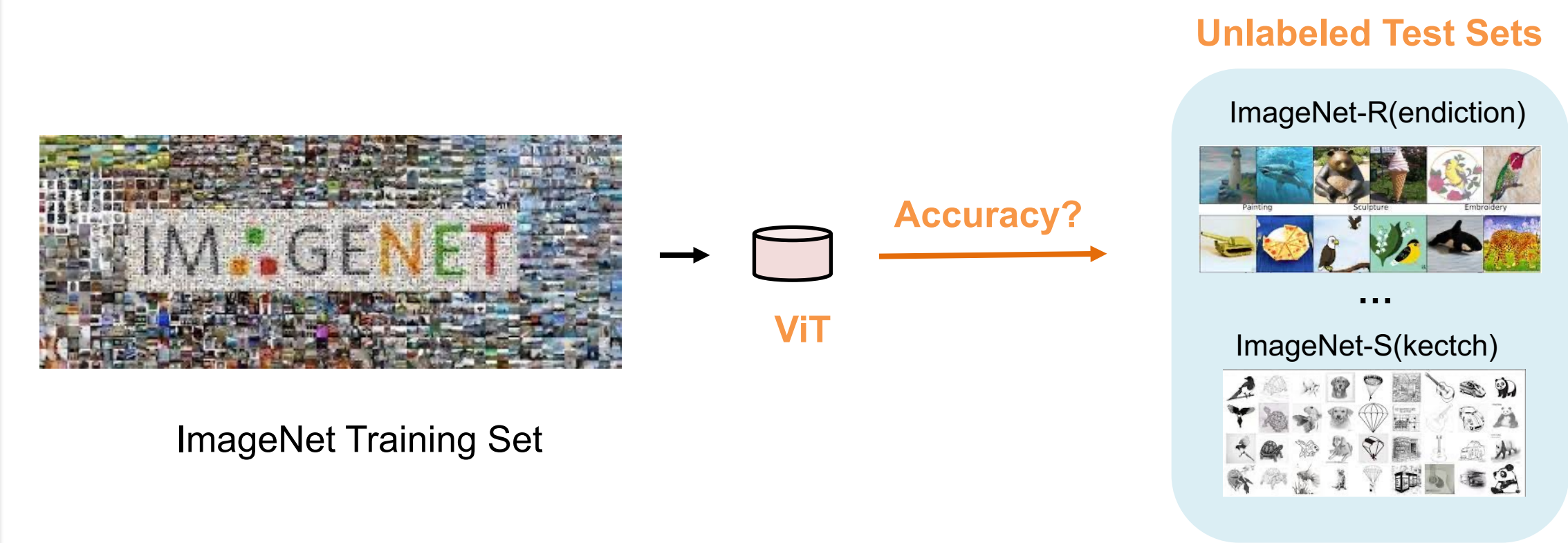# Confidence and Dispersity Speak: Characterizing Prediction Matrix for Unsupervised Accuracy Estimation

