# **Tunable Domain Adaptation Using Unfolding**

Anonymous authors Paper under double-blind review

## Abstract

Machine learning models often struggle to generalize across domains with varying data distributions, such as differing noise levels, leading to degraded performance. Traditional strategies like personalized training, which trains separate models per domain, and joint training, which uses a single model for all domains, have significant limitations in flexibility and effectiveness. To address this, we propose two novel domain adaptation methods for regression tasks based on interpretable unrolled networks—deep architectures inspired by iterative optimization algorithms. These models leverage the functional dependence of select tunable parameters on domain variables, enabling controlled adaptation during inference. Our methods include Parametric Tunable-Domain Adaptation (P-TDA), which uses known domain parameters for dynamic tuning, and Data-Driven Tunable-Domain Adaptation (DD-TDA), which infers domain adaptation directly from input data. We validate our approach on compressed sensing problems involving noise-adaptive sparse signal recovery and domain-adaptive gain calibration, demonstrating improved or comparable performance to domain-specific models while surpassing joint training baselines. This work highlights the potential of unrolled networks for effective, interpretable domain adaptation in regression settings.

# 1 Introduction

Machine learning models typically learn input-output mappings by optimizing parameters to minimize prediction errors on paired data. Their ability to generalize effectively to unseen data hinges on the assumption that both training and test data are drawn from similar distributions. However, this assumption often breaks down when data varies due to domain-specific factors, such as noise levels, leading to performance degradation. To ensure robust model performance in such settings, it becomes critical to account for this variability during training.

A straightforward strategy to handle domain shifts is *Personalized Training (PT)*, where separate models are trained for each domain. While this yields strong performance, it becomes impractical when the domain parameters vary continuously, and it also requires precise domain information during inference. On the other hand, *Joint Training (JT)* consolidates data from all domains into a single model, eliminating the need for domain-specific information. However, JT typically underperforms compared to PT, as it attempts to learn a unified representation across varying domains, which may not capture domain-specific nuances effectively.

To address the issues arising in the PT and the JT approaches, *transfer learning* and *domain adaptation* techniques have been proposed, often outperforming both PT and JT. In classification tasks, these methods commonly involve retraining only the final layers of a model, under the assumption that the feature extractor is domain-invariant Yosinski et al. (2014); He et al. (2016); Devlin et al. (2019); Pan & Yang (2010). However, identifying which layers should serve as feature extractors is non-trivial Tseng et al. (2023). In He et al. (2024), the authors describe a gradual domain adaptation method using optimal transport to improve classification performance across large domain shifts. In Farahani et al. (2021), Farahani et al. describe domain adaptation methods that address domain shifts by aligning features between domains, learning deep domain-invariant representations, iteratively refining models with pseudo-labels on target data, and using adversarial training to make features indistinguishable across domains. These require separate domain-specific data for progressive retraining during adaptation. But, all these approaches do not extend well to regression tasks—which is the primary focus of this work.

Our goal is to develop domain adaptation methods suitable for regression problems where the underlying assumption is that the domains are parametric. While Deep Neural Networks (DNN) and Convolutional Neural Networks (CNN)-based methods have demonstrated strong performance, their non-interpretable nature obscures the internal mechanisms of the network, making it difficult to determine which parameters should be adjusted for domain adaptation.

In this paper, we address these limitations by adopting *interpretable unrolled networks* Monga et al. (2021) deep learning architectures derived by unrolling iterative optimization algorithms into finite-layer neural networks. These models combine the structure of classical optimization with the flexibility of deep learning, leading to improved interpretability. When trained on data with domain variability, the learnable parameters of unrolled networks often exhibit functional dependence on the domain parameter. Importantly, only a subset of these parameters—referred to as *tunable parameters*—tend to vary significantly with the domain, while the rest remain relatively invariant. Because the network is based on known algorithms, it allows us to uncover the functional relationship between the domain and its corresponding parameters. This allows us to model tunable parameters as functions of the domain parameter, enabling controlled adaptation during inference.

With an objective of designing a single network that effectively adapts to multiple domains, we propose two novel domain adaptation methods:

- Parametric Tunable-Domain Adaptation (P-TDA): This method assumes that domain parameters are either known for each data instance or that there exists a predefined function linking domain parameters to the tunable parameters of the unrolled network. The model is trained to minimize prediction error across domains by adjusting its parameters based on domain-specific characteristics. During inference, the model dynamically adapts by tuning the relevant parameters according to the known domain value.
- Data-Driven Tunable-Domain Adaptation (DD-TDA): In cases where neither the domain parameters nor the parametric relationship is known, this method learns a function that maps input data to domain-specific tunable parameters. During inference, the model simultaneously predicts outputs and adapts its parameters in a self-supervised fashion, using only the input data.

We demonstrate the effectiveness of these methods through two applications in *Compressed Sensing (CS)* Donoho (2006), a domain where the objective is to recover sparse signals from a limited number of measurements. The Iterative Shrinkage-Thresholding Algorithm (ISTA) Daubechies et al. (2004) is a commonly used, interpretable, mathematically derived method for solving the CS problem that can be naturally converted into an unrolled network. Building on this, unrolled optimization methods like Learned-ISTA (LISTA) Gregor & LeCun (2010)—a learned, unrolled variant of ISTA—have gained popularity for their balance of speed and interpretability, modeling each ISTA iteration as a separate layer in the network.

- Noise-Adaptive LISTA: In our first application, we focus on sparse signal recovery under varying noise conditions. Prior work Donoho & Johnstone (1994) has shown a strong relationship between the noise standard deviation and the soft-thresholding parameter in LISTA. Here, the *domain parameter* is the noise standard deviation (assuming Gaussian noise), and the *tunable parameter* is the soft-threshold parameter. Both P-TDA and DD-TDA are applied in this context. P-TDA leverages known noise levels to train a black-box model for dynamically adjusting the soft-thresholding parameter. In contrast, DD-TDA infers the relationship directly from data, learning how to adapt without access to the noise level during inference. Several experiments revealed that our methods perform at least as well as PT—and often better—in terms of the defined hit rate metrics, while also achieving lower estimation error compared to the JT method.
- Domain-Adaptive Gain Calibration: In our second application, we extend the proposed framework to a different compressed sensing problem involving multiplicative distortions modeled via unknown sensor gains. Here, the domain variability arises due to the changes in sensor gains. Using explicit gain information (P-TDA) or data directly (DD-TDA), we demonstrate that both P-TDA and DD-TDA can efficiently adjust to such gain-induced variability by learning gain-compensated network

parameters. Our results demonstrated that the suggested adjustable models beat joint training baselines in both structured and random gain cases, while achieving performance competitive with domain-specific models.

The paper is organized as follows. In Section 2, we formulate the general learning problem in a domainparameterized setting. In Section 3, we present the proposed domain-adaptation method through unrollingbased networks. Two applications of the proposed approach for noise adaptation and gain calibration are discussed in Sections 4 and 5, respectively, followed by conclusions.

