# 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
shivvrat.arya@utdallas.edu

**Tahrima Rahman**
Department of Computer Science
The University of Texas at Dallas
tahrima.rahman@utdallas.edu

**Vibhav Gogate**
Department of Computer Science
The University of Texas at Dallas
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 [39, 40] 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, PGMs 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., $\mathbf{x}$).

---

[1]Arya 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 likelihood[2] 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:

$$\mathrm{MPE}(\mathbf{Q}, \mathbf{e}) = \operatorname*{argmax}_{\mathbf{q}} p_{\mathcal{M}}(\mathbf{q}|\mathbf{e}) = \operatorname*{argmax}_{\mathbf{q}} \{\log p_{\mathcal{M}}(\mathbf{q}, \mathbf{e})\} \tag{1}$$

It is known that the MPE task is NP-hard in general and even hard to approximate [9, 11, 36, 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. Additionally, we develop an output encoding scheme that specifies the number of NN output nodes required for the given PM and enables the recovery of 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 $\mathrm{MPE}(\mathbf{Q}, \mathbf{e})$. 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 $\mathrm{MPE}(\mathbf{Q}, \mathbf{e})$, 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

---

[2]or 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 $\mathrm{MPE}(\mathbf{Q}, \mathbf{e})$, let $\mathbf{q}^c \in [0, 1]^{|\mathbf{Q}|}$ 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 $\log \mathrm{p}_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$ is maximized, namely, $-\log \mathrm{p}_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$ is minimized. Thus, a natural loss function that we can use is $-\log \mathrm{p}_{\mathcal{M}}(\mathbf{q}, \mathbf{e})$. Unfortunately, the NN outputs a continuous vector $\mathbf{q}^c$ and as a result $\mathrm{p}_{\mathcal{M}}(\mathbf{q}^c, \mathbf{e})$ 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 $\ell(\mathbf{q}, \mathbf{e}) = -\log \mathrm{p}_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$, which is a function from $\{0, 1\}^n \to \mathbb{R}$ 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 $\ell^c(\mathbf{q}^c, \mathbf{e}) : [0, 1]^n \to \mathbb{R}$ such that $\ell^c$ coincides with $\ell$ on $\{0, 1\}^n$. For PGMs and PCs, $\ell$ is a multi-linear function and $\ell^c$ 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, $\ell$ 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 $\ell^c$ is that it can be evaluated and differentiated in polynomial time. Moreover, when $\ell$ is defined by either a neural network (in NAMs) or a multilinear function (in BNs, MNs and PCs), the minimum value of $\ell^c$ over the domain $[0, 1]^n$ is less than or equal to the minimum value of the original function $\ell$ over the discrete domain $\{0, 1\}^n$. Formally,

**Proposition 1.** *Let $l(\mathbf{q}, \mathbf{e}) : \{0, 1\}^n \to \mathbb{R}$ be either a neural network or a multilinear function, and let $l^c(\mathbf{q}^c, \mathbf{e}) : [0, 1]^n \to \mathbb{R}$ 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 ($\ell_E$), governed by $\alpha > 0$.

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

This penalty encourages discrete solutions by preferring $q_j^c$ values close to 0 or 1, where $\alpha$ modulates the trade-off. Setting $\alpha$ to 0 yields the continuous approximation; conversely, an $\alpha$ value of $\infty$ 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 $\alpha$ 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 $\{1, \ldots, n\}$ 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.
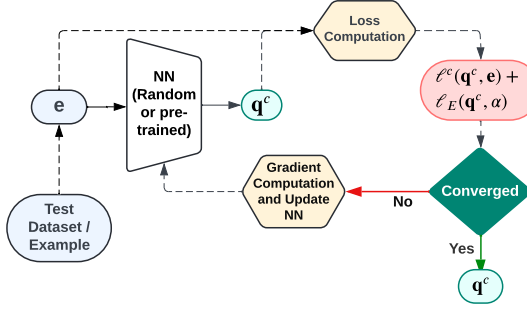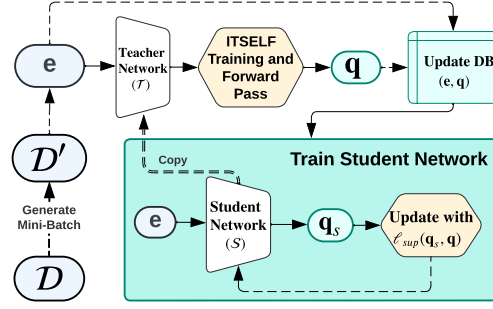
Figure 1: ITSELF Training Procedure



Figure 2: One Training Epoch for $\mathcal{GUIDE}$

### 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. 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), as detailed in Figure 1. 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, 38, 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 $\alpha$ 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 outputs[3]. 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 [35], all of which have the same loss function value. Thus, optimizing and regularizing using the self-supervised loss is difficult compared to a 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.

---

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

**Algorithm 1** <u>GU</u>ided <u>I</u>terative <u>D</u>ual <u>LE</u>arning with Self-supervised Teacher ($\mathcal{GUIDE}$)

---

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:             ▷ 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:    $\mathcal{T} \leftarrow \mathcal{S}$                          ▷ Initialize $\mathcal{T}$ with $\mathcal{S}$ for the next epoch
14: **end for**

---

## 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, 37, 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 ITSELF algorithm, 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.

## 4.2 Training Procedure

Our proposed training procedure, which we call $\mathcal{GUIDE}$, 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 ITSELF algorithm (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 ($\ell_{sup}$) 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). Figure 2 illustrates a single training epoch of GUIDE.

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 $\mathcal{GUIDE}$ 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 PCs and NAMs. For the purpose of evaluating PGMs, 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 NAMs, we trained Masked Autoencoder for Distribution Estimation (MADE) models using PyTorch, following the approach in Germain et al. [17]. For Markov Networks (MNs), a specific type of PGM, we applied Gibbs sampling to generate 8,000, 1,000, and 1,000 samples for the training, testing, and validation sets, respectively. The query ratio ($qr$), defined as the fraction of variables in the query set, was varied across the set $\{0.1, 0.3, 0.5, 0.7, 0.8, 0.9\}$ for each probabilistic model (PM).

## 5.2 Baseline Methods and Evaluation Criteria

**PC**s - We used three polynomial-time baseline methods from the probabilistic circuits and probabilistic graphical models literature as 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 $j - 1$.

We further evaluated 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].

**NAM**s - As a baseline, we used the stochastic hill-climbing search (HC) algorithm. Following a procedure similar to that used for PCs, we conducted a 60-second hill-climbing search for each test example, with query variables initialized randomly and setting evidence variables according to the values in the given example.

**PGM**s - We employed the distributed AND/OR Branch and Bound (AOBB) method [39] as a baseline, using the implementation outlined in Otten [40]. Since AOBB is an anytime algorithm, we set a 60-second time limit for inference per test example.

**Neural Baselines** - Arya et al. [4] introduced Self-Supervised learning based MMAP solver for PCs (SSMP), training a neural network to handle queries on a fixed variable partition within PCs. We extend this approach to address the any-MPE task in PMs (see Section 3.2), using a single network to answer any-MPE queries as an additional neural baseline.

**Evaluation Criteria** - We evaluated competing approaches based on log-likelihood (LL) scores, calculated as $\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$, and inference times for given evidence $\mathbf{e}$ and query output $\mathbf{q}$. Higher log-likelihood scores indicate better performance, while shorter inference times are preferable.

## 5.3 Neural Network-Based Approaches

We implemented two neural network training protocols for each PM and query ratio: SSMP and $\mathcal{GUIDE}$. Each model was trained for 20 epochs following the training procedure outlined by Arya et al. [4] for SSMP. Both protocols employed two distinct inference strategies, thus forming four neural-based variants. In the first strategy, we performed a single forward pass through the network to estimate the values of query variable, as specified by Arya et al. [4]. The second strategy utilized our novel test-time optimization-based ITSELF approach for inference. The ITSELF optimization terminates after 100 iterations or upon loss convergence for both PCs and PGMs. For NAMs, we increase the limit to 1,000 iterations while keeping the convergence criterion.

We standardized network architectures for PMs across all experiments. For PCs, we used 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]. We optimized all models using Adam [25] and implemented them 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 $\mathcal{GUIDE}$ methods—on 20 TPM datasets as shown in the contingency table in figure 3a (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. We present similar contingency tables for Hill Climbing Search over PCs (Fig. 3b), NAMs (Fig. 3c), and PGMs (Fig. 3d) to benchmark the proposed methods against the baselines.

The contingency table for PC (Fig. 3a) shows that methods incorporating ITSELF consistently outperform both polynomial-time and traditional neural baselines, as indicated by the dark blue cells in the corresponding rows. Notably, $\mathcal{GUIDE}$ + 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 $\mathcal{GUIDE}$. In contrast, the polynomial-time baseline MAX is better than both SSMP and $\mathcal{GUIDE}$ (as 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 $\mathcal{GUIDE}$ + ITSELF using a heatmap in Figure 4a. The y-axis presents datasets by variable count and the x-axis represents query ratio. Each cell displays the percentage difference in mean LL scores between the methods, calculated as %Diff. $= 100 \times (ll_{nn} - ll_{max})/|ll_{max}|$. The heatmap shows that $\mathcal{GUIDE}$ + ITSELF achieves performance comparable to MAX for small query sets. As the problem complexity increases with an increase in query set size, our method consistently outperforms MAX across all datasets, except for NLTCS and Tretail, as highlighted by the green cells. In the 12 cases where $\mathcal{GUIDE}$ + ITSELF underperforms, the performance gap remains minimal, as indicated by the limited number of red cells in the heatmap.

| | MAX | ML | Seq | SSMP | $\mathcal{GUIDE}$ | SSMP ITSELF | $\mathcal{GUIDE}$ ITSELF |
|---|---|---|---|---|---|---|---|
| MAX | 0 | 79 | 94 | 81 | 71 | 39 | 12 |
| ML | 41 | 0 | 43 | 55 | 49 | 15 | 5 |
| Seq | 26 | 64 | 0 | 44 | 43 | 16 | 7 |
| SSMP | 39 | 65 | 76 | 0 | 24 | 4 | 0 |
| $\mathcal{GUIDE}$ | 49 | 71 | 77 | 71 | 0 | 8 | 0 |
| SSMP ITSELF | 81 | 105 | 104 | 104 | 99 | 0 | 7 |
| $\mathcal{GUIDE}$ ITSELF | 108 | 115 | 113 | 107 | 107 | 94 | 0 |

(a) PCs: Initial Solutions

| | MAX | ML | Seq | SSMP | $\mathcal{GUIDE}$ | SSMP ITSELF | $\mathcal{GUIDE}$ ITSELF |
|---|---|---|---|---|---|---|---|
| MAX | 0 | 65 | 81 | 66 | 62 | 46 | 21 |
| ML | 35 | 0 | 49 | 41 | 27 | 22 | 4 |
| Seq | 21 | 45 | 0 | 29 | 27 | 15 | 4 |
| SSMP | 36 | 55 | 71 | 0 | 23 | 17 | 1 |
| $\mathcal{GUIDE}$ | 39 | 71 | 73 | 49 | 0 | 24 | 0 |
| SSMP ITSELF | 53 | 73 | 82 | 70 | 61 | 0 | 12 |
| $\mathcal{GUIDE}$ ITSELF | 79 | 88 | 94 | 91 | 89 | 72 | 0 |

(b) PCs: Hill-Climbing

| | HC | SSMP | $\mathcal{GUIDE}$ | SSMP ITSELF | $\mathcal{GUIDE}$ ITSELF |
|---|---|---|---|---|---|
| HC | 0 | 46 | 38 | 36 | 12 |
| SSMP | 34 | 0 | 7 | 13 | 4 |
| $\mathcal{GUIDE}$ | 42 | 66 | 0 | 34 | 16 |
| SSMP ITSELF | 44 | 61 | 40 | 0 | 6 |
| $\mathcal{GUIDE}$ ITSELF | 68 | 72 | 60 | 70 | 0 |

(c) NAMs

| | AOBB | SSMP | $\mathcal{GUIDE}$ | SSMP ITSELF | $\mathcal{GUIDE}$ ITSELF |
|---|---|---|---|---|---|
| AOBB | 0 | 7 | 6 | 6 | 5 |
| SSMP | 9 | 0 | 3 | 4 | 0 |
| $\mathcal{GUIDE}$ | 9 | 12 | 0 | 15 | 3 |
| SSMP ITSELF | 10 | 12 | 0 | 0 | 3 |
| $\mathcal{GUIDE}$ ITSELF | 11 | 16 | 12 | 13 | 0 |

(d) PGMs

Figure 3: MPE method comparison across PMs. Blue shows row superiority, red shows column superiority; darker shades indicate larger values.

8

Figure 3b further analyzes the performance of our proposed methods against various baselines as initialization strategies for Hill Climbing Search. This comparison evaluates the effectiveness of ITSELF and $\mathcal{GUIDE}$ in enhancing *anytime methods* compared to conventional heuristic initialization approaches. Notably, methods incorporating ITSELF provide superior initialization for local search-based algorithms.

**Results for NAMs**: The contingency table in Figure 3c presents our evaluation of several methods for NAMs, including HC and two neural network approaches, SSMP and $\mathcal{GUIDE}$, each tested with two inference schemes. We evaluated these methods on 20 TPM datasets, creating 80 test sets for the MPE task using 20 MADEs across four query ratios ($qr$).

The $\mathcal{GUIDE}$ + ITSELF approach demonstrates superior performance compared to both baseline methods and other neural inference schemes, aligning with observations from PC. While HC outperforms SSMP, both $\mathcal{GUIDE}$ and the combination of SSMP-based training with ITSELF-based inference surpass HC, highlighting their advantages over the baseline.

The heatmaps in Figure 4b further highlight the superior performance of $\mathcal{GUIDE}$ + 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 combination of $\mathcal{GUIDE}$-based learning with ITSELF-based inference consistently outperforms the baseline across most datasets, with exceptions only in the Mushrooms, Connect 4, and Retail. Overall, the $\mathcal{GUIDE}$ + ITSELF approach significantly enhances the quality of the MPE solutions in NAM models.

**Results for PGMs**: The contingency table in 3d compares the performance of AOBB and four neural-network-based methods on PGMs across four high-treewidth networks. For this evaluation, we generated 16 test datasets for the MPE task using four PGMs across four query ratios ($qr$).

Consistent with results from previous PMs, methods using IT-SELF for inference consistently outperform the baseline methods AOBB and SSMP across most scenarios. Both $\mathcal{GUIDE}$ and SSMP outperform AOBB in at least 50 percent of the tests. The supplementary material presents comparisons against exact solutions, conducted on less complex probabilistic models where ground truth computation remains tractable.

**Does a teacher-student-based network outperform a single network trained with the self-supervised loss? ($\mathcal{GUIDE}$ vs. SSMP):**

This analysis aims to evaluate the performance of $\mathcal{GUIDE}$ against traditional neural network training methods used in SSMP across different PMs and inference schemes. Using traditional inference scheme (i.e., one forward pass through the network), $\mathcal{GUIDE}$ 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-based inference, $\mathcal{GUIDE}$ maintains this advantage, achieving higher quality solutions in more than 75%, 85%, and 80% of cases for PCs, NAMs, and PGMs, respectively. Therefore, models trained using $\mathcal{GUIDE}$ are consistently superior to those trained with SSMP for the *any-MPE* task.



(a) PC: $\mathcal{GUIDE}$ + ITSELF vs. MAX



(b) NAM: $\mathcal{GUIDE}$ + ITSELF vs. HC

Figure 4: Heatmaps showing LL % Differences. Top: PC; Bottom: NAM. Green cells: our method is better. Darker shades indicate larger values.

**Does inference time optimization improve performance? (One-Pass vs. Multi-Pass):**

In this analysis, we compare the performance of the single-pass inference method to that of the proposed multi-pass inference method (ITSELF). ITSELF combined with SSMP training outperforms the other methods in over 85% cases for PC, and more than 75% for NAM and PGM models. When used on models trained with $\mathcal{GUIDE}$, ITSELF demonstrates even better results, achieving superior performance in nearly 90% of PC cases and 75% for both NAMs and PGMs. Overall, $\mathcal{GUIDE}$ with ITSELF inference emerges as the most effective method across all experiments. Empirical evidence consistently demonstrates ITSELF's superiority over single-pass inference across PMs.

The inference time analysis, detailed in the supplementary material, compares computational efficiency across methods using the natural logarithm of execution time in microseconds. Neural network-based approaches with traditional inference demonstrate the fastest performance across all PMs, as they only require a single forward pass to compute query variable values. For MADE, models trained with $\mathcal{GUIDE}$ and ITSELF are the next most efficient. In PGMs, $\mathcal{GUIDE}$ + ITSELF ranks third, followed by SSMP + ITSELF. For PCs, MAX is marginally faster than both $\mathcal{GUIDE}$ + ITSELF and SSMP + ITSELF, while ML and Seq have the longest computational times. In general, models trained with $\mathcal{GUIDE}$ achieve shorter inference times than those trained with the self-supervised loss (SSMP), as they require fewer ITSELF iterations due to more effective initial training.

**Summary:** Our experiments demonstrate that $\mathcal{GUIDE}$ + ITSELF outperforms both polynomial-time and neural-based baselines across various PMs, as evidenced by higher log-likelihood scores. Notably, ITSELF demonstrates significant advantages over traditional single-pass inference in addressing the complex *any-MPE* query task within probabilistic models, emphasizing the importance of Inference Time Optimization. Furthermore, the superior performance of models trained with $\mathcal{GUIDE}$ compared to SSMP highlights the effectiveness of the dual network approach, which improves initial model quality and establishes 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 $\mathcal{GUIDE}$ 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, neural autoregressive models, and probabilistic graphical models, 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 with constraints; training neural networks with losses from multiple probabilistic models to embed their inference mechanisms; boosting performance by developing advanced encoding strategies for similar tasks; implementing sophisticated neural architectures tailored to probabilistic models; etc.

## Acknowledgements

## 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. Technical report, University of California, Los Angeles, 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 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.

[36] Denis Deratani Mauá and Cassio P. de Campos. Approximation complexity of maximum A posteriori inference in sum-product networks. *CoRR*, abs/1703.06045, 2017.

[37] 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.

[38] 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.

[39] L. Otten and R. Dechter. A case study in complexity estimation: Towards parallel branch-and-bound over graphical models. *Conference on Uncertainty in Artificial Intelligence*, 2012.

[40] Lars Otten. DAOOPT: Sequential and distributed AND/OR branch and bound for MPE problems. 2012. URL https://github.com/lotten/daoopt.

[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 Uria, 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.

# A  Experimental Setup

## A.1  Datasets and Models

Table 1 summarizes the datasets and the probabilistic circuits trained on them. We use the same datasets for both PCs [8] and NAMs [17, 50]. The selection includes both smaller datasets, such as NLTCS and MSNBC, and larger datasets with over 1000 nodes.

For Markov networks, we utilize high treewidth grid networks, specifically grid40x40.f2.wrap, grid40x40.f5.wrap, grid40x40.f10.wrap, and grid40x40.f15.wrap. Each model contains 4800 variables and 1600 factors.

Table 1: Summary of datasets used with their respective numbers of variables and nodes in probabilistic circuits.

| Dataset | Number of Variables | Number of Nodes in PC |
|---|---|---|
| NLTCS | 16 | 125 |
| MSNBC | 17 | 46 |
| KDDCup2k | 64 | 274 |
| Plants | 69 | 3737 |
| Audio | 100 | 348 |
| Jester | 100 | 274 |
| Netflix | 100 | 400 |
| Accidents | 111 | 1178 |
| Mushrooms | 112 | 902 |
| Connect 4 | 126 | 2128 |
| Retail | 135 | 359 |
| RCV-1 | 150 | 519 |
| DNA | 180 | 1855 |
| Book | 500 | 1628 |
| WebKB | 839 | 3154 |
| Reuters-52 | 889 | 7348 |
| 20 NewsGroup | 910 | 2467 |
| Movie reviews | 1001 | 2567 |
| BBC | 1058 | 3399 |
| Ad | 1556 | 9666 |

## A.2  Hyperparameters Details

Our experimental framework was designed to ensure consistency and efficiency across all conducted experiments. For NAM's, we used MADE, training the model with two hidden layers of 512 and 1024 units, respectively, using the hyperparameters from Germain et al. [17].

For neural network-based solvers, the mini-batch size was set to 512 samples, and a learning rate decay strategy, reducing the rate by 0.9 upon loss plateauing, was implemented to improve training efficiency. Optimal hyperparameters were identified via extensive 5-fold cross-validation.

In discrete loss scenarios, the hyperparameter $\alpha$ played a pivotal role. We systematically explored the optimal $\alpha$ value across the range $0.001, 0.01, 0.1, 1, 10, 100, 1000$ for neural-based models, including ITSELF and $\mathcal{GUIDE}$. Notably, higher $\alpha$ values better constrain outputs to binary, thereby facilitating near-optimal results.

# B  Extending the Current Approach to Other Data Types and Inference Tasks

The current approach can be extended to support both multi-valued discrete and continuous variables, broadening its utility in diverse scenarios.

For multi-valued discrete variables, the method can be adapted by implementing a multi-class, multi-output classification head. Each query variable is represented by a softmax output node, which provides soft evidence by generating probabilistic distributions across multiple discrete values.

To incorporate continuous variables, we introduce a linear activation function in the output layer. The loss function, specifically the multi-linear representation of the PM, is modified to accommodate continuous neural network outputs. For example, in Probabilistic Circuits that use Gaussian distributions, continuous values can be directly integrated into the loss function, facilitating gradient-based backpropagation.

These extensions primarily involve adjusting the network's output layer and refining the self-supervised loss function represented by the PM. Notably, other elements of our approach, including the ITSELF and GUIDE procedures, remain unchanged.

Our approach further extends to additional inference tasks over probabilistic models, including marginal MAP and constrained most probable explanation (CMPE) tasks. However, the scalability of this approach depends on the computational efficiency of evaluating the loss function for each inference task. When this evaluation becomes computationally infeasible, the proposed method—training a neural network to answer queries over probabilistic models—may itself become infeasible. For example, performing marginal MAP inference over NAMs and PGMs requires repeated evaluations of the loss function associated with the marginal MAP task and its gradient during training. This iterative process, essential for updating the neural network's parameters, can become prohibitively resource-intensive due to the high computational demands of evaluating the marginal MAP loss over these probabilistic models.

## C  A Comparative Analysis of Performance of ITSELF for Different Pre-Training Methods

This section evaluates the performance of models initialized through various techniques—random initialization, SSMP, and $\mathcal{GUIDE}$. Each plot represents the loss for a distinct test example, with the x-axis denoting the number of ITSELF iterations and the y-axis showing the Negative Log Likelihood (NLL) scores. Lower NLL values signify better solutions. Through this empirical assessment, we compare the impact of different pre-training methods on model performance.

Figures 5 to 28 present the plots for NAMs. The plots for PCs are shown in Figures 29 to 67. Figures 68 to 78 illustrate the plots for PGMs. We selected the following datasets for PCs and NAMs: DNA, RCV-1, Reuters-52, Netflix, WebKB, Audio, Moviereview, and Jester. For PGMs, we used all the datasets presented in the main paper. Each plot consists of two sections. The left section presents the Negative Log-Likelihood Loss for 1000 iterations for all methods. The right section contains two sub-plots: the top sub-plot displays the zoomed-in losses for the first 200 iterations, while the bottom sub-plot shows the zoomed-in losses for the last 200 iterations.

We randomly initialize the parameters for the random model and perform 1000 iterations of ITSELF for inference. For the two pre-trained models (SSMP and $\mathcal{GUIDE}$), we update the top $N$ layers, where $N$ is the number of layers corresponding to that loss curve, and fix the remaining bottom layers. We extract features by passing the input through these fixed layers and then train the parameters of the top $N$ layers. We again perform 1000 iterations of ITSELF for inference. For NAMs and PGMs, we use neural networks with up to one hidden layer, while for PCs, we employ models with up to three hidden layers.

From the plots for the three Probabilistic Modelss (PMs), we observe that models pre-trained using the proposed $\mathcal{GUIDE}$ training scheme generally have a better starting point for ITSELF, indicated by a lower loss, compared to all other models. Across a wide array of datasets, PGMs, and query percentages, the $\mathcal{GUIDE}$ method consistently converges to a lower or equivalent loss compared to other models. Remarkably, it sometimes achieves a loss value that is less than half of the nearest competing model. Furthermore, the losses for $\mathcal{GUIDE}$ are typically more stable than those of other initialization. In some scenarios, all models achieve a similar final loss, although models initialized with SSMP and those randomly initialized may experience oscillations in their loss values.

Models pre-trained using the traditional self-supervised loss (SSMP) typically have better or similar starting points than randomly initialized models. However, models pre-trained using the SSMP method might converge to a worse loss than their $\mathcal{GUIDE}$ pre-trained counterparts.

In most cases, convergence is rapid, even with a reduced learning rate of $10^{-4}$ compared to the experiments shown in the main paper. Most methods converge within 200 to 300 iterations, although some may still oscillate during the later iterations of ITSELF.

Figure 5: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the DNA Dataset at a Query Ratio of 0.5.



Figure 6: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the DNA Dataset at a Query Ratio of 0.7.



