
A Neural Network Approach for Efficiently Answering Most Probable Explanation Queries in Probabilistic Models

Shivvrat Arya

Department of Computer Science
The University of Texas at Dallas
Dallas, TX 75080
shivvrat.arya@utdallas.edu

Tahrima Rahman

Department of Computer Science
The University of Texas at Dallas
Dallas, TX 75080
tahrima.rahman@utdallas.edu

Vibhav Gogate

Department of Computer Science
The University of Texas at Dallas
Dallas, TX 75080
vibhav.gogate@utdallas.edu

Abstract

We propose a novel neural networks based approach to efficiently answer arbitrary Most Probable Explanation (MPE) queries—a well-known NP-hard task—in large probabilistic models such as Bayesian and Markov networks, probabilistic circuits, and neural auto-regressive models. By arbitrary MPE queries, we mean that there is no predefined partition of variables into evidence and non-evidence variables. The key idea is to distill all MPE queries over a given probabilistic model into a neural network and then use the latter for answering queries, eliminating the need for time-consuming inference algorithms that operate directly on the probabilistic model. We improve upon this idea by incorporating inference-time optimization with self-supervised loss to iteratively improve the solutions and employ a teacher-student framework that provides a better initial network, which in turn, helps reduce the number of inference-time optimization steps. The teacher network utilizes a self-supervised loss function optimized for getting the exact MPE solution, while the student network learns from the teacher’s near-optimal outputs through supervised loss. We demonstrate the efficacy and scalability of our approach on various datasets and a broad class of probabilistic models, showcasing its practical effectiveness.

1 Introduction

Probabilistic representations such as Probabilistic Circuits (PCs) [8], graphical models [26] such as Bayesian Networks (BNs) and Markov Networks (MNs), and Neural Autoregressive Models (NAMs) [50] are widely used to model large, multi-dimensional probability distributions. However, they face a significant challenge: as the complexity of these distributions increases, solving practically relevant NP-hard inference tasks such as finding the Most Probable Explanation (MPE) via exact inference techniques [35, 36] becomes increasingly difficult and time-consuming. In particular, although various exact and approximate solvers exist for the MPE task in PCs, BNs and MNs, exact solvers are often too slow for practical use, and approximate solvers tend to lack the necessary accuracy, particularly in autoregressive models that currently rely on slow hill-climbing/beam search methods.

In recent work, Arya et al. [4] proposed a method to overcome the limitations of existing approximate methods by using neural networks (NNs) to solve the MPE task in PCs.¹ Their method draws inspiration from the learning to optimize literature [12, 15, 29, 42, 55]. Given a PC and a *predefined partition of variables into query and evidence sets*, the core idea is to train a NN that takes an assignment to the evidence variables as input and outputs the most likely assignment to the query variables w.r.t. the distribution defined by the PC. Arya et al. suggest using either supervised or self-supervised learning techniques to train the NN; the former requires access to exact inference schemes, while the latter does not and is therefore more practical.

In this paper, we address a more general and complex version of the MPE task than the one considered by Arya et al. Specifically, we assume that there is *no predefined partition of the variables into evidence and query sets*, which we refer to as the **any-MPE** task. The complexity of the any-MPE task arises from the exponential increase in the number of input configurations, compounded by the exponential number of possible divisions of variables into evidence and query sets. Furthermore, our method applies to a broad class of probabilistic models, including BNs, MNs and NAMs, whereas Arya et al.’s method is limited to PCs. In addition, Arya et al.’s method does not fully exploit the capabilities of self-supervision, and the benefits of combining supervised and self-supervised loss functions.

This paper presents a novel approach that uses a NN for solving the any-MPE task in a broad class of probabilistic models (PMs) and achieves technical advancements in three key aspects:

1. Efficient MPE Inference via Encoding Scheme and Loss Function: We introduce a new encoding scheme that tailors the NN architecture to the specific structure of the input PM. This scheme not only delineates the input and output nodes for the NN but also establishes a methodology for setting input values and extracting the MPE solution from the NN’s outputs. Furthermore, we propose a tractable, and differentiable self-supervised loss function, enabling efficient training.

2. Inference Time Optimization with ITSELF: We introduce a novel inference technique called Inference Time Self Supervised Training (ITSELF). This technique iteratively refines the MPE solution during the inference process itself. It utilizes gradient descent (*back-propagation*) to update the NN’s parameters using our proposed self-supervised loss, leading to continual (anytime) improvement towards near-optimal solutions. ITSELF fully utilizes the power of our self-supervised loss, as it does not require labeled data or an external MPE solver.

3. Two-Phase Pre-training with Teacher-Student Architecture: To address challenges associated with self-supervised learning and ITSELF, we propose a two-phase pre-training strategy that leverages a teacher-student architecture. Self-supervised learning can suffer from overfitting and requires careful regularization. Additionally, ITSELF, especially with random initializations, might necessitate a substantial number of gradient updates to converge on optimal solutions. Our approach addresses these issues using the following methodology: (i) The teacher network first overfits the training data using ITSELF and (ii) The student network is then trained using supervised loss functions (e.g., binary cross-entropy) by treating the teacher network’s output as pseudo-labels. This supervised training phase improves and regularizes the parameter learning process of the student network. It also provides a robust starting point for ITSELF, significantly reducing the required optimization steps and leading to substantial performance gains.

Finally, we conduct a detailed experimental comparison of our method with existing approaches on several types of PMs such as PCs, BNs, MNs and NAMs. Our results demonstrate that our method surpasses state-of-the-art approximate inference techniques in terms of both accuracy and speed.

2 Background and Motivation

Without loss of generality, we use binary variables which take values from the set $\{0, 1\}$. We denote a random variable by an uppercase letter (e.g., X), and a value assigned to it by the corresponding lowercase letter (e.g., x). We denote a set of random variables by a bold uppercase letter (e.g., \mathbf{X}) and an assignment of values to all variables in the set by the corresponding bold lowercase letter (e.g., \mathbf{x}).