### 2 Problem Formulation

Consider a learning problem from a set of paired data points  $\mathcal{D} = \{\mathbf{y}_n, \mathbf{x}_n\}_{n=1}^N$  where  $\mathbf{y}_n \in \mathbb{C}^{N_y}$  is input or feature to a network to be learned and  $\mathbf{x}_n \in \mathbb{C}^{N_x}$  is the corresponding output or label. Here, N represents the sample count of the dataset. The objective is to learn a function or a network  $f_{\alpha,\beta} : \mathbb{C}^{N_y} \to \mathbb{C}^{N_x}$ , with the learnable parameters  $(\alpha, \beta)$ , such that the estimate  $f_{\alpha,\beta}(\mathbf{y}_n)$  is close to the true output  $\mathbf{x}_n$  in a predefined distance metric. The goal is typically achieved by fine-tuning the learnable parameters by solving the following optimization problem:

$$\min_{\boldsymbol{\alpha},\boldsymbol{\beta}} \sum_{(\mathbf{y}_n, \mathbf{x}_n) \in \mathcal{D}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha},\boldsymbol{\beta}}(\mathbf{y}_n)\right),\tag{1}$$

where  $d(\cdot, \cdot)$  is a distance metric in  $\mathbb{C}^{N_x}$ . The parameters learned through the process are optimal for the training data  $\mathcal{D}$ . It is expected that for any unseen data pairs, the learned network will provide fairly accurate results, given that the unseen or test data is similar to the training data. The similarity is generally ascertained by assuming that the train and test data are samples from the same probability distribution.

In many practical applications, the data set could be parameterized as  $\mathcal{D}_{\theta}$  where  $\theta$  is the parameter vector whose entries could be either discrete or continuous values. For example,  $\theta$  could represent the noise level in the features/input. Since the data changes with  $\theta$ , an important question arises on how to train the network in this scenario. For the simplicity of the discussion, let us assume that  $\theta$  can take on one of the following J values:  $\{\theta_j\}_{j=1}^J$ . Then we can have the following training options.

The first is to train individual networks for each data set  $\mathcal{D}_{\theta_j}$ . Such PT approach results in the learnable parameters  $\{\alpha_j, \beta_j\}$  as a result of the following optimization problem

$$\arg\min_{\boldsymbol{\alpha},\boldsymbol{\beta}} \sum_{(\mathbf{y}_n,\mathbf{x}_n)\in\mathcal{D}_{\boldsymbol{\theta}_j}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha},\boldsymbol{\beta}}(\mathbf{y}_n)\right).$$
(2)

The PT method requires training J networks, which may not be practically feasible, especially when  $\theta$  is a continuous variable. In addition, during inference, the precise value of  $\theta$  should be known such that the corresponding network is used. Any mismatch may amplify the inference error.

To reduce the number of networks to be trained, a JT strategy could be used. Here, a single network is trained for all the data as

$$\{\boldsymbol{\alpha}_{\text{joint}}, \boldsymbol{\beta}_{\text{joint}}\} = \arg\min_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \sum_{j=1}^{J} \sum_{(\mathbf{y}_n, \mathbf{x}_n) \in \mathcal{D}_{\boldsymbol{\theta}_j}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha}, \boldsymbol{\beta}}(\mathbf{y}_n)\right).$$
(3)

Unlike PT, in the JT approach, knowledge of the precise value of  $\theta$  is not required. However, this approach may have an inferior performance compared with PT due to combined training. The reduction in performance depends on how the data parameter  $\theta$  is related to the data and hence, varies for different models.

An alternative and, perhaps, a more suitable solution to the problem is to apply methods from transfer learning and domain adaptation. In transfer learning, a general-purpose pre-trained network, let us say, on data set  $\mathcal{D}_{\theta_1}$ , is retrained for a new data set, say  $\mathcal{D}_{\theta_2}$  Yosinski et al. (2014); He et al. (2016); Devlin et al. (2019); Pan & Yang (2010). While retraining, most layers are fixed, and the remaining ones, usually the last few, are relearned for the new data. For example, let  $f_{\alpha_1,\beta_1}$  is the network learned on  $\mathcal{D}_{\theta_1}$ . Then, during fine-tuning or retraining, the parameters  $\alpha_1$  are kept fixed, and the parameters  $\beta$  are learned. The pretraining and fine-tuning steps for the example considered here are given below.

Pre-training: 
$$\{\boldsymbol{\alpha}_1, \boldsymbol{\beta}_1\} = \arg\min_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \sum_{(\mathbf{y}_n, \mathbf{x}_n) \in \mathcal{D}_{\boldsymbol{\theta}_1}} d(\mathbf{x}_n, f_{\boldsymbol{\alpha}, \boldsymbol{\beta}}(\mathbf{y}_n)).$$
  
Fine-tuning:  $\{\boldsymbol{\beta}_2\} = \arg\min_{\boldsymbol{\beta}} \sum_{(\mathbf{y}_n, \mathbf{x}_n) \in \mathcal{D}_{\boldsymbol{\theta}_2}} d(\mathbf{x}_n, f_{\boldsymbol{\alpha}_1, \boldsymbol{\beta}}(\mathbf{y}_n)).$ 

The fine-tuning-based method to adapt a network for new data is typically applied in classification tasks where the first few layers are assumed to act as feature extractors. The last few layers, which are retrainable, work as classifiers. The features could be the same for different data sets; hence, only the classifier must be tuned. However, in conventional DNNs/CNNs, it is not straightforward to differentiate between feature extractor layers and classification layers. The method may also not be directly applicable to regression tasks, which are of interest in this work.

In the previous approach, once the network is pre-trained on data  $\mathcal{D}_{\theta_1}$ , the data is no longer used during fine-tuning for dataset  $\mathcal{D}_{\theta_2}$ . Note that the data from  $\mathcal{D}_{\theta_1}$  and  $\mathcal{D}_{\theta_2}$  have different distributions, and in the fine-tuning approach, there is no emphasis on distribution alignment. On the other hand, in the supervised domain-adaptation methods, data from all the distributions are used simultaneously to train a network, as in equation 3, with an additional term that explicitly reduces distribution differences between the data Long et al. (2013; 2015); Motiian et al. (2017). Though these methods seem suitable for the problem considered in this paper, the approaches do not consider the parametric relationship among the data, which is prevalent in many applications. Further, most notable works in transfer learning and supervised-domain-adaptation focus on classification tasks, and it is unclear how to extend them to regression problems, except a few works Tolooshams et al. (2021); Lu et al. (2019); Li et al. (2023); Zhang et al. (2024); Chen et al. (2022); Pardoe & Stone (2010). Again, these works treat each domain independently and do not consider any underlying parametric relationship.

