

# Gradient-Based Learning of Discrete Structured Measurement Operators for Signal Recovery

Jonathan Sauder  
EPFL  
jonathan.sauder@epfl.ch

Martin Genzel  
Utrecht University  
m.genzel@uu.nl

Peter Jung  
TU Berlin  
peter.jung@tu-berlin.de

## Abstract

Countless signal processing applications include the reconstruction of signals from few indirect linear measurements. The design of effective measurement operators is typically constrained by the underlying hardware and physics, posing a challenging and often even discrete optimization task. While the potential of gradient-based learning via the unrolling of iterative recovery algorithms has been demonstrated, it has remained unclear how to leverage this technique when the set of admissible measurement operators is structured and discrete. We tackle this problem by combining unrolled optimization with Gumbel reparametrizations, which enable the computation of low-variance gradient estimates of categorical random variables. Our approach is formalized by GLODISMO (Gradient-based Learning of DIScrete Structured Measurement Operators). This novel method is easy-to-implement, computationally efficient, and extendable due to its compatibility with automatic differentiation. We empirically demonstrate the performance and flexibility of GLODISMO in several prototypical signal recovery applications, verifying that the learned measurement matrices outperform conventional designs based on randomization as well as discrete optimization baselines.

## I. INTRODUCTION

Linear measurement operators adhering a specific structure due to physical constraints of the observation process or the hardware of the measuring device are ubiquitous in signal processing. Corresponding structural constraints appear in numerous relevant applications, like magnetic resonance imaging with masked Fourier matrices [Lustig et al., 2007], communication and remote sensing tasks [Zhu and Bamler, 2010], single-pixel imaging [Duarte et al., 2008], compressed sensing with expander-graphs [Foucart and Rauhut, 2013], or pooling matrices for group testing [Petersen et al., 2020]. The optimal design of such measurement operators—which typically reside in a discrete subset—to improve the performance of downstream tasks poses great computational challenges. While it is often easy to create a suitable random mask, it is not obvious how to optimize the measurement matrix in a way that is both efficient and respects the feasibility constraints. Classical approaches commonly use discrete optimization to find such sets, as no gradients can be directly computed.

On the other hand, gradient-based optimization via backpropagation through massive nonlinear computational graphs has shown impressive results in the field of machine learning. In order to enable these techniques even in non-differentiable settings, recent work has developed a variety of approaches for computing low-variance estimates of gradients. In this context, a promising concept is given by Gumbel reparametrizations [Jang et al., 2016, Maddison et al., 2016]. These allow for estimating the gradients of categorical random variables by simply adding component-wise independent and identically distributed (i.i.d.) noise from a Gumbel distribution to the logits. This easy-to-implement method has proven successful in many machine learning applications, including graph learning [Krawczuk et al., 2021], discrete representation learning [Kusner and Hernández-Lobato, 2016], and neural architecture search [Xie et al., 2018]. However, Gumbel reparametrizations have found little attention in the signal processing and inverse problems communities so far, except for magnetic resonance imaging [Huijben et al., 2019, 2020, Bahadir et al., 2019, 2020]. The present work aims at using Gumbel reparametrizations to fuse gradient-based learning with the design of structured measurement operators that are constrained to a discrete set.

The actual signal recovery problem is typically solved by convex optimization methods. It is well-established that the computational graph of common iterative optimization schemes can be unrolled to obtain a neural network that can be readily backpropagated through [Gregor and LeCun, 2010]. This strategy allows the solver to adapt to

(training) data and has proven very powerful, especially in reducing the number of iterations for convex programs by orders of magnitudes [Gregor and LeCun, 2010, Liu and Chen, 2019, Behrens et al., 2020]. Importantly, unrolling also enables the computation of gradients with respect to the measurement operator and other parameters of the reconstruction algorithm. However, it has remained unclear how this technique can be leveraged when the underlying parameter set is structured and discrete. This important shortcoming has made unrolled optimization infeasible for the design of practicable measurement operators. In this paper, we present a novel approach to tackle this problem. Our main contributions can be summarized as follows:

- We propose an efficient and easy-to-implement method for Gradient-based Learning of DIscrete Structured Measurement Operators (GLODISMO), which combines the concepts of unrolled optimization and Gumbel reparametrizations.
- We successfully apply GLODISMO to several prototypical signal recovery tasks, namely single-pixel imaging, compressed sensing with left- $d$ -regular graphs, and pooling matrices for group testing. In particular, conventional baselines relying on randomization and discrete optimization are significantly outperformed.
- GLODISMO provides a flexible learning framework, rather than a rigid algorithm. For example, we demonstrate that our basic approach (Algorithm 1) can be easily extended to produce measurement operators that enable fast transforms. Moreover, the compatibility with automatic differentiation unlocks many potential applications beyond signal recovery.

## II. BACKGROUND & RELATED WORK

### A. Linear Inverse Problems & Compressed Sensing

In linear inverse problems, the basic task is to recover an unknown target signal  $x \in \mathbb{R}^n$  from indirect observations of the form  $y = \Phi x$  (mostly also contaminated by additive noise), where  $\Phi \in \mathbb{R}^{m \times n}$  is a known measurement matrix. The number of measurements  $m$  is usually much smaller than the ambient dimension of the signal  $n$ , making the inverse problem ill-posed and only solvable under prior knowledge about the underlying signals.

A prominent class of such inverse problems is given by compressed sensing, in which the signal is assumed to be sparse. The seminal papers in this field [Candes et al., 2006, Donoho, 2006] have shown that robust and stable recovery can be achieved with computationally tractable algorithms if the measurement matrix fulfills appropriate conditions. For example, when  $x$  is  $s$ -sparse, i.e.,  $x \in \Sigma_s^n := \{x \in \mathbb{R}^n \mid \|x\|_0 \leq s\}$ , unique reconstruction is possible via convex optimization, given that  $\Phi$  fulfills the null space property [Foucart and Rauhut, 2013]. The convex program that is to be solved corresponds to a LASSO problem [Tibshirani, 1996] with hyperparameter  $\lambda$ :

$$\min_{\hat{x}} \|\Phi \hat{x} - y\|_2^2 + \lambda \|\hat{x}\|_1. \quad (1)$$

When the signal  $x$  is not sparse by itself but with respect to a basis transform  $\Psi \in \mathbb{R}^{n \times n}$ , the reparametrization  $\bar{x} = \Psi^* x$  can be incorporated into the reconstruction algorithm, so that the problem (1) yields the synthesis formulation:

$$\min_{\bar{x}} \|\Phi \Psi^* \bar{x} - y\|_2^2 + \lambda \|\bar{x}\|_1.$$

For natural images, useful sparsifying transforms include wavelets [Haar, 1911, Rao, 2002, Daubechies, 1988], shearlets [Guo et al., 2006], discrete cosine transforms [Ahmed et al., 1974], and total variation synthesis [Candes and Guo, 2002].

As mentioned above, the success of compressed sensing particularly relies on desirable properties of the measurement matrix  $\Phi$  (or  $\Phi \Psi^*$ ), such as the null space property. On the theoretical side, it has been shown that various random matrices with  $m \in O(s \log(n/s))$  satisfy such criteria with high probability. Prototypical examples are Gaussian matrices with  $\Phi_{ij} \sim \mathcal{N}(0, 1/m)$  or Bernoulli matrices whose entries take values in  $\{1/\sqrt{m}, -1/\sqrt{m}\}$  with equal probability [Foucart and Rauhut, 2013]. Guarantees are also available for randomly subsampled Fourier matrices [Rudelson and Vershynin, 2008], which is accompanied by highly efficient matrix-vector multiplications due to the fast Fourier transform. While the randomization of measurement operators is key to the theoretical analysis of sparse recovery problems, this technique is neither necessary nor practicable in most real-world applications.

A widely-used class of algorithms for solving LASSO-type problems are gradient-based methods. For example, an iterative proximal scheme of the following type can be used to solve (1):

$$\hat{x}^{(t+1)} = \text{prox}_{\lambda \|\cdot\|_1} \left( \hat{x}^{(t)} + \gamma \nabla_{\hat{x}^{(t)}} (\|y - \Phi \hat{x}^{(t)}\|_2^2) \right), \quad (2)$$

where  $t = 0, 1, \dots, T-1$  and the proximal operator of the  $\lambda$ -weighted  $\ell_1$ -norm corresponds to element-wise soft thresholding:

$$\text{prox}_{\lambda \|\cdot\|_1}(v) := \text{sign}(v) \cdot \max\{|v| - \lambda, 0\}, \quad v \in \mathbb{R}. \quad (3)$$

The algorithm described in (2) is known as the Iterative Soft Thresholding Algorithm (ISTA) [Daubechies et al., 2004]. Another popular method is based on Iterative Hard Thresholding (IHT) [Blumensath and Davies, 2009], in which the proximal operator in (3) is replaced by a projection onto the set of  $s$ -sparse vectors  $\Sigma_s^n$ . This corresponds to a hard-thresholding operator, which sets all entries of a vector to zero except for those with the  $s$  largest absolute values.