¹Arya et al. [4] developed a NN-based method for solving the *marginal maximum-a-posteriori* (MMAP) task in PCs. In this paper, we focus on the MPE task, also sometimes referred to as the full MAP task, which is a special case of MMAP. Our method can be easily extended for solving the MMAP problem in PCs and tractable graphical models. For simplicity of exposition, we concentrate on the MPE task in this paper.

Throughout the paper when we use the term probabilistic models (PMs), we are referring to a broad class of probabilistic models in which computing the likelihood² of an assignment to all variables in the model can be done in polynomial (preferably linear) time in the size of the model. This class includes, among others, Bayesian and Markov networks collectively called probabilistic graphical models (PGMs) [26], smooth and decomposable probabilistic circuits (PCs) [8], and neural autoregressive models (NAMs) such as NADE [50] and MADE [17].

We are interested in solving the most probable explanation (MPE) task in PMs, namely the task of finding the most likely assignment to all unobserved (non-evidence) variables given observations (evidence). Formally, let \mathcal{M} denote a probabilistic model defined over a set of variables \mathbf{X} that represents the distribution $p_{\mathcal{M}}(\mathbf{x})$. We categorize the variables \mathbf{X} into evidence $\mathbf{E} \subseteq \mathbf{X}$ and query $\mathbf{Q} \subseteq \mathbf{X}$ groups, ensuring that $\mathbf{E} \cap \mathbf{Q} = \emptyset$ and $\mathbf{E} \cup \mathbf{Q} = \mathbf{X}$. Then, given an assignment \mathbf{e} to the set of evidence variables \mathbf{E} , the MPE task can be formulated as:

$$\text{MPE}(\mathbf{Q}, \mathbf{e}) = \underset{\mathbf{q}}{\operatorname{argmax}} p_{\mathcal{M}}(\mathbf{q}|\mathbf{e}) = \underset{\mathbf{q}}{\operatorname{argmax}} \{\log p_{\mathcal{M}}(\mathbf{q}, \mathbf{e})\} \quad (1)$$

It is known that the MPE task is NP-hard in general and even hard to approximate [9, 11, 38, 41, 44].

Motivation: The goal of this paper is to develop a method that trains a NN for a given PM and, at test time, serves as an approximate MPE solver for any-MPE query posed over the PM. By any-MPE, we mean that the NN can take an assignment to an arbitrary subset of variables (evidence) as input and output the most likely assignments to the remaining (query) variables. Recently, Arya et al. [4] proposed a NN-based solution for solving the MPE task in PCs under the constraint that the partition of the variables into evidence and query sets *is known before training the NN*. This constraint is highly restrictive because, for generative models, it is unlikely that such a partition of variables is known in advance. In such cases, one would typically train a discriminative model rather than a generative one. Unlike Arya et al.’s method, our approach yields an any-MPE solver. Additionally, Arya et al.’s approach has several limitations in that it does not fully exploit the benefits of self-supervision during inference time and requires the use of relatively large NNs to achieve good performance in practice. Our proposed approach, described next, addresses these limitations.

3 A Self-Supervised Neural Approximator for any-MPE

In this section, we develop a neural network (NN) based approach for solving the *any-MPE* task. Specifically, given a PM, we develop an input encoding (see Section 3.1) that determines the number of input nodes of the NN and sets their values for the given MPE query. We also develop an output encoding that determines the number of output nodes for the given PM and allows us to recover the MPE solution from the outputs. For training the NN, we introduce a tractable and differentiable self-supervised loss function (see Section 3.2), whose global minima aligns with the MPE solutions to efficiently learn the parameters of the NN given *unlabeled data*.

3.1 An Encoding For any-MPE Instances

Since NNs require fixed-sized inputs and outputs, we introduce input and output encodings that generate fixed-length input and output vectors for each PM from a given MPE problem instance $\text{MPE}(\mathbf{Q}, \mathbf{e})$. To encode the input, for each variable $X_i \in \mathbf{X}$, we associate two input nodes in the NN, denoted by \hat{X}_i and \bar{X}_i . Thus for a PM having n (namely, $|\mathbf{X}| = n$) variables, the corresponding NN has $2n$ input nodes. Given a query $\text{MPE}(\mathbf{Q}, \mathbf{e})$, we set the values of the input nodes as follows: (1) If $X_i \in \mathbf{E}$ and $X_i = 0$ is in \mathbf{e} , then we set $\hat{X}_i = 0$ and $\bar{X}_i = 1$; (2) If $X_i \in \mathbf{E}$ and $X_i = 1$ is in \mathbf{e} , then we set $\hat{X}_i = 1$ and $\bar{X}_i = 0$; and (3) If $X_i \in \mathbf{Q}$ then we set $\hat{X}_i = 0$ and $\bar{X}_i = 0$. (The assignment $\hat{X}_i = 1$ and $\bar{X}_i = 1$ is not used.) It is easy to see that the input encoding described above yields an *injective* mapping between the set of all possible MPE queries over the given PM and the set $\{0, 1\}^{2n}$. This means that each unique MPE query (\mathbf{Q}, \mathbf{e}) will yield a unique 0-1 input vector of size $2n$.

The output of the neural network comprises of n nodes with sigmoid activation, where each output node is associated with a variable $X_i \in \mathbf{X}$. We ignore the outputs corresponding to the evidence variables and define a loss function over the outputs corresponding to the query variables in the set \mathbf{Q} . The MPE solution can be reconstructed from the output nodes of the NN by thresholding the

²or a value proportional to it such as the unnormalized probability in Markov networks.

output nodes corresponding to the query variables appropriately (e.g., if the value of the output node is greater than 0.5, then the query variable is assigned the value 1; otherwise it is assigned to 0).

3.2 A Self-Supervised Loss Function for any-MPE

Since the output nodes of our proposed NN use sigmoid activation, each output is continuous and lies in the range $[0, 1]$. Given an MPE query $\text{MPE}(\mathbf{Q}, \mathbf{e})$, let $\mathbf{q}^c \in [0, 1]^{|Q|}$ denote the (continuous) *Most Probable Explanation* (MPE) assignment predicted by the NN. In MPE inference, given \mathbf{e} , we want to find an assignment \mathbf{q} such that $\log p_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$ is maximized, namely, $-\log p_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$ is minimized. Thus, a natural loss function that we can use is $-\log p_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$. Unfortunately, the NN outputs a continuous vector \mathbf{q}^c and as a result $p_{\mathcal{M}}(\mathbf{q}^c, \mathbf{e})$ is not defined.