# 3 Domain-Adaptation Through Unrolling

In this section, we discuss the proposed method for adapting a network from the data. In the fine-tuning approach discussed in the previous section, only  $\beta$  is learned for the new data. As discussed, in DNNs/CNNs, though there is intuitive reasoning about it, it has not been proven theoretically why only a certain part of the networks should be tuned. Additionally, the layers of a DNN or a CNN are not interpretable, which makes it challenging to prove the tuning part. Unlike DNNs/CNNs, unrolling (or unfolding)-based methods, which are deep learning models that are derived by unfolding an iterative optimization algorithm into a finite number of layers, are interpretable Gregor & LeCun (2010); Monga et al. (2021). These networks incorporate domain-specific knowledge from optimization techniques and leverage the power of deep learning for efficient and interpretable learning. By noting that the parametric relationship between the domains stems from physical models that generate data and unrolling networks are domain and model-specific, we show that fine-tuning can be made interpretable and efficient.

To elaborate on the interpretable domain adaptation via unrolling framework, let  $f_{\alpha,\beta}(\cdot)$  be the unfolded network. Unlike standard DNNs/CNNs, the learnable parameters  $\{\alpha, \beta\}$  in this network are model-specific and are interpretable. When the network is trained on parametric data  $\mathcal{D}_{\theta}$ , the network's parameters  $\{\alpha, \beta\}$ are function of  $\theta$ . In the following, we assume that among these,  $\alpha$  changes negligibly with  $\theta$ , whereas,  $\beta$ strongly depends on  $\theta$ . Mathematically, there exist a function  $g(\cdot)$  such that

$$\boldsymbol{\beta} = g_{\boldsymbol{\varphi}}(\boldsymbol{\theta}),\tag{4}$$

where  $\varphi$  are parameters of the function. The function  $g(\cdot)$  can be either derived from the physics of the generating model or learned from the available data. The assumption of the existence of  $g_{\varphi}$  is the same as in the fine-tuning approach. However, the assumption holds in practice for many unrolling networks and parametric data, as shown in the following sections. A P-TDA approach is used for known g and  $\theta$ , to learn

the network as

$$\min_{\boldsymbol{\alpha}} \sum_{j=1}^{J} \sum_{(\mathbf{y}_n, \mathbf{x}_n) \in \mathcal{D}_{\boldsymbol{\theta}_j}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha}, g_{\boldsymbol{\varphi}}(\boldsymbol{\theta}_j)}(\mathbf{y}_n)\right).$$
(5)

The resulting network  $f_{\alpha,g_{\varphi}(\theta)}(\cdot)$  can then be adapted to test data coming from any domain by substituting appropriate  $\theta$ . When  $g_{\varphi}$  is unknown and  $\theta$  is known, it can be learned by solving the problem

$$\min_{\boldsymbol{\alpha},\boldsymbol{\varphi}} \sum_{j=1}^{J} \sum_{(\mathbf{y}_n,\mathbf{x}_n)\in\mathcal{D}_{\boldsymbol{\theta}_j}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha},g_{\boldsymbol{\varphi}}(\boldsymbol{\theta}_j)}(\mathbf{y}_n)\right).$$
(6)

In practice,  $g_{\varphi}$  is realized using a neural network where the inputs are the domain information  $\theta_j$  as well as the measurements  $\mathbf{y}_n$ s. In addition, the parameters to the forward model can be used. We observe that the additional information of the measurements improves the performance of the overall tunable network.

In the case of unknown g and  $\theta$ , we propose the following DD-TDA approach. By noting that the data is implicitly dependent on  $\theta$ , we train a function  $h_{\varphi}(\cdot)$  which is expected to estimate  $\beta$  from the data. Specifically, we optimize the following problem

$$\min_{\boldsymbol{\alpha},\boldsymbol{\varphi}} \sum_{j=1}^{J} \sum_{(\mathbf{y}_n,\mathbf{x}_n)\in\mathcal{D}_{\boldsymbol{\theta}_j}} d\left(\mathbf{x}_n, f_{\boldsymbol{\alpha},h_{\boldsymbol{\varphi}}(\mathbf{y}_n)}(\mathbf{y}_n)\right).$$
(7)

The proposed method results in a single tunable network  $f_{\alpha,h_{\varphi}(\cdot)}(\cdot)$ . Note that during inference, for any test input **y**, the network not only predicts the output but also tunes the network via the function  $h_{\varphi}$ . As in the previous model, the function  $h_{\varphi}$  is realized using a neural network.

In the aforementioned formulation, the network  $f_{\alpha,h_{\varphi}(\cdot)}(\cdot)$  is trained over the domains  $\{\theta_j, j = 1, \dots, J\}$ . The network is expected to adapt well to any data from the trained domains. However, whether it will perform well on an unseen domain  $\theta \notin \{\theta_j, j = 1, \dots, J\}$  depends on the generalization ability of the function  $h_{\varphi}(\cdot)$ . As we will show in the following sections, the generalization ability depends on the task and how the domain parameter depends on the network parameter. In general, the generalization is good for large J, as with any learned network, where the generalization typically improves with the data size.

In the following sections, we consider two practical scenarios of the aforementioned framework. In particular, we discuss domain adaptation for compressive sensing in the next section, followed by adaptation for blind calibration problems.

### 4 Noise-Adaptive Tunable LISTA

This section considers the first problem of domain adaptation through tunable unrolled networks. The problem is to recover sparse vectors from their compressed measurements under different noise conditions. Before detailing the problem and the proposed solution, we first give a quick background of the compressive sensing problem, its solution, and LISTA.

#### 4.1 The Compressed Sensing Problem

Consider a set of noisy linear measurements as

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\eta},\tag{8}$$

where  $\eta \in \mathbb{R}^{N_y}$  is an additive Gaussian noise vector with zero mean and variance  $\sigma^2$ . Here,  $\mathbf{y} \in \mathbb{R}^{N_y}$  representing the measurement signal,  $\mathbf{x} \in \mathbb{R}^{N_x}$  denoting the unknown signal to be estimated, and  $\mathbf{A} \in \mathbb{R}^{N_y \times N_x}$  signifying the measurement matrix with  $N_y < N_x$ . The underdetermined system of equations can

be solved with a sparsity assumption on  $\mathbf{x}$ , that is,  $\|\mathbf{x}\|_0 \leq K$  where  $K < N_y$ . A widely used approach to solve the problem is to minimize the following  $\ell_1$ -minimization problem:

$$\min_{\mathbf{y}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \lambda \|\mathbf{x}\|_{1}.$$
(9)

The objective function encourages solutions that balance data fidelity, captured by the first term, and sparsity, controlled by  $\|\mathbf{x}\|_1$ . The balance can be achieved by choosing the hyperparameter  $\lambda$  appropriately. Several practical algorithms, such as the basis pursuit, iterative hard-thresholding, and ISTA Daubechies et al. (2004); Beck & Teboulle (2009), have been developed to solve the optimization problem.