Weijian Deng[1]  Yumin Suh[2]  Stephen Gould[1]  Liang Zheng[1]
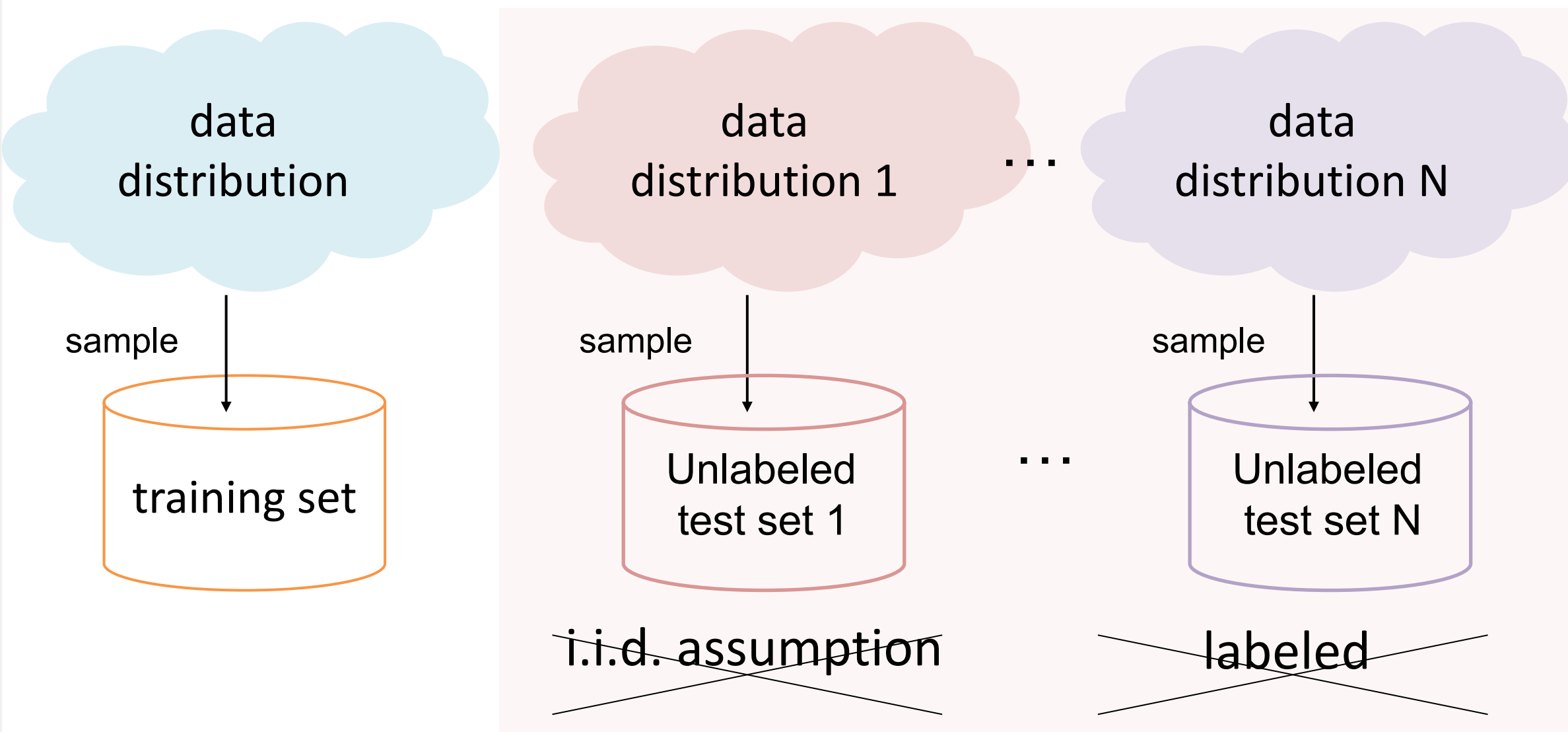
[1]Australian National University   [2]NEC Labs America

Australian National University

NEC Laboratories America

ICML International Conference On Machine Learning — 40 Years

## Unsupervised Accuracy Estimation

- **Definition**: given a trained model, the goal is to estimate its accuracy on various test datasets **without labels**