Next, we describe how to solve the above problem by leveraging the following property of the class of PMs that we consider in this paper—specifically BNs, MNs, PCs and NAMs. In these PMs, the function $\ell(\mathbf{q}, \mathbf{e}) = -\log p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$, which is a function from $\{0, 1\}^n \rightarrow \mathbb{R}$ is either a multi-linear polynomial or a neural network, and can be computed in linear time in the size of the PM. To facilitate the use of continuous outputs, we define a loss function $\ell^c(\mathbf{q}^c, \mathbf{e}) : [0, 1]^n \rightarrow \mathbb{R}$ such that ℓ^c coincides with ℓ on $\{0, 1\}^n$. For PGMs and PCs, ℓ is a multi-linear function and ℓ^c is obtained by substituting each occurrence of a discrete variable $q_i \in \mathbf{q}$ with the corresponding continuous variable $q_i^c \in \mathbf{q}^c$ where $q_i^c \in [0, 1]$. In NAMs, ℓ is a NN and we can perform a similar substitution—we substitute each binary input q_i in the NN with a continuous variable $q_i^c \in [0, 1]$. This substitution transforms the discrete NN into a continuous function while preserving its functional form.

An important property of ℓ^c is that it can be evaluated and differentiated in polynomial time. Moreover, when ℓ is defined by either a neural network (in NAMs) or a multilinear function (in BNs, MNs and PCs), the minimum value of ℓ^c over the domain $[0, 1]^n$ is less than or equal to the minimum value of the original function ℓ over the discrete domain $\{0, 1\}^n$. Formally,

Proposition 1. *Let $\ell(\mathbf{q}, \mathbf{e}) : \{0, 1\}^n \rightarrow \mathbb{R}$ be either a neural network or a multilinear function, and let $\ell^c(\mathbf{q}^c, \mathbf{e}) : [0, 1]^n \rightarrow \mathbb{R}$ be its continuous extension obtained by substituting each binary input q_i with a continuous variable $q_i^c \in [0, 1]$. Then,*

$$\min_{\mathbf{q}^c \in [0, 1]^n} \ell^c(\mathbf{q}^c, \mathbf{e}) \leq \min_{\mathbf{q} \in \{0, 1\}^n} \ell(\mathbf{q}, \mathbf{e})$$

Following Arya et al. [4], we propose to improve the quality of the loss function by tightening the lower bound given in proposition 1 with an entropy-based penalty (ℓ_E), governed by $\alpha > 0$.

$$\ell_E(\mathbf{q}^c, \alpha) = -\alpha \sum_{j=1}^{|Q|} [q_j^c \log(q_j^c) + (1 - q_j^c) \log(1 - q_j^c)] \quad (2)$$

This penalty encourages discrete solutions by preferring q_j^c values close to 0 or 1, where α modulates the trade-off. Setting α to 0 yields the continuous approximation; conversely, an α value of ∞ results exclusively in discrete outcomes. From proposition 1 and by using the theory of Lagrange multipliers, we can show that for any $\alpha > 0$, the use of the entropy penalty yields a tighter lower bound:

Proposition 2.

$$\min_{\mathbf{q}^c \in [0, 1]^n} \ell^c(\mathbf{q}^c, \mathbf{e}) \leq \min_{\mathbf{q}^c \in [0, 1]^n} \ell^c(\mathbf{q}^c, \mathbf{e}) + \ell_E(\mathbf{q}^c, \alpha) \leq \min_{\mathbf{q} \in \{0, 1\}^n} \ell(\mathbf{q}, \mathbf{e})$$

How to use the Loss Function: Given a PM defined over n variables, we can use the self-supervised loss function $\ell^c(\mathbf{q}^c, \mathbf{e}) + \ell_E(\mathbf{q}^c, \alpha)$ (treating α as a hyper-parameter) to train any neural network (NN) architecture that has $2n$ input nodes and n output nodes. This trained NN can then be used to answer any arbitrary MPE query posed over the PM. The training data for the neural network consists of assignments (evidence \mathbf{e}) to a subset of the variables. Each training example can be generated using the following three-step process. We first sample a full assignment \mathbf{x} to all variables in the PM using techniques like Gibbs sampling or perfect sampling for tractable distributions such as PCs and BNs. Second, we choose an integer k uniformly at random from the range $\{1, \dots, n\}$ and designate k randomly selected variables as evidence variables \mathbf{E} , and the remaining $n - k$ as query variables \mathbf{Q} . Finally, we project the full assignment \mathbf{x} on \mathbf{E} . The primary advantage of using the self-supervised loss function is that it eliminates the need for access to a dedicated MPE solver to provide supervision during training; gradient-based training of the neural network provides the necessary supervision.

3.3 Inference-Time Neural Optimization using Self-Supervised Loss

At a high level, assuming that the NN is over-parameterized, if we use the self-supervised loss and repeatedly run (stochastic) gradient updates over the NN for a given dataset, theoretical results [2, 13] as well as prior experimental work [46, 56] suggest that the parameters of the NN will converge to a point near the global minimum of the self-supervised loss function. This means that through gradient updates, the network will find a near-optimal MPE assignment for each training example. This strategy of performing gradient updates over the NN can also be used *during inference (test) time to iteratively improve the MPE solution*, thereby maximizing the benefits of self-supervision.

Specifically, at test time, given a test dataset (or example), we initialize the NN either randomly or using a pre-trained model and then run gradient-based updates over the NN iteratively until convergence. The gradient is computed w.r.t. to the self-supervised loss function $\ell^c(\mathbf{q}^c, \mathbf{e}) + \ell_E(\mathbf{q}^c, \alpha)$. We call the resulting algorithm ITSELF (Inference Time Optimization using SELF-Supervised Loss). The performance of ITSELF typically improves with each iteration until the loss converges.