Most of the aforementioned algorithms are iterative and, starting from an initial estimate, refine the estimate in each iteration. For example, let  $\mathbf{x}^k$  be the current estimate of the sparse signal at iteration k,  $\mathbf{x}^{k+1}$  is the updated estimate for the next iteration. The update rule in ISTA is given as

$$\mathbf{x}^{k+1} = \mathcal{S}_{\beta} \left( \mathbf{W}_1 \mathbf{y} + \mathbf{W}_2 \mathbf{x}^k \right), \tag{10}$$

with

$$\mathbf{W}_1 = \frac{1}{\kappa} \mathbf{A}^T, \quad \mathbf{W}_2 = \mathbf{I}_{N_x} - \frac{1}{\kappa} \mathbf{A}^T \mathbf{A}, \quad \text{and} \quad \beta = \frac{\lambda}{\kappa}, \tag{11}$$

where  $\mathbf{I}_{N_x}$  is the  $N_x \times N_x$  identity matrix and  $\mathcal{S}_{\beta}$  is an elementwise soft-thresholding operation with threshold parameter  $\beta = \frac{\lambda}{\kappa}$  where  $\kappa$  is greater than the largest eigenvalue of  $\mathbf{A}^T \mathbf{A}$ .

While ISTA and other algorithms work well in practice, their convergence is slow. In addition, the methods are data-independent. On the other hand, given a set of data with some underlying structure (for example, drawn from the same distribution), the algorithms should be able to be fine-tuned to improve the performance and convergence. This is precisely the LISTA Gregor & LeCun (2010) algorithm where the matrices  $\alpha = \{\mathbf{W}_1, \mathbf{W}_2\}$  together with the thresholding parameter  $\beta$  are learned from the data by keeping the number of iterations fixed. The LISTA is shown to converge in fewer iterations than the original ISTA.

Before we discuss the proposed noise-adaptive LISTA, we quickly summarize a few results on adaptive LISTA approaches and compare them with the current one.

### 4.2 Adaptive LISTA

In the original LISTA framework Gregor & LeCun (2010), the network is trained to learn the parameters  $\{\mathbf{W}_1, \mathbf{W}_2, \beta\}$  (cf. equation 11) using data generated from a specific distribution and a known measurement matrix **A**. Changes in **A** lead to changes in the distribution, affecting the measurements  $\mathbf{y}_n$  and necessitating network retraining. To address this, Aberdam et al. Aberdam et al. (2021) introduced an adaptive LISTA algorithm where matrices  $\mathbf{W}_{1,2}$  are decomposed into two parts: a fixed component dependent on **A** and a learnable component similar to traditional LISTA. This enables training on data generated by multiple choices of **A**, optimizing the network accordingly. At inference, the network adapts to new **A** before prediction, demonstrating improved performance across varying **A** values with theoretical guarantees. However, this method requires precise domain parameter knowledge (here **A**) during inference, and the matrix decompositions introducing domain dependence are less straightforward, leaving room for alternative approaches Chen et al. (2018); Liu & Chen (2019); Chen et al. (2021).

While our focus differs, our DD-TDA framework does not necessitate prior domain parameter knowledge. In cases where it is known, we optimize its use akin to the P-TDA approach.

#### 4.3 Noise-Level-Specific Domains

In this problem, noise variance parameterized the domains, and the domain-specific measurements are given as

$$\mathcal{D}_{\theta_j = \sigma_j} = \{ \mathbf{y}_{ij}, \mathbf{x}_i \}_{i=1}^N \ \forall \ j = \{1, 2, \dots, J\},$$
(12)

			-
Training approach	D1(6  dB)	D2(16  dB)	D3(31  dB)
	NMSE (in dB), $HR(\%)$		
Joint train	<b>-6</b> .60, 54.22	-12.34, 76.72	-13.33, 79.81
Retrain last layer	<b>-6.72</b> , 53.73	-12.64, 77.87	-13.79, 81.17
Retrain last 3 layers	-7.10, 53.93	-13.92, 83.40	-16.17, 89.24

Table 1: Test results of retrained DNNs for domain adaptation

where

$$\mathbf{y}_{ij} = \mathbf{A}\mathbf{x}_i + \boldsymbol{\eta}_j, \quad \boldsymbol{\eta}_j \sim \mathcal{N}(\mathbf{0}, \sigma_j^2 \mathbf{I}_{N_y}), \tag{13}$$

and  $\mathbf{I}_{N_y}$  is an identity matrix of dimension  $N_y \times N_y$ 

Specifically, for each domain, the noise levels are different while the rest of the measurement model remains fixed. The measurement model is practically relevant, as the noise levels can vary significantly. As discussed in Section 2, both the PT and JT approaches have limitations. Here, we consider the problem of fine-tuning the LISTA model to adapt it to different noise domains.

For ease of discussion, the different domains, based on the noise levels, are also characterized by average SNR levels defined as

$$SNR_{j} = \frac{1}{N} \sum_{i=1}^{N} 10 \log_{10} \left( \frac{\frac{1}{N_{y}} \| \mathbf{A} \mathbf{x}_{i} \|_{2}^{2}}{\sigma_{j}^{2}} \right).$$
(14)

### 4.4 A Retraining Approach

Before introducing our noise-tunable LISTA framework, we first examine the effectiveness of a retrainingbased method that utilizes a DNN to address the noise-tunability challenge. The experiment is conducted on a dataset with three domains with SNRs of 6 dB, 16 dB, and 32 dB. Data is generated following equation 12 and equation 13 (see the data generation part of Section 4.5 for details). We start by training a fully connected neural network on a combined dataset containing samples from all domains. The network consists of three hidden layers with 256, 512, and 256 neurons, respectively, each followed by ReLU activations. For domain-specific adaptation, we perform staged retraining by progressively unfreezing deeper layers of the network on the target domain data. For retraining, we use domain-specific data that was previously included in a mixed-domain dataset during the initial training phase. Initially, only the final layer is retrained; then, we freeze the first hidden layer and fine-tune the remaining layers, including the output layer. Results in Table 1 show that retraining just the final layer does not yield significant improvements across varying noise levels. However, as more layers are made trainable, performance on the test data, measured in terms of NMSE and HR (defined in equation 16 and equation 17), improves consistently. These findings suggest that adapting only the final layer is insufficient for effective tunability. Instead, identifying the optimal subset of layers to retrain may require a trial-and-error approach, which can be task-specific and assumes access to domain-specific data for adaptation.

