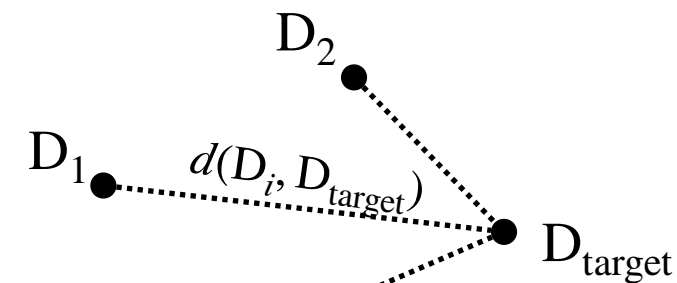


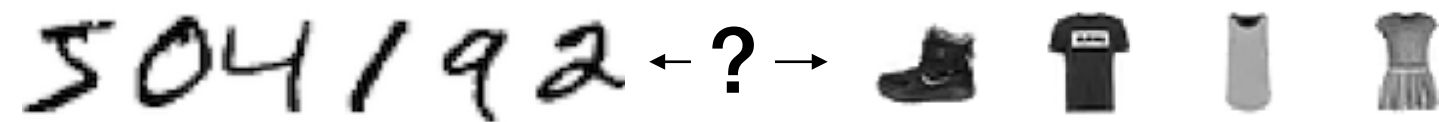
## Distances between Datasets

- Key in various settings: transfer/meta-learning, etc



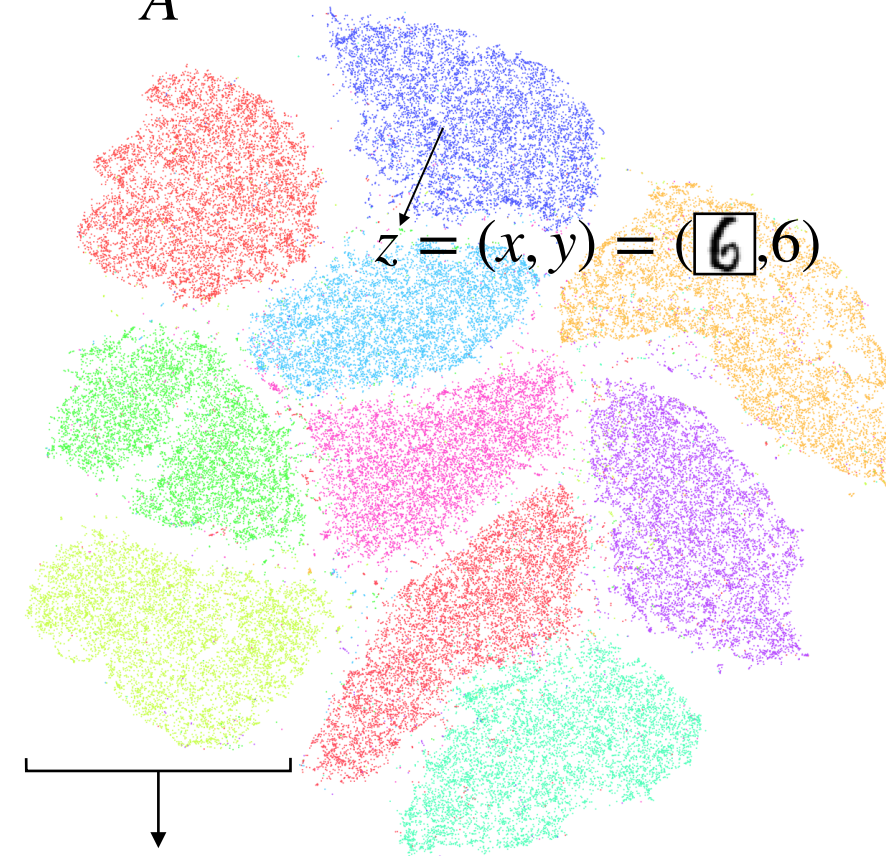
$D_3$  → not best choice for pre-training!

