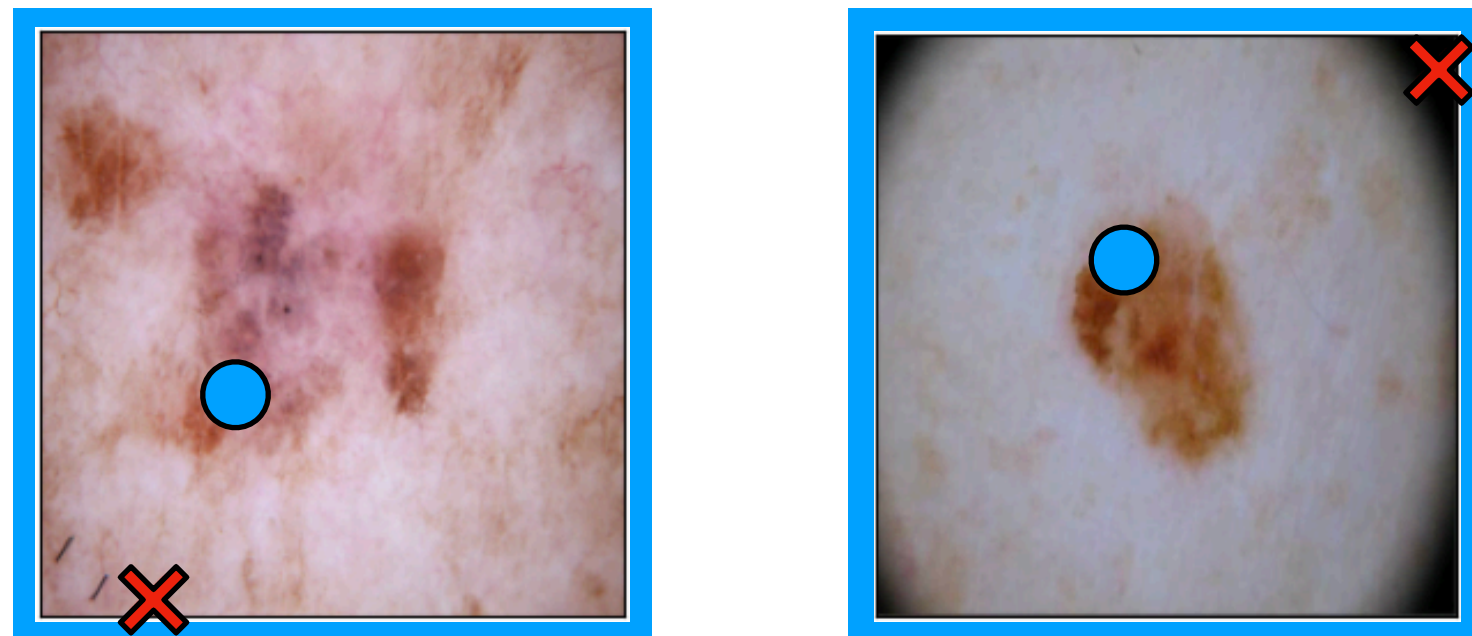
# Visualization





# Flexible for different types of annotation

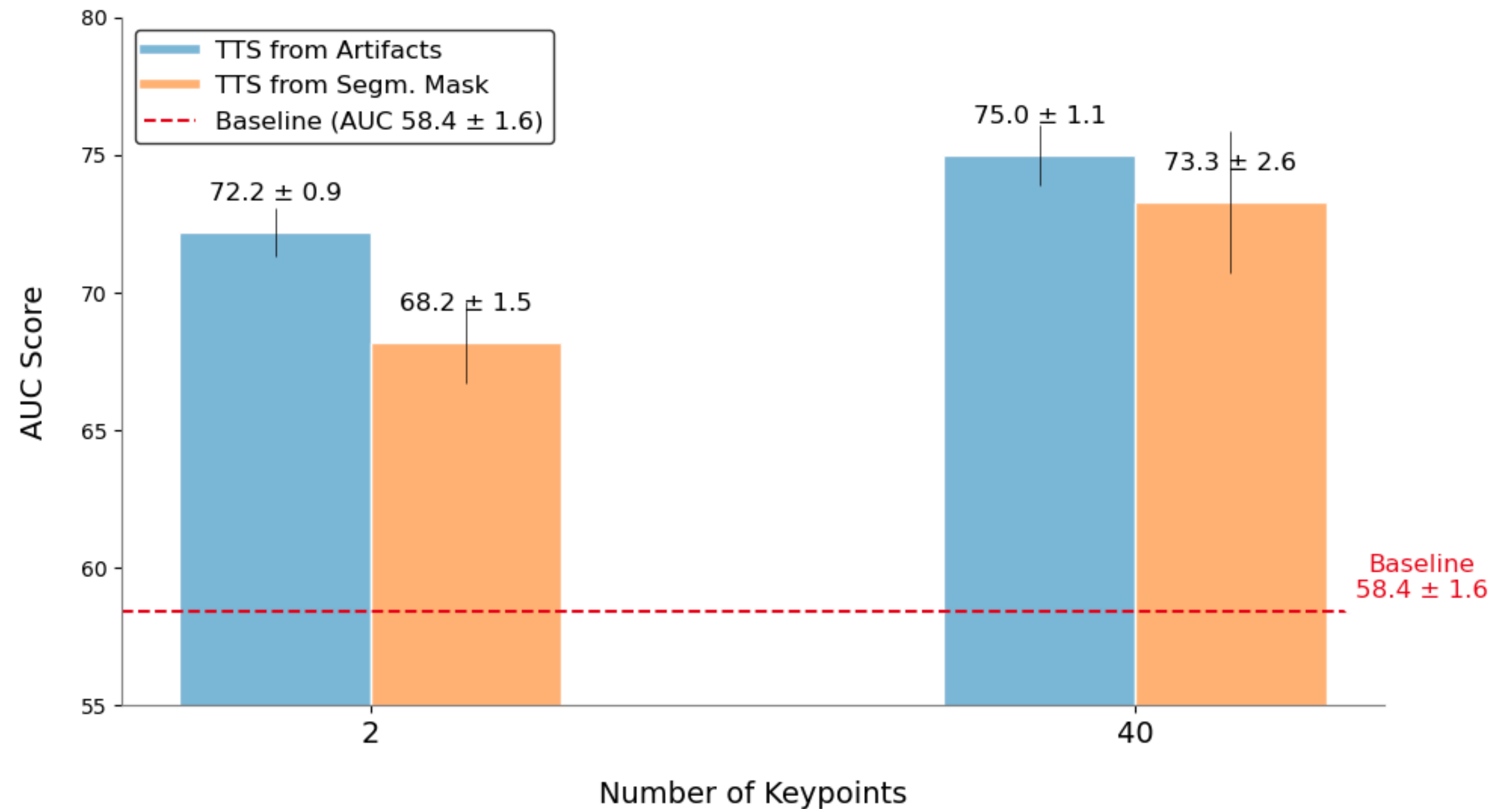
Keypoints on Artifacts



Keypoints from Segmentation Mask



Comparison of TTS Methods for Different Number of Keypoints



# Limitations

- How to adapt this solution to Vision Transformers?
- How to deal with biases uniformly spread across the image? (e.g., different acquisition devices.)

