

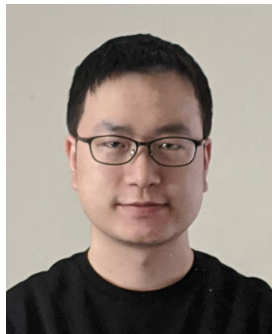


AIFARMS

Artificial Intelligence for Future Agricultural
Resilience, Management, and Sustainability



Distribution-Informed Neural Networks for Domain Adaptation Regression



Jun Wu



Jingrui He



Sheng Wang



Kaiyu Guan



Elizabeth Ainsworth

University of Illinois Urbana-Champaign
{junwu3, jingrui, sheng12, kaiyug, ainswort}@illinois.edu

Domain Adaptation Regression

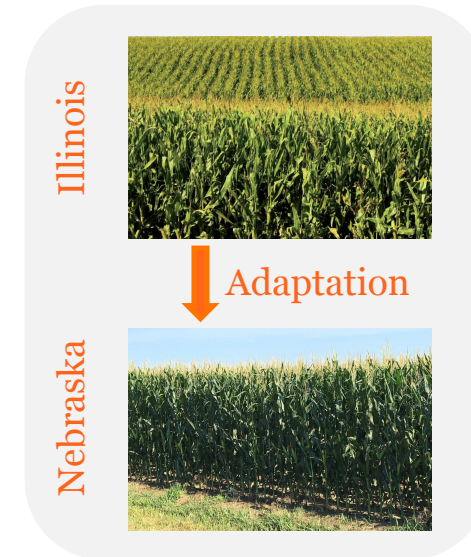
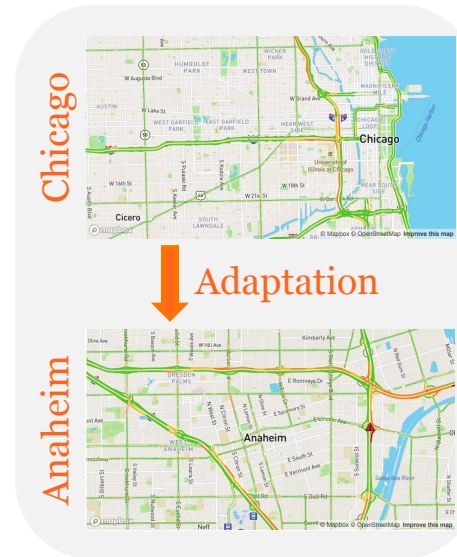
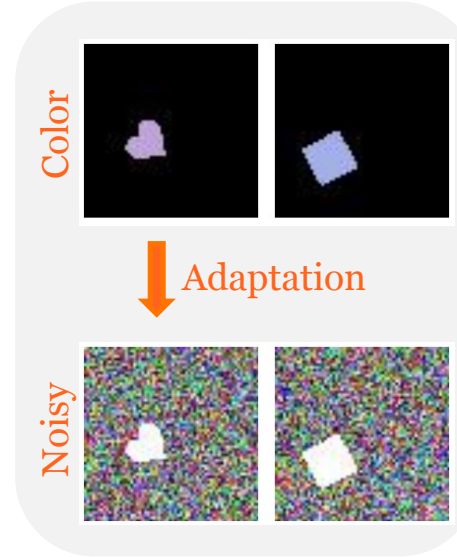


□ Problem definition

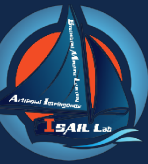
- Input: A source domain and a target domain
- Output: Prediction function on the target domain

□ Applications

- **Computer vision:** Object localization
- **Natural language processing:** Sentiment analysis
- **Graph mining:** Traffic flow prediction
- **Agriculture analysis:** Plant phenotyping



Distribution-Informed Neural Networks



□ Motivation

- For standard neural network $f(\cdot)$ on a single domain

$$f(x_2) \approx f(x_3) \text{ if } x_2 \approx x_3 \text{ (only source domain)}$$