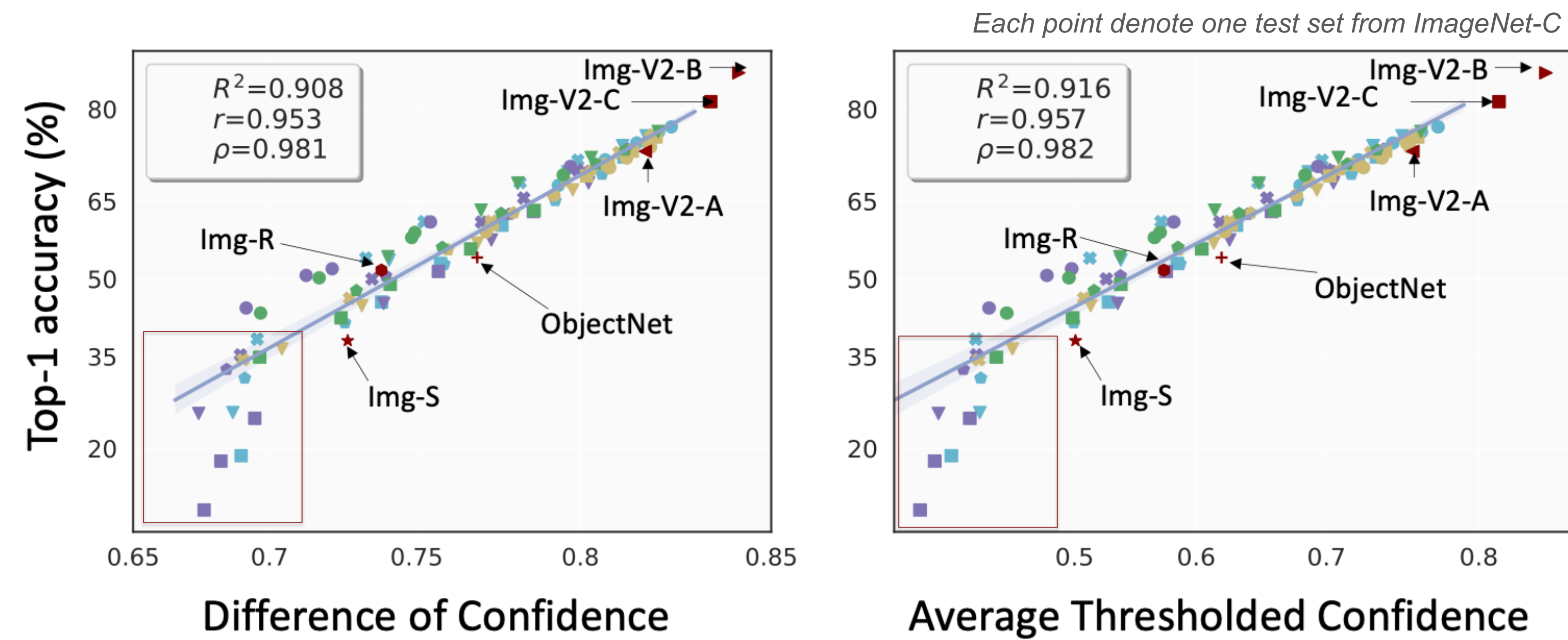ImageNet Training Set → ViT → Accuracy? → Unlabeled Test Sets: ImageNet-R(endiction) … ImageNet-S(ketch)

**Real-world evaluation**: 1) the distributions of test sets are often *different* from that of training set (*no i.i.d*); 2) test labels are *unavailable* or *expensive to obtain*.



data distribution → sample → training set

data distribution 1 → sample → Unlabeled test set 1 … data distribution N → sample → Unlabeled test set N

i.i.d. assumption     labeled

In-distribution accuracy may only be a weak predictor of performance on out-of-distribution data;
**Evaluation without labels and under distribution shifts**

## Prediction Confidence

- **Confidence** reflects whether the individual prediction is certain
  Existing methods (*e.g.*, DoC and ATC) explore such information