# Takeaways

- Consider evaluating your models' robustness on trap sets
- TTS improves robustness across different levels of bias
- TTS is effective even with a single pair of keypoints
- TTS is flexible to different types of annotations

## Code, Data & Paper:

<https://github.com/alceubissoto/skin-tts>

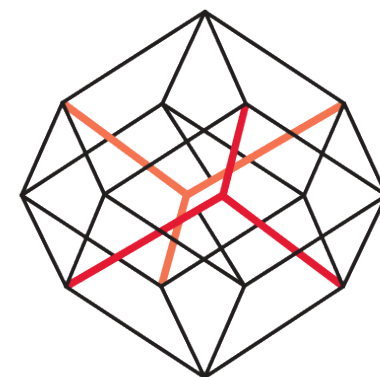
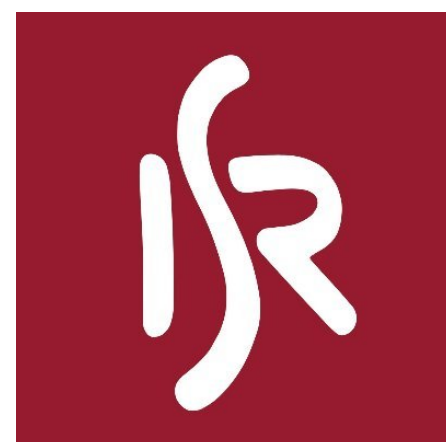
# Thank you!

