

A Call to Reflect on Evaluation Practices for Failure Detection in Image Classification

Paul F. Jaeger^{1,2}, Carsten T. Lueth^{1,2},
Lukas Klein^{1,2,3}, Till Bungert^{1,2}

¹Interactive Machine Learning Group, German Cancer Research Center (DKFZ), Heidelberg, Germany

²Helmholtz Imaging, DKFZ, Heidelberg, Germany

³Institute for Machine Learning, ETH Zürich, Zürich, Switzerland

Paper:



Benchmark:

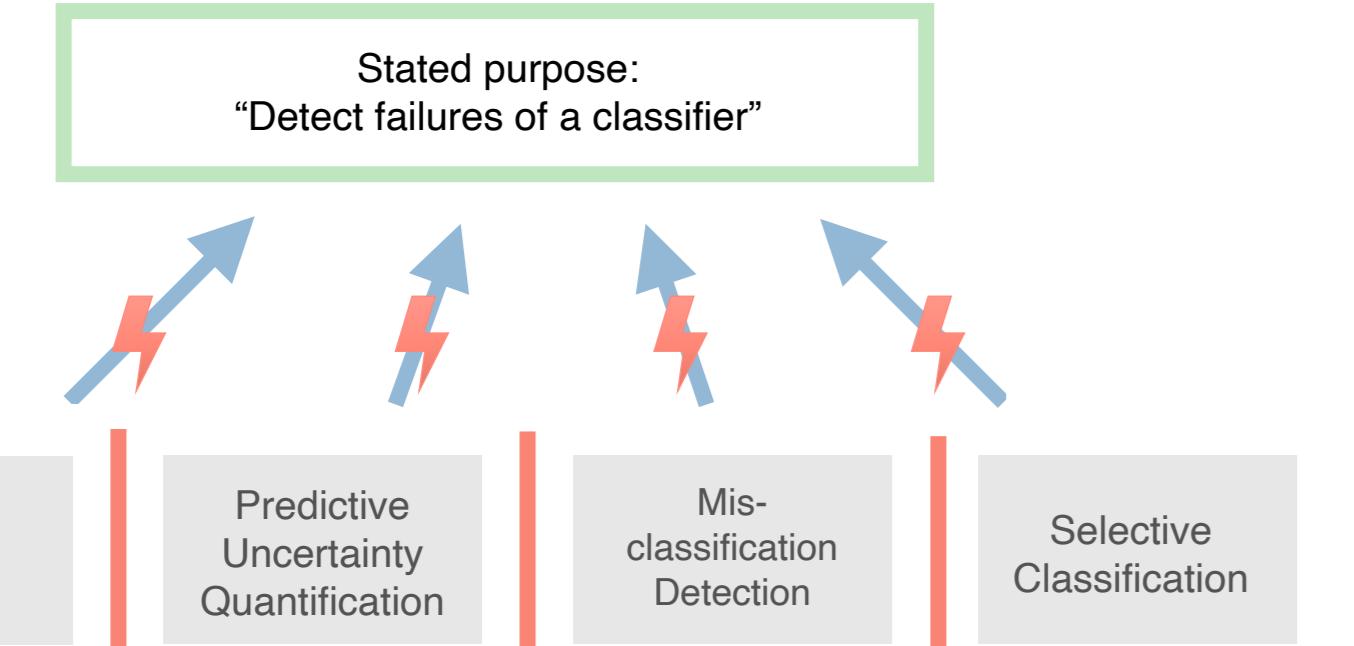


HELMHOLTZ
IMAGING

IML
Interactive Machine Learning Group
at Helmholtz Imaging and the DKFZ

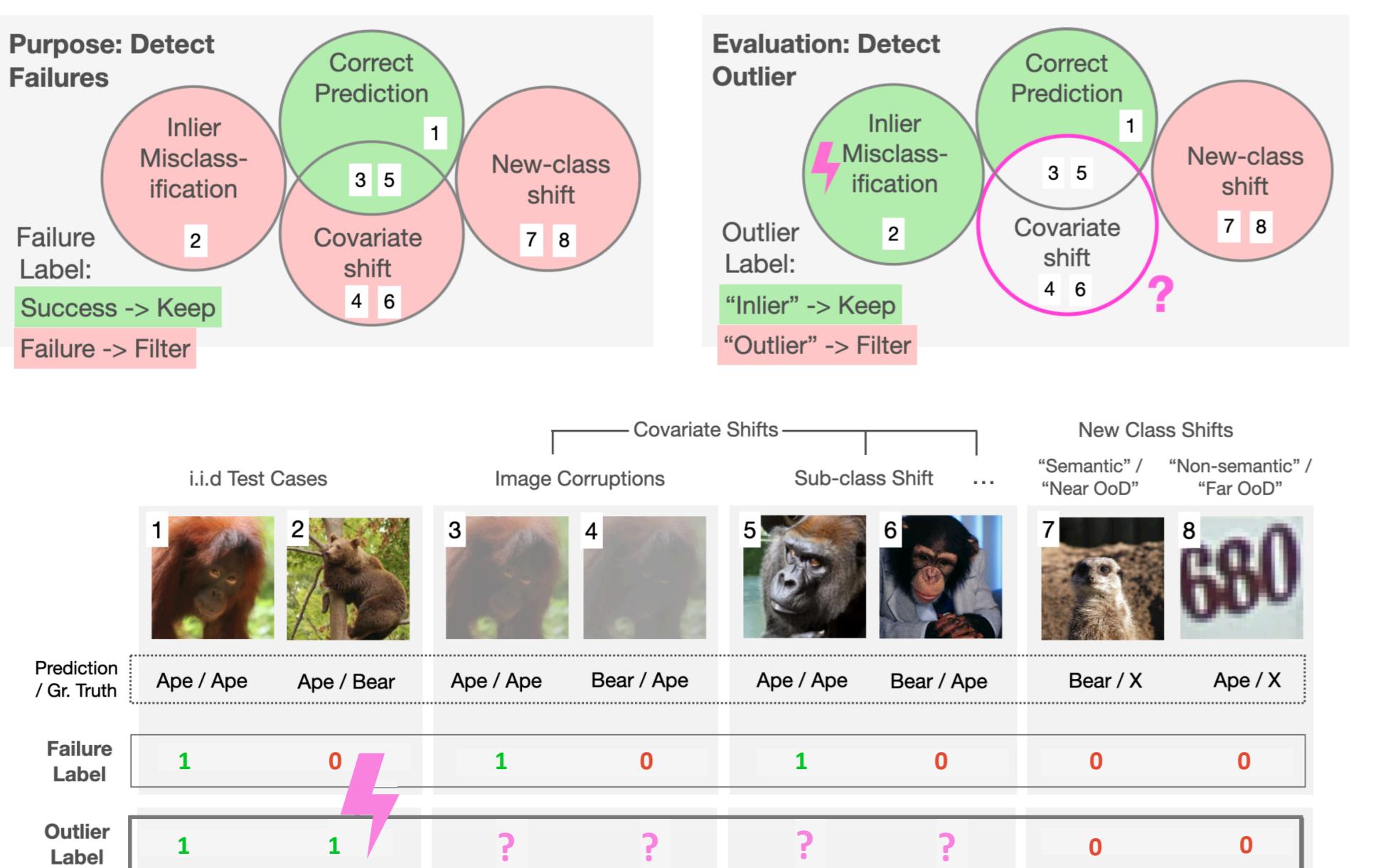
dkfz.
GERMAN
CANCER RESEARCH CENTER
IN THE HELMHOLTZ ASSOCIATION

