

Model-Based Deep Learning

Nir Shlezinger, Jay Whang, Yonina C. Eldar, and Alexandros G. Dimakis

Abstract—Signal processing, communications, and control have traditionally relied on classical statistical modeling techniques. Such model-based methods utilize mathematical formulations that represent the underlying physics, prior information and additional domain knowledge. Simple classical models are useful but sensitive to inaccuracies and may lead to poor performance when real systems display complex or dynamic behavior. On the other hand, purely data-driven approaches that are model-agnostic are becoming increasingly popular as data sets become abundant and the power of modern deep learning pipelines increases. Deep neural networks (DNNs) use generic architectures which learn to operate from data, and demonstrate excellent performance, especially for supervised problems. However, DNNs typically require massive amounts of data and immense computational resources, limiting their applicability for some scenarios.

In this article we present the leading approaches for studying and designing model-based deep learning systems. These are methods that combine principled mathematical models with data-driven systems to benefit from the advantages of both approaches. Such model-based deep learning methods exploit both partial domain knowledge, via mathematical structures designed for specific problems, as well as learning from limited data. Among the applications detailed in our examples for model-based deep learning are compressed sensing, digital communications, and tracking in state-space models. Our aim is to facilitate the design and study of future systems on the intersection of signal processing and machine learning that incorporate the advantages of both domains.

I. INTRODUCTION

Traditional signal processing is dominated by algorithms that are based on simple mathematical models which are hand-designed from domain knowledge. Such knowledge can come from statistical models based on measurements and understanding of the underlying physics, or from fixed deterministic representation of the particular problem at hand. These domain-knowledge-based processing algorithms, which we refer to henceforth as *model-based methods*, carry out inference based on knowledge of the underlying model relating the observations at hand and the desired information. Model-based methods do not rely on data to learn their mapping, though data is often used to estimate a small number of parameters. Fundamental techniques like the Kalman filter and message passing algorithms belong to the class of model-based methods. Classical statistical models rely on simplifying assumptions (e.g., linear systems, Gaussian and independent noise, etc.) that make inference tractable, understandable and computationally efficient. On the other hand, simple models

frequently fail to represent nuances of high-dimensional complex data and dynamic variations.

The incredible success of deep learning, e.g., on vision [1], [2] as well as challenging games such as Go [3] and Starcraft [4], has initiated a general data-driven mindset. It is currently prevalent to replace simple principled models with purely data-driven pipelines, trained with massive labeled data sets. In particular, deep neural networks (DNNs) can be trained in a supervised way end-to-end to map inputs to predictions. The benefits of data-driven methods over model-based approaches are twofold: First, purely-data-driven techniques do not rely on analytical approximations and thus can operate in scenarios where analytical models are not known. Second, for complex systems, data-driven algorithms are able to recover features from observed data which are needed to carry out inference [5]. This is sometimes difficult to achieve analytically, even when complex models are perfectly known.

The computational burden of training and utilizing highly-parametrized DNNs, as well as the fact that massive data sets are typically required to train such DNNs to learn a desirable mapping, may constitute major drawbacks in various signal processing, communications, and control applications. This is particularly relevant for hardware-limited devices, such as mobile phones, unmanned aerial vehicles, and Internet of Things (IOT) systems, which are often limited in their ability to utilize highly-parametrized DNNs [6], and require adapting to dynamic conditions. Furthermore, DNNs are commonly utilized as black-boxes; understanding how their predictions are obtained and characterizing confidence intervals tends to be quite challenging. As a result, deep learning does not yet offer the interpretability, flexibility, versatility, and reliability of model-based methods [7].

The limitations associated with model-based methods and black-box deep learning systems gave rise to a multitude of techniques for combining signal processing and machine learning to benefit from both approaches. These methods are application-driven, and are thus designed and studied in light of a specific task. For example, the combination of DNNs and model-based compressed sensing (CS) recovery algorithms was shown to facilitate sparse recovery [8], [9] as well as enable CS beyond the domain of sparse signals [10], [11]; Deep learning was used to empower regularized optimization methods [12], [13], while model-based optimization contributed to the design of DNNs for such tasks [14]; Digital communication receivers used DNNs to learn to carry out and enhance symbol detection and decoding algorithms in a data-driven manner [15]–[17], while symbol recovery methods enabled the design of model-aware deep receivers [18]–[21]. The proliferation of hybrid model-based/data-driven systems, each designed for a unique task, motivates establishing a concrete systematic framework for combining domain knowledge in the

N. Shlezinger is with the School of ECE, Ben-Gurion University of the Negev, Be'er-Sheva, Israel (e-mail: nirshl@bgu.ac.il). J. Whang is with the Department of CS, University of Texas at Austin, Austin, TX (e-mail: jaywhang@cs.utexas.edu). Y. C. Eldar is with the Faculty of Math and CS, Weizmann Institute of Science, Rehovot, Israel (e-mail: yonina@weizmann.ac.il). A. G. Dimakis is with the Department of ECE, University of Texas at Austin, Austin, TX (e-mail: dimakis@austin.utexas.edu).

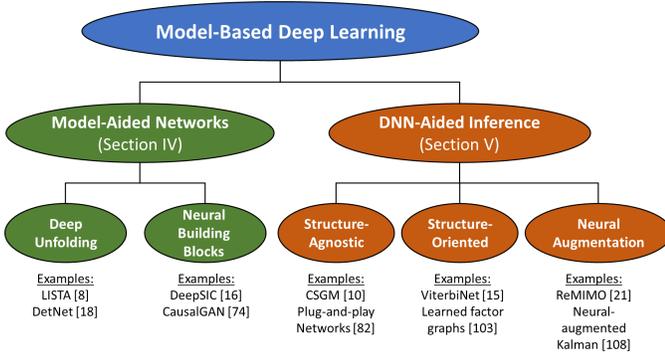


Fig. 1: Division of model-based deep learning techniques into categories and sub-categories.

form of model-based methods and deep learning, which is the focus of this article.

In this article we review leading strategies for designing systems whose operation combines domain knowledge and data via model-based deep learning in a tutorial fashion. To that aim, we present a unified framework for studying hybrid model-based/data-driven systems, without focusing on a specific application, while being geared towards families of problems typically studied in the signal processing literature. The proposed framework divides systems combining model-based signal processing and deep learning into two main strategies: The first category includes DNNs whose architecture is specialized to the specific problem using model-based methods, referred to here as *model-aided networks*. The second one, which we refer to as *DNN-aided inference*, consists of techniques in which inference is carried out by a model-based algorithm whose operation is augmented with deep learning tools. This integration of model-agnostic deep learning tools allows one to use model-based inference algorithms while having access only to partial domain knowledge. Based on this division, we provide concrete guidelines for studying, designing, and comparing model-based deep learning systems. An illustration of the proposed division into categories and sub-categories is depicted in Fig. 1.

We begin by discussing the high level concepts of model-based, data-driven, and hybrid schemes. Since we focus on DNNs as the current leading data-driven technique, we briefly review basic concepts in deep learning, ensuring that the tutorial is accessible to readers without background in deep learning. We then elaborate on the fundamental strategies for combining model-based methods with deep learning. For each such strategy, we present a few concrete implementation approaches in a systematic manner, including established approaches such as deep unfolding, which was originally proposed in 2010 by Gregor and LeCun [8], as well as recently proposed model-based deep learning paradigms such as DNN-aided inference [22] and neural augmentation [23]. For each approach we formulate system design guidelines for a given problem; provide detailed examples from the recent literature; and discuss its properties and use-cases. Each of our detailed examples focuses on a different application in signal processing, communications, and control, demonstrating the breadth and the wide variety of applications that can

benefit from such hybrid designs. We conclude the article with a summary and a qualitative comparison of model-based deep learning approaches, along with a description of some future research topics and challenges. We aim to encourage future researchers and practitioners with a signal processing background to study and design model-based deep learning.

This overview article focuses on strategies for designing architectures whose operation combines deep learning with model-based methods, as illustrated in Fig. 1. These strategies can also be integrated into existing mechanisms for incorporating model-based domain knowledge in the selection of the task for which data-driven systems are applied, as well as in the generation and manipulation of the data. An example of a family of such mechanisms for using model-based knowledge in the selection of the application and the data is the learning-to-optimize framework, which is the focus of growing attention in the context of wireless networks design [24]–[26]; this framework advocates the usage of pre-trained DNNs for realizing fast solvers for complex optimization problems which rely on objectives and constraints formulated based on domain knowledge, along with the usage of model-based generated data for offline training. An additional related family is that of channel autoencoders, which integrate mathematical modelling of random communication channels as layers of deep autoencoders to design channel codes [27], [28] and compression mechanisms [29].

The rest of this article is organized as follows: Section II discusses the concepts of model-based methods as compared to data-driven schemes, and how they give rise to the model-based deep learning paradigm. Section III reviews some basics of deep learning. The main strategies for designing model-based deep learning systems, i.e., model-aided networks and DNN-aided inference, are detailed in Sections IV–V, respectively. Finally, we provide a summary and discuss some future research challenges in Section VI.

II. MODEL-BASED VERSUS DATA-DRIVEN INFERENCE

We begin by reviewing the main conceptual differences between model-based and data-driven inference. To that aim, we first present a mathematical formulation of a generic inference problem. Then we discuss how this problem is tackled from a purely model-based perspective as well as from a purely data-driven one, where for the latter we focus on deep learning as a family of generic data-driven approaches. We then formulate the notion of model-based deep learning based upon these distinct strategies.

A. Inference Systems

The term *inference* refers to the ability to conclude based on evidence and reasoning. While this generic definition can refer to a broad range of tasks, we focus in our description on systems which estimate or make predictions based on a set of observed variables. In this wide family of problems, the system is required to map an input variable $\mathbf{x} \in \mathcal{X}$ into a prediction of a label variable $\mathbf{s} \in \mathcal{S}$, denoted $\hat{\mathbf{s}}$, where \mathcal{X}

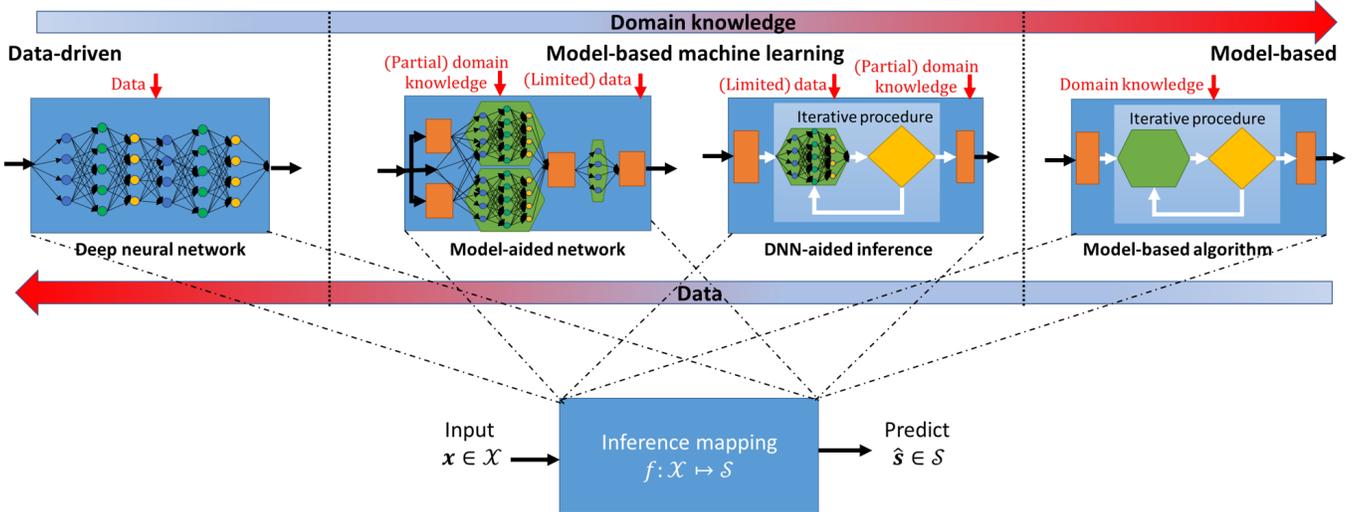


Fig. 2: Illustration of model-based versus data-driven inference. The red arrows correspond to computation performed before the particular inference data is received.

and \mathcal{S} are referred to as the *input space* and the *label space*, respectively. An inference rule can thus be expressed as

$$f : \mathcal{X} \mapsto \mathcal{S} \quad (1)$$

and the space of inference mappings is denoted by \mathcal{F} . We use $l(\cdot)$ to denote a cost measure defined over $\mathcal{F} \times \mathcal{X} \times \mathcal{S}$, dictated by the specific task [30, Ch. 2]. The fidelity of an inference mapping is measured by the *risk function*, also known as the generalization error, given by $\mathbb{E}_{\mathbf{x}, \mathbf{s} \sim p_{\mathbf{x}, \mathbf{s}}} \{l(f, \mathbf{x}, \mathbf{s})\}$, where $p_{\mathbf{x}, \mathbf{s}}$ is the underlying statistical model relating the input and the label. The goal of both model-based methods and data-driven schemes is to design the inference rule $f(\cdot)$ to minimize the risk for a given problem. The main difference between these strategies is what information is utilized to tune $f(\cdot)$.

B. Model-Based Methods

Model-based algorithms, also referred to as *hand-designed methods* [31], set their inference rule, i.e., tune f in (1) to minimize the risk function, based on domain knowledge. The term *domain knowledge* typically refers to prior knowledge of the underlying statistics relating the input \mathbf{x} and the label \mathbf{s} . In particular, an analytical mathematical expression describing the underlying model, i.e., $p_{\mathbf{x}, \mathbf{s}}$, is required. Model-based algorithms can provably implement the risk minimizing inference mapping, e.g., the maximum a-posteriori probability (MAP) rule. While computing the risk minimizing rule is often computationally prohibitive, various model-based methods approximate this rule at controllable complexity, and in some cases also provably approach its performance. This is typically achieved using iterative methods comprised of multiple stages, where each stage involves generic mathematical manipulations and model-specific computations.

Model-based methods do not rely on data to learn their mapping, as illustrated in the right part of Fig. 2, though data is often used to estimate unknown model parameters. In practice, accurate knowledge of the statistical model relating the observations and the desired information is typically unavailable,

and thus applying such techniques commonly requires imposing some assumptions on the underlying statistics, which in some cases reflects the actual behavior, but may also constitute a crude approximation of the true dynamics. In the presence of inaccurate model knowledge, either as a result of estimation errors or due to enforcing a model which does not fully capture the environment, the performance of model-based techniques tends to degrade. This limits the applicability of model-based schemes in scenarios where, e.g., $p_{\mathbf{x}, \mathbf{s}}$ is unknown, costly to estimate accurately, or too complex to express analytically.

C. Data-Driven Schemes

Data-driven systems learn their mapping from data. In a supervised setting, data is comprised of a training set consisting of n_t pairs of inputs and their corresponding labels, denoted $\{(\mathbf{x}_t, \mathbf{s}_t)\}_{t=1}^{n_t}$. Data-driven schemes do not have access to the underlying distribution, and thus cannot compute the risk function. As a result, the inference mapping is typically tuned based on an empirical risk function, referred henceforth as *loss function*, which for an inference mapping f is given by

$$\mathcal{L}(f) = \frac{1}{n_t} \sum_{t=1}^{n_t} l(f, \mathbf{x}_t, \mathbf{s}_t). \quad (2)$$

Since one can usually form an inference rule which minimizes the empirical loss (2) by memorizing the data, i.e., overfit, data-driven schemes often restrict the domain of feasible inference rules [30, Ch. 2]. A leading strategy in data-driven systems, upon which deep learning is based, is to assume some highly-expressive generic parametric model on the mapping in (1), while incorporating optimization mechanisms to avoid overfitting and allow the resulting system to infer reliably with new data samples. In such cases, the inference rule is dictated by a set of parameters denoted θ , and thus the system mapping is written as f_θ .

The conventional application of deep learning implements f_θ using a DNN architecture, where θ represent the weights

of the network. Such highly-parametrized networks can effectively approximate any Borel measurable mapping, as follows from the universal approximation theorem [32, Ch. 6.4.1]. Therefore, by properly tuning their parameters using sufficient training data, as we elaborate in Section III, one should be able to obtain the desirable inference rule.

Unlike model-based algorithms, which are specifically tailored to a given scenario, purely-data-driven methods are model-agnostic, as illustrated in the left part of Fig. 2. The unique characteristics of the specific scenario are encapsulated in the learned weights. The parametrized inference rule, e.g., the DNN mapping, is generic and can be applied to a broad range of different problems. While standard DNN structures are highly model-agnostic and are commonly treated as black boxes, one can still incorporate some level of domain knowledge in the selection of the specific network architecture. For instance, when the input is known to exhibit temporal correlation, architectures based on recurrent neural networks (RNNs) [33] or attention mechanisms [34] are often preferred. Alternatively, in the presence of spatial patterns, one may utilize convolutional layers [35]. An additional method to incorporate domain knowledge into a black box DNN is by pre-processing of the input via, e.g., feature extraction.

The generic nature of data-driven strategies induces some drawbacks. Broadly speaking, learning a large number of parameters requires a massive data set to train on. Even when a sufficiently large data set is available, the resulting training procedure is typically lengthy and involves high computational burden. Finally, the black-box nature of the resulting mapping implies that data-driven systems in general lack interpretability, making it difficult to provide performance guarantees and insights into the system operation.

D. Model-Based Deep Learning

Completely separating existing literature into model-based versus data-driven is a subjective and debatable task. Instead, we focus on some approaches which clearly lie in the middle ground to give a useful overview of the landscape. The considered families of methods incorporate domain knowledge in the form of an established model-based algorithm which is suitable for the problem at hand, while combining capabilities to learn from data via deep learning techniques.

Model-based deep learning schemes tune their mapping of the input \mathbf{x} based on both data, e.g., a labeled training set $\{(\mathbf{x}_t, \mathbf{s}_t)\}_{t=1}^{n_t}$, as well as some domain knowledge, such as partial knowledge of the underlying distribution $p_{\mathbf{x}, \mathbf{s}}$. Such hybrid data-driven model-aware systems can typically learn their mappings from smaller training sets compared to purely model-agnostic DNNs, and commonly operate without full accurate knowledge of the underlying model upon which model-based methods are based.

Most existing techniques for implementing inference rules in a hybrid model-based/data-driven fashion are designed for a specific application, i.e., to solve a given problem rather than formulate a systematic methodology. Nonetheless, one can identify a common rationale for categorizing existing schemes in a systematic manner that is not tailored to a specific

scenario. In particular, model-based deep learning techniques can be divided into two main strategies, as illustrated in Fig. 2. These strategies may each be further specialized to various different tasks, as we show in the sequel. The first of the two, which we refer to as *model-aided networks*, utilizes DNNs for inference; however, rather than using conventional DNN architectures, here a specific DNN tailored for the problem at hand is designed by following the operation of suitable model-based methods. The second strategy, which we call *DNN-aided inference systems*, uses conventional model-based methods for inference; however, unlike purely model-based schemes, here specific parts of the model-based algorithm are augmented with deep learning tools, allowing the resulting system to implement the algorithm while learning to overcome partial or mismatched domain knowledge from data.

The systematic categorization of model-based deep learning methodologies can facilitate the study and design of future techniques in different and diverse application areas. One may also propose schemes which combine aspects from both categories, building upon the understanding of the characteristics and gains of each approach, discussed in the sequel. Since both strategies rely on deep learning tools, we first provide a brief overview of key concepts in deep learning in the following section, after which we elaborate on model-aided networks and DNN-aided inference in Sections IV and V, respectively.

III. BASICS OF DEEP LEARNING

Here, we cover the basics of deep learning required to understand the DNN-based components in the model-based/data-driven approaches discussed later. Our aim is to equip the reader with necessary foundations upon which our formulations of model-based deep learning systems are presented.

As discussed in Subsection II-C, in deep learning, the target mapping is constrained to take the form of a parametrized function $f_{\theta} : \mathcal{X} \rightarrow \mathcal{S}$. In particular, the inference mapping belongs to a fixed family of functions \mathcal{F} specified by a predefined DNN architecture, which is represented by a specific choice of the parameter vector θ . Once the function class \mathcal{F} and loss function \mathcal{L} are defined, where the latter is dictated by the training data (2) while possibly including some regularization on θ , one may attempt to find the function which minimizes the loss within \mathcal{F} , i.e.,

$$\theta^* = \arg \min_{f_{\theta} \in \mathcal{F}} \mathcal{L}(f_{\theta}). \quad (3)$$

A common challenge in optimizing based on (3) is to guarantee that the inference mapping learned using the data-based loss function rather than the model-based risk function will not overfit and be able to generalize, i.e., infer reliably from new data samples. Since the optimization in (3) is carried out over θ , we write the loss as $\mathcal{L}(\theta)$ for brevity.

The above formulation naturally gives rise to three fundamental components of deep learning: the DNN *architecture* that defines the function class \mathcal{F} ; the task-specific *loss* function $\mathcal{L}(\theta)$; and the *optimizer* that dictates how to search for the optimal f_{θ} within \mathcal{F} . Therefore, our review of the basics of deep learning commences with a description of the fundamental architecture and optimizer components in Subsection III-A.

We then present several representative tasks along with their corresponding typical loss functions in Subsection III-B.

A. Deep Learning Preliminaries

The formulation of the parametric empirical risk in (3) is not unique to deep learning, and is in fact common to numerous machine learning schemes. The strength of deep learning, i.e., its ability to learn accurate complex mappings from large data sets, is due to its use of DNNs to enable a highly-expressive family of function classes \mathcal{F} , along with dedicated optimization algorithms for tuning the parameters from data. In the following we discuss the high level notion of DNNs, followed by a description of how they are optimized.

a) *Neural Network Architecture*: DNNs implement parametric functions comprised of a sequence of differentiable transformations called *layers*, whose composition maps the input to a desired output. Specifically, a DNN f_{θ} consisting of k layers $\{h_1, \dots, h_k\}$ maps the input \mathbf{x} to the output $\hat{\mathbf{s}} = f_{\theta}(\mathbf{x}) = h_k \circ \dots \circ h_1(\mathbf{x})$, where \circ denotes function composition. Since each layer h_i is itself a parametric function, the parameters set of the entire network f_{θ} is the union of all of its layers' parameters, and thus f_{θ} denotes a DNN with parameters θ . The *architecture* of a DNN refers to the specification of its layers $\{h_i\}_{i=1}^k$.

A generic formulation which captures various parametrized layers is that of an affine transformation, i.e., $h(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ whose parameters are $\{\mathbf{W}, \mathbf{b}\}$. For instance, in *fully-connected (FC)* layers, also referred to as *dense* layers, one can optimize $\{\mathbf{W}, \mathbf{b}\}$ to take any value. Another extremely common affine transform layer is *convolutional* layers. Such layers apply a set of discrete convolutional kernels to signals that are possibly comprised of multiple channels, e.g., tensors. The vector representation of their output can be written as an affine mapping of the form $\mathbf{W}\mathbf{x} + \mathbf{b}$, where \mathbf{x} is the vectorization of the input, and \mathbf{W} is constrained to represent multiple channels of discrete convolutions [32, Ch. 9]. These convolutional neural networks (CNNs) are known to yield a highly parameter-efficient mapping that captures important invariances such as translational invariance in image data.