Our proposed method, ITSELF, is closely related to test-time training approaches which are widely used to solve problems in deep learning [1, 10, 19, 30–32, 40, 49, 51, 57]. Our method differs from these previous approaches in that the global minima of our proposed self-supervised loss correspond to the MPE solutions, provided that the penalty α is sufficiently large.

4 Supervised Knowledge Transfer from ITSELF

A drawback of our self-supervised loss function is that, unlike supervised loss functions such as binary cross entropy, it is a non-convex function of the NN outputs³. As a result, it has a significantly larger number of local minima compared to the supervised loss function, but also a potentially exponential number of global minima, because an MPE problem can have multiple optimal solutions [37], all of which have the same loss function value. Thus, optimizing and regularizing using self-supervised loss is difficult compared to supervised loss, especially when the number of training examples is large.

Moreover, our experiments show that large datasets necessitate large, over-parameterized neural networks (NNs) to achieve near-optimal MPE solutions for all examples. However, when the training data is limited and the NN is sufficiently over-parameterized, our preliminary findings, along with theoretical and empirical results from prior studies [3, 6, 23, 27, 28], suggest that the NN is more likely to approach the global optima. Specifically, with a reasonably sized NN and a small dataset, the algorithm ITSELF tends to yield near-optimal MPE solutions. A further challenge with ITSELF is that even for small datasets, achieving convergence from a random initialization requires numerous iterations of gradient descent, rendering the training process inefficient and slow.

4.1 Teacher-Student Strategy

To address these challenges (using small datasets with ITSELF; designing better initialization for it; and using non-convex loss functions for training), we propose a two-network teacher-student strategy [7, 16, 20–22, 24, 39, 47, 52–54], where we have two networks with the same structure that are trained via mini-batch gradient updates. The teacher network is overfitted to the mini-batch using our self-supervised loss via the algorithm ITSELF, and the student network is subsequently trained with a supervised loss function such as binary cross entropy. By overfitting the teacher network via ITSELF on the mini-batch, we ensure that it finds near-optimal MPE assignments for all (unlabeled) examples in the mini-batch and eventually over the whole training dataset.

The student network then learns from the teacher’s outputs, using them as soft labels in a supervised learning framework. This transfer of knowledge mitigates the optimization difficulties associated with the non-convex self-supervised loss, allowing the student network to achieve faster convergence and better generalization with a more manageable model size. Additionally, this strategy reduces the need for severe over-parameterization and extensive training iterations for the teacher network because it is operating on a smaller dataset. It also helps achieve better initialization for ITSELF.

³Note that we are referring to convexity with respect to the outputs, not the parameters of the NN.

Algorithm 1 GUIDed Iterative Dual LEarning with Self-supervised Teacher (*GUIDE*)

```
1: Input: Training data  $\mathcal{D}$ , teacher  $\mathcal{T}$  and student  $\mathcal{S}$  having the same structure
2: Output: Trained student network  $\mathcal{S}$ 
3:    $\triangleright$  Database  $DB$  stores the best MPE assignment and loss value for each example in  $\mathcal{D}$ 
4: Initialize: Randomly initialize  $\mathcal{T}$ ,  $\mathcal{S}$ , and  $DB$ 
5: for each epoch do
6:   Sample a mini-batch  $\mathcal{D}'$  from  $\mathcal{D}$ 
7:   Update the parameters of  $\mathcal{T}$  using the algorithm ITSELF (self-supervised loss) with Dataset  $\mathcal{D}'$ 
8:   for each example  $\mathbf{e}_i$  in  $\mathcal{D}'$  do
9:     Make a forward-pass over  $\mathcal{T}$  to get an MPE assignment  $\mathbf{q}_i$  for  $\mathbf{e}_i$ 
10:    Update the entry in  $DB$  for  $\mathbf{e}_i$  with  $\mathbf{q}_i$  if it has a lower loss value than the current entry
11:   end for
12:   Update the parameters of  $\mathcal{S}$  using the mini-batch  $\mathcal{D}'$  and labels from  $DB$  and a supervised loss
13:    $\mathcal{T} \leftarrow \mathcal{S}$   $\triangleright$  Initialize  $\mathcal{T}$  with  $\mathcal{S}$  for the next epoch
14: end for
```

4.2 Training Procedure

Our proposed training procedure, which we call *GUIDE*, is detailed in Algorithm 1. The algorithm trains a two-network system comprising a teacher network (\mathcal{T}) and a student network (\mathcal{S}) with the same structure. The goal is to train the student network using a combination of self-supervised and supervised learning strategies. The algorithm takes as input the training data \mathcal{D} , along with the teacher and student networks, \mathcal{T} and \mathcal{S} , respectively and outputs a trained network \mathcal{S} . A database (DB) is utilized to store the best MPE assignment and corresponding loss value for each example in \mathcal{D} . The parameters of \mathcal{T} and \mathcal{S} , and the entries in DB , are randomly initialized at the start.

In each epoch, a mini-batch \mathcal{D}' is sampled from the training data \mathcal{D} . The parameters of the teacher network \mathcal{T} are then updated using the algorithm ITSELF (which uses a self-supervised loss), applied to the mini-batch \mathcal{D}' (the mini-batch helps address large data issues associated with ITSELF). For each example \mathbf{e}_i in \mathcal{D}' , we perform a forward-pass over \mathcal{T} to obtain an MPE assignment \mathbf{q}_i . The database DB is subsequently updated with \mathbf{q}_i if it has a lower loss value than the current entry for \mathbf{e}_i .

Following this, the parameters of the student network \mathcal{S} are updated using the mini-batch \mathcal{D}' , the labels from DB , and a supervised loss function such as Binary Cross Entropy or $L2$ loss. Finally, the parameters of the teacher network \mathcal{T} are reinitialized with the updated parameters of the student network \mathcal{S} to prepare for the next epoch (addressing the initialization issue associated with ITSELF).

Thus, at a high level, Algorithm 1 leverages the strengths of both self-supervised and supervised learning to improve training efficiency and reduce the model complexity, yielding a student network \mathcal{S} . Moreover, at test time, the student network can serve as an initialization for ITSELF.