Figure 7: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the DNA Dataset at a Query Ratio of 0.9.
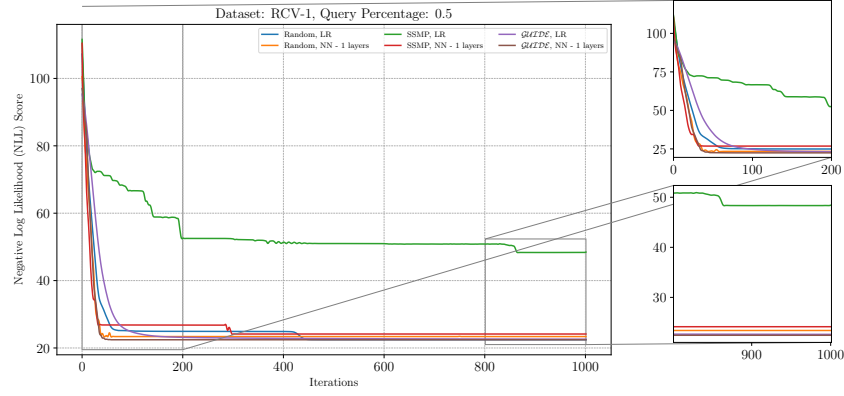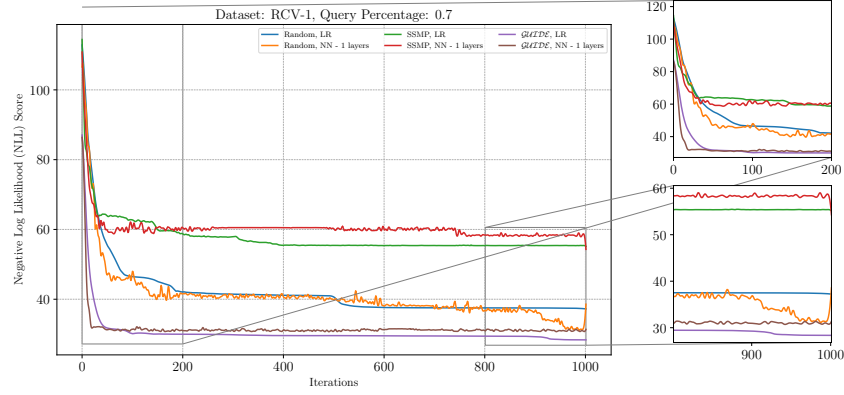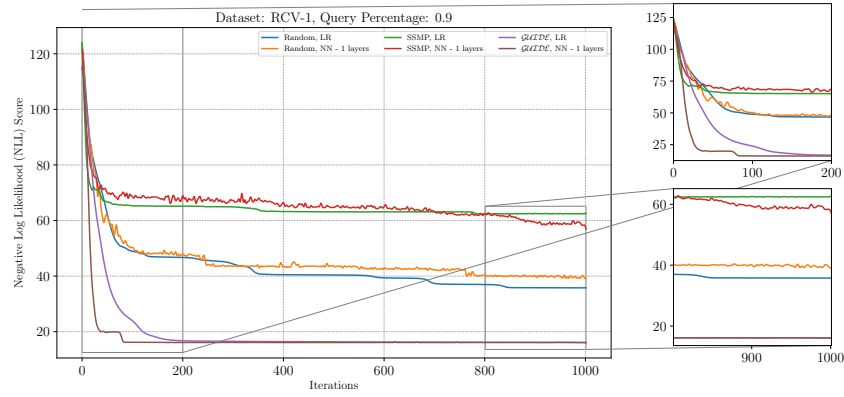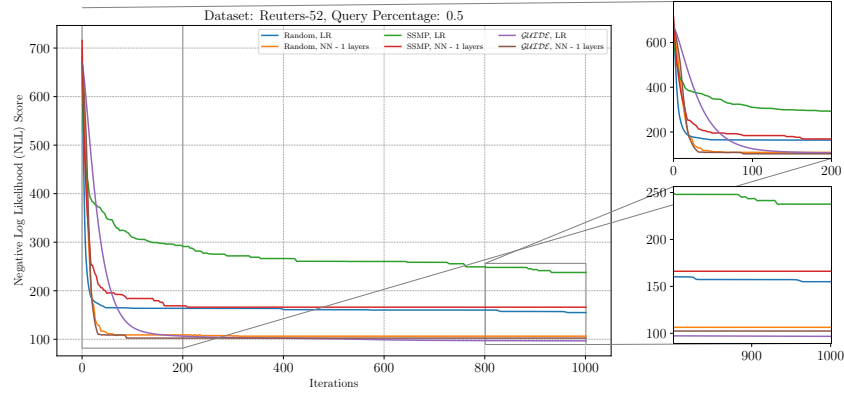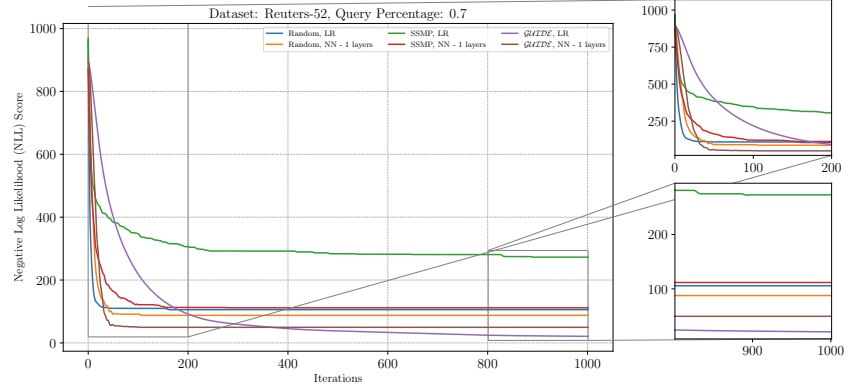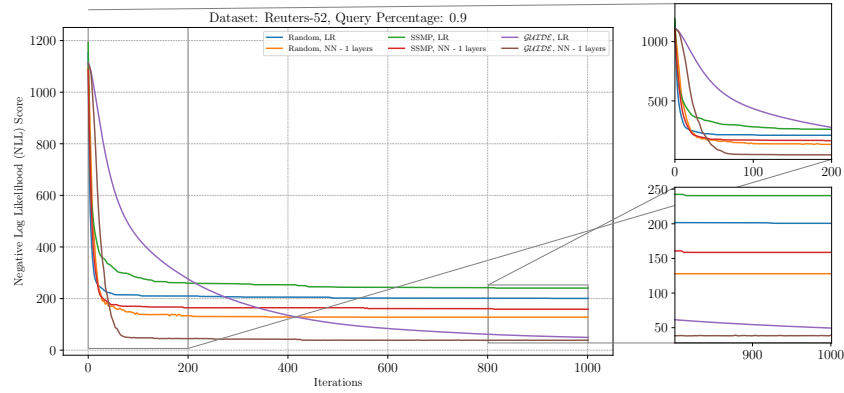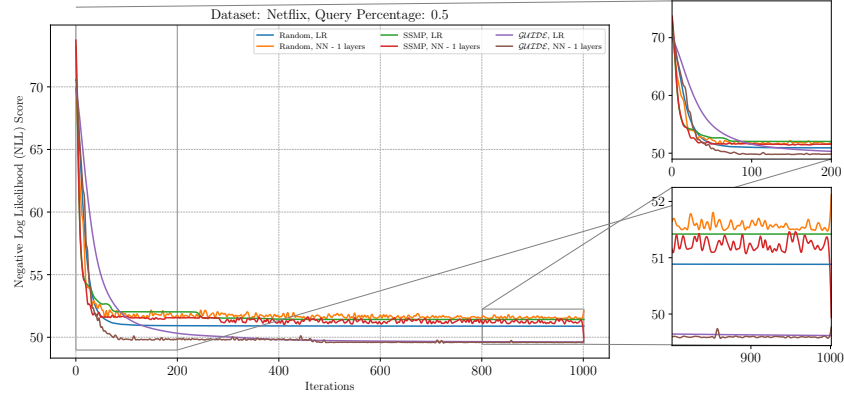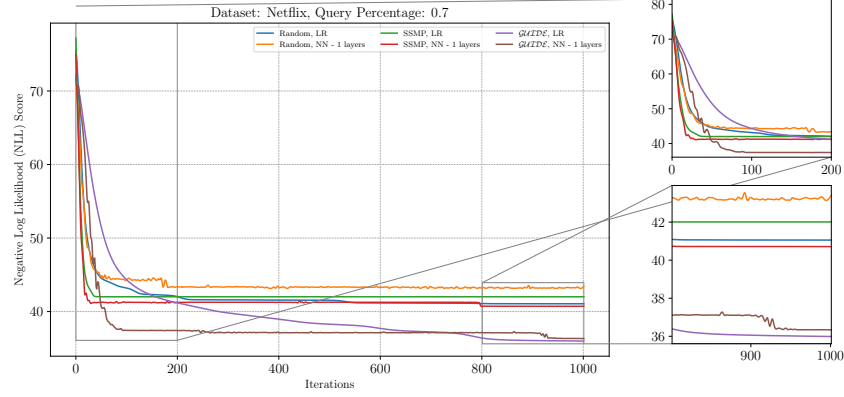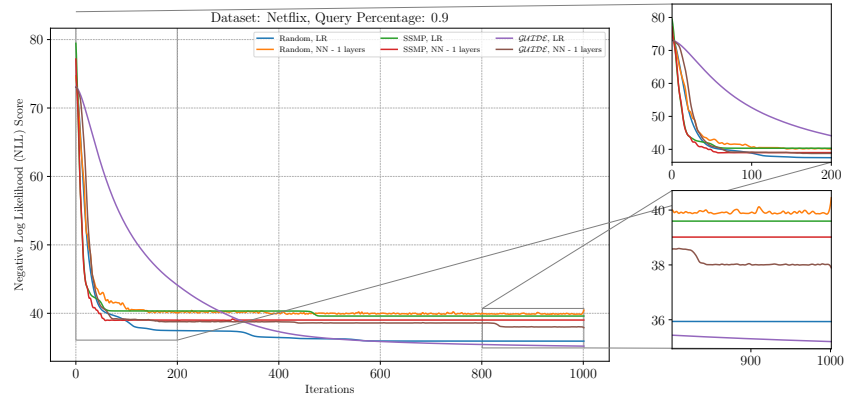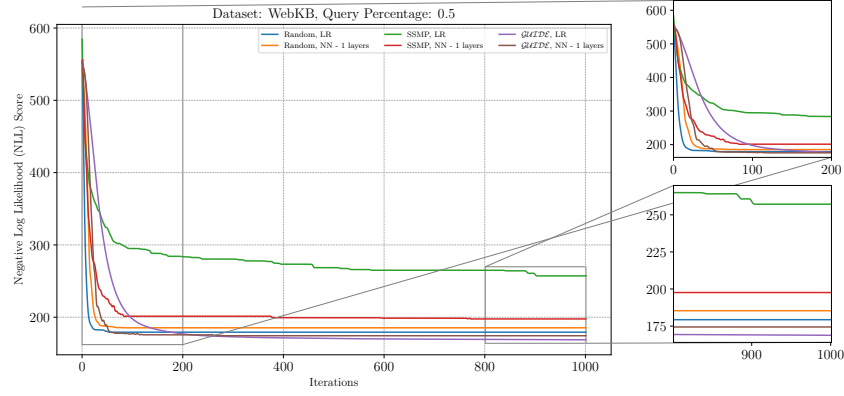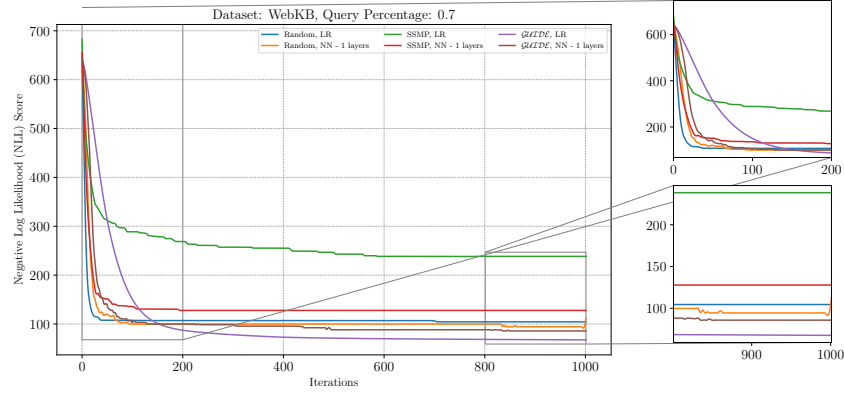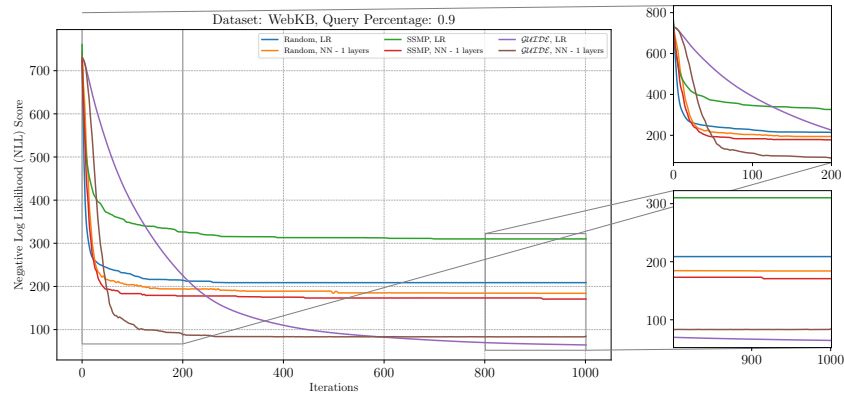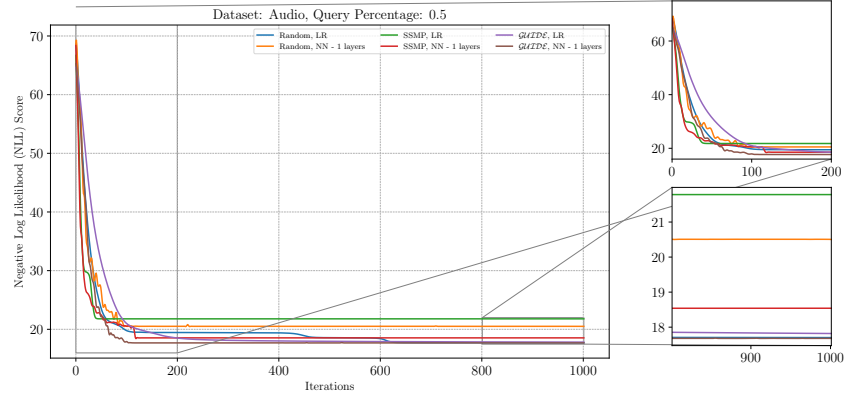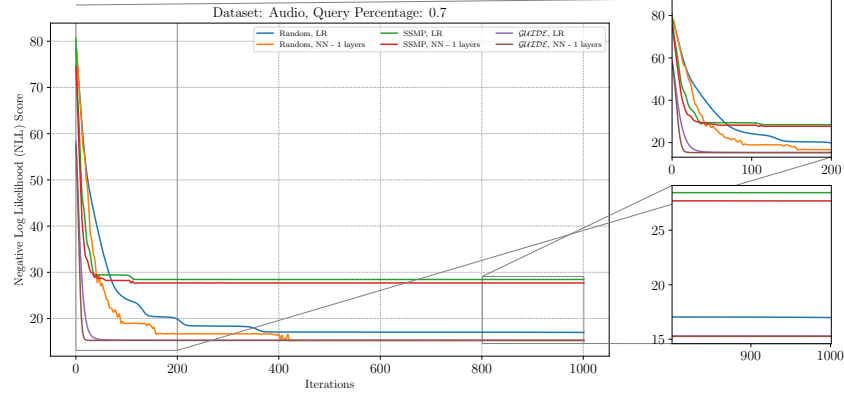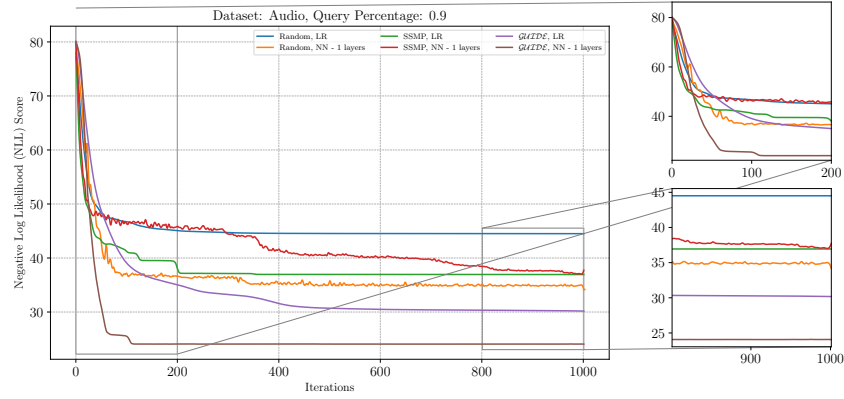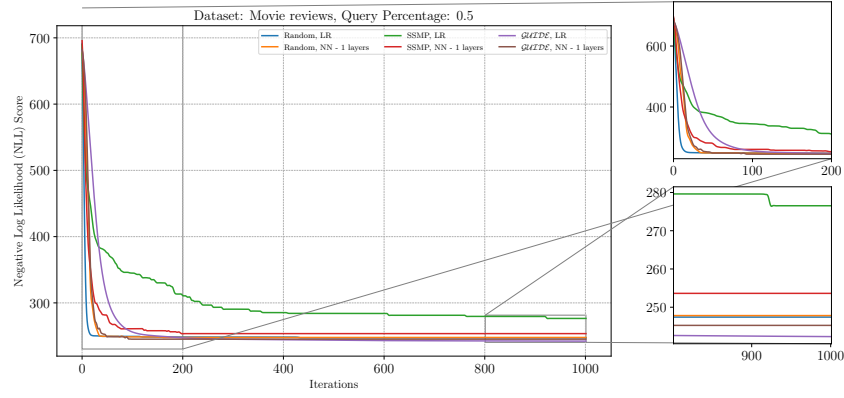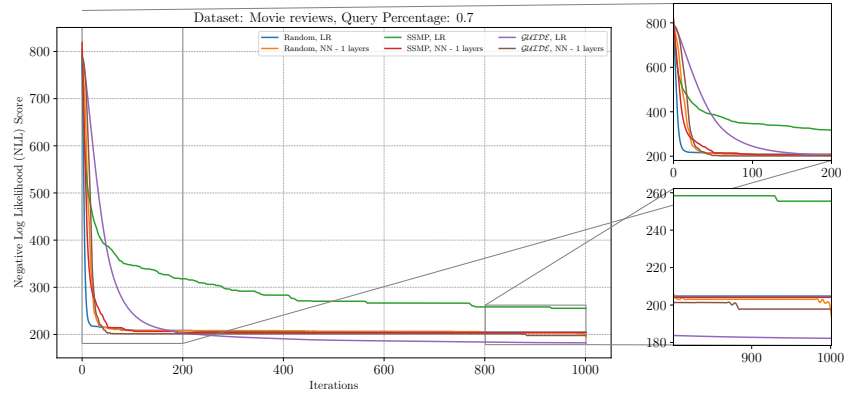
Figure 8: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the RCV-1 Dataset at a Query Ratio of 0.5.



Figure 9: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the RCV-1 Dataset at a Query Ratio of 0.7.
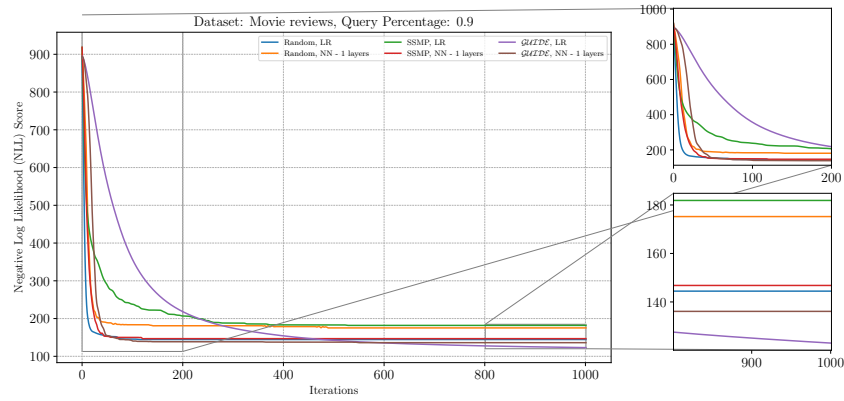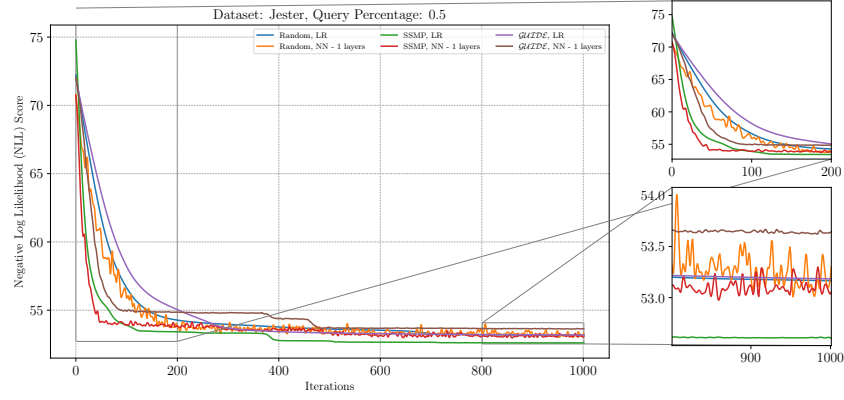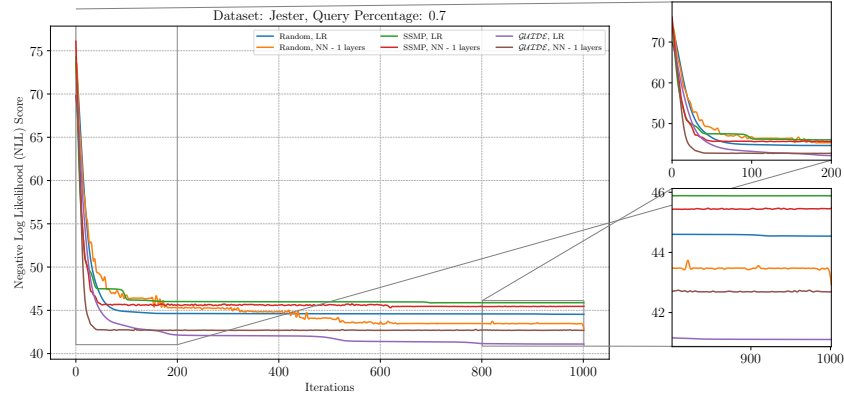


Figure 10: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the RCV-1 Dataset at a Query Ratio of 0.9.

Figure 11: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Reuters-52 Dataset at a Query Ratio of 0.5.



Figure 12: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Reuters-52 Dataset at a Query Ratio of 0.7.
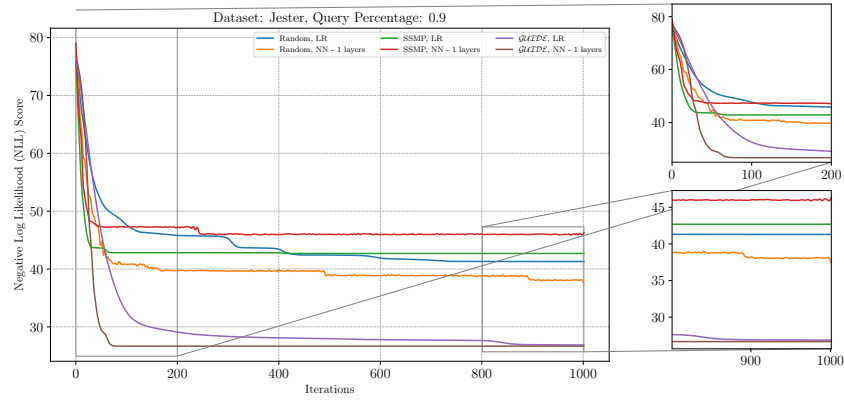


Figure 13: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Reuters-52 Dataset at a Query Ratio of 0.9.

Figure 14: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Netflix Dataset at a Query Ratio of 0.5.
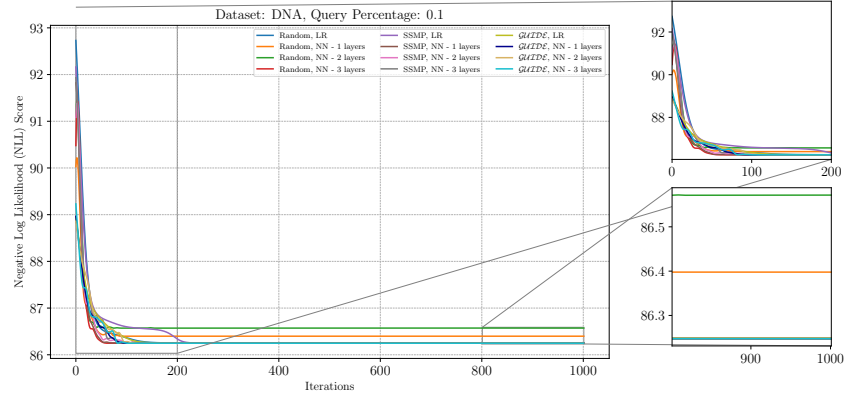


Figure 15: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Netflix Dataset at a Query Ratio of 0.7.



Figure 16: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Netflix Dataset at a Query Ratio of 0.9.

Figure 17: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the WebKB Dataset at a Query Ratio of 0.5.
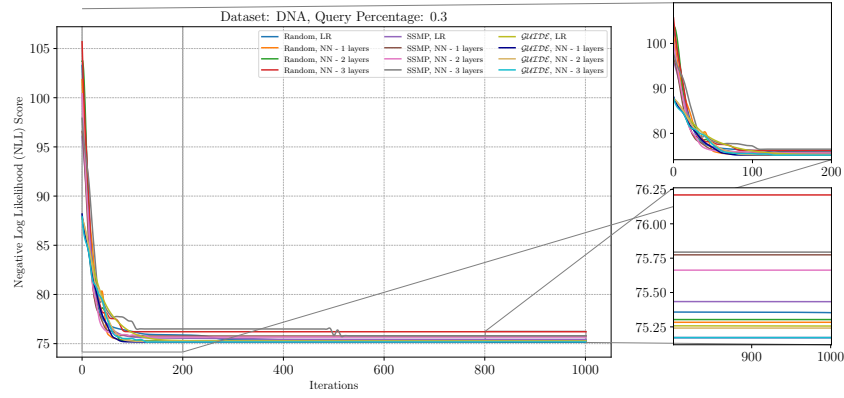


Figure 18: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the WebKB Dataset at a Query Ratio of 0.7.



Figure 19: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the WebKB Dataset at a Query Ratio of 0.9.

Figure 20: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Audio Dataset at a Query Ratio of 0.5.
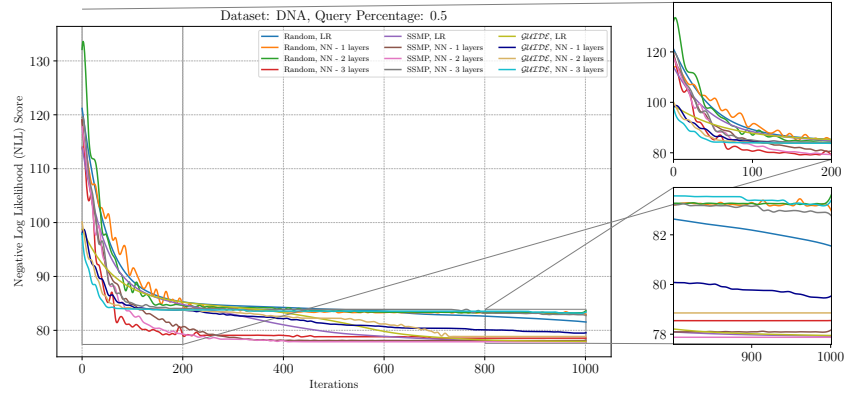


Figure 21: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Audio Dataset at a Query Ratio of 0.7.
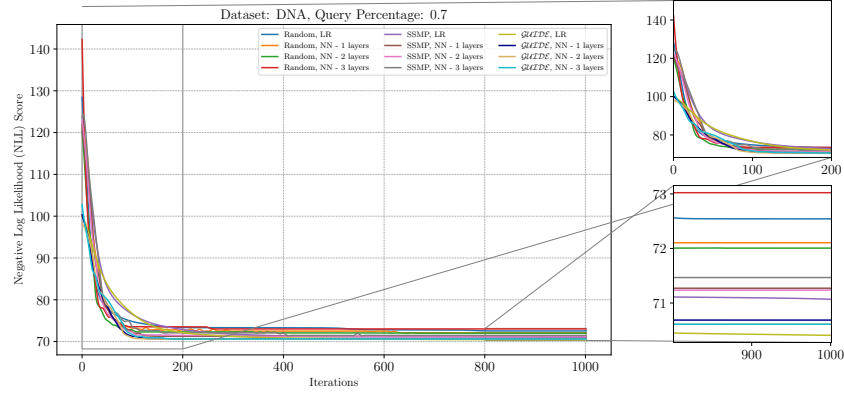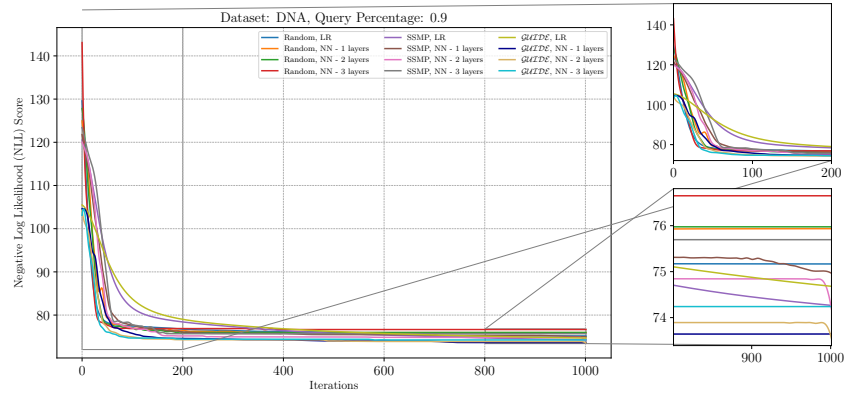


Figure 22: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Audio Dataset at a Query Ratio of 0.9.

Figure 23: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Movie reviews Dataset at a Query Ratio of 0.5.



Figure 24: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Movie reviews Dataset at a Query Ratio of 0.7.



Figure 25: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Movie reviews Dataset at a Query Ratio of 0.9.

Figure 26: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Jester Dataset at a Query Ratio of 0.5.
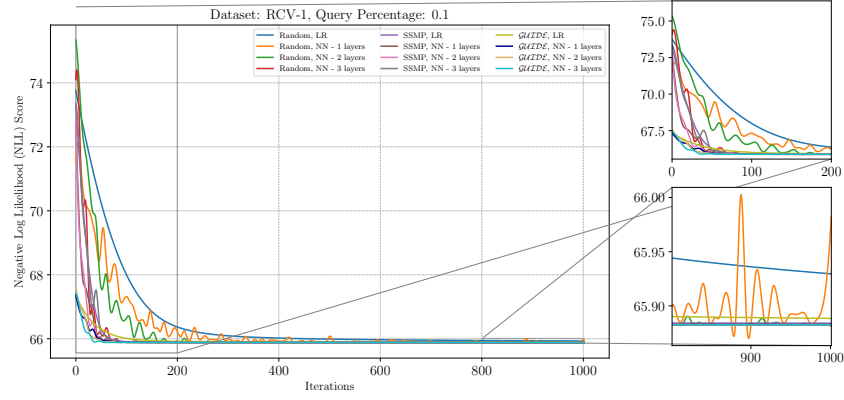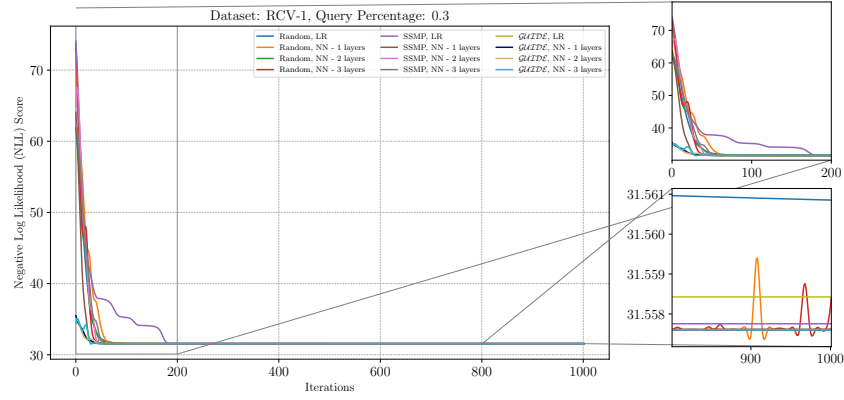


Figure 27: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Jester Dataset at a Query Ratio of 0.7.



Figure 28: Analysis of ITSELF Loss Across Various Pre-Trained Models for NAMs on the Jester Dataset at a Query Ratio of 0.9.

24

Figure 29: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the DNA Dataset at a Query Ratio of 0.1.



Figure 30: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the DNA Dataset at a Query Ratio of 0.3.
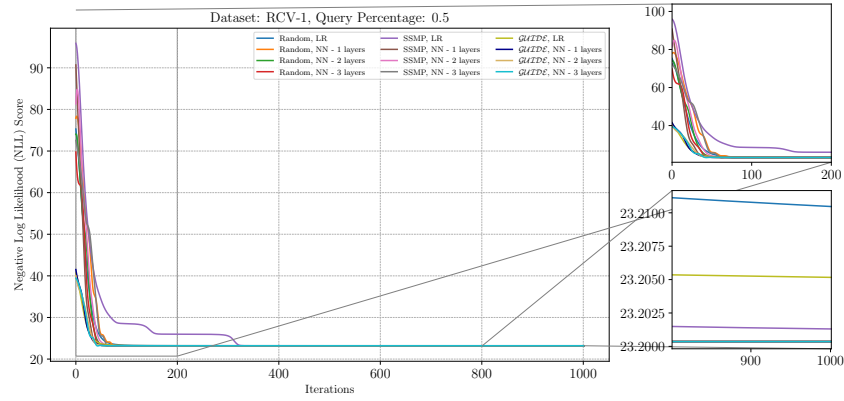
25

Figure 31: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the DNA Dataset at a Query Ratio of 0.5.
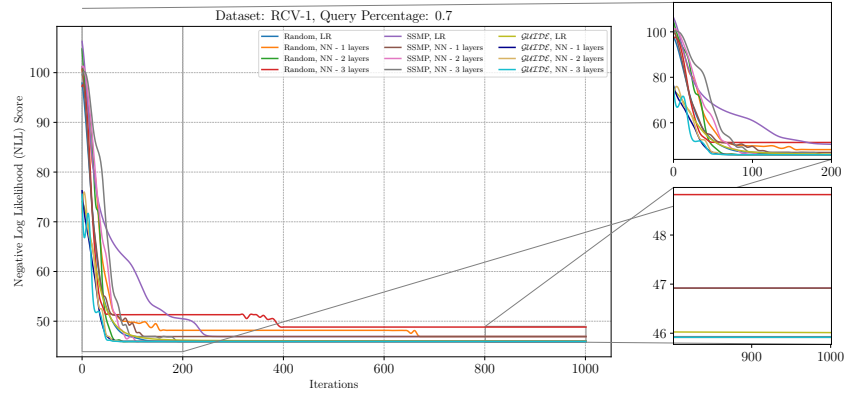


Figure 32: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the DNA Dataset at a Query Ratio of 0.7.



Figure 33: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the DNA Dataset at a Query Ratio of 0.9.
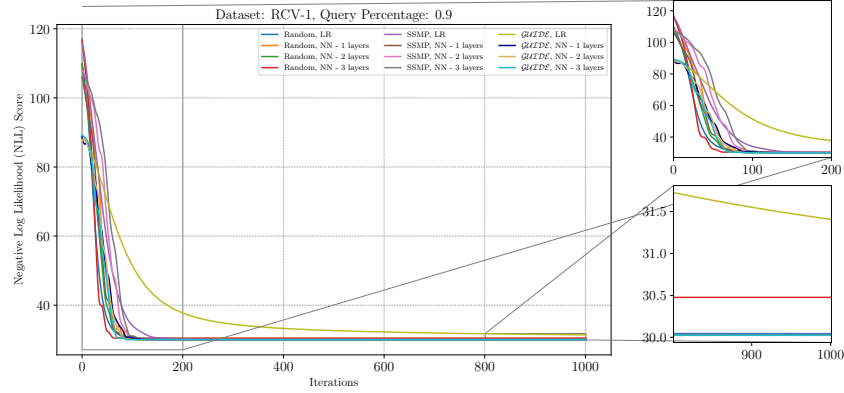
Figure 34: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the RCV-1 Dataset at a Query Ratio of 0.1.



Figure 35: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the RCV-1 Dataset at a Query Ratio of 0.3.



Figure 36: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the RCV-1 Dataset at a Query Ratio of 0.5.

Figure 37: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the RCV-1 Dataset at a Query Ratio of 0.7.