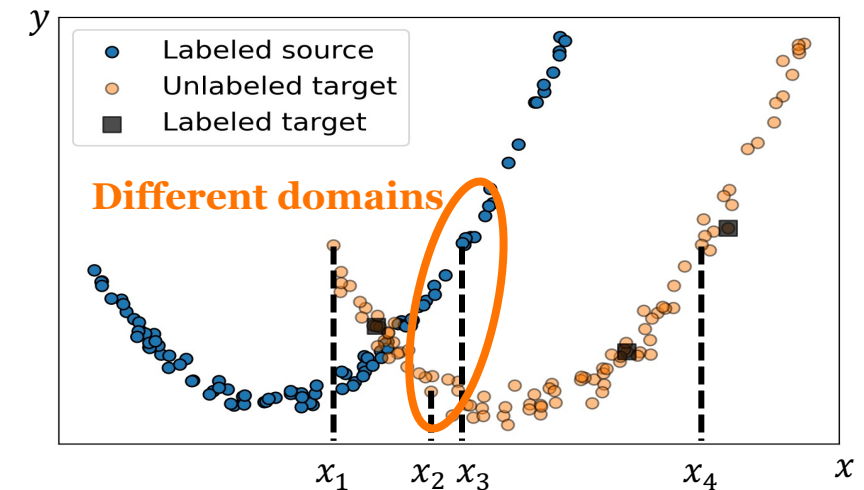
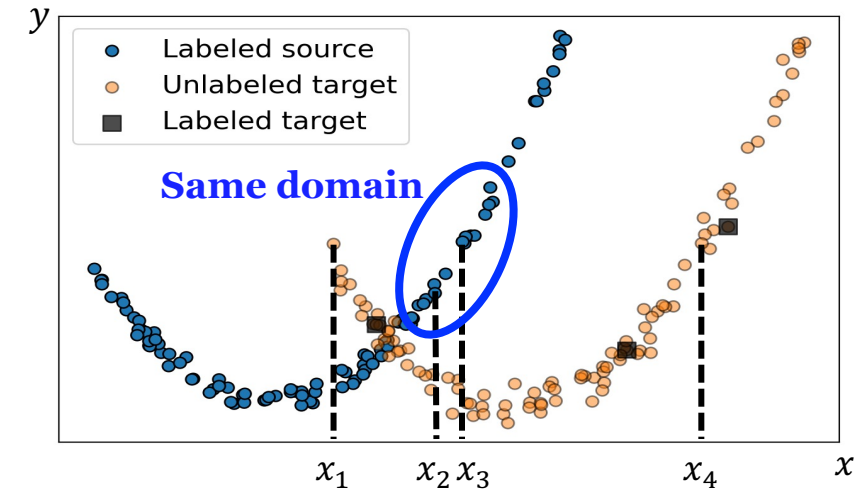
- Limitations on **heterogeneous domains**

$$f(x_2) \neq f(x_3) \text{ for } x_2 \approx x_3 \text{ (heterogeneous case)}$$

- x_2 from target domain, and x_3 from source domain

- Implication:

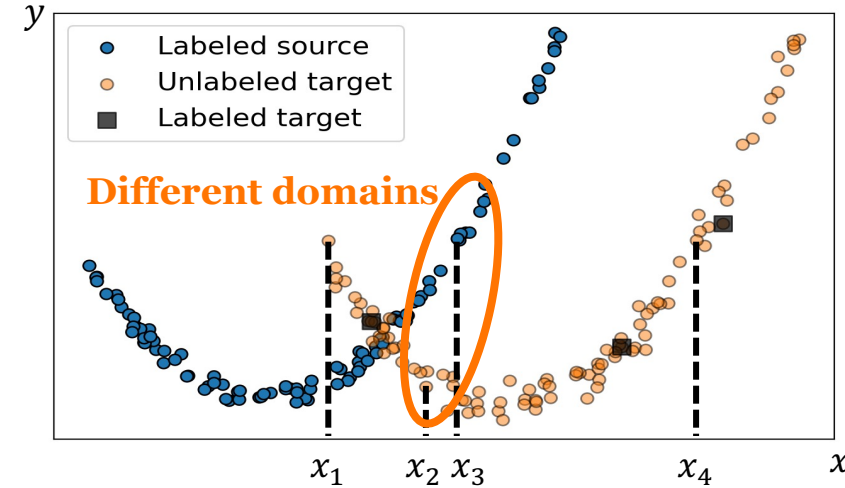
- Input-output relationship varies among different domains
- It cannot be directly captured by a standard neural network



Distribution-Informed Neural Networks

Formal definition

- Distribution-aware input-output relationship
 - $\tilde{f}: (x, \mathbb{P}) \rightarrow y$ where $x \in \mathbb{P}$
 - $\tilde{f}(x_2, \mathbb{P}_{tgt}) \neq \tilde{f}(x_3, \mathbb{P}_{src})$ for $x_2 \approx x_3$ (heterogeneous case)
- Distribution-informed neural network (DINO)



