

# Optimization For Robustness Evaluation Beyond $l_p$ Metrics



Hengyue Liang, Tiancong Chen, Buyun Liang, Ying Cui, Tim Mitchell, Ju Sun



## Robustness Evaluation Problems

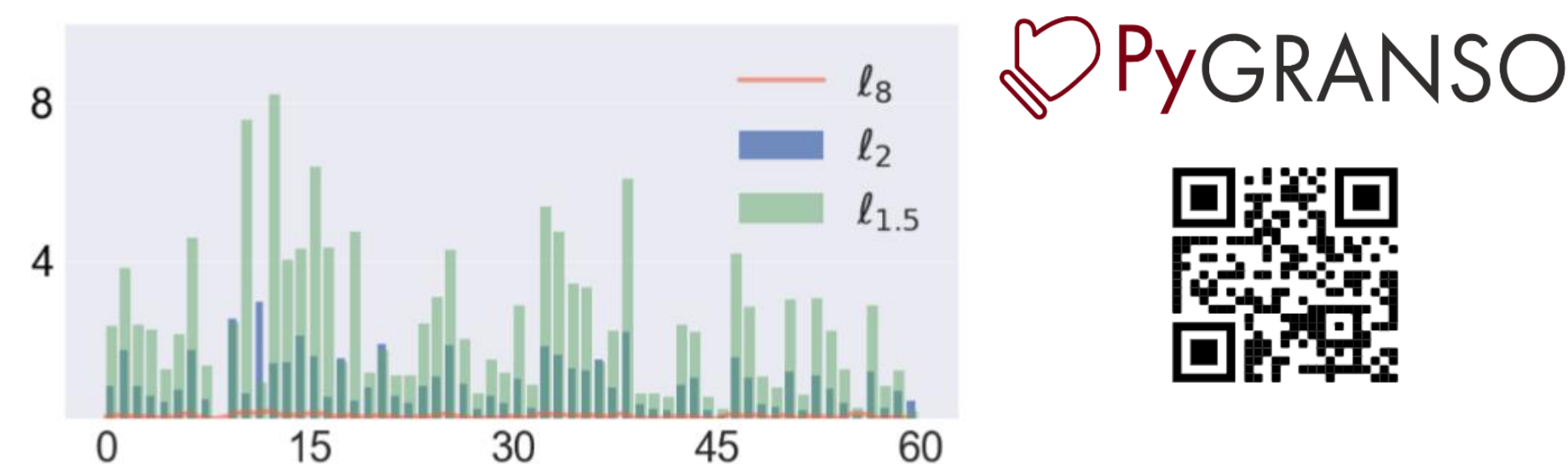
Current (adversarial) robustness evaluation of neural networks are generally formulated as solving the following two forms of **constrained optimization problems**:

- Finding an **adversarial perturbation** via max form:
 
$$\max_{x'} \ell(y, f_{\theta}(x'))$$
 s.t.  $d(x, x') \leq \epsilon, \quad x' \in [0, 1]^n$
- Finding the **robustness radius** via min form:
 
$$\min_{x'} d(x, x')$$
 s.t.  $\max_{i \neq y} f_{\theta}^i(x') \geq f_{\theta}^y(x'), \quad x' \in [0, 1]^n$