### B. Unrolled Optimization in Linear Inverse Problems

It is well-established that many iterative optimization schemes can be viewed as a neural network via unrolling the computational graph of a finite number of iterations  $T$ . Once unrolled, the tunable parameters  $\theta$  of a given reconstruction algorithm  $f_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^n$  can be fit to a (training) dataset by minimizing a loss function  $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  via (stochastic) gradient-based optimization:

$$\min_{\theta} \mathbb{E}_x \left[ \mathcal{L}(f_\theta(\Phi x), x) \right]. \quad (4)$$

In the context of linear inverse problems, a prominent example of unrolled optimization is Learned ISTA (LISTA) [Gregor and LeCun, 2010], where ISTA is unrolled and the involved matrices, step-sizes, and thresholds are learned in an end-to-end fashion via (4). Such a data-driven approach can reduce the number of iterations required for convergence by orders of magnitude. LISTA has inspired a line of research on unrolled optimization algorithms for compressed sensing. In Analytic LISTA (ALISTA), for example, the involved matrices are computed analytically and remain fixed, while only the step-sizes and thresholds are learned [Liu and Chen, 2019]. This significantly reduces the number of trainable parameters, and theoretical guarantees known for ISTA are retained, but the performance of a fully learned LISTA is still matched. Going one step further, Neurally Augmented ALISTA (NA-ALISTA) equips ALISTA with a recurrent neural network that enables adaptive step-sizes and thresholds during reconstruction, depending on the target vector [Behrens et al., 2020].

While compressed sensing theory focuses on recovery guarantees for generic sparse vectors, real-world signals often carry additional structure, for which random measurement matrices may not be optimal. This shortcoming has inspired the search for data-driven optimization of measurement operators. Unrolled optimization provides a natural approach in this regard, since  $\Phi$  can be directly included in the above end-to-end training procedure, i.e., the objective in (4) is also minimized over  $\Phi$ . This is equivalent to learning an autoencoder neural network in which the encoder (a single linear layer) becomes the measurement matrix. In the context of compressed sensing, this has been explored by Wu et al. [2019], where a measurement matrix (initialized with Gaussian entries) is learned by backpropagating through a fixed number of steps of subgradient descent. This scheme has been also investigated in a compressive learning setting by Adler et al. [2016], where a neural network classifier is learned with a linear first layer. Both works have demonstrated that end-to-end learning can provide decent measurement matrices when there are no additional constraints on the involved operator. Similarly, LISTA and its variants use end-to-end learning on unconstrained tuning parameter sets, initialized via Gaussian random variables. However, in real-world applications, the measurement matrix often must follow specific structures imposed by the underlying hardware and physics of the problem. In magnetic resonance imaging, for example, the measurement operator is a subsampled Fourier matrix [Lustig et al., 2007], whereas in single-pixel imaging [Duarte et al., 2008], only binary matrices are admissible. In conclusion, simply including the measurement matrix into end-to-end schemes like (4) does not respect important constraints and therefore limits its practicability.

### C. Gumbel Reparametrizations

Backpropagating gradients through immensely large computational graphs has shown to be highly effective in the training of deep neural networks with hundreds of layers [Tai et al., 2017] and billions of learnable parameters [Brown et al., 2020]. However, when discrete nodes are included in the computational graphs, gradients cannot be computed directly and have to be estimated. Formally, we consider a computational graph including a discrete random variable  $v$  taking values in  $\{1, \dots, a\}$  (in one-hot encoding) for some  $a \in \mathbb{N}$ , where its unnormalized log-probabilities are denoted by  $\varphi = [\varphi_1, \dots, \varphi_a]^T \in \mathbb{R}^a$ . The value of  $v$  is then passed through a deterministic, differentiable function  $f$ . The goal is to compute gradients with respect to the probability parameters  $\varphi$ . As  $v$  is discrete, it is not possible to directly backpropagate through  $v$ . However, the gradient with respect to  $\varphi$  can be estimated via sampling from  $v$ :

$$\nabla_{\varphi} \mathbb{E}_v [f(v)].$$

A simple yet effective way to sample from such a categorical random variable is to use the Gumbel-Max trick [Gumbel, 1954, Maddison et al., 2014], according to which a realization of  $v$  can be obtained by

$$v_i = \text{one-hot} \left( \arg \max_{j \in \{1, \dots, a\}} (g_j + \varphi_j) \right)_i, \quad i = 1, \dots, a, \quad (5)$$

where  $g_j \sim \text{Gumbel}(0, 1)$  follows a Gumbel distribution, which is equivalent to a twice negative-log transformed uniform distribution, i.e.,  $g_j = -\log(-\log(u_j))$  with  $u_j \sim U(0, 1)$ .

The Gumbel-softmax trick [Jang et al., 2016], also known as the Concrete distribution [Maddison et al., 2016], replaces the argmax operator in (5) by a softmax operator (with temperature  $\tau$ ), which ensures differentiability everywhere:

$$v_i = \frac{\exp((g_i + \varphi_i)/\tau)}{\sum_{j=1}^a \exp((g_j + \varphi_j)/\tau)}, \quad i = 1, \dots, a. \quad (6)$$

Note that in the limit  $\tau \rightarrow 0$ , the Gumbel-softmax trick based on (6) reduces to the non-differentiable argmax in (5). This modification allows for the backpropagation through discrete random variables by adding element-wise i.i.d. noise, outperforming previous approaches to estimating the gradient of discrete nodes, such as the straight-through estimator [Bengio et al., 2013], or score-function-based estimators [Williams, 1992, Mnih and Gregor, 2014] in a series of supervised learning tasks. For  $\tau > 0$ , the Gumbel-softmax is a continuous relaxation of the (one-hot encoded) discrete random variable. In cases where true discreteness is needed, it is possible to use the argmax operator in the forward pass and the softmax in the backward pass of backpropagation. This strategy is known as the straight-through Gumbel softmax estimator. It can be directly extended to taking multiple samples without replacement from a categorical distribution, namely by selecting the top- $K$  values instead of only the largest one. This was first highlighted in a blog post by Vieira [2014], connecting Gumbel reparametrizations to weighted reservoir sampling.

Gumbel reparametrizations have found various applications in the field of machine learning, ranging from learned discrete representations [Jang et al., 2016, Kusner and Hernández-Lobato, 2016], over graph generation [Krawczuk et al., 2021], to neural architecture search [Xie et al., 2018]. Recently, they have been also applied in the context of compressive magnetic resonance imaging [Huijben et al., 2019, 2020], where a single Gumbel top- $K$  operation is employed to select a fixed number of rows from a Fourier matrix; see also [Bahadir et al., 2019, 2020]. This can be seen as a special case of the GLODISMO framework (see Section III), when learning only a single global mask. Having said that, our actual main concern is to allow for additional structural constraints, which is not addressable with existing methods. Moreover, the aforementioned works rely on off-the-shelf feedforward neural networks for recovery, which are impractical for most applications (where convex optimization algorithms are ubiquitous). We close this important feasibility gap by focusing on unrolled algorithms instead.<sup>1</sup>

<sup>1</sup>Note that this can be viewed as a recurrent neural network in (4), since the matrix  $\Phi$  appears in every iteration (layer) of the unrolled algorithm, cf. (2).

### III. METHOD

In this section, we present our method GLODISMO (Gradient-based Learning of DIcrete Structured Measurement Operators). The main objective is to combine the unrolling approach of (4) with learning a measurement operator  $\Phi$ :

$$\min_{\theta, \Phi} \mathbb{E}_x \left[ \mathcal{L}(f_\theta(\Phi x), x) \right]. \quad (7)$$

Our key idea is to impose discrete structural constraints on  $\Phi$  by Gumbel reparametrizations to enable gradient estimation.

More specifically, we model  $\Phi \in \mathbb{R}^{m \times n}$  as a matrix-valued discrete variable over a certain probability distribution that is parametrized by a learnable parameter  $\varphi \in \mathbb{R}^{m \times n}$ . To this end, we consider the index set  $\mathcal{I} := \{1, \dots, m\} \times \{1, \dots, n\}$  and let  $\mathcal{P}(\mathcal{I}) = \{I_1, \dots, I_l\}$  be a partition<sup>2</sup> of  $\mathcal{I}$ , whose purpose is to capture the structural constraints on  $\Phi$ . As such,  $\Phi$  follows a joint distribution over random variables supported on the index subsets  $I_1, \dots, I_l$ , each of which is independently obtained from a Gumbel reparametrization. At this point the parameter  $\varphi$  comes into play: The notation  $\varphi[I_i]$  refers to those entries of  $\varphi$  indexed by  $I_i$ . To obtain  $\Phi$ , for each  $I_i$ , we add element-wise i.i.d. Gumbel noise to  $\varphi[I_i]$  and then select the indices corresponding to the top- $d_i$  values. This procedure is used in the forward pass to sample  $\Phi$  and to compute the loss in (7). In the backward pass, we use the gradient of the softmax with temperature  $\tau$  instead of the hard top- $d_i$ . We refer to the above procedure as GLODISMO. A pseudo-code implementation is provided in Algorithm 1, assuming a software framework capable of automatic differentiation. For the sake of simplicity, the described method is limited to learning binary matrices, but an extension to larger alphabets is straightforward.