- How to deal with datasets with different label sets?



- Ideally: model agnostic, sound theoretical footing

$D_A$



Labels represented as distribution over features  $v_y = \mathcal{N}(\mu_y, \Sigma_y)$

## The Optimal Transport Dataset Distance

Distance between **feature/label pairs**:

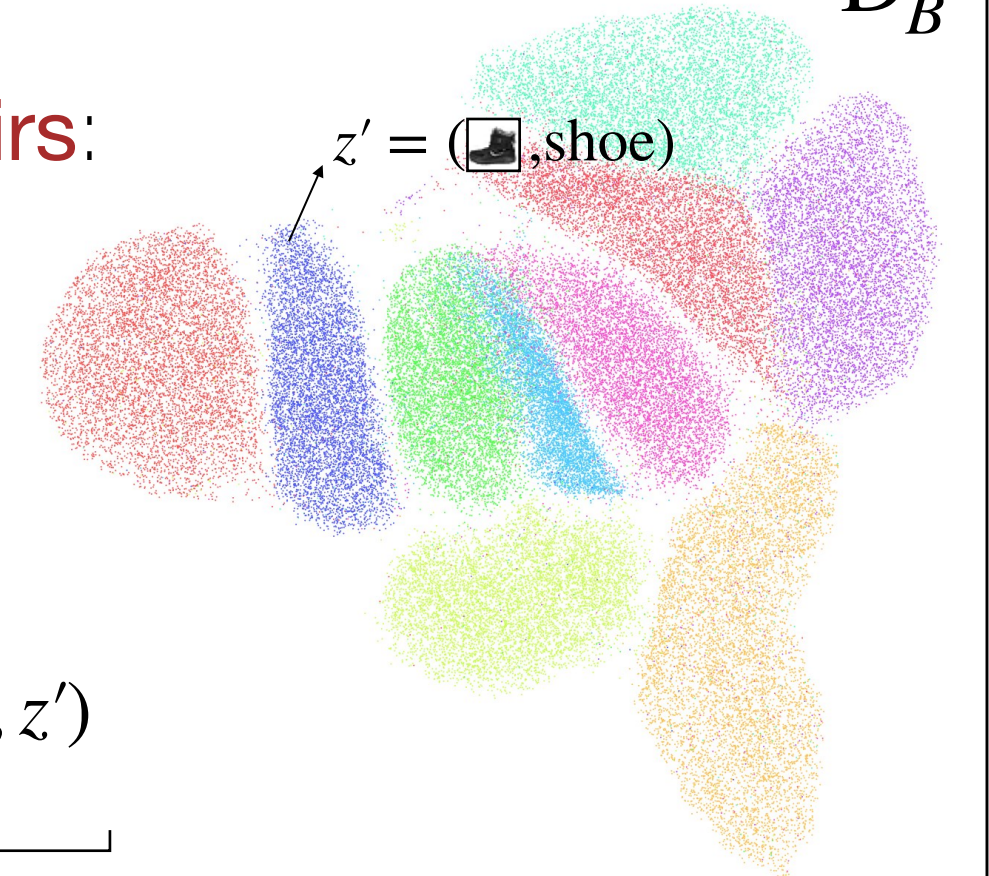
$$d(z, z') = \left( d(x, x')^p + \underbrace{W_p^p(v_y, v_{y'})}_{\text{Wasserstein distance}} \right)^{1/p}$$

Distance between **datasets**:

$$d_{OT}(D_A, D_B) = \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} d(z, z')^p d\pi(z, z')$$