While many commonly used layers are affine, DNNs rely on the inclusion of *non-linear* layers. If all the layers of a DNN were affine, the composition of all such layers would also be affine, and thus the resulting network would only represent affine functions. For this reason, layers in a DNN are interleaved with *activation functions*, which are simple non-linear functions applied to each dimension of the input separately. Activations are often fixed, i.e., their mapping is not parametric and is thus not optimized in the learning process. Some notable examples of widely-used activation functions include the rectified linear unit (ReLU) defined as $\text{ReLU}(x) = \max\{x, 0\}$ and the sigmoid $\sigma(x) = (1 + \exp(-x))^{-1}$.

b) *Choice of Optimizer*: Given a DNN architecture and a loss function $\mathcal{L}(\theta)$, finding a globally optimal θ that minimizes \mathcal{L} is a hopelessly intractable task, especially at the scale of millions of parameters or more. Fortunately, recent success of deep learning has demonstrated that gradient-based optimization methods work surprisingly well despite their inability

to find global optima. The simplest such method is *gradient descent*, which iteratively updates the parameters:

$$\theta_{q+1} = \theta_q - \eta_q \nabla_{\theta} \mathcal{L}(\theta_q) \quad (4)$$

where η_q is the *step size* that may change as a function of the step count q . Since the gradient $\nabla_{\theta} \mathcal{L}(\theta_q)$ is often too costly to compute over the entire training data, it is estimated from a small number of randomly chosen samples (i.e., a mini-batch). The resulting optimization method is called *mini-batch stochastic gradient descent* and belongs to the family of stochastic first-order optimizers.

Stochastic first-order optimization techniques are well-suited for training DNNs because their memory usage grows only linearly with the number of parameters, and they avoid the need to process the entire training data at each step of optimization. Over the years, numerous variations of stochastic gradient descent have been proposed. Many modern optimizers such as RMSProp [36] and Adam [37] use statistics from previous parameter updates to adaptively adjust the step size for each parameter separately (i.e., for each dimension of θ).

B. Common Deep Learning Tasks

As detailed above, the data-driven nature of deep learning is encapsulated in the dependence of the loss function on the training data. Thus, the loss function not only implicitly defines the task of the resulting system, but also dictates what kind of data is required. Based on the requirements placed on the training data, problems in deep learning largely fall under three different categories: supervised, semi-supervised, and unsupervised. Here, we define each category and list some example tasks as well as their typical loss functions.

a) *Supervised Learning*: In supervised learning, the training data consists of a set of input-label pairs $\{(\mathbf{x}_t, \mathbf{s}_t)\}_{t=1}^{n_t}$, where each pair takes values in $\mathcal{X} \times \mathcal{S}$. As discussed in Subsection II-C, the goal is to recover a mapping f_{θ} which minimizes the risk function, i.e., the generalization error. This is done by optimizing the DNN mapping f_{θ} using the data-based empirical loss function $\mathcal{L}(\theta)$ (2). This setting encompasses a wide range of problems including regression, classification, and structured prediction, through a judicious choice of the loss function. Below we review commonly used loss functions for classification and regression tasks.

1) *Classification*: Perhaps one of the most widely-known success stories of DNNs, classification (image classification in particular) has remained a core benchmark since the introduction of AlexNet [38]. In this setting, we are given a training set $\{(\mathbf{x}_t, \mathbf{s}_t)\}_{t=1}^{n_t}$ containing input-label pairs, where each \mathbf{x}_t is a fixed-size input, e.g., an image, and \mathbf{s}_t is the one-hot encoding of the class. Such one-hot encoding of class c can be viewed as a probability vector for a K -way categorical distribution, with $K = |\mathcal{S}|$, with all probability mass placed on class c .

The DNN mapping f_{θ} for this task is appropriately designed to map an input \mathbf{x}_t to the probability vector $\hat{\mathbf{s}}_t \triangleq f(\mathbf{x}_t) = \langle \hat{s}_{t,1}, \dots, \hat{s}_{t,K} \rangle$, where $\hat{s}_{t,c}$ denotes the c -th component of $\hat{\mathbf{s}}_t$. This parametrization allows for the model to return a soft decision in the form of a categorical distribution over the classes.

A natural choice of loss function for this setting is the *cross-entropy* loss, defined as

$$\mathcal{L}_{\text{CE}}(\boldsymbol{\theta}) = \frac{1}{n_t} \sum_{t=1}^{n_t} \sum_{c=1}^K s_{t,c} (-\log \hat{s}_{t,c}). \quad (5)$$

For a sufficiently large set of i.i.d. training pairs, the empirical cross entropy loss approaches the expected cross entropy measure, which is minimized when the DNN output matches the true conditional distribution $p_{s|x}$. Consequently, minimizing the cross-entropy loss encourages the DNN output to match the ground truth label, and its mapping closely approaches the true underlying posterior distribution when properly trained.

The formulation of the cross entropy loss (5) implicitly assumes that the DNN returns a valid probability vector, i.e., $\hat{s}_{t,c} \geq 0$ and $\sum_{c=1}^K \hat{s}_{t,c} = 1$. However, there is no guarantee that this will be the case, especially at the beginning of training when the parameters of the DNN are more or less randomly initialized. To guarantee that the DNN mapping yields a valid probability distribution, classifiers typically employ the softmax function (e.g., on top of the output layer), given by:

$$\text{Softmax}(\mathbf{x}) = \left\langle \frac{\exp(x_1)}{\sum_{i=1}^d \exp(x_i)}, \dots, \frac{\exp(x_d)}{\sum_{i=1}^d \exp(x_i)} \right\rangle$$

where x_i is the i th entry of \mathbf{x} . Due to the exponentiation followed by normalization, the output of the softmax function is guaranteed to be a valid probability vector. In practice, one can compute the softmax function of the network outputs when evaluating the loss function, rather than using a dedicated output layer.

2) *Regression*: Another task where DNNs have been successfully applied is regression, where one attempts to predict continuous variables instead of categorical ones. Here, the labels $\{s_t\}$ in the training data represent some continuous value, e.g., in \mathbb{R} or some specified range $[a, b]$.

Similar to the usage of softmax layers for classification, an appropriate final activation function σ is needed, depending on the range of the variable of interest. For example, when regressing on a strictly positive value, a common choice is $\sigma(x) = \exp(x)$ or the softplus activation $\sigma(x) = \log(1 + \exp(x))$, so that the range of the network f_θ is constrained to be the positive reals. When the output is to be limited to an interval $[a, b]$, then one may use the mapping $\sigma(x) = a + (b - a)(1 + \tanh(x))/2$.

Arguably the most common loss function for regression tasks is the empirical mean-squared error (MSE), i.e.,

$$\mathcal{L}_{\text{MSE}}(\boldsymbol{\theta}) = \frac{1}{n_t} \sum_{t=1}^{n_t} (s_t - \hat{s}_t)^2. \quad (6)$$

b) *Unsupervised Learning*: In unsupervised learning, we are only given a set of examples $\{\mathbf{x}_t\}_{t=1}^{n_t}$ without labels. Since there is no label to predict, unsupervised learning algorithms are often used to discover interesting patterns present in the given data. Common tasks in this setting include clustering, anomaly detection, generative modeling, and compression.

1) *Generative models*: One goal in unsupervised learning of a generative model is to train a *generator* network $G_\theta(\mathbf{z})$ such

that the *latent variables* \mathbf{z} , which follow a simple distribution such as standard Gaussian, are mapped into samples obeying a distribution similar to that of the training data [32, Ch. 20]. For instance, one can train a generative model to map Gaussian vectors into images of human faces. A popular type of DNN-based generative model that tries to achieve this goal is *generative adversarial network (GAN)* [39], which has shown remarkable success in many domains.

GANs learn the generative model by employing a *discriminator* network D_φ to assess the generated samples, thus avoiding the need to mathematically handcraft a loss measure quantifying their quality. The parameters $\{\boldsymbol{\theta}, \varphi\}$ of the two networks are learned via *adversarial* training, where $\boldsymbol{\theta}$ and φ are updated in an alternating manner. The two networks G_θ and D_φ “compete” against each other to achieve opposite goals: G_θ tries to fool the discriminator, whereas D_φ tries to reliably distinguish real examples from the fake ones made by the generator.

Various methods have been proposed to train generative models in this adversarial fashion, including, e.g., the Wasserstein GAN [40], [41]; the least-squares GAN [42]; the Hinge GAN [43]; and the relativistic average GAN [44]. For simplicity, in the following we describe the original GAN formulation of [39]. Here, $D_\varphi : \mathcal{X} \rightarrow [0, 1]$ is a binary classifier trained to distinguish real examples \mathbf{x}_t from the fake examples generated by G_θ , and the GAN loss function is the minmax loss. The loss is optimized in an alternating fashion by tuning the discriminator φ to minimize $\mathcal{L}_D(\cdot)$ for a given generator $\boldsymbol{\theta}$, followed by a corresponding optimization of the generator based on its loss $\mathcal{L}_G(\cdot)$. These loss measures are given by

$$\begin{aligned} \mathcal{L}_D(\varphi|\boldsymbol{\theta}) &= \frac{-1}{2n_t} \sum_{t=1}^{n_t} \log D_\varphi(\mathbf{x}_t) + \log(1 - D_\varphi(G_\theta(\mathbf{z}_t))), \\ \mathcal{L}_G(\boldsymbol{\theta}|\varphi) &= \frac{-1}{n_t} \sum_{t=1}^{n_t} \log \log D_\varphi(G_\theta(\mathbf{z}_t)). \end{aligned}$$

Here, the latent variables $\{\mathbf{z}_t\}$ are drawn from its known prior distribution for each mini-batch.

Among currently available deep generative models, GANs achieve the best sample quality at an unprecedented resolution. For example, the current state-of-the-art model StyleGAN2 [45] is able to generate high-resolution (1024×1024) images that are nearly indistinguishable from real photos to a human observer. That said, GANs do come with several disadvantages as well. The adversarial training procedure is known to be unstable, and many tricks are necessary in practice to train a large GAN. Also because GANs do not offer any probabilistic interpretation, it is difficult to objectively evaluate the quality of a GAN.

2) *Autoencoders*: Another well-studied task in unsupervised learning is the training of an *autoencoder*, which has many uses such as dimensionality reduction and representation learning. An autoencoder consists of two neural networks: an *encoder* $f_{\text{enc}} : \mathcal{X} \mapsto \mathcal{Z}$ and a *decoder* $f_{\text{dec}} : \mathcal{Z} \mapsto \mathcal{X}$, where \mathcal{Z} is some predefined latent space. The primary goal of an autoencoder is to reconstruct a signal \mathbf{x} from itself by mapping it through $f_{\text{dec}} \circ f_{\text{enc}}$.

The task of autoencoding may seem pointless at first; indeed one can trivially recover \mathbf{x} by setting $\mathcal{Z} = \mathcal{X}$ and $f_{\text{enc}}, f_{\text{dec}}$ to be identity functions. The interesting case is when one imposes constraints which limit the ability of the network to learn the identity mapping [32, Ch. 14]. One way to achieve this is to form an undercomplete autoencoder, where the latent space \mathcal{Z} is restricted to be lower-dimensional than \mathcal{X} , e.g., $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Z} = \mathbb{R}^m$ for some $m < n$. This constraint forces the encoder to map its input into a more compact representation, while retaining enough information so that the reconstruction is as close to the original input as possible. Additional mechanisms for preventing an autoencoder from learning the identity mapping include imposing a regularizing term on the latent representation, as done in sparse autoencoders and contractive autoencoders, or alternatively, by distorting the input to the system, as carried out by denoising autoencoders [32, Ch. 14.2]. A common metric used to measure the quality of reconstruction is the MSE loss. Under this setting, we obtain the following loss function for training

$$\mathcal{L}_{\text{MSE}}(f_{\text{enc}}, f_{\text{dec}}) = \frac{1}{n_t} \sum_{t=1}^{n_t} \|\mathbf{x}_t - f_{\text{dec}}(f_{\text{enc}}(\mathbf{x}_t))\|_2^2. \quad (7)$$

c) *Semi-Supervised Learning*: Semi-supervised learning lies in the middle ground between the above two categories, where one typically has access to a large amount of unlabeled data and a small set of labeled data. The goal is to leverage the unlabeled data to improve performance on some supervised task to be trained on the labeled data. As labeling data is often a very costly process, semi-supervised learning provides a way to quickly learn desired inference rules without having to label all of the available unlabeled data points.

Various approaches have been proposed in the literature to utilize unlabeled data for a supervised task, see the detailed survey [46]. One such common technique is to guess the missing labels, while integrating dedicated mechanisms to boost confidence [47]. This can be achieved by, e.g., applying the DNN to various augmentations of the unlabeled data [48], while combining multiple regularization terms for encouraging consistency and low-entropy of the guessed labels [49], as well as training a teacher DNN using the available labeled data to produce guessed labels [50].

IV. MODEL-AIDED NETWORKS

Model-aided networks implement model-based deep learning by using model-aware algorithms to design deep architectures. Broadly speaking, model-aided networks implement the inference system using a DNN, similar to conventional deep learning. Nonetheless, instead of applying generic off-the-shelf DNNs, the rationale here is to tailor the architecture specifically for the scenario of interest, based on a suitable model-based method. By converting a model-based algorithm into a model-aided network, that learns its mapping from data, one typically achieves improved inference speed, as well as overcome partial or mismatched domain knowledge. In particular, model-aided networks can learn missing model parameters, such as channel matrices [19], dictionaries [51], and noise covariances [52], as part of the learning procedure.

Alternatively, it can be used to learn a surrogate model for which the resulting inference rule best matches the training data [53].

Model-aided networks obtain dedicated DNN architectures by identifying structures in a model-based algorithm one would have utilized for the problem given full domain knowledge and sufficient computational resources. Such structures can be given in the form of an iterative representation of the model-based algorithm, as exploited by *deep unfolding* detailed in Subsection IV-A, or via a block diagram algorithmic representation, which *neural building blocks* rely upon, as presented in Subsection IV-B. The dedicated neural network is then formulated as a discriminative architecture [54], [55] whose trainable parameters, intermediate mathematical manipulations, and interconnections follow the operations of the model-based algorithm, as illustrated in Fig. 3.

In the following we describe these methodologies in a systematic manner. In particular, our presentation of each approach commences with a high level description and generic design outline, followed by one or two concrete model-based deep learning examples from the literature, and concludes with a summarizing discussion. For each example, we first detail the system model and model-based algorithm from which it originates. Then, we describe the hybrid model-based/data-driven system by detailing its architecture and training procedure, as well as present some representative quantitative results.

A. Deep Unfolding

Deep unfolding [56], also referred to as *deep unrolling*, converts an iterative algorithm into a DNN by designing each layer to resemble a single iteration. Deep unfolding was originally proposed by Greger and LeCun in [8], where a deep architecture was designed to learn to carry out the iterative soft thresholding algorithm (ISTA) for sparse recovery. Deep unfolded networks have since been applied in various applications in image denoising [57], [58], sparse recovery [9], [31], [59], dictionary learning [51], [60], communications [18], [19], [61]–[64], ultrasound [65], and super resolution [66]–[68]. A recent review can be found in [7].

Design Outline: The application of deep unfolding to design a model-aided deep network is based on the following steps:

- 1) Identify an iterative optimization algorithm which is useful for the problem at hand. For instance, recovering a sparse vector from its noisy projections can be tackled using ISTA, unfolded into LISTA in [8].
- 2) Fix a number of iterations in the optimization algorithm.
- 3) Design the layers to imitate the free parameters of each iteration in a trainable fashion.
- 4) Train the overall resulting network end-to-end.

The selection of the free parameters to learn in the third step determines the resulting trainable architecture. One can set these parameters to be the hyperparameters of the iterative optimizer (such as step-size), thus leveraging data to automatically determine parameters typically selected by hand [53]. Alternatively, the architecture may be designed to learn the parameters of the objective optimized in each iteration, thus achieving a more abstract family of inference rules compared

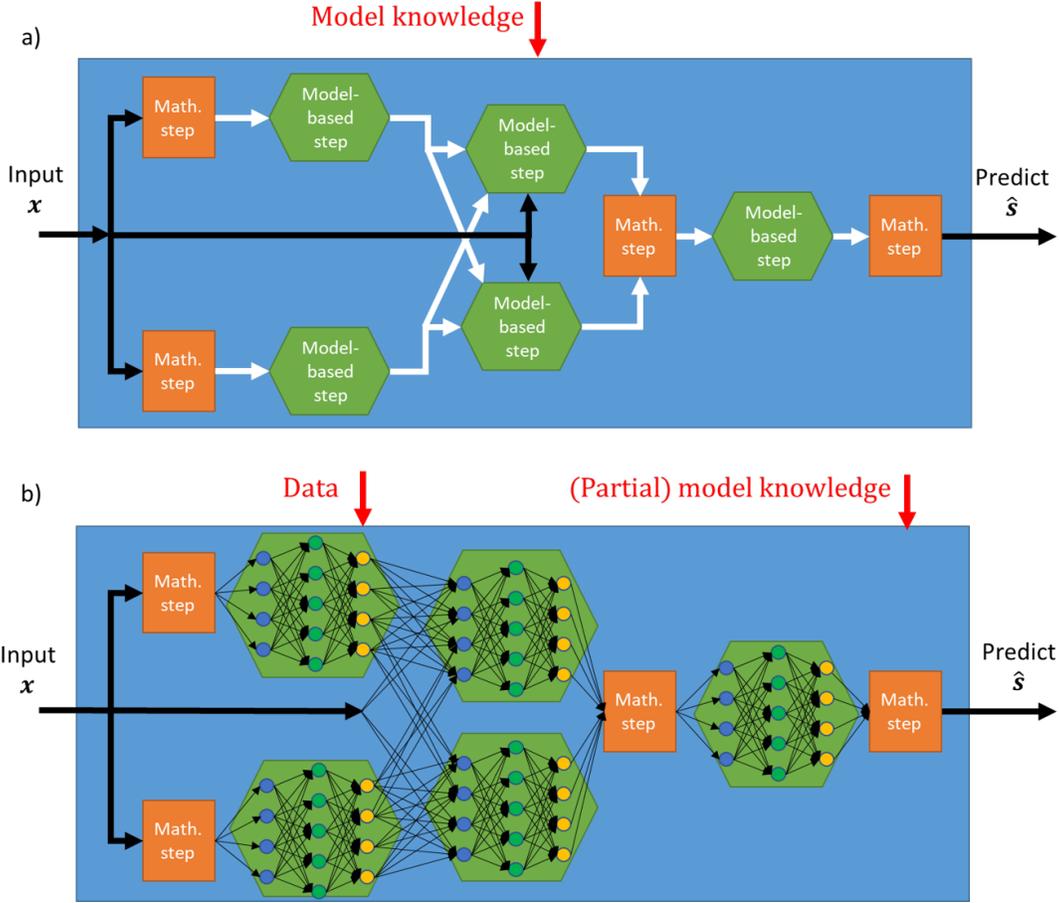


Fig. 3: Model-aided DNN illustration: *a)* a model-based algorithm comprised of a series of model-aware computations and generic mathematical steps; *b)* A DNN whose architecture and inter-connections are designed based on the model-based algorithm. Here, data can be used to train the overall architecture end-to-end, typically requiring the intermediate mathematical steps to be either differentiable or well-approximated by a differentiable mapping.

with the original iterative algorithm, or even convert the operation of each iteration into a trainable neural architecture. We next demonstrate how this rationale is translated into concrete architectures, using two examples: the first is the DetNet system of [18] which unfolds projected gradient descent optimization; the second is the unfolded dictionary learning for Poisson image denoising proposed in [51].

Example 1: Deep Unfolded Projected Gradient Descent: Projected gradient descent is a simple and common iterative algorithm for tackling constrained optimization. While the projected gradient descent method is quite generic and can be applied in a broad range of constrained optimization setup, in the following we focus on its implementation for symbol detection in linear memoryless multiple-input multiple-output (MIMO) Gaussian channels. In such cases, where the constraint follows from the discrete nature of digital communication symbols, the iterative projected gradient descent gives rise to the DetNet architecture proposed in [18] via deep unfolding.

a) System Model: Consider the problem of symbol detection in linear memoryless MIMO Gaussian channels. The task is to recover the K -dimensional vector \mathbf{s} from the $N \times 1$

observations \mathbf{x} , which are related via:

$$\mathbf{x} = \mathbf{H}\mathbf{s} + \mathbf{w}. \quad (8)$$

Here, \mathbf{H} is a known deterministic $N \times K$ channel matrix, and \mathbf{w} consists of N i.i.d Gaussian random variables (RVs). For our presentation we consider the case in which the entries of \mathbf{s} are symbols generated from a binary phase shift keying (BPSK) constellation in a uniform i.i.d. manner, i.e., $\mathcal{S} = \{\pm 1\}^K$. In this case, the MAP rule given an observation \mathbf{x} becomes the minimum distance estimate, given by

$$\hat{\mathbf{s}} = \arg \min_{\mathbf{s} \in \{\pm 1\}^K} \|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2. \quad (9)$$

b) Projected Gradient Descent: While directly solving (9) involves an exhaustive search over the 2^K possible symbol combinations, it can be tackled with affordable computational complexity using the iterative projected gradient descent algorithm. This method, whose derivation is detailed in Appendix A, is summarized as Algorithm 1, where $\mathcal{P}_{\mathcal{S}}(\cdot)$ denotes the projection operator into \mathcal{S} , which for BPSK constellations is the element-wise sign function.

c) Unfolded DetNet: DetNet unfolds the projected gradient descent iterations, repeated until convergence in Algo-

Algorithm 1: Projected gradient descent for system model (8)

Init: Fix step-size $\eta > 0$. Set initial guess \hat{s}_0
1 for $q = 0, 1, \dots$ **do**
2 | Update $\hat{s}_{q+1} = \mathcal{P}_S(\hat{s}_q - \eta \mathbf{H}^T \mathbf{x} + \eta \mathbf{H}^T \mathbf{H} \hat{s}_q)$.
3 end
Output: Estimate $\hat{s} = s_q$.

rithm 1, into a DNN, which learns to carry out this optimization procedure from data. To formulate DetNet, we first fix a number of iterations Q . Next, a DNN with Q layers is designed, where each layer imitates a single iteration of Algorithm 1 in a trainable manner.

Architecture: DetNet builds upon the observation that the update rule in Step 2 of Algorithm 1 consists of two stages: gradient descent computation, i.e., gradient step $\hat{s}_q - \eta \mathbf{H}^T \mathbf{x} + \eta \mathbf{H}^T \mathbf{H} \hat{s}_q$; and projection, namely, applying $\mathcal{P}_S(\cdot)$. Therefore, each unfolded iteration is represented as two sub-layers: The first sub-layer learns to compute the gradient descent stage by treating the step-size as a learned parameter and applying an FC layer with ReLU activation to the obtained value. For iteration index q , this results in

$$z_q = \text{ReLU}(\mathbf{W}_{1,q}((\mathbf{I} + \delta_{2,q} \mathbf{H}^T \mathbf{H})\hat{s}_{q-1} - \delta_{1,q} \mathbf{H}^T \mathbf{x}) + \mathbf{b}_{1,q})$$

in which $\{\mathbf{W}_{1,q}, \mathbf{b}_{1,q}, \delta_{1,q}, \delta_{2,q}\}$ are learnable parameters. The second sub-layer learns the projection operator by approximating the sign operation with a soft sign activation preceded by an FC layer, leading to

$$\hat{s}_q = \text{soft sign}(\mathbf{W}_{2,q} z_q + \mathbf{b}_{2,q}). \quad (10)$$

Here, the learnable parameters are $\{\mathbf{W}_{2,q}, \mathbf{b}_{2,q}\}$. The resulting deep network is depicted in Fig. 4, in which the output after Q iterations, denoted \hat{s}_Q , is used as the estimated symbol vector by taking the sign of each element.