*Each point denote one test set from ImageNet-C*



$R^2=0.908$
$r=0.953$
$\rho=0.981$

$R^2=0.916$
$r=0.957$
$\rho=0.982$

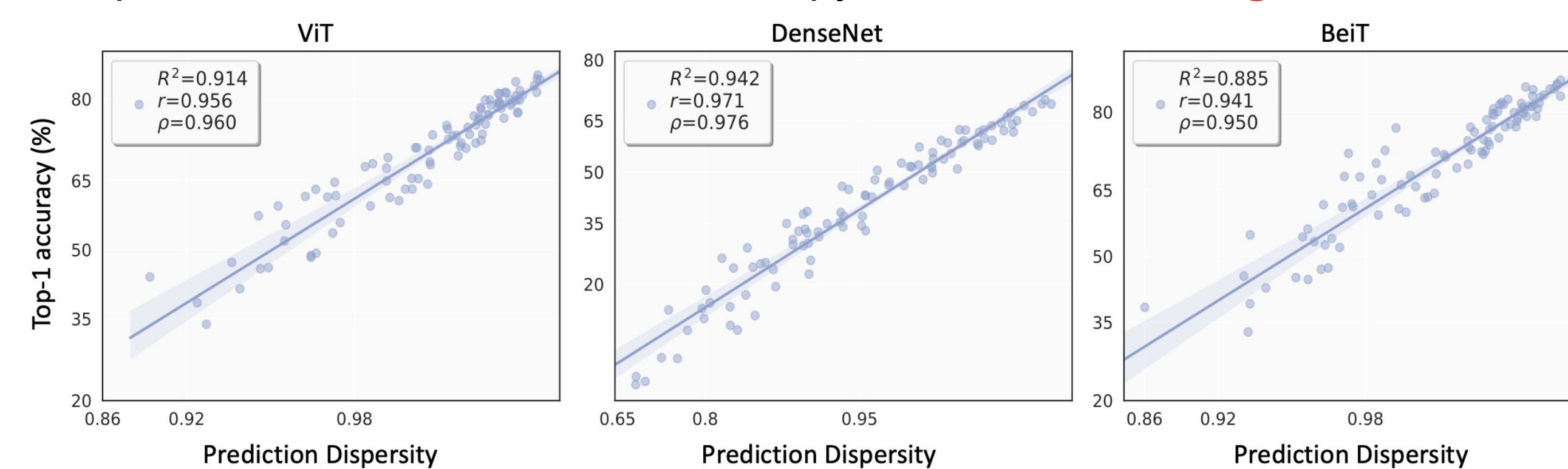Difference of Confidence    Average Thresholded Confidence

- **Confidence may be a weak indicator**
  Predication score-based methods **cannot well capture** the test sets in the **low-accuracy region** (bottom-left area of the above correlation figure)

## Prediction Dispersity

- **Dispersity** indicates how the predictions are distributed across all categories
  **Prediction Dispersity Score**: we first calculate the histogram of the number of the predicted class and then use entropy to measure **the degree of balance**
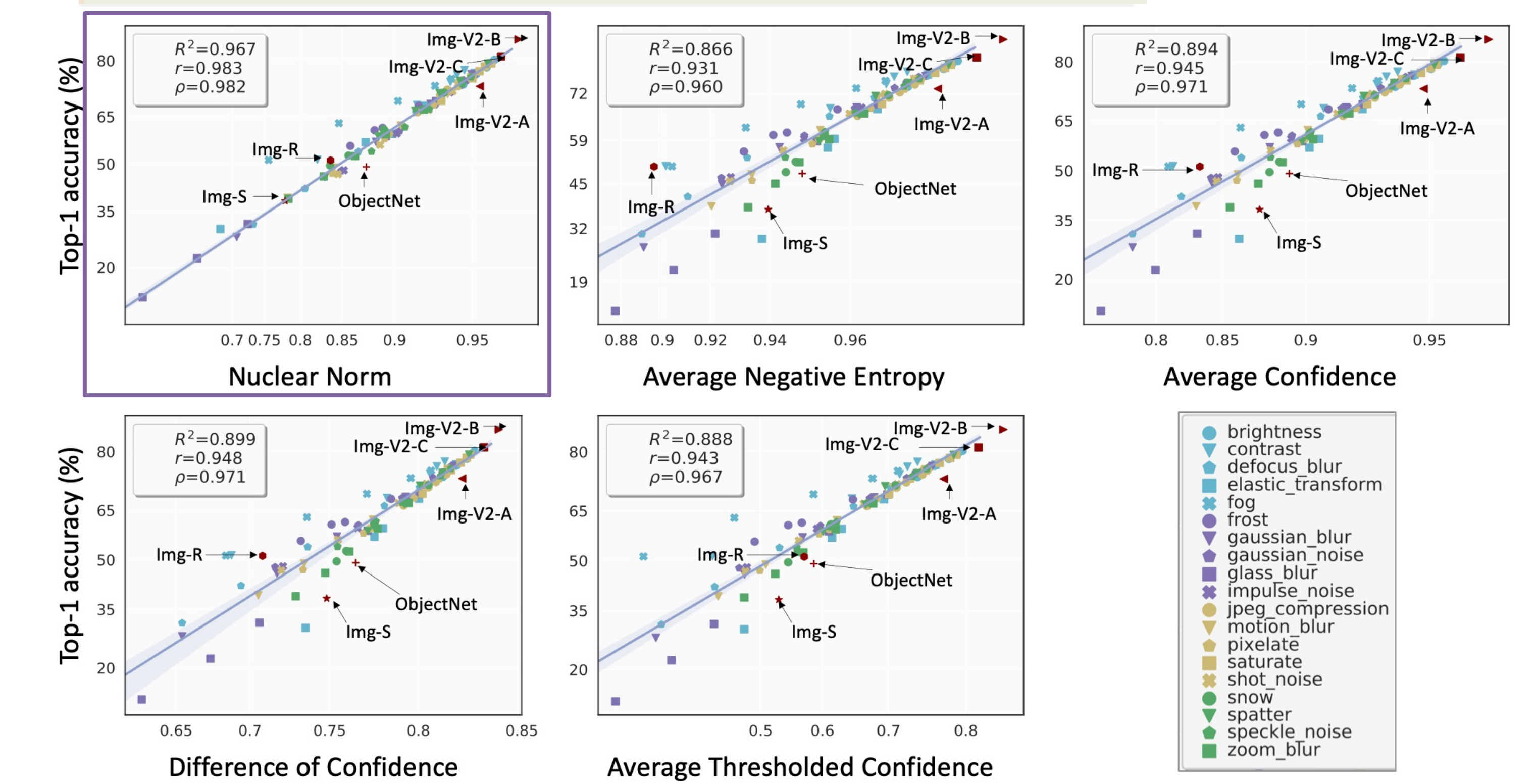