5 Experiments

This section evaluates the ITSELF method (see section 3.3), the *GUIDE* teacher-student training method (see section 4) and the method that uses only self-supervised training, which we call SSMP (see section 3.2). We benchmark these against various baselines, including neural network-based and traditional polynomial-time algorithms that directly operate on the probabilistic model. We begin by detailing our experimental framework, including competing methods, evaluation metrics, neural network architectures, and datasets.

5.1 Datasets and Graphical Models

We used twenty binary datasets extensively used in tractable probabilistic models literature [5, 18, 34, 50]—referred to as TPM datasets—for evaluating Probabilistic Circuits (PC)s and Neural Auto-Regressive Model (NAM). For the purpose of evaluating Probabilistic Graphical Models (PGM)s, we utilized high treewidth models from previous UAI inference competitions [14].

To train Sum Product Networks (SPNs), our choice of PCs, we employed the DeeProb-kit library [33], with SPN sizes ranging from 46 to 9666 nodes. For NAM, Masked Autoencoder for Distribution

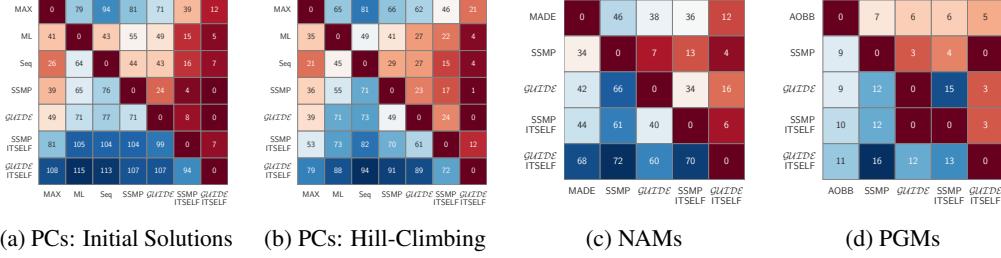


Figure 1: Contingency Tables: Comparing Methods for MPE Problems Across PMs

Estimation (MADE) models were trained using PyTorch as described in [17]. In the case of Markov Networks (MNs), a specific category of PGMs, we utilize Gibbs Sampling, generating 8000, 1000, and 1000 examples for the training, testing, and validation sets, respectively. The query ratio (qr) is defined as the fraction of variables in the query set. For each PM, we varied qr with values from the set $\{0.1, 0.3, 0.5, 0.7, 0.8, 0.9\}$.

5.2 Baseline Methods and Evaluation Criteria

PCs - We employed three polynomial-time baseline methods from the PC and PGM literature as initial benchmarks [41, 45]. **MAX Approximation (MAX)** [45] transforms sum nodes into max nodes. During the upward pass, max nodes output the highest weighted value from their children. The downward pass, starting from the root, selects the child with the highest value at each max node and includes all children of product nodes. **Maximum Likelihood Approximation (ML)** [41] computes the marginal distribution $p_{\mathcal{M}}(Q_i|\mathbf{e})$ for each variable $Q_i \in \mathbf{Q}$, setting Q_i to its most likely value. **Sequential Approximation (Seq)** [41] iteratively assigns query variables according to an order o . At each step j , it selects the j -th query variable Q_j in o and assigns to it a value q_j such that $p_{\mathcal{M}}(q_j|\mathbf{e}, \mathbf{y})$ is maximized, where \mathbf{y} is an assignment of values to all query variables from 1 to $j - 1$. Our study further assessed the impact of initializing stochastic hill climbing searches using solutions from all baseline approaches and our proposed methods for MPE inference, conducting 60-second searches for each MPE problem in our experiments, as detailed in Park and Darwiche [41].

NAMs - As a baseline, we used the stochastic hill-climbing search (HC). Similarly to the PC procedure, for each test example, we conducted a 60-second hill-climbing search, initializing query variables randomly and setting evidence variables based on the example.

PGMs - As a baseline we utilize the AND/OR Branch-and-Bound (AOBB)[36], employing the implementation described in Marinescu [35]. Given that AOBB is an anytime scheme, for each test example, we designated a 60-second time limit for the inference process.

Neural Baselines - Arya et al. [4] introduced Self-Supervised learning based MMAP solver for PCs (SSMP), which trains a neural network to handle queries on a fixed partition of variables on PCs. We proposed an extension of this approach for solving the any-MPE task in PMs (see section 3.2), where a single neural network is trained to answer any-MPE query. We use it as another neural baseline.

Evaluation Criteria - Competing approaches were compared using log-likelihood (LL) scores, calculated as $\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$, and inference times for a given evidence \mathbf{e} and query output \mathbf{q} .

5.3 Neural Network-Based Approaches

For each PM and query ratio, we implemented two neural network training protocols: SSMP and *GUIDE*. We subjected each model to 20 training epochs, adhering to the established procedures for SSMP as delineated by Arya et al. [4]. Both protocols employed two distinct inference strategies, thus forming four neural-based variants. The first strategy consisted of a single forward pass through the network to estimate query variable values, as specified by Arya et al. [4]. The second strategy utilized our novel test-time optimization-based *ITSELF* approach for inference. For *ITSELF*, we undertook 100 optimization iterations or terminated earlier upon achieving loss convergence, applicable to both PCs and PGMs. In cases involving NAMs, the optimization iterations extended to 1,000, or concluded upon convergence.

We standardized network architectures for PMs across all experiments. For PCs, we used a fully connected Neural Networks (NN) with three hidden layers (128, 256, 512 nodes). For NAMs and PGMs, a single hidden layer of 512 nodes was employed. All hidden layers featured ReLU activation, while the output layers used sigmoid functions with dropout for regularization [48]. Optimization was performed using the Adam optimizer [25], and models were implemented in PyTorch [43] on an NVIDIA A40 GPU.

Results for PCs: We compare methods—including three polynomial-time baselines, neural network-based SSMP, and our **ITSELF** and *GUIDE* methods—on 20 TPM datasets as shown in the contingency table in figure 1a (detailed results in the supplementary materials). We generated 120 test datasets for the MPE task using 20 PCs across 6 query ratios (qr). Each cell (i, j) in the table represents how often (out of 120) the method in row i outperformed the method in column j based on average log-likelihood scores. Any difference between 120 and the combined frequencies of cells (i, j) and (j, i) indicates cases where the compared methods achieved similar scores.