Training: Let $\theta = \{(\mathbf{W}_{1,q}, \mathbf{W}_{2,q}, \mathbf{b}_{1,q}, \mathbf{b}_{2,q}, \delta_{1,q}, \delta_{2,q})\}_{q=1}^Q$ be the trainable parameters of DetNet¹. To tune θ , the overall network is trained end-to-end to minimize the empirical weighted ℓ_2 norm loss over its intermediate layers, given by

$$\mathcal{L}(\theta) = \frac{1}{n_t} \sum_{t=1}^{n_t} \sum_{q=1}^Q \log(q) \|s_t - \hat{s}_q(\mathbf{x}_t; \theta)\|^2 \quad (11)$$

where $\hat{s}_q(\mathbf{x}_t; \theta)$ is the output of the q th layer of DetNet with parameters θ and input \mathbf{x}_t . This loss measure accounts for the interpretable nature of the unfolded network, in which the output of each layer is a further refined estimate of s .

Quantitative Results: The experiments reported in [18] indicate that, when provided sufficient training examples, DetNet outperforms leading MIMO detection algorithms based on approximate message passing and semi-definite relaxation. It is also noted in [18] that the unfolded network requires an order of magnitude less layers compared to the number of

¹The formulation of DetNet in [18] includes an additional sub-layer in each iteration intended to further lift its input into higher dimensions and introduce additional trainable parameters, as well as reweighing of the outputs of subsequent layers. As these operations do not follow directly from unfolding projected gradient descent, they are not included in the description here.

iterations required by the model-based optimizer to converge. This gain is shown to be translated into reduced run-time during inference, particularly when processing batches of data in parallel. In particular, it is reported in [18, Tbl. 1] that DetNet successfully detects a batch of 1000 channel outputs in a 60×30 static MIMO channel at run-time which is three times faster than that required by approximate message passing to converge, and over 80 times faster than semi-definite relaxation.

Example 2: Deep Unfolded Dictionary Learning: DetNet exemplifies how deep unfolding can be used to realize rapid implementations of exhaustive optimization algorithms that typically require a very large amount of iterations to converge. However, DetNet requires full domain knowledge, i.e., it assumes the system model follows (8) and that the channel parameters \mathbf{H} are known. An additional benefit of deep unfolding is its ability to learn missing model parameters along with the overall optimization procedure, as we illustrate in the following example proposed in [51], which focuses on dictionary learning for Poisson image denoising. Similar examples where channel knowledge is not required in deep unfolding can be found in, e.g., [19], [57], [64]

a) *System Model:* Consider the problem of reconstructing an image $\mu \in \mathbb{R}^N$ from its noisy measurements $\mathbf{x} \in \mathbb{R}^N$. The image is corrupted by Poisson noise, namely, $p_{\mathbf{x}|\mu}$ is a multivariate Poisson distribution with mutually independent entries and mean μ . Furthermore, it is assumed that for the clean image μ , it holds that $\log(\mu)$ (taken element-wise) can be written as

$$\log(\mu) = \mathbf{H}s. \quad (12)$$

In (12), \mathbf{H} , referred to as the *dictionary*, is an unknown block-Toeplitz matrix (representing a convolutional dictionary), while s is an unknown sparse vector.

b) *Proximal Gradient Mapping:* The recovery of the clean image μ is tackled by alternating optimization [69]. In each iteration, one first recovers s for a fixed \mathbf{H} , after which s is set to be fixed and \mathbf{H} is estimated. The resulting iterative algorithm, whose detailed derivation is given in Appendix B, is summarized as Algorithm 2. Here, $\eta > 0$ is the step size; $\mathbf{1}$ is the all-ones vector; b is a threshold parameter; and \mathcal{T}_b is the soft-thresholding operator, also referred to as the shrinkage operator, applied element-wise and is given by $\mathcal{T}_b(x) = \text{sign}(x) \max\{|x| - b, 0\}$. Furthermore, the optimization variable \mathbf{H} in Step 2 is constrained to be block-Toeplitz.

c) *Deep Convolutional Exponential-Family Autoencoder:* The hybrid model-based/data-driven architecture entitled deep convolutional exponential-family autoencoder (DCEA) architecture proposed in [51] unfolds the proximal gradient iterations in Step 5 of Algorithm 2. By doing so, it avoids the need to learn the dictionary \mathbf{H} by alternating optimization, as it is implicitly learned from data in the training procedure of the unfolded network.

Architecture: DCEA treats the two-step convolutional sparse coding problem as an autoencoder, where the encoder computes the sparse vector s by unfolding Q proximal gradient iterations as in Step 5 of Algorithm 2. The decoder then converts \hat{s} produced by the encoder into a recovered clean

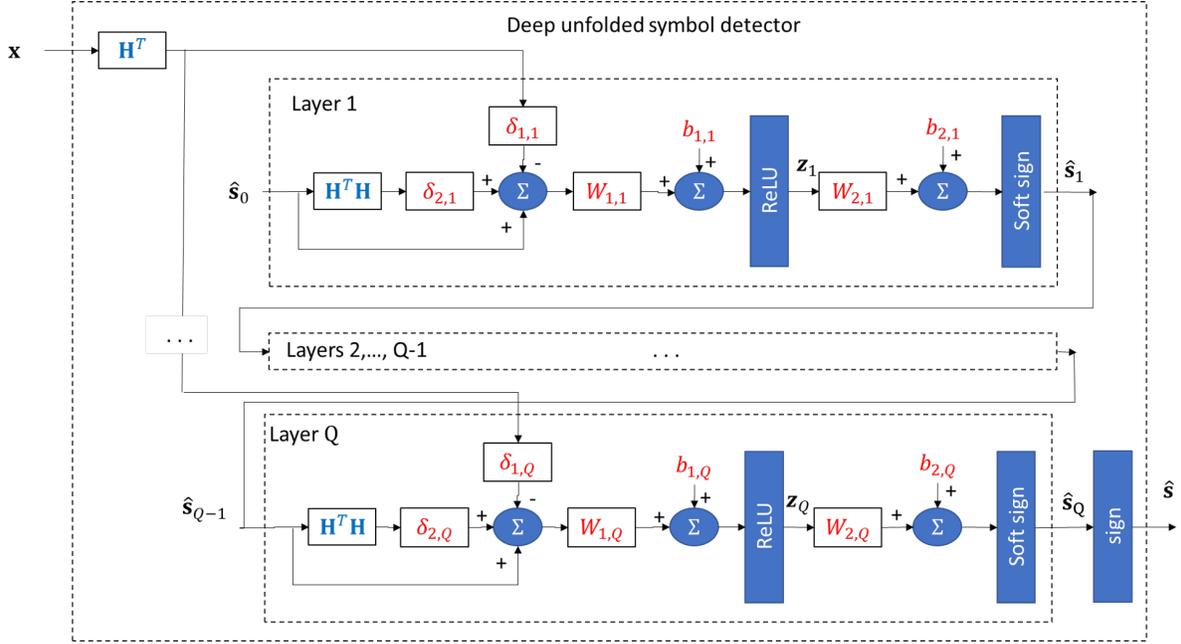


Fig. 4: DetNet illustration. Parameters in red fonts are learned in training, while those in blue fonts are externally provided.

Algorithm 2: Alternating image recovery and dictionary learning for system model (12)

Init: Fix step-size $\eta > 0$. Set initial guess s_0

- 1 **for** $l = 0, 1, \dots$ **do**
- 2 Update $\hat{H}_l = \arg \min_{\mathbf{H}} \mathbf{1}^T \exp(\mathbf{H} s_l) - x^T \mathbf{H} s_l$.
- 3 Set $\hat{s}_0 = s_l$.
- 4 **for** $q = 0, 1, \dots$ **do**
- 5 Update
- 6 $\hat{s}_{q+1} = \mathcal{T}_b(\hat{s}_q + \eta \mathbf{H}^T (x - \exp(\hat{H}_l \hat{s}_q)))$.
- 7 **end**
- 8 Set $s_{l+1} = \hat{s}_q$.
- 9 **end**

Output: Estimate clean image via $\hat{\mu} = \exp(\hat{H}_l s_l)$.

image $\hat{\mu}$.

In particular, [51] proposed two implementations of DCEA. The first, referred to as DCEA-C, directly implements Q proximal gradient iterations followed by the decoding step which computes $\hat{\mu}$, where both the encoder and the decoder use the same value of the dictionary matrix \mathbf{H} . This is replaced with a convolutional layer and is learned via end-to-end training along with the thresholding parameters, bypassing the need to explicitly recover the dictionary for each image, as in Step 2 of Algorithm 2. The second implementation, referred to DCEA-UC, decouples the convolution kernels of the encoder and the decoder, and lets the encoder carry out Q iterations of the form

$$\hat{s}_{q+1} = \mathcal{T}_b \left(\hat{s}_q + \eta \mathbf{W}_2^T (x - \exp(\mathbf{W}_1 \hat{s}_q)) \right). \quad (13)$$

Here, \mathbf{W}_1 and \mathbf{W}_2 are convolutional kernels which are not

constrained to be equal to \mathbf{H} used by the decoder². An illustration of the resulting architecture is depicted in Fig. 5.

Training: The parameters of DCEA are $\theta = \{\mathbf{H}, \mathbf{b}\}$ for DCEA-C, and $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{H}, \mathbf{b}\}$ for DCEA-UC. The vector $\mathbf{b} \in \mathbb{R}^C$ is comprised of the thresholding parameters used at each channel. When applied for Poisson image denoising, DCEA is trained in a supervised manner using the MSE loss, namely, a set of n_t clean images $\{\mu_t\}_{t=1}^{n_t}$ are used along with their Poisson noisy version $\{x_t\}_{t=1}^{n_t}$. By letting $f_\theta(\cdot)$ denote the resulting mapping of the unfolded network, the loss function is formulated as

$$\mathcal{L}(\theta) = \frac{1}{n_t} \sum_{t=1}^{n_t} \|\mu_t - f_\theta(x_t)\|^2. \quad (14)$$

Quantitative Results: The experimental results reported in [51] evaluated the ability of the unfolded DCEA-C and DCEA-UC in recovering images corrupted with different levels of Poisson noise. An example of an image denoised by the unfolded system is depicted in Fig. 6. It was noted in [51] that the proposed approach allows to achieve similar and even improved results to those of purely data-driven techniques based on black-box CNNs [70]. However, the fact that the denoising system is obtained by unfolding the model-based optimizer in Step 5 of Algorithm 2 allows this performance to be achieved while utilizing 3%–10% of the overall number of trainable parameters as those used by the conventional CNN.

Discussion: Deep unfolding incorporates model-based domain knowledge to obtain a dedicated DNN design which resembles an iterative optimization algorithm. Compared to conventional DNNs, unfolded networks are typically inter-

²The architecture proposed in [51] is applicable for various exponential-family noise signals. Particularly for Poisson noise, an additional exponential linear unit was applied to $x - \exp(\mathbf{W}_1 \hat{s}_q)$ which was empirically shown to improve the convergence properties of the network.

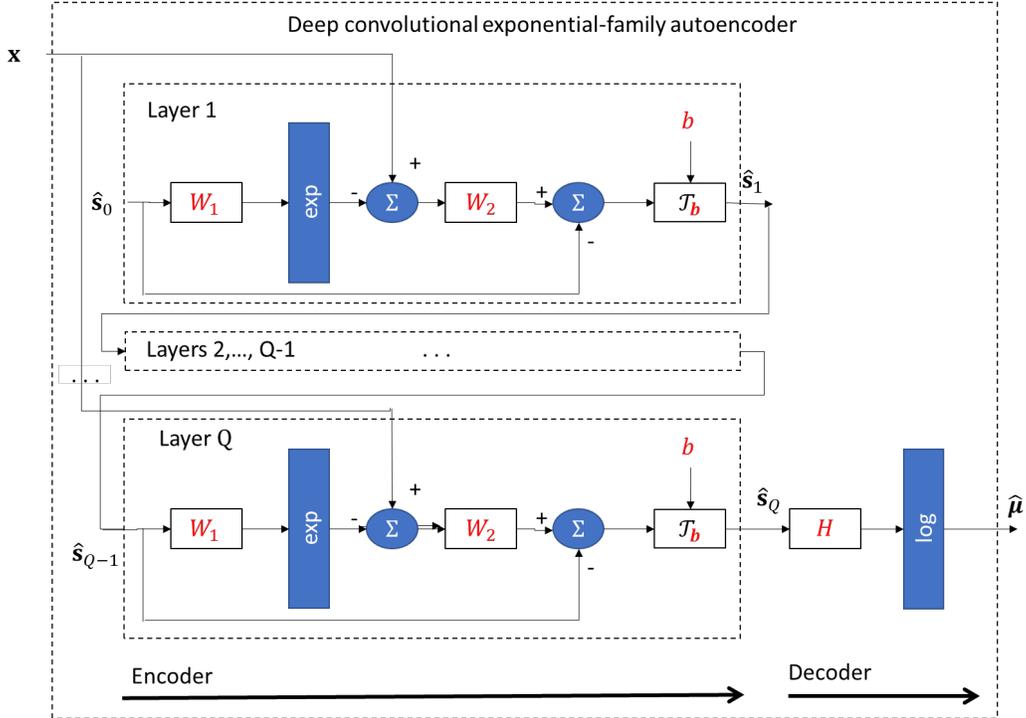


Fig. 5: DCEA illustration. Parameters in red fonts are learned in training, while those in blue fonts are externally provided.

pretable, and tend to have a smaller number of parameters, and can thus be trained more quickly [7], [61]. A key advantage of deep unfolding over model-based optimization is in inference speed. For instance, unfolding projected gradient descent iterations into DetNet allows to infer with much fewer layers compared to the number of iterations required by the model-based algorithm to converge. Similar observations have been made in various unfolded algorithms [58], [66].

One of the key properties of unfolded networks is their reliance on knowledge of the model describing the setup (though not necessarily on its parameters). For example, one must know that the image is corrupted by Poisson noise to formulate the iterative procedure in Algorithm 2 unfolded into DCEA, or that the observations obey a linear Gaussian model to unfold the projected gradient descent iterations into DetNet. However, the parameters of this model, e.g., the matrix H in (8) and (12), can be either provided based on domain knowledge, as done in DetNet, or alternatively, learned in the training procedure, as carried out by DCEA. The model-awareness of deep unfolding has its advantages and drawbacks. When the model is accurately known, deep unfolding essentially incorporates it into the DNN architecture, as opposed to conventional black-box DNNs which must learn it from data. However, this approach does not exploit the model-agnostic nature of deep learning, and thus may lead to degraded performance when the true relationship between the measurements and the desired quantities deviates from the model assumed in design. Nonetheless, training an unfolded network designed with a mismatched model using data corresponding to the true underlying scenario typically yields more accurate inference compared to the model-based iterative algorithm with the same model-mismatch, as the unfolded

network can learn to compensate for this mismatch [64].

B. Neural Building Blocks

Neural building blocks is an alternative approach to design model-aided networks, which can be treated as a generalization of deep unfolding. It is based on representing a model-based algorithm, or alternatively prior knowledge of an underlying statistical model, as an interconnection of distinct building blocks. Neural building blocks implement a DNN comprised of multiple sub-networks. Each module learns to carry out the specific computations of the different building blocks constituting the model-based algorithm, as done in [16], [71]–[73], or to capture a known statistical relationship, as in [74].

Neural building blocks are designed for scenarios which are tackled using algorithms with flow diagram representations, that can be captured as a sequential and parallel interconnection of building blocks. In particular, deep unfolding can be obtained as a special case of neural building blocks, where the original algorithm is an iterative optimizer, such that the building blocks are interconnected in a sequential fashion, and implemented using a single layer. However, the generalization of neural building blocks compared to deep unfolding is not encapsulated merely in its ability to implement non-sequential interconnections between algorithmic building blocks in a learned fashion, but also in the identification of the specific task of each block, as well as the ability to convert known statistical relationships such as causal graphs into dedicated DNN architectures.

Design Outline: The application of neural building blocks to design a model-aided deep network is based on the following steps:

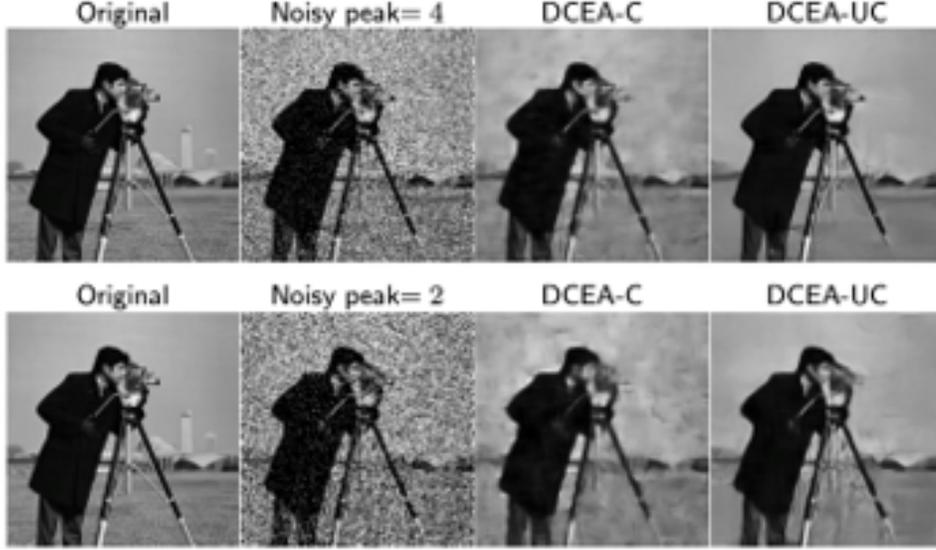


Fig. 6: Illustration of an image corrupted by different levels of Poisson noise and the resulting denoised images produced by the unfolded DCEA-C and DCEA-UC. Figure reproduced from [51] with authors' permission.

- 1) Identify an algorithm or a flow-chart structure which is useful for the problem at hand, and can be decomposed into multiple building blocks.
- 2) Identify which of these building blocks should be learned from data, and what is their concrete task.
- 3) Design a dedicated neural network for each building block capable of learning to carry out its specific task.
- 4) Train the overall resulting network, either in an end-to-end fashion or by training each building block network individually.

We next demonstrate how one can design a model-aided network comprised of neural building blocks. Our example focuses on symbol detection in flat MIMO channels, where we consider the data-driven implementation of the iterative soft interference cancellation (SIC) scheme of [75], which is the DeepSIC algorithm proposed in [16].

Example 3: DeepSIC for MIMO Detection: Iterative SIC [75] is a MIMO detection method suitable for linear Gaussian channels, i.e., the same channel models as that described in the example of DetNet in Subsection IV-A. DeepSIC is a hybrid model-based/data-driven implementation of the iterative SIC scheme [16]. However, unlike its model-based counterpart, and alternative deep MIMO receivers [18], [19], [61], DeepSIC is not particularly tailored for linear Gaussian channels, and can be utilized in various flat MIMO channels. We formulate DeepSIC by first reviewing the model-based iterative SIC, and present DeepSIC as its data-driven implementation.

a) Iterative Soft Interference Cancellation: The iterative SIC algorithm proposed in [75] is a MIMO detection method that combines multi-stage interference cancellation with soft decisions. The detector operates iteratively, where in each iteration, an estimate of the conditional probability mass function (PMF) of s_k , which is the k th entry of \mathbf{s} , given the observed \mathbf{x} , is generated for every symbol $k \in \{1, 2, \dots, K\} := \mathcal{K}$. Each PMF, which is an $|\mathcal{S}| \times 1$ vector denoted $\hat{\mathbf{p}}_k^{(q)}$ at the

q th iteration, is computed using the corresponding estimates of the interfering symbols $\{s_l\}_{l \neq k}$ obtained in the previous iteration. Iteratively repeating this procedure refines the PMF estimates, allowing to accurately recover each symbol from the output of the last iteration. This iterative procedure is illustrated in Fig. 7(a) and summarized as Algorithm 3, whose derivation is detailed in Appendix C. Algorithm 3 is detailed for linear Gaussian models as in (8), assuming that the noise \mathbf{w} has variance σ_w^2 . We use \mathbf{h}_l to denote the l th column of \mathbf{H} , while $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$.

Algorithm 3: Iterative SIC for system model (8)

Init: Set initial PMFs guess $\{\hat{\mathbf{p}}_k^{(0)}\}_{k=1}^K$

- 1 **for** $q = 0, 1, \dots$ **do**
- 2 For each $k \in \mathcal{K}$, compute expected values $e_k^{(q)}$ and variance $v_k^{(q)}$ from $\hat{\mathbf{p}}_k^{(q)}$.
- 3 *Interference cancellation:* For each $k \in \mathcal{K}$ compute

$$\mathbf{z}_k^{(q+1)} = \mathbf{x} - \sum_{l \neq k} \mathbf{h}_l e_l^{(q)}.$$
- 4 *Soft decoding:* For each $k \in \mathcal{K}$, estimate $\hat{\mathbf{p}}_k^{(q+1)}$ as the PMF of s_k given $\mathbf{z}_k^{(q+1)}$, assuming that

$$\mathbf{z}_k^{(q+1)} | s_k \sim \mathcal{N}(\mathbf{h}_k s_k, \sigma_w^2 \mathbf{I}_K + \sum_{l \neq k} v_l^{(q)} \mathbf{h}_l \mathbf{h}_l^T).$$
- 5 **end**

Output: Estimate $\hat{\mathbf{s}}$ by setting each \hat{s}_k as the symbol maximizing the estimated PMF $\hat{\mathbf{p}}_k^{(q)}$.

b) DeepSIC: Iterative SIC is specifically designed for linear channels of the form (8). In particular, the interference cancellation Step 3 of Algorithm 3 requires the contribu-

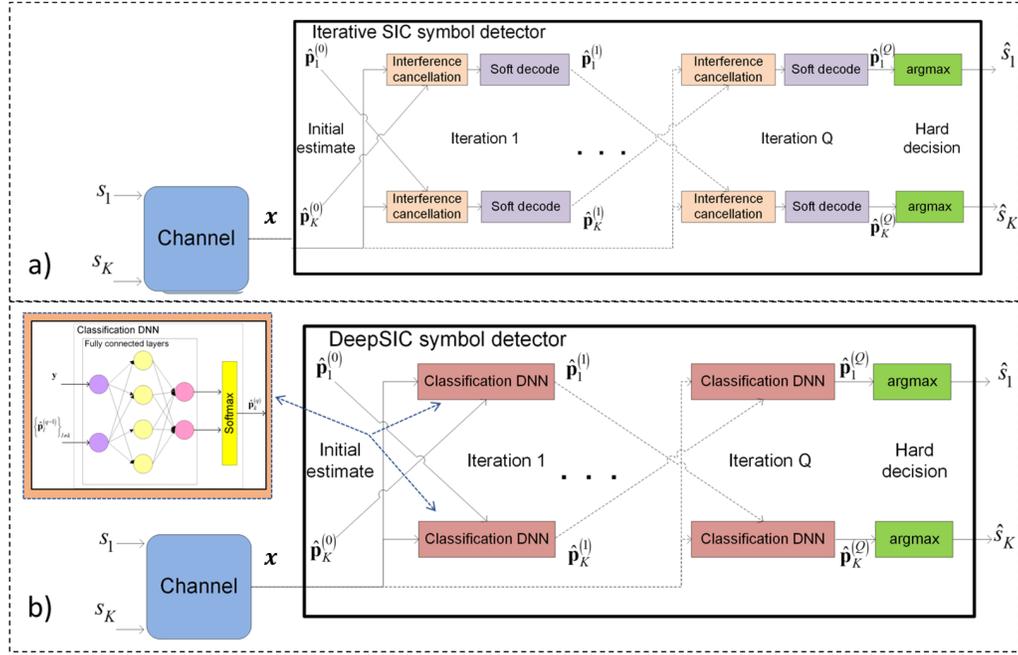


Fig. 7: Iterative SIC illustration: a) model-based method; b) DeepSIC.

tion of the interfering symbols to be additive. Furthermore, it requires accurate complete knowledge of the underlying statistical model, i.e., of (8). DeepSIC proposed in [16] learns to implement the iterative SIC from data as a set of neural building blocks, thus circumventing these limitations of its model-based counterpart.