Figure 38: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the RCV-1 Dataset at a Query Ratio of 0.9.
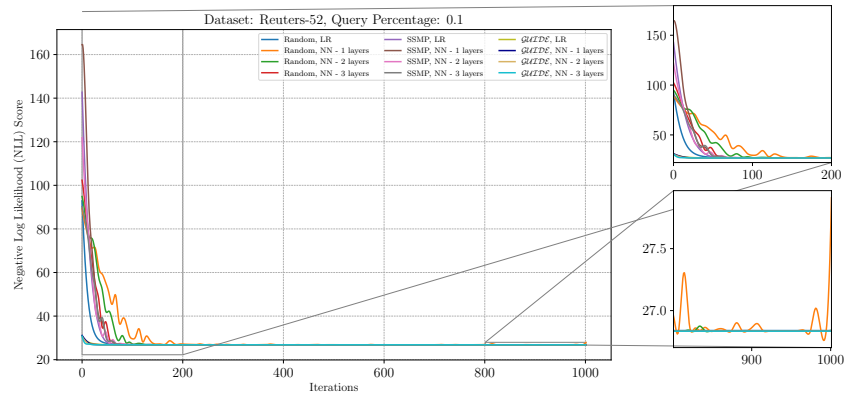


Figure 39: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Reuters-52 Dataset at a Query Ratio of 0.1.
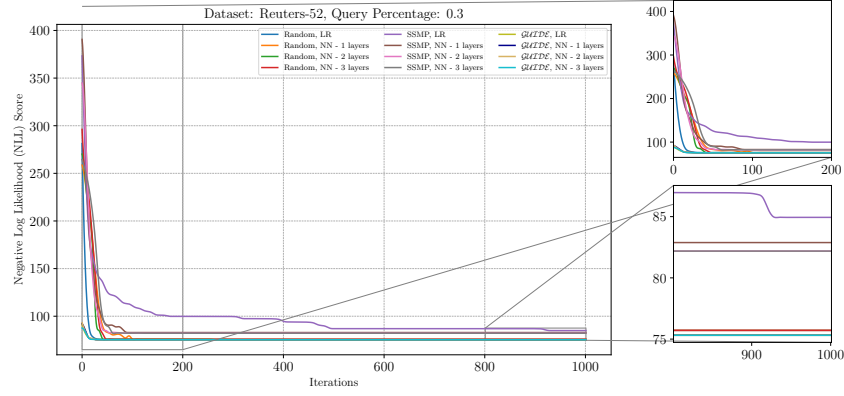
Figure 40: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Reuters-52 Dataset at a Query Ratio of 0.3.
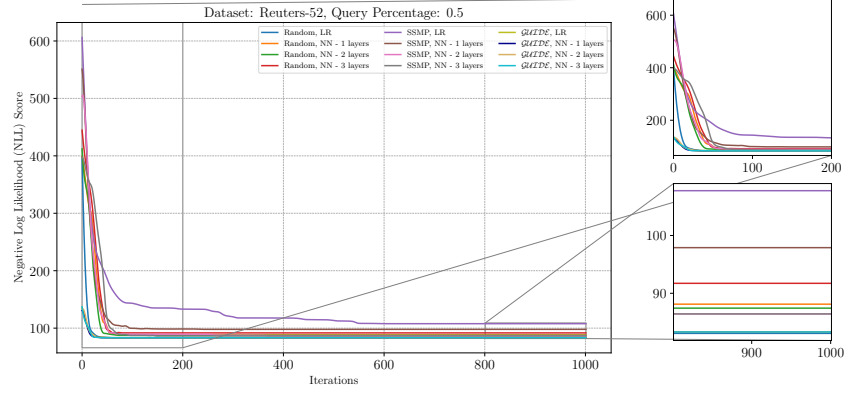


Figure 41: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Reuters-52 Dataset at a Query Ratio of 0.5.
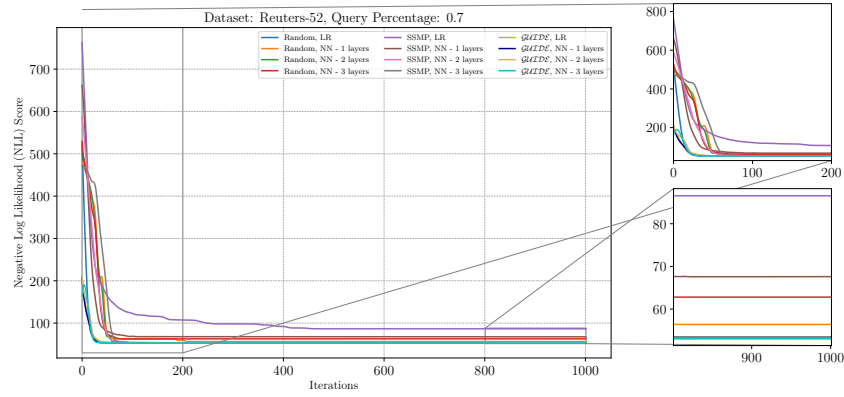


Figure 42: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Reuters-52 Dataset at a Query Ratio of 0.7.
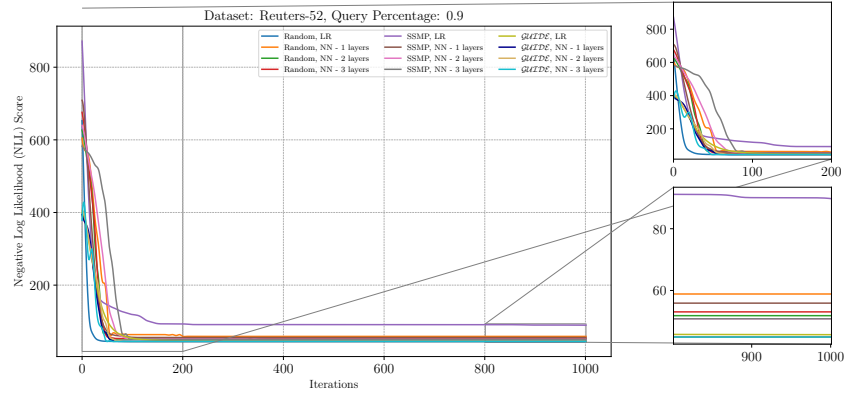
Figure 43: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Reuters-52 Dataset at a Query Ratio of 0.9.



Figure 44: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Netflix Dataset at a Query Ratio of 0.1.



Figure 45: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Netflix Dataset at a Query Ratio of 0.3.
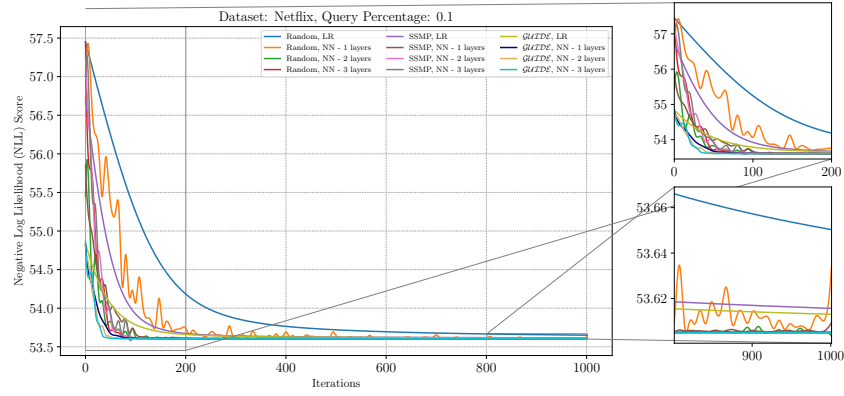
Figure 46: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Netflix Dataset at a Query Ratio of 0.5.
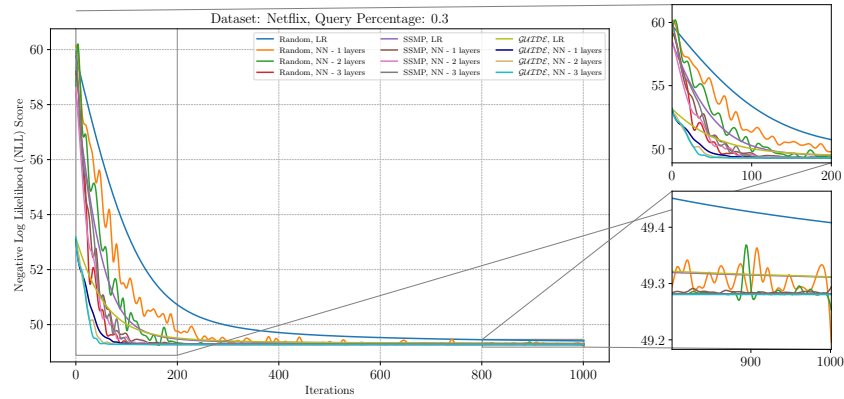


Figure 47: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Netflix Dataset at a Query Ratio of 0.7.
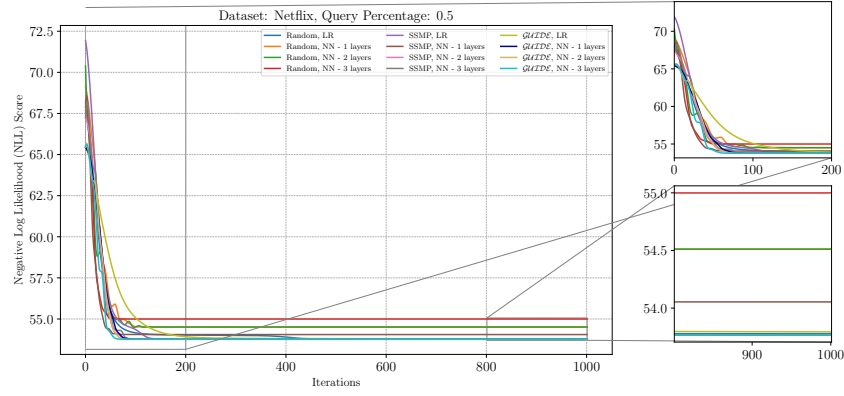


Figure 48: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Netflix Dataset at a Query Ratio of 0.9.

Figure 49: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the WebKB Dataset at a Query Ratio of 0.1.



Figure 50: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the WebKB Dataset at a Query Ratio of 0.3.



Figure 51: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the WebKB Dataset at a Query Ratio of 0.5.

Figure 52: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the WebKB Dataset at a Query Ratio of 0.7.



Figure 53: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the WebKB Dataset at a Query Ratio of 0.9.



Figure 54: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Audio Dataset at a Query Ratio of 0.1.
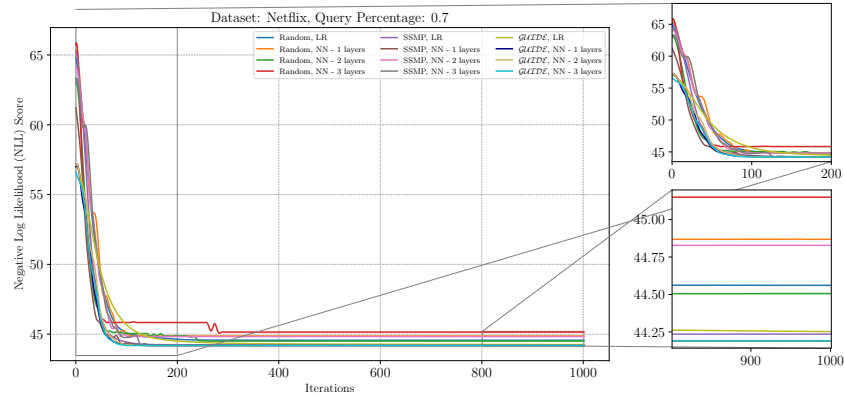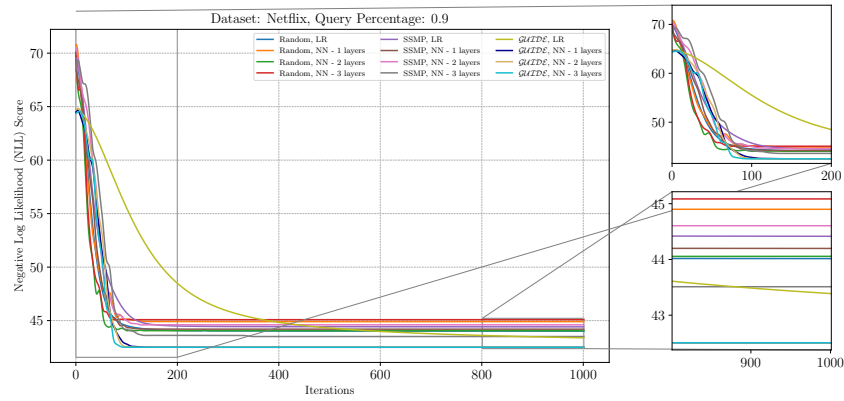
Figure 55: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Audio Dataset at a Query Ratio of 0.3.



Figure 56: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Audio Dataset at a Query Ratio of 0.5.



Figure 57: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Audio Dataset at a Query Ratio of 0.7.

Figure 58: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Audio Dataset at a Query Ratio of 0.9.
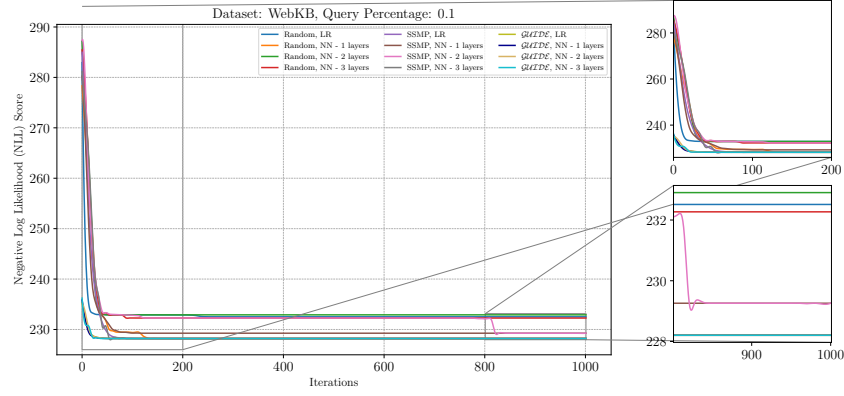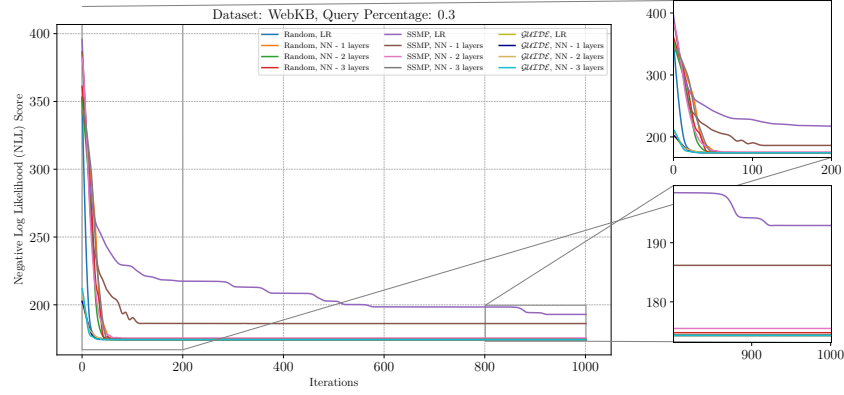


Figure 59: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Movie reviews Dataset at a Query Ratio of 0.1.



Figure 60: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Movie reviews Dataset at a Query Ratio of 0.3.

Figure 61: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Movie reviews Dataset at a Query Ratio of 0.5.



Figure 62: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Movie reviews Dataset at a Query Ratio of 0.7.



Figure 63: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Movie reviews Dataset at a Query Ratio of 0.9.

Figure 64: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Jester Dataset at a Query Ratio of 0.1.



Figure 65: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Jester Dataset at a Query Ratio of 0.3.
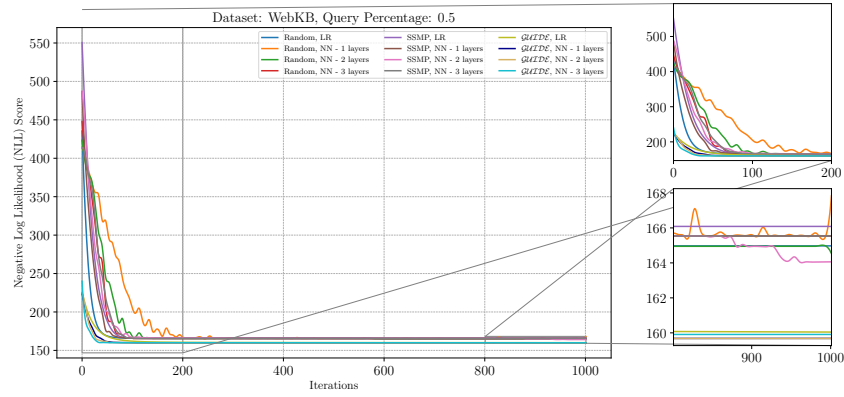


Figure 66: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Jester Dataset at a Query Ratio of 0.5.
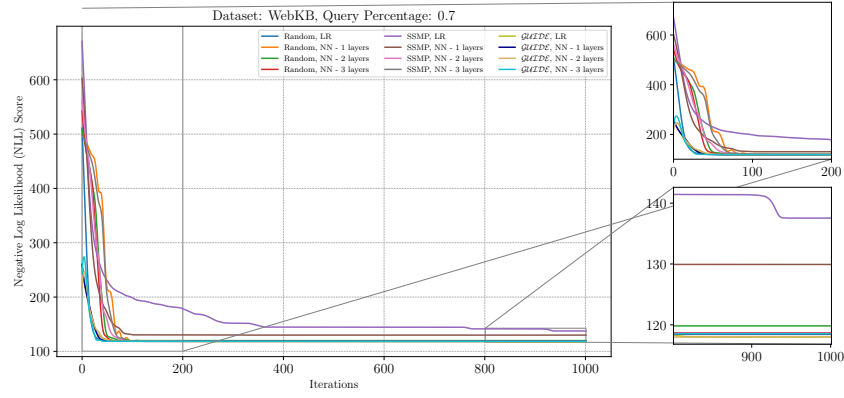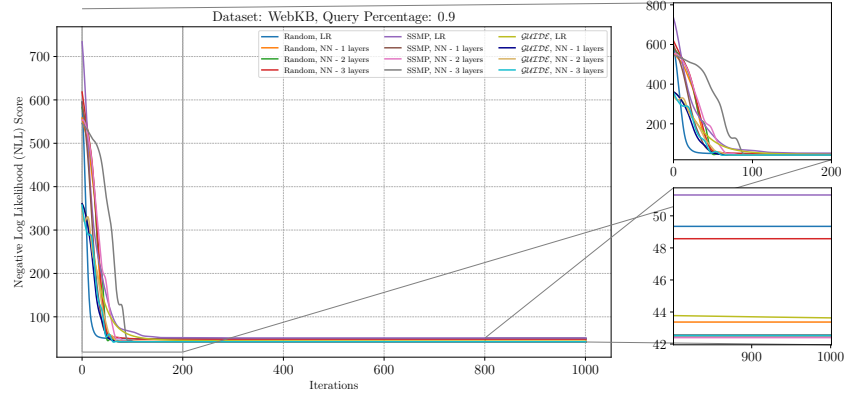
Figure 67: Analysis of ITSELF Loss Across Various Pre-Trained Models for PCs on the Jester Dataset at a Query Ratio of 0.7.
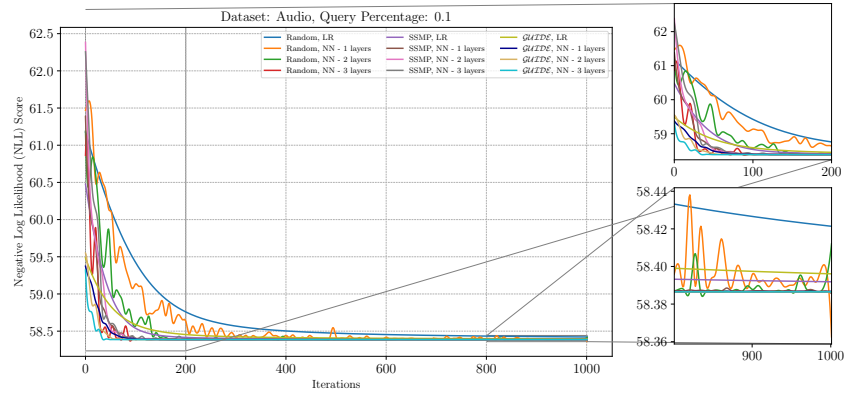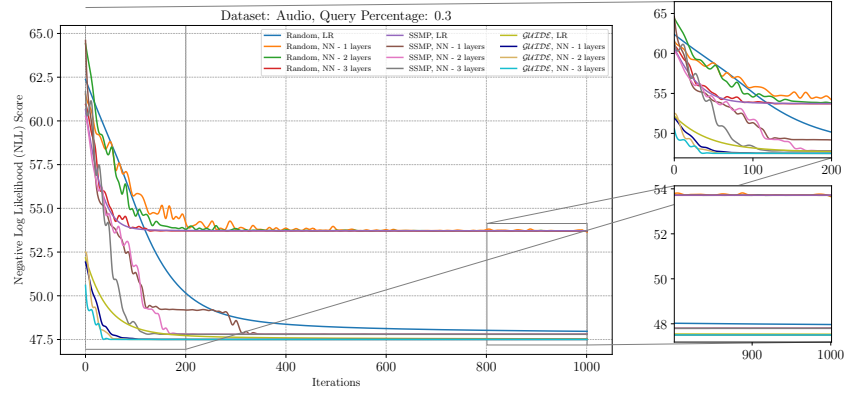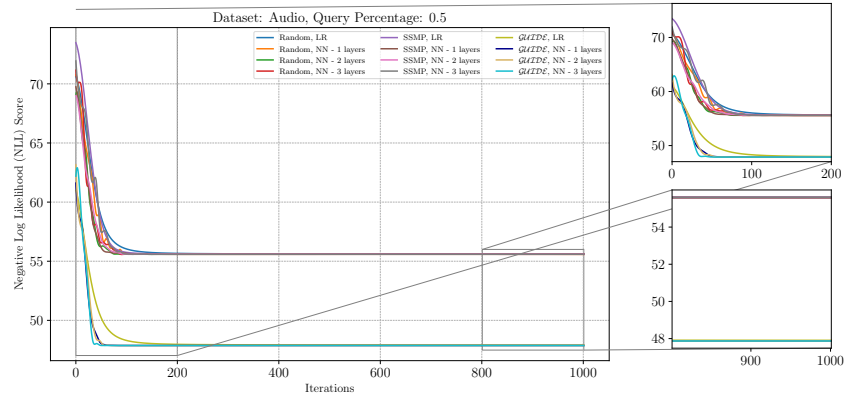
Figure 68: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f10.wrap Dataset at a Query Ratio of 0.7.



Figure 69: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f10.wrap Dataset at a Query Ratio of 0.9.

Figure 70: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f15.wrap Dataset at a Query Ratio of 0.5.



Figure 71: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f15.wrap Dataset at a Query Ratio of 0.7.



Figure 72: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f15.wrap Dataset at a Query Ratio of 0.9.

Figure 73: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f2.wrap Dataset at a Query Ratio of 0.5.
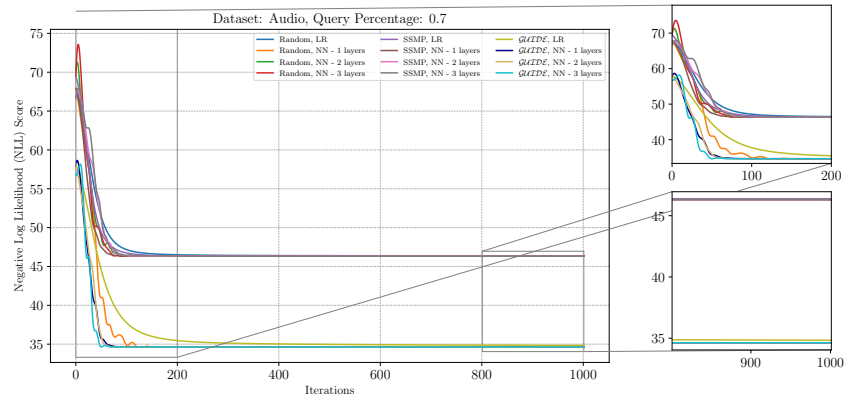


Figure 74: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f2.wrap Dataset at a Query Ratio of 0.7.



Figure 75: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f2.wrap Dataset at a Query Ratio of 0.9.

Figure 76: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f5.wrap Dataset at a Query Ratio of 0.5.



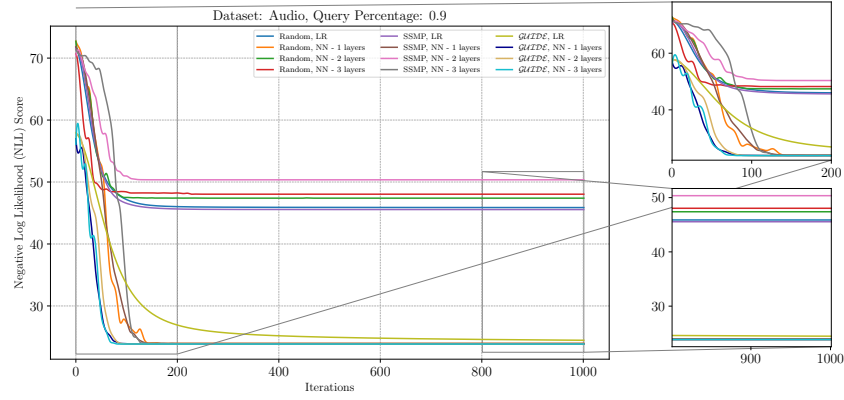Figure 77: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f5.wrap Dataset at a Query Ratio of 0.7.



Figure 78: Analysis of ITSELF Loss Across Various Pre-Trained Models for PGMs on the grid40x40.f5.wrap Dataset at a Query Ratio of 0.9.
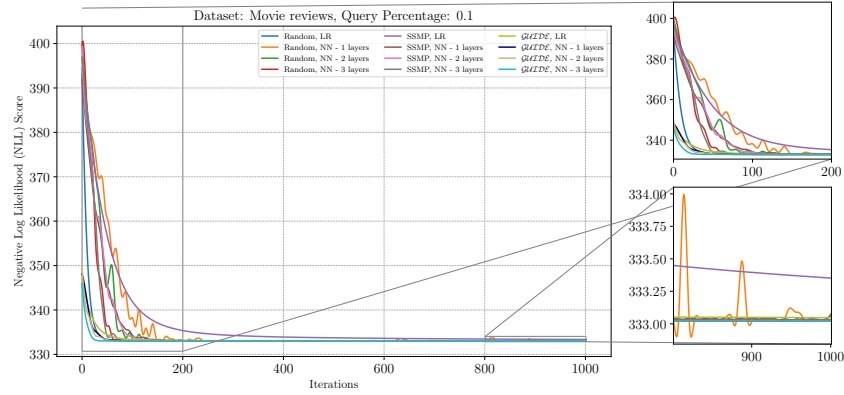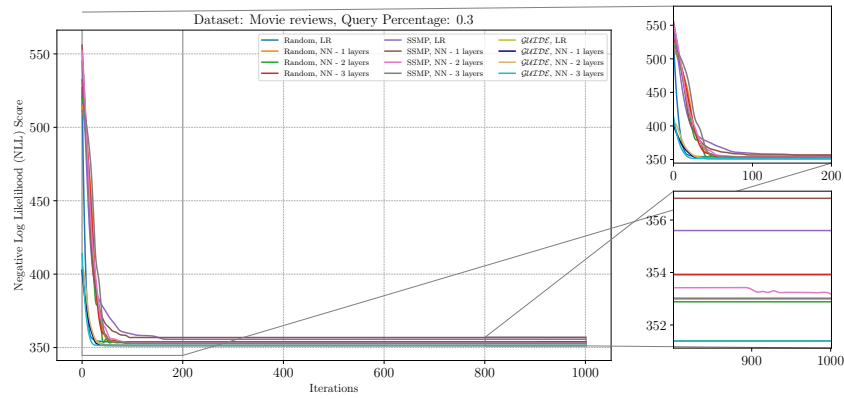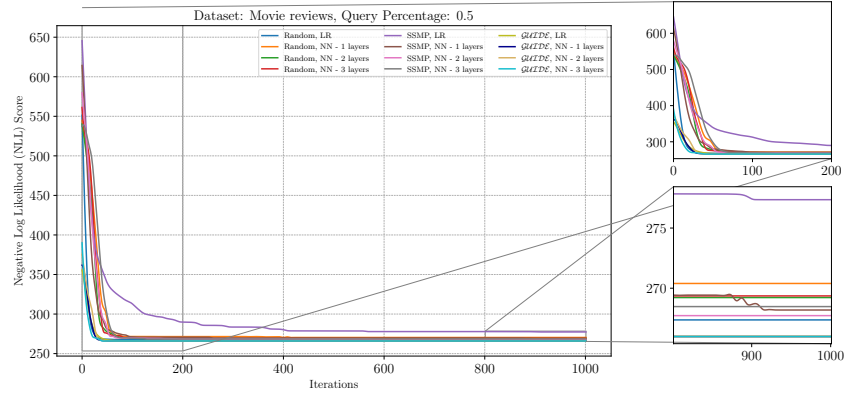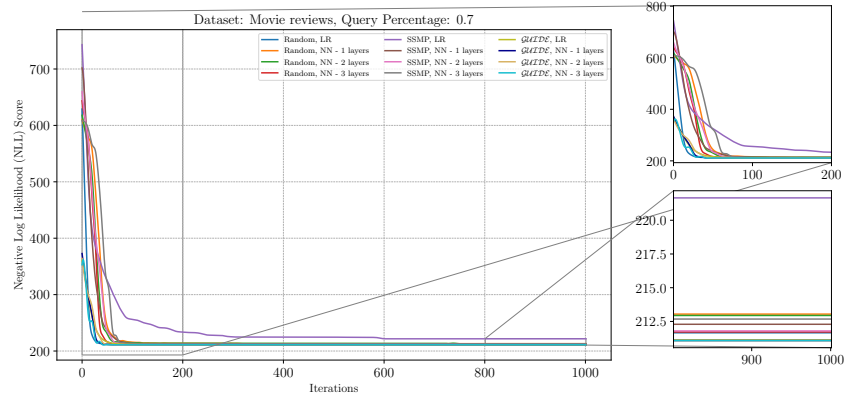
# D Inference Time Comparison



Figure 79: Heatmap depicting the inference time for MADE on a logarithmic microsecond scale, where a lighter color denotes shorter (more favorable) durations.



Figure 80: Heatmap depicting the inference time for PC on a logarithmic microsecond scale, where a lighter color denotes shorter (more favorable) durations.

We present the inference times for all baselines and proposed methods in Figures 79 to 81. Figure 79 details the inference times for MADE, while Figures 80 and 81 respectively illustrate the times for PCs and PGMs. This comparison facilitates a direct evaluation of the computational efficiency across different models.

Figure 81: Heatmap depicting the inference time for PGM on a logarithmic microsecond scale, where a lighter color denotes shorter (more favorable) durations.

Each cell displays the natural logarithm of the time, measured in microseconds, for each method and dataset. Lighter colors indicate lower values. Notably, inferences using SSMP and $\mathcal{GUIDE}$ require the shortest time, as these methods necessitate only a single forward pass through the neural network to obtain the values for the query variables.