### 4.5 Proposed Noise-Adaptive LISTA (NA-LISTA)

Unlike the previous DNN-based methods, we show that LISTA is tunable to noise levels. Specifically, we note that among the learnable parameters of LISTA,  $\boldsymbol{\alpha} = \{\mathbf{W}_1, \mathbf{W}_2\}$  and  $\beta$ , the latter one is a strong function of the noise variance  $\sigma$  for a fixed  $\mathbf{A}$ , that is, we have that  $\beta = g(\sigma)$ . There are several choices of g Galatsanos & Katsaggelos (1992) and one of the widely used ones is given as Donoho & Johnstone (1994)

$$\beta = \sigma \sqrt{2\log N_x} / \kappa. \tag{15}$$



Figure 1: Layers/iterations of different LISTA architectures: (a) Conventional LISTA, (b), (c) NA-LISTA.



Figure 2: The network architechture for  $g_{\varphi}$  with output dimensions of each layer in red. The  $h_{\varphi}$  network architecture is same with only **y** and **A** as inputs. Here, BN stands for batch normalization.

The above formulation implies that the higher the noise variance, the larger the soft thresholding is. Since soft-thresholding operation promotes sparsity, larger  $\beta$  generates sparser estimates. The choice in equation 15 requires precise knowledge of  $\sigma$ , which may not be available in practice. An alternative approach is cross-validation Golub et al. (1979) where  $\beta$  is learned from validation data. The learned parameter is not optimal if the validation data consists of measurements with various noise levels.

To address the limitations mentioned above, we apply the proposed P-TDA (assuming a known  $\sigma$ ) and DD-TDA methods. One iteration or layer of these two methods, including the conventional LISTA, are shown in Fig. 1. In DD-TDA,  $\{\mathbf{A}, \mathbf{y}, \sigma\}$  are used to tune the network by altering  $\beta$ , whereas, in DD-TDA framework, only  $\{\mathbf{A}, \mathbf{y}\}$  are given as input to  $h_{\varphi}$  to tune  $\beta$ . In the following, we discuss the data generation and results.

**Data Generation and Network Architecture**: To evaluate the proposed NA-LISTA, we generated a synthetic dataset as per equation 12–equation 13. The matrix **A** was kept fixed in all the experiments, and its entries were sampled from a standard normal distribution followed by columnwise normalization. Further, we set  $N_x = 100$ ,  $N_y = 30$ , and K = 3. The entries of the noise vector were zero-mean Gaussian random variables with domain-specific standard deviations. Specifically, noise variance  $\sigma_j$  simulates J distinct SNR levels for robustness assessment. The total number of samples generated was split into train, validation, and test sets in the ratio of 56%, 24%, and 20%, respectively.

Since samples from all the domains were used to train and test the JT, P-TDA. and DD-TDA methods, in order to maintain dataset size uniformity, we used the same number of samples for each PT model. For example, for the J domains used for training, we generated  $N_{train}$  (training set slice of N) examples from each domain, and used  $N_{train} \times J$  training samples for JT, P-TDA, and DD-TDA methods. Whereas, we generated  $N_{train} \times J$  examples for training each domain-specific PT model. In the following experiments, we consider J = 3 and N = 43,000.

Fig. 2 depicts the network architecture used to realize the functions  $g_{\varphi}$  and  $h_{\varphi}$  for P-TDA and DD-TDA methods. The dimensions of the outputs of each layer are marked.

LISTA networks with the proposed  $g_{\varphi}$  and  $h_{\varphi}$  networks were jointly trained in an end-to-end fashion using backpropagation, and the mean-squared error (MSE) loss was used as the cost function. The MSE was measured by averaging the Euclidean distance  $\|\mathbf{x}^* - \hat{\mathbf{x}}\|_2^2$  over the training batch, where  $\mathbf{x}^*$  and  $\hat{\mathbf{x}}$  denote the ground truth and predicted sparse signals, respectively. The network parameters were initialized as



(d) Domains with block bit tanget. D1 ( $\sigma = 0.03$ , SNR = 16 dB), and D3( $\sigma = 0.005$ , SNR = 32 dB).

(b) Domains with narrow SNR ranges: D1 ( $\sigma = 0.12$ , SNR = 4 dB), D2 ( $\sigma = 0.06$ , SNR = 10 dB), and D3( $\sigma = 0.035$ , SNR = 15 dB).

Figure 3: Comparison of the proposed methods with the JT and PT methods for NA-LISTA for different SNR ranges. PT's performances are better than JT's for each domain. The P-TDA/DD-TDA has comparable performance to PT, demonstrating the noise-tunability aspect.

mentioned in equation 11, where  $\lambda$  was initialized to 0.1 and  $\kappa$  was 1.001 times the maximum absolute eigenvalue of  $\mathbf{A}^T \mathbf{A}$ .

**Performance Metrics**: Model performance was evaluated using Normalized Mean Squared Error (NMSE) and Hit Rate (HR). Let  $\mathbf{x}^*$  be the ground truth signal and  $\hat{\mathbf{x}}$  be the predicted sparse signal. The NMSE contribution by each test sample is given by

NMSE = 
$$\frac{\|\mathbf{x}^* - \hat{\mathbf{x}}\|_2^2}{\|\mathbf{x}^*\|_2^2}$$
. (16)

For our experiments, the average NMSE over the entire test set was computed and reported in decibels.

The HR, with a relative error tolerance t, is defined as

$$HR = \frac{\sum_{n=1}^{N_x} \mathbb{1}(\mathbf{x}^*(n) \neq 0) \mathbb{1}(|\mathbf{x}^*(n) - \hat{\mathbf{x}}(n)| \le t \, \mathbf{x}^*(n))}{K}.$$
(17)

Here, n indexes the elements of the vectors  $\mathbf{x}^*$  and  $\mathbf{\hat{x}}$ . HR measures the fraction of nonzero entries in  $\mathbf{\hat{x}}$  whose estimates fall within a relative error tolerance of t, set to 0.3 in our experiments. For our experiments, the average HR (in %) over the entire test set is reported.

**Experiments**: We conducted three evaluations. The first one focused on a broad range of SNRs, another covers a lower SNR range, and a third for testing generalization by evaluating on unseen SNR levels. In each of these, our proposed DD-TDA and P-TDA methods were compared with the baseline JT and PT methods. For both JT and PT, conventional LISTA networks were trained.