Architecture: The iterative SIC algorithm can be viewed as a set of interconnected basic building blocks, each implementing the two stages of interference cancellation and soft decoding, as illustrated in Fig. 7(a). While the block diagram in Fig. 7(a) is ignorant of the underlying channel model, the basic building blocks are model-dependent. Although each of these basic building blocks consists of two sequential procedures which are completely channel-model-based, the purpose of these computations is to carry out a classification task. In particular, the k th building block of the q th iteration, $k \in \mathcal{K}$, produces $\hat{p}_k^{(q)}$, which is an estimate of the conditional PMF of s_k given \mathbf{x} based on $\{\hat{p}_l^{(q-1)}\}_{l \neq k}$. Such computations are naturally implemented by classification DNNs, e.g., FC networks with a softmax output layer. Embedding these conditional PMF computations into the iterative SIC block diagram in Fig. 7(a) yields the overall receiver architecture depicted in Fig. 7(b).

A major advantage of using classification DNNs as the basic building blocks in Fig. 7(b) stems from their ability to accurately compute conditional distributions in complex non-linear setups without requiring a-priori knowledge of the channel model and its parameters. Consequently, when these building blocks are trained to properly implement their classification task, the receiver essentially realizes iterative SIC for arbitrary channel models in a data-driven fashion.

Training: In order for DeepSIC to reliably implement symbol detection, its building block classification DNNs must be properly trained. Two possible training approaches are

considered based on a labeled set of n_t samples $\{(s_t, \mathbf{x}_t)\}_{t=1}^{n_t}$:

(i) **End-to-end training:** The first approach jointly trains the entire network, i.e., all the building block DNNs. Since the output of the deep network is the set of PMFs $\{\hat{p}_k^{(Q)}\}_{k=1}^K$, the sum cross entropy loss is used. Let θ be the network parameters, and $\hat{p}_k^{(Q)}(\mathbf{x}, \alpha; \theta)$ be the entry of $\hat{p}_k^{(Q)}$ corresponding to $s_k = \alpha$ when the input to the network parameterized by θ is \mathbf{x} . The sum cross entropy loss is

$$\mathcal{L}(\theta) = \frac{1}{n_t} \sum_{t=1}^{n_t} \sum_{k=1}^K -\log \hat{p}_k^{(Q)}(\mathbf{x}_t, (s_t)_k; \theta). \quad (15)$$

Training the interconnection of DNNs in Fig. 7(b) end-to-end based on (15) jointly updates the coefficients of all the $K \cdot Q$ building block DNNs. For a large number of symbols, i.e., large K , training so many parameters simultaneously is expected to require a large labeled set.

(ii) **Sequential training:** The fact that DeepSIC is implemented as an interconnection of neural building blocks, implies that each block can be trained with a reduced number of training samples. Specifically, the goal of each building block DNN does not depend on the iteration index: The k th building block of the q th iteration outputs a soft estimate of s_k for each iteration q . Therefore, each building block DNN can be trained individually, by minimizing the conventional cross entropy loss. To formulate this objective, let $\theta_k^{(q)}$ represent the parameters of the k th DNN at iteration q , and write $\hat{p}_k^{(q)}(\mathbf{x}, \{\hat{p}_l^{(q-1)}\}_{l \neq k}, \alpha; \theta_k^{(q)})$ as the entry of $\hat{p}_k^{(q)}$ corresponding to $s_k = \alpha$ when the DNN parameters are $\theta_k^{(q)}$ and its inputs are \mathbf{x} and $\{\hat{p}_l^{(q-1)}\}_{l \neq k}$. The cross entropy loss is

$$\mathcal{L}(\theta_k^{(q)}) = \frac{-1}{n_t} \sum_{t=1}^{n_t} \log \hat{p}_k^{(q)}(\tilde{\mathbf{x}}_t, \{\hat{p}_{t,l}^{(q-1)}\}_{l \neq k}, (\tilde{s}_t)_k; \theta_k^{(q)}) \quad (16)$$

where $\{\hat{\mathcal{P}}_{t,l}^{(q-1)}\}$ represent the estimated PMFs associated with \mathbf{x}_i computed at the previous iteration. The problem with training each DNN individually is that the soft estimates $\{\hat{\mathcal{P}}_{t,l}^{(q-1)}\}$ are not provided as part of the training set. This challenge can be tackled by training the DNNs corresponding to each layer in a sequential manner, where for each layer the outputs of the trained previous iterations are used as the soft estimates fed as training samples.

Quantitative Results: Two experimental studies of DeepSIC taken from [16] are depicted in Fig. 8. These results compare the symbol error rate (SER) achieved by DeepSIC which learns to carry out $Q = 5$ SIC iterations from $n_t = 5000$ labeled samples. In particular, Fig. 8a considers a Gaussian channel of the form (8) with $K = N = 32$, resulting in MAP detection being computationally infeasible, and compares DeepSIC to the model-based iterative SIC as well as the data-driven DetNet [18]. Fig. 8b considers a Poisson channel, where \mathbf{x} is related to \mathbf{s} via a multivariate Poisson distribution, for which schemes requiring a linear Gaussian model such as the iterative SIC algorithm are not suitable. The ability to use DNNs as neural building blocks to carry out their model-based algorithmic counterparts in a robust and model-agnostic fashion is demonstrated in Fig. 8. In particular, it is demonstrated that DeepSIC approaches the SER values of the iterative SIC algorithm in linear Gaussian channels, while notably outperforming it in the presence of model mismatch, as well as when applied in non-Gaussian setups. It is also observed in Fig. 8a that the resulting architecture of DeepSIC can be trained with smaller data sets compared to alternative data-driven receivers, such as DetNet.

Discussion: The main rationale in designing DNNs as interconnected neural building blocks is to facilitate learned inference by preserving the structured operation of a model-based algorithm applicable for the problem at hand given full domain knowledge. As discussed earlier, this approach can be treated as an extension of deep unfolding, allowing to exploit additional structures beyond a sequential iterative operation. The generalization of deep unfolding into a set of learned building blocks opens additional possibilities in designing model-aided networks.

First, the treatment of the model-based algorithm as a set of building blocks with concrete tasks allows a DNN architecture designed to comply with this structure not only to learn to carry out the original model-based method from data, but also to robustify it and enable its application in diverse new scenarios. This follows since the block diagram structure of the algorithm may be ignorant of the specific underlying statistical model, and only rely upon a set of generic assumptions, e.g., that the entries of the desired vector \mathbf{s} are mutually independent. Consequently, replacing these building blocks with dedicated DNNs allows to exploit their model-agnostic nature, and thus the original algorithm can now be learned to be carried out in complex environments. For instance, DeepSIC can be applied to non-linear channels, owing to the implementation of the building blocks of the iterative SIC algorithm using generic DNNs, while the model-based algorithm is limited to setups of the form (8).

In addition, the division into building blocks gives rise

to the possibility to train each block separately. The main advantage in doing so is that a smaller training set is expected to be required, though in the horizon of a sufficiently large amount of training, end-to-end training is likely to yield a more accurate model as its parameters are jointly optimized. For example, in DeepSIC, sequential training uses the n_t input-output pairs to train each DNN individually. Compared to the end-to-end training that utilizes the training samples to learn the complete set of parameters, which can be quite large, sequential training uses the same data set to learn a significantly smaller number of parameters, reduced by a factor of $K \cdot Q$, multiple times. This indicates that the ability to train the blocks individually is expected to require much fewer training samples, at the cost of a longer learning procedure for a given training set, due to its sequential operation, and possible performance degradation as the building blocks are not jointly trained. In addition, training each block separately facilitates adding and removing blocks, when such operations are required in order to adapt the inference rule.

V. DNN-AIDED INFERENCE

DNN-aided inference is a family of model-based deep learning algorithms in which DNNs are incorporated into model-based methods. As opposed to model-aided networks discussed in Section IV, where the resultant system is a deep network whose architecture imitates the operation of a model-based algorithm, here inference is carried out using a traditional model-based method, while some of the intermediate computations are augmented by DNNs. The main motivation of DNN-aided inference is to exploit the established benefits of model-based methods, in terms of performance, complexity, and suitability for the problem at hand. Deep learning is incorporated to mitigate sensitivity to inaccurate model knowledge, facilitate operation in complex environments, and enable application in new domains. An illustration of a DNN-aided inference system is depicted in Fig. 9.

DNN-aided inference is particularly suitable for scenarios in which one only has access to partial domain knowledge. In such cases, the available domain knowledge dictates the algorithm utilized, while the part that is not available or is too complex to model analytically is tackled using deep learning. We divide our description of DNN-aided inference schemes into three main families of methods: The first, referred to as *structure-agnostic DNN-aided inference* detailed in Subsection V-A, utilizes deep learning to capture structures in the underlying data distribution, e.g., to represent the domain of natural images. This DNN is then utilized by model-based methods, allowing them to operate in a manner which is invariant to these structures. The family of *structure-oriented DNN-aided inference* schemes, detailed in Subsection V-B, utilizes model-based algorithms to exploit a known tractable statistical structure, such as an underlying Markovian behavior of the considered signals. In such methods, deep learning is incorporated into the structure-aware algorithm, thereby capturing the remaining portions of the underlying model as well as mitigating sensitivity to uncertainty. Finally, in Subsection V-C, we discuss *neural augmentation* methods,

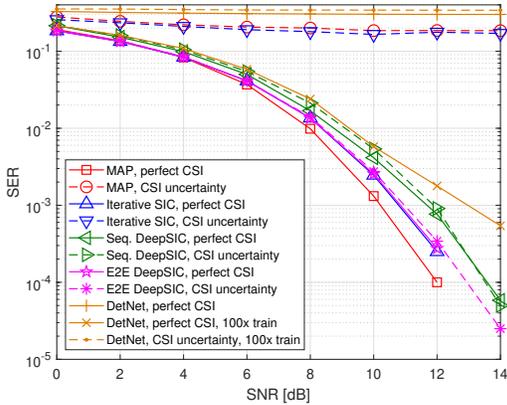
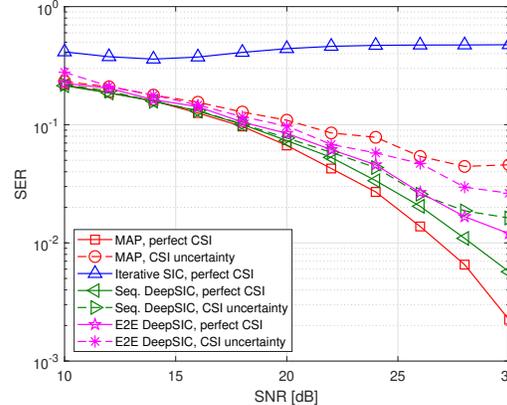
(a) 32×32 Gaussian channel.(b) 4×4 Poisson channel.

Fig. 8: Experimental results from [16] of DeepSIC compared to the model-based iterative SIC, the model-based MAP (when feasible) and the data-driven DetNet of [18] (when applicable). *Perfect CSI* implies that the system is trained and tested using samples from the same channel, while under *CSI uncertainty* they are trained using samples from a set of different channels.

which are tailored to robustify model-based processing in the presence of inaccurate knowledge of the parameters of the underlying model. Here, inference is carried out using a model-based algorithm based on its available domain knowledge, while a deep learning system operating in parallel is utilized to compensate for errors induced by model inaccuracy. Our description of these methodologies in Subsections V-A-V-C follows the same systematic form used in Section IV, where each approach is detailed by a high-level description; design outline; one or two concrete examples; and a summarizing discussion.

A. Structure-Agnostic DNN-Aided Inference

The first family of DNN-aided inference utilizes deep learning to implicitly learn structures and statistical properties of the signal of interest, in a manner that is amenable to model-based optimization. These inference systems are particularly relevant for various inverse problems in signal processing, including denoising, sparse recovery, deconvolution, and super resolution [76]. Tackling such problems typically involves imposing some structure on the signal domain. This prior knowledge is then incorporated into a model-based optimization procedure, such as alternating direction method of multipliers (ADMM) [77], fast iterative shrinkage and thresholding algorithm [78], and primal-dual splitting [79], which recover the desired signal with provable performance guarantees.

Traditionally, the prior knowledge encapsulating the structure and properties of the underlying signal is represented by a handcrafted regularization term or constraint incorporated into the optimization objective. For example, a common model-based strategy used in various inverse problems is to impose sparsity in some given dictionary, which facilitates CS-based optimization. Deep learning brings forth the possibility to avoid such explicit constraint, thereby mitigating the detrimental effects of crude, handcrafted approximation of the true underlying structure of the signal, while enabling optimization with implicit data-driven regularization. This can

be implemented by incorporating deep denoisers as learned proximal mappings in iterative optimization, as carried out by plug-and-play networks³ [13], [14], [80]–[85]. DNN-based priors can also be used to enable, e.g., CS beyond the domain of sparse signals [10], [11].

Design Outline: Designing structure-agnostic DNN-aided systems can be carried out via the following steps:

- 1) Identify a suitable optimization procedure, given the domain knowledge for the signal of interest.
- 2) The specific parts of the optimization procedure which rely on complicated and possibly analytically intractable domain knowledge are replaced with a DNN.
- 3) The integrated data-driven module can either be trained separately from the inference system, possibly in an unsupervised manner as in [10], or alternatively, the complete inference system is trained end-to-end [12].

We next demonstrate how these steps are carried out in two examples: CS over complicated domains, where deep generative networks are used for capturing the signal domain [10]; and plug-and-play networks, which augment ADMM with a DNN to bypass the need to express a proximal mapping.

Example 4: Compressed Sensing using Generative Models: CS refers to the task of recovering some unknown signal from (possibly noisy) lower-dimensional observations. The mapping that transforms the input signal into the observations is known as the *forward operator*. In our example, we focus on the setting where the forward operator is a particular linear function that is known at the time of signal recovery.

The main challenge in CS is that there could be (potentially infinitely) many signals that agree with the given observations. Since such a problem is underdetermined, it is necessary to make some sort of structural assumptions on the unknown

³The term *plug-and-play* typically refers to the usage of an image denoiser as proximal mapping in regularized optimization [80]. As this approach can also utilize model-based denoisers, we use the term *plug-and-play networks* for such methods with DNN-based denoisers.

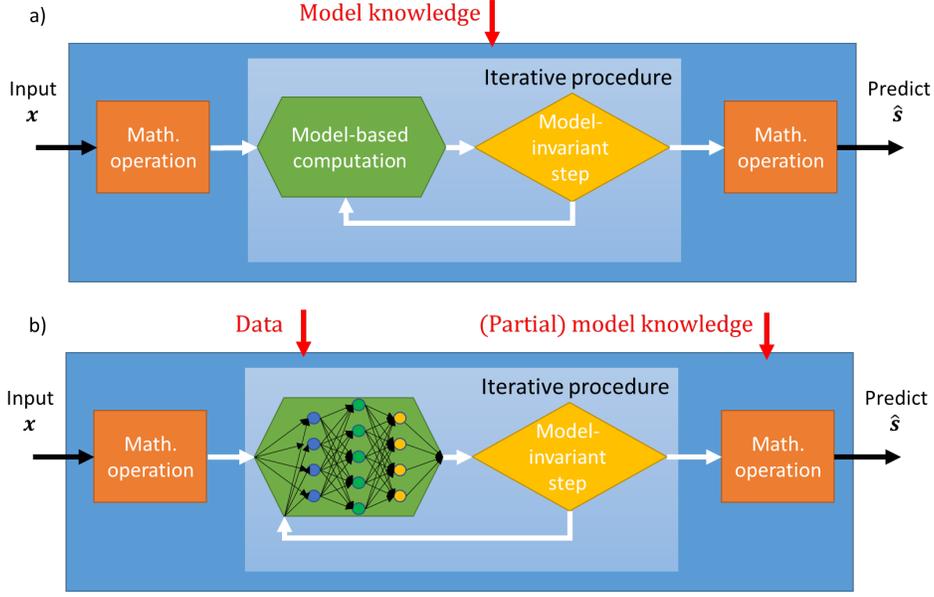


Fig. 9: DNN-aided inference illustration: *a)* a model-based algorithm comprised of multiple iterations with intermediate model-based computations; *b)* A data-driven implementation of the algorithm, where the specific model-based computations are replaced with dedicated learned deep models. Here, one can possibly use data to train the internal DNNs individually, or to train the overall inference mapping end-to-end as a discriminative learning model [54], [55], typically requiring the intermediate mathematical steps to be either differentiable or well-approximated by a differentiable mapping.

signal to identify the most plausible one. A classic assumption is that the signal is *sparse* in some known basis.

a) System Model: We consider the problem of noisy CS, where we wish to reconstruct an unknown N -dimensional signal s^* from the following observations

$$\mathbf{x} = \mathbf{H}\mathbf{s}^* + \mathbf{w} \quad (17)$$

where \mathbf{H} is an $M \times N$ matrix, modeled as random Gaussian matrix with entries $\mathbf{H}_{ij} \sim \mathcal{N}(0, 1/M)$, with $M < N$, and \mathbf{w} is an $M \times 1$ noise vector.

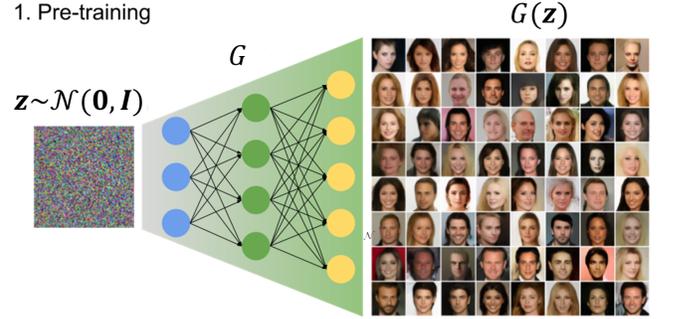
b) Sparsity-based CS: We next focus on a particular technique as a representative example of model-based CS. We rely here on the assumption that s^* is sparse, and seek to recover s^* from \mathbf{x} by solving the ℓ_1 relaxed LASSO objective

$$\mathcal{L}_{\text{LASSO}}(\mathbf{s}) \triangleq \|\mathbf{H}\mathbf{s} - \mathbf{x}\|_2^2 + \lambda\|\mathbf{s}\|_1. \quad (18)$$

While the derivation above assumes that s^* is sparse, the LASSO objective can also be used when s^* is sparse in some dictionary \mathbf{B} , e.g., in the wavelet domain, and the detailed formulation is given in Appendix D.

c) DNN-Aided Compressed Sensing: In a data-driven approach, we aim to replace the sparsity prior with a learned DNN. The following description is based on [10], which proposed to use a deep generative prior. Specifically, we replace the explicit sparsity assumption on true signal s^* , with a requirement that it lies in the range of a pre-trained generator network $G: \mathbb{R}^l \rightarrow \mathbb{R}^N$ (e.g., the generator network of a GAN).

Pre-training: To implement deep generative priors, one first has to train a generative network G to map a latent vector \mathbf{z} into a signal \mathbf{s} which lies in the domain of interest. A major advantage of employing a DNN-based prior in this setting is



$$\hat{\mathbf{z}} = \operatorname{argmin}_{\mathbf{z}} \|\mathbf{H}\mathbf{G}(\mathbf{z}) - \mathbf{x}\|_2^2 + \lambda\|\mathbf{z}\|_2^2; \quad \hat{\mathbf{s}} = \mathbf{G}(\hat{\mathbf{z}})$$

Fig. 10: High-level overview of CS with a DNN-based prior. The generator network G is pre-trained to map Gaussian latent variables to plausible signals in the target domain. Then signal recovery is done by finding a point in the range of G that minimizes reconstruction error via gradient-based optimization over the latent variable.

that generator networks are agnostic to how they are used and can be *pre-trained* and reused for multiple downstream tasks. The pre-training thus follows the standard unsupervised training procedure, as discussed, e.g., in Subsection III-B for GANs. In particular, the work [10] trained a Deep convolutional GAN [86] on the CelebA data set [87], to represent 64×64 color images of human faces, as well as a variational autoencoder (VAE) [88] for representing handwritten digits in 28×28 grayscale form based on the MNIST data set [89].

Architecture: Once a pre-trained generator network G is available, it can be incorporated as an alternative prior for the

inverse model in (17). The key intuition behind this approach is that the range of G should only contain *plausible* signals. Thus one can replace the handcrafted sparsity prior with a data-driven DNN prior G by constraining our signal recovery to the range of G .

One natural way to impose this constraint is to perform the optimization in the latent space to find z whose image $G(z)$ matches the observations. This is carried out by minimizing the following loss function in the latent space of G :

$$\mathcal{L}(z) = \|HG(z) - x\|_2^2. \quad (19)$$

Because the above loss function involves a highly non-convex function G , there is no closed-form solution or guarantee for this optimization problem. However the loss function is differentiable with respect to z , so it can be tackled using conventional gradient-based optimization techniques. Once a suitable latent z is found, the signal is recovered as $G(z)$.

In practice, [10] reports that incorporating an ℓ_2 regularizer on z helps. This is possibly due to the Gaussian prior assumption for the latent variable, as the density of z is proportional to $\exp(-\|z\|_2^2)$. Therefore, minimizing $\|z\|_2^2$ is equivalent to maximizing the density of z under the Gaussian prior. This has the effect of avoiding images that are extremely unlikely under the Gaussian prior even if it matches the observation well. The final loss includes this regularization term:

$$\mathcal{L}_{CS}(z) = \|HG(z) - x\|_2^2 + \lambda\|z\|_2^2 \quad (20)$$

where λ is a regularization coefficient.

In summary, DNN-aided CS replaces the constrained optimization over the complex input signal with tractable optimization over the latent variable z , which follows a known simple distribution. This is achieved using a pre-trained DNN-based prior G to map it into the domain of interest. Inference is performed by minimizing \mathcal{L}_{CS} in the latent space of G . An illustration of the system operation is depicted in Fig. 10.