Optimal Transport distance:  $\approx$  min-cost matching

$D_B$

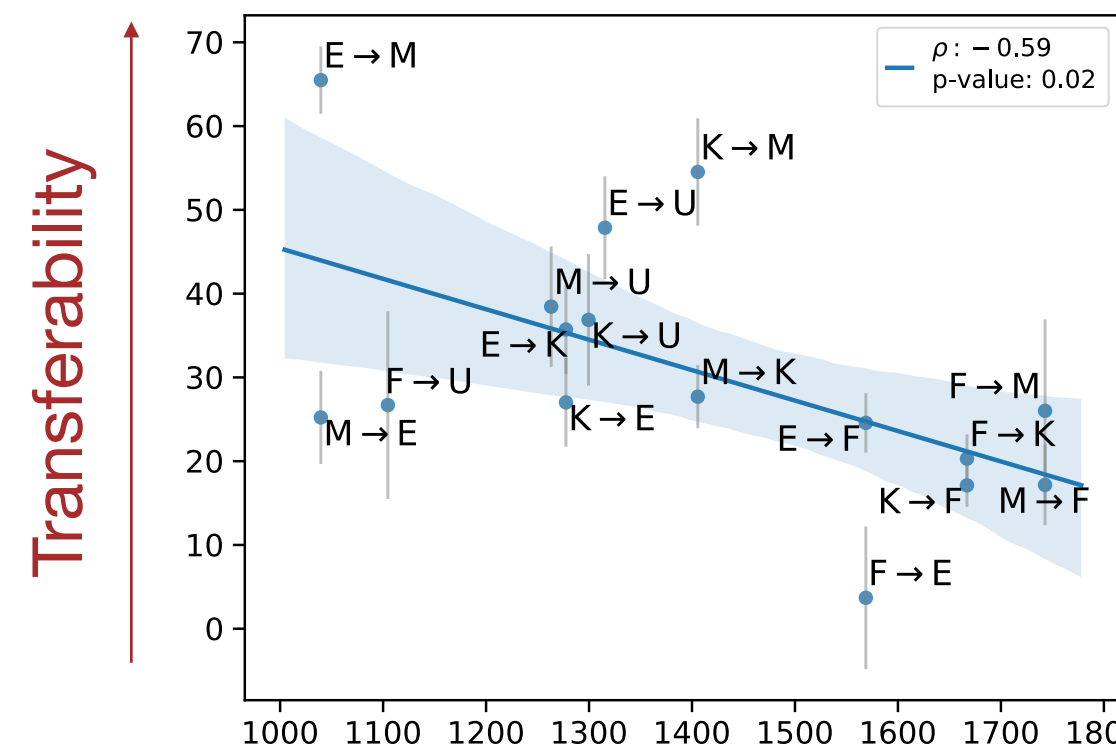


## Properties of OTDD

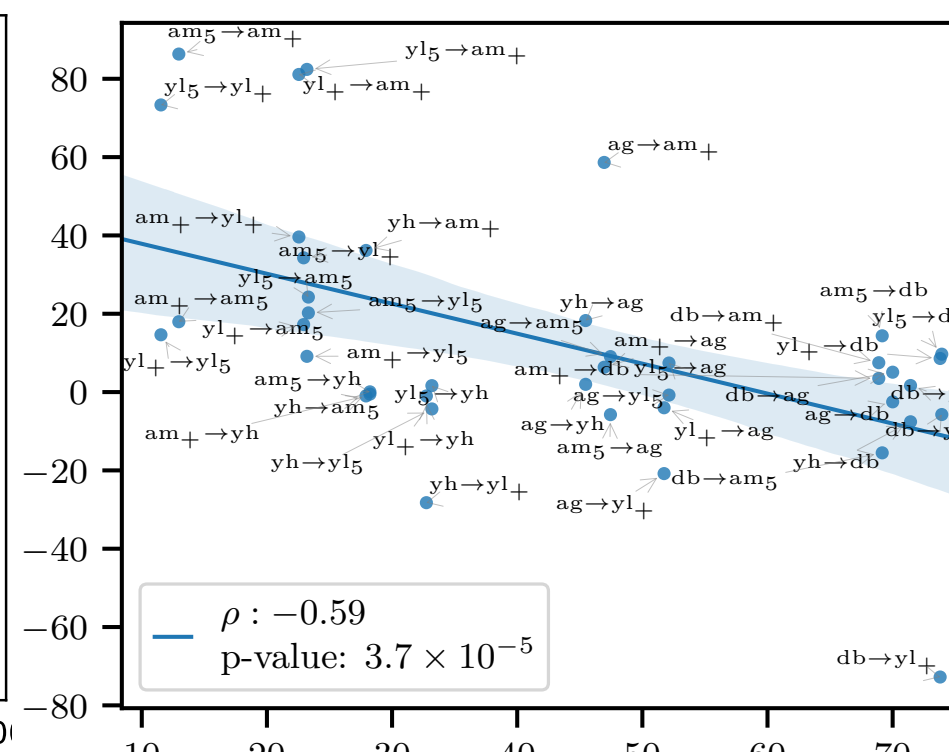
- A **true metric** in the space of measures over feature-label pairs  $\mathcal{P}(\mathcal{X} \times \mathcal{Y})$
- Can be estimated from finite samples
- **Efficient computation** through a Gaussian approximation + online moment estimation