DINO:

$$\begin{aligned} \tilde{f}(x, \mathbb{P}) &:= f_{\theta}(x) \cdot g_{w_g}(\mathbb{P}|x) \\ &= (\phi_{\theta^{<L}}(x)^T w) \cdot (\Phi_x(\mathbb{P})^T w_g) = w^T (\phi_{\theta^{<L}}(x) \Phi_x(\mathbb{P})^T) w_g \end{aligned}$$

Input representation learning

A fully-connected NN: $f_{\theta}(x) = \phi_{\theta^{<L}}(x)^T w$

$\theta^{<L}$: Parameters of the first $L - 1$ layers
 w : Parameters of the output layer

Input-oriented distribution representation learning

Infinitely-wide $f_{\theta}(\cdot)$ \rightarrow NNGP kernel space $K_x \rightarrow \Phi_x(\mathbb{P}) = \sum_{i=1}^n \beta_{x, \tilde{x}_i} \langle \cdot, \tilde{x}_i \rangle_{K_x} \rightarrow g_{w_g}(\mathbb{P}|x) = \Phi_x(\mathbb{P})^T w_g$



□ Observation at model initialization

- DINO is a Gaussian process with adaptive NNGP kernel

Under random initialization, when the network width goes to infinity, we have $\tilde{f}(\cdot) \sim \mathcal{N}(0, K^{DA})$ with

$$K^{DA}((x, \mathbb{P}), (x', \mathbb{P}')) = K_{\mathcal{X}}(x, x') \cdot K_{\mathcal{P}|\mathcal{X}}(\mathbb{P}, \mathbb{P}' | x, x')$$

where $K_{\mathcal{X}}(\cdot, \cdot)$ is the NNGP kernel, and $K_{\mathcal{P}|\mathcal{X}}(\cdot, \cdot)$ is a distribution kernel, i.e.,

$$K_{\mathcal{P}|\mathcal{X}}(\mathbb{P}, \mathbb{P}' | x, x') = \sum_{i=1}^n \sum_{j=1}^{n'} \beta_{x, \tilde{x}_i} \beta_{x', \tilde{x}_j} K_{\mathcal{X}}(\tilde{x}_i, \tilde{x}_j)$$

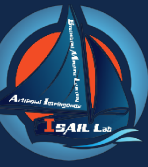
□ Adaptive Gaussian process algorithm

- Prior GP $\tilde{f}(\cdot) \sim \mathcal{N}(0, K^{DA})$
- Prediction function $p(Y | X_*^{tgt}) = \mathcal{N}(\bar{\mu}, \bar{\Sigma})$

$$\bar{\mu} = K^{DA}(X_*^{tgt}, X)C^{-1}Y \quad \bar{\Sigma} = K^{DA}(X_*^{tgt}, X_*^{tgt}) - K^{DA}(X_*^{tgt}, X)C^{-1}K^{DA}(X_*^{tgt}, X)^T$$



Proposed Algorithm: DINO-TRAIN



□ DINO under gradient descent training

- Objective function

$$\mathcal{L}(\theta) = \underbrace{\frac{\alpha}{2n_{src}} \sum_{i=1}^{n_{src}} (\tilde{f}(x_i^{src}, \mathbb{P}^{src}) - y_i^{src})^2 + \frac{1-\alpha}{2n_{tgt}^l} \sum_{j=1}^{n_{tgt}^l} (\tilde{f}(x_j^{tgt}, \mathbb{P}^{tgt}) - y_j^{tgt})^2}_{\text{Supervised loss over labeled examples}} + \underbrace{\frac{\mu}{2} \text{MMD}_{\Theta_{DA}}^2(\mathbb{P}^{src}, \mathbb{P}^{tgt})}_{\text{Empirical MMD-NTK}}$$

Supervised loss over labeled examples

Empirical MMD-NTK

- Empirical Maximum Mean Discrepancy (MMD) over training dynamics
 - Measure the distribution shift during gradient descent training

$$\text{MMD}_{\Theta_{DA}}^2(\mathbb{P}^{src}, \mathbb{P}^{tgt}) = \left\| \frac{1}{n_{src}} \sum_{i=1}^{n_{src}} \nabla_{\theta} \tilde{f}(x_i^{src}, \mathbb{P}^{src}) - \frac{1}{n_{tgt}} \sum_{j=1}^{n_{tgt}} \nabla_{\theta} \tilde{f}(x_j^{tgt}, \mathbb{P}^{tgt}) \right\|_{\mathcal{H}_{DA}}^2$$