**Comparison with the existing methods:** For the broad range we considered three noise levels as D1 ( $\sigma = 0.1$ , SNR = 6 dB), D2 ( $\sigma = 0.03$ , SNR = 16 dB), and D3( $\sigma = 0.005$ , SNR = 32 dB). For the narrow range, the SNRs are given as D1 ( $\sigma = 0.12$ , SNR = 4 dB), D2 ( $\sigma = 0.06$ , SNR = 10 dB), and D3 ( $\sigma = 0.035$ , SNR = 15 dB). For these two scenarios, the NMSEs and HRs of the methods when tested for each domain are shown in Figs. 3(a) and 3(b). Here, PT-D1, PT-D2, and PT-D3, refers to PT methods trained for domains D1, D2, and D3, respectively. We observe that PT methods trained for a specific domain perform well when tested on the same domain. The performance deteriorates when trained and tested on different

rable = c (11) and $p$ values for $r = models$				
Domains (SNR)	D1 (6 dB)	D2 (16 dB)	D3 (32 dB)	
$S(\mathbf{W}_1)$	16.35	16.1	15.55	
$S(\mathbf{W}_2)$	15.07	17.77	15.62	
$\beta$	0.088	0.048	0.037	

Table 2:  $S(\mathbf{W})$  and  $\beta$  values for PT models



Figure 4: Results for generalization experiment. Like before, we can see that DDTDA and PTDA perform much better than JT and as well as PT in the hit rate (HR) metric.

domains. Further, PT gives lower error (better HR) than JT for all the domains, and this observation is more evident in the HRs, where PT shows a performance improvement of 6 - 12% over the JT for D2 and D3. We note that the JT and PT perform similarly on D1, likely because PT-D1 was trained only on highly noisy samples, unlike JT, which saw both clean and noisy data, leading to PT-D1 underperforming relative to expectations.

Comparing the proposed P-TDA and DD-TDA with the PT method, we observe that the proposed approaches are performing on par with the respective PT methods for each domain. Focusing on the DD-TDA method, which does not consider the domain knowledge, and has comparable performance with the PT methods. The only notable performance gain of PT over our methods is that in D3 during the broad SNR range experiment, with an NMSE gain of around 5 dB, likely due to PT-D3 having access to significantly more clean training examples than the DD-TDA method. The results show the tunability aspect of the NA-LISTA for different noise levels.

Having assessed the fact that few network parameters are tunable with the domains, next, we show that the remaining parameters are nearly invariant to the domains.

**Invariance of**  $\alpha$  with  $\theta$ : Here, we asses the variation of the network parameters  $\alpha = {\mathbf{W}_1, \mathbf{W}_2}$  and threshold  $\beta$  with the noise levels  $\theta = \sigma$ . For NA-LISTA, it was assumed that  $\alpha$  does not vary significantly with  $\sigma$  and  $\beta$  changes significantly with domains.

To assess the invariance of  $\alpha$  across domains, we define the normalized Frobenius norm metric. To compare the structure of learned weight matrices across trained models, we define a scale-invariant metric such that it normalizes the Frobenius norm by the maximum absolute value in the matrix. Mathematically, we measured the metric  $S(\mathbf{W}) = \frac{\|\mathbf{W}\|_F}{\max \|\mathbf{W}\|}$  and the values are shown in Table 2. We note negligible changes in the metric for both matrices when the networks are trained for different domains. On the other hand, the values of  $\beta$  have a significant change with domains, justifying the assumption. Next, we discuss the generalization ability of the proposed approaches.



Figure 5: Tunable model for blind calibration gain with **B** as a tuned parameter.

**Generalization:** Building on the previous experiments, which evaluated the model's performance on domains with known SNR levels, we now assess its ability to generalize to unseen domains. Unlike before, the model is tested on SNR levels that were not part of the training set, allowing us to evaluate its adaptability and robustness in entirely new noise conditions. We considered six domains, D1-D6, with average SNRs (in dB) of {4, 10, 15, 20, 26, 32}. All the models, JT, PT, DD-TDA, and P-TDA, were trained on domains D1, D3, and D5 while they were tested on D2, D4, and D6.

Since this experiment evaluates the ability to generalize to unseen domains, we compare our proposed models specifically against the JT method, focusing on their respective generalization capabilities. The results, presented in Fig. 4, demonstrate that our methods outperform the baseline JT approach, especially in the HR metric by around 9-16%. On comparing the P-TDA and DD-TDA methods with PT methods, we observe that PT methods result in a lower error of 3-5 dB for D4 and D6, which are high SNR domains. Whereas, for D2, with SNR of 10 dB, the P-TDA and DD-TDA have comparable performance to the PT approach.

In the following, we consider another application for the P-TDA and DD-TDA methods.

### 5 Domain-Adaptive Sparse Gain Calibration

Consider the measurements of the form

$$\mathbf{y} = \operatorname{diag}(\mathbf{c})\mathbf{A}\mathbf{x} + \boldsymbol{\eta},\tag{18}$$

where, except  $\mathbf{c} \in \mathbb{R}^{N_y}$ , the rest of the parameters have the same dimensions and characteristics as in equation 8. Here  $\mathbf{c}_i$  is the gain of  $\mathbf{y}_i$ -th measurement. In particularly,  $\mathbf{x}$  is K-sparse. The measurement model is prevalent in several applications, such as time-of-flight imaging Steiger et al. (2008); Mersmann et al. (2013), direction of arrival estimation with unequal sensor gains Paulraj & Kailath (1985), and more. An unknown  $\mathbf{c}$  amounts to a sparse-blind calibration problem Schulke et al. (2013). In Tolooshams et al. (2023), a data-driven approach, based on LISTA, was proposed to solve for an unknown but fixed  $\mathbf{c}$ . However, the network has to vary as  $\mathbf{c}$  varies, which is the practice case. Here, we consider the applicability of the proposed tunability approach by considering the domain parameter as  $\mathbf{c}$ . Specifically, the data for the parametric domains are defined as

$$\mathcal{D}_{\boldsymbol{\theta}=\mathbf{c}_j} = \{\mathbf{y}_{ij}, \mathbf{x}_i\}_{i=1}^N \ \forall \ j = \{1, 2, \dots, J\},\tag{19}$$

where

$$\mathbf{y}_{ij} = \operatorname{diag}(\mathbf{c}_j) \mathbf{A} \mathbf{x}_i + \boldsymbol{\eta}_j, \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{N_{\eta}}).$$
(20)

Note that the noise variance is constant across domains, unlike the parametric-domain model in equation 12 and equation 13. However, this is not a limitation; noise  $(\eta)$  and gain (c) could be the domain parameters. Since our objective is to show the tunability aspect of the network as c varies,  $\sigma$  is kept fixed. Next, we discuss the details of the data generation and network architecture, followed by the simulation results.



Figure 6: The network architechture for  $g_{\varphi}$  with output dimensions of each layer in red. The  $h_{\varphi}$  network architecture is the same with only y as inputs.

#### 5.1 The Proposed Approaches and Network Architecture

We first propose an approach that is applicable to both the PT and JT methods. Then, we suggest a modification to make the method domain-adaptive.

There are several possible approaches to solving the sparse-gain calibration problem. For example, one could treat  $\operatorname{diag}(\mathbf{c})\mathbf{A}$  as the measurement matrix and then apply the conventional LISTA algorithm. Alternatively, one could learn the matrices  $\operatorname{diag}(\mathbf{c})$  and  $\mathbf{A}$  separately, as in Tolooshams et al. (2023). The former approach lacks adaptability, as the matrix  $\operatorname{diag}(\mathbf{c})$  is not learned independently. In contrast, the latter approach leads to a more complex network, since the information in  $\operatorname{diag}(\mathbf{c})$  is embedded within each layer of LISTA.