As indicated above, the index partition  $\mathcal{P}(\mathcal{I})$  is used to impose structural constraints on the measurement matrix  $\Phi$ . While a proper choice of  $\mathcal{P}(\mathcal{I})$  is problem-specific, it is often naturally given by the rows or columns of  $\Phi$ , and the number of selected elements  $d_i$  is equal for each  $i \in \{1, \dots, l\}$  (see Section IV for examples). In such cases, vectorized, and therefore computationally efficient softmax and top- $K$  operations over one of the axes can be applied. In Algorithm 1,  $\varphi \in \mathbb{R}^{m \times n}$  has  $m \cdot n$  learnable parameters, which must be stored in memory during training time. This is identical to the cost of learning a dense sensing matrix without constraints. At test time, a single mask is sampled after adding Gumbel noise and performing the top- $K$ , and then kept fixed. This means that there is no additional computational cost compared to using a conventional (randomly chosen, but fixed) matrix satisfying the constraints. Hence, our method is feasible for training and comes at no additional expenses during testing or deployment.

---

**Algorithm 1 (GLODISMO)** Learning a binary matrix with  $d_i$  ones per set  $I_i$  of the partition.

---

**Input:** signal  $x$  (training data), temperature  $\tau$ , top- $K$ -keeps  $d_1, \dots, d_l \in \mathbb{N}$ , differentiable reconstruction algorithm  $f_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^n$ , index partition  $\mathcal{P}(\mathcal{I}) = \{I_1, \dots, I_l\}$

**Learnable Parameters:** Parameters of measurement matrix  $\varphi \in \mathbb{R}^{m \times n}$ , parameters of reconstruction algorithm  $\theta$

```

1:  $G \sim_{i.i.d.} \text{Gumbel}(0, 1)^{m \times n}$ 
2: for  $i \in \{1, \dots, l\}$  do
3:    $\text{logits} := (\varphi[I_i] + G[I_i]) / \tau$ 
4:    $\text{probs} := \text{softmax}(\text{logits})$ 
5:    $\text{hard} := \text{topk}(\text{probs}, d_i)$ 
6:    $\Phi[I_i] := \text{hard.detach}() + \text{probs} - \text{probs.detach}()$ 
7: end for
8:  $y := \Phi x$ 
9:  $\text{loss} := \mathcal{L}(f_\theta(y), x)$ 
10:  $\text{loss.backward}()$ 
```

---

<sup>2</sup>That means the sets  $I_1, \dots, I_l$  are non-empty and pairwise disjoint such that  $\bigcup_{i \in \{1, \dots, l\}} I_i = \mathcal{I}$ .

Note that the binary matrix computed in Algorithm 1 can be used like any other parameter in a framework with automatic differentiation, since the gradient with respect to  $\varphi$  is well-defined. In particular, this allows for simultaneous optimization of  $\Phi$  and other tunable parameters  $\theta$ , which is an important advantage over discrete optimization techniques. More generally, GLODISMO should be seen as an extendable learning framework rather than a rigid algorithm. For example, Algorithm 1 can be easily modified to learn any measurement operator that is constructed by differentiable transforms of one or even multiple binary matrices (which are structured in the sense that their entries are subselected from specific partitions). We demonstrate this flexibility in the context of fast transforms (see Appendix C) and learning super-pixel masks for single pixel imaging (see Appendix D). A full evaluation of potential generalizations is beyond the scope of this paper and left to future research.

#### IV. EXPERIMENTS

In this section, we demonstrate the effectiveness of GLODISMO in a selection of prototypical signal recovery applications. We refer the reader to Appendix A for more details on the implementation of all considered methods.

*Baselines:* We compare our method to random matrices (fixed after a single random draw), which is a standard baseline for all considered problem setups. Moreover, we benchmark against two baseline discrete optimization algorithms: a greedy approach as well as Simulated Annealing [Pincus, 1970, Romeijn and Smith, 1994], referred to as “Greedy” and “SimAn” in our experiments, respectively. Note that both methods are incompatible with reconstruction algorithms that include any (real-valued) learnable parameters.

*Error metrics:* As common in many signal recovery tasks, we report the normalized mean squared error (NMSE) and normalized mean absolute error (NMAE) between the ground truth signal  $x$  and the estimated reconstruction  $\hat{x}$ , which are defined as (in dB-scale):

$$\text{NMSE}(x, \hat{x}) := 10 \log_{10} \left( \frac{\mathbb{E}_x[\|\hat{x} - x\|_2^2]}{\mathbb{E}_x[\|x\|_2^2]} \right),$$

$$\text{NMAE}(x, \hat{x}) := 10 \log_{10} \left( \frac{\mathbb{E}_x[\|\hat{x} - x\|_1]}{\mathbb{E}_x[\|x\|_1]} \right).$$

##### A. Application: Single Pixel Imaging

The fact that conventional visible light cameras with millions of pixels are cheaply available stems from the fortunate coincidence that the wavelengths of visible light are in the same order of magnitude as the wavelengths that silicon responds to [Mackenzie, 2009]. For imaging beyond visible light, cameras composed of many pixels may be prohibitively expensive due to the difficulty in manufacturing such pixels. A remarkable application of compressed sensing is single pixel imaging [Duarte et al., 2008], in which the original image is reconstructed from compressive measurements acquired by a spatial light modulator or a digital micromirror device that collects light onto a single pixel using a series of masks. For each of the  $m$  measurements, the pixel intensity is measured while using a distinct mask. Mathematically, this process can be expressed as observing the unknown (vectorized) image  $x \in \mathbb{R}^n$  through a binary measurement matrix  $\Phi \in \mathbb{R}^{m \times n}$ , and then using a reconstruction algorithm to estimate

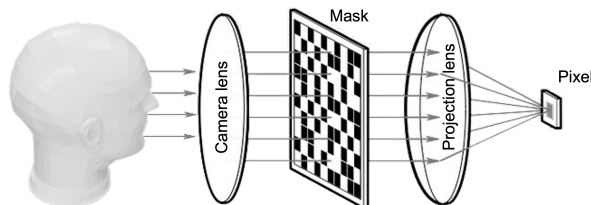


Fig. 1: Schematic visualization of a single-pixel imaging setup; image adapted from Bacca et al. [2019].

the original image from  $m$  sums of pixels. As natural images are often sparse with respect to a certain transform domain, compressed sensing can be employed for the reconstruction. An example of a single-pixel imaging setup is visualized in Figure 1. While theoretical results highlight that random masks have favorable reconstruction properties for generic sparse signals, this may not be the case when the physics and hardware constraints of the measurement process or additional image structures are taken into account. Examples are the modeling of diffraction effects or implicit signal structure that is extracted in a data-driven fashion.

GLODISMO is well suited to incorporate such aspects. For our case study, we partition the indices of  $\varphi \in \mathbb{R}^{m \times n}$  into the set of row vectors and use Algorithm 1 to learn adaptive masks for the MNIST dataset ( $n = 784$ ) [LeCun et al., 1998]. The squared error  $\mathcal{L}(\hat{x}, x) = \|\hat{x} - x\|_2^2$  is used as loss function during training, and bi-orthogonal 2.2 wavelets with one level as sparsifying transform. We unroll IHT, which has no additional learnable parameters, and NA-ALISTA, in which an LSTM-network [Hochreiter and Schmidhuber, 1997] is employed to adaptively predict step-sizes and thresholds. In a first experiment, we demonstrate that GLODISMO enables much better image reconstructions from far fewer measurements. For this, we fix the number of ones per row (i.e., the number of "on" pixels per mask) to 32. We run IHT and NA-ALISTA for  $T = 20$  iterations and compare the reconstruction performance for random pixel masks and those learned via GLODISMO. The number of measurements is varied from 10 to 500. Additive i.i.d. Gaussian noise with a signal-to-noise ratio of 40dB is added to the measurement vector  $y$ .

The results are reported in Figure 2. With only 10 measurements, IHT with a learned  $\Phi$  is able to match the performance of IHT with a random  $\Phi$  of 200 measurements. As IHT itself has no learnable parameters, this difference must stem from the learned mask. Both the greedy and Simulated Annealing baselines lead to only marginally better results than random matrices. As the hard-thresholding in IHT removes all but the top 50 absolute values of the reconstruction (in the wavelet basis), it saturates at around -10dB NMSE for learned  $\Phi$ , which is essentially perfect reconstruction of the top 50 wavelet coefficients under 40dB Gaussian noise. On the other hand, NA-ALISTA, which is based on soft-thresholding, is able to reach a significantly higher reconstruction quality. Similarly to IHT, NA-ALISTA with a learned mask significantly outperforms its counterpart with a random mask. Overall, Figure 2 highlights that both learning  $\Phi$  and learning the parameters of the reconstruction algorithm enable more accurate and faster reconstructions; accordingly, this enables a significant reduction of the number of required measurements. The effectiveness of GLODISMO is further underpinned when fixing the number of measurements and instead evaluating the improvement in convergence speed. For  $m = 50$  and  $m = 200$  measurements, we compare the convergence of the aforementioned algorithms in Figure 3. Even after a single iteration of using a learned  $\Phi$  with either IHT or NA-ALISTA, the reconstruction is better than a random  $\Phi$  with a thousand iterations of IHT or 20 iterations of NA-ALISTA. Furthermore, a learned  $\Phi$  leads to a much higher accuracy after only a few iterations when compared to IHT with a random  $\Phi$  after thousands of iterations. This also holds true for  $m = 200$ , but NA-ALISTA outperforms IHT for learned  $\Phi$ . Again, both the greedy and Simulated Annealing baseline are only insignificantly better than random matrices. These results verify that GLODISMO can be used seamlessly within the algorithm unrolling framework to speed up convergence of signal recovery.