For MADE, the subsequent fastest method employs a model trained with $\mathcal{GUIDE}$ and conducts inference using ITSELF, outperforming the approach that uses SSMP for training. This advantage stems from the reduced number of ITSELF iterations required by $\mathcal{GUIDE}$, benefiting from a more effectively trained model. In PGMs, a similar pattern emerges with $\mathcal{GUIDE}$ + ITSELF as the next fastest method, followed by SSMP + ITSELF. For PCs, MAX ranks as the next fastest, closely followed by the $\mathcal{GUIDE}$ + ITSELF and SSMP + ITSELF methods. Finally, the ML and Seq methods display the highest inference times.

Thus, if you require a highly efficient method capable of performing inference in a fraction of a millisecond, $\mathcal{GUIDE}$ is the optimal choice. It outperforms the baseline for both MADE and PGMs. However, if higher log-likelihood scores are necessary, $\mathcal{GUIDE}$ + ITSELF would be suitable, as it generally surpasses the baselines in speed and performance across various scenarios.

# E  Gap Analysis For PGM

Table 2 presents the log-likelihood score gap between the neural network methods (SSMP, $\mathcal{GUIDE}$, SSMP + ITSELF, $\mathcal{GUIDE}$ + ITSELF) and exact solutions. These exact solutions are obtained using AOBB, which provides near-optimal results for smaller datasets. For each approach M, the gap is calculated as the relative difference between the score of the near-optimal solution (determined by AOBB) and the score achieved by M. This approach is feasible due to the use of small datasets, allowing identification of exact solutions.

The final column highlights the neural-based approach achieving the best performance for each dataset and query ratio combination. Notably, $\mathcal{GUIDE}$ and ITSELF consistently surpass other neural baselines across almost all dataset-query pairs. This analysis provides a comprehensive assessment of the proposed methods relative to exact solutions on small datasets, enabling a direct comparison of their effectiveness.

44

Table 2: Gap Between AOBB And Other Methods.

| Method | Query Ratio | SSMP | $\mathcal{GUIDE}$ | SSMP + ITSELF | $\mathcal{GUIDE}$ + ITSELF | Best Method |
|---|---|---|---|---|---|---|
| **Grids-17** | 0.900 | 0.082 | 0.060 | 0.069 | 0.075 | $\mathcal{GUIDE}$ |
| **Grids-17** | 0.800 | 0.051 | 0.038 | 0.040 | 0.035 | $\mathcal{GUIDE}$ + ITSELF |
| **Grids-17** | 0.700 | 0.042 | 0.030 | 0.034 | 0.016 | $\mathcal{GUIDE}$ + ITSELF |
| **Grids-17** | 0.500 | 0.026 | 0.024 | 0.024 | 0.007 | $\mathcal{GUIDE}$ + ITSELF |
| **Grids-18** | 0.900 | 0.081 | 0.062 | 0.071 | 0.102 | $\mathcal{GUIDE}$ |
| **Grids-18** | 0.700 | 0.033 | 0.027 | 0.024 | 0.015 | $\mathcal{GUIDE}$ + ITSELF |
| **Grids-18** | 0.500 | 0.020 | 0.018 | 0.018 | 0.006 | $\mathcal{GUIDE}$ + ITSELF |
| **Grids-18** | 0.800 | 0.054 | 0.035 | 0.045 | 0.037 | $\mathcal{GUIDE}$ |
| **Segmentation-14** | 0.500 | 0.032 | 0.032 | 0.032 | 0.004 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-14** | 0.900 | 0.045 | 0.014 | 0.014 | 0.005 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-14** | 0.800 | 0.051 | 0.024 | 0.024 | 0.006 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-14** | 0.700 | 0.029 | 0.029 | 0.029 | 0.005 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-15** | 0.800 | 0.046 | 0.002 | 0.002 | 0.002 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-15** | 0.500 | 0.003 | 0.003 | 0.003 | 0.000 | $\mathcal{GUIDE}$ + ITSELF |
| **Segmentation-15** | 0.900 | 0.675 | 0.255 | 0.433 | 0.305 | $\mathcal{GUIDE}$ |
| **Segmentation-15** | 0.700 | 0.003 | 0.003 | 0.003 | 0.002 | $\mathcal{GUIDE}$ + ITSELF |

# F  Log Likelihood Scores Comparison

This section compares log-likelihood scores across baselines, SSMP, SSMP + ITSELF, $\mathcal{GUIDE}$ and $\mathcal{GUIDE}$ + ITSELF for all datasets and PMs. The log likelihood plots for NAMs are depicted in Figures 82 to 101, while those for PCs are illustrated in Figures 102 to 121. Each bar represents the mean log likelihood score of the corresponding method, with tick marks indicating the mean ± standard deviation. Higher values in these scores signify better performance by the method, considering they represent log likelihood scores.

## F.1  Scores for NAM

Figures 82 to 101 present the log likelihood scores for NAMs, illustrating the performance of ITSELF inference and $\mathcal{GUIDE}$ training relative to other baselines. The heatmaps and contingency tables discussed in the main paper corroborate the superior performance of $\mathcal{GUIDE}$ + ITSELF. These visual representations allow for a comprehensive understanding of the performance of our methods and baseline approaches, including HC and SSMP, across various datasets and query ratios.



Figure 82: Log-Likelihood Scores on NLTCS  for NAM. Higher Scores Indicate Better Performance.

## F.2  Scores for PCs

Analyzing Figures 102 to 121, which focuses on PCs, reveals similar patterns. The neural-based methods significantly outperform the MAX baseline. Among these, $\mathcal{GUIDE}$ + ITSELF surpasses all other polynomial-time baselines and neural methods in over 80 percent of the experiments. This demonstrates that ITSELF substantially enhances the chances of approaching optimal solutions by performing test-time optimization. When comparing traditional inference with ITSELF, ITSELF consistently proves superior. Moreover, $\mathcal{GUIDE}$ outperforms the other neural-based training methods (SSMP).

Figure 83: Log-Likelihood Scores on   for NAM. Higher Scores Indicate Better Performance.



Figure 84: Log-Likelihood Scores on KDDCup2k  for NAM. Higher Scores Indicate Better Performance.



Figure 85: Log-Likelihood Scores on Plants  for NAM. Higher Scores Indicate Better Performance.

Figure 86: Log-Likelihood Scores on Audio for NAM. Higher Scores Indicate Better Performance.



Figure 87: Log-Likelihood Scores on Jester for NAM. Higher Scores Indicate Better Performance.



Figure 88: Log-Likelihood Scores on Netflix for NAM. Higher Scores Indicate Better Performance.

Figure 89: Log-Likelihood Scores on Accidents  for NAM. Higher Scores Indicate Better Performance.



Figure 90: Log-Likelihood Scores on Mushrooms  for NAM. Higher Scores Indicate Better Performance.



Figure 91: Log-Likelihood Scores on Connect 4  for NAM. Higher Scores Indicate Better Performance.

Figure 92: Log-Likelihood Scores on RCV-1 for NAM. Higher Scores Indicate Better Performance.



Figure 93: Log-Likelihood Scores on Retail for NAM. Higher Scores Indicate Better Performance.



Figure 94: Log-Likelihood Scores on DNA for NAM. Higher Scores Indicate Better Performance.

Figure 95: Log-Likelihood Scores on Movie reviews for NAM. Higher Scores Indicate Better Performance.



Figure 96: Log-Likelihood Scores on Book for NAM. Higher Scores Indicate Better Performance.



Figure 97: Log-Likelihood Scores on WebKB for NAM. Higher Scores Indicate Better Performance.

Figure 98: Log-Likelihood Scores on Reuters-52 for NAM. Higher Scores Indicate Better Performance.



Figure 99: Log-Likelihood Scores on 20 NewsGroup for NAM. Higher Scores Indicate Better Performance.



Figure 100: Log-Likelihood Scores on Ad for NAM. Higher Scores Indicate Better Performance.

Figure 101: Log-Likelihood Scores on BBC for NAM. Higher Scores Indicate Better Performance.



Figure 102: Log-Likelihood Scores on NLTCS for PCs. Higher Scores Indicate Better Performance.



Figure 103: Log-Likelihood Scores on for PCs. Higher Scores Indicate Better Performance.

Figure 104: Log-Likelihood Scores on KDDCup2k for PCs. Higher Scores Indicate Better Performance.



Figure 105: Log-Likelihood Scores on Plants for PCs. Higher Scores Indicate Better Performance.



Figure 106: Log-Likelihood Scores on Audio for PCs. Higher Scores Indicate Better Performance.

Figure 107: Log-Likelihood Scores on Jester for PCs. Higher Scores Indicate Better Performance.



Figure 108: Log-Likelihood Scores on Netflix for PCs. Higher Scores Indicate Better Performance.



Figure 109: Log-Likelihood Scores on Accidents for PCs. Higher Scores Indicate Better Performance.

Figure 110: Log-Likelihood Scores on Mushrooms for PCs. Higher Scores Indicate Better Performance.



Figure 111: Log-Likelihood Scores on Connect 4 for PCs. Higher Scores Indicate Better Performance.



Figure 112: Log-Likelihood Scores on RCV-1 for PCs. Higher Scores Indicate Better Performance.

Figure 113: Log-Likelihood Scores on Retail for PCs. Higher Scores Indicate Better Performance.



Figure 114: Log-Likelihood Scores on DNA for PCs. Higher Scores Indicate Better Performance.



Figure 115: Log-Likelihood Scores on Movie reviews for PCs. Higher Scores Indicate Better Performance.

Figure 116: Log-Likelihood Scores on Book for PCs. Higher Scores Indicate Better Performance.



Figure 117: Log-Likelihood Scores on WebKB for PCs. Higher Scores Indicate Better Performance.



Figure 118: Log-Likelihood Scores on Reuters-52 for PCs. Higher Scores Indicate Better Performance.

Figure 119: Log-Likelihood Scores on 20 NewsGroup for PCs. Higher Scores Indicate Better Performance.



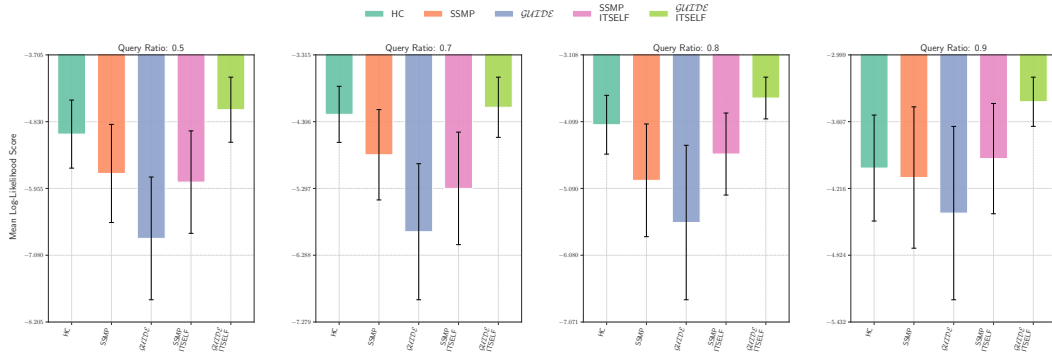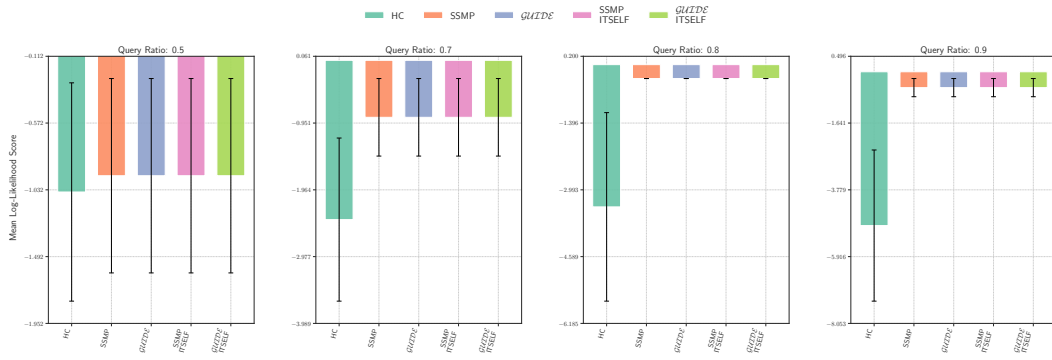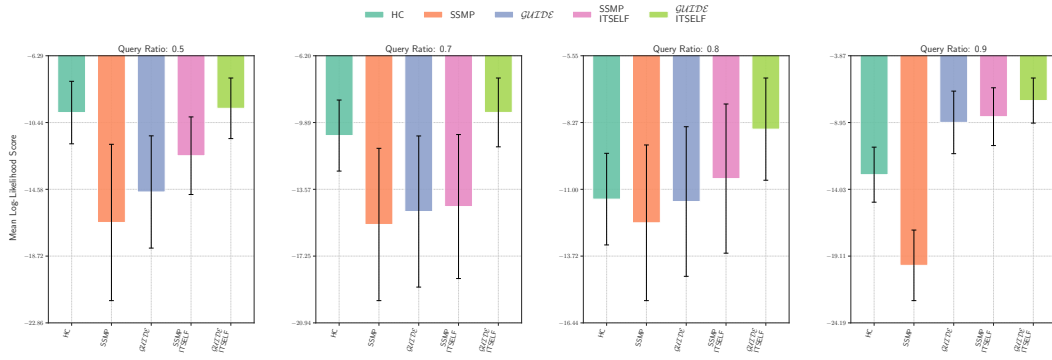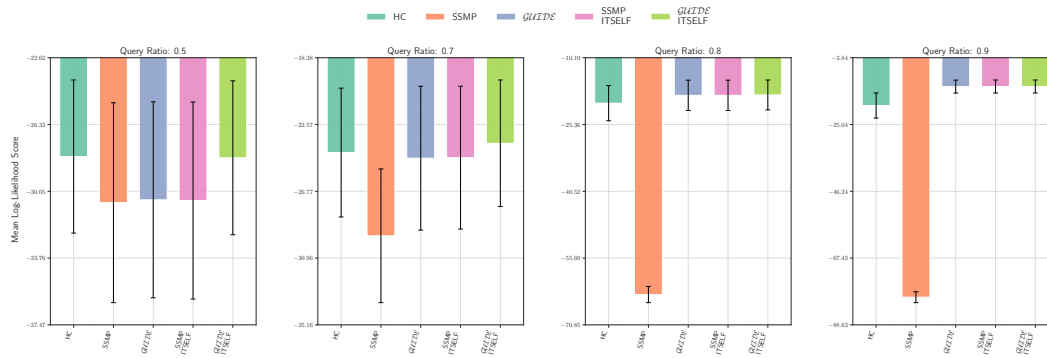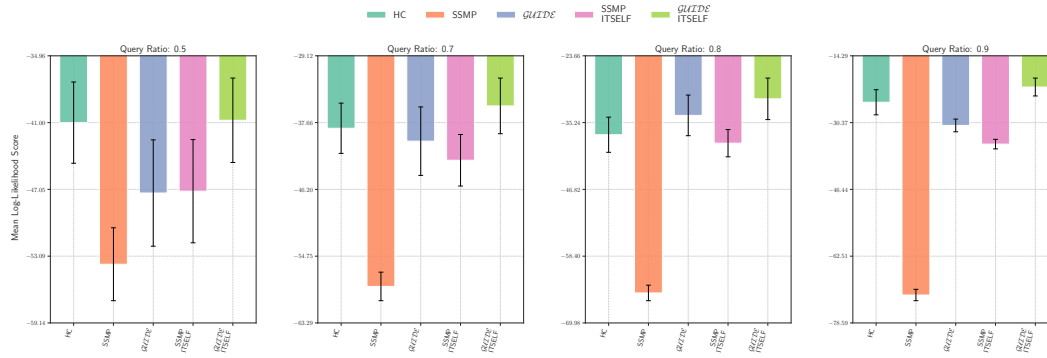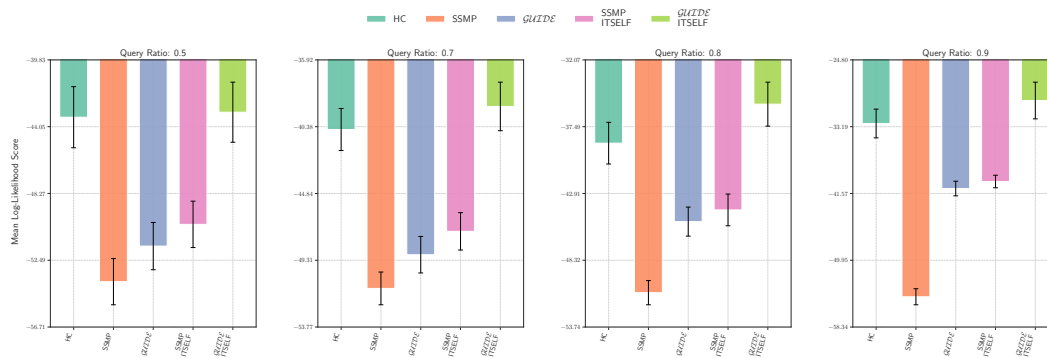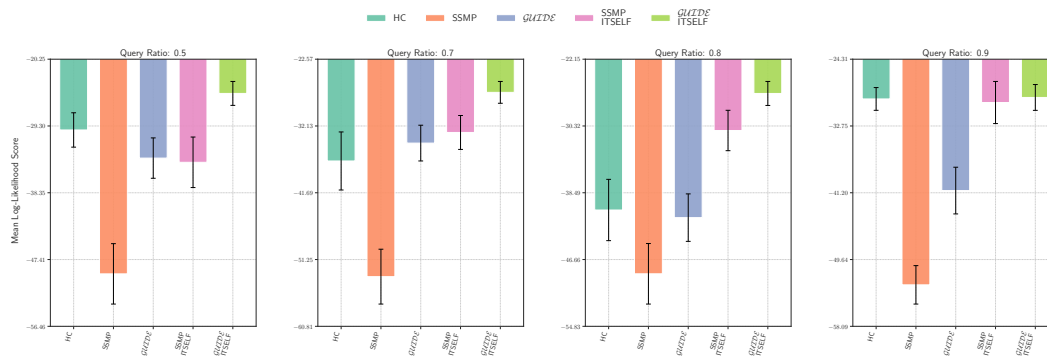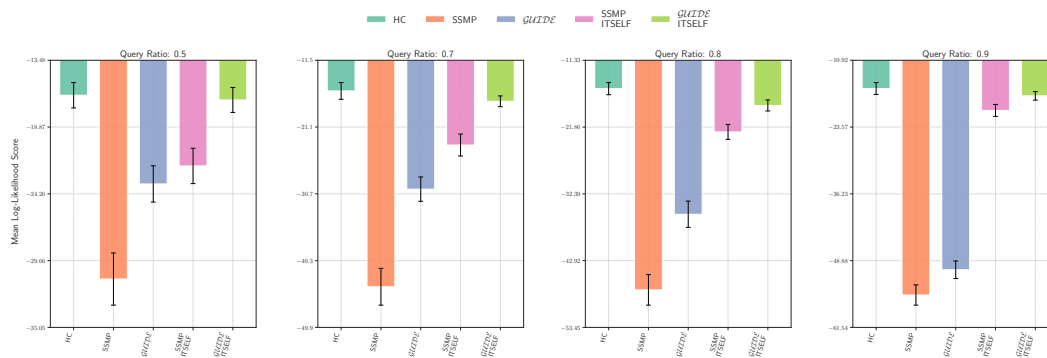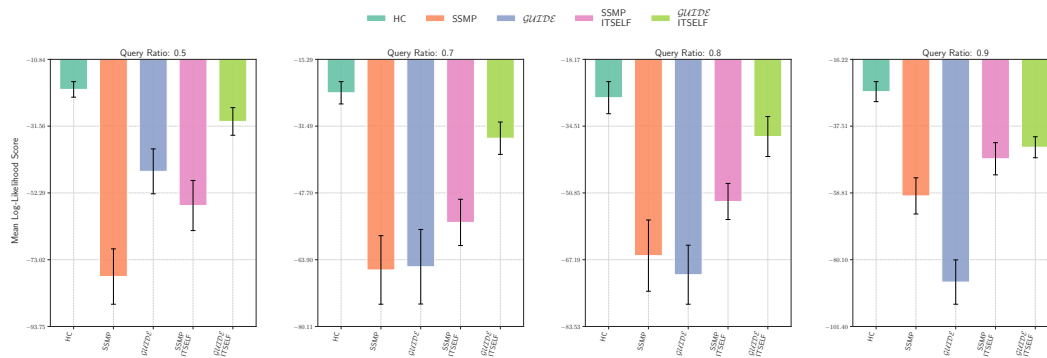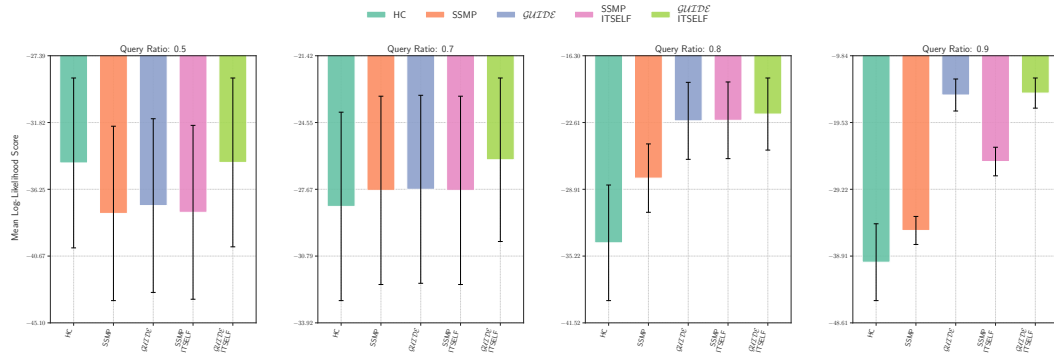Figure 120: Log-Likelihood Scores on Ad for PCs. Higher Scores Indicate Better Performance.



Figure 121: Log-Likelihood Scores on BBC for PCs. Higher Scores Indicate Better Performance.

# 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.

# Neural Network Approximators for Marginal MAP in Probabilistic Circuits

**Shivvrat Arya, Tahrima Rahman, Vibhav Gogate**

The University of Texas at Dallas
{shivvrat.arya, tahrima.rahman, vibhav.gogate}@utdallas.edu

## Abstract

Probabilistic circuits (PCs) such as sum-product networks efficiently represent large multi-variate probability distributions. They are preferred in practice over other probabilistic representations, such as Bayesian and Markov networks, because PCs can solve marginal inference (MAR) tasks in time that scales linearly in the size of the network. Unfortunately, the most probable explanation (MPE) task and its generalization, the marginal maximum-a-posteriori (MMAP) inference task remain NP-hard in these models. Inspired by the recent work on using neural networks for generating near-optimal solutions to optimization problems such as integer linear programming, we propose an approach that uses neural networks to approximate MMAP inference in PCs. The key idea in our approach is to approximate the cost of an assignment to the query variables using a continuous multilinear function and then use the latter as a loss function. The two main benefits of our new method are that it is self-supervised, and after the neural network is learned, it requires only linear time to output a solution. We evaluate our new approach on several benchmark datasets and show that it outperforms three competing linear time approximations: max-product inference, max-marginal inference, and sequential estimation, which are used in practice to solve MMAP tasks in PCs.

## Introduction

Probabilistic circuits (PCs) (Choi, Vergari, and Van den Broeck 2020) such as sum-product networks (SPNs) (Poon and Domingos 2011), arithmetic circuits (Darwiche 2003), AND/OR graphs (Dechter and Mateescu 2007), cutset networks (Rahman, Kothalkar, and Gogate 2014), and probabilistic sentential decision diagrams (Kisa et al. 2014) represent a class of *tractable probabilistic models* which are often used in practice to compactly encode a large multi-dimensional joint probability distribution. Even though all of these models admit linear time computation of marginal probabilities (MAR task), only some of them (Vergari et al. 2021; Peharz 2015), specifically those without any latent variables or having specific structural properties, e.g., cutset networks, selective SPNs (Peharz et al. 2016), AND/OR graphs having small contexts, etc., admit tractable most

probable explanation (MPE) inference[1].

However, none of these expressive PCs can efficiently solve the *marginal maximum-a-posteriori (MMAP) task* (Peharz 2015; Vergari et al. 2021), a task that combines MAR and MPE inference. More specifically, the distinction between MPE and MMAP tasks is that, given observations over a subset of variables (evidence), the MPE task aims to find the most likely assignment to all the non-evidence variables. In contrast, in the MMAP task, the goal is to find the most likely assignment to a subset of non-evidence variables known as the query variables, while marginalizing out non-evidence variables that are not part of the query. The MMAP problem has numerous real-world applications, especially in health care, natural language processing, computer vision, linkage analysis and diagnosis where hidden variables are present and need to be marginalized out (Bioucas-Dias and Figueiredo 2016; Kiselev and Poupart 2014; Lee, Marinescu, and Dechter 2014; Ping, Liu, and Ihler 2015).

In terms of computational complexity, both MPE and MMAP tasks are at least NP-hard in SPNs, a popular class of PCs (Peharz 2015; Conaty, de Campos, and Mauá 2017). Moreover, it is also NP-hard to approximate MMAP in SPNs to $2^{n^\delta}$ for fixed $0 \le \delta < 1$, where $n$ is the input size (Conaty, de Campos, and Mauá 2017; Mei, Jiang, and Tu 2018). It is also known that the MMAP task is much harder than the MPE task and is NP-hard even on models such as cutset networks and AND/OR graphs that admit linear time MPE inference (Park and Darwiche 2004; de Campos 2011).

To date, both exact and approximate methods have been proposed in literature for solving the MMAP task in PCs. Notable exact methods include branch-and-bound search (Mei, Jiang, and Tu 2018), reformulation approaches which encode the MMAP task as other combinatorial optimization problems with widely available solvers (Mauá et al. 2020) and circuit transformation and pruning techniques (Choi, Friedman, and Van den Broeck 2022). These methods can be quite slow in practice and are not applicable when fast, real-time inference is desired. As a result, approximate approaches that require only a few passes over the PC are often used in practice. A popular approximate approach is to

---

[1]The MPE inference task is also called full maximum-a-posteriori (full MAP) inference in literature. In this paper, we adopt the convention of calling it MPE.

compute an MPE solution over both the query and unobserved variables and then project the MPE solution over the query variables (Poon and Domingos 2011; Rahman, Jin, and Gogate 2019). Although this approach can provide fast answers at query time, it often yields MMAP solutions that are far from optimal.

In this paper, we propose to address the limitations of existing approximate methods for MMAP inference in PCs by using neural networks (NNs), leveraging recent work in the *learning to optimize* literature (Li and Malik 2016; Fioretto, Mak, and Hentenryck 2020; Donti, Rolnick, and Kolter 2020; Zamzam and Baker 2020; Park and Hentenryck 2023). In particular, several recent works have shown promising results in using NNs to solve both constrained and unconstrained optimization problems (see Park and Hentenryck (2023) and the references therein).

The high-level idea in these works is the following: given data, train NNs, either in a supervised or self-supervised manner, and then use them at test time to predict high-quality, near-optimal solutions to future optimization problems. A number of reasons have motivated this idea of *learning to optimize* using NNs: 1) NNs are good at approximating complex functions (distributions), 2) once trained, they can be faster at answering queries than search-based approaches, and 3) with ample data, NNs can learn accurate mappings of inputs to corresponding outputs. This has led researchers to employ NNs to approximately answer probabilistic inference queries such as MAR and MPE in Bayesian and Markov networks (Yoon et al. 2019; Cui et al. 2022). To the best of our knowledge, there is no prior work on solving MMAP in BNs, MNs, or PCs using NNs.

This paper makes the following contributions. First, we propose to learn a neural network (NN) approximator for solving the MMAP task in PCs. Second, by leveraging the tractability of PCs, we devise a loss function whose gradient can be computed in time that scales linearly in the size of the PC, allowing fast gradient-based algorithms for learning NNs. Third, our method trains an NN in a self-supervised manner without having to rely on pre-computed solutions to arbitrary MMAP problems, thus circumventing the need to solve intractable MMAP problems in practice. Fourth, we demonstrate via a large-scale experimental evaluation that our proposed NN approximator yields higher quality MMAP solutions as compared to existing approximate schemes.

## Preliminaries

We use upper case letters (e.g., $X$) to denote random variables and corresponding lower case letters (e.g., $x$) to denote an assignment of a value to a variable. We use bold upper case letters (e.g., $\mathbf{X}$) to denote a set of random variables and corresponding bold lower case letters (e.g., $\mathbf{x}$) to denote an assignment of values to all variables in the set. Given an assignment $\mathbf{x}$ to all variables in $\mathbf{X}$ and a variable $Y \in \mathbf{X}$, let $\mathbf{x}_Y$ denote the projection of $\mathbf{x}$ on $Y$. We assume that all random variables take values from the set $\{0, 1\}$; although note that it is easy to extend our method to multi-valued variables.

## Probabilistic Circuits

A probabilistic circuit (PC) $\mathcal{M}$ (Choi, Vergari, and Van den Broeck 2020) defined over a set of variables $\mathbf{X}$ represents a joint probability distribution over $\mathbf{X}$ using a rooted directed acyclic graph. The graph consists of three types of nodes: internal sum nodes that are labeled by $+$, internal product nodes that are labeled by $\times$, and leaf nodes that are labeled by either $X$ or $\neg X$ where $X \in \mathbf{X}$. Sum nodes represent conditioning, and an edge into a sum node $n$ from its child node $m$ is labeled by a real number $\omega(m, n) > 0$. Given an internal node (either a sum or product node) $n$, let $\mathtt{ch}(n)$ denote the set of children of $n$. We assume that each sum node $n$ is normalized and satisfies the following property: $\sum_{m \in \mathtt{ch}(n)} \omega(m, n) = 1$.

In this paper, we focus on a class of PCs which are *smooth and decomposable* (Choi, Vergari, and Van den Broeck 2020; Vergari et al. 2021). Examples of such PCs include sum-product networks (Poon and Domingos 2011; Rahman and Gogate 2016b), mixtures of cutset networks (Rahman, Kothalkar, and Gogate 2014; Rahman and Gogate 2016a), and arithmetic circuits obtained by compiling probabilistic graphical models (Darwiche 2003). These PCs admit tractable marginal inference, a key property that we leverage in our proposed method.

**Definition 1.** *We say that a sum or a product node $n$ is defined over a variable $X$ if there exists a directed path from $n$ to a leaf node labeled either by $X$ or $\neg X$. A PC is **smooth** if each sum node is such that its children are defined over the same set of variables. A PC is **decomposable** if each product node is such that its children are defined over disjoint subsets of variables.*

**Example 1.** *Figure 1(a) shows a smooth and decomposable probabilistic circuit defined over $\mathbf{X} = \{X_1, \ldots, X_4\}$.*

## Marginal Inference in PCs

Next, we describe how to compute the probability of an assignment to a subset of variables in a smooth and decomposable PC. This task is called the marginal inference (MAR) task. We begin by describing some additional notation.