Status Quo: Inconsistent landscape of confidence scoring



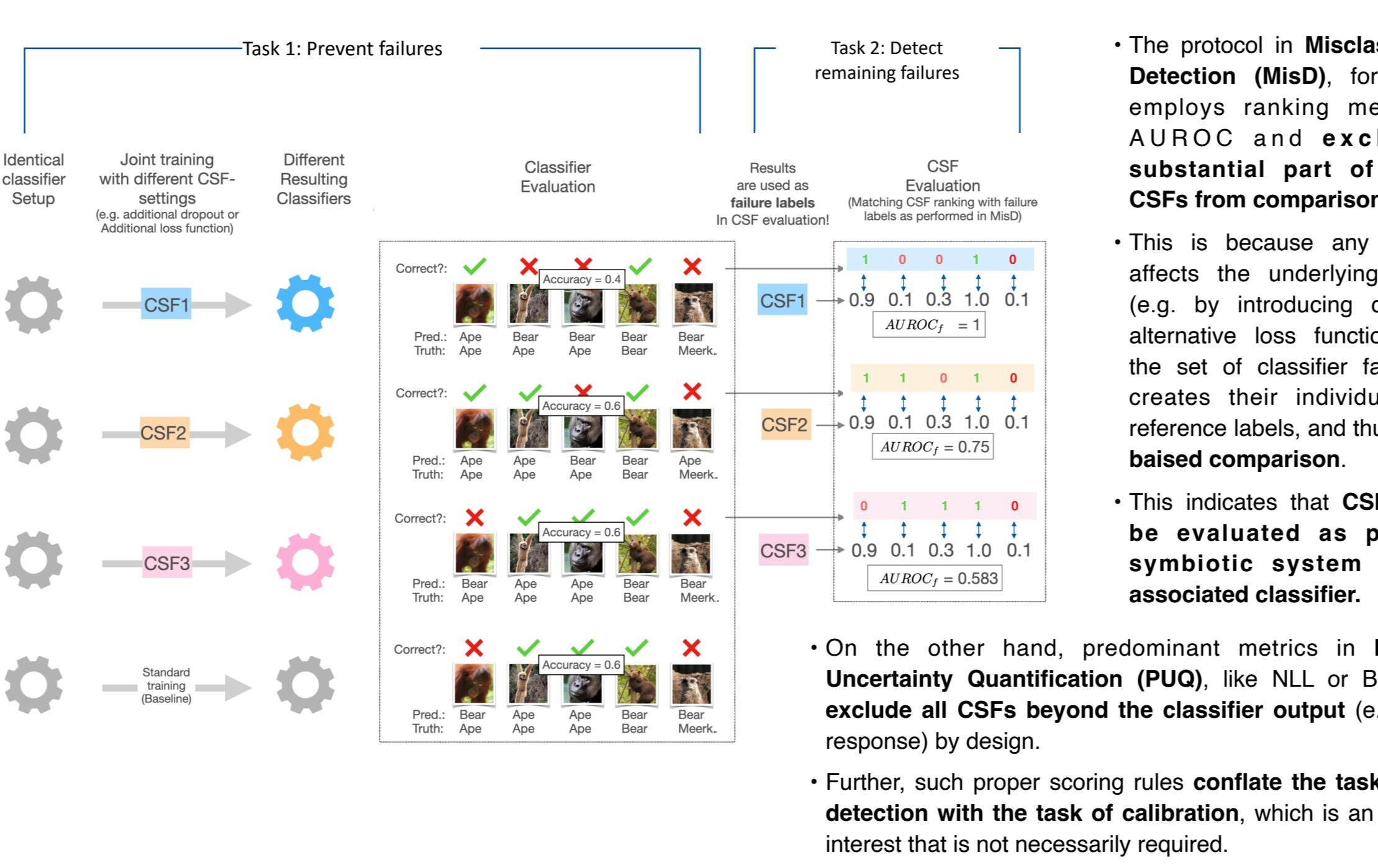
- Motivational statements in many publications of these four research areas (OoD-D, PPU, MisD, SC) indicate that all respective studies address the same goal of using confidence scoring functions (CSFs) to detect failures of a classifier.
- In our work, we make the case that the evaluation in these fields is substantially flawed for two reasons:
 1. There is a discrepancy between this stated purpose and the way methods are evaluated.
 2. All these fields are currently mostly siloed meaning that there is no cross-comparison of methods, although they address the same goal.

Pitfall 1: OoD-Detection often deviates from its stated purpose



- In Out-of-distribution Detection, indeed, the majority of publications states failure detection in a well-defined classification task as their purpose.
- In the OoD-D task protocol, however, an outlier label is employed instead of the failure status of the classifier, only aiming to determine whether cases are subject to a new-class shift or not.
- Thus, for one, a CSF is rewarded for giving high confidence to all inliers (purple lightning), including failures, but perhaps even more concerning, the subjective outlier label is not clearly defined on the covariate shifts (purple question marks).
- It could be argued that these more subtle shifts where the image label is preserved are the more realistic and thus more relevant ones.