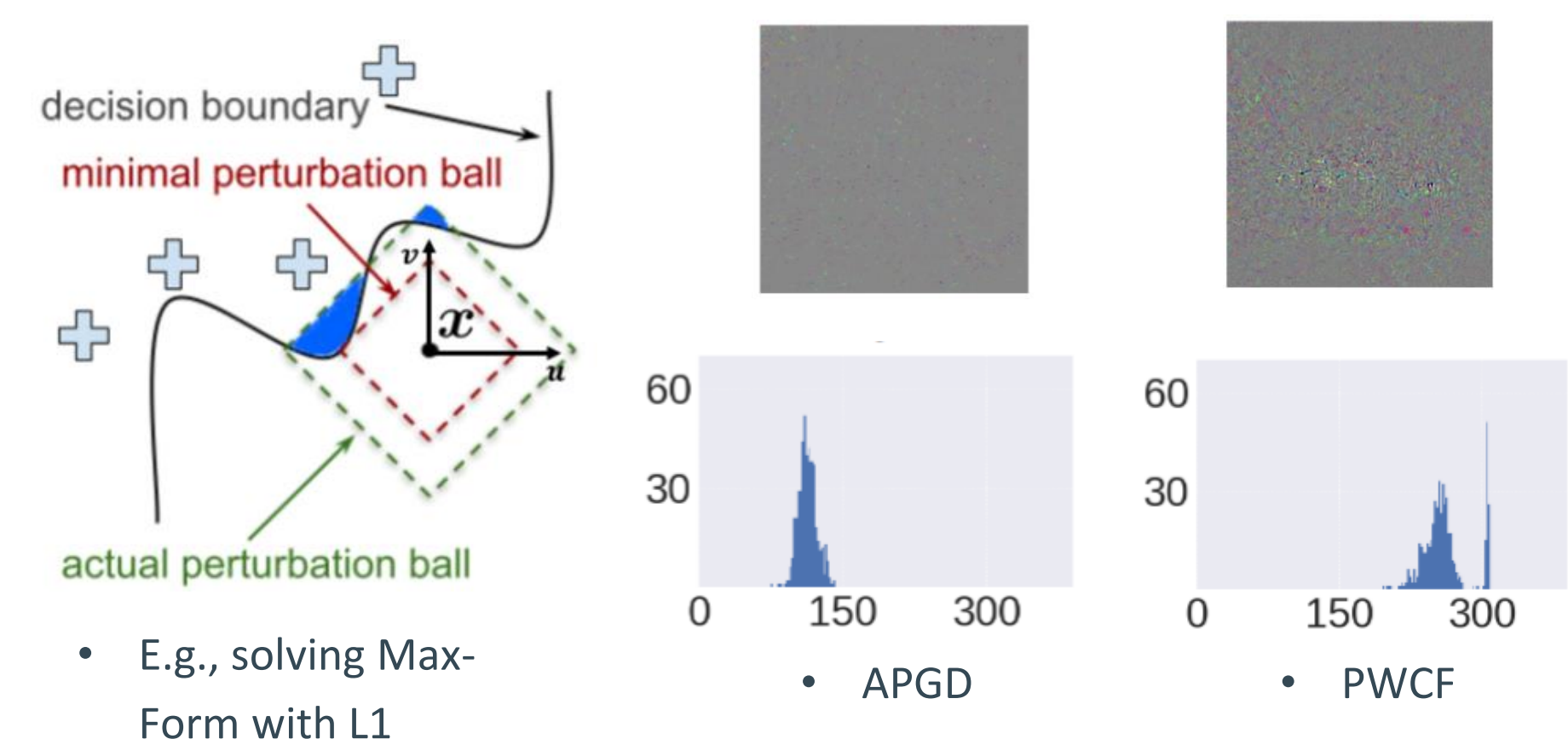
## PyGRANSO with Constraint-Folding (PWCF)

- General purposed non-linear optimization solver
- Can handle non-smooth functions
- With GPU-acceleration --- Deep Learning OK

- Can solve both formulations with general distance metrics with high quality. E.g., min form with L8 and L1.5 distances below:



## Solution patterns depending on solvers used

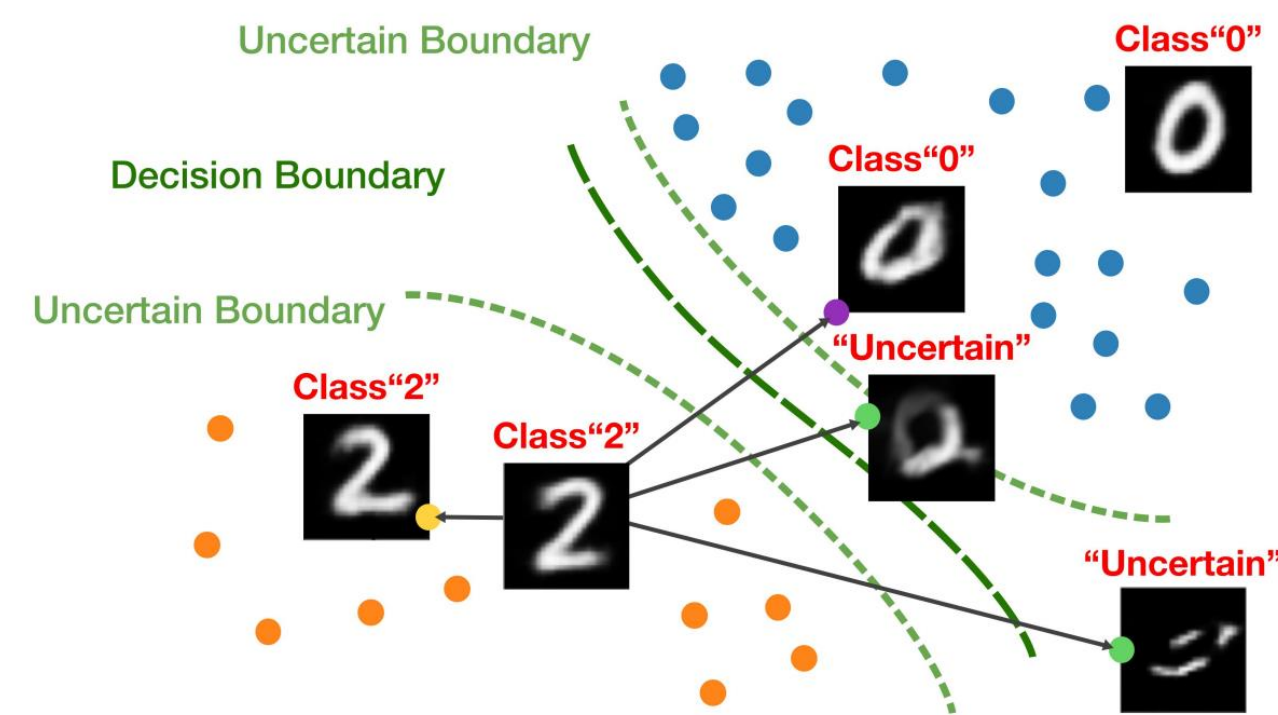


What this may imply:

- Current robustness evaluation may be insufficient and misleading
- Abs robustness may be hard to achieve

[1] Liang H, Liang B, Peng L, Cui Y, Mitchell T, Sun J. Optimization and optimizers for adversarial robustness. arXiv preprint arXiv:2303.13401. 2023 Mar 23.

## Uncertainty-aware Boldness



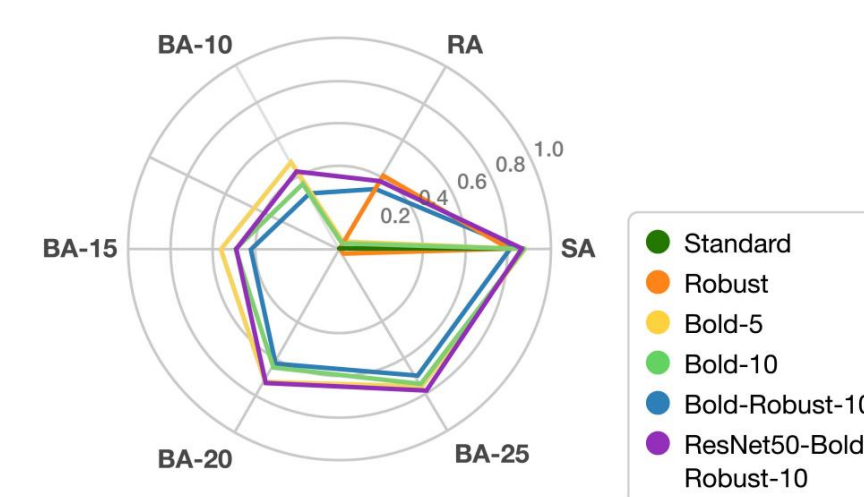
A classifier should achieve good: 1) generalizability 2) robustness 3) uncertainty-aware Boldness.

Introduce a new evaluation metric towards reliability: boldness accuracy (BA)

Existing models, including robust models, are not uncertainty-aware bold



Formulate robustness and uncertainty-aware boldness as min-max optimization problem, improve the overall performance