Quantitative Results: To showcase the efficacy of the data-driven prior at capturing complex high-dimensional signal domains, we present the evaluation of its performance as reported in [10]. The baseline model used for comparison is based on directly solving the LASSO loss (18). For CelebA, we formulate the LASSO objective in the discrete cosine transform (DCT) and the wavelet (WVT) basis, and minimize it via coordinate descent.

The first task is the recovery of handwritten digit images from low-dimensional projections corrupted by additive Gaussian noise. The reconstruction error is evaluated for various numbers of observations M . The results are depicted in Fig. 11. We clearly see the benefit of using a data-driven deep prior in Fig. 11, where the VAE-based methods (labeled VAE and VAE+REG) show notable performance gain compared to the sparsity prior for small number of measurements. Implicitly imposing a sparsity prior via the LASSO objective outperforms the deep generative priors as the number of observations approaches the dimension of the signal. One explanation for this behavior is that the pre-trained generator G does not perfectly model the MNIST digit distribution and may not actually contain the ground truth signal in its range.

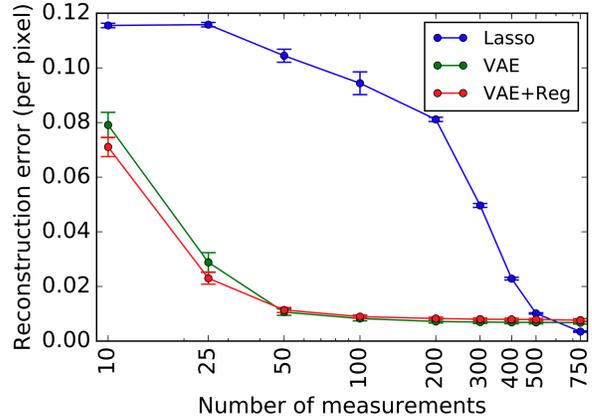


Fig. 11: Experimental result for noisy CS on the MNIST data set. Reproduced from [10] with the authors' permission.



Fig. 12: Visualization of the recovered signals from noisy CS on the CelebA data set. Reproduced from [10] with the authors' permission.

As such, its reconstruction error may never be exactly zero regardless of how many observations are given. The LASSO objective, on the other hand, does not suffer from this issue and is able to make use of the extra observations available.

The ability of deep generative priors to facilitate recovery from compressed measurements is also observed in Fig. 12, which qualitatively evaluates GAN-based CS recovery on the CelebA data set. This experiment uses $M = 500$ noisy measurements (out of $N = 12288$ total dimensions). As shown in Fig. 12, in this low-measurement regime, the data-driven prior again provides much more reasonable samples.

Example 5: Plug-and-Play Networks for Image Restoration:

The above example of DNN-aided CS allows to carry out regularized optimization over complex domains while using deep learning to avoid regularizing explicitly. This is achieved via deep priors, where the domain of interest is captured by a generative network. An alternative strategy, referred to as plug-

and-play networks, applies deep denoisers as learned proximal mappings. Namely, instead of using DNNs to evaluate the regularized objective as in [10], one uses DNNs to carry out an optimization procedure which relies on this objective without having to express the desired signal domain. In the following we exemplify the application of plug-and-play networks for image restoration using ADMM optimization [80].

a) System Model: We again consider the linear inverse problem formulated in (17) in which the additive noise w is comprised of i.i.d. mutually independent Gaussian entries with zero mean and variance σ_w^2 . However, unlike the setup considered in the previous example, the sensing matrix \mathbf{H} is not assumed to be random, and can be any fixed matrix dictated by the underlying setup.

The recovery of the desired signal s can be obtained via the MAP rule, which is given by

$$\begin{aligned} \hat{s} &= \arg \min_s -\log p(s|\mathbf{x}) \\ &= \arg \min_s -\log p(\mathbf{x}|s) - \log p(s) \\ &= \arg \min_s \frac{1}{2} \|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2 + \phi(s) \end{aligned} \quad (21)$$

where $\phi(s)$ is a regularization term which equals $-\sigma_w^2 \log p(s)$, with possibly some additive constant that does not affect the minimization in (21).

b) Alternating Direction Method of Multipliers: The regularized optimization problem which stems from the MAP rule in (21) can be solved using ADMM [77]. ADMM introduces two auxiliary variables, denoted \mathbf{v} and \mathbf{u} , and is given by the iterative procedure in Algorithm 4, whose derivation is detailed in Appendix E. In Step 2, we defined $f(\mathbf{v}) \triangleq \frac{1}{2} \|\mathbf{x} - \mathbf{H}\mathbf{v}\|^2$, while the proximal mapping of some function $g(\cdot)$ used in Steps 2-3 is defined as

$$\text{prox}_g(\mathbf{v}) := \arg \min_{\mathbf{z}} \left(g(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{v}\|_2^2 \right). \quad (22)$$

The ADMM algorithm is illustrated in Fig. 13(a).

Algorithm 4: ADMM

Init: Fix $\alpha > 0$. Initialize $\mathbf{u}^{(0)}, \mathbf{v}^{(0)}$ randomly
1 for $q = 0, 1, \dots$ **do**
2 Update $\hat{\mathbf{s}}_{q+1} = \text{prox}_{\alpha f}(\mathbf{v}_q + \mathbf{u}_q)$.
3 Update $\mathbf{v}_{q+1} = \text{prox}_{\alpha \phi}(\mathbf{s}_{q+1} + \mathbf{u}_q)$.
4 Update $\mathbf{u}_{q+1} = \mathbf{u}_q + (\hat{\mathbf{s}}_{q+1} - \mathbf{v}_{q+1})$.
5 end
Output: Estimate $\hat{\mathbf{s}} = \hat{\mathbf{s}}_q$.

c) Plug-and-Play ADMM: The key challenge in implementing the ADMM iterations stems from the computation of the proximal mapping in Step 3. In particular, while one can evaluate Step 2 in closed-form, as shown in Appendix E, computing Step 3 of Algorithm 4 requires explicit knowledge of the prior $\phi(\cdot)$, which is often not available. Furthermore, even when one has a good approximation of $\phi(\cdot)$, computing the proximal mapping in Step 3 may still be extremely challenging to carry out analytically.

However, the proximal mapping in Step 3 of Algorithm 4 is invariant of the task and the data. In particular, it is the solution to the problem of MAP denoising $\hat{\mathbf{s}}_{q+1} + \mathbf{u}_q$ assuming the noise-free signal has prior $\phi(\cdot)$ and the noise is Gaussian with variance α . Now, denoisers are common DNN models, and are known to operate reliably on signal domains with intractable priors (e.g., natural images) [81]. One can thus implement ADMM optimization without having to specify the prior $\phi(\cdot)$ by replacing Step 3 of Algorithm 4 with a DNN denoiser [80], as illustrated in Fig. 13. Specifically, the proximal mapping is replaced with a DNN-based denoiser f_{θ} , such that

$$\mathbf{v}_{q+1} = f_{\theta}(\hat{\mathbf{s}}_{q+1} + \mathbf{u}_q; \alpha_q) \quad (23)$$

where α_q denotes the noise level to which the denoiser is tuned. This noise level can either be fixed to represent that used during training, or alternatively, one can use flexible DNN-based denoiser in which, e.g., the noise level is provided as an additional input [90].

Quantitative Results: As an illustrative example of the quantitative gains on plug-and-play networks we consider the setup of cardiac magnetic resonance imaging image reconstruction reported in [80]. The proximal mapping here is replaced with a five-layer CNN with residual connection operating on spatiotemporal volumetric patches. The CNN is trained offline to denoise clean images manually corrupted by Gaussian noise. The experimental results reported in Fig. 14 demonstrate that the introduction of deep denoisers notably improves both the performance and the convergence rate of the iterative optimizer compared to utilizing model-based approaches for approximating the proximal mapping.

Discussion: Using deep learning to strengthen regularized optimization builds upon the model-agnostic nature of DNNs. Traditional optimization methods rely on mathematical expressions to capture the structure of the solution one is looking for, inevitably inducing model mismatch in domains which are extremely challenging to describe analytically. The ability of deep learning to learn complex mappings without relying on domain knowledge is exploited here to bypass the need for explicit regularization. The need to learn to capture the domain of interest facilitates using pre-trained networks, thus reducing the dependency on massive amounts of labeled data. For instance, deep generative priors use DNN architectures that are trained in an unsupervised manner, and thus rely only on unlabeled data, e.g., natural images. Such unlabeled samples are typically more accessible and easy to aggregate compared to labeled data, e.g., tagged natural images. One can often utilize off-the-shelf pre-trained DNNs when such network exist for domains related to the ones over which optimization is carried out, with possible adjustments to account for the subtleties of the problem by transfer learning.

Finally, while our description of DNN-aided regularized optimization relies on model-based iterative optimizers which utilize a deep learning module, one can also incorporate deep learning into the optimization procedure. For instance, the iterative optimization steps can be unfolded into a DNN, as in, e.g., [12]. This approach allows to benefit from both the ability of deep learning to implicitly represent complex domains, as well as the inference speed reduction of deep

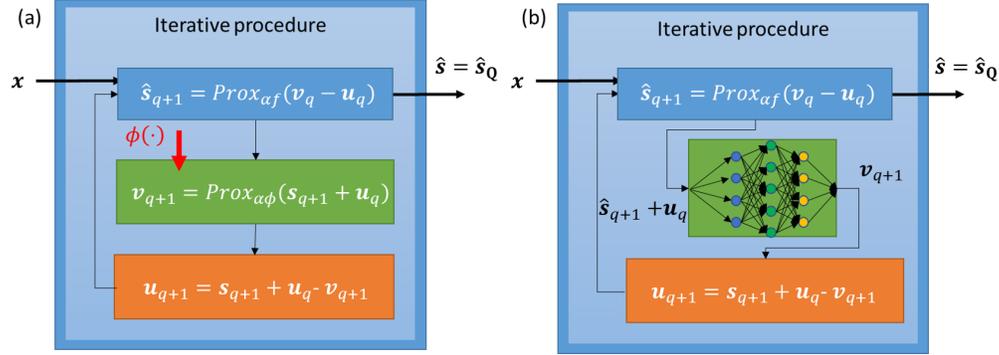


Fig. 13: Illustration of (a) ADMM algorithm compared to (b) plug-and-play ADMM network.

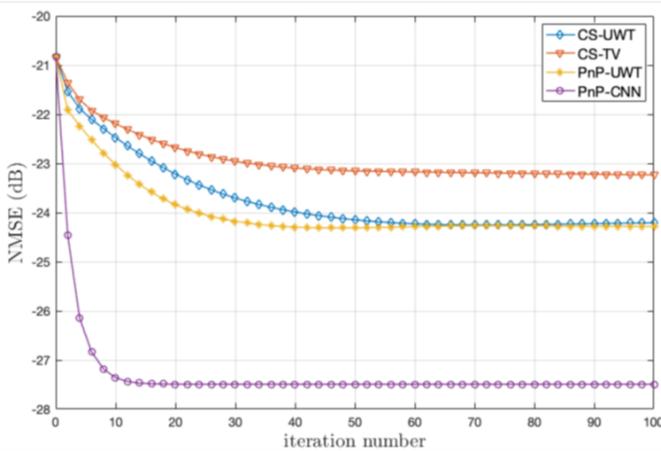


Fig. 14: Normalized MSE versus iteration for the recovery of cardiac MRI images. Here, plug-and-play networks using a CNN denoiser (PnP-CNN) is compared to the model-based strategies of computing the proximal mapping by imposing as prior sparsity in the undecimated wavelet domain (PnP-UWT), as well as CS with a similar constraint (CS-UWT) and with total-variation prior (CS-TV). Figure reproduced from [80] with authors' permission.

unfolding along with its robustness to uncertainty and errors in the model parameters assumed to be known. Nonetheless, the fact that the iterative optimization must be learned from data in addition to the structure of the domain of interest implies that larger amounts of labeled data are required to train the system, compared to using the model-based optimizer.

B. Structure-Oriented DNN-Aided Inference

The family of structure-oriented DNN-aided inference algorithms utilize model-based methods designed to exploit an underlying statistical structure, while integrating DNNs to enable operation without additional explicit characterization of this model. The types of structures exploited in the literature can come in the form of an a-priori known factorizable distribution, such as causality and finite memory in communication channels [15], [22], [91]; it can follow from an established approximation of the statistical behavior, such as modelling of images as conditional random fields [92]–[94]; follow from physical knowledge of the system operation [95]–[97]; or arise due to the distributed nature of the problem, as in [98].

The main advantage in accounting for such statistical structures stems from the availability of various model-based methods, tailored specifically to exploit these structures to facilitate accurate inference at reduced complexity. Many of these algorithms, such as the Kalman filter and its variants [99, Ch. 7], which build upon an underlying state-space structure, or the Viterbi algorithm [100], which exploits the presence of a hidden Markov model, can be represented as special cases of the broad family of factor graph methods. Consequently, our main example used for describing structure-oriented DNN-aided inference focuses on the implementation of message passing over data-driven factor graphs.

Design Outline: Structure-oriented DNN-aided algorithms utilize deep learning not for the overall inference task, but for robustifying and relaxing the model-dependence of established model-based inference algorithms designed specifically for the structure induced by the specific problem being solved. The design of such DNN-aided hybrid inference systems consists of the following steps:

- 1) A proper inference algorithm is chosen based on the available knowledge of the underlying statistical structure. The domain knowledge is encapsulated in the selection of the algorithm which is learned from data.
- 2) Once a model-based algorithm is selected, we identify its model-specific computations, and replace them with dedicated compact DNNs.
- 3) The resulting DNNs are either trained individually, or the overall system can be trained in an end-to-end manner.

We next demonstrate how these steps are translated in a hybrid model-based/data-driven algorithm, using the example of learned factor graph inference for Markovian sequences proposed in [91], [101].

Example 6: Learned Factor Graphs: Factor graph methods, such as the sum-product (SP) algorithm, exploit the factorization of a joint distribution to efficiently compute a desired quantity [102]. The application of the SP algorithm for distributions which can be represented as non-cyclic factor graphs, such as Markovian models, allows computing the MAP rule, an operation whose burden typically grows exponentially with the label space dimensionality, with complexity that only grows linearly with it. While the following description focuses on Markovian stationary time sequences, it can be extended to various forms of factorizable distributions.

a) *System Model*: We consider the recovery of a time series $\{s_i\}$ taking values in a finite set \mathcal{S} from an observed sequence $\{x_i\}$ taking values in a set \mathcal{X} . The subscript i denotes the time index. The joint distribution of $\{s_i\}$ and $\{x_i\}$ obeys an l th-order Markovian stationary model, $l \geq 1$. Consequently, when the initial state $\{s_i\}_{i=-l}^0$ is given, the joint distribution of $\mathbf{x} = [x_1, \dots, x_t]^T$ and $\mathbf{s} = [s_1, \dots, s_t]^T$ satisfies

$$p(\mathbf{x}, \mathbf{s}) = \prod_{i=1}^t p(x_i | \mathbf{s}_{i-l}^i) p(s_i | \mathbf{s}_{i-l}^{i-1}) \quad (24)$$

for any fixed sequence length $t > 0$, where we write $\mathbf{s}_i^j \triangleq [s_i, s_{i+1}, \dots, s_j]^T$ for $i < j$.

b) *The Sum-Product Algorithm*: When the joint distribution of \mathbf{s} and \mathbf{x} is a-priori known and can be computed, the inference rule that minimizes the error probability for each time instance is the MAP detector,

$$\hat{s}_i(\mathbf{x}) = \arg \max_{s_i \in \mathcal{S}} p(s_i | \mathbf{x}) \quad (25)$$

for each $i \in \{1, \dots, t\} \triangleq \mathcal{T}$. This rule can be efficiently approached when (24) holds using the SP algorithm [102]. The SP algorithm represents the joint distribution (24) and computes the posterior distribution by message passing over this graph, as illustrated in Fig. 15(a). The resulting procedure, detailed further in Appendix F, is summarized as Algorithm 5, where we define $\mathbf{s}_i \triangleq \mathbf{s}_{i-l+1}^i \in \mathcal{S}^l$, and the function

$$f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1}) \triangleq p(x_i | \mathbf{s}_i, \mathbf{s}_{i-1}) p(\mathbf{s}_i | \mathbf{s}_{i-1}). \quad (26)$$

Algorithm 5 approaches the MAP detector in (25) with complexity that only grows linearly with t .

Algorithm 5: The SP algorithm for system model (24)

Init: Fix an initial forward message $\vec{\mu}_{\mathbf{s}_0}(\mathbf{s}) = 1$ and a final backward message $\overleftarrow{\mu}_{\mathbf{s}_t}(\mathbf{s}) \equiv 1$.

1 **for** $i = t - 1, t - 2, \dots, 1$ **do**

2 For each $\mathbf{s}_i \in \mathcal{S}^l$, compute backward message

$$\overleftarrow{\mu}_{\mathbf{s}_i}(\mathbf{s}_i) = \sum_{\mathbf{s}_{i+1}} f(x_{i+1}, \mathbf{s}_{i+1}, \mathbf{s}_i) \overleftarrow{\mu}_{\mathbf{s}_{i+1}}(\mathbf{s}_{i+1}).$$

3 **end**

4 **for** $i = 1, 2, \dots, t$ **do**

5 For each $\mathbf{s}_i \in \mathcal{S}^l$, compute forward message

$$\vec{\mu}_{\mathbf{s}_i}(\mathbf{s}_i) = \sum_{\mathbf{s}_{i-1}} f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1}) \vec{\mu}_{\mathbf{s}_{i-1}}(\mathbf{s}_{i-1}).$$

6 Estimate

$$\hat{s}_i = \arg \max_{s_i \in \mathcal{S}} \sum_{\mathbf{s}_{i-1} \in \mathcal{S}^l} \vec{\mu}_{\mathbf{s}_{i-1}}(\mathbf{s}_{i-1}) f(x_i, [s_{i-l+1}, \dots, s_i], \mathbf{s}_{i-1}) \times \overleftarrow{\mu}_{\mathbf{s}_i}([s_{i-l+1}, \dots, s_i]).$$

7

8 **end**

Output: $\hat{\mathbf{s}}^t = [\hat{s}_1, \dots, \hat{s}_t]^T$

c) *Learned Factor Graphs*: Learned factor graphs enable learning to implement MAP detection from labeled data. It

utilizes partial domain knowledge to determine the structure of the factor graph, while using deep learning to compute the function nodes without having to explicitly specify their computations. Finally, it carries out the SP method for inference over the resulting learned factor graph.

Architecture: For Markovian relationships, the structure of the factor graph is that illustrated in Fig. 15(a) regardless of the specific statistical model. Furthermore, the stationarity assumption implies that the complete factor graph is encapsulated in the single function $f(\cdot)$ (26) regardless of the block size t . Building upon this insight, DNNs can be utilized to learn the mapping carried out at the function node separately from the inference task. The resulting learned stationary factor graph is then used to recover $\{s_i\}$ by message passing, as illustrated in Fig. 15(b). As learning a single function node is expected to be a simpler task compared to learning the overall inference method for recovering \mathbf{s} from \mathbf{x} , this approach allows using relatively compact DNNs, which can be learned from a relatively small data set.

Training: In order to learn a stationary factor graph from samples, one must only learn its function node, which here boils down to learning $p(x_i | \mathbf{s}_{i-l}^i)$ and $p(\mathbf{s}_i | \mathbf{s}_{i-l}^{i-1})$ by (26). Since \mathcal{S} is finite, the transition probability $p(\mathbf{s}_i | \mathbf{s}_{i-l}^{i-1})$ can be learned via a histogram.

For learning the distribution $p(x_i | \mathbf{s}_{i-l}^i)$, it is noted that

$$p(x_i | \mathbf{s}_i) = p(\mathbf{s}_i | x_i) p(x_i) (p(\mathbf{s}_i))^{-1}. \quad (27)$$

A parametric estimate of $p(\mathbf{s}_i | x_i)$, denoted $\hat{P}_\theta(\mathbf{s}_i | x_i)$, is obtained for each $\mathbf{s}_i \in \mathcal{S}^{l+1}$ by training classification networks with softmax output layers to minimize the cross entropy loss. As the SP mapping is invariant to scaling $f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1})$ with some factor which does not depend on the $\mathbf{s}_i, \mathbf{s}_{i-1}$, one can set $p(x_i) \equiv 1$ in (27), and use the result to obtain a scaled value of the function node, which, as discussed above, does not affect the inference mapping.

Quantitative Results: As a numerical example of learned factor graphs for Markovian models, we consider a scenario of symbol detection over causal stationary communication channels with finite memory, reproduced from [91]. Fig. 16 depicts the numerically evaluated SER achieved by applying the SP algorithm over a factor graph learned from $n_t = 5000$ labeled samples, for channels with memory $l = 4$. The results are compared to the performance of model-based SP, which requires complete knowledge of the underlying statistical model, as well as the sliding bidirectional RNN detector proposed in [103] for such setups, which utilizes a conventional DNN architecture that does not explicitly account for the Markovian structure. Fig. 16a considers a Gaussian channel, while in Fig. 16b the conditional distribution $p(x_i | \mathbf{s}_{i-l}^i)$ represents a Poisson distribution. Fig. 16 demonstrates the ability of learned factor graphs to enable accurate message passing inference in a data-driven manner, as the performance achieved using learned factor graphs approaches that of the SP algorithm, which operates with full knowledge of the underlying statistical model. The numerical results also demonstrate that combining model-agnostic DNNs with model-aware inference notably improves robustness to model uncertainty compared to

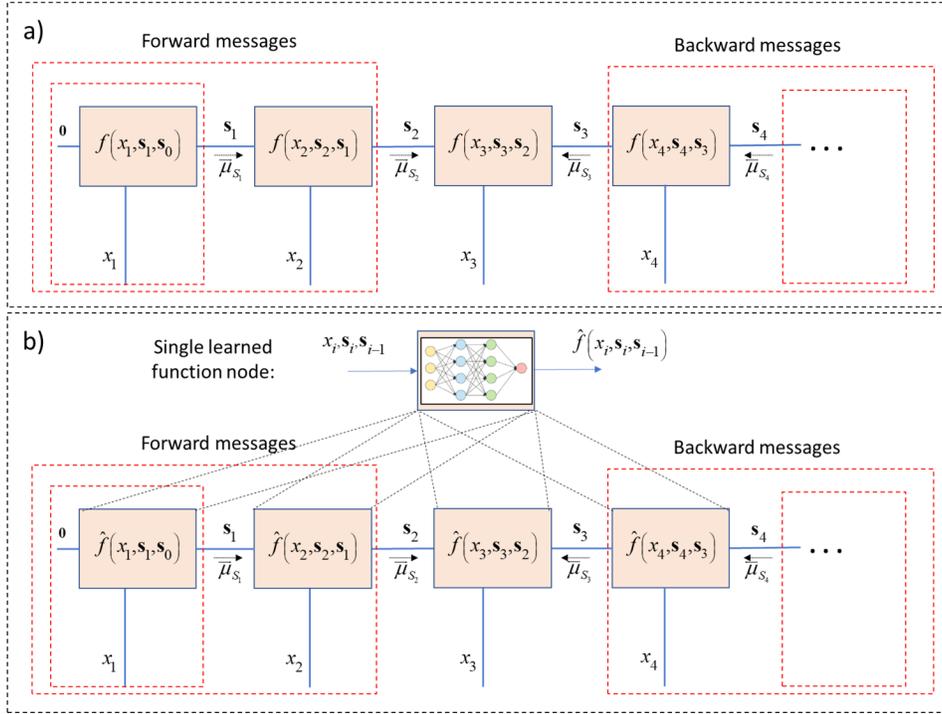


Fig. 15: Illustration of the SP method for Markovian sequences using a) the true factor graph; and b) a learned factor graph.

applying SP with the inaccurate model. Furthermore, it also observed that explicitly accounting for the Markovian structure allows to achieve improved performance compared to utilizing black-box DNN architectures such as the sliding bidirectional RNN detector, with limited data sets for training.