## Application: Predicting Transferability

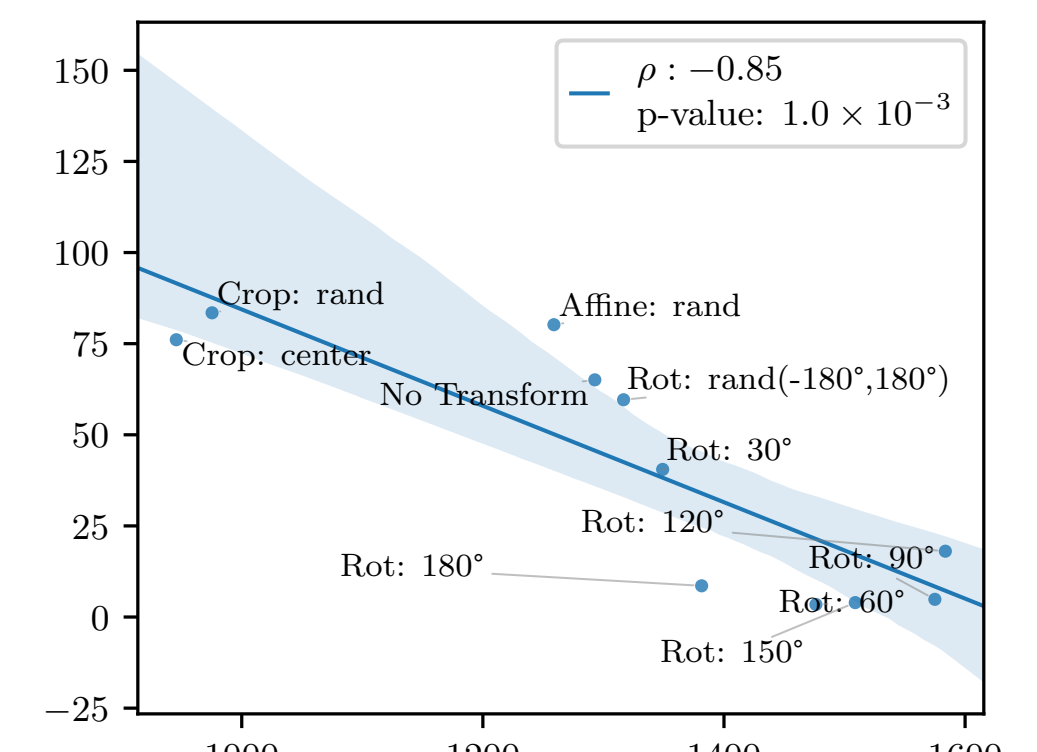
\*NIST Datasets



Text Datasets



MNIST + Augmentations



OT Dataset Distance