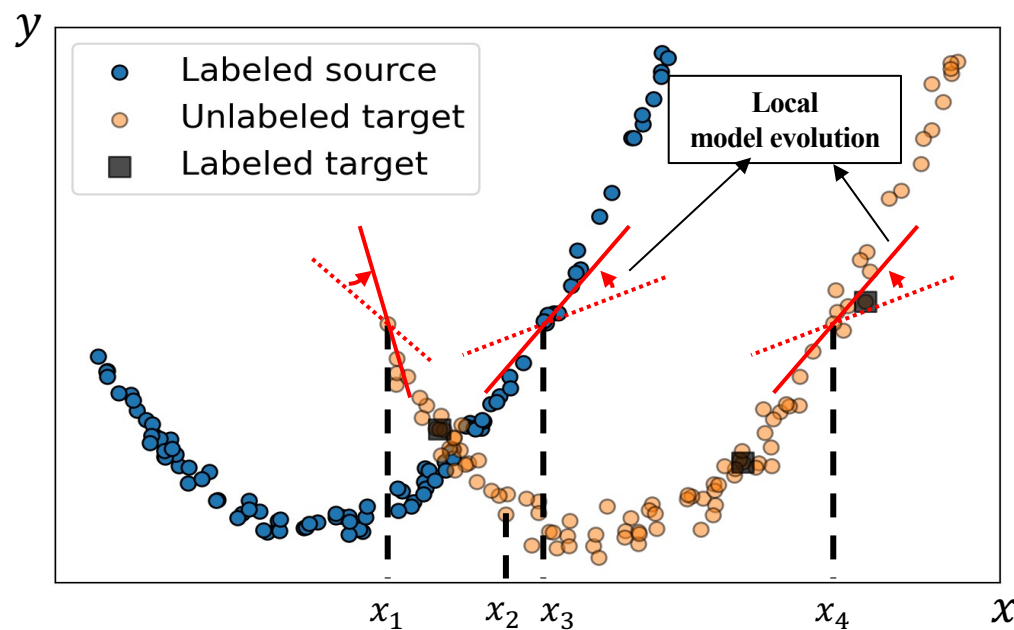
Proposed Algorithm: DINO-TRAIN



□ Intuition behind MMD-NTK

- Unified distribution discrepancy estimator
- Local model evolvment

	Estimator	Characterization
MMD-RBF [Long et al. 2015]	<ul style="list-style-type: none">▪ Two-stage discrepancy;▪ Heuristic division: feature extractor and label predictor	<ul style="list-style-type: none">▪ Simple model output
MMD-NTK	<ul style="list-style-type: none">▪ Unified estimator over training dynamics of $\tilde{f}(\cdot)$	<ul style="list-style-type: none">▪ Local model evolvment



□ Convergence

There exists $\eta^* \in \mathbb{R}_+$ such that for the the infinitely-wide DINO $\tilde{f}(\cdot)$ trained under gradient flow with learning rate $\eta < \eta^*$, the prediction function $\lim_{t \rightarrow \infty} \tilde{f}_{\theta_t}(X_*^{tgt})$ **converges to a Gaussian process** with

$$\begin{aligned} \mu &= \Theta_{DA}(X_*^{tgt}, X) \Theta_{DA}(X, X)^{-1} Y \\ \Sigma &= K^{DA}(X_*^{tgt}, X_*^{tgt}) + \Theta_{DA}(X_*^{tgt}, X) \Theta_{DA}(X, X)^{-1} K^{DA} \Theta_{DA}(X, X_*^{tgt}) - (\Theta_{DA}(X_*^{tgt}, X) \Theta_{DA}(X, X)^{-1} K^{DA} (X, X_*^{tgt}) + h.c.) \end{aligned}$$

□ Generalization

For any $\delta > 0$, with probability $1 - \delta$, the expected error in the target domain can be bounded by

$$\epsilon_{tgt}(\tilde{f}) \leq \underbrace{\frac{\alpha}{n_{src}} \sum_{i=1}^{n_{src}} (\tilde{f}(x_i^{src}, \mathbb{P}^{src}) - y_i^{src})^2 + \frac{1 - \alpha}{n_{tgt}^l} \sum_{j=1}^{n_{tgt}^l} (\tilde{f}(x_j^{tgt}, \mathbb{P}^{tgt}) - y_j^{tgt})^2}_{\text{DINO-TRAIN empirically minimizes the error bound}} + 8\alpha M_0 \cdot \text{MMD}_{\Theta_{DA}}(\mathbb{P}^{src}, \mathbb{P}^{tgt}) + \Omega$$