Discussion: The integration of deep learning into structure-oriented model-based algorithms allows to exploit the model-agnostic nature of DNNs while explicitly accounting for available structural domain knowledge. Consequently, structure-oriented DNN-aided inference is most suitable for setups in which structured domain knowledge naturally follows from established models, while the subtleties of the complete statistical knowledge may be challenging to accurately capture analytically. Such structural knowledge is often present in various problems in signal processing and communications. For instance, modelling communication channels as causal finite-memory systems, as assumed in the above quantitative example, is a well-established representation of many physical channels. The availability of established structures in signal processing related setups makes structure-oriented DNN-aided inference a candidate approach to facilitate inference in such scenarios in a manner which is ignorant of the possibly intractable subtleties of the problem, by learning to account for them implicitly from data.

The fact that DNNs are used to learn an intermediate computation rather than the complete predication rule, facilitates the usage of relatively compact DNNs. This property can be exploited to implement learned inference on computationally limited devices, as was done in [97] for DNN-aided velocity tracking in autonomous racing cars. An additional consequence is that the resulting system can be trained using scarce data sets. One can exploit the fact that the system can

be trained using small training sets to, e.g., enable online adaptation to temporal variations in the statistical model based on some feedback on the correctness of the inference rule. This property was exploited in [104] to facilitate online training of DNN-aided receivers in coded communications.

A DNN integrated into a structure-oriented model-based inference method can be either trained individually, i.e., independently of the inference task, or in an end-to-end fashion. The first approach typically requires less training data, and the resulting trained DNN can be combined with various inference algorithms. For instance, the learned function node used to carry out SP inference in the above example can also be integrated into the Viterbi algorithm as done in [15]. Alternatively, the learned modules can be tuned end-to-end by formulating their objective as that of the overall inference algorithm, and backpropagating through the model-based computations, see, e.g., [94]. Learning in an end-to-end fashion facilitates overcoming inaccuracies in the assumed structures, possibly by incorporating learned methods to replace the generic computations of the model-based algorithm, at the cost of requiring larger volumes of data for training purposes.

C. Neural Augmentation

The DNN-aided inference strategies detailed in Subsections V-A and V-B utilize model-based algorithms to carry out inference, while replacing explicit domain-specific computations with dedicated DNNs. An alternative approach, referred to as *neural augmentation*, utilizes the complete model-based algorithm for inference, i.e., without embedding deep learning into its components, while using an external DNN for correcting some of its intermediate computations [21], [23], [105], [106]. An illustration of this approach is depicted in Fig. 17.

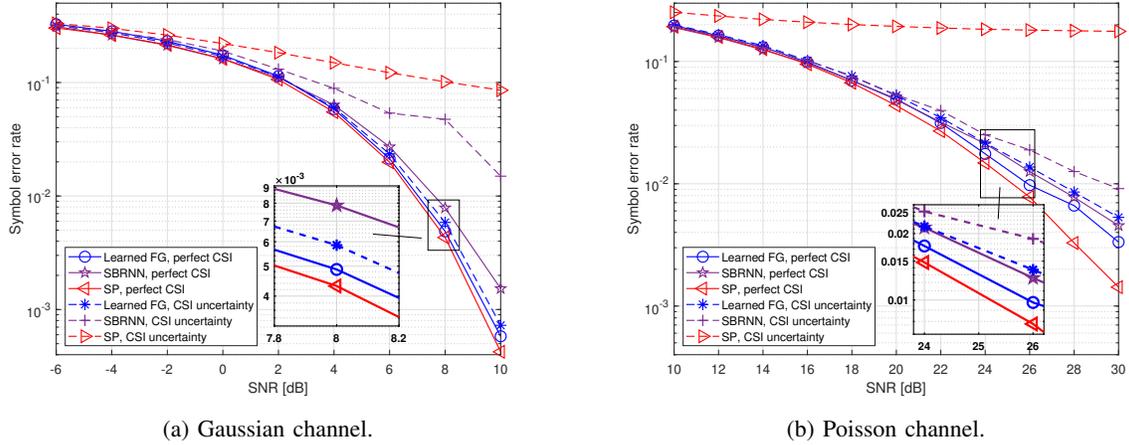


Fig. 16: Experimental results from [91] of learned factor graphs (Learned FG) compared to the model-based SP algorithm and the data-driven sliding bidirectional RNN (SBRNN) of [103]. *Perfect CSI* implies that the system is trained and tested using samples from the same channel, while under *CSI uncertainty* they are trained using samples from a set of different channels.

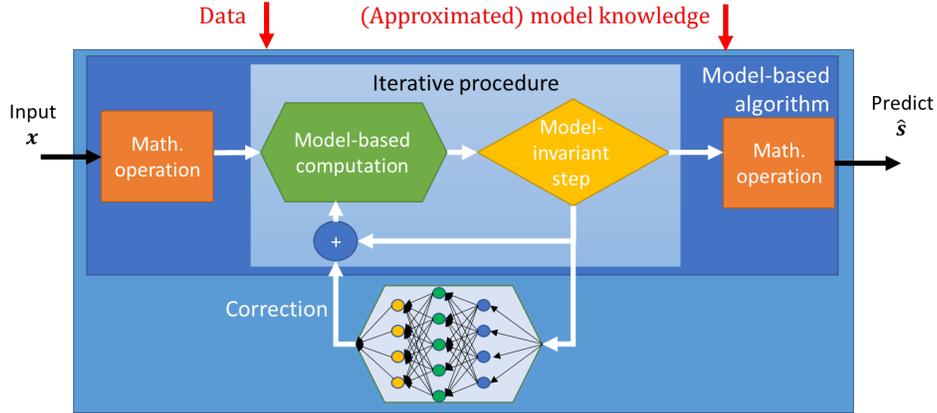


Fig. 17: Neural augmentation illustration.

The main advantage in utilizing an external DNN for correcting internal computations stems from its ability to notably improve the robustness of model-based methods to inaccurate knowledge of the underlying model parameters. Since the model-based algorithm is individually implemented, one must possess the complete domain knowledge it requires, and thus the external correction DNN allows the resulting system to overcome inaccuracies in this domain knowledge by learning to correct them from data. Furthermore, the learned correction term incorporated by neural augmentation can improve the performance of model-based algorithms in scenarios where they are sub-optimal, as detailed in the example in the sequel.

Design Outline: The design of neural-augmented inference systems is comprised of the following steps:

- 1) Choose a suitable iterative optimization algorithm for the problem of interest, and identify the information exchanged between the iterations, along with the intermediate computations used to produce this information.
- 2) The information exchanged between the iterations is updated with a correction term learned by a DNN. The DNN is designed to combine the same quantities used by

the model-based algorithm, only in a learned fashion.

- 3) The overall hybrid model-based/data-driven system is trained in an end-to-end fashion, where one can consider not only the algorithm outputs in the loss function, but also the intermediate outputs of the internal iterations.

We next demonstrate how these steps are carried out in order to augment Kalman smoothing, as proposed in [105].

Example 7: Neural-Augmented Kalman Smoothing: The DNN-aided Kalman smoother proposed in [105] implements state estimation in environments characterized by state-space models. Here, neural augmentation does not only to robustify the smoother in the presence of inaccurate model knowledge, but also improves its performance in non-linear setups, where variants of the Kalman algorithm, such as the extended Kalman method, may be sub-optimal [99, Ch. 7].

a) System Model: Consider a linear Gaussian state-space model. Here, one is interested in recovering a sequence of t state RVs $\{\mathbf{s}_i\}_{i=1}^t$ taking values in a continuous set from an observed sequence $\{\mathbf{x}_i\}_{i=1}^t$. The observations are related to the desired state sequence via

$$\mathbf{x}_i = \mathbf{H}\mathbf{s}_i + \mathbf{r}_i \quad (28a)$$

while the state transition takes the form

$$\mathbf{s}_i = \mathbf{F}\mathbf{s}_{i-1} + \mathbf{w}_i. \quad (28b)$$

In (28), \mathbf{r}_i and \mathbf{w}_i obey an i.i.d. zero-mean Gaussian distributions with covariance \mathbf{R} and \mathbf{W} , respectively, while \mathbf{H} and \mathbf{F} are known linear mappings.

We focus on scenarios where the state-space model in (28) that is available to the inference system, is an inaccurate approximation of the true underlying dynamics. For such scenarios, one can apply Kalman smoothing, which is known to achieve minimal MSE recovery when (28) holds, while introducing a neural augmentation correction term [105].

b) Kalman Smoothing: The Kalman smoother computes the minimal MSE estimate of each s_i given a realization of $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_t]^T$. Its procedure is comprised of forward and backward message passing, exploiting the Markovian structure of the state-space model to operate at complexity which only grows linearly with t . In particular, by writing $\mathbf{s} = [s_1, \dots, s_t]^T$, one way to implement such smoothing to approach the minimal MSE estimate involves applying gradient descent optimization on the joint log likelihood function, i.e., by iterating over

$$\mathbf{s}^{(q+1)} = \mathbf{s}^{(q)} + \eta \nabla_{\mathbf{s}^{(q)}} \log p(\mathbf{x}, \mathbf{s}^{(q)}) \quad (29)$$

where $\eta > 0$ is a step-size. Leveraging the state-space model (28), one can implement gradient descent iterations as message passing, via the procedure summarized in Algorithm 6, whose detailed formulation is given in Appendix G.

Algorithm 6: Smoothing via iterative gradient descent

Init: Fix step-size $\eta > 0$. Set initial guess $\hat{\mathbf{s}}^{(0)}$

1 **for** $q = 0, 1, \dots$ **do**

2 **for** $i = 1, \dots, t$ **do**

3 Compute messages

$$\mu_{\mathbf{S}_{i-1} \rightarrow \mathbf{S}_i}^{(q)} = -\mathbf{W}^{-1} \left(\mathbf{s}_i^{(q)} - \mathbf{F}\mathbf{s}_{i-1}^{(q)} \right),$$

$$\mu_{\mathbf{S}_{i+1} \rightarrow \mathbf{S}_i}^{(q)} = \mathbf{F}^T \mathbf{W}^{-1} \left(\mathbf{s}_{i+1}^{(q)} - \mathbf{F}\mathbf{s}_i^{(q)} \right),$$

$$\mu_{\mathbf{X}_i \rightarrow \mathbf{S}_i}^{(q)} = \mathbf{H}^T \mathbf{R}^{-1} \left(\mathbf{x}_i - \mathbf{H}\mathbf{s}_i^{(q)} \right).$$

4 Update gradient step via

$$\hat{\mathbf{s}}_i^{(q+1)} = \hat{\mathbf{s}}_i^{(q)} + \eta \left(\mu_{\mathbf{S}_{i-1} \rightarrow \mathbf{S}_i}^{(q)} + \mu_{\mathbf{S}_{i+1} \rightarrow \mathbf{S}_i}^{(q)} + \mu_{\mathbf{X}_i \rightarrow \mathbf{S}_i}^{(q)} \right).$$

5 **end**

6 **end**

Output: Estimate $\hat{\mathbf{s}} = \hat{\mathbf{s}}^{(q)}$.

c) Neural-Augmented Kalman Smoothing: The gradient descent formulation in (29) is evaluated by the messages in Step 3 of Algorithm 6, which in turn rely on accurate knowledge of the state-space model (28). To facilitate operation with inaccurate model knowledge due to, e.g., (28) being a linear approximation of a non-linear setup, one can introduce neural

augmentation to learn to correct inaccurate computations of the log-likelihood gradients. This is achieved by using an external DNN to map the messages in Step 3 into a correction term, denoted $\epsilon^{(q+1)}$.

Architecture: The learned mapping of the messages (28) into a correction term operates in the form of a graph neural network (GNN) [107]. This is implemented by maintaining an internal node variable for each variable in Step 3 of Algorithm 6, denoted $h_{\mathbf{s}_i}^{(q)}$ for each $\mathbf{s}_i^{(q)}$ and $h_{\mathbf{x}_i}$ for each \mathbf{x}_i , as well as internal message variables $m_{\mathbf{V}_n \rightarrow \mathbf{S}_i}^{(q)}$ for each message computed by the model-based Algorithm 6. The node variables $h_{\mathbf{s}_i}^{(q)}$ are updated along with the model-based smoothing algorithm iterations as estimates of their corresponding variables, while the variables $h_{\mathbf{x}_i}$ are obtained once from \mathbf{x} via a neural network. The GNN then maps the messages produced by the model-based Kalman smoother into its internal messages via a neural network $f_e(\cdot)$ which operates on the corresponding node variables, i.e.,

$$m_{\mathbf{V}_n \rightarrow \mathbf{S}_i}^{(q)} = f_e \left(h_{\mathbf{v}_n}^{(q)}, h_{\mathbf{s}_i}^{(q)}, \mu_{\mathbf{V}_n \rightarrow \mathbf{S}_i}^{(q)} \right) \quad (30)$$

where $h_{\mathbf{x}_n}^{(q)} \equiv h_{\mathbf{x}_n}$ for each q . These messages are then combined and forwarded into a gated recurrent unit (GRU), which produces the refined estimate of the node variables $\{h_{\mathbf{s}_i}^{(q+1)}\}$ based on their corresponding messages (30). Finally, each updated node variable $h_{\mathbf{s}_i}^{(q+1)}$ is mapped into its corresponding error term $\epsilon_i^{(q+1)}$ via a fourth neural network, denoted $f_d(\cdot)$.

The correction terms $\{\epsilon_i^{(q+1)}\}$ aggregated into the vector $\epsilon^{(q+1)}$ are used to update the log-likelihood gradients, resulting in the update equation (29) replaced with

$$\mathbf{s}^{(q+1)} = \mathbf{s}^{(q)} + \eta \left(\nabla_{\mathbf{s}^{(q)}} \log p(\mathbf{x}, \mathbf{s}^{(q)}) + \epsilon^{(q+1)} \right). \quad (31)$$

The overall architecture is illustrated in Fig. 18.

Training: Let θ be the parameters of the GNN in Fig. 18. The hybrid system is trained end-to-end to minimize the empirical weighted ℓ_2 norm loss over its intermediate layers, where the contribution of each iteration to the overall loss increases as the iterative procedure progresses. In particular, letting $\{(\mathbf{s}_t, \mathbf{x}_t)\}_{t=1}^{n_t}$ be the training set, the loss function used to train the neural-augmented Kalman smoother is given by

$$\mathcal{L}(\theta) = \frac{1}{n_t} \sum_{t=1}^{n_t} \sum_{q=1}^Q \frac{q}{Q} \|\mathbf{s}_t - \hat{\mathbf{s}}_q(\mathbf{x}_t; \theta)\|^2 \quad (32)$$

where $\hat{\mathbf{s}}_q(\mathbf{x}_t; \theta)$ is the estimate produced by the q th iteration, i.e., via (31), with parameters θ and input \mathbf{x}_t .

Quantitative Results: The experiment whose results are depicted in Fig. 19 considers a non-linear state-space model described by the Lorenz attractor equations, which describe atmospheric convection via continuous-time differential equations. The state space model is approximated as a discrete-time linear one by replacing the dynamics with their j th order Taylor series. Fig. 19 demonstrates the ability of neural augmentation to improve model-based inference. It is observed that introducing the DNN-based correction term allows the system to learn to overcome the model inaccuracy, and achieve an error which decreases with the amount of available training data. It is also observed that the hybrid approach of combining

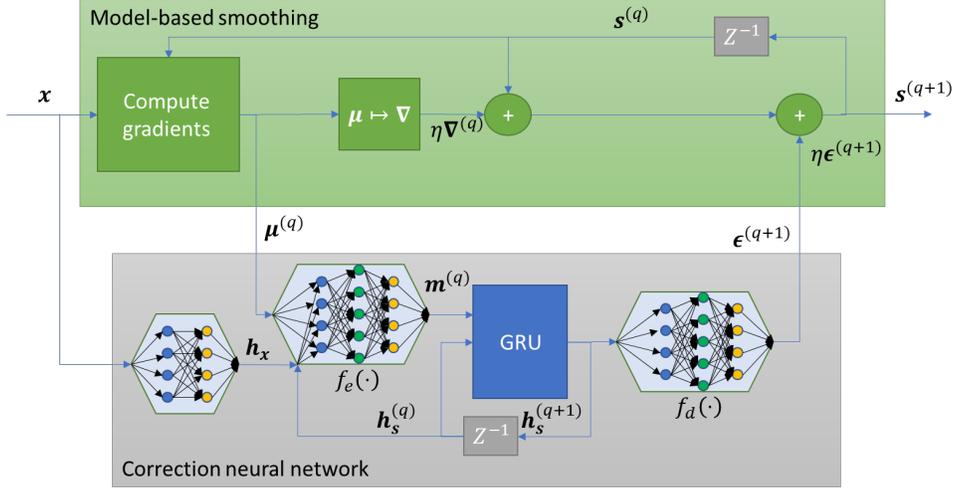


Fig. 18: Neural augmented Kalman smoother illustration. Blocks marked with Z^{-1} represent a single iteration delay.

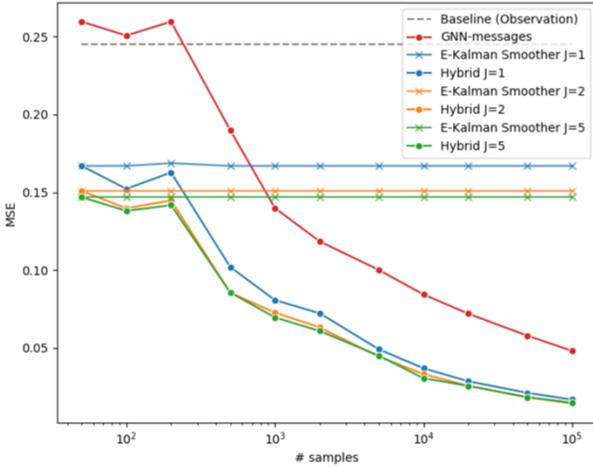


Fig. 19: MSE versus data set size for the neural-augmented Kalman smoother (Hybrid) compared to the model-based extended Kalman smoother (E-Kalman) and a solely data-driven GNN, for various linearizations of state-space models (represented by the index j). Figure reproduced from [105] with authors' permission.

model-based inference and deep learning enables accurate inference with notably reduced volumes of training data, as the individual application of the GNN for state estimation, which does not explicitly account for the available domain knowledge, requires much more training data to achieve similar accuracy as that of the neural-augmented Kalman smoother.

Discussion: Neural augmentation implements hybrid model-based/data-driven inference by utilizing two individual modules – a model-based algorithm and a DNN – with each capable of inferring on its own. The rationale here is to benefit from both approaches by interleaving the iterative operation of the modules, and specifically by utilizing the data-driven component to learn to correct the model-based algorithm, rather than produce individual estimates. This approach thus conceptually differs from the DNN-aided inference strategies discussed in Subsections V-A and V-B, where a DNN is

integrated into a model-based algorithm.

The fact that neural augmentation utilizes individual model-based and data-driven modules reflects on its requirements and use cases. First, one must possess full domain knowledge, or at least an approximation of the true model, in order to implement model-based inference. For instance, the neural-augmented Kalman smoother requires full knowledge of the state-space model (28), or at least an approximation of this analytical closed-form model as used in the quantitative example, in order to compute the exchanged messages in Algorithm 6. Additionally, the presence of an individual DNN module implies that relatively large amounts of data are required in order to train it. Nonetheless, the fact that this DNN only produces a correction term which is interleaved with the model-based algorithm operation implies that the amount of training data required to achieve a given accuracy is notably smaller compared to that required when using solely the DNN for inference. For instance, the quantitative example of the neural augmented Kalman smoother demonstrate that it requires 10 – 20 times less samples compared to that required by the individual GNN to achieve similar MSE results.

VI. CONCLUSIONS AND FUTURE CHALLENGES

In this article, we presented a mapping of methods for combining domain knowledge and data-driven inference via model-based deep learning in a tutorial manner. We noted that hybrid model-based/data-driven systems can be categorized into model-aided networks, which utilize model-based algorithms to design DNN architectures, and DNN-aided inference, where deep learning is integrated into traditional model-based methods. We detailed representative design approaches for each strategy in a systematic manner, along with design guidelines and concrete examples. To conclude this overview, we first summarize the key advantages of model-based deep learning in Subsection VI-A. Then, we present guidelines for selecting a design approach for a given application in Subsection VI-B, intended to facilitate the derivation of future hybrid data-driven/model-based systems. Finally, we review some future research challenges in Subsection VI-C.

A. Advantages of Model-Based Deep Learning

The combination of traditional handcrafted algorithms with emerging data-driven tools via model-based deep learning brings forth several key advantages. Compared to purely model-based schemes, the integration of deep learning facilitates inference in complex environments, where accurately capturing the underlying model in a closed-form mathematical expression may be infeasible. For instance, incorporating DNN-based implicit regularization was shown to enable CS beyond its traditional domain of sparse signals, as discussed in Subsection V-A, while the implementation of the SIC method as an interconnection of neural building blocks enables its operation in non-linear setups, as demonstrated in Subsection IV-B. The model-agnostic nature of deep learning also allows hybrid model-based/data-driven inference to achieve improved resiliency to model uncertainty compared to inferring solely based on domain knowledge. For example, augmenting model-based Kalman smoothing with a GNN was shown in Subsection V-C to notably improve its performance when the state-space model does not fully reflect the true dynamics, while the usage of learned factor graphs for SP inference was demonstrated to result in improved robustness to model uncertainty in Subsection V-B. Finally, the fact that hybrid systems learn to carry out part of their inference based on data allows to infer with reduced delay compared to the corresponding fully model-based methods, as demonstrated by deep unfolding in Subsection IV-A.

Compared to utilizing conventional DNN architectures for inference, the incorporation of domain knowledge via a hybrid model-based/data-driven design results in systems which are tailored for the problem at hand. As a result, model-based deep learning systems require notably less data in order to learn an accurate mapping, as demonstrated in the comparison of learned factor graphs and the sliding bidirectional RNN system in the quantitative example in Subsection V-B, as well as the comparison between the neural augmented Kalman smoother and the GNN state estimator in the corresponding example in Subsection V-C. This property of model-based deep learning systems enables quick adaptation to variations in the underlying statistical model, as shown in [104]. Finally, a system combining DNNs with model-based inference often provides the ability to analyze its resulting predictions, yielding interpretability and confidence which are commonly challenging to obtain with conventional black-box deep learning.