**Alceu Bissoto** `alceubissoto@ic.unicamp.br`

**Catarina Barata** `ana.c.fidalgo.barata@tecnico.ulisboa.pt`

**Eduardo Valle** `dovalle@dca.fee.unicamp.br`

**Sandra Avila** `sandra@ic.unicamp.br`



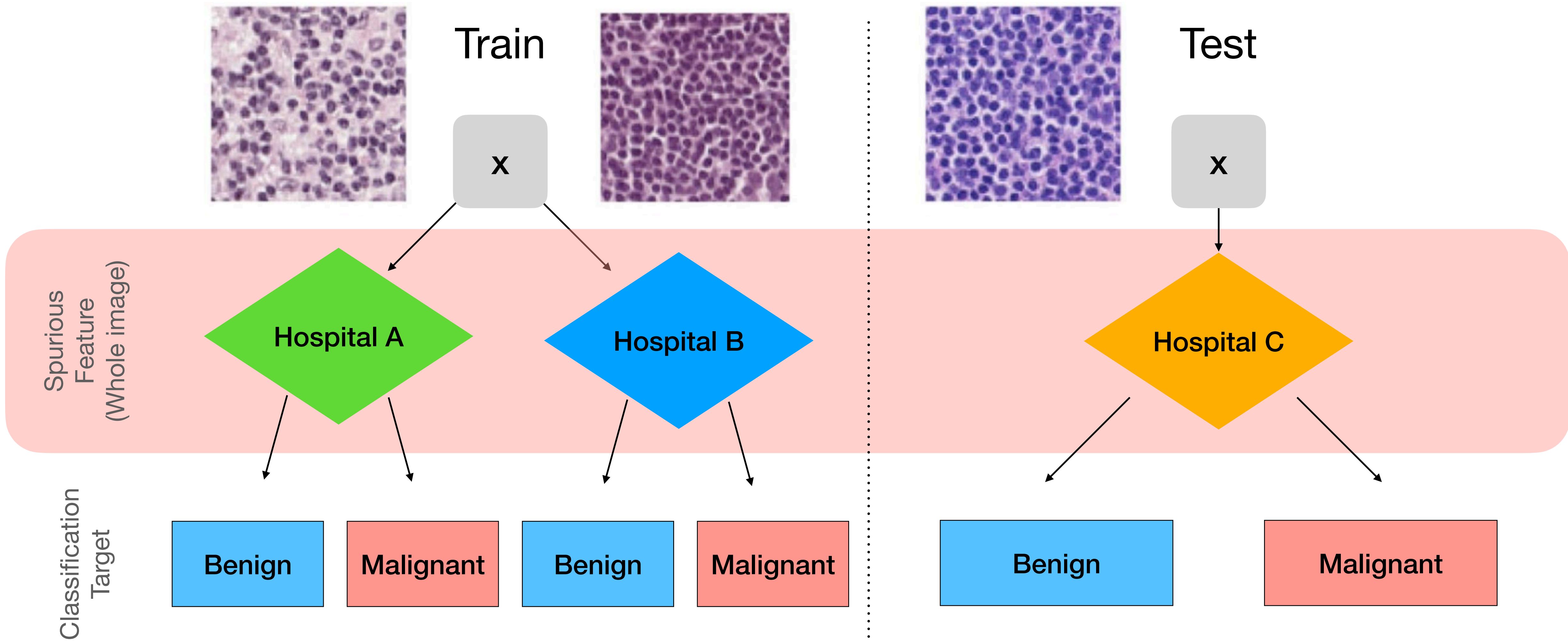
recod



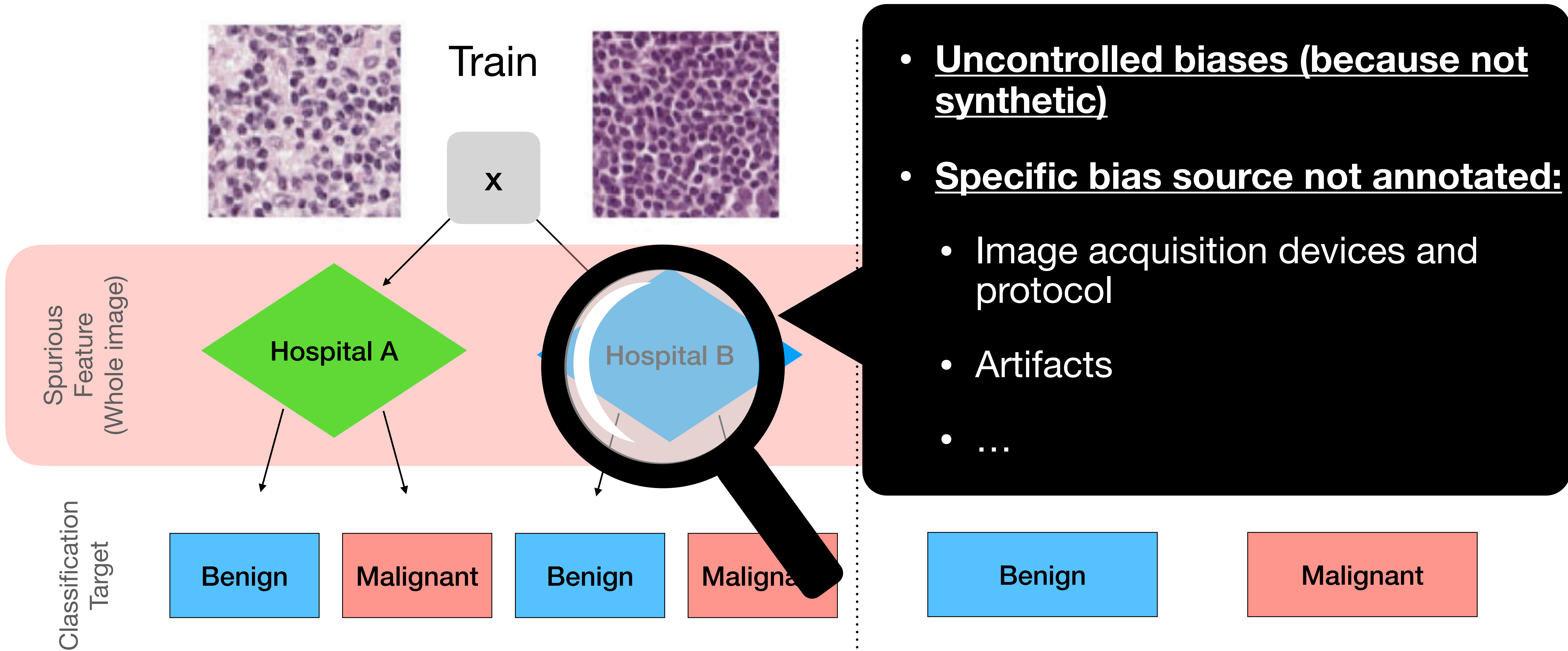
ISIC Workshop @ MICCAI 2023