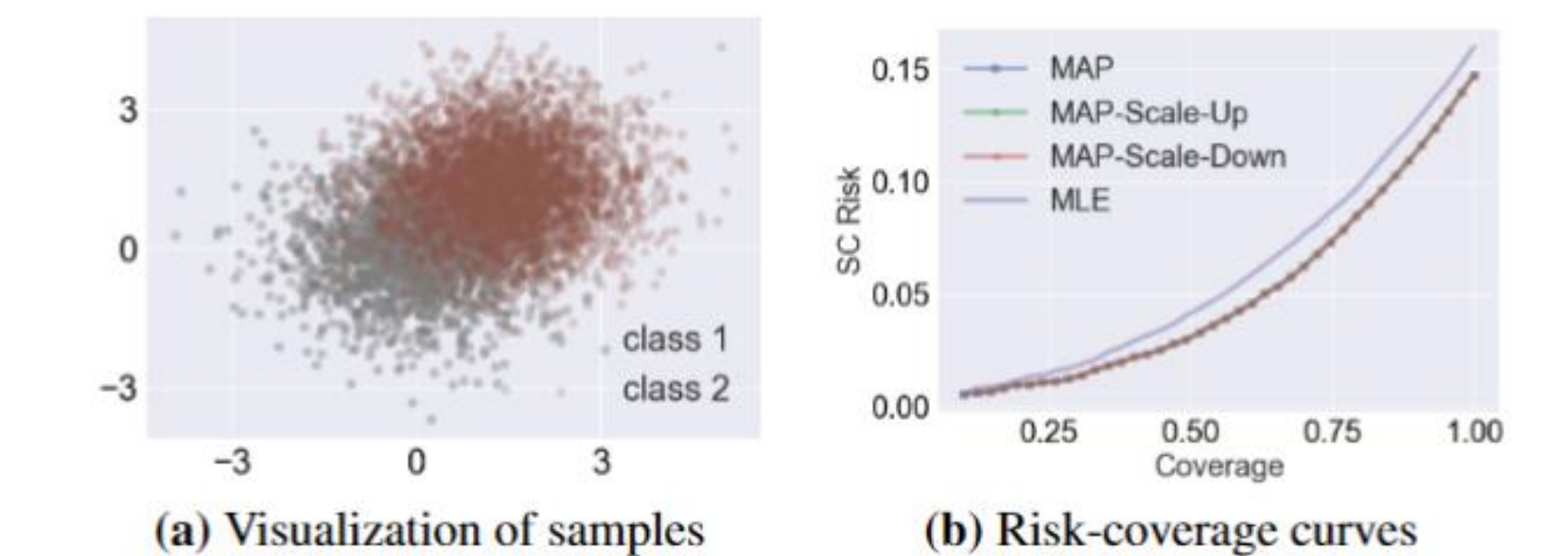
## Selective Classification (SC)

- Selectively making predictions to avoid excessive errors.
- Beneficial to deploy the imperfect AI models to practical applications with high-stakes requirements

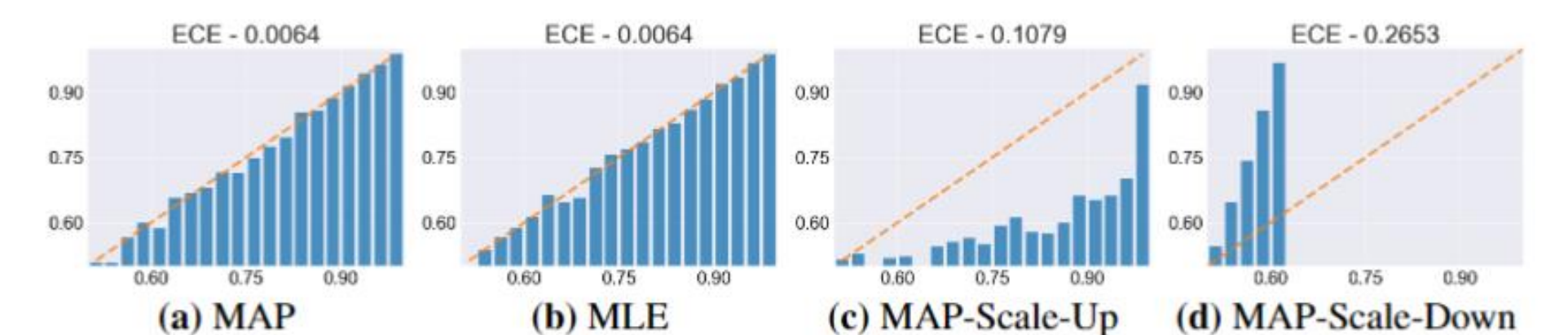
E.g., an AI-powered medical diagnosis assistant can make confident and correct predictions on its own, saving a significant amount of doctors' labor, while turning unconfident cases to doctor.