DINO-TRAIN empirically minimizes the error bound

□ Data sets

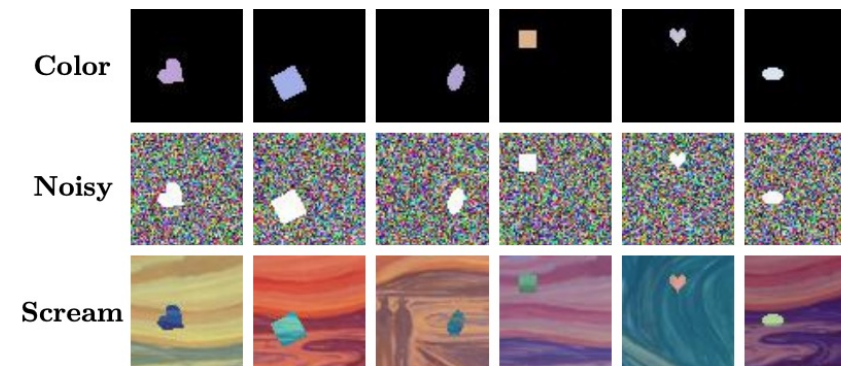
- dSprites
- MPI3D
- Plant Phenotyping

□ Metric

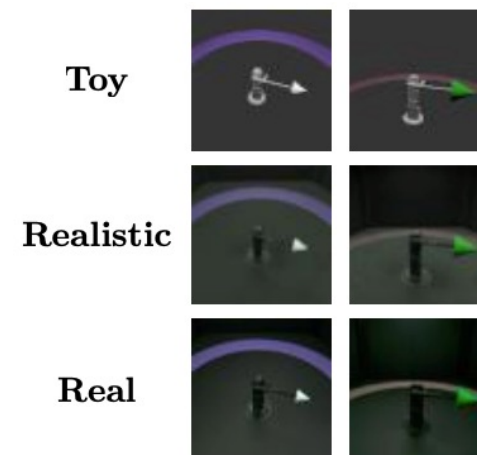
- MAE: Mean Absolute Error

□ Baseline

- Plain Gaussian process: NNGP and NTKGP
- Adaptive Gaussian process: AT-GP and TL-NTK



Examples of dSprites



Examples of MPI3D

Results



Methods	C → N	C → S	N → C	N → S	S → C	S → N	Avg.
NNGP [34]	2.041 \pm 0.001	1.823 \pm 0.001	0.445 \pm 0.002	0.624 \pm 0.001	0.197 \pm 0.002	0.459 \pm 0.002	0.932
NTKGP [25]	1.345 \pm 0.002	1.227 \pm 0.000	0.323 \pm 0.002	0.529 \pm 0.004	0.248 \pm 0.001	0.425 \pm 0.002	0.683
AT-GP [7]	0.194 \pm 0.005	0.259 \pm 0.002	0.104 \pm 0.001	0.252 \pm 0.005	0.118 \pm 0.003	0.189 \pm 0.006	0.186
TL-NTK [38]	0.164 \pm 0.001	0.231 \pm 0.000	0.124 \pm 0.005	0.242 \pm 0.002	0.125 \pm 0.001	0.197 \pm 0.004	0.181
DINO-INIT (ours)	0.128 \pm 0.001	0.233 \pm 0.003	0.114 \pm 0.002	0.227 \pm 0.002	0.112 \pm 0.001	0.181 \pm 0.005	0.166
DINO-TRAIN (ours)	0.127 \pm 0.002	0.240 \pm 0.003	0.127 \pm 0.000	0.243 \pm 0.000	0.128 \pm 0.001	0.194 \pm 0.001	0.177