Pitfall 2: Current evaluation metrics lead to biased and incomplete comparison



- The protocol in Misclassification Detection (MisD), for instance, employs ranking metrics like AUCROC and excludes a substantial part of relevant CSFs from comparison

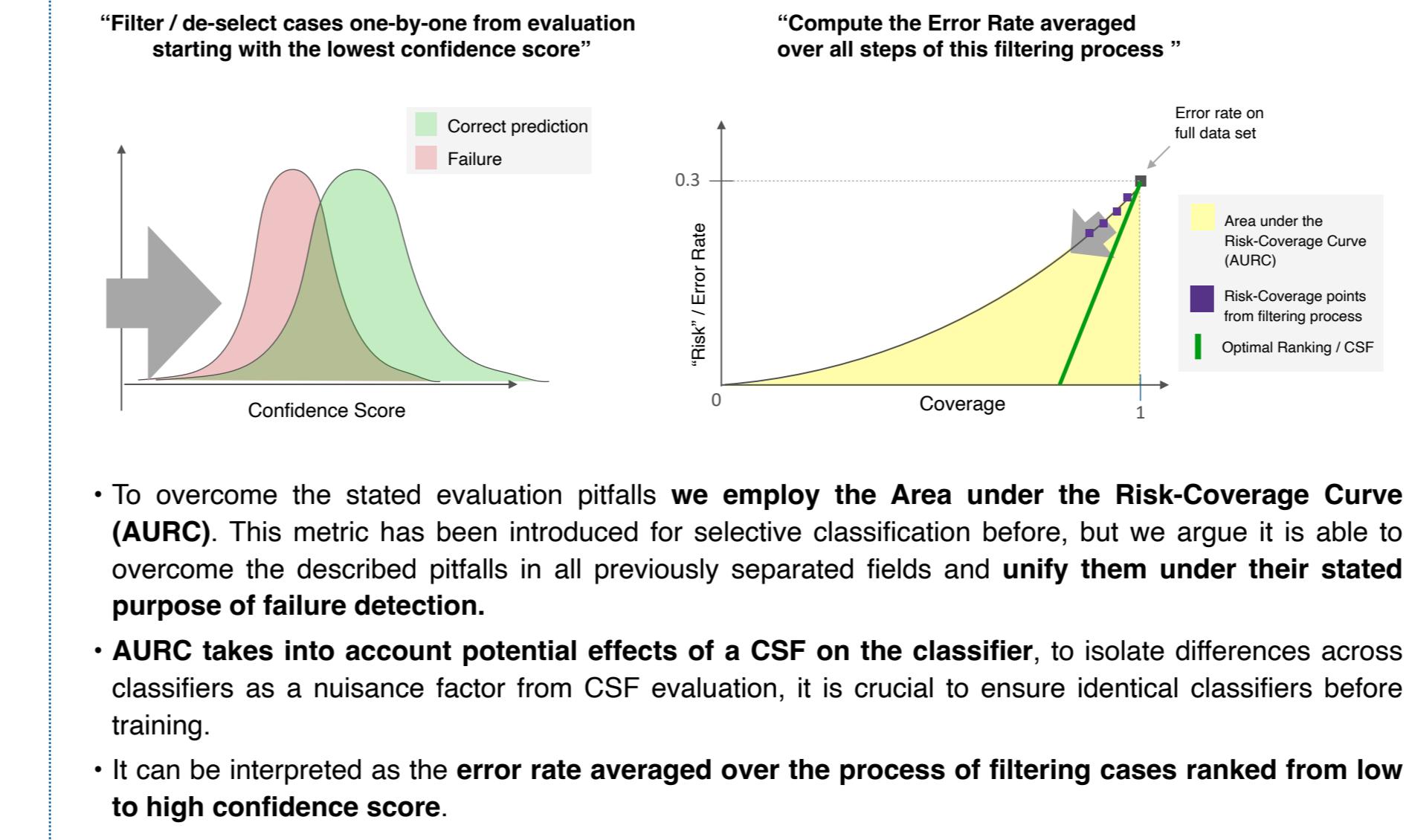
- This is because any CSF that affects the underlying classifier (e.g. by introducing dropout or alternative loss functions) alters the set of classifier failures, i.e. creates their individual set of reference labels, and thus leads to biased comparison.

- This indicates that CSFs should be evaluated as part of a symbiotic system with the associated classifier.