# Spurious Features vs. Generalization



# Spurious Features vs. Generalization

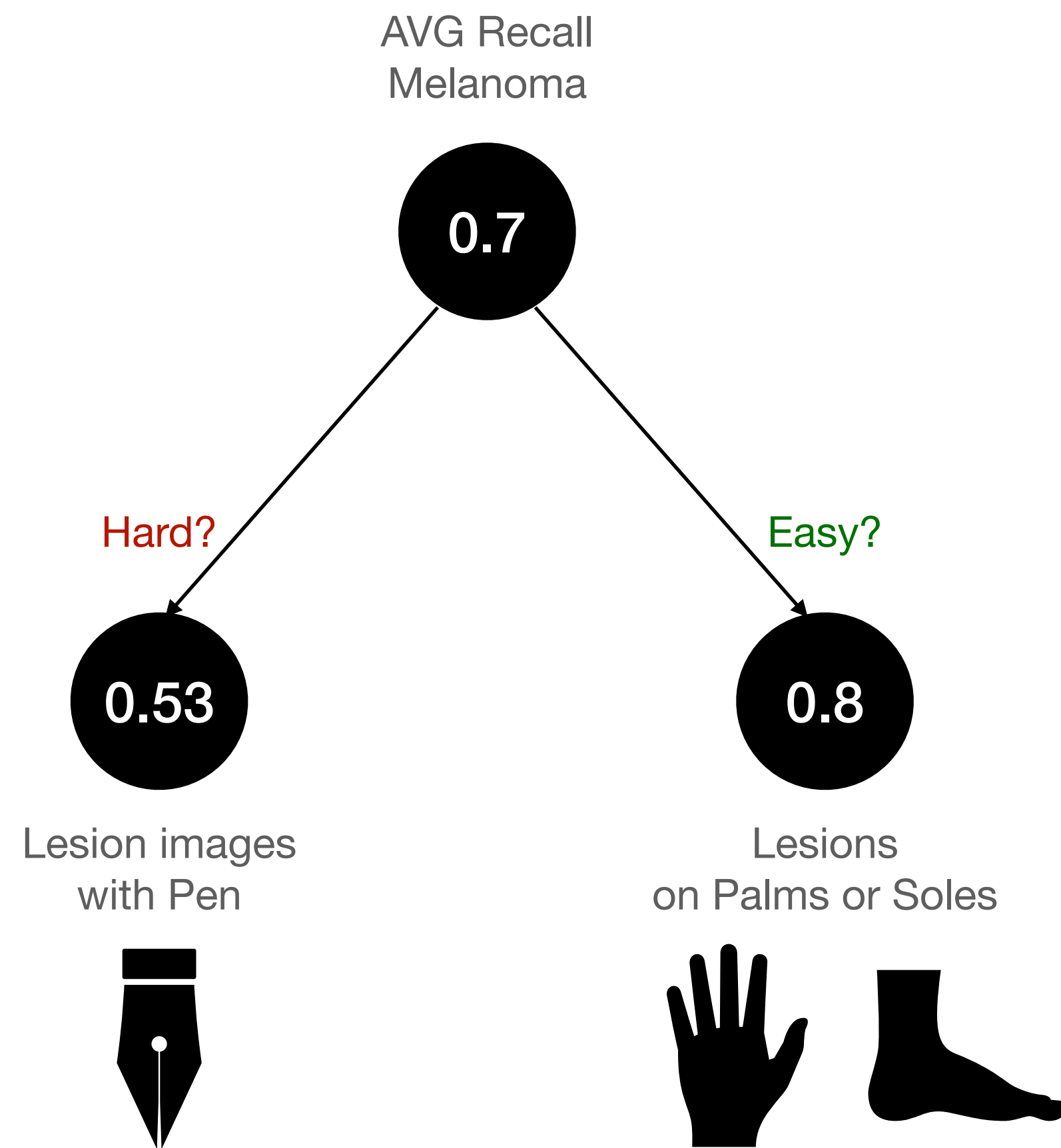