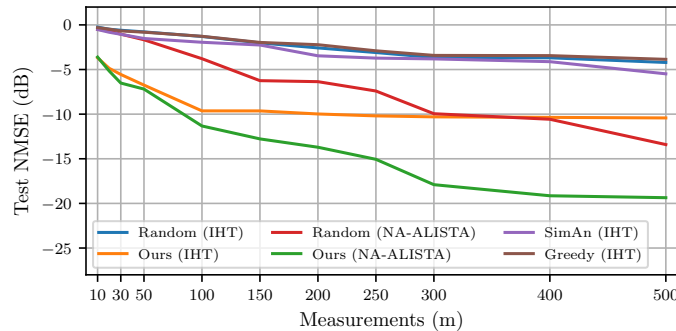


Fig. 2: NMSE of reconstruction with single pixel imaging masks as a function of the number of measurements  $m$ .

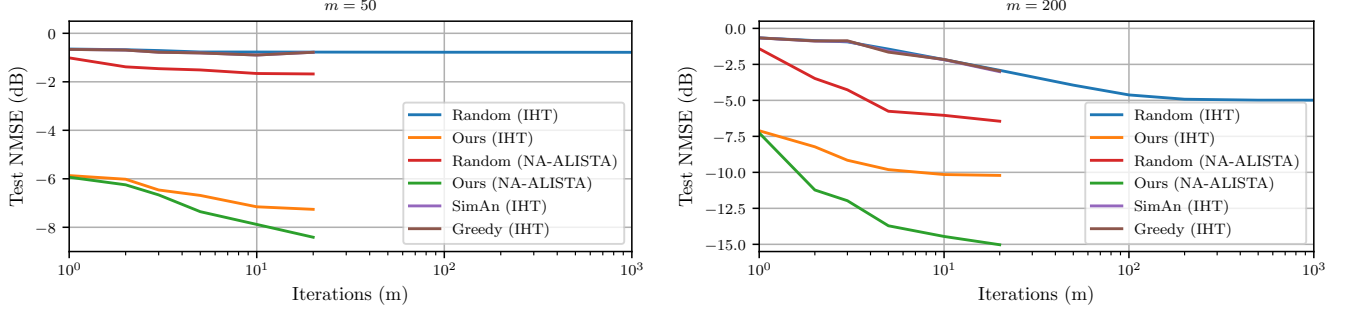


Fig. 3: NMSE of image reconstruction with single pixel imaging masks as a function of the number of iterations  $t$  for  $m = 50$  (left) and  $m = 200$  measurements (right).

Finally, to visualize how changing the mask influences the reconstruction quality beyond a quantitative metric like NMSE, Figure 4 shows a random  $\Phi$  and a learned  $\Phi$  (in combination with unrolled IHT), as well as individual reconstructions from the test set. The four images in subfigure (a) correspond to the top four rows of the matrix which is displayed in subfigure (c) (rows resized to image shape), i.e., the first four light patterns that are projected onto the scene. No structure can be seen in the matrix or the masks, and the samples reconstructed with 20 IHT-iterations in subfigure (e) are of poor quality. The right part of Figure 4 displays the learned counterparts. The four shown masks clearly exhibit additional structure, and so does the full matrix  $\Phi$ . The digits of the reconstructed samples are clearly readable.

*Further experiments:* We have conducted further experiments on the single-pixel imaging setup to highlight possible extensions of GLODISMO within the automatic differentiation framework. More specifically, we show that GLODISMO can be used to learn fast-transform matrices in Appendix C as well as super-pixel masks in Appendix D.

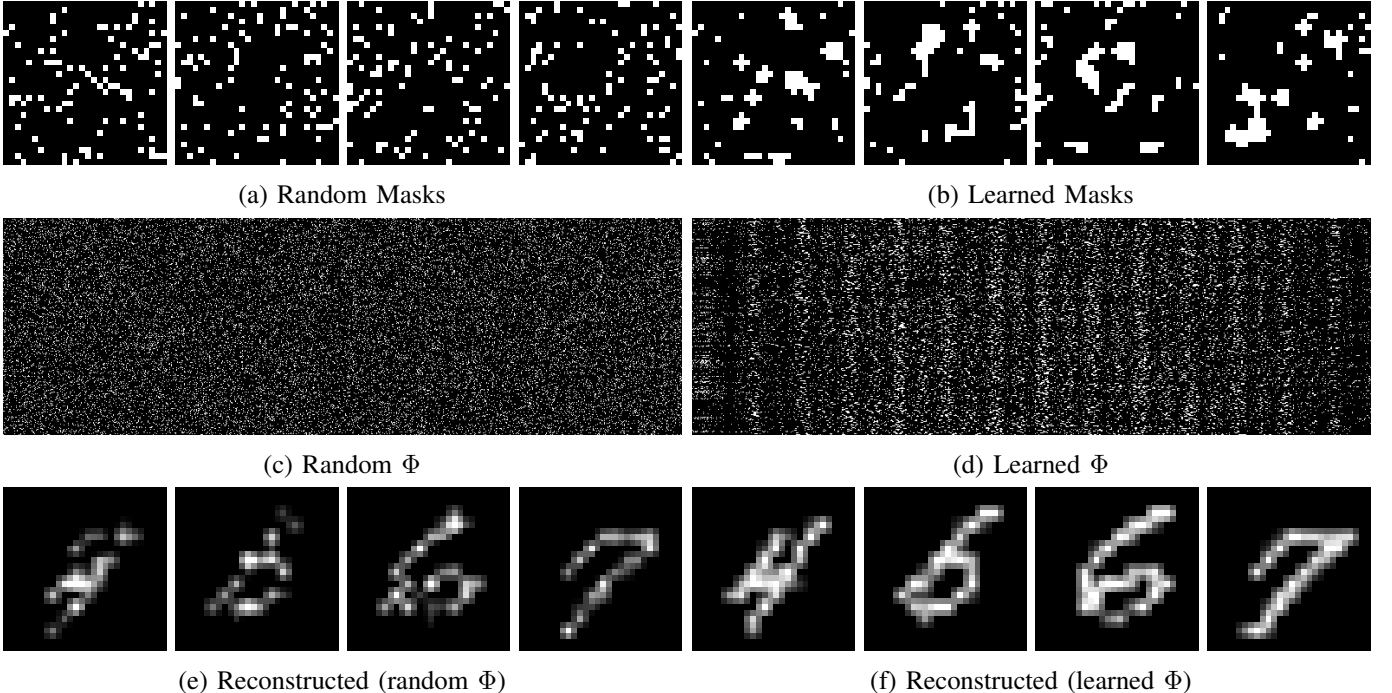


Fig. 4: Single pixel imaging masks for  $m = 250$ , where  $d = 90$  ones per row are selected (i.e., ones per mask). The reconstructions are obtained by  $T = 20$  IHT-iterations.



### B. Application: Compressed Sensing with Left- $d$ -Regular Graphs

Another important class of measurement matrices suitable for compressed sensing is given by adjacency matrices of lossless expanders [Foucart and Rauhut, 2013]. In this setting,  $\Phi$  corresponds to a 0-1 adjacency matrix of a bipartite graph, connecting the signal  $x$  (left vertices) to the measurements  $y$  (right vertices). A graph is called left- $d$ -regular if every left vertex has exactly  $d$  connected right vertices. In matrix notation, this translates into  $\Phi$  having exactly  $d$  ones per column. An example is visualized in Figure 5. Extremely large random left- $d$ -regular graphs with  $d \in O(n/s)$  and  $m \in O(s \log(n/s))$  satisfy the restricted expansion property with high probability, which is a sufficient recovery criterion [Foucart and Rauhut, 2013]. However, for smaller values of  $m$  and  $n$ , such random graphs are unlikely to enjoy this property and therefore may not be suited for compressed sensing. In this regime, we can employ GLODISMO to learn left- $d$ -regular graphs. Invoking Algorithm 1, we partition the entries of the measurement matrix into its column vectors, and select  $d = 7$  ones per column. As unrolled method, we consider Iterative Hard Thresholding for expanders (E-IHT), see Foucart and Rauhut [2013]:

$$\hat{x}^{(t+1)} = \mathcal{H}_s \left( \hat{x}^{(t)} + \mathcal{M}(y - \Phi \hat{x}^{(t)}) \right).$$

Here,  $\mathcal{H}_s : \mathbb{R}^n \rightarrow \Sigma_s^n$  denotes the hard thresholding operator, which sets all but the  $s$  largest absolute coefficients to 0. The nonlinear function  $\mathcal{M} : \mathbb{R}^m \rightarrow \mathbb{R}^m$  denotes the median operator associated with  $\Phi$ , i.e., the  $j$ -th coefficient of  $\mathcal{M}(y)$  is given by  $\text{median}\{y_i : i \in R(j)\}$ , where  $R(j)$  is the set of right vertices connected to the left vertex  $j$  (which are  $d$ -many for left- $d$ -regular graphs). Similarly to E-IHT, we adapt NA-ALISTA to expanders by replacing the adjoint  $\Phi^T$  in each iteration with the median operator  $\mathcal{M}$ , which we call E-NA-ALISTA. As common in expander theory, we use  $\mathcal{L}(x, \hat{x}) = \|x - y\|_1$  as loss function. The experiments in this section are conducted with heavy-tailed noise (student t-distributed with 1 degree of freedom) and a signal-to-noise ratio of 40dB.