Given a PC $\mathcal{M}$ defined over $\mathbf{X}$, let $\mathcal{S}$, $\mathcal{P}$ and $\mathcal{L}$ denote the set of sum, product and leaf nodes of $\mathcal{M}$ respectively. Let $\mathbf{Q} \subseteq \mathbf{X}$. Given a node $m$ and an assignment $\mathbf{q}$, let $v(m, \mathbf{q})$ denote the value of $m$ given $\mathbf{q}$. Given a leaf node $n$, let $\mathtt{var}(n)$ denote the variable associated with $n$ and let $l(n, \mathbf{q})$ be a function, which we call *leaf function*, that is defined as follows. $l(n, \mathbf{q})$ equals 0 if any of the following two conditions are satisfied: (1) the label of $n$ is $Q$ where $Q \in \mathbf{Q}$ and $\mathbf{q}$ contains the assignment $Q = 0$; and (2) if the label of $n$ is $\neg Q$ and $\mathbf{q}$ contains the assignment $Q = 1$. Otherwise, it is equal to 1. Intuitively, the leaf function assigns all leaf nodes that are inconsistent with the assignment $\mathbf{q}$ to 0 and the remaining nodes, namely those that are consistent with $\mathbf{q}$ and those that are not part of the query to 1.

Under this notation, and given a leaf function $l(n, \mathbf{q})$, the marginal probability of any assignment $\mathbf{q}$ w.r.t $\mathcal{M}$, and denoted by $\mathrm{p}_{\mathcal{M}}(\mathbf{q})$ can be computed by performing the follow-

Figure 1: (a) An example smooth and decomposable PC. The figure also shows value computation for answering the query $p_{\mathcal{M}}(X_3 = 1, X_4 = 0)$. The values of the leaf, sum, and product nodes are given in parentheses on their bottom, top, and left, respectively. The value of the root node is the answer to the query. (b) QPC obtained from the PC given in (a) for query variables $\{X_3, X_4\}$. For simplicity, here, we use an MMAP problem without any evidence. This is because a given evidence can be incorporated into the PC by appropriately setting the leaf nodes. We also show value computations for the following leaf initialization: $X_3^c = 0.99, \neg X_3^c = 0.01, X_4^c = 0.05, \neg X_4^c = 0.95$ and all other leaves are set to 1.

ing recursive *value computations*:

$$v(n, \mathbf{q}) = \begin{cases} l(n, \mathbf{q}) & \text{if } n \in \mathcal{L} \\ \sum_{m \in \text{ch}(n)} \omega(m, n) v(m, \mathbf{q}) & \text{if } n \in \mathcal{S} \\ \prod_{m \in \text{ch}(n)} v(m, \mathbf{q}) & \text{if } n \in \mathcal{P} \end{cases} \quad (1)$$

Let $r$ denote the root node of $\mathcal{M}$. Then, the probability of $\mathbf{q}$ w.r.t. $\mathcal{M}$, denoted by $p_{\mathcal{M}}(\mathbf{q})$ equals $v(r, \mathbf{q})$. Note that if $\mathbf{Q} = \mathbf{X}$, then $v(r, \mathbf{x})$ denotes the probability of the joint assignment $\mathbf{x}$ to all variables in the PC. Thus

$$v(r, \mathbf{q}) = \sum_{\mathbf{y} \in \{0,1\}^{|\mathbf{Y}|}} v(r, (\mathbf{q}, \mathbf{y}))$$

where $\mathbf{Y} = \mathbf{X} \setminus \mathbf{Q}$ and the notation $(\mathbf{q}, \mathbf{y})$ denotes the *composition* of the assignments to $\mathbf{Q}$ and $\mathbf{Y}$ respectively.

Since the recursive value computations require only one bottom-up pass over the PC, MAR inference is tractable or linear time in smooth and decomposable PCs.

**Example 2.** *Figure 1(a) shows bottom-up, recursive value computations for computing the probability of the assignment $(X_3 = 1, X_4 = 0)$ in our running example. Here, the leaf nodes $\neg X_3$ and $X_4$ are assigned to 0 and all other leaf nodes are assigned to 1. The number in parentheses at the top, left, and bottom of each sum, product and leaf nodes respectively shows the value of the corresponding node. The value of the root node equals $p_{\mathcal{M}}(X_3 = 1, X_4 = 0)$.*

## Marginal Maximum-a-Posteriori (MMAP) Inference in PCs

Given a PC $\mathcal{M}$ defined over $\mathbf{X}$, let $\mathbf{E} \subseteq \mathbf{X}$ and $\mathbf{Q} \subseteq \mathbf{X}$ denote the set of evidence and query variables respectively such that $\mathbf{E} \cap \mathbf{Q} = \emptyset$. Let $\mathbf{H} = \mathbf{X} \setminus (\mathbf{Q} \cup \mathbf{E})$ denote the set of hidden variables. Given an assignment $\mathbf{e}$ to the evidence variables (called evidence), the MMAP task seeks to

find an assignment $\mathbf{q}$ to $\mathbf{Q}$ such that the probability of the assignment $(\mathbf{e}, \mathbf{q})$ w.r.t. $\mathcal{M}$ is maximized. Mathematically,

$$\text{MMAP}(\mathbf{Q}, \mathbf{e}) = \underset{\mathbf{q}}{\arg\max} \, p_{\mathcal{M}}(\mathbf{e}, \mathbf{q}) \quad (2)$$

$$= \underset{\mathbf{q}}{\arg\max} \sum_{\mathbf{h} \in \{0,1\}^{|\mathbf{H}|}} p_{\mathcal{M}}(\mathbf{e}, \mathbf{q}, \mathbf{h}) \quad (3)$$

If $\mathbf{H} = \emptyset$ (namely $\mathbf{Q}$ is the set of non-evidence variables), then MMAP corresponds to the most probable explanation (MPE) task. It is known that both MMAP and MPE tasks are at least NP-hard in smooth and decomposable PCs (Park and Darwiche 2004; de Campos 2011; Peharz 2015), and even NP-hard to approximate (Conaty, de Campos, and Mauá 2017; Mei, Jiang, and Tu 2018).

A popular approach to solve the MMAP task in PCs is to replace the sum ($\sum$) operator with the max operator during bottom-up, recursive value computations and then performing a second top-down pass to find the assignment (Poon and Domingos 2011).

## A Neural Optimizer for MMAP in PCs

In this section, we introduce a learning-based approach using deep neural networks (NNs) to approximately solve the MMAP problem in PCs. Formally, the NN represents a function $f_\theta(.)$ that is parameterized by $\theta$, and takes an assignment $\mathbf{e}$ over the evidence variables as input and outputs an assignment $\mathbf{q}$ over the query variables. Our goal is to design *generalizable, continuous loss functions* for updating the parameters of the NN such that once learned, at test time, given an assignment $\mathbf{e}$ to the evidence variables as input, the NN outputs near-optimal solutions to the MMAP problem.

In this paper, we assume that the sets of evidence $\left(\mathbf{E} = \{E_i\}_{i=1}^N\right)$ and query $\left(\mathbf{Q} = \{Q_j\}_{j=1}^M\right)$ variables are known *a priori* and do not change at both training and test

time. We leave as future work the generalization of our approach that can handle variable length, arbitrarily chosen evidence, and query sets. Also, note that our proposed method does not depend on the particular NN architecture used, and we only require that each output node is a continuous quantity in the range $[0, 1]$ and uses a differentiable activation function (e.g., the sigmoid function).

We can learn the parameters of the given NN either in a *supervised* manner or in a *self-supervised* manner. However, the supervised approach is impractical, as described below.

In the supervised setting, we assume that we are given training data $\mathcal{D} = \{\langle \mathbf{e}_1, \mathbf{q}_1^* \rangle, \ldots, \langle \mathbf{e}_d, \mathbf{q}_d^* \rangle\}$, where each input $\mathbf{e}_i$ is an assignment to the evidence variables, and each (label) $\mathbf{q}_i^*$ is an optimal solution to the corresponding MMAP task, namely $\mathbf{q}_i^* = \text{MMAP}(\mathbf{Q}, \mathbf{e}_i)$. We then use supervised loss functions such as the mean-squared-error (MSE) $\sum_{i=1}^{d} \|\mathbf{q}_i^* - \mathbf{q}_i^c\|_2^2 / d$ and the mean-absolute-error (MAE) $\sum_{i=1}^{d} \|\mathbf{q}_i^* - \mathbf{q}_i^c\|_1 / d$ where $\mathbf{q}_i^c$ is the predicted assignment (note that $\mathbf{q}_i^c$ is continuous), and standard gradient-based methods to learn the parameters. Although supervised approaches allow us to use simple-to-implement loss functions, they are *impractical* if the number of query variables is large because they require access to the exact solutions to several intractable MMAP problems[2]. We therefore propose to use a *self-supervised approach*.

## A Self-Supervised Loss Function for PCs

In the self-supervised setting, we need access to training data in the form of assignments to the evidence variables, i.e., $\mathcal{D}' = \{\mathbf{e}_1, \ldots, \mathbf{e}_d\}$. Since smooth and decomposable PCs admit perfect sampling, these assignments can be easily sampled from the PC via top-down AND/OR sampling (Gogate and Dechter 2012). The latter yields an assignment $\mathbf{x}$ over all the random variables in the PC. Then we simply project $\mathbf{x}$ on the evidence variables $\mathbf{E}$ to yield a training example $\mathbf{e}$. Because each training example can be generated in time that scales linearly with the size of the PC, in practice, our proposed self-supervised approach is likely to have access to much larger number of training examples compared to the supervised approach.

Let $\mathbf{q}^c$ denote the MMAP assignment predicted by the NN given evidence $\mathbf{e} \in \mathcal{D}'$ where $\mathbf{q}^c \in [0, 1]^M$. In MMAP inference, given $\mathbf{e}$, we want to find an assignment $\mathbf{q}$ such that $\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$ is maximized, namely, $-\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$ is minimized. Thus, a natural loss function that we can use is $-\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$. Unfortunately, the NN outputs a continuous vector $\mathbf{q}^c$ and as a result $p_{\mathcal{M}}(\mathbf{e}, \mathbf{q}^c)$ is not defined. Therefore, we cannot use $-\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q}^c)$ as a loss function.

One approach to circumvent this issue is to use a threshold (say $0.5$) to convert each continuous quantity in the range $[0,1]$ to a binary one. A problem with this approach is that the threshold function is not differentiable.

Therefore, we propose to construct a smooth, differentiable loss function that given $\mathbf{q}^c = (q_1^c, \ldots, q_M^c)$ ap-

proximates $-\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$ where $\mathbf{q} = (q_1 = [q_1^c > 0.5], \ldots, q_M = [q_M^c > 0.5])$ and $[q_i^c > 0.5]$ is an indicator function which is 1 if $q_i^c > 0.5$ and 0 otherwise. The key idea in our approach is to construct a new PC, which we call *Query-specific PC* (QPC) by replacing all binary leaf nodes associated with the query variables in the original PC, namely those labeled by $Q$ and $\neg Q$ where $Q \in \mathbf{Q}$, with continuous nodes $Q^c \in [0, 1]$ and $\neg Q^c \in [0, 1]$. Then our proposed loss function is obtained using value computations (at the *root node* of the QPC) via a simple modification of the *leaf function* of the PC. At a high level, our new leaf function assigns each leaf node labeled by $Q_j^c$ such that $Q_j \in \mathbf{Q}$ to its corresponding estimate $q_j^c$, obtained from the NN and each leaf node labeled by $\neg Q_j^c$ such that $Q_j \in \mathbf{Q}$ to $1 - q_j^c$.

Formally, for the QPC, we propose to use leaf function $l'(n, (\mathbf{e}, \mathbf{q}^c))$ defined as follows:

1. If the label of $n$ is $Q_j^c$ such that $Q_j \in \mathbf{Q}$ then $l'(n, (\mathbf{e}, \mathbf{q}^c)) = q_j^c$.

2. If $n$ is labeled by $\neg Q_j^c$ such that $Q_j \in \mathbf{Q}$ then $l'(n, (\mathbf{e}, \mathbf{q}^c)) = 1 - q_j^c$.

3. If $n$ is labeled by $E_k$ such that $E_k \in \mathbf{E}$ and the assignment $E_k = 0$ is in $\mathbf{e}$ then $l'(n, (\mathbf{e}, \mathbf{q}^c)) = 0$.

4. If $n$ is labeled by $\neg E_k$ such that $E_k \in \mathbf{E}$ and the assignment $E_k = 1$ is in $\mathbf{e}$ then $l'(n, (\mathbf{e}, \mathbf{q}^c)) = 0$.

5. If conditions (1)-(4) are not met then $l'(n, (\mathbf{e}, \mathbf{q}^c)) = 1$.

The value of each node $n$ in the QPC, denoted by $v'(n, (\mathbf{e}, \mathbf{q}^c))$ is given by a similar recursion to the one given in Eq. (1) for PCs, except that the leaf function $l(n, \mathbf{q})$ is replaced by the new (continuous) leaf function $l'(n, (\mathbf{e}, \mathbf{q}^c))$. Formally, $v'(n, (\mathbf{e}, \mathbf{q}^c))$ is given by

$$
v'(n, (\mathbf{e}, \mathbf{q}^c))
$$
$$
= \begin{cases} l'(n, (\mathbf{e}, \mathbf{q}^c)) & \text{if } n \in \mathcal{L} \\ \sum_{m \in \text{ch}(n)} \omega(m, n) v'(m, (\mathbf{e}, \mathbf{q}^c)) & \text{if } n \in \mathcal{S} \\ \prod_{m \in \text{ch}(n)} v'(n, (\mathbf{e}, \mathbf{q}^c)) & \text{if } n \in \mathcal{P} \end{cases} \quad (4)
$$

Let $r$ denote the root node of $\mathcal{M}$, then we propose to use $-\ln v'(r, (\mathbf{e}, \mathbf{q}^c))$ as a loss function.

**Example 3.** *Figure 1(b) shows the QPC corresponding to the PC shown in Figure 1(a). We also show value computations for the assignment ($X_3^c = 0.99, X_4^c = 0.05$).*

## Tractable Gradient Computation

Our proposed loss function is smooth and continuous because by construction, it is a negative logarithm of a *multilinear function* over $\mathbf{q}^c$. Next, we show that the partial derivative of the function w.r.t. $q_j^c$ can be computed in linear time in the size of the QPC[3]. More specifically, in order to compute the partial derivative of QPC with respect to $\mathbf{q}_j^c$, we simply have to use a new leaf function which is identical to $l'$ except that if the label of a leaf node $n$ is $Q_j^c$ then we set its value to 1 (instead of $q_j^c$) and if it is $\neg Q_j^c$ then we set its value $-1$ (instead of $1 - q_j^c$). We then perform bottom-up recursive

---

[2]Note that the training data used to train the NN in the supervised setting is different from the training data used to learn the PC. In particular, in the data used to train the PC, the assignments to the query variables $\mathbf{Q}$ may not be optimal solutions of $\text{MMAP}(\mathbf{Q}, \mathbf{e})$.

[3]Recall that $q_j^c$ is an output node of the NN and therefore backpropagation over the NN can be performed in time that scales linearly with the size of the NN and the QPC

value computations over the QPC and the value of the root node is the partial derivative of the QPC with respect to $\mathbf{q}_j^c$. In summary, it is straight-forward to show that:

**Proposition 1.** *The gradient of the loss function* $-\ln v'(r, (\mathbf{e}, \mathbf{q}^c))$ *w.r.t.* $\mathbf{q}_j^c$ *can be computed in time and space that scales linearly with the size of $\mathcal{M}$.*

**Example 4.** *The partial derivative of the QPC given in figure 1(b) w.r.t. $x_3^c$ given $(X_3^c = 0.99, X_4^c = 0.05)$ can be obtained by setting the leaf nodes $X_3^c$ to 1 and $\neg X_3^c$ to $-1$, assigning all other leaves to the values shown in Figure 1(b) and then performing value computations. After the value computation phase, the value of the root node will equal the partial derivative of the QPC w.r.t. $x_3^c$.*

### Improving the Loss Function

As mentioned earlier, our proposed loss function is a continuous approximation of the discrete function $-\ln v(r, (\mathbf{e}, \mathbf{q}))$ where $\mathbf{q} = (q_1 = [q_1^c > 0.5], \ldots, q_M = [q_M^c > 0.5])$ and the difference between the two is minimized iff $\mathbf{q} = \mathbf{q}^c$. Moreover, since the set of continuous assignments includes the discrete assignments, it follows that:

$$\min_{\mathbf{q}^c} \left\{ -\ln v'(r, (\mathbf{e}, \mathbf{q}^c)) \right\} \leq \min_{\mathbf{q}} \left\{ -\ln v(r, (\mathbf{e}, \mathbf{q})) \right\} \quad (5)$$

Since the right-hand side of the inequality given in (5) solves the MMAP task, we can improve our loss function by tightening the lower bound. This can be accomplished using an entropy-based penalty, controlled by a hyper-parameter $\alpha > 0$, yielding the loss function

$$\ell(\mathbf{q}^c) = -\ln v'(r, (\mathbf{e}, \mathbf{q}^c)) -$$

$$\alpha \sum_{j=1}^{M} q_j^c \log(q_j^c) + (1 - q_j^c) \log(1 - q_j^c) \quad (6)$$

The second term in the expression given above is minimized when each $q_j^c$ is closer to 0 or 1 and is maximized when $q_j^c = 0.5$. Therefore, it encourages 0/1 (discrete) solutions. The hyperparameter $\alpha$ controls the magnitude of the penalty. When $\alpha = 0$, the above expression finds an assignment based on the continuous approximation $-\ln v'(r, (\mathbf{e}, \mathbf{q}^c))$. On the other hand, when $\alpha = \infty$ then only discrete solutions are possible yielding a non-smooth loss function. $\alpha$ thus helps us trade the smoothness of our proposed loss function with its distance to the true loss.

## Experiments

In this section, we describe and analyze the results of our comprehensive experimental evaluation for assessing the performance of our novel Self-Supervised learning based MMAP solver for PCs, referred to as SSMP hereafter. We begin by describing our experimental setup including competing methods, evaluation criteria, as well as NN architectures, datasets, and PCs used in our study.

### Competing Methods

We use *three polytime baseline methods* from the PC and probabilistic graphical models literature (Park and Darwiche 2004; Poon and Domingos 2011). We also compared the impact of using the solutions computed by the three baseline schemes as well our method SSMP as initial state for stochastic hill climbing search.

**Baseline 1: MAX Approximation (Max).** In this scheme (Poon and Domingos 2011), the MMAP assignment is derived by substituting sum nodes with max nodes. During the upward pass, a max node produces the maximum weighted value from its children instead of their weighted sum. Subsequently, the downward pass begins from the root and iteratively selects the highest-valued child of a max node (or one of them), along with all children of a product node.

**Baseline 2: Maximum Likelihood Approximation (ML)** (Park and Darwiche 2004) For each variable $Q \in \mathbf{Q}$, we first compute the marginal distribution $p_{\mathcal{M}}(Q|\mathbf{e})$ and then set $Q$ to $\arg\max_{j \in \{0,1\}} p_{\mathcal{M}}(Q = j|\mathbf{e})$.

**Baseline 3: Sequential Approximation (Seq)** In this scheme (Park and Darwiche 2004), we assign the query variables one by one until no query variables remain unassigned. At each step, we choose an unassigned query variable $Q_j \in \mathbf{Q}$ that maximizes the probability $p_{\mathcal{M}}(q_j|\mathbf{e}, \mathbf{y})$ for one of its values $q_j$ and assign it to $q_j$ where $\mathbf{y}$ represents the assignment to the previously considered query variables.

**Stochastic Hill Climbing Search.** We used the three baselines and our SSMP method as the initial state in stochastic hill climbing search for MMAP inference described in (Park and Darwiche 2004). The primary goal of this experiment is to assess whether our scheme can assist local search-based *anytime methods* in reaching better solutions than other heuristic methods for initialization. In our experiments, we ran stochastic hill climbing for 100 iterations for each MMAP problem.

### Evaluation Criteria

We evaluated the performance of the competing schemes along two dimensions: log-likelihood scores and inference times. Given evidence $\mathbf{e}$ and query answer $\mathbf{q}$, the log-likelihood score is given by $\ln p_{\mathcal{M}}(\mathbf{e}, \mathbf{q})$.

### Datasets and Probabilistic Circuits

We use twenty-two widely used binary datasets from the tractable probabilistic models' literature (Lowd and Davis 2010; Haaren and Davis 2012; Larochelle and Murray 2011; Bekker et al. 2015) (we call them TPM datasets) as well as the binarized MNIST (Salakhutdinov and Murray 2008), EMNIST (Cohen et al. 2017) and CIFAR-10 (Krizhevsky, Nair, and Hinton 2009) datasets. We used the DeeProb-kit library (Loconte and Gala 2022) to learn a sum-product network (our choice of PC) for each dataset. The number of nodes in these learned PCs ranges from 46 to 22027.

For each PC and each test example in the 22 TPM datasets, we generated two types of MMAP instances: MPE instances in which $\mathbf{H}$ is empty and MMAP instances in which $\mathbf{H}$ is not empty. We define query ratio, denoted by $qr$, as the fraction of variables that are part of the query set. For MPE, we selected $qr$ from $\{0.1, 0.3, 0.5, 0.6, 0.7, 0.8, 0.9\}$, and for MMAP, we replaced $0.9$ with $0.4$ to avoid small $\mathbf{H}$ and $\mathbf{E}$. For generating MMAP instances, we used 50% of

| | Initial | | | | | | | | Hill Climbing Search | | | | | | | |
| | MPE | | | | MMAP | | | | MPE | | | | MMAP | | | |
| | Max | SSMP | ML | Seq | Max | SSMP | ML | Seq | Max | SSMP | ML | Seq | Max | SSMP | ML | Seq |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Max | 0 | 64 | 33 | 14 | 0 | 46 | 23 | 10 | 0 | 40 | 13 | 9 | 0 | 27 | 10 | 16 |
| SSMP | 88 | 0 | 96 | 77 | 97 | 0 | 102 | 82 | 93 | 0 | 99 | 87 | 98 | 0 | 100 | 86 |
| ML | 6 | 49 | 0 | 15 | 3 | 34 | 0 | 10 | 19 | 37 | 0 | 14 | 12 | 26 | 0 | 17 |
| Seq | 105 | 63 | 105 | 0 | 117 | 53 | 117 | 0 | 85 | 44 | 82 | 0 | 89 | 39 | 90 | 0 |

Table 1: Contingency tables for competing methods across MPE and MMAP Problems, including initial and Hill Climbing Search comparisons. Highlighted values represent results for SSMP.

the remaining variables as evidence variables (and for MPE instances all remaining variables are evidence variables).

For the MNIST, EMNIST, and CIFAR-10 datasets, we used $qr = 0.7$ and generated MPE instances only. More specifically, we used the top 30% portion of the image as evidence, leaving the bottom 70% portion as query variables. Also, in order to reduce the training time for PCs, note that for these datasets, we learned a PC for each class, yielding a total of ten PCs for each dataset.

**Neural Network Optimizers**

For each PC and query ratio combination, we trained a corresponding neural network (NN) using the loss function described in the previous section. Because we have 22 TPM datasets and 7 query ratios for them, we trained 154 NNs for the MPE task and 154 for the MMAP task. For the CIFAR-10, MNIST and EMNIST datasets, we trained 10 NNs, one for each PC (recall that we learned a PC for each class).

Because our learning method does not depend on the specific choice of neural network architectures, we use a fixed neural network architecture across all experiments: fully connected with four hidden layers having 128, 256, 512, and 1024 nodes respectively. We used ReLU activation in the hidden layers, sigmoid in the output layer, dropout for regularization (Srivastava et al. 2014) and Adam optimizer (Kingma and Ba 2017) with a standard learning rate scheduler for 50 epochs. All NNs were trained using PyTorch (Paszke et al. 2019) on a single NVIDIA A40 GPU. We select a value for the hyperparameter $\alpha$ used in our loss function (see equation (6)) via 5-fold cross validation.

**Results on the TPM Datasets**

We summarize our results for the competing schemes (3 baselines and SSMP) on the 22 TPM datasets using the first two contingency tables given in Table 1, one for MPE and one for MMAP. Detailed results are provided in the supplementary material. Recall that we generated 154 test datasets each for MPE and MMAP (22 PCs × 7 $qr$ values). In all contingency tables, the number in the cell $(i, j)$ equals the number of times (out of 154) that the scheme in the $i$-th row was better in terms of average log-likelihood score than the scheme in the $j$-th column. The difference between 154 and the sum of the numbers in the cells $(i, j)$ and $(j, i)$ equals the number of times the scheme in the $i$-th row and $j$-th column had identical log-likelihood scores.

From the MPE contingency table given in Table 1, we observe that SSMP is superior to Max, ML, and Seq approximations. The Seq approximation is slightly better than

the Max approximation, and ML is the worst-performing scheme. For the harder MMAP task, we see a similar ordering among the competing schemes (see Table 1) with SSMP dominating other schemes. In particular, SSMP outperforms the Max and ML approximations in almost two-thirds of the cases and the Seq method in more than half of the cases.

We also investigate the effectiveness of SSMP and other baseline approaches when employed as initialization strategies for Hill Climbing Search. These findings are illustrated in the last two contingency tables given in Table 1. Notably, SSMP outperforms all other baseline approaches in nearly two-thirds of the experiments for both MPE and MMAP tasks. These results demonstrate that SSMP can serve as an effective initialization technique for anytime local search-based algorithms.



(a) MPE      (b) MMAP

Figure 2: Heat map showing the % difference in log-likelihood scores between SSMP and Max approximation. Blue represents Max's superiority (negative values) and red indicates SSMP better performance (positive values).

In Figure 2, via a heat-map representation, we show a more detailed performance comparison between SSMP and the Max approximation, which is a widely used baseline for MPE and MMAP inference in PCs. In the heat-map representation, the y-axis represents the datasets (ordered by the number of variables), while the x-axis shows the query ratio. The values in each cell represent the percentage difference between the mean log-likelihood scores of SSMP and

| | CIFAR | | | | MNIST | | | | EMNIST | | | |
|------|-----|------|----|-----|-----|------|----|-----|-----|------|----|-----|
| | Max | SSMP | ML | Seq | Max | SSMP | ML | Seq | Max | SSMP | ML | Seq |
| Max | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 5 |
| SSMP | 9 | 0 | 9 | 9 | 9 | 0 | 9 | 9 | 7 | 0 | 7 | 7 |
| ML | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 5 |
| Seq | 7 | 0 | 7 | 0 | 9 | 1 | 9 | 0 | 3 | 1 | 3 | 0 |

Table 2: Contingency tables comparing competing methods for MPE on CIFAR, MNIST and EMNIST datasets. Highlighted values represent results for SSMP.

the Max approximation. Formally, let $ll_{ssmp}$ and $ll_{max}$ denote the mean LL scores of SSMP and Max approximation respectively, then the percentage difference is given by

$$\%\text{Diff.} = \frac{ll_{ssmp} - ll_{max}}{|ll_{max}|} \times 100 \qquad (7)$$

From the heatmap for MPE given in Figure 2(a), we observe that SSMP is competitive with the Max approximation when the size of the query set is small. However, as the number of query variables increases, signaling a more challenging problem, SSMP consistently outperforms or has similar performance to the Max method across all datasets, except for accidents, pumsb-star, and book.

The heatmaps for MMAP are illustrated in Figure 2(b). We see a similar trend as the one for MPE; SSMP remains competitive with the Max approximation, particularly when the number of query variables is small. While SSMP outperforms (with some exceptions) the Max approximation when the number of query variables is large.

Finally, we present inference times in the supplement. On average SSMP requires in the order of 7-10 micro-seconds for MMAP inference on an A40 GPU. The Max approximation takes 7 milli-seconds (namely, SSMP is almost 1000 times faster). In comparison, as expected, the Seq and ML approximations are quite slow, requiring roughly 400 to 600 milliseconds to answer MPE and MMAP queries. In the case of our proposed method (SSMP), during the inference process, the size of the SPN holds no relevance; its time complexity is linear in the size of the neural network. On the contrary, for the alternative methods, the inference time is intricately dependent on the size of the SPN.

### Results on the CIFAR-10 Dataset

We binarized the CIFAR-10 dataset using a variational autoencoder having 512 bits. We then learned a PC for each of the 10 classes; namely, we learned a PC conditioned on the class variable. As mentioned earlier, we randomly set 70% of the variables as query variables. The contingency table for CIFAR-10 is shown in Table 2. We observe that SSMP dominates all competing methods while the Seq approximation is the second-best performing scheme (although note that Seq is computationally expensive).

### Results on the MNIST and EMNIST Datasets

Finally, we evaluated SSMP on the image completion task using the Binarized MNIST (Salakhutdinov and Murray 2008) and the EMNIST datasets (Cohen et al. 2017). As mentioned earlier, we used the top 30% of the image as evidence and estimated the bottom 70% by solving the MPE

task over PCs using various competing methods. The contingency tables for the MNIST and EMNIST datasets are shown in Table 2. We observe that on the MNIST dataset, SSMP is better than all competing schemes on 9 out of the 10 PCs, while it is inferior to all on one of them. On the EMNIST dataset, SSMP is better than all competing schemes on 7 out of the 10 PCs and inferior to all on one of the PCs. Detailed results on the image datasets, including qualitative comparisons, are provided in the supplement.

In summary, we find that, on average, our proposed method (SSMP) is better than other baseline MPE/MMAP approximations in terms of log-likelihood score. Moreover, it is substantially better than the baseline methods when the number of query variables is large. Also, once learned from data, it is also significantly faster than competing schemes.

### Conclusion and Future Work

In this paper, we introduced a novel self-supervised learning algorithm for solving MMAP queries in PCs. Our contributions comprise a neural network approximator and a self-supervised loss function which leverages the tractability of PCs for achieving scalability. Notably, our method employs minimal hyperparameters, requiring only one in the discrete case. We conducted a comprehensive empirical evaluation across various benchmarks; specifically, we experimented with 22 binary datasets used in tractable probabilistic models community and three classic image datasets, MNIST, EMNIST, and CIFAR-10. We compared our proposed neural approximator to polytime baseline techniques and observed that it is superior to the baseline methods in terms of log-likelihood scores and is significantly better in terms of computational efficiency. Additionally, we evaluated how our approach performs when used as an initialization scheme in stochastic hill climbing (local) search and found that it improves the quality of solutions output by anytime local search schemes. Our empirical results clearly demonstrated the efficacy of our approach in both accuracy and speed.

Future work includes compiling PCs to neural networks for answering more complex queries that involve constrained optimization; developing sophisticated self-supervised loss functions; learning better NN architecture for the given PC; generalizing our approach to arbitrarily chosen query and evidence subsets; etc.

### Acknowledgements

# References

Bekker, J.; Davis, J.; Choi, A.; Darwiche, A.; and Van den Broeck, G. 2015. Tractable Learning for Complex Probability Queries. In Cortes, C.; Lawrence, N.; Lee, D.; Sugiyama, M.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

Bioucas-Dias, J.; and Figueiredo, M. 2016. Bayesian Image Segmentation Using Hidden Fields: Supervised, Unsupervised, and Semi-Supervised Formulations. In *2016 24th European Signal Processing Conference (EUSIPCO)*, 523–527.

Choi, Y.; Friedman, T.; and Van den Broeck, G. 2022. Solving marginal map exactly by probabilistic circuit transformations. In *International Conference on Artificial Intelligence and Statistics*, 10196–10208. PMLR.

Choi, Y.; Vergari, A.; and Van den Broeck, G. 2020. Probabilistic Circuits: A Unifying Framework for Tractable Probabilistic Models. Technical report, University of California, Los Angeles.