B. Choosing a Model-Based Deep Learning Strategy

The aforementioned gains of model-based deep learning are shared at some level by all the different approaches presented in Sections IV-V. However, each strategy is focused on exploiting a different advantage of hybrid model-based/data-driven inference, particularly in the context of signal processing oriented applications. Consequently, to complement the mapping of model-based deep learning strategies and facilitate the implementation of future application-specific hybrid systems, we next enlist the main considerations one should take into account when seeking to combine model-based methods with data-driven tools for a given problem.

Step 1: Domain knowledge and data characterization:

First, one must ensure the availability of the two key ingredients in model-based deep learning, i.e., domain knowledge and data. The former corresponds to what is known *a priori* about the problem at hand, in terms of statistical models and established assumptions, as well as what is unknown, or is based on some approximation that is likely to be inaccurate. The latter addresses the amount of labeled and unlabeled samples one possesses in advance for the considered problem, as well as whether or not they reflect the scenario in which the system is requested to infer in practice.

Step 2: Identifying a model-based method: Based on the available domain knowledge, the next step is to identify a suitable model-based algorithm for the problem. This choice should rely on the portion of the domain knowledge which is *available*, and not on what is *unknown*, as the latter can be compensated for by integration of deep learning tools. This stage must also consider the requirements of the inference system in terms of performance, complexity, and real-time operation, as these are encapsulated in the selection of the algorithm. The identification of a model-based algorithm, combined with the availability of domain knowledge and data, should also indicate whether model-based deep learning mechanisms are required for the application of interest.

Step 3: Implementation challenges: Having identified a suitable model-based algorithm, the selection of the approach to combine it with deep learning should be based on the understanding of its main implementation challenges. Some representative issues and their relationship with the recommended model-based deep learning approaches include:

- 1) Missing domain knowledge - model-based deep learning can implement the model-based inference algorithm when parts of the underlying model are unknown, or alternatively, too complex to be captured analytically, by harnessing the model-agnostic nature of deep learning. In this case, the selection of the implementation approach depends on the format of the identified model-based algorithm: When it builds upon some known structures via, e.g., message passing based inference, structure-oriented DNN-aided inference detailed in Subsection V-B can be most suitable as means of integrating DNNs to enable operation with missing domain knowledge. Similarly, when the missing domain knowledge can be represented as some complex search domain, or alternatively, an unknown and possibly intractable regularization term, structure-agnostic DNN-aided inference detailed in Subsection V-A can typically facilitate optimization with implicitly learned regularizers. Finally, when the algorithm can be represented as an interconnection of model-dependent building blocks, one can maintain the overall flow of the algorithm while operating in a model-agnostic manner via neural building blocks, as discussed in Subsection IV-B.
- 2) Inaccurate domain knowledge - model-based algorithms are typically sensitive to inaccurate knowledge of the underlying model and its parameters. In such cases, where one has access to a complete description of the underlying model up to some uncertainty, model-based deep learning

can robustify the model-based algorithm and learn to achieve improved accuracy. A candidate approach to robustify model-based processing is by adding a learned correction term via neural augmentation, as detailed in Subsection V-C. Alternatively, when the model-based algorithm takes an iterative form, improved resiliency can be obtained by unfolding the algorithm into a DNN, as discussed in Subsection IV-A, as well as use robust optimization in unfolding [108].

- 3) Inference speed - model-based deep learning can learn to implement iterative inference algorithms, which typically require a large amount of iterations to converge, with reduced inference speed. This is achieved by designing model-aided networks, typically via deep unfolding (see Subsection IV-A) or neural building blocks (see Subsection IV-B). The fact that model-aided networks learn their iterative computations from data allows the resulting system to infer reliably with a much smaller number of iteration-equivalent layers, compared to the iterations required by the model-based algorithm. Alternatively, when the delaying aspect is an internal lengthy computation, one can improve run-time by replacing it with a fixed run-time DNNs via DNN-aided inference, as shown in, e.g., [109].

The aforementioned implementation challenges constitute only a partial list of the considerations one should account for when selecting a model-based deep learning design approach. Additional considerations include computational capabilities during both training as well as inference; the need to handle variations in the statistical model, which in turn translate to a possible requirement to periodically re-train the system; and the quantity and the type of available data. Nonetheless, the above division provides systematic guidelines which one can utilize and possibly extend when seeking to implement an inference system relying on both data and domain knowledge. Finally, we note that some of the detailed model-based deep learning strategies can be combined, and thus one can select more than a single design approach. For instance, one can interleave DNN-aided inference via implicitly learned regularization and/or priors, with deep unfolding of the iterative optimization algorithm, as discussed in Subsection V-A.

C. Future Research Directions

We end by discussing a few representative unexplored research aspects of model-based deep learning:

Performance Guarantees: One of the key strengths of model-based algorithms is their established theoretical performance guarantees. In particular, the analytical tractability of model-based methods implies that one can quantify their expected performance as a function of the parameters of underlying statistical or deterministic models. For conventional deep learning, such performance guarantees are very challenging to characterize, and deeper theoretical understanding is a crucial missing component. The combination of deep learning with model-based structure increases interpretability thus possibly leading to theoretical guarantees. Theoretical guarantees improve the reliability of hybrid model-based/data-driven systems, as well as improve performance. For example,

some preliminary theoretical results were identified for specific model-based deep learning methods, such as the convergence analysis of the unfolded LISTA in [110] and of plug-and-play networks in [82].

Deep Learning Algorithms: Improving model interpretability and incorporating human knowledge is crucial for artificial intelligence development. Model-based deep learning can constitute a systematic framework to incorporate domain knowledge into data-driven systems, and can thus give rise to new forms of deep learning algorithms. For instance, while our description of the methodologies in Sections IV-V systematically commences with a model-based algorithm which is then augmented into a data-aided design via deep learning techniques, one can also envision algorithms in which model-based algorithms are utilized to improve upon an existing DNN architecture. Alternatively, one can leverage model-based techniques to propose interpretable DNN architectures which follow traditional model-based methods to account for domain knowledge.

Collaborative Model-Based Deep Learning: The increasing demands for accessible and personalized artificial intelligence give rise to the need to operate DNNs on edge devices such as smartphones, sensors, and autonomous cars [6]. The limited computational and data resources of edge devices make model-based deep learning strategies particularly attractive for edge intelligence. Latency considerations and privacy constraints for mobile and sensitive data are further driving research in distributed training (e.g., through the framework of federated learning [111], [112]) and collaborative inference [113]. Combining model-based structures with federated learning and distributed inference remains as interesting research directions.

Unexplored Applications: The increasing interest in hybrid model-based/data-driven deep learning methods is motivated by the need for robustness and structural understanding. Applications falling under the broad family of signal processing, communications, and control problems are natural candidates to benefit due to the proliferation of established model-based algorithms. We believe that model-based deep learning can contribute to the development of technologies such as IOT networks, autonomous systems, and wireless communications.

APPENDIX

A. Detailed Formulation of Project Gradient Descent (Example 1, Section IV)

Projected gradient descent iteratively refines its estimate by taking a gradient step with respect to unconstrained objective, followed by projection into the constrained set of the optimization variable. For the system model in (8), this operation at iteration index $q + 1$ is obtained recursively as

$$\begin{aligned} \hat{\mathbf{s}}_{q+1} &= \mathcal{P}_{\mathcal{S}} \left(\hat{\mathbf{s}}_q - \eta \frac{\partial \|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2}{\partial \mathbf{s}} \Big|_{\mathbf{s}=\hat{\mathbf{s}}_q} \right) \\ &= \mathcal{P}_{\mathcal{S}} \left(\hat{\mathbf{s}}_q - \eta \mathbf{H}^T \mathbf{x} + \eta \mathbf{H}^T \mathbf{H} \hat{\mathbf{s}}_q \right) \end{aligned} \quad (\text{A.1})$$

where η is the step size, and $\hat{\mathbf{s}}_0$ is an initial guess.

B. Detailed Formulation of Proximal Gradient Method (Example 2, Section IV)

The recovery of the clean image $\boldsymbol{\mu}$ which can be represented using a convolutional dictionary from the noisy observations \boldsymbol{x} can be formulated as a convolutional sparse coding problem:

$$\begin{aligned} (\hat{\boldsymbol{s}}, \hat{\boldsymbol{H}}) &= \arg \min_{\boldsymbol{s}, \boldsymbol{H}} -\log p_{\boldsymbol{x}|\boldsymbol{\mu}}(\boldsymbol{x}|\boldsymbol{\mu} = \boldsymbol{H}\boldsymbol{s}) + \lambda \|\boldsymbol{s}\|_1 \\ &= \arg \min_{\boldsymbol{s}, \hat{\boldsymbol{H}}} \mathbf{1}^T \exp(\boldsymbol{H}\boldsymbol{s}) - \boldsymbol{x}^T \boldsymbol{H}\boldsymbol{s} + \lambda \|\boldsymbol{s}\|_1, \end{aligned} \quad (\text{B.1})$$

where the dictionary optimization variable is constrained to be block-Toeplitz. The clean image is then obtained as

$$\hat{\boldsymbol{\mu}} = \exp(\hat{\boldsymbol{H}}\hat{\boldsymbol{s}}). \quad (\text{B.2})$$

Here, $\mathbf{1}$ is the all ones vector, λ is a regularizing term that controls the degree of sparsity, boosted by the usage of the ℓ_1 norm.

Algorithm 2 tackles (B.1) via alternating optimization, where the update equations at iteration of index l are given by

$$\begin{aligned} \hat{\boldsymbol{s}}_{l+1} &= \arg \min_{\boldsymbol{s}} \mathbf{1}^T \exp(\boldsymbol{H}\boldsymbol{s}) - \boldsymbol{x}^T \boldsymbol{H}\boldsymbol{s} + \lambda \|\boldsymbol{s}\|_1, \quad (\text{B.3}) \\ &\text{subject to } \boldsymbol{H} = \hat{\boldsymbol{H}}_l \end{aligned}$$

and

$$\begin{aligned} \hat{\boldsymbol{H}}_{l+1} &= \arg \min_{\boldsymbol{H}} \mathbf{1}^T \exp(\boldsymbol{H}\boldsymbol{s}) - \boldsymbol{x}^T \boldsymbol{H}\boldsymbol{s}, \quad (\text{B.4}) \\ &\text{subject to } \boldsymbol{s} = \hat{\boldsymbol{s}}_{l+1}. \end{aligned}$$

The ℓ_1 regularized optimization problem (B.3) can be tackled for a given \boldsymbol{H} and index l via proximal gradient descent iterations. This optimizer involves multiple iterations, indexed $q = 0, 1, 2, \dots$, of the form

$$\hat{\boldsymbol{s}}_{q+1} = \mathcal{T}_b \left(\hat{\boldsymbol{s}}_q + \eta \boldsymbol{H}^T (\boldsymbol{x} - \exp(\boldsymbol{H}\hat{\boldsymbol{s}}_q)) \right). \quad (\text{B.5})$$

The threshold parameter b is dictated by the regularization parameter λ .

C. Detailed Formulation of Iterative Soft Interference Cancellation (Example 3, Section IV)

To formulate the iterative SIC algorithm, we consider the Gaussian MIMO channel in (8). Each iteration of the iterative SIC algorithm indexed q generates K distribution vectors over the set of possible symbols \mathcal{S} . The PMFs are denoted by the vectors $\hat{\boldsymbol{p}}_k^{(q)}$ of size $|\mathcal{S}| \times 1$, where $k \in \mathcal{K}$. These vectors are computed from the observed \boldsymbol{x} as well as the distribution vectors obtained at the previous iteration, $\{\hat{\boldsymbol{p}}_k^{(q-1)}\}_{k=1}^K$. The entries of $\hat{\boldsymbol{p}}_k^{(q)}$ are estimates of the distribution of s_k for each possible symbol in \mathcal{S} , given the observed \boldsymbol{x} and assuming that the interfering symbols $\{s_l\}_{l \neq k}$ are distributed via $\{\hat{\boldsymbol{p}}_l^{(q-1)}\}_{l \neq k}$. Every iteration consists of two steps, carried out in parallel for each user: *Interference cancellation*, and *soft decoding*. Focusing on the k th user and the q th iteration, the interference cancellation stage first computes the expected values and variances of $\{s_l\}_{l \neq k}$ based on the estimated PMF $\{\hat{\boldsymbol{p}}_l^{(q-1)}\}_{l \neq k}$. The contribution of the interfering symbols from \boldsymbol{x} is then canceled by replacing them with $\{e_l^{(q-1)}\}$ and

subtracting their resulting term. Letting \boldsymbol{h}_l be the l th column of \boldsymbol{H} , the interference canceled channel output is given by

$$\boldsymbol{z}_k^{(q)} = \boldsymbol{x} - \sum_{l \neq k} \boldsymbol{h}_l e_l^{(q-1)}. \quad (\text{C.1})$$

Substituting the channel output \boldsymbol{x} into (C.1), the realization of the interference canceled $\boldsymbol{z}_k^{(q)}$ is obtained.

To implement soft decoding, it is assumed that $\boldsymbol{z}_k^{(q)} = \boldsymbol{h}_k s_k + \tilde{\boldsymbol{w}}_k^{(q)}$, where the interference plus noise term $\tilde{\boldsymbol{w}}_k^{(q)}$ obeys a zero-mean Gaussian distribution, independent of s_k , with covariance $\boldsymbol{\Sigma}_k^{(q)} = \sigma_w^2 \boldsymbol{I}_K + \sum_{l \neq k} v_l^{(q-1)} \boldsymbol{h}_l \boldsymbol{h}_l^T$, where σ_w^2 is the noise variance. Combining this assumption with (C.1), while writing the set of possible symbols as $\mathcal{S} = \{\alpha_j\}_{j=1}^{|\mathcal{S}|}$, the conditional distribution of $\boldsymbol{z}_k^{(q)}$ given $s_k = \alpha_j$ is multivariate Gaussian with mean $\boldsymbol{h}_k \alpha_j$ and covariance $\boldsymbol{\Sigma}_k^{(q)}$. The conditional PMF of s_k given \boldsymbol{x} is approximated from the conditional distribution of $\boldsymbol{z}_k^{(q)}$ given s_k via Bayes theorem, assuming that the marginal PMF of each s_k is uniform over \mathcal{S} , this estimated conditional distribution is computed as

$$\begin{aligned} (\hat{\boldsymbol{p}}_k^{(q)})_j &= \\ &= \frac{\exp \left\{ -\frac{1}{2} \left(\boldsymbol{z}_k^{(q)} - \boldsymbol{h}_k \alpha_j \right)^T \left(\boldsymbol{\Sigma}_k^{(q)} \right)^{-1} \left(\boldsymbol{z}_k^{(q)} - \boldsymbol{h}_k \alpha_j \right) \right\}}{\sum_{\alpha_{j'} \in \mathcal{S}} \exp \left\{ -\frac{1}{2} \left(\boldsymbol{z}_k^{(q)} - \boldsymbol{h}_k \alpha_{j'} \right)^T \left(\boldsymbol{\Sigma}_k^{(q)} \right)^{-1} \left(\boldsymbol{z}_k^{(q)} - \boldsymbol{h}_k \alpha_{j'} \right) \right\}} \end{aligned}$$

After the final iteration, the symbols are decoded by maximizing the estimated PMFs for each $k \in \mathcal{K}$, i.e., via

$$\hat{s}_k = \alpha_{\hat{j}}, \quad \hat{j} = \arg \max_j \left(\hat{\boldsymbol{p}}_k^{(Q)} \right)_j, \quad (\text{C.2})$$

and the overall estimate is set to $\hat{\boldsymbol{s}} = [\hat{s}_1, \dots, \hat{s}_K]$.

D. Detailed Formulation of Sparsity-Based CS (Example 4, Section V)

Consider the case where \boldsymbol{s}^* is sparse in some dictionary \boldsymbol{B} , e.g., in the wavelet domain, such that $\boldsymbol{s}^* = \boldsymbol{B}\boldsymbol{c}^*$ where $\|\boldsymbol{c}^*\|_0 = l$ with $l \ll N$. In this case, the goal is to find the sparsest \boldsymbol{c} such that $\boldsymbol{s} = \boldsymbol{B}\boldsymbol{c}$ agrees with the noisy observations:

$$\begin{aligned} &\text{minimize } \|\boldsymbol{c}\|_0 \\ &\text{subject to } \|\boldsymbol{H}\boldsymbol{B}\boldsymbol{c} - \boldsymbol{x}\|_2 \leq \epsilon, \end{aligned}$$

where ϵ is a noise threshold. Since one can define $\tilde{\boldsymbol{H}} := \boldsymbol{H}\boldsymbol{B}$, we henceforth focus on the setting where \boldsymbol{B} is the identity matrix, and the optimization variable of the above ℓ_0 norm optimization problem is \boldsymbol{s} .

Although the above problem is NP-hard, [114], [115] showed that it suffices to minimize the ℓ_1 relaxed LASSO objective in (18). The formulation (18) is convex, and for Gaussian \boldsymbol{A} with $l = \|\boldsymbol{s}^*\|_0$ and $M = \Theta(l \log \frac{N}{l})$, the unique minimizer of $\mathcal{L}_{\text{LASSO}}$ is equal to \boldsymbol{s}^* with high probability.

E. Detailed Formulation of ADMM (Example 5, Section V)

ADMM tackles the optimization problem in (21) by utilizing variable splitting. Namely, it introduces an additional

auxiliary variable \mathbf{v} in order to decouple the regularizer $\phi(\mathbf{s})$ from the likelihood term $\|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2$. The resulting formulation of (21) is expressed as

$$\hat{\mathbf{s}} = \arg \min_{\mathbf{s}} \min_{\mathbf{v}} \frac{1}{2} \|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2 + \phi(\mathbf{v}), \quad (\text{E.1})$$

$$\text{subject to } \mathbf{v} = \mathbf{s}. \quad (\text{E.2})$$

The problem (E.1) is then solved by formulating the augmented Lagrangian (which introduces an additional optimization variable \mathbf{u}) and solving it in an alternating fashion. This results in the following update equations for the q th iteration [82]

$$\hat{\mathbf{s}}_{q+1} = \arg \min_{\mathbf{s}} \frac{\alpha}{2} \|\mathbf{x} - \mathbf{H}\mathbf{s}\|^2 + \frac{1}{2} \|\mathbf{s} - (\mathbf{v}_q - \mathbf{u}_q)\|^2, \quad (\text{E.3a})$$

$$\mathbf{v}_{q+1} = \arg \min_{\mathbf{v}} \alpha \phi(\mathbf{v}) + \frac{1}{2} \|\mathbf{v} - (\hat{\mathbf{s}}_{q+1} + \mathbf{u}_q)\|^2, \quad (\text{E.3b})$$

$$\mathbf{u}_{q+1} = \mathbf{u}_q + (\hat{\mathbf{s}}_{q+1} - \mathbf{v}_{q+1}). \quad (\text{E.3c})$$

Here, $\alpha > 0$ is an optimization hyperparameter. Steps (E.3a) and (E.3b) are the proximal mappings with respect to the functions $\alpha\phi(\cdot)$ and $\alpha f(\cdot)$, respectively, with $f(\mathbf{v}) \triangleq \frac{1}{2} \|\mathbf{x} - \mathbf{H}\mathbf{v}\|^2$. Step (E.3c) represents a gradient ascent iteration.

For brevity, in Algorithm 4 we write (E.3a) as $\hat{\mathbf{s}}_{q+1} = \text{Prox}_{\alpha f}(\mathbf{v}_q - \mathbf{u}_q)$ and (E.3b) as $\mathbf{v}_{q+1} = \text{Prox}_{\alpha\phi}(\hat{\mathbf{s}}_{q+1} + \mathbf{u}_q)$ in Algorithm 4. In particular, it is noted that (E.3a) equals $\mathbf{s}_{q+1} = (\alpha \mathbf{H}^T \mathbf{H} + \mathbf{I})^{-1} (\alpha \mathbf{H}^T \mathbf{x} + (\mathbf{v}_q - \mathbf{u}_q))$.

F. Detailed Formulation of Sum-Product Method (Example 6, Section V)

To formulate the SP method, the factorizable distribution (24) is first represented as a factor graph. To that aim, we recall the definitions of the vector variable $\mathbf{s}_i = \mathbf{s}_{i-l+1}^i \in \mathcal{S}^l$, and the function $f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1})$ in (26). When \mathbf{s}_i is a shifted version of \mathbf{s}_{i-1} , (26) coincides with $p(x_i | \mathbf{s}_{i-l}^i) p(\mathbf{s}_i | \mathbf{s}_{i-l}^i)$, and equals zero otherwise. Using (26), the joint distribution $p(\mathbf{x}, \mathbf{s})$ in (24) can be written as

$$p(\mathbf{x}, \mathbf{s}) = \prod_{i=1}^t f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1}). \quad (\text{F.1})$$

The factorizable expression of the joint distribution (F.1) implies that it can be represented as a factor graph with t function nodes $\{f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1})\}$, in which $\{\mathbf{s}_i\}_{i=2}^{t-1}$ are edges while the remaining variables are half-edges.

Using its factor graph representation, one can compute the joint distribution of \mathbf{s} and \mathbf{x} by recursive message passing along its factor graph as illustrated in Fig. 15(a). In particular,

$$p(\mathbf{s}_k, \mathbf{s}_{k+1}, \mathbf{x}) = \vec{\mu}_{\mathbf{s}_k}(\mathbf{s}_k) f(x_{k+1}, \mathbf{s}_{k+1}, \mathbf{s}_k) \overleftarrow{\mu}_{\mathbf{s}_{k+1}}(\mathbf{s}_{k+1}) \quad (\text{F.2})$$

where the forward path messages satisfy

$$\vec{\mu}_{\mathbf{s}_i}(\mathbf{s}_i) = \sum_{\mathbf{s}_{i-1}} f(x_i, \mathbf{s}_i, \mathbf{s}_{i-1}) \vec{\mu}_{\mathbf{s}_{i-1}}(\mathbf{s}_{i-1}) \quad (\text{F.3})$$

for $i = 1, 2, \dots, k$. Similarly, the backward messages are

$$\overleftarrow{\mu}_{\mathbf{s}_i}(\mathbf{s}_i) = \sum_{\mathbf{s}_{i+1}} f(x_{i+1}, \mathbf{s}_{i+1}, \mathbf{s}_i) \overleftarrow{\mu}_{\mathbf{s}_{i+1}}(\mathbf{s}_{i+1}) \quad (\text{F.4})$$

for $i = t-1, t-2, \dots, k+1$.

The ability to compute the joint distribution in (F.2) via message passing allows to obtain the MAP detector in (25) with complexity that only grows linearly with t . This is achieved by noting that the MAP estimate satisfies

$$\hat{\mathbf{s}}_i(\mathbf{x}) = \arg \max_{\mathbf{s}_i \in \mathcal{S}} \sum_{\mathbf{s}_{i-1} \in \mathcal{S}^l} \vec{\mu}_{\mathbf{s}_{i-1}}(\mathbf{s}_{i-1}) f(x_i, [\mathbf{s}_{i-l+1}, \dots, \mathbf{s}_i], \mathbf{s}_{i-1}) \times \overleftarrow{\mu}_{\mathbf{s}_i}([\mathbf{s}_{i-l+1}, \dots, \mathbf{s}_i]) \quad (\text{F.5})$$

for each $i \in \mathcal{T}$, where the summands can be computed recursively, result in Algorithm 5. It is noted that when the block size l is large, the messages may tend to zero, and are thus commonly scaled [116], e.g., $\overleftarrow{\mu}_{\mathbf{s}_i}(\mathbf{s})$ is replaced with $\gamma_i \overleftarrow{\mu}_{\mathbf{s}_i}(\mathbf{s})$ for some scale factor which does not depend on \mathbf{s} , and thus does not affect the MAP rule.