# Evaluation - Generalization and Spurious Correlations

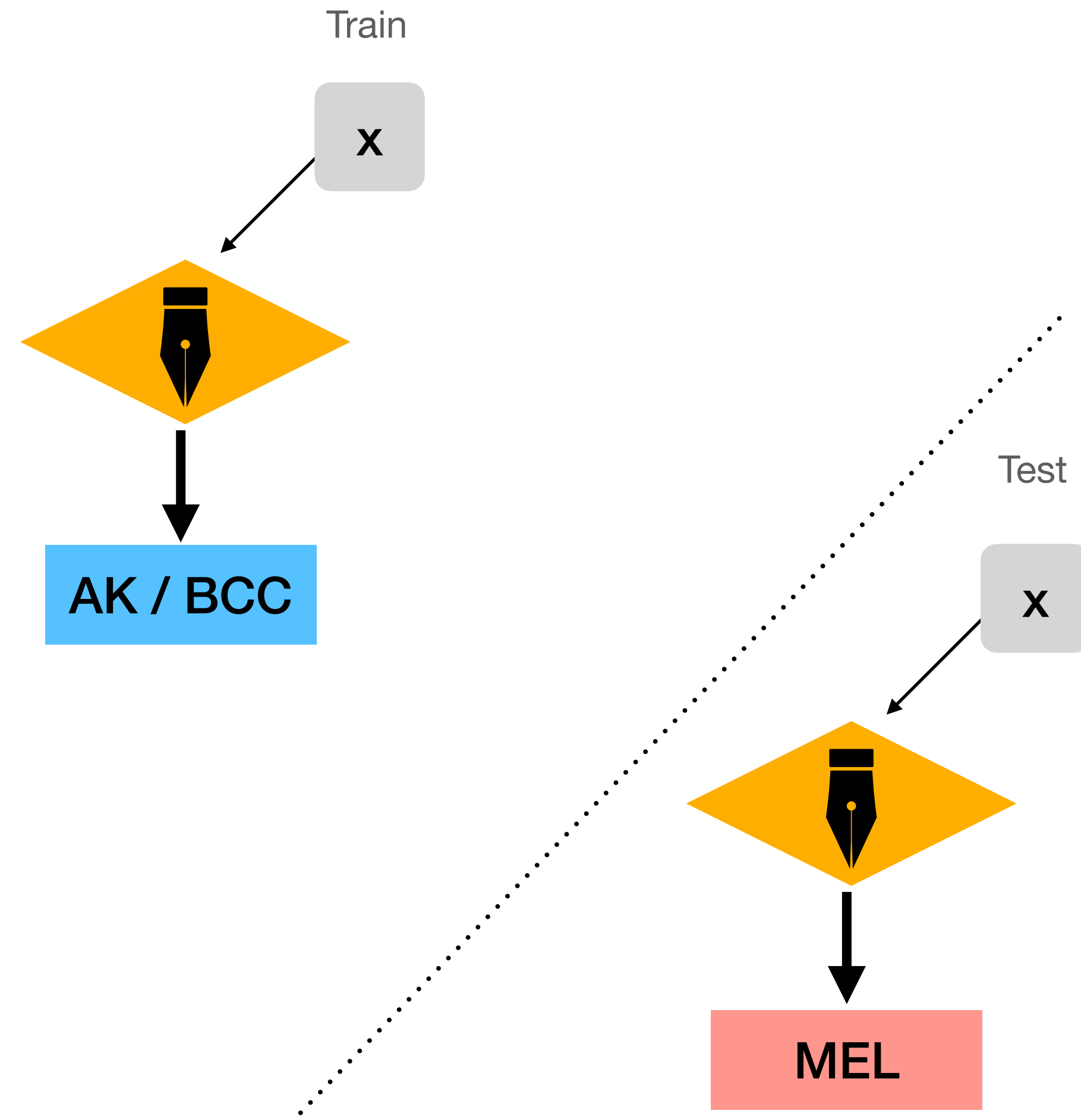
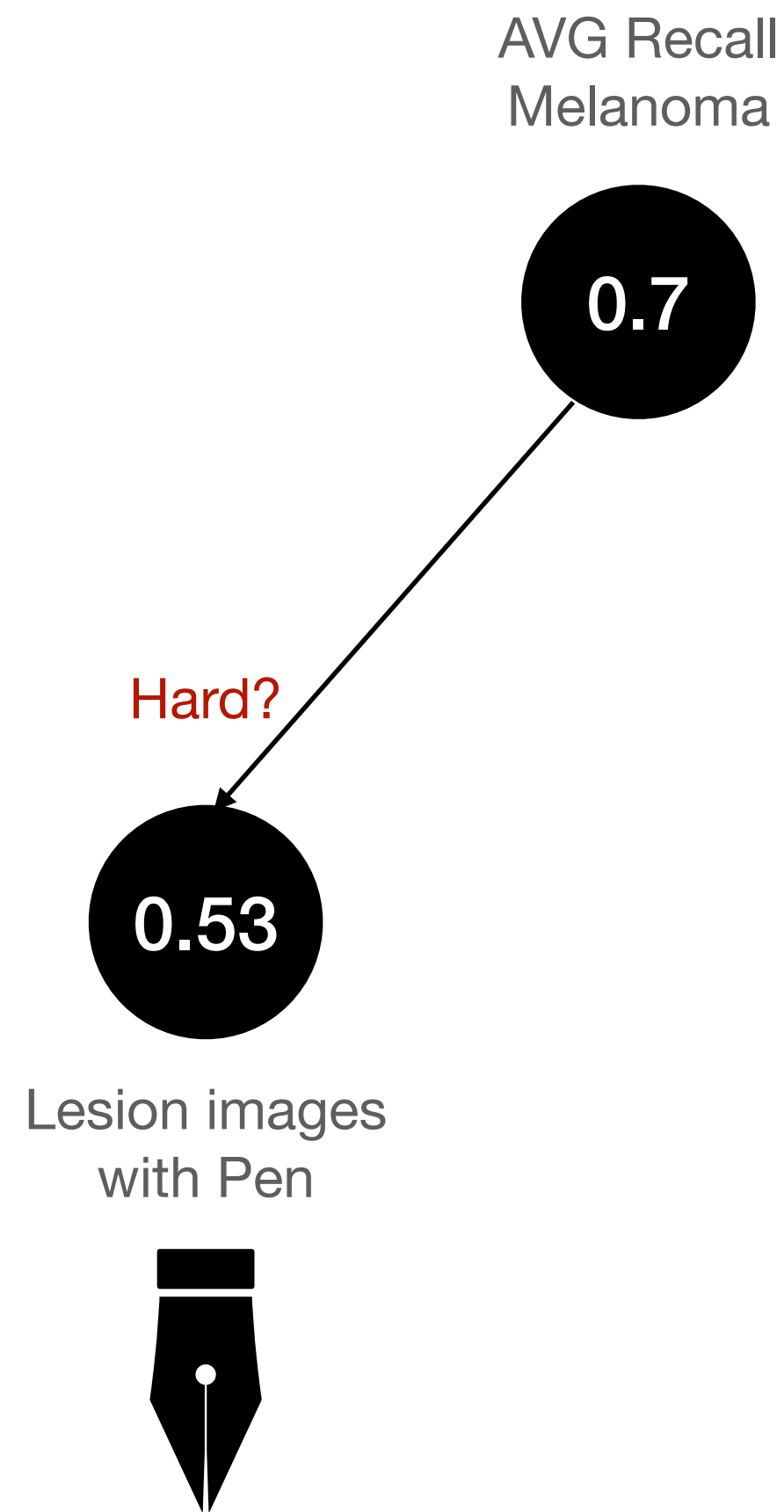
## Subgroup Evaluation