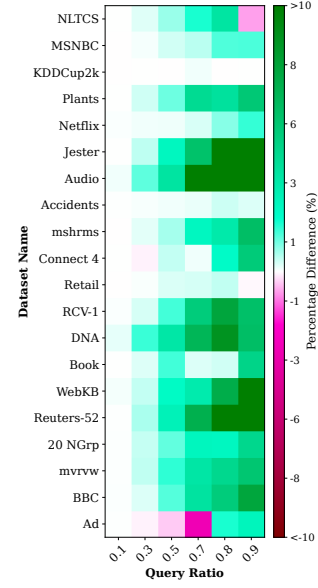
The contingency table for PC shows that the **ITSELF** methods outperform polynomial-time and traditional neural baselines. Specifically, *GUIDE* + **ITSELF** is superior to all the other methods in almost two-thirds of the 120 cases, while SSMP + **ITSELF** is better than both SSMP and *GUIDE* without **ITSELF**. However, the polynomial-time baseline MAX is better than both SSMP and *GUIDE* (note that these methods were used in Arya et al. [4]), highlighting **ITSELF**'s significant role in boosting model performance for the complex *any-MPE* task.

We compare MAX and *GUIDE* + **ITSELF** using a heatmap in Figure 2a. The y-axis presents datasets by variable count and the x-axis by query ratio. Each cell displays the percentage difference in mean LL scores between the methods, calculated as $\%Diff. = 100 \times (l_{nn} - l_{max}) / |l_{max}|$. From the heatmap in Figure 2a for PCs, we observe that *GUIDE* + **ITSELF** exhibit performance comparable to the MAX approximation when the query set size is small. As the problem complexity increases with an increase in query set size, our new method consistently outperforms MAX across all datasets, except for NLTCS and Tretail. Even when *GUIDE* + **ITSELF** underperforms compared to MAX (noted in 12 instances out of 120) the performance gap is very small, evidenced by sparse red cells in the heatmaps.

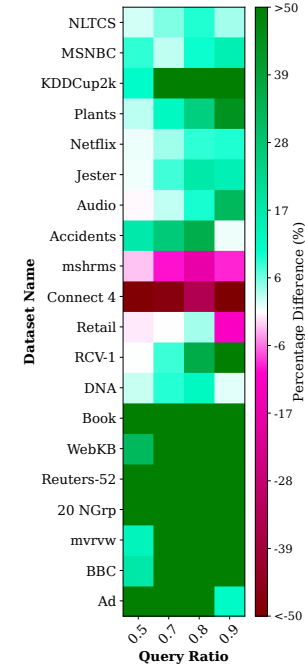
Further analysis in Figure 1b compares the performance of our proposed methods with various baselines as initialization strategies for Hill Climbing Search. The goal is to assess whether **ITSELF** and *GUIDE* enhance *anytime methods* to outperform other heuristic initialization techniques, with methods utilizing **ITSELF** proving superior in initializing local search-based algorithms.

Results for NAMs: Our evaluation, detailed in the contingency table in Figure 1c, examines the performance of several methods for NAM, including HC and two neural network methods, SSMP and *GUIDE*, each with two inference schemes. We tested four neural-based schemes on 20 TPM datasets, generating 80 test datasets for the MPE task using 20 MADEs across four query ratios (qr).

Similar to the results observed with PC, the *GUIDE* + **ITSELF** method exhibits superior performance over the baselines and other neural inference schemes. HC outperforms SSMP, while *GUIDE* and SSMP + **ITSELF** are superior to HC.



(a) PC: *GUIDE* + **ITSELF** vs. MAX



(b) NAM: *GUIDE* + **ITSELF** vs. HC

Figure 2: Heatmaps showing LL Percentage Differences: Top—PC; Bottom—MADE.

The heatmaps in Figure 2b highlight the superior performance of *GUIDE* + *ITSELF* for NAMs, particularly in larger datasets where it outperforms the HC baseline by over 50% in most cases, as indicated by the dark green cells. The integration of *GUIDE*-based learning with *ITSELF*-based inference consistently outperforms the baseline across most scenarios, with exceptions only in the Mushrooms, Connect 4, and Retail datasets. Thus, the *GUIDE* + *ITSELF* approach significantly enhances MPE query answering in NAM models.

Results for PGMs: In our evaluation, as detailed in the contingency table in 1d, we assessed the efficacy of various methods for PGMs, including AOBB and four neural-network-based methods, across four high-treewidth networks (details of these networks are provided in the supplement). To accomplish this, we constructed 16 test datasets for the MPE task by employing four PGMs across four query ratios (qr).

Similar to the outcomes observed with previous PMs, the methods employing *ITSELF* for inference consistently demonstrate enhanced performance compared to the baseline methods AOBB and SSMP in most scenarios. Both *GUIDE* and SSMP outperform AOBB in at least 50 percent of the tests.

Is Teacher Student Training Better than A Single Network Trained with Self-Supervised Loss? (SSMP versus *GUIDE*): In our comparative analysis between SSMP and *GUIDE* across different models, we aim to evaluate the performance of *GUIDE* against the traditional neural network training methods used in SSMP. Using traditional inference schemes (i.e., one forward pass through the network), *GUIDE* consistently outperforms SSMP, demonstrating its superiority in 60% of scenarios for PCs, more than 80% for NAM models, and 75% for PGM models. When employing *ITSELF* for inference over the trained models, *GUIDE* continues to outperform SSMP, achieving better results in more than 75%, 85%, and 80% for PCs, NAMs, and PGMs, respectively.

Does Inference Time Optimization help? (One Pass versus Multiple Passes): Comparative analyses reveal that *ITSELF* consistently outperforms traditional single forward pass inference across various PMs. *ITSELF* with SSMP training outperforms the other methods in over 85% of PC cases, and more than 75% for NAM and PGM models. When trained using *GUIDE*, *ITSELF* demonstrates even better results, achieving superior performance in nearly 90% of PC cases and 75% for both NAMs and PGMs. *GUIDE* with the *ITSELF* inference scheme emerges as the best method across all our experiments.