Cohen, G.; Afshar, S.; Tapson, J.; and van Schaik, A. 2017. EMNIST: An Extension of MNIST to Handwritten Letters. arxiv:1702.05373.

Conaty, D.; de Campos, C. P.; and Mauá, D. D. 2017. Approximation Complexity of Maximum A Posteriori Inference in Sum-Product Networks. In Elidan, G.; Kersting, K.; and Ihler, A., eds., *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence (UAI)*. AUAI Press.

Cui, Z.; Wang, H.; Gao, T.; Talamadupula, K.; and Ji, Q. 2022. Variational Message Passing Neural Network for Maximum-A-Posteriori (MAP) Inference. In Cussens, J.; and Zhang, K., eds., *Uncertainty in Artificial Intelligence, Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI 2022, 1-5 August 2022, Eindhoven, the Netherlands*, volume 180 of *Proceedings of Machine Learning Research*, 464–474. PMLR.

Darwiche, A. 2003. A differential approach to inference in Bayesian networks. *Journal of the ACM (JACM)*, 50(3): 280–305.

de Campos, C. P. 2011. New Complexity Results for MAP in Bayesian Networks. *IJCAI International Joint Conference on Artificial Intelligence*, 2100–2106.

Dechter, R.; and Mateescu, R. 2007. AND/OR search spaces for graphical models. *Artificial intelligence*, 171(2-3): 73–106.

Donti, P. L.; Rolnick, D.; and Kolter, J. Z. 2020. DC3: A Learning Method for Optimization with Hard Constraints. In *International Conference on Learning Representations*.

Fioretto, F.; Mak, T. W. K.; and Hentenryck, P. V. 2020. Predicting AC Optimal Power Flows: Combining Deep Learning and Lagrangian Dual Methods. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01): 630–637.

Gogate, V.; and Dechter, R. 2012. Importance sampling-based estimation over AND/OR search spaces for graphical models. *Artificial Intelligence*, 184-185: 38–77.

Haaren, J. V.; and Davis, J. 2012. Markov Network Structure Learning: A Randomized Feature Generation Approach. *Proceedings of the AAAI Conference on Artificial Intelligence*, 26(1): 1148–1154.

Kingma, D. P.; and Ba, J. 2017. Adam: A Method for Stochastic Optimization. arxiv:1412.6980.

Kisa, D.; Van den Broeck, G.; Choi, A.; and Darwiche, A. 2014. Probabilistic sentential decision diagrams. In *Fourteenth International Conference on the Principles of Knowledge Representation and Reasoning*.

Kiselev, I.; and Poupart, P. 2014. POMDP Planning by Marginal-MAP Probabilistic Inference in Generative Models. In *Proceedings of the 2014 AAMAS Workshop on Adaptive Learning Agents*.

Krizhevsky, A.; Nair, V.; and Hinton, G. 2009. Learning multiple layers of features from tiny images. Technical report, Technical Report, University of Toronto.

Larochelle, H.; and Murray, I. 2011. The Neural Autoregressive Distribution Estimator. In Gordon, G.; Dunson, D.; and Dudík, M., eds., *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, 29–37. Fort Lauderdale, FL, USA: PMLR.

Lee, J.; Marinescu, R.; and Dechter, R. 2014. Applying Marginal MAP Search to Probabilistic Conformant Planning: Initial Results. In *Statistical Relational Artificial Intelligence, Papers from the 2014 AAAI Workshop, Québec City, Québec, Canada, July 27, 2014*, volume WS-14-13 of *AAAI Technical Report*. AAAI.

Li, K.; and Malik, J. 2016. Learning to Optimize. arXiv:1606.01885.

Loconte, L.; and Gala, G. 2022. DeeProb-kit: a Python Library for Deep Probabilistic Modelling.

Lowd, D.; and Davis, J. 2010. Learning Markov Network Structure with Decision Trees. In *2010 IEEE International Conference on Data Mining*, 334–343. IEEE. ISBN 978-1-4244-9131-5.

Mauá, D. D.; Reis, H. R.; Katague, G. P.; and Antonucci, A. 2020. Two reformulation approaches to maximum-a-posteriori inference in sum-product networks. In *International Conference on Probabilistic Graphical Models*, 293–304. PMLR.

Mei, J.; Jiang, Y.; and Tu, K. 2018. Maximum a posteriori inference in sum-product networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.

Park, J. D.; and Darwiche, A. 2004. Complexity Results and Approximation Strategies for MAP Explanations. *J. Artif. Int. Res.*, 21(1): 101–133.

Park, S.; and Hentenryck, P. V. 2023. Self-Supervised Primal-Dual Learning for Constrained Optimization. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(4): 4052–4060.

Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Peharz, R. 2015. *Foundations of sum-product networks for probabilistic modeling*. Ph.D. thesis, PhD thesis, Medical University of Graz.

Peharz, R.; Gens, R.; Pernkopf, F.; and Domingos, P. 2016. On the latent variable interpretation in sum-product networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(10): 2030–2044.

Ping, W.; Liu, Q.; and Ihler, A. T. 2015. Decomposition Bounds for Marginal MAP. In Cortes, C.; Lawrence, N.; Lee, D.; Sugiyama, M.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

Poon, H.; and Domingos, P. 2011. Sum-Product Networks: A New Deep Architecture. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, 337–346. AUAI Press.

Rahman, T.; and Gogate, V. 2016a. Learning Ensembles of Cutset Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1).

Rahman, T.; and Gogate, V. 2016b. Merging Strategies for Sum-Product Networks: From Trees to Graphs. In *Proceedings of the Thirty-Second Conference Conference on Uncertainty in Artificial Intelligence*, 617–626.

Rahman, T.; Jin, S.; and Gogate, V. 2019. Look ma, no latent variables: Accurate cutset networks via compilation. In *International Conference on Machine Learning*, 5311–5320. PMLR.

Rahman, T.; Kothalkar, P.; and Gogate, V. 2014. Cutset networks: A simple, tractable, and scalable approach for improving the accuracy of chow-liu trees. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2014, Nancy, France, September 15-19, 2014. Proceedings, Part II 14*, 630–645. Springer.

Salakhutdinov, R.; and Murray, I. 2008. On the quantitative analysis of deep belief networks. In *Proceedings of the 25th international conference on Machine learning*, 872–879. ACM.

Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15(56): 1929–1958.

Vergari, A.; Choi, Y.; Liu, A.; Teso, S.; and Van den Broeck, G. 2021. A Compositional Atlas of Tractable Circuit Operations for Probabilistic Inference. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 13189–13201. Curran Associates, Inc.

Yoon, K.; Liao, R.; Xiong, Y.; Zhang, L.; Fetaya, E.; Urtasun, R.; Zemel, R.; and Pitkow, X. 2019. Inference in probabilistic graphical models by graph neural networks. In *2019 53rd Asilomar Conference on Signals, Systems, and Computers*, 868–875. IEEE.

Zamzam, A. S.; and Baker, K. 2020. 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)*, 1–6.

# Deep Dependency Networks and Advanced Inference Schemes for Multi-Label Classification

**Shivvrat Arya**
The University of Texas at Dallas

**Yu Xiang**
The University of Texas at Dallas

**Vibhav Gogate**
The University of Texas at Dallas

## Abstract

We present a unified framework called deep dependency networks (DDNs) that combines dependency networks and deep learning architectures for multi-label classification, with a particular emphasis on image and video data. The primary advantage of dependency networks is their ease of training, in contrast to other probabilistic graphical models like Markov networks. In particular, when combined with deep learning architectures, they provide an intuitive, easy-to-use loss function for multi-label classification. A drawback of DDNs compared to Markov networks is their lack of advanced inference schemes, necessitating the use of Gibbs sampling. To address this challenge, we propose novel inference schemes based on local search and integer linear programming for computing the most likely assignment to the labels given observations. We evaluate our novel methods on three video datasets (Charades, TACoS, Wetlab) and three image datasets (MS-COCO, PASCAL VOC, NUS-WIDE), comparing their performance with (a) basic neural architectures and (b) neural architectures combined with Markov networks equipped with advanced inference and learning techniques. Our results demonstrate the superiority of our new DDN methods over the two competing approaches.

## 1 INTRODUCTION

In this paper, we focus on the multi-label classification (MLC) task and more specifically, on its two notable instantiations, multi-label action classification (MLAC) for videos and multi-label image classification (MLIC). At a high level, given a pre-defined set of labels (or actions) and a test example (video or image), the goal is to assign each test example to a subset of labels. It is well known that MLC is notoriously difficult because, in practice, the labels are often correlated, and thus, predicting them independently may lead to significant errors. Therefore, most advanced methods explicitly model the relationship or dependencies between the labels, using either probabilistic techniques (Antonucci et al., 2013; Di Mauro et al., 2016; Guo and Xue, 2013; Tan et al., 2015; Wang et al., 2008, 2014) or non-probabilistic/neural methods (Chen et al., 2019a,b; Kong et al., 2013; Liu et al., 2022, 2021; Nguyen et al., 2021; Papagiannopoulou et al., 2015; Qu et al., 2021; Wang et al., 2021a,b; Weng et al., 2023; Zhou et al., 2023).

To this end, motivated by approaches that combine probabilistic graphical models (PGMs) with neural networks (NNs) (Johnson et al., 2016; Krishnan et al., 2015), we jointly train a hybrid model, termed *deep dependency networks* (DDNs) (Guo and Weng, 2020), as illustrated in Figure 1. In a Deep Dependency Network (DDN), a conditional dependency network sits on top of a neural network. The underlying neural network transforms input data (e.g., an image) into a feature set. The dependency network (Heckerman et al., 2000) then utilizes these features to establish a local conditional distribution for each label, considering not only the features but also the other labels. Thus, at a high level, a DDN is a *neuro-symbolic model* where the neural network extracts features from data and the dependency network acts as a symbolic counterpart, learning the weighted constraints between the labels.

However, a limitation of DDNs is that they rely on naive techniques such as Gibbs sampling and mean-field inference for probabilistic reasoning and lack advanced probabilistic inference techniques (Lowd, 2012; Lowd and Shamaei, 2011). This paper addresses these limitations by introducing sophisticated inference schemes tailored for the Most Probable Explanation (MPE) task in DDNs, which involves finding the most likely assignment to the unobserved variables given observations. In essence, a solution to the MPE task, when applied to a probabilistic model defined over labels and observed variables, effectively solves the multi-label classification problem.
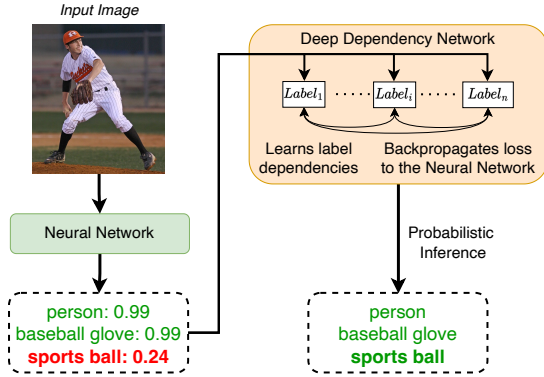
Figure 1: Illustrating improvements from our new inference schemes for DDNs. The DDN learns relationships between labels, and the inference schemes reason over them to accurately identify concealed objects, such as **sports ball**.

More specifically, we propose two new methods for MPE inference in DDNs. Our first method uses a random-walk based local search algorithm. Our second approach uses a piece-wise approximation of the log-sigmoid function to convert the non-linear MPE inference problem in DDNs into an integer linear programming problem. The latter can then be solved using off-the-shelf commercial solvers such as Gurobi (Gurobi Optimization, LLC, 2023).

We evaluate DDNs equipped with our new MPE inference schemes and trained via joint learning on three video datasets: Charades (Sigurdsson et al., 2016), TACoS (Regneri et al., 2013), and Wetlab (Naim et al., 2014), and three image datasets: MS-COCO (Lin et al., 2014), PASCAL VOC 2007 (Everingham et al., 2010), and NUS-WIDE (Chua et al., 2009). We compare their performance to two categories of models: (a) basic neural networks (without dependency networks) and (b) hybrids of Markov random fields (MRFs), an undirected PGM, and neural networks equipped with sophisticated reasoning and learning algorithms. Specifically, we employ three advanced approaches: (1) iterative join graph propagation (IJGP) (Mateescu et al., 2010), a type of generalized Belief propagation method (Yedidia et al., 2000) for marginal inference, (2) integer linear programming (ILP) based techniques for computing most probable explanations (MPE) and (3) a well-known structure learning method based on logistic regression with $\ell_1$-regularization (Lee et al., 2006; Wainwright et al., 2006) for pairwise MRFs.

Via a detailed experimental evaluation, we found that, generally speaking, the MRF+NN hybrids outperform NNs, as measured by metrics such as Jaccard index and subset accuracy, especially when advanced inference methods such as IJGP are employed. Additionally, DDNs, when equipped with our novel MILP-based MPE inference approach, of-

ten outperform both MRF+NN hybrids and NNs. This enhanced performance of DDNs with advanced MPE solvers is likely attributed to their superior capture of dense label interdependencies, a challenge for MRFs. Notably, MRFs rely on sparsity for efficient inference and learning.

## 2 PRELIMINARIES

A **log-linear model** or a **Markov random field** (MRF), denoted by $\mathcal{M}$, is an undirected probabilistic graphical model (Koller and Friedman, 2009) that is widely used in many real-world domains for representing and reasoning about uncertainty. MRFs are defined as a triple $\langle \mathbf{X}, \mathcal{F}, \Theta \rangle$ where $\mathbf{X} = \{X_1, \ldots, X_n\}$ is a set of Boolean random variables, $\mathcal{F} = \{f_1, \ldots, f_m\}$ is a set of features such that each feature $f_i$ (we assume that a feature is a Boolean formula) is defined over a subset $\mathbf{D}_i$ of $\mathbf{X}$, and $\Theta = (\theta_1, \ldots, \theta_m)$ are real-valued weights or parameters, namely $\forall \theta_i \in \Theta; \theta_i \in \mathbb{R}$ such that each feature $f_i$ is associated with a parameter $\theta_i$. $\mathcal{M}$ represents the following probability distribution:

$$P(\mathbf{x}) = \frac{1}{Z(\Theta)} \exp \left\{ \sum_{i=1}^{m} \theta_i f_i (\mathbf{x}_{\mathbf{D}_i}) \right\} \quad (1)$$

where $\mathbf{x}$ is an assignment of values to all variables in $\mathbf{X}$, $\mathbf{x}_{\mathbf{D}_i}$ is the projection of $\mathbf{x}$ on the variables $\mathbf{D}_i$ of $f_i$, $f_i(\mathbf{x}_{\mathbf{D}_i})$ is an *indicator function* that equals 1 when $\mathbf{x}_{\mathbf{D}_i}$ evaluates $f_i$ to True and is 0 otherwise, and $Z(\Theta)$ is the normalization constant called the *partition function*.

We focus on three tasks over MRFs: (1) structure learning; (2) posterior marginal inference; and (3) finding the most likely assignment to all the non-evidence variables given evidence (this task is often called most probable explanation or MPE inference in short). All of these tasks are at least NP-hard in general, and therefore approximate methods are often preferred over exact ones in practice.

A popular and fast method for structure learning is to learn binary pairwise MRFs (MRFs in which each feature is defined over at most two variables) by training an $\ell_1$-regularized logistic regression classifier for each variable given all other variables as features (Lee et al., 2006; Wainwright et al., 2006). $\ell_1$-regularization induces sparsity in that it encourages many weights to take the value zero. All non-zero weights are then converted into conjunctive features. Each conjunctive feature evaluates to True if both variables are assigned the value 1 and to False otherwise. Popular approaches for posterior marginal inference are the Gibbs sampling algorithm and generalized Belief propagation (Yedidia et al., 2000) techniques such as Iterative Join Graph Propagation (Mateescu et al., 2010). For MPE inference, a popular approach is to encode it as an integer linear programming (ILP) problem (Koller and Friedman, 2009) and then use off-the-shelf approaches such as Gurobi Optimization, LLC (2023) to solve the ILP.

**Dependency Networks (DNs)** (Heckerman et al., 2000) represent the joint distribution using a set of local conditional probability distributions, one for each variable. Each conditional distribution defines the probability of a variable given all of the others. A DN is consistent if there exists a joint probability distribution $P(\mathbf{x})$ such that all conditional distributions $P_i(x_i|\mathbf{x}_{-i})$ where $\mathbf{x}_{-i}$ is the projection of $\mathbf{x}$ on $\mathbf{X} \setminus \{X_i\}$, are conditional distributions of $P(\mathbf{x})$.

A DN is learned from data by learning a classifier (e.g., logistic regression, multi-layer perceptron, etc.) for each variable, and thus DN learning is embarrassingly parallel. However, because the classifiers are independently learned from data, we often get an inconsistent DN. It has been conjectured (Heckerman et al., 2000) that most DNs learned from data are almost consistent in that only a few parameters need to be changed in order to make them consistent.

The most popular inference method over DNs is *fixed-order* Gibbs sampling (Liu, 2008). If the DN is consistent, then its conditional distributions are derived from a joint distribution $P(\mathbf{x})$, and the stationary distribution (namely the distribution that Gibbs sampling converges to) will be the same as $P(\mathbf{x})$. If the DN is inconsistent, then the stationary distribution of Gibbs sampling will be inconsistent with the conditional distributions.

## 3 TRAINING DEEP DEPENDENCY NETWORKS

In this section, we detail the training process for our proposed Deep Dependency Network model (Guo and Weng, 2020), the hybrid framework for multi-label action classification in videos and multi-label image classification. Two key components are trained jointly: a neural network and a conditional dependency network. The neural network is responsible for extracting high-quality features from video segments or images, while the dependency network models the relationships between these features and their corresponding labels, supplying the neural network with gradient information regarding these relationships.

### 3.1 Framework

Let $\mathbf{V}$ denote the set of random variables corresponding to the pixels and $\mathbf{v}$ denote the RGB values of the pixels in a frame or a video segment. Let $\mathbf{E}$ denote the (continuous) output nodes of a neural network which represents a function $\mathcal{N} : \mathbf{v} \mapsto \mathbf{e}$, that takes $\mathbf{v}$ as input and outputs an assignment $\mathbf{e}$ to $\mathbf{E}$. Let $\mathbf{X} = \{X_1, \ldots, X_n\}$ denote the set of indicator variables representing the labels. For simplicity, we assume that $|\mathbf{E}| = |\mathbf{X}| = n$. Given $(\mathbf{V}, \mathbf{E}, \mathbf{X})$, a deep dependency network (DDN) is a pair $\langle \mathcal{N}, \mathcal{D} \rangle$ where $\mathcal{N}$ is a neural network that maps $\mathbf{V} = \mathbf{v}$ to $\mathbf{E} = \mathbf{e}$ and $\mathcal{D}$ is a conditional dependency network (Guo and Gu, 2011) that models $P(\mathbf{x}|\mathbf{e})$ where $\mathbf{e} = \mathcal{N}(\mathbf{v})$.



Figure 2: Illustration of Dependency Network for multi-label video classification. The NN takes video clips (frames) as input and outputs the features $e_1, e_2, ..., e_n$ (denoted by red colored nodes). These features are then used by the sigmoid output ($\sigma_1, \ldots, \sigma_n$) of the dependency layer to model the local conditional distributions.

The conditional dependency network represents the distribution $P(\mathbf{x}|\mathbf{e})$ using a collection of local conditional distributions $P_i(x_i|\mathbf{x}_{-i}, \mathbf{e})$, one for each label $X_i$, where $\mathbf{x}_{-i} = \{x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n\}$.

Thus, a DDN is a discriminative model and represents the conditional distribution $P(\mathbf{x}|\mathbf{v})$ using several local conditional distributions $P(x_i|\mathbf{x}_{-i}, \mathbf{e})$ and makes the following conditional independence assumptions $P(x_i|\mathbf{x}_{-i}, \mathbf{v}) = P(x_i|\mathbf{x}_{-i}, \mathbf{e})$ where $\mathbf{e} = \mathcal{N}(\mathbf{v})$. Figure 2 demonstrates the DDN architecture.

### 3.2 Learning

We employ the conditional pseudo log-likelihood loss (CPLL) (Besag, 1975) in order to jointly train the two components of DDN, drawing inspiration from the DDN training approach outlined by Guo and Weng (2020). For each training example $(\mathbf{v}, \mathbf{x})$, we send the video/image through the neural network to obtain a new representation $\mathbf{e}$ of $\mathbf{v}$. In the dependency layer, we learn a classifier for each label $X_i$ to model the conditional distribution $P_i(x_i|\mathbf{x}_{-i}, \mathbf{e})$. More precisely, with the representation of the training instance $(\mathbf{e}, \mathbf{x})$, each sigmoid output of the dependency layer indexed by $i$ and denoted by $\sigma_i$ (see figure 2) uses $X_i$ as the class variable and $(\mathbf{E} \cup \mathbf{X}_{-i})$ as the attributes. The joint training of the model is achieved through CPLL applied to the outputs of the dependency layer ($\sigma_i$'s).

Let $\Theta$ represent the parameter set of the DDN. We employ gradient-based optimization methods (e.g., backpropagation), to minimize the Conditional Pseudo-Likelihood

(CPLL) loss function given below

$$\mathcal{L}(\Theta, \mathbf{v}, \mathbf{x}) = -\sum_{i=1}^{n} \log P_i(x_i | \mathbf{e} = \mathcal{N}(\mathbf{v}), \mathbf{x}_{-i}; \Theta) \quad (2)$$

## 4 MPE INFERENCE IN DDNs

Unlike a conventional discriminative model, such as a neural network, in a DDN, we cannot predict the output labels by simply making a forward pass over the network. This is because each sigmoid output $\sigma_i$ (which yields a probability distribution over $X_i$) of the dependency layer requires an assignment $\mathbf{x}_{-i}$ to all labels except $x_i$ and $\mathbf{x}_{-i}$ is not available at prediction time. Consequently, it is imperative to employ specialized techniques for obtaining output labels in multilabel classification tasks within the context of DDNs.

Given a DDN representing a distribution $P(\mathbf{x}|\mathbf{e})$ where $\mathbf{e} = \mathcal{N}(\mathbf{v})$, the multilabel classification task can be solved by finding the most likely assignment to all the unobserved variables based on a set of observed variables. This task is also called the most probable explanation (MPE) task. Formally, we seek to find $\mathbf{x}^*$ such that:

$$\mathbf{x}^* = \arg\max_{\mathbf{x}} P(\mathbf{x} \mid \mathbf{e}) \quad (3)$$

Next, we present three algorithms for solving the MPE task.

### 4.1 Gibbs Sampling

To perform MPE inference, we first send the video (or frame) through the neural network to yield an assignment $\mathbf{e}$ to all variables in $\mathbf{E}$. Then given $\mathbf{e}$, we generate $N$ samples $(\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)})$ via Gibbs sampling (see for example Koller and Friedman (2009)), a classic MCMC technique. These samples can then be used to estimate the marginal probability distribution of each label $X_i$ using the following mixture estimator (Liu, 2008):

$$\hat{P}_i(x_i|\mathbf{v}) = \frac{1}{N}\sum_{j=1}^{N} P_i\left(x_i \mid \mathbf{x}_{-i}^{(j)}, \mathbf{e}\right) \quad (4)$$

Given an estimate of the marginal distribution $\hat{P}_i(x_i|\mathbf{v})$ for each variable $X_i$, we can estimate the MPE assignment using the standard max-marginal approximation:

$$\max_{\mathbf{x}} P(\mathbf{x}|\mathbf{e}) \approx \prod_{i=1}^{n} \max_{x_i} \hat{P}(x_i|\mathbf{e})$$

In other words, we can construct an approximate MPE assignment by finding the value $x_i$ for each variable that maximizes $\hat{P}_i(x_i|\mathbf{e})$.

### 4.2 Local Search Based Methods

Local search algorithms (see for example Selman et al. (1993)) systematically examine the solution space by making localized adjustments, with the objective of either finding an optimal solution or exhausting a pre-established time limit. They offer a viable approach for solving the MPE inference task in DDNs. These algorithms, through their structured exploration and score maximization, are effective in identifying near-optimal label configurations.

In addressing the MPE inference within DDNs, we define the objective function for local search as $\sum_{i=1}^{n} \log\left(P_i(x_i \mid \mathbf{x}_{-i}, \mathbf{e})\right)$, where $\mathbf{e}$ is the evidence provided by $\mathcal{N}$. We propose to use two distinct local search strategies for computing the MPE assignment: random walk and greedy local search.

In *random walk (RW)*, the algorithm begins with a random assignment to the labels and at each iteration, flips a random label (changes the value of the label from a 1 to a 0 or a 0 to a 1) to yield a new assignment. At termination, the algorithm returns the assignment with the highest score explored during the random walk. In *greedy local search*, the algorithm begins with a random assignment to the labels, and at each iteration, with probability $p$ flips a random label to yield a new assignment and with probability $1 - p$ flips a label that yields the maximum improvement in the score. At termination, the algorithm returns the assignment with the highest score explored during its execution.

### 4.3 Multi-Linear Integer Programming

In this section, we present a novel approach for MPE inference in DDNs by formulating the problem as a Second-Order Multi-Linear Integer Programming task. Specifically, we show that the task of maximizing the scoring function $\sum_{i=1}^{n} \log\left(P_i(x_i \mid \mathbf{x}_{-i}, \mathbf{e})\right)$ is equivalent to the following optimization problem (derivation is provided in the supplementary material):

$$\underset{\mathbf{x}, \mathbf{z}}{\text{maximize}} \sum_{i=1}^{n} \left(x_i \log P_i + (1 - x_i) \log(1 - P_i)\right)$$

subject to:

$$P_i = \sigma(z_i), \quad \forall i \in \{1, \ldots, n\}$$

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \ldots, n\}$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \ldots, n\}$$

$$(5)$$

Here, $P_i = P_i(x_i|\mathbf{x}_{-i}, \mathbf{e})$ and $w_{ij}$'s and $v_{ik}$'s denote the weights associated with $\mathbf{e}$ and $\mathbf{x}$ for $P_i$, respectively. The bias term for each $P_i$ is denoted by $b_i$. In this context, $z_i$ represents the values acquired prior to applying the sigmoid activation function. Substituting the constraint

$P_i(x_i|\mathbf{x}_{-i}, \mathbf{e}) = \sigma(z_i) = \frac{1}{1+e^{(-z_i)}}$ in the objective and simplifying, we get

$$\underset{\mathbf{x},\mathbf{z}}{\text{maximize}} \sum_{i=1}^{n} x_i z_i - \log(1 + e^{z_i}) \quad (6)$$

subject to:

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \dots, n\} \quad (7)$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\} \quad (8)$$

The optimal value for $\mathbf{x}$ corresponds to the solution of this optimization problem. The second term in the objective specified in equation 6 comprises logarithmic and exponential functions, which are non-linear, thereby making it a *non-linear optimization problem*. To address this non-linearity, we propose the use of piece-wise linear approximations(Asghari et al., 2022; Geißler et al., 2012; Kumar, 2007; Li et al., 2022; Lin et al., 2013; Rovatti et al., 2014) for these terms. Further details about the piece-wise linear approximation can be found in the supplement. Let $g(z_i)$ represent the piece-wise linear approximation of $\log(1 + e^{z_i})$, then the optimization problem can be expressed as:

$$\underset{\mathbf{x},\mathbf{z}}{\text{maximize}} \sum_{i=1}^{n} x_i z_i - g(z_i)$$

subject to:

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \dots, n\} \quad (9)$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\}$$

The optimization problem given in equation 9 is an integer multilinear program of order 2 because it includes terms of the form $x_i x_j$ where both $x_i$ and $x_j$ take values from the set $\{0, 1\}$. Since $x_i x_j$ corresponds to a "logical and" between two Boolean variables, all such expressions can be easily encoded as linear constraints yielding an integer linear program (ILP) (see the supplementary material for an example).

In our experiments, we solved the ILP using Gurobi (Gurobi Optimization, LLC, 2023), which is an anytime ILP solver. Another useful feature of the ILP formulation is that we can easily incorporate prior knowledge (e.g., an image may not have more than ten objects/labels) into the ILP under the assumption that the knowledge can be reliably modeled using linear constraints.

## 5 EXPERIMENTAL EVALUATION

In this section, we evaluate the proposed methods on two multi-label classification tasks: (1) multi-label activity classification using three video datasets; and (2) multi-label image classification using three image datasets. We begin by describing the datasets and metrics, followed by the experimental setup, and conclude with the results. All models were implemented utilizing PyTorch and were trained and evaluated on a machine with an NVIDIA A40 GPU and an Intel(R) Xeon(R) Silver 4314 CPU.

### 5.1 Datasets and Metrics

We evaluated our algorithms on three video datasets: (1) Charades (Sigurdsson et al., 2016); (2) TACoS (Regneri et al., 2013); and (3) Wetlab (Naim et al., 2015). Charades dataset comprises videos of people performing daily indoor activities while interacting with various objects. We adopted the train-test split instructions given in PySlowFast (Fan et al., 2020) with 7,986 training and 1,863 validation videos. TACoS consists of third-person videos of a person cooking in a kitchen. The dataset comes with hand-annotated labels of actions, objects, and locations for each video frame. From the complete set of these labels, we selected 28 labels resulting in 60,313 training frames and 9,355 test frames across 17 videos. Wetlab features experimental activities in labs, consisting of five training videos (100,054 frames) and one test video (11,743 frames) with 57 distinct labels.

For multi-label image classification (MLIC), we examined: (1) MS-COCO (Lin et al., 2014); (2) PASCAL VOC 2007 (Everingham et al., 2010); and (3) NUS-WIDE (Chua et al., 2009). MS-COCO, a well-known dataset for detection and segmentation, comprises 122,218 labeled images with an average of 2.9 labels per image. We used its 2014 version. NUS-WIDE is a real-world web image dataset that contains 269,648 images from Flickr. Each image has been manually annotated with a subset of 81 visual classes that include objects and scenes. PASCAL VOC 2007 contains 5,011 train-validation and 4,952 test images, with each image labeled with one or more of the 20 available object classes.

We follow the instructions provided in (Qu et al., 2021) to do the train-test split for NUS-WIDE and PASCAL VOC. We evaluated the performance on the TACoS, Wetlab, MS-COCO, NUS-WIDE, and VOC datasets using Subset Accuracy (SA), Jaccard Index (JI), Hamming Loss (HL), Macro F1 Score (Macro F1), Micro F1 Score (Micro F1) and F1 Score (F1). With the exception of Hamming Loss, superior performance is indicated by higher scores in all considered metrics. We omit SA for the Charades dataset due to the infeasibility of achieving reasonable scores given its large label space. Additionally, Hamming Loss (HL) has been excluded for the MS-COCO dataset as the perfor-

mance of all methods is indistinguishable.

Given that the focus in MPE inference is on identifying the most probable label configurations rather than individual label probabilities, the use of Mean Average Precision (mAP) as a performance metric is not applicable to our study. Therefore, our primary analysis relies on SA, JI, HL, Macro F1, Micro F1, and F1. Precision-based metrics are detailed in the supplementary material.

## 5.2 Experimental Setup and Methods

We used three types of architectures in our experiments: (1) Baseline neural networks, which are specific to each dataset; (2) neural networks augmented with MRFs, which we will refer to as deep random fields or DRFs in short; and (3) a dependency network on top of the neural networks called deep dependency networks (DDNs).

