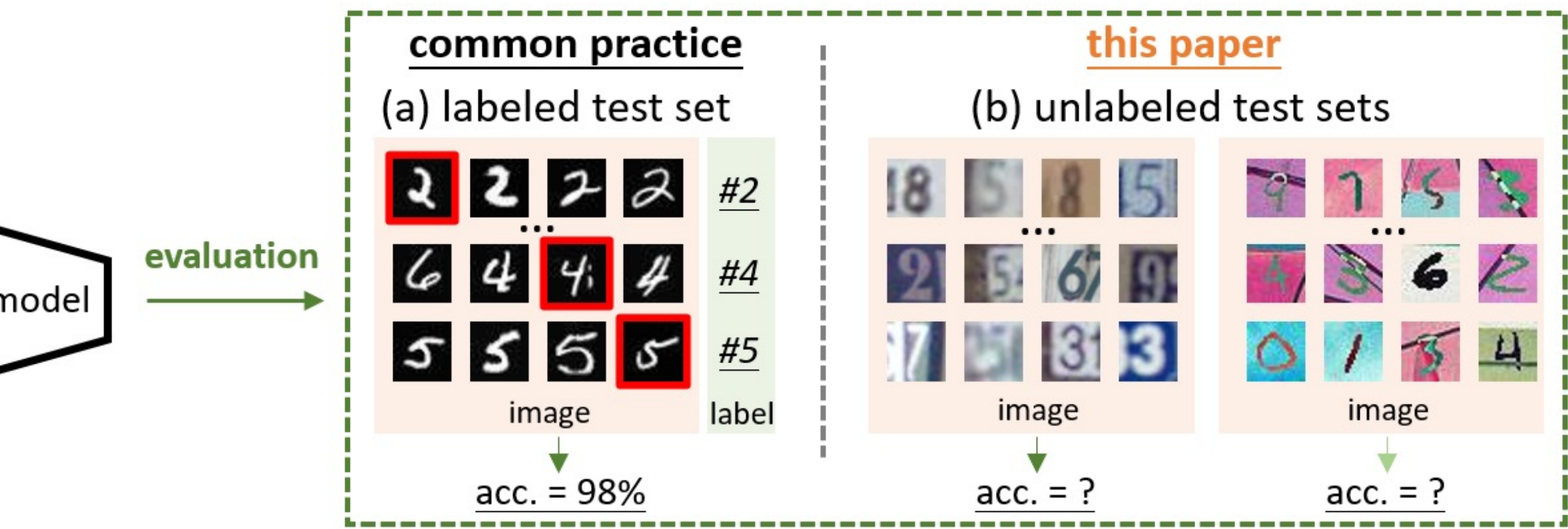


## Automatic Model Evaluation

Given a trained classifier, the overall goal is to estimate its accuracy on various test datasets **without labels**