Finally, we present inference times in the supplement. For all the PMs, the neural based method that utilizes traditional inference are the quickest. For MADE, the method employing a *GUIDE*-trained model with *ITSELF* is the second fastest. Similarly, in PGMs, the *GUIDE* + *ITSELF* method emerges as the third fastest, after SSMP + *ITSELF*; in PCs, MAX leads, slightly ahead of both *GUIDE* + *ITSELF* and SSMP + *ITSELF*, while ML and Seq record the longest inference times.

Summary: Our experiments demonstrate that *GUIDE* + *ITSELF* outperforms both polynomial-time and neural-based baselines across various PMs, as evidenced by higher log-likelihood scores. Specifically, *ITSELF* exceeds the capabilities of traditional forward pass inference in addressing the challenging task of answering *any-MPE* queries within a probabilistic model, highlighting the necessity of Inference Time Optimization. Additionally, the superiority of models trained with *GUIDE* over SSMP underscores the effectiveness of using a dual network approach that utilizes two loss functions to enhance the initial model quality and provide an optimal starting point for *ITSELF*.

6 Conclusion and Future Work

We introduced novel methods for answering Most Probable Explanation (MPE) queries in probabilistic models. Our approach employs self-supervised loss functions to represent MPE objectives, enabling tractable loss and gradient computations during neural network training. We also proposed a new inference time optimization technique, *ITSELF*, which iteratively improves the solution to the MPE problem via gradient updates. Additionally, we introduced a dual-network-based strategy that combines supervised and unsupervised training which we call *GUIDE* to provide better initialization for *ITSELF* and addressing various challenges associated with self-supervised training. Our method was tested on various benchmarks, including probabilistic circuits, MADE, and PGMs, using 20 binary datasets and high tree-width networks. It outperformed polytime baselines and other neural

methods, substantially in some cases. Additionally, it improved the effectiveness of stochastic hill climbing (local) search strategies.

Future work includes solving complex queries in Probabilistic Models (PM)s with constraints; training neural networks with losses from multiple PMs to embed their inference mechanisms; boosting performance by developing advanced encoding strategies for similar tasks; implementing sophisticated neural architectures tailored to PMs; etc.

References

- [1] Ferran Alet, Maria Bauza, Kenji Kawaguchi, Nurullah Giray Kuru, Tomás Lozano-Pérez, and Leslie Kaelbling. Tailoring: Encoding inductive biases by optimizing unsupervised objectives at prediction time. In *Advances in Neural Information Processing Systems*, volume 34, pages 29206–29217. Curran Associates, Inc., 2021.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019.
- [3] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 322–332. PMLR, 2019.
- [4] Shivvrat Arya, Tahrima Rahman, and Vibhav Gogate. Neural Network Approximators for Marginal MAP in Probabilistic Circuits. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(10):10918–10926, March 2024. ISSN 2374-3468. doi: 10.1609/aaai.v38i10.28966.
- [5] Jessa Bekker, Jesse Davis, Arthur Choi, Adnan Darwiche, and Guy Van den Broeck. Tractable learning for complex probability queries. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [6] Lénaïc Chizat and Francis R. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 3040–3050, 2018.
- [7] Jang Hyun Cho and Bharath Hariharan. On the Efficacy of Knowledge Distillation. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4793–4801, October 2019. doi: 10.1109/ICCV.2019.00489.
- [8] YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. oct 2020.
- [9] Gregory F. Cooper. The computational complexity of probabilistic inference using bayesian belief networks. *Artificial Intelligence*, 42(2-3):393–405, March 1990. ISSN 00043702. doi: 10.1016/0004-3702(90)90060-D.
- [10] Mohammad Zalbagi Darestani, Jiayu Liu, and Reinhard Heckel. Test-time training can close the natural distribution shift performance gap in deep learning based compressed sensing. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 4754–4776. PMLR, 17–23 Jul 2022.
- [11] Cassio P de Campos. New Complexity Results for MAP in Bayesian Networks. *IJCAI International Joint Conference on Artificial Intelligence*, pages 2100–2106, 01 2011. doi: 10.5591/978-1-57735-516-8/IJCAI11-351.
- [12] Priya L. Donti, David Rolnick, and J. Zico Kolter. DC3: A learning method for optimization with hard constraints. In *International Conference on Learning Representations*, October 2020.

- [13] S. Du, Xiyu Zhai, B. Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *International Conference on Learning Representations*, 2018.
- [14] G. Elidan and A. Globerson. *The 2010 UAI Approximate Inference Challenge*. 2010.
- [15] Ferdinando Fioretto, Terrence W. K. Mak, and Pascal Van Hentenryck. Predicting AC Optimal Power Flows: Combining Deep Learning and Lagrangian Dual Methods. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):630–637, April 2020. ISSN 2374-3468. doi: 10.1609/aaai.v34i01.5403.
- [16] Tommaso Furlanello, Zachary Chase Lipton, M. Tschannen, L. Itti, and Anima Anandkumar. Born Again Neural Networks. In *International Conference on Machine Learning*, May 2018.
- [17] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 881–889, Lille, France, 07–09 Jul 2015. PMLR.
- [18] Jan Van Haaren and Jesse Davis. Markov Network Structure Learning: A Randomized Feature Generation Approach. *Proceedings of the AAAI Conference on Artificial Intelligence*, 26(1): 1148–1154, 2012. ISSN 2374-3468. doi: 10.1609/aaai.v26i1.8315.
- [19] Moritz Hardt and Yu Sun. Test-time training on nearest neighbors for large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [20] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A Comprehensive Overhaul of Feature Distillation. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1921–1930, October 2019. doi: 10.1109/ICCV.2019.00201.
- [21] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge Transfer via Distillation of Activation Boundaries Formed by Hidden Neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3779–3787, July 2019. doi: 10.1609/aaai.v33i01.33013779.
- [22] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015.
- [23] Arthur Jacot, Clément Hongler, and Franck Gabriel. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8580–8589, 2018.
- [24] Jangho Kim, Seonguk Park, and Nojun Kwak. Paraphrasing Complex Network: Network Compression via Factor Transfer. *ArXiv*, February 2018.
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [26] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [27] Jaehoon Lee, S. Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Narain Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Neural Information Processing Systems*, 2020.
- [28] Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(12):124002, December 2020. ISSN 1742-5468. doi: 10.1088/1742-5468/abc62b.
- [29] Ke Li and Jitendra Malik. Learning to optimize. *International Conference on Learning Representations*, 2016.