**Neural Networks**. We choose four different types of neural networks, and they act as a baseline for the experiments and as a feature extractor for DRFs and DDNs. Specifically, we experimented with: (1) 2D CNNs, (2) 3D CNNs, (3) transformers, and (4) CNNs with attention module and graph attention networks (GAT) (Veličković et al., 2018). This helps us show that our proposed method can improve the performance of a wide variety of neural architectures, even those which model label relationships, because unlike the latter, it performs probabilistic inference.

For the Charades dataset, we use the PySlowFast (Fan et al., 2020) implementation of the SlowFast Network (Feichtenhofer et al., 2019) (a state-of-the-art 3D CNN for video classification). For TACoS and Wetlab datasets, we use InceptionV3 (Szegedy et al., 2016), a state-of-the-art 2D CNN model for image classification. For the MS-COCO dataset, we used Query2Label (Q2L) (Liu et al., 2021), which uses transformers to pool class-related features. Q2L also learns label embeddings from data to capture the relationships between the labels. Finally, we used the multi-layered semantic representation network (MSRN) (Qu et al., 2021) for NUS-WIDE and PASCAL VOC. MSRN also models label correlations and learns semantic representations at multiple convolutional layers. We adopt pre-trained models and hyper-parameters from existing repositories for Charades, MS-COCO, NUS-WIDE, and PASCAL VOC. For TaCoS and Wetlab datasets, we fine-tuned an InceptionV3 model that was pre-trained on the ImageNet dataset.

**Deep Random Fields (DRFs)**. As a baseline, we used a model that combines MRFs with neural networks. This DRF model is similar to DDN except that we use an MRF instead of a DN to compute $P(\mathbf{x}|\mathbf{e})$. We trained the MRFs generatively; namely, we learned a joint distribution $P(\mathbf{x}, \mathbf{e})$, which can be used to compute $P(\mathbf{x}|\mathbf{e})$ by instantiating evidence. We chose generative learning be-

cause we learned the structure of the MRFs from data, and discriminative structure learning is slow in practice (Koller and Friedman, 2009).

For inference over MRFs, we used Gibbs sampling (GS), Iterative Join Graph Propagation (IJGP) (Mateescu et al., 2010), and Integer Linear Programming (ILP) methods. Thus, three versions of DRFs corresponding to the inference scheme were used. We refer to these schemes as DRF-GS, DRF-ILP, and DRF-IJGP, respectively. Note that IJGP and ILP are advanced schemes, and we are unaware of their use for multi-label classification. Our goal is to test whether advanced inference schemes help improve the performance of deep random fields.

**Deep Dependency Networks (DDNs)**. We trained the Deep Dependency Networks (DDNs) using the joint learning loss described in equation 2. We examined four unique inference methods for DDNs: (1) DDN-GS, employing Gibbs Sampling; (2) DDN-RW, leveraging a random walk local search; (3) DDN-Greedy, implementing a greedy local search; and (4) DDN-ILP, utilizing Integer Linear Programming to optimize the objective given in equation 9. It is noteworthy that, until now, only DDN-GS has been used for inference in Dependency Networks. DDN-RW and DDN-Greedy, which are general-purpose local search techniques, are our proposals for MPE inference on Dependency Networks. Lastly, DDN-ILP introduces a novel approach using optimization techniques with the objective of enhancing MPE inference on dependency networks.

**Hyperparameters.** For DRFs, in order to learn a sparse structure (using the logistic regression with $\ell_1$ regularization method of Wainwright et al. (2006)), we increased the regularization constant associated with the $\ell_1$ regularization term until the number of neighbors of each node in $\mathcal{G}$ is bounded between 2 and 10. We enforced this sparsity constraint in order to ensure that the inference schemes, specifically IJGP and ILP, are accurate and the model does not overfit the training data. IJGP, ILP, and GS are anytime methods; for them, we used a time-bound of one minute per example.

For DDNs, we used $\ell_1$ regularization for the dependency layer. Through cross-validation, we selected the regularization constants from the set $\{0.1, 0.01, 0.001\}$. In the context of joint learning, the learning rates of the DDN model were adjusted using an extended version of the learning rate scheduler from PySlowFast (Fan et al., 2020). We impose a strict time constraint of 60 seconds per example for DDN-GS, DDN-RW, DDN-Greedy, and DDN-ILP. The DDN-ILP method is solved using Gurobi Optimization, LLC (2023) and utilizes accurate piece-wise linear approximations for non-linear functions, subject to a predefined error tolerance of $0.001$.

Table 1: Comparison of our methods with the feature extractor for MLAC. The best/second best values are bold/underlined.

| Method | Charades | | | | | TACoS | | | | | | Wetlab | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | JI | HL↓ | Macro F1 | Micro F1 | F1 | SA | JI | HL↓ | Macro F1 | Micro F1 | F1 | SA | JI | HL↓ | Macro F1 | Micro F1 | F1 |
| SlowFast | 0.29 | **0.052** | 0.32 | <u>0.45</u> | <u>0.42</u> | | | | | | | | | | | | |
| InceptionV3 | | | | | | 0.40 | 0.61 | 0.082 | 0.26 | 0.48 | 0.44 | 0.35 | 0.64 | 0.017 | 0.24 | 0.76 | 0.72 |
| DRF - GS | 0.22 | 0.194 | 0.16 | 0.31 | 0.28 | 0.47 | 0.65 | 0.044 | 0.54 | 0.75 | 0.70 | 0.35 | 0.52 | 0.058 | 0.22 | 0.69 | 0.68 |
| DRF - ILP | 0.31 | 0.067 | 0.18 | 0.21 | 0.22 | 0.51 | 0.65 | **0.030** | 0.42 | 0.73 | 0.66 | 0.60 | 0.73 | 0.014 | <u>0.27</u> | <u>0.81</u> | 0.77 |
| DRF - IJGP | <u>0.32</u> | <u>0.054</u> | 0.19 | 0.29 | 0.26 | 0.44 | <u>0.70</u> | 0.043 | 0.54 | **0.83** | **0.78** | 0.58 | 0.74 | <u>0.014</u> | 0.24 | **0.82** | 0.77 |
| DDN - GS | 0.31 | <u>0.054</u> | 0.30 | 0.43 | 0.41 | 0.54 | 0.69 | 0.041 | 0.52 | 0.74 | 0.68 | <u>0.63</u> | <u>0.78</u> | **0.011** | **0.28** | 0.80 | <u>0.79</u> |
| DDN - RW | 0.30 | 0.056 | 0.13 | 0.24 | 0.22 | 0.46 | 0.64 | 0.042 | 0.59 | 0.77 | 0.75 | 0.46 | 0.63 | 0.025 | 0.22 | 0.66 | 0.68 |
| DDN - Greedy | 0.31 | 0.056 | <u>0.33</u> | 0.44 | 0.41 | <u>0.56</u> | 0.69 | <u>0.040</u> | <u>0.61</u> | <u>0.80</u> | 0.75 | 0.55 | 0.68 | <u>0.014</u> | **0.28** | 0.80 | <u>0.79</u> |
| DDN - ILP | **0.33** | **0.052** | **0.36** | **0.47** | **0.44** | **0.63** | **0.72** | <u>0.040</u> | **0.62** | 0.79 | <u>0.76</u> | **0.65** | **0.76** | <u>0.014</u> | **0.28** | <u>0.81</u> | **0.80** |

Table 2: Comparison of our methods with the feature extractor for MLIC. The best/second best values are bold/underlined.

| Method | MS-COCO | | | | | NUS-WIDE | | | | | | PASCAL-VOC | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SA | JI | Macro F1 | Micro F1 | F1 | SA | JI | HL↓ | Macro F1 | Micro F1 | F1 | SA | JI | HL↓ | Macro F1 | Micro F1 | F1 |
| Q2L | 0.51 | 0.80 | **0.86** | <u>0.86</u> | **0.88** | | | | | | | | | | | | |
| MSRN | | | | | | 0.31 | <u>0.64</u> | **0.015** | **0.56** | <u>0.73</u> | **0.71** | 0.71 | 0.85 | 0.015 | 0.89 | 0.90 | 0.91 |
| DRF - GS | 0.35 | 0.69 | 0.81 | 0.82 | 0.79 | 0.28 | 0.55 | 0.020 | 0.23 | 0.63 | 0.62 | 0.73 | 0.83 | 0.021 | 0.75 | 0.85 | 0.83 |
| DRF - ILP | <u>0.54</u> | <u>0.82</u> | 0.84 | 0.83 | 0.82 | <u>0.32</u> | 0.59 | <u>0.016</u> | 0.28 | 0.70 | 0.67 | 0.76 | 0.88 | 0.019 | 0.81 | 0.87 | 0.87 |
| DRF - IJGP | **0.55** | <u>0.82</u> | <u>0.85</u> | 0.85 | 0.85 | <u>0.32</u> | 0.63 | 0.017 | 0.32 | 0.71 | <u>0.69</u> | 0.76 | 0.87 | 0.019 | 0.77 | 0.86 | 0.86 |
| DDN - GS | **0.55** | <u>0.82</u> | **0.86** | **0.87** | **0.88** | **0.33** | 0.62 | <u>0.016</u> | <u>0.52</u> | 0.71 | <u>0.69</u> | 0.76 | 0.87 | 0.008 | 0.94 | 0.94 | <u>0.96</u> |
| DDN - RW | 0.53 | 0.81 | 0.82 | 0.85 | 0.86 | 0.25 | 0.56 | <u>0.016</u> | 0.26 | 0.70 | 0.68 | 0.82 | 0.89 | **0.006** | 0.94 | 0.94 | 0.93 |
| DDN - Greedy | **0.55** | <u>0.82</u> | <u>0.85</u> | <u>0.86</u> | <u>0.87</u> | 0.30 | 0.61 | <u>0.016</u> | 0.29 | 0.72 | <u>0.69</u> | <u>0.86</u> | <u>0.91</u> | <u>0.007</u> | <u>0.95</u> | <u>0.95</u> | <u>0.96</u> |
| DDN - ILP | **0.55** | **0.83** | <u>0.85</u> | <u>0.86</u> | **0.88** | **0.33** | **0.65** | **0.015** | 0.51 | **0.74** | **0.71** | **0.89** | **0.95** | **0.006** | **0.96** | **0.96** | **0.97** |

## 5.3 Results

We compare the baseline neural networks with three versions of DRFs and four versions of DDNs using the six metrics and six datasets given in Section 5.1. The results are presented in tables 1 and 2.

**Comparison between baseline neural network and DRFs**. We observe that IJGP and ILP outperform the baseline neural networks (which include transformers for some datasets) in terms of JI, SA, and HL on four out of the six datasets. IJGP typically outperforms GS and ILP on JI. In terms of F1 metrics, IJGP and ILP methods tend to perform better than Gibbs Sampling (GS) approaches, but they are less effective than baseline NNs. ILP's superiority in SA—a metric that scores 1 for an exact label match and 0 otherwise—can be attributed to its precise most probable explanation (MPE) inference. An accurate MPE inference, when paired with a precise model, is likely to achieve high SA scores. Note that getting a higher SA is much harder in datasets having large number of labels. Specifically, SA does not distinguish between models that predict *almost* correct labels and completely incorrect outputs. We observe that advanced inference schemes, particularly IJGP and ILP, are superior on average to GS.

**Comparison between baseline neural networks and DDNs**. Our study has shown that the DDN model with the proposed MILP based inference method is superior to the baseline neural networks in four out of six datasets, with notable enhancements in SA and JI metrics, such as an 18%, 23%, and 30% improvement in SA for the PASCAL-VOC, TACoS, and Wetlab datasets, respectively. Moreover, the MILP-based method either significantly surpasses or matches all alternative inference strategies for DDNs.

While Gibbs Sampling-based inference is typically better than the Random Walk based approach in DDNs, its performance against the greedy sampling method varies.

We observe that on the MLIC task, the DDN with advanced MPE inference outperforms Q2L and MSRN, even though both Q2L and MSRN model label correlations. This suggests that DDNs are either able to uncover additional relationships between labels during the learning phase or better reason about them during the inference phase or both. In particular, both Q2L and MSRN do not use MPE inference to predict the labels because they do not explicitly model the joint probability distribution over the labels.

Figure 3 displays the images and their predicted labels by both Q2L and DDN-ILP on the MS-COCO dataset. Our method not only adds the labels omitted by Q2L but also eliminates several incorrect predictions. In the first two images, our approach rectifies label omissions by Q2L, conforming to the ground truth. In the third image, our method removes erroneous predictions. The last image illustrates a case where DDN underperforms compared to Q2L by failing to identify a ground-truth label. Additional instances are available in the appendix.

**Comparison between DRFs and DDNs**. Based on our observations, we found that jointly trained DDNs in conjunction with the proposed inference method consistently lead to superior performance compared to the top-performing DRFs across all datasets. Nonetheless, in certain situations, DRFs that employ advanced inference strategies produce results that closely match those of DDNs for both JI and SA, making DRFs a viable option, particularly when limited GPU resources are available for training and optimization of both JI and SA is prioritized.

| Image |  |  |  |  |
|---|---|---|---|---|
| Ground Truth | person, sports ball, tennis racket, chair, clock | chair, potted plant, tv, remote, book, vase | person, bottle, pizza, clock | bird, skateboard, couch |
| Q2L | person (1.00), tennis racket (1.00), chair (0.96), clock (0.95), **[sports ball (0.33)]** | chair (0.99), potted plant (0.99), tv (1.00), book (1.00), vase (0.96), [**remote (0.25)**] | person (1.00), bottle (0.99), pizza (1.00), **potted plant (0.74)**, clock (0.77), **vase (0.80),** | bird (1.00), skateboard (1.00), couch (0.73) |
| DDN | person, sports ball, tennis racket, chair, clock | chair, potted plant, tv, remote, book, vase | person, bottle, pizza, clock | bird, skateboard |

Figure 3: Comparison of labels predicted by Q2L (Liu et al., 2021) and our DDN-ILP scheme on the MS-COCO dataset. Labels in bold represent the difference between the predictions of the two methods, assuming that a threshold of 0.5 is used (i.e., every label whose probability > 0.5 is considered a predicted label). Due to the MPE focus in DDN-ILP, only label configurations are generated, omitting corresponding probabilities. The first three column shows examples where DDN improves over Q2L, while the last column (outlined in red) shows an example where DDN is worse than Q2L.

In summary, the empirical results indicate that Deep Dependency Networks (DDNs) equipped with advanced inference strategies consistently outperform conventional neural networks and MRF+NN hybrids in multi-label classification tasks. Furthermore, these advanced inference methods surpass traditional sampling-based techniques for DDNs. The superior performance of both DDNs and DRFs utilizing advanced inference techniques supports the value of such mechanisms and suggests that further advancement in this area has the potential to unlock additional capabilities within these models.

## 6 RELATED WORK

A large number of methods have been proposed that train PGMs and NNs jointly. For example, Zheng et al. (2015) proposed to combine conditional random fields (CRFs) and recurrent neural networks (RNNs), Arnab et al. (2016); Larsson et al. (2017, 2018); Schwing and Urtasun (2015) showed how to combine CNNs and CRFs, Chen et al. (2015) proposed to use densely connected graphical models with CNNs, and Johnson et al. (2016) combined latent graphical models with neural networks. The combination of PGMs and NNs has also been applied to improve performance on a wide variety of real-world tasks. Notable examples include human pose estimation (Liang et al., 2018; Song et al., 2017; Tompson et al., 2014; Yang et al., 2016), semantic labeling of body parts (Kirillov et al., 2016), stereo estimation (Knöbelreiter et al., 2017), language understanding (Yao et al., 2014), face sketch synthesis (Zhu et al., 2021) and crowd-sourcing aggregation

(Li et al., 2021)). These hybrid models have also been used for solving a range of computer vision tasks such as semantic segmentation (Arnab et al., 2018; Guo and Dou, 2021), image crowd counting (Han et al., 2017), visual relationship detection (Yu et al., 2022), modeling for epileptic seizure detection (Craley et al., 2019), face sketch synthesis (Zhang et al., 2020), semantic image segmentation (Chen et al., 2018; Lin et al., 2016), 2D Hand-pose Estimation (Kong et al., 2019), depth estimation from a single monocular image (Liu et al., 2015), animal pose tracking (Wu et al., 2020) and pose estimation (Chen and Yuille, 2014).

To date, dependency networks have been used to solve various tasks such as collective classification (Neville and Jensen, 2003), binary classification (Gámez et al., 2006, 2008), multi-label classification (Guo and Gu, 2011), part-of-speech tagging (Tarantola and Blanc, 2002), relation prediction (Figueiredo et al., 2021), text classification (Guo and Weng, 2020) and collaborative filtering (Heckerman et al., 2000). However, DDNs have traditionally been restricted to Gibbs sampling (Heckerman et al., 2000) and mean-field inference (Lowd and Shamaei, 2011), showing limited compatibility with advanced probabilistic inference methods (Lowd, 2012). This study marks the inaugural attempt to incorporate advanced inference methods for the MPE task in DDNs, utilizing jointly trained networks for MLC scenarios.

# 7  CONCLUSION AND FUTURE WORK

More and more state-of-the-art methods for challenging applications of computer vision tasks usually use deep neural networks. Deep neural networks are good at extracting features in vision tasks like image classification, video classification, object detection, image segmentation, and others. Nevertheless, for more complex tasks involving multi-label classification, these methods cannot model crucial information like inter-label dependencies and infer about them. In this work, we present novel inference algorithms for Deep Dependency Networks (DDNs), a powerful neuro-symbolic approach, that exhibit consistent superiority compared to traditional neural network baselines, sometimes by a substantial margin. These algorithms offer an improvement over existing Gibbs Sampling-based inference schemes used for DDNs, without incurring significant computational burden. Importantly, they have the capacity to infer inter-label dependencies that are commonly overlooked by baseline techniques utilizing transformers, attention modules, and Graph Attention Networks (GAT). By formulating the inference procedure as an optimization problem, our approach permits the integration of domain-specific constraints, resulting in a more knowledgeable and focused inference process. In particular, our optimization-based approach furnishes a robust and computationally efficient mechanism for inference, well-suited for handling intricate multi-label classification tasks.

Avenues for future work include: applying the setup described in the paper to other multi-label classification tasks in computer vision, natural language understanding, and speech recognition; converting DDNs to MRFs for better inference (Lowd, 2012); exploring and validating the neuro-symbolic benefits of DDNs such as improved accuracy in predictions and enhanced interpretability of model decisions; etc.

# ACKNOWLEDGMENTS

## References

A. Antonucci, G. Corani, D. D. Mauá, and S. Gabaglio. An ensemble of bayesian networks for multilabel classification. In *Twenty-third international joint conference on artificial intelligence*, 2013.

A. Arnab, S. Jayasumana, S. Zheng, and P. H. S. Torr. Higher order conditional random fields in deep neural networks. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision – ECCV 2016*, pages 524–540, Cham, 2016. Springer International Publishing.

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 Processing Magazine*, 35(1):37–52, Jan. 2018. ISSN 1558-0792. doi: 10.1109/MSP.2017.2762355.

M. Asghari, A. M. Fathollahi-Fard, S. M. J. Mirzapour Al-e-hashem, and M. A. Dulebenets. Transformation and Linearization Techniques in Optimization: A State-of-the-Art Survey. *Mathematics*, 10(2):283, Jan. 2022. ISSN 2227-7390. doi: 10.3390/math10020283.

J. Besag. Statistical Analysis of Non-Lattice Data. *The Statistician*, 24:179–195, 1975.

L.-C. Chen, A. Schwing, A. Yuille, and R. Urtasun. Learning deep structured models. In *International Conference on Machine Learning*, pages 1785–1794. PMLR, 2015.

L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, Apr. 2018. ISSN 0162-8828, 2160-9292. doi: 10.1109/TPAMI.2017.2699184.

T. Chen, M. Xu, X. Hui, H. Wu, and L. Lin. Learning Semantic-Specific Graph Representation for Multi-Label Image Recognition, Aug. 2019a.

X. Chen and A. L. Yuille. Articulated Pose Estimation by a Graphical Model with Image Dependent Pairwise Relations. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.

Z.-M. Chen, X.-S. Wei, X. Jin, and Y. Guo. Multi-Label Image Recognition with Joint Class-Aware Map Disentangling and Label Correlation Embedding. *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 622–627, July 2019b. doi: 10.1109/ICME.2019.00113.

T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng. NUS-WIDE: A real-world web image database from National University of Singapore. In *Proceedings of the ACM International Conference on Image and Video Retrieval*, pages 1–9, Santorini, Fira Greece, July 2009. ACM. ISBN 978-1-60558-480-5. doi: 10.1145/1646396.1646452.

J. Craley, E. Johnson, and A. Venkataraman. Integrating Convolutional Neural Networks and Probabilistic Graphical Modeling for Epileptic Seizure Detection in Multichannel EEG. In A. C. S. Chung, J. C. Gee, P. A.

Yushkevich, and S. Bao, editors, *Information Processing in Medical Imaging*, volume 11492, pages 291–303. Springer International Publishing, Cham, 2019. ISBN 978-3-030-20350-4 978-3-030-20351-1. doi: 10.1007/978-3-030-20351-1_22. Series Title: Lecture Notes in Computer Science.

N. Di Mauro, A. Vergari, and F. Esposito. Multi-label classification with cutset networks. In A. Antonucci, G. Corani, and C. P. Campos, editors, *Proceedings of the Eighth International Conference on Probabilistic Graphical Models*, volume 52 of *Proceedings of Machine Learning Research*, pages 147–158, Lugano, Switzerland, 06–09 Sep 2016. PMLR.

M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010. ISSN 0920-5691, 1573-1405. doi: 10.1007/s11263-009-0275-4.

H. Fan, Y. Li, B. Xiong, W.-Y. Lo, and C. Feichtenhofer. Pyslowfast. https://github.com/facebookresearch/slowfast, 2020.

C. Feichtenhofer, H. Fan, J. Malik, and K. He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, October 2019.

L. F. d. Figueiredo, A. Paes, and G. Zaverucha. Transfer learning for boosted relational dependency networks through genetic algorithm. In *International Conference on Inductive Logic Programming*, pages 125–139. Springer, 2021.

J. A. Gámez, J. L. Mateo, and J. M. Puerta. Dependency networks based classifiers: learning models by using independence test. In *Third European Workshop on Probabilistica Graphical Models (PGM06)*, pages 115–122. Citeseer, 2006.

J. A. Gámez, J. L. Mateo, T. D. Nielsen, and J. M. Puerta. Robust classification using mixtures of dependency networks. In *Proceedings of the Fourth European Workshop on Probabilistic Graphical Models*, pages 129–136, 2008.

B. Geißler, A. Martin, A. Morsi, and L. Schewe. Using Piecewise Linear Functions for Solving MINLPs. In J. Lee and S. Leyffer, editors, *Mixed Integer Nonlinear Programming*, The IMA Volumes in Mathematics and Its Applications, pages 287–314, New York, NY, 2012. Springer. ISBN 978-1-4614-1927-3. doi: 10.1007/978-1-4614-1927-3_10.

I. Griva, S. G. Nash, and A. Sofer. *Linear and Nonlinear Optimization*. SIAM, Society for Industrial and Applied Mathematics, Philadelphia, 2. ed edition, 2009. ISBN 978-0-89871-661-0.

Q. Guo and Q. Dou. Semantic Image Segmentation based on SegNetWithCRFs. *Procedia Computer Science*, 187:300–306, 2021. ISSN 18770509. doi: 10.1016/j.procs.2021.04.066.

X. Guo and Y. Weng. Deep Dependency Network for Multi-label Text Classification. In Y. Peng, Q. Liu, H. Lu, Z. Sun, C. Liu, X. Chen, H. Zha, and J. Yang, editors, *Pattern Recognition and Computer Vision*, Lecture Notes in Computer Science, pages 298–309, Cham, 2020. Springer International Publishing. ISBN 978-3-030-60636-7. doi: 10.1007/978-3-030-60636-7_25.

Y. Guo and S. Gu. Multi-label classification using conditional dependency networks. In *Twenty-Second International Joint Conference on Artificial Intelligence*, pages 1300–1305, 2011.

Y. Guo and W. Xue. Probabilistic multi-label classification with sparse feature learning. In *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence*, IJCAI '13, page 1373–1379. AAAI Press, 2013. ISBN 9781577356332.

Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023. URL https://www.gurobi.com.

K. Han, W. Wan, H. Yao, L. Hou, and School of Communication and Information Engineering, Shanghai University Institute of Smart City, Shanghai University 99 Shangda Road, BaoShan District, Shanghai 200444, China. Image Crowd Counting Using Convolutional Neural Network and Markov Random Field. *Journal of Advanced Computational Intelligence and Intelligent Informatics*, 21(4):632–638, July 2017. ISSN 1883-8014, 1343-0130. doi: 10.20965/jaciii.2017.p0632.

D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. Kadie. Dependency networks for inference, collaborative filtering, and data visualization. *Journal of Machine Learning Research*, 1(Oct):49–75, 2000.

M. J. Johnson, D. Duvenaud, A. B. Wiltschko, S. R. Datta, and R. P. Adams. Composing graphical models with neural networks for structured representations and fast inference. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 2954–2962, 2016.

A. Kirillov, D. Schlesinger, S. Zheng, B. Savchynskyy, P. H. S. Torr, and C. Rother. Joint Training of Generic CNN-CRF Models with Stochastic Optimization, 2016.

P. Knöbelreiter, C. Reinbacher, A. Shekhovtsov, and T. Pock. End-to-end training of hybrid cnn-crf models for stereo. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1456–1465, 2017. doi: 10.1109/CVPR.2017.159.

D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. Adaptive computation and machine learning. MIT Press, 2009. ISBN 9780262013192.

D. Kong, Y. Chen, H. Ma, X. Yan, and X. Xie. Adaptive graphical model network for 2d handpose estimation. In *BMVC*, 2019.

X. Kong, B. Cao, and P. S. Yu. Multi-label classification by mining label and instance correlations from heterogeneous information networks. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 614–622, Chicago Illinois USA, Aug. 2013. ACM. ISBN 978-1-4503-2174-7. doi: 10.1145/2487575.2487577.

R. G. Krishnan, U. Shalit, and D. Sontag. Deep kalman filters. *stat*, 1050:25, 2015.

M. Kumar. Converting some global optimization problems to mixed integer linear problems using piecewise linear approximations. Master's thesis, University of Missouri–Rolla, 2007.

M. Larsson, J. Alvén, and F. Kahl. Max-Margin Learning of Deep Structured Models for Semantic Segmentation. In P. Sharma and F. M. Bianchi, editors, *Image Analysis*, Lecture Notes in Computer Science, pages 28–40, Cham, 2017. Springer International Publishing. ISBN 978-3-319-59129-2. doi: 10.1007/978-3-319-59129-2_3.

M. Larsson, A. Arnab, F. Kahl, S. Zheng, and P. Torr. A Projected Gradient Descent Method for CRF Inference allowing End-To-End Training of Arbitrary Pairwise Potentials. Technical Report arXiv:1701.06805, arXiv, Jan. 2018. arXiv:1701.06805 [cs] type: article.

S.-i. Lee, V. Ganapathi, and D. Koller. Efficient structure learning of markov networks using $l_1$-regularization. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006.

S.-Y. Li, S.-J. Huang, and S. Chen. Crowdsourcing aggregation with deep Bayesian learning. *Science China Information Sciences*, 64(3):130104, Mar. 2021. ISSN 1674-733X, 1869-1919. doi: 10.1007/s11432-020-3118-7.

Z. Li, Y. Zhang, B. Sui, Z. Xing, and Q. Wang. FPGA Implementation for the Sigmoid with Piecewise Linear Fitting Method Based on Curvature Analysis. *Electronics*, 11(9):1365, Jan. 2022. ISSN 2079-9292. doi: 10.3390/electronics11091365.

G. Liang, X. Lan, J. Wang, J. Wang, and N. Zheng. A Limb-Based Graphical Model for Human Pose Estimation. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 48(7):1080–1092, July 2018. ISSN 2168-2232. doi: 10.1109/TSMC.2016.2639788. Conference Name: IEEE Transactions on Systems, Man, and Cybernetics: Systems.

G. Lin, C. Shen, A. Van Den Hengel, and I. Reid. Efficient piecewise training of deep structured models for seman-

tic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3194–3203, 2016.

M.-H. Lin, J. G. Carlsson, D. Ge, J. Shi, and J.-F. Tsai. A Review of Piecewise Linearization Methods. *Mathematical Problems in Engineering*, 2013:e101376, Nov. 2013. ISSN 1024-123X. doi: 10.1155/2013/101376.

T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.

B. Liu, X. Liu, H. Ren, J. Qian, and Y. Wang. Text multi-label learning method based on label-aware attention and semantic dependency. *Multimedia Tools and Applications*, 81(5):7219–7237, Feb. 2022. ISSN 1573-7721. doi: 10.1007/s11042-021-11663-9.

F. Liu, C. Shen, and G. Lin. Deep convolutional neural fields for depth estimation from a single image. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5162–5170, 2015.

J. Liu. *Monte Carlo strategies in scientific computing*. Springer Verlag, New York, Berlin, Heidelberg, 2008. ISBN 0-387-95230-6.

S. Liu, L. Zhang, X. Yang, H. Su, and J. Zhu. Query2Label: A Simple Transformer Way to Multi-Label Classification, July 2021.

D. Lowd. Closed-form learning of markov networks from dependency networks. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*, pages 533–542, 2012.

D. Lowd and A. Shamaei. Mean Field Inference in Dependency Networks: An Empirical Study. *Proceedings of the AAAI Conference on Artificial Intelligence*, 25(1):404–410, Aug. 2011. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v25i1.7936.

R. Mateescu, K. Kask, V. Gogate, and R. Dechter. Join-graph propagation algorithms. *Journal of Artificial Intelligence Research*, 37:279–328, 2010.

I. Naim, Y. Song, Q. Liu, H. Kautz, J. Luo, and D. Gildea. Unsupervised alignment of natural language instructions with video segments. *Proceedings of the AAAI Conference on Artificial Intelligence*, 28(1), Jun. 2014. doi: 10.1609/aaai.v28i1.8939.

I. Naim, Y. C. Song, Q. Liu, L. Huang, H. Kautz, J. Luo, and D. Gildea. Discriminative unsupervised alignment of natural language instructions with corresponding video segments. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 164–174, Denver, Colorado, May–June

2015. Association for Computational Linguistics. doi: 10.3115/v1/N15-1017.

J. Neville and D. Jensen. Collective classification with relational dependency networks. In *Workshop on Multi-Relational Data Mining (MRDM-2003)*, page 77, 2003.

H. D. Nguyen, X.-S. Vu, and D.-T. Le. Modular Graph Transformer Networks for Multi-Label Image Classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):9092–9100, May 2021. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v35i10.17098.

C. Papagiannopoulou, G. Tsoumakas, and I. Tsamardinos. Discovering and Exploiting Deterministic Label Relationships in Multi-Label Learning. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 915–924, Sydney NSW Australia, Aug. 2015. ACM. ISBN 978-1-4503-3664-2. doi: 10.1145/2783258.2783302.

X. Qu, H. Che, J. Huang, L. Xu, and X. Zheng. Multi-layered Semantic Representation Network for Multi-label Image Classification, June 2021.