ViT: $R^2=0.914$, $r=0.956$, $\rho=0.960$
DenseNet: $R^2=0.942$, $r=0.971$, $\rho=0.976$
BeiT: $R^2=0.885$, $r=0.941$, $\rho=0.950$

## Nuclear Norm

- **Nuclear norm** is effective in characterizing both confidence and dispersity

Prediction Matrix $P \in \mathbb{R}^{N_t \times K}$ ($N_t$ test samples, and $K$ classes)

Nuclear Norm: the sum of singular values of prediction matrix



Nuclear Norm: $R^2=0.967$, $r=0.983$, $\rho=0.982$
Average Negative Entropy: $R^2=0.866$, $r=0.931$, $\rho=0.960$
Average Confidence: $R^2=0.894$, $r=0.945$, $\rho=0.971$
Difference of Confidence: $R^2=0.899$, $r=0.948$, $\rho=0.971$
Average Thresholded Confidence: $R^2=0.888$, $r=0.943$, $\rho=0.967$

Legend: brightness, contrast, defocus_blur, elastic_transform, fog, frost, gaussian_blur, gaussian_noise, glass_blur, impulse_noise, jpeg_compression, motion_blur, pixelate, saturate, shot_noise, snow, spatter, speckle_noise, zoom_blur

Nuclear norm exhibits the highest correlation strength with OOD accuracy

## Potential Direction

1) Other methods are stable under class imbalance;
2) Nuclear Norm **is resistant to moderate class imbalance;**
3) Nuclear Norm **is less effective under severe class imbalance.**

If we have **prior knowledge** about the imbalanced class distribution, we can expect class predictions to follow it rather than a uniform one



$m=0.1$, $m=0.2$, $m=0.4$, $m=0.6$, $m=0.8$, $m=1.0$

Nuclear Norm: $R^2=0.881$, $\rho=0.918$
Average Negative Entropy: $R^2=0.950$, $\rho=0.980$
Average Confidence: $R^2=0.960$, $\rho=0.982$