Another alternative, when  $\mathbf{c}$  is known, is to solve the optimization problem in equation 8 using LISTA, after normalizing the measurements via  $\hat{\mathbf{y}} = \text{diag}(\mathbf{c}^{-1})\mathbf{y}$ . However, this may significantly amplify noise in measurements corresponding to small values of  $\mathbf{c}$  — a well-known issue in deconvolution problems. Moreover, the exact values of  $\mathbf{c}$  may not be known in practice.

We address these issues by proposing an approach based on Wiener filtering for deconvolution, combined with a learning strategy. Specifically, instead of using the normalized measurements diag( $\mathbf{c}^{-1}$ ) $\mathbf{y}$ , we learn a vector  $\mathbf{b}$  jointly with the LISTA parameters, such that diag( $\mathbf{b}$ ) $\mathbf{y}$  is used as the input to LISTA. As in Wiener filtering, the entries of the learned  $\mathbf{b}$  are expected to be small when the SNR for the corresponding measurement is low. This approach mitigates the problems associated with directly using  $\mathbf{c}^{-1}$ . Furthermore, deconvolution is performed only once, after which the subsequent layers follow the conventional LISTA structure (cf. Fig. 5(a)). The network can be trained either jointly across all domains or individually per domain to learn the parameters  $\mathbf{b}, \mathbf{W}_{1,2}$  and the soft-thresholding parameter  $\boldsymbol{\lambda}$ .

Since the domains are parameterized by the calibration vector  $\mathbf{c}$ , enabling adaptability or tunability requires identifying the parameters that vary significantly with  $\mathbf{c}$ . In Fig. 5(a), the matrix  $\mathbf{B}$  is introduced to compensate for the gain factor, while the rest of the network retains the structure of the traditional LISTA. Moreover, in the noise-free case, if one were to solve the problem using conventional ISTA,  $\mathbf{B}$  would act as  $\mathbf{c}^{-1}$ . Hence, the calibration gains  $\boldsymbol{\beta} = \mathbf{B}$  are expected to vary significantly with  $\mathbf{c}$ , in contrast to the parameters  $\boldsymbol{\alpha} = \mathbf{W}_{1,2}, \lambda$ , which are relatively invariant.

Motivated by this insight, we propose the P-TDA and DD-TDA methods, as shown in Figs. 5(b) and (c), respectively. In the P-TDA method, **B** is learned from both the known **c** and the measurements using the network  $g_{\varphi}()$ . In contrast, in the DD-TDA approach, **B** is inferred solely from the measurements using the network  $h_{\varphi}()$ .

**Data Generation**: For the experiment, we considered two sets of parametric domain datasets following the model in equation 19. The first set uses structured gains, where  $\mathbf{c}$  follows a sinusoidal distribution. Specifically, the entries of  $\mathbf{c}$  for each domain were generated as

$$\mathbf{c}_{i}(n) = \left| a \sin\left(2\pi f_{i} n + \phi_{i}\right) + b \right|. \tag{21}$$

The variation in  $\mathbf{c}_j$  across domains is governed by changes in  $\{f_j, \phi_j\}$  and can therefore be controlled. In contrast, when  $\mathbf{c}_j(n)$  was randomly generated, the resulting measurements vary significantly across domains.

In the experiments, we fixed a = 0.5 and b = 0.6, and sampled  $f_j, \phi_j$  uniformly at random from the interval (0, 1]. For the random gain scenario, each  $\mathbf{c}_j(n)$  was sampled independently from a uniform distribution over



Figure 7: Comparison of the proposed methods with the JT and PT methods for the sparse gain calibration problem. The P-TDA/DD-TDA has comparable performance to PT, demonstrating the domain-adaptability aspect.

[0.1, 1.3]. In both the structured and random gain cases, we considered three domains (J = 3) and generated N = 43,000 examples per domain, following the procedure in Section 4.5.

The training procedures for JT, PT, P-TDA, and DD-TDA were identical to those described in Section 4. For the P-TDA and DD-TDA approaches, the network architectures of  $g_{\varphi}()$  and  $h_{\varphi}()$  are shown in Fig. 6.

Using the setup, we evaluated the methods in terms of NMSE and HR. The results are shown in Figs. 7(a) and 7(b). The following observations can be made. For both gain models, the JT approach results in NMSEs that are 5–10 dB higher than those achieved by the corresponding PT approach. In the PT method, training and testing on different domains yields substantially higher errors (more than 20 dB NMSE) compared to NA-LISTA (cf. Figs. 3(a) and 3(b)). These cross-domain errors are particularly severe when measurement changes are significant across domains due to variations in the domain parameter.

To quantify domain similarity, we compute the cosine similarities among the vectors  $\{\mathbf{c}_j\}_{j=1}^3$ . Let  $\rho_{i,j}$  denote the similarity between  $\mathbf{c}_i$  and  $\mathbf{c}_j$ . For the structured and random gain experiments, the similarity values are  $\{\rho_{1,2}, \rho_{2,3}, \rho_{3,1}\} = \{0.73, 0.76, 0.99\}$  and  $\{0.59, 0.78, 0.82\}$ , respectively. Due to the high similarity between domains 1 and 3 in the structured gain case  $(\rho_{1,3} = 0.99)$ , the NMSE degradation when testing PT-D1 on D3 and vice versa is limited to around 5 dB. In contrast, for less similar pairs —for example, testing PT-D2 on D3— the NMSE degradation exceeds 20 dB.

Comparing the proposed P-TDA and DD-TDA methods to PT, we find that the NMSE and HR values remain consistent across domains. These results demonstrate that a single network can adapt to multiple domains, eliminating the need to train separate networks for each one.

Next, we examine the generalization capability of our proposed methods.

**Generalization:** To assess the domain generalization capability of the proposed approaches for the calibration problem, we trained the models on 45 different domains (D1-D45) and then tested them on five unseen domains (D46-D50). For the JT, P-TDA, and DD-TDA methods, we generated N = 10,000 examples from each domain, which amounts to a total of 500,000 examples for training and testing. Whereas, for the PT models, we generated 50000 examples per domain. The datasets were divided into test, validation, and train sets in proportions of 10%, 20%, and 70%, respectively. In all these cases, the gains were generated using equation 21.



Figure 8: Generalization performance of models, across domains with structured calibration gains. The JT, P-TDA, and DD-TDA models were trained on 45 domains and tested on 5 unseen domains.