We employ E-IHT and E-NA-ALISTA with  $m = 250$  measurements on synthetic sparse vectors with  $n = 784$  and an expected sparsity of  $s = 40$ . The support of the synthetic data is generated via i.i.d. Bernoulli random variables, while the non-zero coefficients are normally distributed. Figure 6 reports the NMAE on the test set as the training progresses; note that E-IHT itself has no learnable parameters. GLODISMO clearly improves the reconstruction performance for both E-IHT and E-NA-ALISTA. This is remarkable from a compressed sensing perspective, since there is no signal structure beyond generic sparsity, and yet the learned matrices empirically outperform their random counterparts.

To evaluate if GLODISMO can speed up solving inverse problems with left- $d$ -regular graphs, we report the NMAE as a function of the algorithm iterations in Figure 7. E-IHT with a random left- $d$ -regular graph converges much more slowly than in the case of learning: the former takes about 2000 iterations to reach the same NMAE level as the latter with only 50 iterations. E-NA-ALISTA generally performs better than E-IHT and can be also further accelerated by GLODISMO. Similar to single pixel imaging, we observe in both Figures 6 and 7 that the greedy and Simulated Annealing baselines only slightly improve upon random designs.

Finally, an interesting question is whether the measurement matrices learned by GLODISMO overfit to the unrolled algorithm in use. To this end, we evaluate whether a matrix learned by E-IHT also works well for E-NA-ALISTA and vice versa (in the same setup as before). The results are reported in Figures 8 and 9. We find that using a matrix previously learned by E-IHT as a fixed matrix for E-NA-ALISTA leads to a slightly better performance

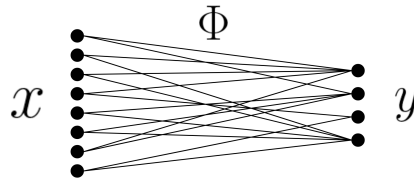


Fig. 5: Example of a linear measurement process using a left- $d$ -regular bipartite graph with  $d = 2$ , where the signal space is 8-dimensional and the measurement space 4-dimensional; image adapted from Foucart and Rauhut [2013, p. 436].

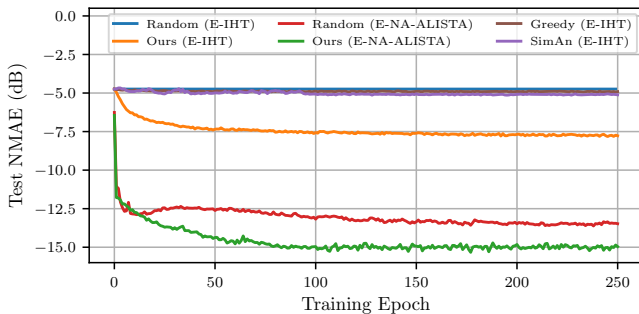


Fig. 6: NMAE of compressed sensing with left- $d$ -regular graphs on synthetic data as a function of the training progress.

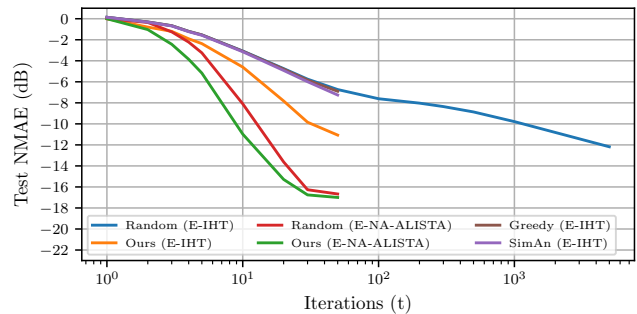


Fig. 7: NMAE of compressed sensing with left- $d$ -regular graphs on synthetic data as a function of the number of iterations  $t$ .

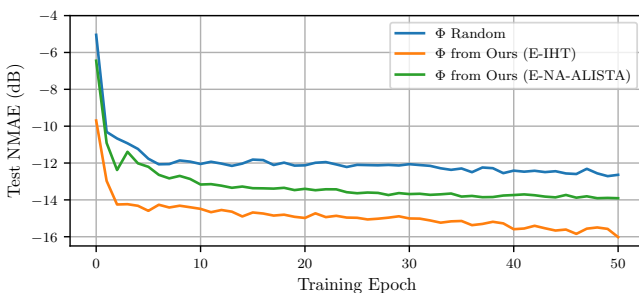


Fig. 8: Comparison of training E-NA-ALISTA with a fixed left- $d$ -regular graph, which is randomly drawn (blue), previously learned by E-IHT with  $T = 20$  iterations (orange), and previously learned by E-NA-ALISTA with  $T = 20$  iterations (green).

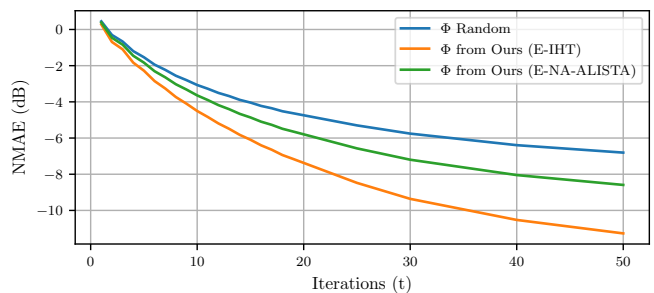


Fig. 9: Comparison of E-IHT (with varying number of iterations  $t$ ) using a fixed left- $d$ -regular graph, which is randomly drawn (blue), previously learned by E-IHT with  $T = 20$  iterations (orange), and previously learned by E-NA-ALISTA with  $T = 20$  iterations (green).

than learning it directly with E-NA-ALISTA. In the other direction, using a matrix learned by E-NA-ALISTA for E-IHT works better than a random matrix, but it is outperformed by a matrix specifically learned by E-IHT. These observations appear plausible to us, since E-NA-ALISTA learns the measurement operator and algorithm parameters simultaneously, whereas E-IHT itself has no tuning parameters.

### C. Application: Pooling Matrices for Group Testing

Another prototypical inverse problem with discrete structured measurement operators is posed by group testing, which is a technique that can be used to reduce the number of tests to identify infectious diseases. Conceptually, the measurement matrix describes which specimen is pooled into which test, while the actual recovery problem is conveniently addressable by Non-Negative Least Absolute Deviation (NNLAD) algorithms [Petersen et al., 2020]. An analysis of GLODISMO in this setup is presented in Appendix E.

## V. CONCLUSION & FUTURE WORK

In this work, we have proposed GLODISMO, an efficient and extendable method for learning discrete structured measurement operators for signal recovery, based on a fusion of unrolled optimization and Gumbel reparametrizations. The effectiveness and flexibility of our approach has been empirically demonstrated in several prototypical applications, thereby significantly outperforming discrete optimization baselines and random matrix designs. There are many more potential applications of GLODISMO, which could be explored in future work. In particular, it is evident to evaluate our methodology in a large-scale setting, with much higher ambient dimensions (both in the signal and measurement domain) and larger datasets. Furthermore, future research could involve an extension to compressive classification [Duarte et al., 2007, Davenport et al., 2007], where the measurement operator would be optimized to classify signals based on the measurements, instead of reconstructing them.

**Acknowledgements:** PJ was partially funded by the German Federal Ministry of Education and Research (BMBF) in the framework of the international future AI lab "AI4EO – Artificial Intelligence for Earth Observation: Reasoning, Uncertainties, Ethics and Beyond" (Grant number: 01DD20001).