- [30] Yushu Li, Xun Xu, Yongyi Su, and Kui Jia. On the Robustness of Open-World Test-Time Training: Self-Training with Dynamic Prototype Expansion. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11802–11812, Paris, France, October 2023. IEEE. ISBN 9798350307184. doi: 10.1109/ICCV51070.2023.01087.
- [31] Huan Liu, Zijun Wu, Liangyan Li, Sadaf Salehkalaibar, Jun Chen, and Keyan Wang. Towards Multi-domain Single Image Dehazing via Test-time Training. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5821–5830, New Orleans, LA, USA, June 2022. IEEE. ISBN 978-1-66546-946-3. doi: 10.1109/CVPR52688.2022.00574.
- [32] Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre Alahi. TTT++: When Does Self-Supervised Test-Time Training Fail or Thrive? In *Advances in Neural Information Processing Systems*, volume 34, pages 21808–21820. Curran Associates, Inc., 2021.
- [33] Lorenzo Loconte and Gennaro Gala. DeeProb-kit: a python library for deep probabilistic modelling, 2022. URL <https://github.com/deeprob-org/deeprob-kit>.
- [34] Daniel Lowd and Jesse Davis. Learning Markov Network Structure with Decision Trees. In *2010 IEEE International Conference on Data Mining*, pages 334–343. IEEE, 2010. ISBN 978-1-4244-9131-5. doi: 10.1109/ICDM.2010.128.
- [35] Radu Marinescu. Daoopt: Distributed and/or optimization-uai’10 inference competition. 2010. URL <https://github.com/lotten/daoopt>.
- [36] Radu Marinescu and Rina Dechter. AND/OR Branch-and-Bound search for combinatorial optimization in graphical models. *Artificial Intelligence*, 173(16-17):1457–1491, November 2009. ISSN 00043702. doi: 10.1016/j.artint.2009.07.003.
- [37] Radu Marinescu and Rina Dechter. Counting the optimal solutions in graphical models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/fc2e6a440b94f64831840137698021e1-Paper.pdf.
- [38] Denis Deratani Mauá and Cassio P. de Campos. Approximation complexity of maximum A posteriori inference in sum-product networks. *CoRR*, abs/1703.06045, 2017.
- [39] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved Knowledge Distillation via Teacher Assistant. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5191–5198, April 2020. doi: 10.1609/aaai.v34i04.5963.
- [40] David Osowiechi, G. A. V. Hakim, Mehrdad Noori, Milad Cheraghalikhani, Ismail Ben Ayed, and Christian Desrosiers. Tttflow: Unsupervised test-time training with normalizing flow. *IEEE Workshop/Winter Conference on Applications of Computer Vision*, 2022. doi: 10.48550/arXiv.2210.11389.
- [41] James D. Park and Adnan Darwiche. Complexity results and approximation strategies for map explanations. *J. Artif. Int. Res.*, 21(1):101–133, feb 2004. ISSN 1076-9757.
- [42] Seonho Park and Pascal Van Hentenryck. Self-Supervised Primal-Dual Learning for Constrained Optimization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(4):4052–4060, June 2023. ISSN 2374-3468. doi: 10.1609/aaai.v37i4.25520.
- [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [44] Robert Peharz. *Foundations of Sum-Product Networks for Probabilistic Modeling*. PhD thesis, PhD thesis, Medical University of Graz, 2015.

- [45] Hoifung Poon and Pedro Domingos. Sum-Product Networks: A New Deep Architecture. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pages 337–346. AUAI Press, 2011.
- [46] Tiancheng Qin, S. Rasoul Etesami, and Cesar A Uribe. Faster convergence of local SGD for over-parameterized models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- [47] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, C. Gatta, and Yoshua Bengio. FitNets: Hints for Thin Deep Nets. *CoRR*, December 2014.
- [48] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014. ISSN 1533-7928.
- [49] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-Time Training with Self-Supervision for Generalization under Distribution Shifts. In *Proceedings of the 37th International Conference on Machine Learning*, pages 9229–9248. PMLR, November 2020.
- [50] Benigno Uribe, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. Neural autoregressive distribution estimation. *Journal of Machine Learning Research*, 17(205):1–37, 2016.
- [51] Dequan Wang, Evan Shelhamer, Shaoteng Liu, B. Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *International Conference on Learning Representations*, 2021.
- [52] Chenglin Yang, Lingxi Xie, Chi Su, and Alan L. Yuille. Snapshot Distillation: Teacher-Student Optimization in One Generation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2854–2863, June 2019. doi: 10.1109/CVPR.2019.00297.
- [53] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A Gift from Knowledge Distillation: Fast Optimization, Network Minimization and Transfer Learning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7130–7138, July 2017. doi: 10.1109/CVPR.2017.754.
- [54] Sergey Zagoruyko and N. Komodakis. Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer. *ArXiv*, November 2016.
- [55] Ahmed S. Zamzam and Kyri Baker. Learning Optimal Solutions for Extremely Fast AC Optimal Power Flow. In *2020 IEEE International Conference on Communications, Control, and Computing Technologies for Smart Grids (SmartGridComm)*, pages 1–6, November 2020. doi: 10.1109/SmartGridComm47815.2020.9303008.
- [56] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [57] Wentao Zhu, Yufang Huang, Daguang Xu, Zhen Qian, Wei Fan, and Xiaohui Xie. Test-time training for deformable multi-scale image registration. *IEEE International Conference on Robotics and Automation*, 2021.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: [\[NA\]](#)

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: [\[NA\]](#)

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Our paper improves the speed, accuracy and scalability of MPE inference algorithms. Since algorithms already exist for solving MPE tasks, we do not perceive any negative or positive societal impacts beyond what currently exists.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: [\[NA\]](#)

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: [\[NA\]](#)

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[NA\]](#)

Justification: [\[NA\]](#)

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.