Results on dSprites

Methods	RL → RC	RL → T	RC → RL	RC → T	T → RL	T → RC	Avg.
NNGP [34]	0.313 \pm 0.001	0.438 \pm 0.004	0.356 \pm 0.005	0.515 \pm 0.008	0.367 \pm 0.001	0.324 \pm 0.004	0.386
NTKGP [25]	0.396 \pm 0.001	0.365 \pm 0.001	0.200 \pm 0.007	0.390 \pm 0.003	0.390 \pm 0.000	0.354 \pm 0.003	0.349
AT-GP [7]	0.214 \pm 0.011	0.209 \pm 0.002	0.227 \pm 0.010	0.198 \pm 0.002	0.236 \pm 0.000	0.249 \pm 0.000	0.222
TL-NTK [38]	0.206 \pm 0.004	0.200 \pm 0.002	0.213 \pm 0.000	0.197 \pm 0.000	0.226 \pm 0.001	0.218 \pm 0.000	0.210
DINO-INIT (ours)	0.204 \pm 0.001	0.185 \pm 0.006	0.207 \pm 0.003	0.182 \pm 0.004	0.218 \pm 0.001	0.212 \pm 0.001	0.201
DINO-TRAIN (ours)	0.193 \pm 0.001	0.194 \pm 0.003	0.207 \pm 0.003	0.188 \pm 0.002	0.226 \pm 0.001	0.218 \pm 0.001	0.204

Results on MPI3D

Lower is better



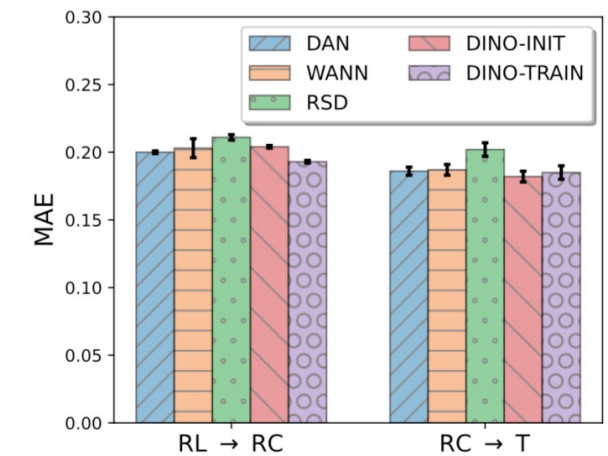
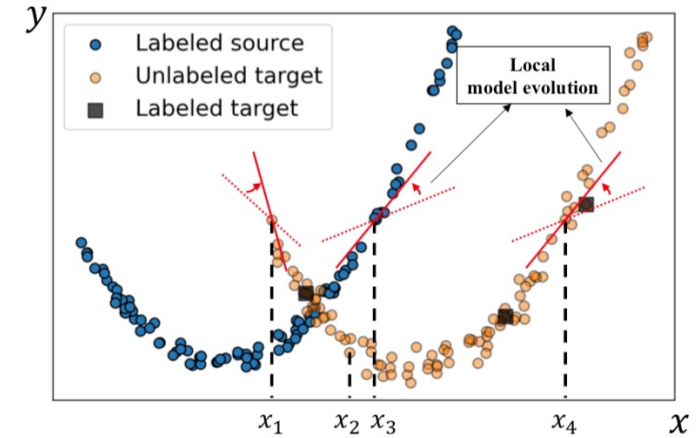
Methods	M → MU	MU → M
NNGP [34]	0.562 \pm 0.001	0.672 \pm 0.010
NTKGP [25]	0.562 \pm 0.004	0.702 \pm 0.010
AT-GP [7]	0.308 \pm 0.006	0.593 \pm 0.025
TL-NTK [38]	0.316 \pm 0.008	0.488 \pm 0.027
DINO-INIT (ours)	0.316 \pm 0.007	0.645 \pm 0.017
DINO-TRAIN (ours)	0.314 \pm 0.009	0.443 \pm 0.030

Results on Plant Phenotyping

Conclusion



- ❑ **Problem:** Domain adaptation regression
- ❑ **Algorithms:** Distribution-informed neural networks
 - **DINO-INIT:** Gaussian process with adaptive NNGP kernel
 - **DINO-TRAIN:** Discrepancy minimization using MMD in the NTK-induced RKHS
- ❑ **Evaluation:** Effectiveness on adaptive regression tasks
 - Object localization
 - Plant phenotyping



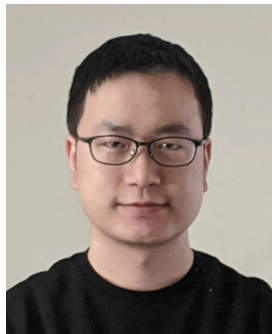


AIFARMS

Artificial Intelligence for Future Agricultural
Resilience, Management, and Sustainability



Distribution-Informed Neural Networks for Domain Adaptation Regression



Jun Wu



Jingrui He



Sheng Wang



Kaiyu Guan



Elizabeth Ainsworth

University of Illinois Urbana-Champaign
{junwu3, jingrui, sheng12, kaiyug, ainswort}@illinois.edu