## Calibrated confidence $\neq$ selection confidence

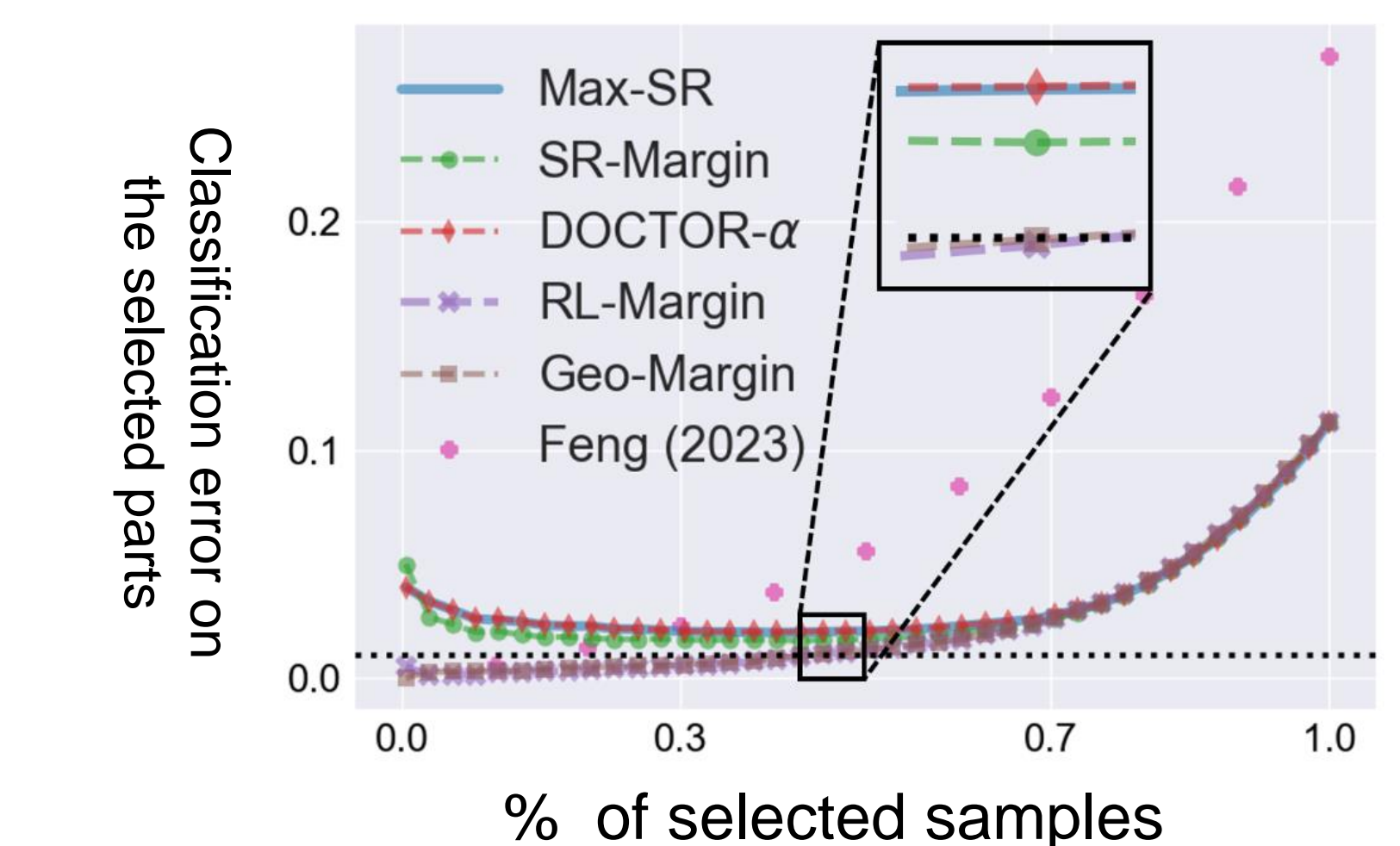
- A simple example and SC performance of 4 score functions



- Their empirical calibration assessment



## "Margin" is better than maximum softmax scores



[N] Liang H, Peng L, Sun J. Toward Effective Post-Training Selective Classification for High-Stakes Applications. In preparation for Neurips 2023.