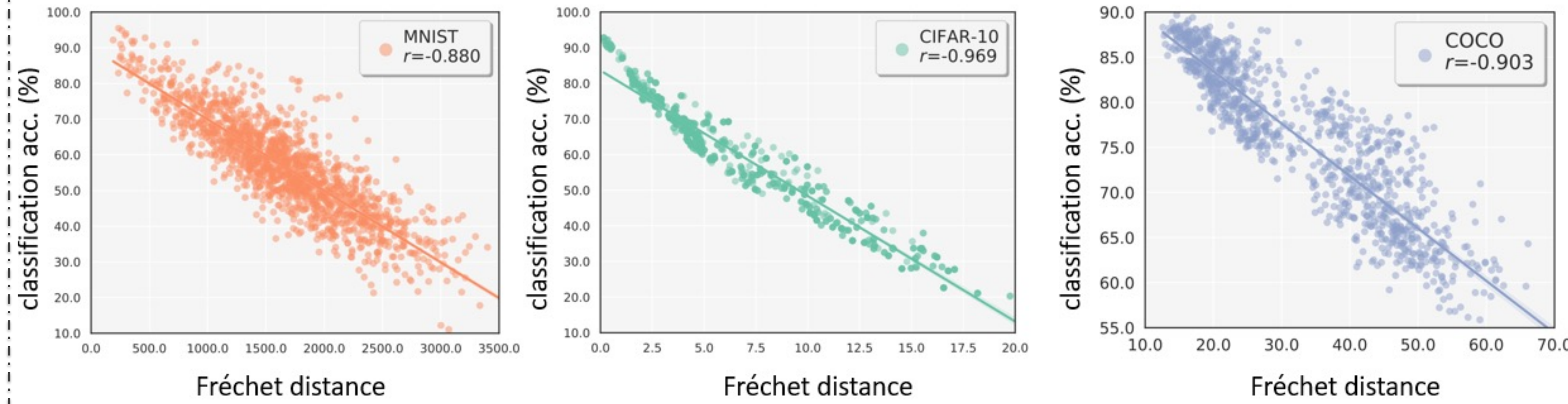
## Motivation

- Distribution shift degrades classifier accuracy  
Reduce the shift to make model generalizable to the target domain
- Meta-dataset (dataset of datasets)  
It contains many datasets from different distributions  
**Data synthesis:** 1) image transform and 2) background change



## Correlation Study

- Negative Linear Correlation between Test Accuracy and Distribution Shift



Classifier tends to have a **low accuracy** on the sample set which has a **high distribution shift** from the training set

## Dataset-level Regression

Predict classifier accuracy from **distribution-related statistics** on an unlabeled test set

- Linear regression

$$\text{Model: } a_{\text{linear}} = A_{\text{linear}}(\mathbf{f}) = w_1 f_{\text{linear}} + w_0$$

$$\text{Feature: } f_{\text{linear}} = \text{FD}(\mathcal{D}_{\text{ori}}, \mathcal{D}) = \|\mu_{\text{ori}} - \mu\|_2^2 + \text{Tr}(\Sigma_{\text{ori}} + \Sigma - 2(\Sigma_{\text{ori}}\Sigma))^{\frac{1}{2}}$$

**Fréchet distance (FD)** measures the difference between training and test distributions

- Neural network regression

$$\text{Model: } a_{\text{neural}} = A_{\text{neural}}(\mathbf{f}_{\text{neural}})$$