#### REFERENCES

- A. Adler, M. Elad, and M. Zibulevsky. Compressed learning: A deep neural network approach. *arXiv preprint arXiv:1610.09615*, 2016.
- N. Ahmed, T. Natarajan, and K. R. Rao. Discrete cosine transform. *IEEE transactions on Computers*, 100(1):90–93, 1974.
- S. Augustin, T. Szollmann, P. Jung, and H.-W. Huebers. Breaking imaging limits using dithering masks in 0.35 terahertz single-pixel imaging. *arXiv preprint arXiv:1711.02995*, 2017.
- J. Bacca, C. V. Correa, E. Vargas, S. Castillo, and H. Arguello. Compressive classification from single pixel measurements via deep learning. In *IEEE 29th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6, 2019.
- C. D. Bahadir, A. V. Dalca, and M. R. Sabuncu. Adaptive compressed sensing mri with unsupervised learning. *arXiv preprint arXiv:1907.11374*, 2019.
- C. D. Bahadir, A. Q. Wang, A. V. Dalca, and M. R. Sabuncu. Deep-learning-based optimization of the under-sampling pattern in mri. *IEEE Transactions on Computational Imaging*, 6:1139–1152, 2020.
- F. Behrens, J. Sauder, and P. Jung. Neurally augmented alista. *arXiv preprint arXiv:2010.01930*, 2020.
- Y. Bengio, N. Léonard, and A. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- T. Blumensath and M. E. Davies. Iterative hard thresholding for compressed sensing. *Applied and computational harmonic analysis*, 27(3):265–274, 2009.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- E. J. Candes and F. Guo. New multiscale transforms, minimum total variation synthesis: Applications to edge-preserving image reconstruction. *Signal Processing*, 82(11):1519–1543, 2002.
- E. J. Candes, J. K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- I. Daubechies. Orthonormal bases of compactly supported wavelets. *Communications on pure and applied mathematics*, 41(7):909–996, 1988.
- I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004. doi: <https://doi.org/10.1002/cpa.20042>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpa.20042>.
- M. A. Davenport, M. F. Duarte, M. B. Wakin, J. N. Laska, D. Takhar, K. F. Kelly, and R. G. Baraniuk. The smashed filter for compressive classification and target recognition. In *Computational Imaging V*, volume 6498, page 64980H. International Society for Optics and Photonics, 2007.
- D. L. Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- M. F. Duarte, M. A. Davenport, M. B. Wakin, J. N. Laska, D. Takhar, K. F. Kelly, and R. G. Baraniuk. Multiscale random projections for compressive classification. In *2007 IEEE International Conference on Image Processing*, volume 6, pages VI–161. IEEE, 2007.
- M. F. Duarte, M. A. Davenport, D. Takhar, J. N. Laska, T. Sun, K. F. Kelly, and R. G. Baraniuk. Single-pixel imaging via compressive sampling. *IEEE signal processing magazine*, 25(2):83–91, 2008.
- S. Foucart and H. Rauhut. An invitation to compressive sensing. In *A mathematical introduction to compressive sensing*, pages 1–39. Springer, 2013.
- K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In *Proceedings of the 27th international conference on international conference on machine learning*, pages 399–406, 2010.
- E. J. Gumbel. Statistical theory of extreme values and some practical applications. *NBS Applied Mathematics Series*, 33, 1954.
- K. Guo, G. Kutyniok, and D. Labate. Sparse multidimensional representations using anisotropic dilation and shear operators, 2006.

- V. Guruswami, C. Umans, and S. Vadhan. Unbalanced expanders and randomness extractors from parvaresh–vardy codes. *Journal of the ACM (JACM)*, 56(4):1–34, 2009.
- A. Haar. Zur theorie der orthogonalen funktionensysteme. *Mathematische Annalen*, 71(1):38–53, 1911.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- I. A. Huijben, B. S. Veeling, and R. J. van Sloun. Deep probabilistic subsampling for task-adaptive compressed sensing. In *International Conference on Learning Representations*, 2019.
- I. A. Huijben, B. S. Veeling, and R. J. van Sloun. Learning sampling and model-based signal recovery for compressed sensing mri. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8906–8910. IEEE, 2020.
- E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- I. Krawczuk, P. Abranches, A. Loukas, and V. Cevher. GG-GAN: A geometric graph generative adversarial network, 2021. URL <https://openreview.net/forum?id=qiAxL3Xqx1o>.
- M. J. Kusner and J. M. Hernández-Lobato. Gans for sequences of discrete elements with the gumbel-softmax distribution. *arXiv preprint arXiv:1611.04051*, 2016.
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- J. Liu and X. Chen. Alista: Analytic weights are as good as learned weights in lista. In *International Conference on Learning Representations (ICLR)*, 2019.
- M. Lustig, D. Donoho, and J. M. Pauly. Sparse mri: The application of compressed sensing for rapid mr imaging. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 58(6):1182–1195, 2007.
- D. Mackenzie. Compressed sensing makes every pixel count. *What’s happening in the mathematical sciences*, 7: 114–127, 2009.
- C. J. Maddison, D. Tarlow, and T. Minka. A\* sampling. *arXiv preprint arXiv:1411.0030*, 2014.
- C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- A. Mnih and K. Gregor. Neural variational inference and learning in belief networks. In *International Conference on Machine Learning*, pages 1791–1799. PMLR, 2014.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019.
- H. B. Petersen, B. Bah, and P. Jung. Practical high-throughput, non-adaptive and noise-robust sars-cov-2 testing. *arXiv preprint arXiv:2007.09171*, 2020.
- H. B. Petersen, B. Bah, and P. Jung. Efficient tuning-free l1-regression of nonnegative compressible signals. *Frontiers in Applied Mathematics and Statistics*, 7, 2021. ISSN 2297-4687. doi: 10.3389/fams.2021.615573. URL <https://www.frontiersin.org/article/10.3389/fams.2021.615573>.
- M. Pincus. A monte carlo method for the approximate solution of certain types of constrained optimization problems. *Operations research*, 18(6):1225–1228, 1970.
- R. Rao. Wavelet transforms. *Encyclopedia of imaging science and technology*, 2002.
- H. Rauhut, J. Romberg, and J. A. Tropp. Restricted isometries for partial random circulant matrices. *Applied and Computational Harmonic Analysis*, 32(2):242–254, 2012.
- H. E. Romeijn and R. L. Smith. Simulated annealing for constrained global optimization. *Journal of Global Optimization*, 5(2):101–126, 1994.
- M. Rudelson and R. Vershynin. On sparse reconstruction from fourier and gaussian measurements. *Comm. Pure Appl. Math.*, 61(8):1025–1045, 2008.
- M.-J. Sun, M. P. Edgar, D. B. Phillips, G. M. Gibson, and M. J. Padgett. Improving the signal-to-noise ratio of single-pixel imaging using digital microscanning. *Optics express*, 24(10):10476–10485, 2016.
- Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3147–3155, 2017.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B*

- (*Methodological*), 58(1):267–288, 1996.
- T. Vieira. Gumbel-max trick and weighted reservoir sampling, 2014. URL <http://timvieira.github.io/blog/post/2014/08/01/gumbel-max-trick-and-weighted-reservoir-sampling/>.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- S. Wu, A. Dimakis, S. Sanghavi, F. Yu, D. Holtmann-Rice, D. Storchus, A. Rostamizadeh, and S. Kumar. Learning a compressed sensing measurement matrix via gradient unrolling. In *International Conference on Machine Learning*, pages 6828–6839. PMLR, 2019.
- S. Xie, H. Zheng, C. Liu, and L. Lin. Snas: stochastic neural architecture search. *arXiv preprint arXiv:1812.09926*, 2018.
- X. X. Zhu and R. Bamler. Tomographic sar inversion by  $l_1$ -norm regularization—the compressive sensing approach. *IEEE transactions on Geoscience and Remote Sensing*, 48(10):3839–3846, 2010.

## APPENDIX A IMPLEMENTATION DETAILS

In our experimental setup,<sup>3</sup> we optimize  $\Phi$  and  $\theta$  using the Adam optimizer [Kingma and Ba, 2014] with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  (default in PyTorch [Paszke et al., 2019]) and a fixed learning rate of 0.0002. We use mini-batches of size 512 training samples. One epoch of training with synthetic data is comprised of 50 000 samples, while the test set with 10 000 samples is kept fixed and not seen during training. The parameters  $\varphi$  are initialized using a standard Gumbel distribution. We find that when learning both  $\Phi$  and  $\theta$ , training is more stable when rescaling the Gumbel noise by a factor of 0.001, both for initialization as well as for the Gumbel reparametrizations during training. We keep the softmax temperature in all Gumbel reparametrizations fixed at  $\tau = 1$ , which has been demonstrated to work well in practice [Jang et al., 2016]. For all evaluated algorithms, the learned  $\Phi$  is always initialized with the same random seed, which also yields the random but fixed  $\Phi$  after the top- $K$  operation. Before each training run, the optimal scalar scaling factor of the measurement operator is determined via a grid search using the initial  $\Phi$ . When training on MNIST, the previously determined optimal scalar is multiplied by 0.9 to prevent gradient explosion as more pixels move to the center of the masks.