The NMSEs and HRs, for the JT and PT methods, together with P-TDA and DD-TDA, are shown in Fig. 8. We note that for all the unseen domains, the NMSEs in the PT approaches are 8-15 dB lower than the JT methods. Both the P-TDA and DD-TDA methods result in 10-13 dB lower NMSEs compared to the JT methods. Comparing the P-TDA and the PT methods, we note that the NMSEs of the P-TDA approach are within a  $\pm 3$  dB range of the corresponding PT error for all the unseen domains. Similarly, the HRs of P-TDAs are within  $\pm 4\%$  that of the PT. This shows the generalization ability of the proposed networks.

# 6 Conclusion

In this paper, we propose a systematic framework for interpretable, unrolled architecture-based domainadaptive sparse signal recovery with two tunability strategies: P-TDA, which uses known domain parameters, and DD-TDA, which derives domain-relevant variations from the data itself. We concentrate on two realistic and different types of domain variability: additive noise variance and sensor calibration gains across domains.

In the case of noise-varying domains, our approaches enable the sparse recovery network to adjust its thresholding behavior without requiring domain-specific retraining. We show that P-TDA can successfully normalize the input using known gains in the more difficult blind calibration scenario. When measurements are affected by unknown gain vectors, DD-TDA learns to correct for gain distortions without explicit access to the calibration parameters. The suggested approaches consistently exceed normal JT models and closely resemble or match the performance of domain-specific PT models. In the future, unsupervised domain learning will be used to extend this framework to incorporate more complex models.

# References

- Aviad Aberdam, Alona Golts, and Michael Elad. Ada-lista: Learned solvers adaptive to varying models. IEEE Trans. Pattern Anal. Mach. Intell., 44(12):9222–9235, 2021.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm with application to waveletbased image deblurring. In Proc. IEEE Int. Conf. Acoust., Speech, and Signal Process. (ICASSP), pp. 693–696, 2009. doi: 10.1109/ICASSP.2009.4959678.
- Wei Chen, Yan Li, and Yang Zhou. Self-supervised deep domain-adversarial regression adaptation for remaining useful life prediction. In *IEEE Trans. Ind. Electron.*, pp. 9769904, 2022.
- Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds. *Adv. Neural Info. Process. Syst.*, 31, 2018.

- Xiaohan Chen, Jialin Liu, Zhangyang Wang, and Wotao Yin. Hyperparameter tuning is all you need for LISTA. Adv. Neural Info. Process. Syst., 34:11678–11689, 2021.
- Ingrid Daubechies, Michel Defrise, and Christine De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Commun. Pure Applied Math.*, 57(11):1413–1457, 2004.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proc. North American Chap. Asso. Comput. Linguistics (NAACL-HLT), pp. 4171–4186, 2019.
- David L Donoho and Iain M Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3): 425–455, 1994.
- D.L. Donoho. Compressed sensing. IEEE Trans. Inf. Theory, 52(4):1289–1306, 2006. doi: 10.1109/TIT. 2006.871582.
- Abolfazl Farahani, Sahar Voghoei, Khaled Rasheed, and Hamid R. Arabnia. A brief review of domain adaptation. In Proc. Adv. Data Sci. Info. Engg., pp. 877–894, Cham, 2021. Springer International Publishing. ISBN 978-3-030-71704-9.
- Nikolas P Galatsanos and Aggelos K Katsaggelos. Methods for choosing the regularization parameter and estimating the noise variance in image restoration and their relation. *IEEE Trans. Image Process.*, 1(3): 322–336, 1992.
- Gene H Golub, Michael Heath, and Grace Wahba. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
- Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In Proc. Int. Conf. Machine Learn. (ICML), pp. 399–406, 2010.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proc. IEEE Int. Conf. Comput. Vision and Pattern Recognition (CVPR), pp. 770–778, 2016.
- Yifei He, Haoxiang Wang, Bo Li, and Han Zhao. Gradual domain adaptation: Theory and algorithms. J Machine Learning Res., 25(361):1–40, 2024.
- Hang Li, Yuanjia Wang, and Xiaohui Xie. Transfer learning for high-dimensional linear regression: Prediction via information borrowing. J. Royal Statistical Soc.: Series B, 84(1):149–175, 2023.
- Jialin Liu and Xiaohan Chen. ALISTA: Analytic weights are as good as learned weights in LISTA. In Int. Conf. Learning Representations (ICLR), 2019.
- Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In Proc. Int. Conf. Computer Vision (ICCV), pp. 2200–2207, 2013.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep adaptation networks. In Proc. Int. Conf. on Machine Learning (ICML), pp. 97–105, 2015.
- Yan Lu, Jing Qin, and Yuanjia Wang. Deep domain adaptation for regression. In Proc. Int. Conf. Machine Learning (ICML), pp. 97–105, 2019.
- Sven Mersmann, Alexander Seitel, Michael Erz, Bernd Jähne, Felix Nickel, Markus Mieth, Arianeb Mehrabi, and Lena Maier-Hein. Calibration of time-of-flight cameras for accurate intraoperative surface reconstruction. *Medical Physics*, 40(8):082701, 2013.
- Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Mag. Signal Processing*, 38(2):18–44, 2021.
- Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In Proc. Int. Conf. Computer Vision (ICCV), pp. 5715–5725, 2017.

- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.*, 22(10): 1345–1359, 2010.
- David Pardoe and Peter Stone. Boosting for regression transfer. In Proc. Int. Conf. on Machine Learning (ICML), pp. 863–870, 2010.
- Arogyaswami Paulraj and Thomas Kailath. Direction of arrival estimation by eigenstructure methods with unknown sensor gain and phase. In *Proc. IEEE Int. Conf. Acoust., Speech, and Signal Process. (ICASSP)*, volume 10, pp. 640–643, 1985.
- Christophe Schulke, Francesco Caltagirone, Florent Krzakala, and Lenka Zdeborová. Blind calibration in compressed sensing using message passing algorithms. Adv. Neural Info. Process. Syst., 26, 2013.
- Olivier Steiger, Judith Felder, and Stephan Weiss. Calibration of time-of-flight range imaging cameras. In Proc. Int. Conf. on Image Process., pp. 1968–1971. IEEE, 2008.
- Bahareh Tolooshams, Satish Mulleti, Demba Ba, and Yonina C Eldar. Unrolled compressed blinddeconvolution. *IEEE Trans. Signal Process.*, 71:2118–2129, 2023.
- Behrang Tolooshams, Xue Wang, Xinyu He, Yifan Zhang, and Mathews Jacob. Transfer learning for linear regression: A statistical test of gain. In arXiv preprint arXiv:2102.09504, 2021.
- Huan-Hsin Tseng, Hsin-Yi Lin, Kuo-Hsuan Hung, and Yu Tsao. Interpretations of domain adaptations via layer variational analysis. arXiv preprint arXiv:2302.01798, 2023.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In Proc. Adv. Neural Info. Process. Sys. (NeurIPS), pp. 3320–3328, 2014.
- Yu Zhang, Xinyu Wang, and Xinyu He. Representation transfer learning for semiparametric regression. arXiv preprint arXiv:2406.13197, 2024.