G. Detailed Formulation of Iterative Kalman Smoother (Example 7, Section V)

The state-space model (28) implies that the joint distribution of the state and observations satisfies

$$p(\mathbf{x}, \mathbf{s}) = p(\mathbf{x} | \mathbf{s}) p(\mathbf{s}) = \prod_t p(\mathbf{x}_t | \mathbf{s}_t) p(\mathbf{s}_t | \mathbf{s}_{t-1}). \quad (\text{G.1})$$

Consequently, it holds that

$$\begin{aligned} \frac{\partial}{\partial \mathbf{s}_t} \log p(\mathbf{x}, \mathbf{s}) &= \frac{\partial}{\partial \mathbf{s}_t} \sum_{\tau} \log p(\mathbf{x}_\tau | \mathbf{s}_\tau) + \sum_{\tau} \log p(\mathbf{s}_\tau | \mathbf{s}_{\tau-1}) \\ &= \frac{\partial}{\partial \mathbf{s}_t} \log p(\mathbf{x}_t | \mathbf{s}_t) + \frac{\partial}{\partial \mathbf{s}_t} \log p(\mathbf{s}_t | \mathbf{s}_{t-1}) + \frac{\partial}{\partial \mathbf{s}_t} \log p(\mathbf{s}_{t+1} | \mathbf{s}_t) \\ &= \frac{\partial}{\partial \mathbf{s}_t} (\mathbf{x}_t - \mathbf{H}\mathbf{s}_t)^T \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{H}\mathbf{s}_t) \\ &\quad + \frac{\partial}{\partial \mathbf{s}_t} (\mathbf{s}_t - \mathbf{F}\mathbf{s}_{t-1})^T \mathbf{W}^{-1} (\mathbf{s}_t - \mathbf{F}\mathbf{s}_{t-1}) \\ &\quad + \frac{\partial}{\partial \mathbf{s}_t} (\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t)^T \mathbf{W}^{-1} (\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t) \\ &= \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{H}\mathbf{s}_t) + -\mathbf{W}^{-1} (\mathbf{s}_t - \mathbf{F}\mathbf{s}_{t-1}) \\ &\quad + \mathbf{F}^T \mathbf{W}^{-1} (\mathbf{s}_{t+1} - \mathbf{F}\mathbf{s}_t). \end{aligned} \quad (\text{G.2})$$

Therefore, the t th entry of the log likelihood gradient in (29), abbreviated henceforth as $\nabla_t^{(q)}$, can be obtained as $\nabla_t^{(q)} = \mu_{\mathbf{s}_{t-1} \rightarrow \mathbf{s}_t}^{(q)} + \mu_{\mathbf{s}_{t+1} \rightarrow \mathbf{s}_t}^{(q)} + \mu_{\mathbf{x}_t \rightarrow \mathbf{s}_t}^{(q)}$, where the summands, referred to as messages, are given by

$$\mu_{\mathbf{s}_{t-1} \rightarrow \mathbf{s}_t}^{(q)} = -\mathbf{W}^{-1} (\mathbf{s}_t^{(q)} - \mathbf{F}\mathbf{s}_{t-1}^{(q)}), \quad (\text{G.3a})$$

$$\mu_{\mathbf{s}_{t+1} \rightarrow \mathbf{s}_t}^{(q)} = \mathbf{F}^T \mathbf{W}^{-1} (\mathbf{s}_{t+1}^{(q)} - \mathbf{F}\mathbf{s}_t^{(q)}), \quad (\text{G.3b})$$

$$\mu_{\mathbf{x}_t \rightarrow \mathbf{s}_t}^{(q)} = \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{H}\mathbf{s}_t^{(q)}). \quad (\text{G.3c})$$

The iterative procedure in (29), is repeated until convergence, as stated in Algorithm 6 and the resulting $\mathbf{s}^{(q)}$ is used as the estimate.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, p. 436, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *IEEE International Conference on Computer Vision*, 2015, pp. 1026–1034.

- [3] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of go without human knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [4] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Molloy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama, D. Wünsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps, and D. Silver, "Grandmaster level in StarCraft II using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [5] Y. Bengio, "Learning deep architectures for AI," *Foundations and Trends in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [6] J. Chen and X. Ran, "Deep learning with edge computing: A review," *Proc. IEEE*, vol. 107, no. 8, pp. 1655–1674, 2019.
- [7] V. Monga, Y. Li, and Y. C. Eldar, "Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing," *IEEE Signal Process. Mag.*, vol. 38, no. 2, pp. 18–44, 2021.
- [8] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *International Conference on Machine Learning*, 2010, pp. 399–406.
- [9] 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*, 2019, pp. 6828–6839.
- [10] A. Bora, A. Jalal, E. Price, and A. G. Dimakis, "Compressed sensing using generative models," in *International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 537–546.
- [11] J. Whang, Q. Lei, and A. G. Dimakis, "Compressed sensing with invertible generative models and dependent noise," in *International Conference on Learning Representations*, 2021.
- [12] D. Gilton, G. Ongie, and R. Willett, "Neumann networks for inverse problems in imaging," *IEEE Trans. Comput. Imaging*, vol. 6, pp. 328–343, 2019.
- [13] S. V. Venkatakrisnan, C. A. Bouman, and B. Wohlberg, "Plug-and-play priors for model based reconstruction," in *IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, 2013, pp. 945–948.
- [14] H. K. Aggarwal, M. P. Mani, and M. Jacob, "MoDL: Model-based deep learning architecture for inverse problems," *IEEE Trans. Med. Imag.*, vol. 38, no. 2, pp. 394–405, 2018.
- [15] N. Shlezinger, N. Farsad, Y. C. Eldar, and A. J. Goldsmith, "ViterbiNet: A deep learning based Viterbi algorithm for symbol detection," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3319–3331, 2020.
- [16] N. Shlezinger, R. Fu, and Y. C. Eldar, "DeepSIC: Deep soft interference cancellation for multiuser MIMO detection," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1349–1362, 2021.
- [17] E. Nachmani, E. Marciano, L. Lugosch, W. J. Gross, D. Burshtein, and Y. Be'ery, "Deep learning methods for improved decoding of linear codes," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, pp. 119–131, 2018.
- [18] N. Samuel, T. Diskin, and A. Wiesel, "Learning to detect," *IEEE Trans. Signal Process.*, vol. 67, no. 10, pp. 2554–2564, 2019.
- [19] H. He, C.-K. Wen, S. Jin, and G. Y. Li, "Model-driven deep learning for MIMO detection," *IEEE Trans. Signal Process.*, vol. 68, pp. 1702–1715, 2020.
- [20] M. Khani, M. Alizadeh, J. Hoydis, and P. Fleming, "Adaptive neural signal detection for massive MIMO," *IEEE Trans. Wireless Commun.*, vol. 19, no. 8, pp. 5635–5648, 2020.
- [21] K. Pratik, B. D. Rao, and M. Welling, "RE-MIMO: Recurrent and permutation equivariant neural MIMO detection," *IEEE Trans. Signal Process.*, vol. 69, pp. 459–473, 2020.
- [22] N. Farsad, N. Shlezinger, A. J. Goldsmith, and Y. C. Eldar, "Data-driven symbol detection via model-based machine learning," *Communications in Information and Systems*, vol. 20, no. 3, pp. 283–317, 2020.
- [23] V. G. Satorras and M. Welling, "Neural enhanced belief propagation on factor graphs," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 685–693.
- [24] A. Zappone, M. Di Renzo, M. Debbah, T. T. Lam, and X. Qian, "Model-aided wireless artificial intelligence: Embedding expert knowledge in deep neural networks for wireless system optimization," *IEEE Veh. Technol. Mag.*, vol. 14, no. 3, pp. 60–69, 2019.
- [25] A. Zappone, M. Di Renzo, and M. Debbah, "Wireless networks design in the era of deep learning: Model-based, AI-based, or both?" *IEEE Trans. Commun.*, vol. 67, no. 10, pp. 7331–7376, 2019.
- [26] L. Liang, H. Ye, G. Yu, and G. Y. Li, "Deep-learning-based wireless resource allocation with application to vehicular networks," *Proc. IEEE*, vol. 108, no. 2, pp. 341–356, 2019.
- [27] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Trans. on Cogn. Commun. Netw.*, vol. 3, no. 4, pp. 563–575, 2017.
- [28] H. Kim, Y. Jiang, R. Rana, S. Kannan, S. Oh, and P. Viswanath, "Communication algorithms via deep learning," in *International Conference on Learning Representations*, 2018.
- [29] M. B. Mashhadi, Q. Yang, and D. Gündüz, "Distributed deep convolutional compression for massive MIMO CSI feedback," *IEEE Trans. Wireless Commun.*, vol. 21, no. 4, pp. 2621–2633, 2021.
- [30] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [31] C. Metzler, A. Mousavi, and R. Baraniuk, "Learned D-AMP: Principled neural network based compressive image recovery," in *Advances in Neural Information Processing Systems*, 2017, pp. 1772–1783.
- [32] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [33] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [35] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [36] T. Tieleman and G. Hinton, "Lecture 6.5-RMSProp: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.
- [37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations*, 2015.
- [38] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [39] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
- [40] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 214–223.
- [41] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of Wasserstein GANs," *Advances in neural information processing systems*, vol. 30, 2017.
- [42] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *IEEE international conference on computer vision*, 2017, pp. 2794–2802.
- [43] J. H. Lim and J. C. Ye, "Geometric GAN," *arXiv preprint arXiv:1705.02894*, 2017.
- [44] A. Jolicoeur-Martineau, "The relativistic discriminator: a key element missing from standard GAN," in *International Conference on Learning Representations*, 2019.
- [45] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of stylegan," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8110–8119.
- [46] J. E. Van Engelen and H. H. Hoos, "A survey on semi-supervised learning," *Machine Learning*, vol. 109, no. 2, pp. 373–440, 2020.
- [47] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *International Conference on Machine Learning*, 2013.
- [48] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," in *International Conference on Learning Representations*, 2016.
- [49] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, "MixMatch: A holistic approach to semi-supervised learning," in *Advances in Neural Information Processing Systems*, 2019.
- [50] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, "Self-training with noisy student improves imagenet classification," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10687–10698.

- [51] B. Tolooshams, A. H. Song, S. Temereanca, and D. Ba, "Convolutional dictionary learning based auto-encoders for natural exponential-family distributions," in *International Conference on Machine Learning*, PMLR, 2020, pp. 9493–9503.
- [52] L. Xu and R. Niu, "EKFNNet: Learning system noise statistics from measurement data," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 4560–4564.
- [53] N. Shlezinger, Y. C. Eldar, and S. P. Boyd, "Model-based deep learning: On the intersection of deep learning and optimization," *arXiv preprint arXiv:2205.02640*, 2022.
- [54] A. Ng and M. Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes," *Advances in neural information processing systems*, vol. 14, 2001.
- [55] N. Shlezinger and T. Routtenberg, "Discriminative and generative learning for linear estimation of random signals [lecture notes]," *arXiv preprint arXiv:2206.04432*, 2022.
- [56] J. R. Hershey, J. L. Roux, and F. Weninger, "Deep unfolding: Model-based inspiration of novel deep architectures," *arXiv preprint arXiv:1409.2574*, 2014.
- [57] Y. Li, M. Tofighi, J. Geng, V. Monga, and Y. C. Eldar, "Efficient and interpretable deep blind image deblurring via algorithm unrolling," *IEEE Trans. Comput. Imaging*, vol. 6, pp. 666–681, 2020.
- [58] O. Solomon, R. Cohen, Y. Zhang, Y. Yang, Q. He, J. Luo, R. J. van Sloun, and Y. C. Eldar, "Deep unfolded robust PCA with application to clutter suppression in ultrasound," *IEEE Trans. Med. Imag.*, vol. 39, no. 4, pp. 1051–1063, 2019.
- [59] Y. Cui, S. Li, and W. Zhang, "Jointly sparse signal recovery and support recovery via deep learning with applications in MIMO-based grant-free random access," *IEEE J. Sel. Areas Commun.*, vol. 6, no. 3, pp. 788–803, 2021.
- [60] T. Chang, B. Tolooshams, and D. Ba, "RandNet: deep learning with compressed measurements of images," in *Proc. IEEE MLSP*, 2019.
- [61] A. Balatsoukas-Stimming and C. Studer, "Deep unfolding for communications systems: A survey and some new directions," in *IEEE International Workshop on Signal Processing Systems (SIPS)*, 2019, pp. 266–271.
- [62] S. Takabe, M. Imanishi, T. Wadayama, R. Hayakawa, and K. Hayashi, "Trainable projected gradient detector for massive overloaded MIMO channels: Data-driven tuning approach," *IEEE Access*, vol. 7, pp. 93 326–93 338, 2019.
- [63] Q. Hu, Y. Cai, Q. Shi, K. Xu, G. Yu, and Z. Ding, "Iterative algorithm induced deep-unfolding neural networks: Precoding design for multiuser MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1394–1410, 2021.
- [64] S. Khobahi, N. Shlezinger, M. Soltanalian, and Y. C. Eldar, "LoRD-Net: Low resolution detection network for deep low-resolution receivers," *IEEE Trans. Signal Process.*, vol. 69, pp. 5651–5664, 2021.
- [65] M. Misch, M. A. L. Bell, R. J. van Sloun, and Y. C. Eldar, "Deep learning in medical ultrasound—from image formation to image analysis," *IEEE Trans. Ultrason., Ferroelectr., Freq. Control*, vol. 67, no. 12, pp. 2477–2480, 2020.
- [66] G. Dardikman-Yoffe and Y. C. Eldar, "Learned SPARCOM: Unfolded deep super-resolution microscopy," *Optics Express*, vol. 28, no. 19, pp. 4797–4812, 2020.
- [67] K. Zhang, L. V. Gool, and R. Timofte, "Deep unfolding network for image super-resolution," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3217–3226.
- [68] Y. Huang, S. Li, L. Wang, and T. Tan, "Unfolding the alternating optimization for blind super resolution," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [69] A. Agarwal, A. Anandkumar, P. Jain, and P. Netrapalli, "Learning sparsely used overcomplete dictionaries via alternating minimization," *SIAM Journal on Optimization*, vol. 26, no. 4, pp. 2775–2799, 2016.
- [70] T. Remez, O. Litany, R. Giryes, and A. M. Bronstein, "Class-aware fully convolutional Gaussian and Poisson denoising," *IEEE Trans. Signal Process.*, vol. 27, no. 11, pp. 5707–5722, 2018.
- [71] J. Duan, J. Schlemper, C. Qin, C. Ouyang, W. Bai, C. Biffi, G. Bello, B. Stott, D. P. O'Regan, and D. Rueckert, "VS-Net: Variable splitting network for accelerated parallel MRI reconstruction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 713–722.
- [72] J. P. Merkofer, G. Revach, N. Shlezinger, and R. J. van Sloun, "Deep augmented music algorithm for data-driven DoA estimation," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 3598–3602.
- [73] T. Van Luong, N. Shlezinger, C. Xu, T. M. Hoang, Y. C. Eldar, and L. Hanzo, "Deep learning based successive interference cancellation for the non-orthogonal downlink," *IEEE Trans. Veh. Technol.*, 2022.
- [74] M. Kocaoglu, C. Snyder, A. G. Dimakis, and S. Vishwanath, "Causal-GAN: Learning causal implicit generative models with adversarial training," in *International Conference on Learning Representations*, 2018.
- [75] W.-J. Choi, K.-W. Cheong, and J. M. Cioffi, "Iterative soft interference cancellation for multiple antenna systems," in *Proc. WCNC*, 2000, pp. 304–309.
- [76] G. Ongie, A. Jalal, C. A. Metzler, R. G. Baraniuk, A. G. Dimakis, and R. Willett, "Deep learning techniques for inverse problems in imaging," *IEEE J. Sel. Areas Inform. Theory*, vol. 1, no. 1, pp. 39–56, 2020.
- [77] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [78] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [79] A. Chambolle and T. Pock, "A first-order primal-dual algorithm for convex problems with applications to imaging," *Journal of mathematical imaging and vision*, vol. 40, no. 1, pp. 120–145, 2011.
- [80] R. Ahmad, C. A. Bouman, G. T. Buzzard, S. Chan, S. Liu, E. T. Reehorst, and P. Schniter, "Plug-and-play methods for magnetic resonance imaging: Using denoisers for image recovery," *IEEE Signal Process. Mag.*, vol. 37, no. 1, pp. 105–116, 2020.
- [81] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep CNN denoiser prior for image restoration," in *IEEE conference on computer vision and pattern recognition*, 2017, pp. 3929–3938.
- [82] E. Ryu, J. Liu, S. Wang, X. Chen, Z. Wang, and W. Yin, "Plug-and-play methods provably converge with properly trained denoisers," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5546–5557.
- [83] S. Ono, "Primal-dual plug-and-play image restoration," *IEEE Signal Process. Lett.*, vol. 24, no. 8, pp. 1108–1112, 2017.
- [84] U. S. Kamilov, H. Mansour, and B. Wohlberg, "A plug-and-play priors approach for solving nonlinear imaging inverse problems," *IEEE Signal Process. Lett.*, vol. 24, no. 12, pp. 1872–1876, 2017.
- [85] T. Meinhardt, M. Moller, C. Hazirbas, and D. Cremers, "Learning proximal operators: Using denoising networks for regularizing inverse imaging problems," in *IEEE International Conference on Computer Vision*, 2017, pp. 1781–1790.
- [86] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [87] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [88] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," 2014.
- [89] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [90] K. Zhang, W. Zuo, and L. Zhang, "FFDNet: Toward a fast and flexible solution for cnn-based image denoising," *IEEE Trans. Image Process.*, vol. 27, no. 9, pp. 4608–4622, 2018.
- [91] N. Shlezinger, N. Farsad, Y. C. Eldar, and A. J. Goldsmith, "Data-driven factor graphs for deep symbol detection," in *International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 2682–2687.
- [92] A. Arnab, S. Zheng, S. Jayasumana, B. Romera-Paredes, M. Larsson, A. Kirillov, B. Savchynskyy, C. Rother, F. Kahl, and P. H. Torr, "Conditional random fields meet deep neural networks for semantic segmentation: Combining probabilistic graphical models with deep learning for structured prediction," *IEEE Signal Process. Mag.*, vol. 35, no. 1, pp. 37–52, 2018.
- [93] S. Chandra and I. Kokkinos, "Fast, exact and multi-scale inference for semantic image segmentation with deep Gaussian CRFs," in *European conference on computer vision*. Springer, 2016, pp. 402–418.
- [94] P. Knobelreiter, C. Sormann, A. Shekhovtsov, F. Fraundorfer, and T. Pock, "Belief propagation reloaded: Learning BP-layers for labeling problems," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7900–7909.
- [95] B. Luitjen, R. Cohen, F. J. De Bruijn, H. A. Schmeitz, M. Misch, Y. C. Eldar, and R. J. Van Sloun, "Adaptive ultrasound beamforming using deep learning," *IEEE Trans. Med. Imag.*, vol. 39, no. 12, pp. 3967–3978, 2020.
- [96] G. Revach, N. Shlezinger, X. Ni, A. L. Escoriza, R. J. Van Sloun, and Y. C. Eldar, "KalmanNet: Neural network aided Kalman filtering for

- partially known dynamics,” *IEEE Trans. Signal Process.*, vol. 70, pp. 1532–1547, 2022.
- [97] A. L. Escoriza, G. Revach, N. Shlezinger, and R. J. G. van Sloun, “Data-driven Kalman-based velocity estimation for autonomous racing,” in *IEEE International Conference on Autonomous Systems (ICAS)*, 2021.
- [98] H. Palangi, R. Ward, and L. Deng, “Distributed compressive sensing: A deep learning approach,” *IEEE Trans. Signal Process.*, vol. 64, no. 17, pp. 4504–4518, 2016.
- [99] S. S. Haykin, *Adaptive filter theory*. Pearson Education India, 2005.
- [100] A. Viterbi, “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm,” *IEEE Trans. Inf. Theory*, vol. 13, no. 2, pp. 260–269, 1967.
- [101] N. Shlezinger, N. Farsad, Y. C. Eldar, and A. J. Goldsmith, “Learned factor graphs for inference from stationary time sequences,” *IEEE Trans. Signal Process.*, vol. 70, pp. 366–380, 2021.
- [102] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, “Factor graphs and the sum-product algorithm,” *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 498–519, 2001.
- [103] N. Farsad and A. Goldsmith, “Neural network detection of data sequences in communication systems,” *IEEE Trans. Signal Process.*, vol. 66, no. 21, pp. 5663–5678, 2018.
- [104] T. Raviv, S. Park, N. Shlezinger, O. Simeone, Y. C. Eldar, and J. Kang, “Meta-ViterbiNet: Online meta-learned Viterbi equalization for non-stationary channels,” in *IEEE International Conference on Communications (ICC)*, 2021.
- [105] V. G. Satorras, Z. Akata, and M. Welling, “Combining generative and discriminative models for hybrid inference,” in *Advances in Neural Information Processing Systems*, 2019, pp. 13 802–13 812.
- [106] F. Gao, J. Zhang, and Y. Zhang, “Neural enhanced dynamic message passing,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 10 471–10 482.
- [107] K. Yoon, R. Liao, Y. Xiong, L. Zhang, E. Fetaya, R. Urtasun, R. Zemel, and X. Pitkow, “Inference in probabilistic graphical models by graph neural networks,” in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2019, pp. 868–875.
- [108] W. Pu, C. Zhou, Y. C. Eldar, and M. R. Rodrigues, “REST: Robust learned shrinkage-thresholding network taming inverse problems with model mismatch,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 2885–2889.
- [109] X. Ni, G. Revach, N. Shlezinger, R. J. van Sloun, and Y. C. Eldar, “RT-SNet: Deep learning aided Kalman smoothing,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 5902–5906.
- [110] X. Chen, J. Liu, Z. Wang, and W. Yin, “Theoretical linear convergence of unfolded ISTA and its practical weights and thresholds,” in *Advances in Neural Information Processing Systems*, 2018, pp. 9079–9089.
- [111] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated learning: Challenges, methods, and future directions,” *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [112] T. Gafni, N. Shlezinger, K. Cohen, Y. C. Eldar, and H. V. Poor, “Federated learning: A signal processing perspective,” *IEEE Signal Process. Mag.*, vol. 39, no. 3, pp. 14–41, 2022.
- [113] N. Shlezinger and I. V. Bajic, “Collaborative inference for AI-empowered IoT devices,” *arXiv preprint arXiv:2207.11664*, 2022.
- [114] E. J. Candès, J. Romberg, and T. Tao, “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [115] D. L. Donoho, “Compressed sensing,” *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [116] H.-A. Loeliger, “An introduction to factor graphs,” *IEEE Signal Process. Mag.*, vol. 21, no. 1, pp. 28–41, 2004.