*Discrete optimization baselines:* In each optimization step, the baseline discrete optimization algorithms propose neighbors to the current  $\Phi$ , which are either accepted or rejected. We benchmark against a greedy baseline and a Simulated Annealing (SimAn) [Pincus, 1970, Romeijn and Smith, 1994] baseline. Let the shorthand  $\mathcal{L}_\Phi$  denote the loss over a randomly sampled mini-batch using  $\Phi$  as a measurement operator. The greedy algorithm only accepts neighbors  $\Phi'$  that strictly decrease the loss, i.e., when  $\mathcal{L}_{\Phi'} < \mathcal{L}_\Phi$ . The SimAn baseline can also accept neighbors that increase the loss, with the rule becoming stricter as training progresses. More precisely, SimAn accepts a neighbor  $\Phi'$  if it strictly decreases the loss or if  $\exp((\mathcal{L}_\Phi - \mathcal{L}_{\Phi'})/\tau) < u$ , where  $u \sim U(0, 1)$  is drawn from a uniform distribution and  $\tau$  is a temperature parameters. The acceptance probability depends crucially on  $\tau$ : initially  $\tau$  is chosen larger, such that SimAn explores solutions far from its initial  $\Phi$ . As training progresses,  $\tau$  is steadily decreased. We manually tuned the initial  $\tau$  to accept about 80% of all neighbors in the first epoch and the decay rate such that changes are accepted for at least 100 training epochs. For single pixel imaging this revealed an initial  $\tau$  of 0.0012 and a temperature decay factor of 0.9997. For left- $d$ -regular graphs, a slightly larger  $\tau$  of 0.003 and a temperature decay of 0.9998 was determined to work well. We found that increasing the batch size 10-fold (to stabilize the quantities  $\mathcal{L}_\Phi$  and  $\mathcal{L}_{\Phi'}$ ) did not have a significant effect on the performance of the baseline algorithms. Neighbors to the current  $\Phi$  are chosen as follows: a random subset of the partition defining the structure is chosen. In this subset, the position of a random one and a random zero are swapped. For example, in the case of single-pixel imaging, this amounts to selecting a random row, and swapping a random one with a random zero within that row. This procedure for drawing neighbors ensures that the neighbors always satisfies the constraints. It also fulfills a notion of being the ‘smallest’ step away from  $\Phi$ .

<sup>3</sup>The full code of our implementation is available at <https://github.com/josauder/glodismo>

## APPENDIX B SUPPLEMENTAL EXPERIMENTS

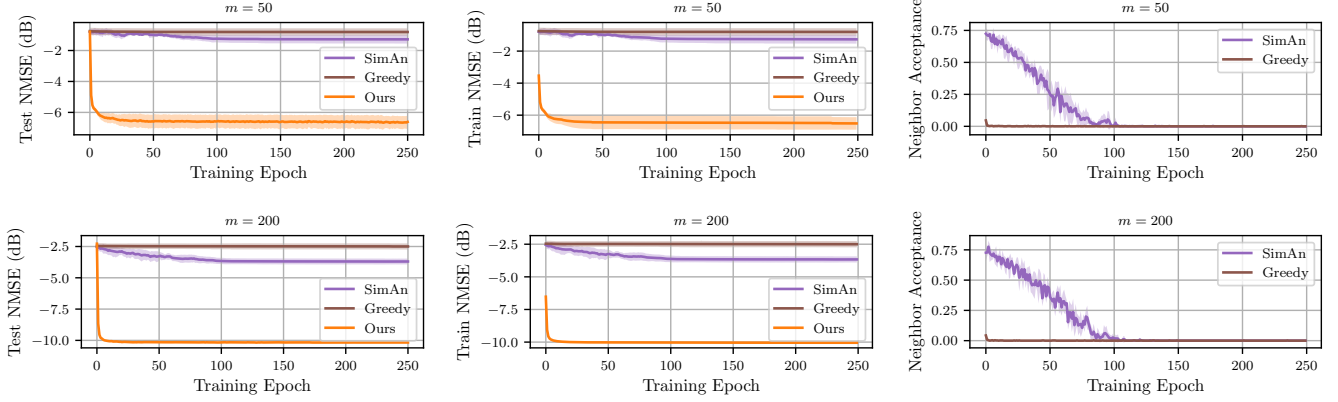


Fig. 10: Comparing the test NMSE (left), train NMSE (center) and the average neighbor acceptance rate of the baseline algorithms throughout a training epoch (right) in the single pixel imaging setup from Section IV with  $m = 50$  (top) and  $m = 200$  (bottom) measurements. The shaded area represents the standard deviation over 10 random seeds. This Figure highlights the training stability of GLODISMO and shows that it is not overly dependent on the random seed. Furthermore, the right plots show that the baselines are reasonably tuned.

### APPENDIX C LEARNING FAST-TRANSFORMABLE STRUCTURED SPARSE MATRICES

In many real-world applications, in particular when  $m$  and  $n$  are large, it is imperative that multiplication with  $\Phi$  can be done via a fast transform algorithm in  $O(n \log n)$  operations instead of a full matrix multiplication which would require  $O(mn)$  operations. In this section we demonstrate how GLODISMO can be easily extended to learn discrete row-wise masked circulant matrices, which are a prime example of a class of matrices for which fast transforms exist.

As circulant matrices are diagonalized by the Discrete Fourier Transform (DFT), a row-wise masked circulant matrix  $\Phi$  can be written as  $\Phi = P_{\Omega} \mathcal{F}^{-1} \Lambda \mathcal{F}$  for a diagonal matrix  $\Lambda$ , DFT matrix  $\mathcal{F}$ , and a row-selection mask  $P_{\Omega}$ . This allows for multiplication in  $O(n \log n)$  with  $\Phi$  by using the only the DFT and point-wise vector multiplication. While it has been shown that certain random row-wise masked circulant matrices fulfill the restricted isometry property with high probability [Rauhut et al., 2012], GLODISMO provides an approach for incorporating structural constraints and training data priors into the measurement operator. Our approach can be easily extended to work with any Toeplitz matrix, for which fast transforms based on the FFT also exist.

In our setup, we use GLODISMO to learn binary row-wise circulant matrices with  $d = 31$  ones per row for single pixel imaging. We employ a Gumbel top- $d$  operation on a vector of length  $n$ , and use this vector to explicitly construct the circulant matrix. With a second Gumbel reparametrization, we learn the row-selection mask  $P_{\Omega}$  that selects  $m$  out of  $n$  rows. Then,  $\Phi$  is obtained by masking the circulant matrix with  $P_{\Omega}$ .

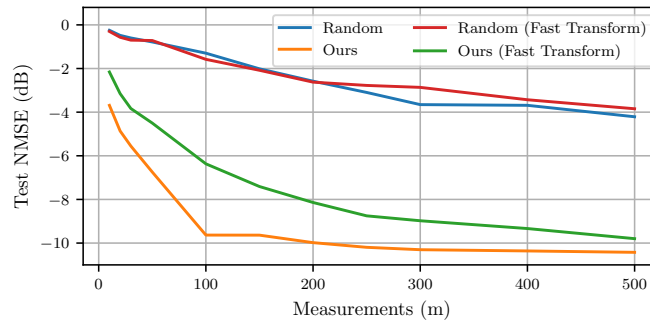


Fig. 11: Comparing the NMSE on the MNIST test set of fast-transformable single pixel imaging matrices with  $d = 32$  to regular single pixel imaging using IHT with  $T = 15$  iterations.

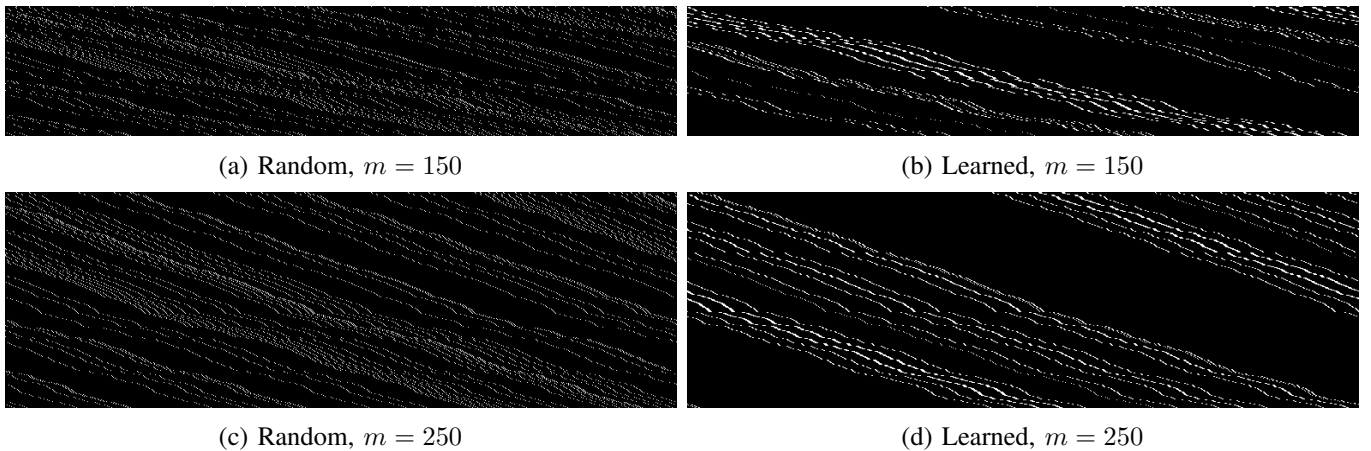


Fig. 12: Row-masked binary circulant single-pixel imaging measurement operators with  $d = 32$  ones per row. The learned fast-transformable operators clearly exhibit structure.



## APPENDIX D

### LEARNING SUPER-PIXEL MASKS

When using single-pixel cameras for imaging at wavelengths which are in the same order of magnitude as the pixel-sizes, strong diffraction effects appear leading to poor image quality [Augustin et al., 2017]. It has been observed that using super-pixel masks, in which multiple adjacent pixels are always selected together, can mitigate this effect [Sun et al., 2016]. We demonstrate that it is possible to learn such super-pixel masks using a Gumbel reparametrization by using a simple modification of our method for learning a pixel mask.