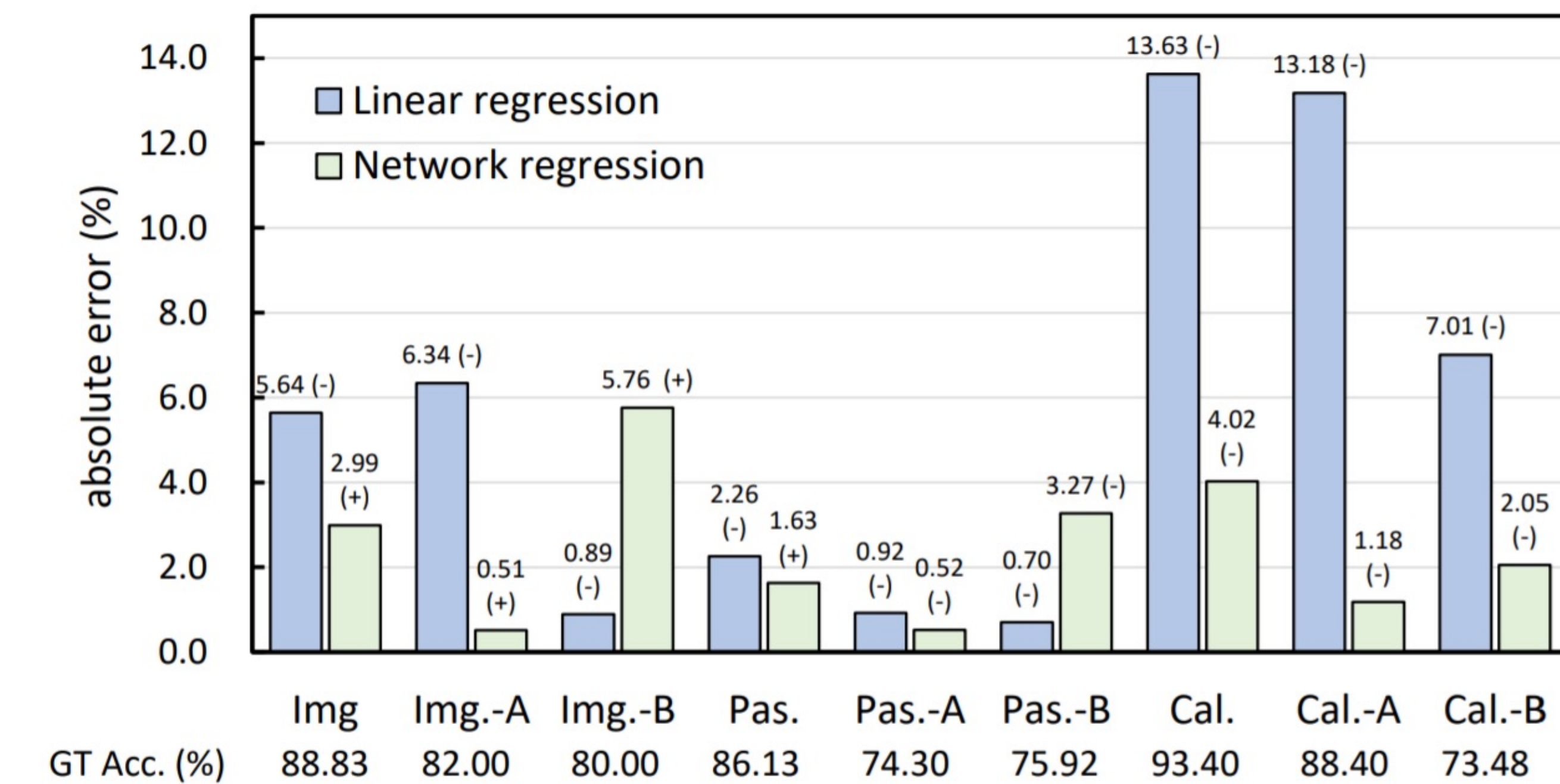
$$\text{Feature: } \mathbf{f}_{\text{neural}} = [f_{\text{linear}}; \mu; \sigma]$$

Network regression uses **mean, co-variance, and FD** to represent each test set

## Experiment

Method	Digits			Natural images			
	SVHN	USPS	RMSE↓	Pascal	Caltech	ImageNet	RMSE↓
Ground-truth accuracy	25.46	64.08	-	86.13	93.40	88.83	-
Predicted score ( $\tau = 0.8$ )	7.97	37.22	22.66	84.32	90.78	86.50	<b>2.28</b>
Predicted score ( $\tau = 0.9$ )	7.03	32.94	25.59	78.61	87.71	81.33	6.96
Linear reg.	26.28	50.14	9.87	83.87	79.77	83.19	8.62
Neural network reg.	27.52	64.11	<b>1.46</b>	87.76	89.39	91.82	<b>3.04</b>

Predicted score-based baseline is sensitive to the threshold;  
**Our regression methods gain promising estimations**



Comparing linear regression and neural network regression when test data undergo new image transformations: A) Cutout and Shear; B) Equalize and ColorTemperature

**Network regression is more robust than linear regression**

## Reference

- A. Torralba and A. A. Efros. Unbiased look at dataset bias. In CVPR, 2011  
Ben-David, Shai, et al. "A theory of learning from different domains." Machine learning, 2010

The code is available at  
<https://weijiandeng.xyz/AutoEval>