M. Regneri, M. Rohrbach, D. Wetzel, S. Thater, B. Schiele, and M. Pinkal. Grounding action descriptions in videos. *Transactions of the Association for Computational Linguistics (TACL)*, 1:25–36, 2013.

R. Rovatti, C. D'Ambrosio, A. Lodi, and S. Martello. Optimistic MILP modeling of non-linear optimization problems. *European Journal of Operational Research*, 239(1):32–45, 2014.

A. G. Schwing and R. Urtasun. Fully Connected Deep Structured Networks. Technical report, arXiv, Mar. 2015.

B. Selman, H. A. Kautz, and B. Cohen. Local search strategies for satisfiability testing. In *Cliques, Coloring, and Satisfiability*, 1993. URL https://api.semanticscholar.org/CorpusID:3215289.

G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *European Conference on Computer Vision*, pages 510–526. Springer, 2016.

J. Song, L. Wang, L. V. Gool, and O. Hilliges. Thin-slicing network: A deep structured model for pose estimation in videos. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5563–5572, 2017.

C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the Inception Architecture for Computer Vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, Las Vegas, NV, USA, June 2016. IEEE. ISBN 978-1-4673-8851-1. doi: 10.1109/CVPR.2016.308.

M. Tan, Q. Shi, A. van den Hengel, C. Shen, J. Gao, F. Hu, and Z. Zhang. Learning graph structure for multi-label

image classification via clique generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4100–4109, June 2015.

C. Tarantola and E. Blanc. Dependency networks and bayesian networks for web mining. *WIT Transactions on Information and Communication Technologies*, 28, 2002.

J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. *CoRR*, abs/1406.2984, 2014.

P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=rJXMpikCZ.

M. J. Wainwright, J. Lafferty, and P. Ravikumar. High-dimensional graphical model selection using $\ell_1$-regularized logistic regression. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006.

H. Wang, M. Huang, and X. Zhu. A generative probabilistic model for multi-label classification. In *2008 Eighth IEEE International Conference on Data Mining*, pages 628–637. IEEE, 2008.

L. Wang, Y. Liu, H. Di, C. Qin, G. Sun, and Y. Fu. Semi-supervised dual relation learning for multi-label classification. *IEEE Transactions on Image Processing*, 30:9125–9135, 2021a.

R. Wang, R. Ridley, X. Su, W. Qu, and X. Dai. A novel reasoning mechanism for multi-label text classification. *Information Processing & Management*, 58(2):102441, Mar. 2021b. ISSN 03064573. doi: 10.1016/j.ipm.2020.102441.

S. Wang, J. Wang, Z. Wang, and Q. Ji. Enhancing multi-label classification by modeling dependencies among labels. *Pattern Recognition*, 47(10):3405–3413, Oct. 2014. ISSN 00313203. doi: 10.1016/j.patcog.2014.04.009.

W. Weng, B. Wei, W. Ke, Y. Fan, J. Wang, and Y. Li. Learning label-specific features with global and local label correlation for multi-label classification. *Applied Intelligence*, 53(3):3017–3033, Feb. 2023. ISSN 1573-7497. doi: 10.1007/s10489-022-03386-7.

A. Wu, E. K. Buchanan, M. Whiteway, M. Schartner, G. Meijer, J.-P. Noel, E. Rodriguez, C. Everett, A. Norovich, E. Schaffer, N. Mishra, C. D. Salzman, D. Angelaki, A. Bendesky, T. I. B. L. The International Brain Laboratory, J. P. Cunningham, and L. Paninski. Deep graph pose: a semi-supervised deep graphical model for improved animal pose tracking. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and

H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6040–6052. Curran Associates, Inc., 2020.

W. Yang, W. Ouyang, H. Li, and X. Wang. End-to-End Learning of Deformable Mixture of Parts and Deep Convolutional Neural Networks for Human Pose Estimation. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3073–3082, June 2016. doi: 10.1109/CVPR.2016.335. ISSN: 1063-6919.

K. Yao, B. Peng, G. Zweig, D. Yu, X. Li, and F. Gao. Recurrent conditional random field for language understanding. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4077–4081, Florence, Italy, May 2014. IEEE. ISBN 978-1-4799-2893-4. doi: 10.1109/ICASSP.2014. 6854368.

J. S. Yedidia, W. Freeman, and Y. Weiss. Generalized Belief Propagation. In *Advances in Neural Information Processing Systems*, volume 13. MIT Press, 2000.

D. Yu, B. Yang, Q. Wei, A. Li, and S. Pan. A probabilistic graphical model based on neural-symbolic reasoning for visual relationship detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10609–10618, 2022.

M. Zhang, N. Wang, Y. Li, and X. Gao. Neural Probabilistic Graphical Model for Face Sketch Synthesis. *IEEE Transactions on Neural Networks and Learning Systems*, 31(7):2623–2637, July 2020. ISSN 2162-237X, 2162-2388. doi: 10.1109/TNNLS.2019.2933590.

S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr. Conditional Random Fields as Recurrent Neural Networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1529–1537, Santiago, Chile, Dec. 2015. IEEE. ISBN 978-1-4673-8391-2. doi: 10.1109/ ICCV.2015.179.

W. Zhou, Z. Xia, P. Dou, T. Su, and H. Hu. Double Attention Based on Graph Attention Network for Image Multi-Label Classification. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 19(1):1–23, Jan. 2023. ISSN 1551-6857, 1551-6865. doi: 10.1145/3519030.

M. Zhu, J. Li, N. Wang, and X. Gao. Learning Deep Patch representation for Probabilistic Graphical Model-Based Face Sketch Synthesis. *International Journal of Computer Vision*, 129(6):1820–1836, June 2021. ISSN 0920-5691, 1573-1405. doi: 10.1007/s11263-021-01442-2.

## Checklist

1. For all models and algorithms presented, check if you include:

   (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]

   (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Not Applicable]

   (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes]

2. For any theoretical claim, check if you include:

   (a) Statements of the full set of assumptions of all theoretical results. [Yes]

   (b) Complete proofs of all theoretical results. [Yes]

   (c) Clear explanations of any assumptions. [Yes]

3. For all figures and tables that present empirical results, check if you include:

   (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes]

   (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes]

   (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]

   (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes]

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:

   (a) Citations of the creator If your work uses existing assets. [Yes]

   (b) The license information of the assets, if applicable. [Not Applicable]

   (c) New assets either in the supplemental material or as a URL, if applicable. [Yes]

   (d) Information about consent from data providers/curators. [Not Applicable]

   (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]

5. If you used crowdsourcing or conducted research with human subjects, check if you include:

   (a) The full text of instructions given to participants and screenshots. [Not Applicable]

   (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]

   (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

## A DERIVING THE OBJECTIVE FUNCTION FOR MOST PROBABLE EXPLANATIONS (MPE)

Let us consider the scoring function we aim to maximize, given as follows:

$$\underset{\mathbf{x}}{\text{maximize}} \sum_{i=1}^{n} \log \left( P_i(x_i \mid \mathbf{x}_{-i}, \mathbf{e}) \right) \tag{10}$$

Given that $x_i$ is binary, namely, $x_i \in \{0, 1\}$, we can express the scoring function as follows:

$$\underset{\mathbf{x}}{\text{maximize}} \sum_{i=1}^{n} \left( x_i \log \left( P_i(X_i = 1 \mid \mathbf{x}_{-i}, \mathbf{e}) \right) + (1 - x_i) \log \left( 1 - P_i(X_i = 1 \mid \mathbf{x}_{-i}, \mathbf{e}) \right) \right) \tag{11}$$

where

$$P_i(X_i = 1 | \mathbf{x}_{-i}, \mathbf{e}) = \sigma \left( \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i \right)$$

Let $p_i = P_i(X_i = 1 | \mathbf{x}_{-i}, \mathbf{e})$ and $z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i$. Then, the MPE task, which involves optimizing the scoring function given above, can be expressed as:

$$\underset{\mathbf{x}, \mathbf{z}}{\text{maximize}} \sum_{i=1}^{n} \left( x_i \log p_i + (1 - x_i) \log(1 - p_i) \right) \tag{12}$$

subject to:

$$p_i = \sigma(z_i), \quad \forall i \in \{1, \dots, n\} \tag{13}$$

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \dots, n\} \tag{14}$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\} \tag{15}$$

The first constraint in the formulation, as expressed by 13, pertains to the sigmoid function employed in the logistic regression module. The subsequent constraint, defined by 14, formalizes the product between the weights and inputs in the Dependency Network (DN). Here, $e_i$ represents evidence values supplied to the DN, while $\mathbf{x}$ represents the decision variables that are subject to optimization. Finally, 15 specifies that the input variables must be constrained to integer values of 0 or 1.

The assembled constraints collectively simulate a forward pass through the network. The objective function is designed to maximize the scoring function pertinent to the Most Probable Explanation (MPE). This framework adeptly aligns the optimization task to maximize the MPE score.

The substitution of the constraint $P_i(x_i|\mathbf{x}_{-i}, \mathbf{e}) = \sigma(z_i) = \frac{1}{1+e^{(-z_i)}}$ into the objective function and the subsequent algebraic simplification result in a simplified formulation as follows:

$$
\begin{aligned}
& x_i \log p_i + (1 - x_i) \log(1 - p_i) \\
&= x_i \log \left( \frac{e^{z_i}}{1 + e^{z_i}} \right) + (1 - x_i) \log \left( 1 - \frac{e^{z_i}}{1 + e^{z_i}} \right) \\
&= x_i \log \left( \frac{e^{z_i}}{1 + e^{z_i}} \right) + (1 - x_i) \log \left( \frac{1}{1 + e^{z_i}} \right) \\
&= x_i \log e^{z_i} - x_i \log(1 + e^{z_i}) - \log(1 + e^{z_i}) + x_i \log(1 + e^{z_i}) \\
&= x_i \log e^{z_i} - \log(1 + e^{z_i}) \\
&= x_i z_i - \log(1 + e^{z_i})
\end{aligned} \tag{16}
$$

Substituting equation 16 in the objective function of our optimization problem, we get

$$\underset{\mathbf{x},\mathbf{z}}{\text{maximize}} \sum_{i=1}^{n} x_i z_i - \log\left(1 + e^{z_i}\right) \tag{17}$$
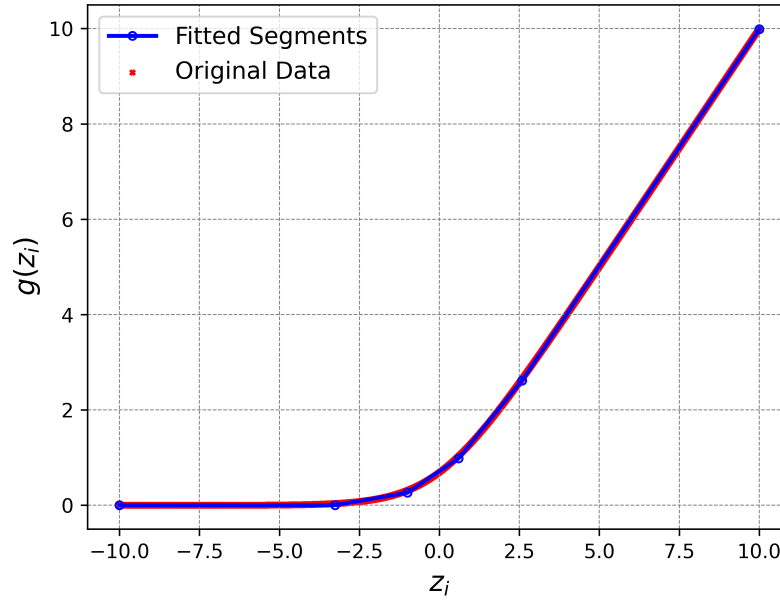
subject to:

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \ldots, n\} \tag{18}$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \ldots, n\} \tag{19}$$

## A.A  Utilizing Piecewise Linear Approximation for Non-linear Functions

Figure SF4: Piece-wise linear approximation of $log(1 + e^{z_i})$



The objective function expressed in equation equation 17 is non-linear due to the inclusion of the term $\log(1 + e^{z_i})$. $f(z)$ can be used as a piecewise approximation for $\log(1 + e^{z_i})$.

$$f(z) = \begin{cases} z & z \gg 1 \\ e^z & z \ll 1 \\ \log(1 + e^z) & \text{otherwise} \end{cases} \tag{20}$$

To obtain a piecewise linear approximation of $\log(1 + e^{z_i})$, a single linear function suffices for $z \gg 1$, while the majority of the linear pieces are utilized for approximating the function near 1. A piecewise linear approximation of the function is detailed in the figure SF4. We employ five segments for the approximation. These piecewise functions can be integrated as linear constraints, facilitating the conversion of the non-linear objective into a linear objective with ease. We also present the piecewise equations for the approximation as follows -

Table ST3: Piecewise Linear Approximation for
$$g(z) \approx \log(1 + e^{z_i})$$

| $z$ | $g(z)$ |
|---|---|
| $]-\infty, -3.257)$ | $0$ |
| $[-3.257, -0.998)$ | $(2^{-4} + 2^{-5} + 2^{-6} + 2^{-8})z + 0.379$ |
| $[-0.998, 0.602)$ | $(2^{-2} + 2^{-3} + 2^{-4} + 2^{-7} + 2^{-8})z + 0.715$ |
| $[0.602, 2.584)$ | $(2^{-1} + 2^{-2} + 2^{-4} + 2^{-7})z + 0.492$ |
| $[2.584, +\infty[$ | $z$ |

Thus after replacing $\log(1 + e^{z_i})$ with its piece-wise approximation ($g(z)$) we get

$$\underset{\mathbf{x}, \mathbf{z}, \alpha}{\text{maximize}} \sum_{i=1}^{n} x_i z_i - g(z_i)$$

subject to:

$$z_i = \sum_{j=1}^{|\mathbf{e}|} w_{ij} e_j + \sum_{\substack{k=1 \\ k \neq i}}^{|\mathbf{x}|} v_{ik} x_k + b_i, \quad \forall i \in \{1, \ldots, n\}$$

$$x_i \in \{0, 1\}, \quad \forall i \in \{1, \ldots, n\}$$

$$
\begin{aligned}
g(z_i) = {} & \alpha_{1i} \times 0 + \alpha_{2i} \times ((2^{-4} + 2^{-5} + 2^{-6} + 2^{-8})z_i + 0.379) + \\
& \alpha_{3i} \times ((2^{-2} + 2^{-3} + 2^{-4} + 2^{-7} + 2^{-8})z_i + 0.715) + \\
& \alpha_{4i} \times ((2^{-1} + 2^{-2} + 2^{-4} + 2^{-7})z_i + 0.492) + \\
& \alpha_{5i} z_i, \quad \forall i \in \{1, \ldots, n\}
\end{aligned}
\tag{21}
$$

$$\sum_{j=1}^{5} \alpha_{ji} = 1, \quad \forall i \in \{1, \ldots, n\}$$

$$\alpha_{ji} \in \{0, 1\}, \quad \forall j \in \{1, 2, 3, 4, 5\}, \forall i \in \{1, \ldots, n\}$$

$$\alpha_{1i} = 1 \Rightarrow z_i < -3.257, \quad \forall i \in \{1, \ldots, n\}$$

$$\alpha_{2i} = 1 \Rightarrow -3.257 \leq z_i < -0.998, \quad \forall i \in \{1, \ldots, n\}$$

$$\alpha_{3i} = 1 \Rightarrow -0.998 \leq z_i < 0.602, \quad \forall i \in \{1, \ldots, n\}$$

$$\alpha_{4i} = 1 \Rightarrow 0.602 \leq z_i < 2.584, \quad \forall i \in \{1, \ldots, n\}$$

$$\alpha_{5i} = 1 \Rightarrow z_i \geq 2.584, \quad \forall i \in \{1, \ldots, n\}$$

The utilization of binary variables, $\alpha_{ji}$, allows for a piece-wise linear approximation of the logarithm of $(1+e^{z_i})$, with these variables serving as selectors to determine the appropriate linear segment based on the value of $z_i$. It is worth noting that the final five constraints are indicator constraints, which can be linearized using the big-M method as described in Griva et al. (2009). Consequently, the problem can be formulated as a mixed-integer multi-linear optimization problem. The mixed-integer multi-linear problem can be further converted to a mixed-integer linear problem (MILP) using the approach describe next.

## B FORMULATING LOGICAL AND BETWEEN BINARY VARIABLES AS LINEAR CONSTRAINTS

Consider an optimization problem involving three binary variables $X_1$, $X_2$, and $X_3$, where $x_1$, $x_2$, and $x_3$ denote their respective assignments. We define an objective function that incorporates products of these binary variables as follows -

$$\max_{x_1, x_2, x_3} x_1 x_2 - x_2 x_3 + x_1 x_3 \tag{22}$$

Although constraints are not incorporated in this instance, the same methodology can be extended to constrained optimization problems. Subsequently, auxiliary variables $z_1$, $z_2$, and $z_3$ are introduced to account for each of the binary products.

The optimization problem can be stated as follows:

$$
\begin{aligned}
\max_{x_1, x_2, x_3} \quad & z_1 - z_2 + z_3 \\
\text{s.t.} \quad & x_1 \wedge x_2 = z_1, \\
& x_2 \wedge x_3 = z_2, \\
& x_1 \wedge x_3 = z_3.
\end{aligned}
\tag{23}
$$

The optimization problem may be expressed by incorporating additional linear constraints to represent each boolean product. The resulting problem is presented as follows:

$$
\begin{aligned}
\max_{x_1, x_2, x_3} \quad & z_1 - z_2 + z_3 \\
\text{s.t.} \quad & x_1 + x_2 - 1 \leq z_1, \\
& z_1 \leq x_1, \\
& z_1 \leq x_2, \\
& x_2 + x_3 - 1 \leq z_2, \\
& z_2 \leq x_2, \\
& z_2 \leq x_3, \\
& x_1 + x_3 - 1 \leq z_3, \\
& z_3 \leq x_1, \\
& z_3 \leq x_3.
\end{aligned}
\tag{24}
$$

The problem formulation outlined in Section A, and more specifically, equation equation 21, includes expressions of the form $x_i x_j$. These terms serve to represent the logical AND operation between binary variables $x_i$ and $x_j$, and are a direct result of the product $x_i \times g_i$. Consequently, this optimization problem is classified as a mixed-integer multi-linear optimization problem.

The formulation detailed in this section enables the capture of these logical AND operations between binary variables through linear constraints, thereby enabling the utilization of Mixed-Integer Linear Programming (MILP) solvers, e.g., Gurobi Optimization, LLC (2023), to find the optimal solution (or an anytime, near optimal solution if a time bound is specified).

## C GIBBS SAMPLING FOR DDNS

---

**Algorithm 1** Gibbs Sampling for DDNs

---

**Input:** video segment/image $\mathbf{v}$, number of samples $N$, DDN $\langle \mathcal{N}, \mathcal{D} \rangle$
**Output:** An estimate of the marginal probability distribution over each label $X_i$ of the DDN given $\mathbf{v}$

1: $\mathbf{e} \leftarrow \mathbb{N}(\mathbf{v})$
2: Randomly initialize $\mathbf{X} = \mathbf{x}^{(0)}$
3: **for** $j = 1$ **to** $N$ **do**
4:     $\pi \leftarrow$ Generate random permutation of $[1, n]$
5:     **for** $i = 1$ **to** $n$ **do**
6:         $x_{\pi(i)}^{(j)} \sim P_{\pi(i)}(x_{\pi(i)} | \mathbf{x}_{\pi(1):\pi(i-1)}^{(j)}, \mathbf{x}_{\pi(i+1):\pi(n)}^{(j-1)}, \mathbf{e})$
7: **for** $i = 1$ **to** $n$ **do**
8:     $\hat{P}_i(x_i | \mathbf{v}) = \frac{1}{N} \sum_{j=1}^{N} P_i(x_i | \mathbf{x}_{-i}^{(j)}, \mathbf{e})$
9: **return** $\{\hat{P}_i(x_i | \mathbf{v}) | i \in \{1, \dots, n\}\}$

---

This section describes the Gibbs Sampling inference procedure for DDNs (see Algorithm 1). Gibbs Sampling serves as an approximation method for the Most Probable Explanation (MPE) inference task in DDNs by offering max-marginals, which can subsequently be utilized to approximate MPE. The inputs to the algorithm are (1) a video segment/image $\mathbf{v}$, (2)

the number of samples $N$ and (3) trained DDN model $\langle \mathcal{N}, \mathcal{D} \rangle$. The algorithm begins (see step 1) by extracting features $\mathbf{e}$ from the video segment/image $\mathbf{v}$ by sending the latter through the neural network $\mathcal{N}$ (which represents the function $\mathbb{N}$). Then in steps 2–8, it generates $N$ samples via Gibbs sampling. The Gibbs sampling procedure begins with a random assignment to all the labels (step 2). Then at each iteration (steps 3–8), it first generates a random permutation $\pi$ over the $n$ labels and samples the labels one by one along the order $\pi$ (steps 5–7). To sample a label indexed by $\pi(i)$ at iteration $j$, we compute $P_{\pi(i)}(x_{\pi(i)} | \mathbf{x}^{(j)}_{\pi(1):\pi(i-1)}, \mathbf{x}^{(j-1)}_{\pi(i+1):\pi(n)}, \mathbf{e})$ from the DN $\mathcal{D}$ where $\mathbf{x}^{(j)}_{\pi(1):\pi(i-1)}$ and $\mathbf{x}^{(j-1)}_{\pi(i+1):\pi(n)}$ denote the assignments to all labels ordered before $x_{\pi(i)}$ at iteration $j$ and the assignments to all labels ordered after $x_{\pi(i)}$ at iteration $j-1$ respectively.

After $N$ samples are generated via Gibbs sampling, the algorithm uses them to estimate (see steps 9–11) the (posterior) marginal probability distribution at each label $X_i$ given $\mathbf{v}$ using the mixture estimator (Liu, 2008). The algorithm terminates (see step 12) by returning these posterior estimates.

## D   PRECISION BASED EVALUATION METRICS

In Tables ST4 and ST5, we present a comprehensive account of precision-focused metrics, specifically Mean Average Precision (mAP) and Label Ranking Average Precision (LRAP). It should be noted that these metrics are not directly applicable to the Most Probable Explanation (MPE) task, as they necessitate scores (probabilities) of the output labels. However, given that SlowFast, InceptionV3, Q2L, MSRN, and DDN-GS can provide such probabilities, their precision scores are anticipated to be elevated. We employed Gibbs Sampling to approximate the MPE, enabling the derivation of marginal probabilities to approximate max-marginals and, thereby, the MPE.

In majority of the instances, the Gibbs Sampling-based inference on DDNs yields performance metrics closely aligned with baseline methods. Variability exists, with some metrics exceeding the baseline in the context of MLAC and falling short in the case of MLIC. Notably, the Integer Linear Programming (ILP)-based approach for DDN outperforms all alternative methods across three datasets when evaluated on LRAP. DRF-based methods and DDN-specific inference techniques generally underperform on these precision metrics. Again, this outcome is anticipated, given that apart from Gibbs Sampling, none of the other techniques can generate probabilistic scores for labels, leading to their inferior performance.

Table ST4: Comparison of our methods with the feature extractor for MLAC - Precision-based metrics. The best/second best values are bold/underlined.

| Method | Charades | | TACoS | | Wetlab | |
|---|---|---|---|---|---|---|
| | mAP | LRAP | mAP | LRAP | mAP | LRAP |
| SlowFast (Fan et al., 2020) | <u>0.39</u> | <u>0.53</u> | | | | |
| InceptionV3 (Szegedy et al., 2016) | | | <u>0.70</u> | 0.81 | <u>0.79</u> | 0.82 |
| DRF - GS | 0.27 | 0.44 | 0.56 | 0.79 | 0.54 | 0.76 |
| DRF - ILP | 0.19 | 0.28 | 0.40 | 0.67 | 0.63 | 0.73 |
| DRF - IJGP | 0.31 | 0.44 | 0.56 | 0.81 | 0.79 | <u>0.85</u> |
| DDN - GS | **0.40** | **0.55** | **0.75** | <u>0.84</u> | **0.84** | **0.87** |
| DDN - RW | 0.19 | 0.28 | 0.49 | 0.66 | 0.63 | 0.77 |
| DDN - Greedy | 0.32 | 0.30 | 0.50 | 0.69 | 0.65 | 0.78 |
| DDN - ILP | 0.36 | <u>0.53</u> | 0.51 | **0.85** | 0.65 | **0.87** |

## E   ANNOTATIONS COMPARISON BETWEEN Q2L AND DDN-MLP-JOINT ON THE MS-COCO DATASET

Table ST6 presents a qualitative evaluation of the label predictions produced by our DDN-ILP inference scheme compared to the baseline Q2L using randomly selected images from the MS-COCO dataset. This analysis provides a perspective on the advantages of DDN-ILP over Q2L, particularly in terms of correcting errors generated by the feature extractor. In the initial set of seven rows, DDN-ILP adds additional correct labels to the Q2L output. For certain images, the objects overlooked by Q2L are inherently challenging to identify. In these instances, the DDN-ILP method, which infers labels based on label relationships, effectively rectifies the limitations of Q2L. DDN-ILP refines Q2L's predictions in the following seven instances by removing incorrect predictions and aligning precisely with the ground truth labels. Finally, we

Table ST5: Comparison of our methods with the feature extractor for MLIC for Precision based metrics. The best/second best values are bold/underlined.

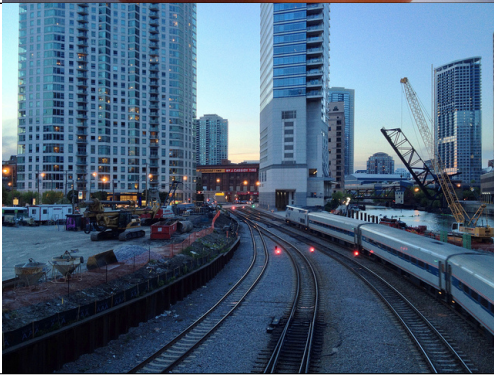| Method | MS-COCO | | NUS-WIDE | | PASCAL-VOC | |
|---|---|---|---|---|---|---|
| | mAP | LRAP | mAP | LRAP | mAP | LRAP |
| Q2L (Liu et al., 2021) | **0.91** | **0.96** | | | | |
| MSRN (Qu et al., 2021) | | | **0.62** | **0.85** | <u>0.96</u> | **0.98** |
| DRF - GS | 0.75 | 0.86 | 0.40 | 0.74 | 0.77 | 0.93 |
| DRF - ILP | 0.74 | 0.82 | 0.25 | 0.59 | 0.81 | 0.88 |
| DRF - IJGP | 0.74 | 0.90 | 0.41 | 0.75 | 0.83 | 0.94 |
| DDN - GS | 0.84 | <u>0.93</u> | <u>0.50</u> | <u>0.82</u> | 0.92 | <u>0.96</u> |
| DDN - RW | 0.74 | 0.82 | 0.24 | 0.61 | 0.91 | 0.94 |
| DDN - Greedy | 0.74 | 0.83 | 0.35 | 0.68 | 0.93 | 0.95 |
| DDN - ILP | <u>0.85</u> | 0.83 | 0.48 | <u>0.82</u> | **0.97** | **0.98** |

examine cases where DDN-ILP's modifications result in inaccuracies (rows have grey background), contrasting with Q2L's correct predictions. These results demonstrate the effectiveness of DDN-ILP in improving the predicted labels relative to the baseline.

Table ST6: Comparison of labels predicted by Q2L (Liu et al., 2021) and our DDN-ILP scheme on the MS-COCO dataset. Labels inside **[]** represent the difference between the predictions of the two methods, assuming that a threshold of 0.5 is used (i.e., every label whose probability $> 0.5$ is considered a predicted label). Due to the MPE focus in DDN-ILP, only label configurations are generated, omitting corresponding probabilities.

| Image | Ground Truth | Q2L | DDN |
|---|---|---|---|
|  | person, car, truck, traffic light, skateboard | person (1.00), car (1.00), skateboard (1.00), **[truck (0.33), traffic light (0.41)]** | person, car, truck, traffic light, skateboard |
|  | knife, spoon, microwave, oven, sink, refrigerator | microwave (1.00), oven (1.00), sink (1.00), refrigerator (1.00), knife (0.98), **[spoon (0.38)]** | knife, spoon, microwave, oven, sink, refrigerator |

| | | | |
|---|---|---|---|
|  | person, skis, snowboard | person (1.00), skis (1.00), **[snowboard (0.20)]** | person, skis, snowboard |
|  | person, truck, suitcase | truck (1.00), suitcase (0.98), **[person (0.24)]** | person, truck, suitcase |
|  | car, bus, truck, cat, dog | car (0.99), truck (0.99), cat (0.96), **[bus (0.37), dog (0.15)]** | car, bus, truck, cat, dog |
|  | fork, sandwich, hot dog, dining table | fork (1.00), hot dog (1.00), dining table (0.90), **[sandwich (0.28)]** | fork, sandwich, hot dog, dining table |

| | | | |
|---|---|---|---|
|  | person, bottle, wine glass, chair, dining table | person (1.00), bottle (1.00), wine glass (1.00), dining table (0.99), **[chair (0.25)]** | person, bottle, wine glass, chair, dining table |
|  | person, sports ball, baseball glove, | person (1.00), baseball glove (1.00), **[sports ball (0.27)]** | person, sports ball, baseball glove, |
|  | person, teddy bear | person (1.00), teddy bear (1.00), **[handbag (0.58), clock (0.96)]** | person, teddy bear |
|  | person, cup, chair, dining table, cell phone, | person (1.00), cup (1.00), chair (1.00), dining table (1.00), cell phone (1.00), **[clock (0.92), potted plant (0.85)]** | person, cup, chair, dining table, cell phone, |

|  | cup, banana, apple | cup (1.00), banana (1.00), apple (1.00), **[spoon (0.90), fork (0.55)]** | cup, banana, apple |
| --- | --- | --- | --- |
|  | cup, banana | banana (1.00), **[bottle (0.90), oven (0.66)]**, cup (0.54) | cup, banana |
|  | train | train (1.00), **[traffic light (0.90), car (0.56), , truck (0.51)]** | train |

| | | | |
|---|---|---|---|
|  | person, bowl, oven | person (1.00), oven (1.00) , bowl (0.99), **[chair (0.85), spoon (0.59)]** | person, bowl, oven |
|  | person, cup, spoon, cake, dining table | cup (1.00), spoon (1.00), cake (1.00), dining table (0.95), person (0.91), **[donut (0.89), fork (0.65)]** | person, cup, spoon, cake, dining table |
|  | bird, skateboard, couch | bird (1.00), skateboard (1.00), couch (0.73) | bird, skateboard, **[couch]** |

| | | | |
|---|---|---|---|
|  | person, bench, suitcase | person (1.00), suitcase (1.00), bench (0.57) | person, suitcase, **[bench]** |
|  | bed, clock | bed (1.00), clock (0.51) | bed, **[clock]** |
|  | bicycle, car, cat, bottle | bicycle (1.00), car (1.00), cat (1.00), bottle (0.77) | bicycle, car, cat, **[bottle]** |
|  | person, bottle, wine glass, fork, bowl, dining table, sandwich | person (1.00), bottle (1.00), wine glass (1.00), fork (1.00), bowl (1.00), dining table (0.99), sandwich (0.52) | person, bottle, wine glass, fork, bowl, dining table, **[sandwich]** |