- On the other hand, predominant metrics in Predictive Uncertainty Quantification (PUQ), like NLL or Brier Score, exclude all CSFs beyond the classifier output (e.g. softmax response) by design.

- Further, such proper scoring rules conflate the task of failure detection with the task of calibration, which is an orthogonal interest that is not necessarily required.

The Area under the Risk-Coverage Curve overcomes previous evaluation pitfalls



- To overcome the stated evaluation pitfalls we employ the Area under the Risk-Coverage Curve (AURC). This metric has been introduced for selective classification before, but we argue it is able to overcome the described pitfalls in all previously separated fields and unify them under their stated purpose of failure detection.
- AURC takes into account potential effects of a CSF on the classifier, to isolate differences across classifiers as a nuisance factor from CSF evaluation, it is crucial to ensure identical classifiers before training.
- It can be interpreted as the error rate averaged over the process of filtering cases ranked from low to high confidence score.

FD-Shifts insights open up new research directions

- "None of the evaluated methods from literature beats the simple Maximum Softmax Response baseline across a realistic range of failure sources."
 - "Prevalent OoD-D methods are only relevant in a narrow range of distribution shifts."
 - "AURC is able to resolve previous obscurities between classifier robustness and CSF performance."
 - "AURC beats ViT on the iWildCam benchmark, indicating interesting transfer-learning issues."
 - "Different types of uncertainty are empirically not distinguishable."
 - "CSFs beyond Maximum Softmax Response yield well-calibrated scores."
 - "The Maximum Softmax Response baseline is disadvantaged by numerical errors in the standard setting."
- Great demand for next generation of robust CSFs!
- Deeper understanding of uncertainty modeling in practice required!
- Research perspective: Calibrated confidence beyond Softmax Response

Easter Egg finding: The Softmax baseline is often disadvantaged by numerical errors in ranking metrics like AUROC

	Round-to-one error rate \downarrow	AURC \downarrow	AUROC \uparrow	Accuracy \uparrow
iWildCam	16bit 32bit 64bit	80.22 69.09 68.80	87.50 92.50	76.01
BREEDS	35.53 42.06 0.00	18.81 12.89 12.86	89.89 92.22 92.22	90.72
CNN	35.53 42.06 0.00	18.81 12.89 12.86	89.89 92.22 92.22	93.99
iCIFAR-10	23.75 23.04 0.00	67.27 71.43 71.39	83.39 87.29 87.29	73.26
CIFAR-10	41.54 14.65 0.00	8.36 6.20 5.67	71.98 93.73 93.73	94.35
SVHN	41.76 17.29 0.00	8.07 4.902 4.850	89.59 92.81 92.87	96.09
iWILDS	44.41 14.91 0.00	22.16 17.78 17.70	75.57 80.35 80.38	62.12
BREEDS	80.19 52.99 0.423	11.43 4.559 1.893	72.65 88.88 94.35	97.92
CNN	80.19 52.99 0.423	11.43 4.559 1.893	72.65 88.88 94.35	97.95
iCIFAR-10	68.65 30.73 0.00	36.27 14.95 14.27	79.29 90.10 90.29	91.62
CIFAR-10	92.16 81.79 1.883	7.614 3.480 0.950	80.58 88.76 90.50	98.76
SVHN	69.02 47.17 0.305	16.94 7.57 5.475	68.75 83.55 88.14	97.30

- Depending on floating point precision, rounding errors occur during the softmax operation thereby losing the ranking information between rounded scores.
- Especially on the ViT classifier, these errors occur at astounding rates leading to substantial ranking performance drops as measured e.g. by AUROC.
- Even at default 32-bit precision, this effect leads to a substantial disadvantage of softmax baselines in all ranking tasks including OoD-Detection.

Hands-on recommendations for evaluating confidence scoring

- State a clear purpose of the confidence scoring function (CSF) and design an evaluation protocol that reflects this purpose.
- If the purpose is failure detection, we recommend AURC as primary metric for method comparison.
- Analogously to classifier robustness, CSFs need to be tested on a wide range of data sets and distribution shifts.
- Compare against all viable solutions addressing the same goal, even if from seemingly separated fields.
- Logits should be cast to 64-bit precision or temperature-scaled prior to the softmax operation for any ranking-related tasks to avoid subpar softmax baselines.