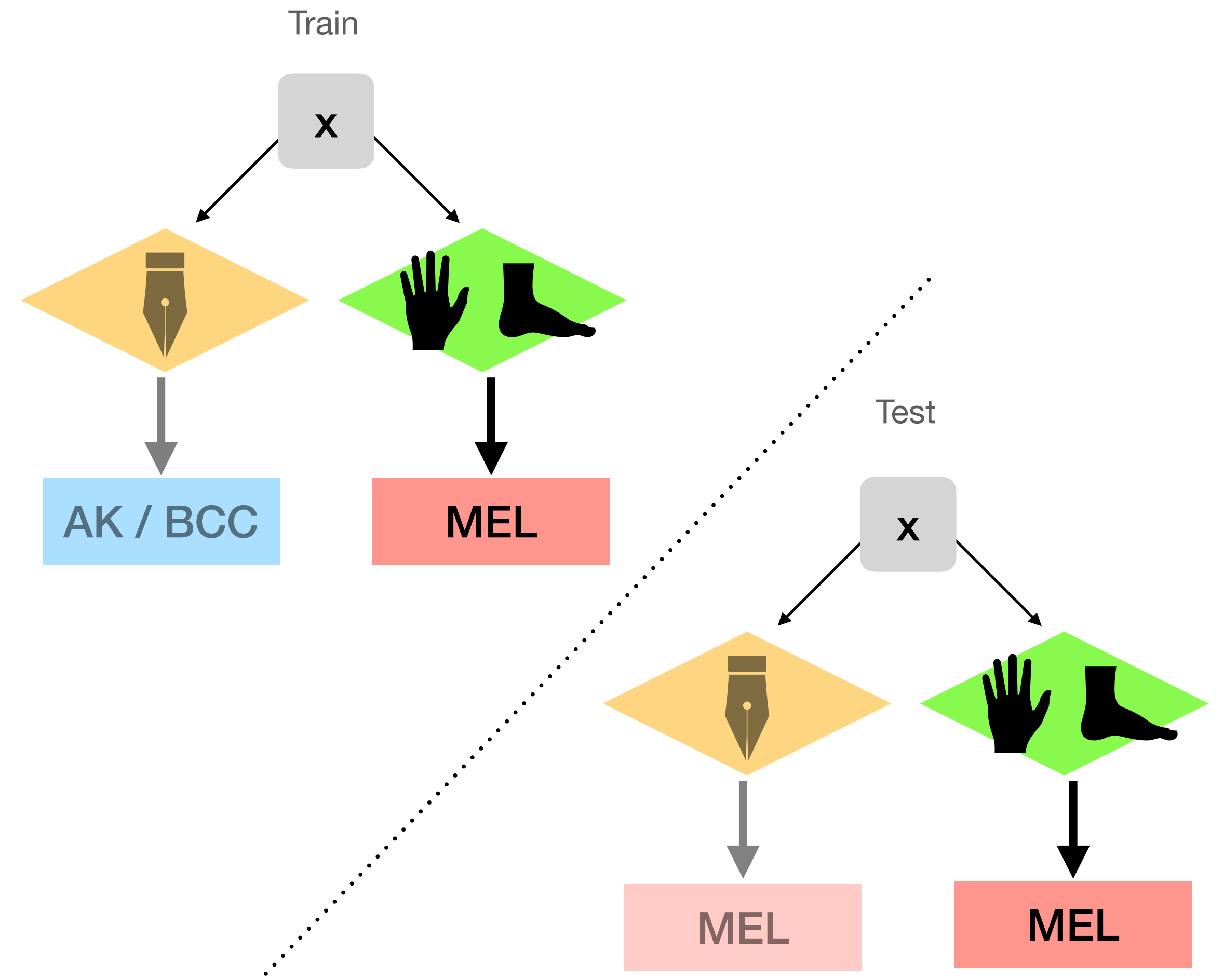
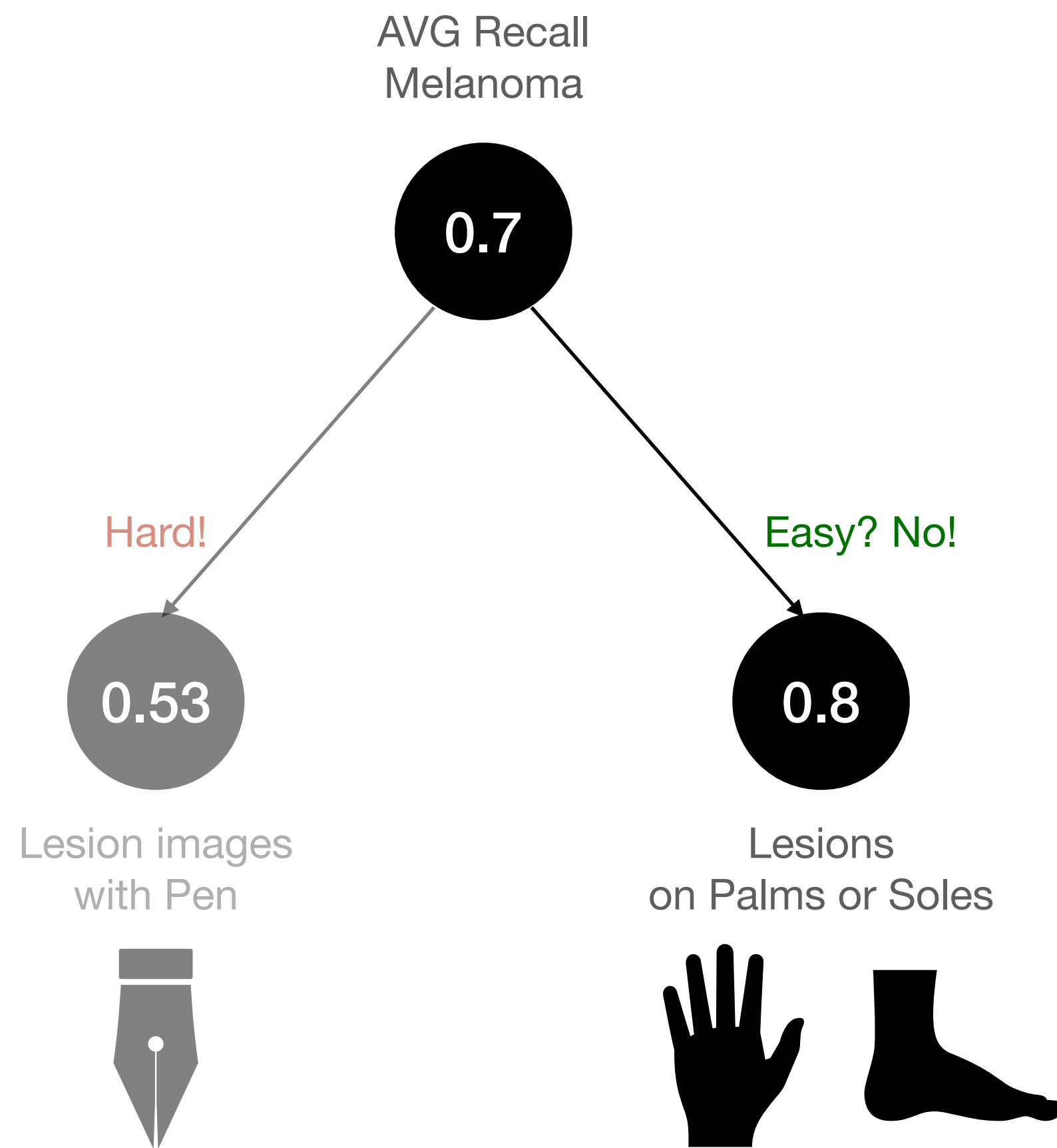
# Evaluation - Generalization and Spurious Correlations

## Subgroup Evaluation



# Evaluation - Generalization and Spurious Correlations

## Subgroup Evaluation



X