The procedure for obtaining a super-pixel mask is as follows: for a fixed number of super-pixels  $d$  with super-pixel height and width of  $\Delta$ , we follow the same procedure as when learning a pixel mask for determining the centers of the super-pixels. Each row of the matrix, corresponding to one mask, is then reshaped to image dimensions and the centers of the super-pixels are then convolved with a 2D convolution kernel of size  $\Delta \times \Delta$  whose values are all ones. Then, the masks are reshaped back to rows of  $\Phi$  and the maximum values of  $\Phi$  are clipped to 1. Computing the gradients with regards to the learned Gumbel parametrization is then done via automatic differentiation.

We employ this technique in the same MNIST single-pixel imaging setup as in Section IV-A. The results are shown in Figure 13. Single-pixel imaging using super-pixel masks learned with GLODISMO outperforms using random super-pixel masks by a significant margin and match the performance of the learned pixel-wise masks. Figure 14 shows how the learned masks adapt to the underlying dataset.

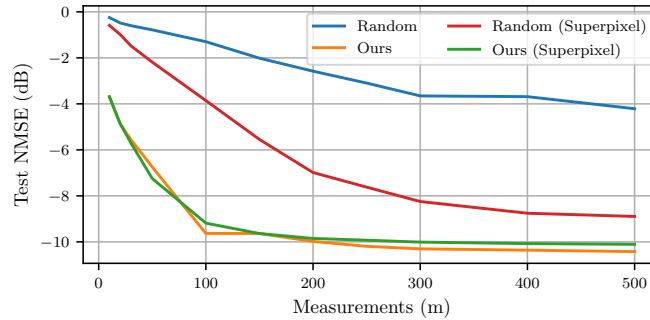


Fig. 13: Comparing the NMSE on the MNIST test set of super-pixel single pixel imaging matrices with  $d = 32$  to regular single pixel imaging using IHT with  $T = 15$  iterations.

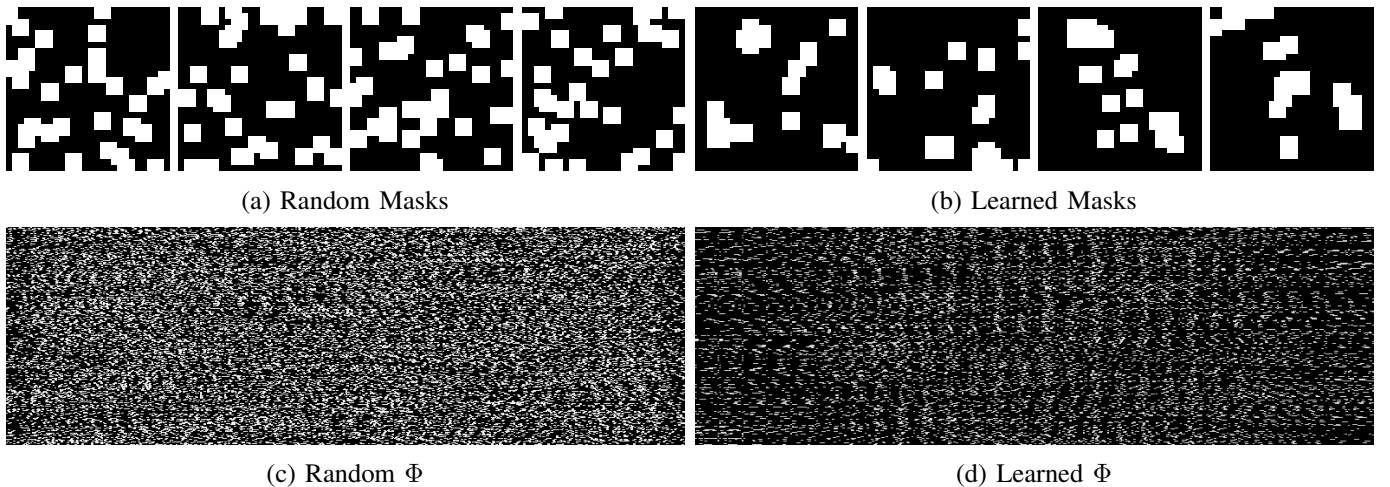


Fig. 14: Super-pixel masks for single pixel imaging with  $d = 32$  super-pixels per mask. The masks of the learned  $\Phi$  clearly exhibit structure adapted to the MNIST dataset.

## APPENDIX E

### LEARNING POOLING MATRICES FOR GROUP TESTING

Consider an array of substances that are to be chemically tested for the possession or absence of a certain rare property. Examples are detecting disease in pandemic situations or detecting tainted products in chemical production chains. It is desirable to detect as many positives with the least amount of tests possible to economise on the tests as well as the personnel conducting the tests. The mathematical field of group testing is concerned with reducing the number of tests needed to faithfully discover all true positives by pooling the substances together. Here we consider a non-adaptive group testing approach in which all tests are conducted in parallel, and therefore the outcomes of the tests are mutually independent. Because usually only a small number of tests are positive, an array of samples to be tested can be modeled as a sparse vector. For  $n$  samples (of which  $s \ll n$  are positive), and  $m$  tests, the  $ij$ -th component of the measurement matrix  $\Phi$  for non-adaptive group testing then describes whether specimen  $j$  is pooled into test  $i$ . As the tests are commonly standardized, it is sensible to model the number  $d$  of specimen per test (i.e. the number of ones per row) to be fixed. The quantity to be measured in the individuals is commonly modeled to be non-negative, as chemical tests in general measure physical quantities. Therefore, the inverse problem of identifying the infected individuals from the pooled tests can be solved using non-negative least absolute deviations (NNLAD) [Petersen et al., 2020], that is

$$\min_{\hat{x} \geq 0} \|\Phi \hat{x} - y\|_1.$$

This convex optimization problem can be solved using a projected subgradient descent method as described in [Petersen et al., 2021], which can be unrolled and back-propagated through out of the box. The crucial question remains how to construct a  $\Phi$  that can faithfully recover the positives from as few tests as possible. It is well-established that random matrices can be employed [Guruswami et al., 2009], when the conditions on  $m, n$  and  $s$  for compressed sensing are satisfied. Here, we follow the same problem size as considered in [Petersen et al., 2020], which assumes  $m = 248, n = 961$ , and  $d = 31$ . The measurement vector  $y$  is contaminated by student-t distributed noise equivalent to a signal-to-noise ratio of 40dB. In these experiments, the size of the support  $s$  of the synthetic sparse vectors is fixed as  $s = 80$  (which is far above the value for which theoretical compressed sensing guarantees hold). More specifically,  $s$  non-zero components are determined randomly without replacement. Mimicking viral load distributions, the non-zero components follow a heavy-tailed Beta distribution, where  $x_i^{\text{nonzero}} \sim_{i.i.d} \text{Beta}(2, 8)$ . We consider an individual to have tested positive if the value of the reconstruction is higher than 0.01 at a given index. The hyperparameters of the NNLAD algorithm were hand-tuned and determined to be  $\sigma = 0.1, \tau = 0.6$ . In group testing, the time taken to conduct the tests is usually orders of magnitudes higher than the time required for solving the recovery problem algorithmically. In particular, the number of iterations taken until convergence is insignificant in group testing. At the same time, learning a pooling matrix with a large number of iterations is infeasible. Therefore, we learn  $\Phi$  using projected subgradient descent for NNLAD [Petersen et al., 2021] with  $T = 200$  iterations, and do inference with 1000 iterations, resulting in a mismatch between the train and test settings. To stabilize GLODISMO when training with 200 iterations, the loss is averaged over the reconstruction after every iteration and not just the final estimate.

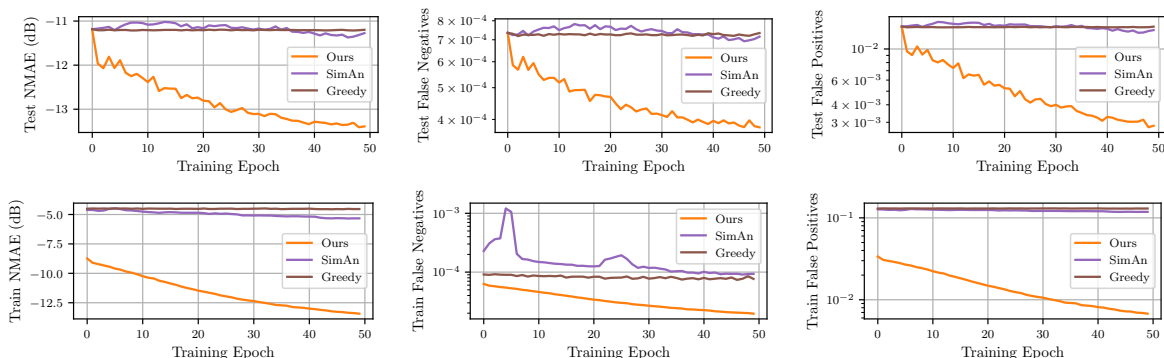


Fig. 15: NMAE (left), false negatives (center) and false positive (right) at training time with  $T = 200$  iterations (bottom) and during inference with  $T = 1000$  iterations (top) of the unrolled NNLAD algorithm. Despite the train/test mismatch, GLODISMO is able to learn a  $\Phi$  that outperforms the baselines.