
Sparse-Matrix Belief Propagation

Reid Bixler

Department of Computer Science
Virginia Tech
Blacksburg, VA 24061

Bert Huang

Department of Computer Science
Virginia Tech
Blacksburg, VA 24061

Abstract

We propose sparse-matrix belief propagation, which executes loopy belief propagation in pairwise Markov random fields by replacing indexing over graph neighborhoods with sparse-matrix operations. This abstraction allows for seamless integration with optimized sparse linear algebra libraries, including those that perform matrix and tensor operations on modern hardware such as graphical processing units (GPUs). The sparse-matrix abstraction allows the implementation of belief propagation in a high-level language (e.g., Python) that is also able to leverage the power of GPU parallelization. We demonstrate sparse-matrix belief propagation by implementing it in a modern deep learning framework (PyTorch), measuring the resulting massive improvement in running time, and facilitating future integration into deep learning models.

1 INTRODUCTION

Belief propagation is a canonical inference algorithm for graphical models such as Markov random fields (MRFs) or Bayesian networks (Pearl, 2014; Wainwright et al., 2008). In graphs with cycles, loopy belief propagation performs approximate inference. Loopy belief propagation passes messages from the variable nodes to their neighbors along the graph structure. These messages are fused to estimate marginal probabilities, also referred to as beliefs. After enough iterations of the algorithm, these beliefs tend to represent a good approximate solution to the actual marginal probabilities. In this paper, we consider pairwise MRFs, which only have unary and pairwise factors.

One drawback of loopy belief propagation is that, though

the algorithm is relatively simple, its implementation requires management of often irregular graph structures. This fact usually results in tedious indexing in software. The algorithm’s message-passing routines can be compiled to be rather efficient, but when implemented in a high-level language, such as those used by data scientists, they can be prohibitively slow. Experts typically resort to writing external software in lower-level, compiled languages such as C++. The implementation of belief propagation (and its variants) as separate, compiled libraries creates a barrier for its integration into high-level data science workflows.

We instead derive loopy belief propagation for pairwise MRFs as a sequence of matrix operations, resulting in sparse-matrix belief propagation. In particular, we use sparse-matrix products to represent the message-passing indexing. The resulting algorithm can then be implemented in a high-level language, and it can be executed using highly optimized sparse and dense matrix operations. Since matrix operations are much more general than loopy belief propagation, they are often built in as primitives in high-level mathematical languages. Moreover, their generality provides access to interfaces that implement matrix operations on modern hardware, such as graphical processing units (GPUs).

In this paper, we describe sparse-matrix belief propagation and analyze its running time, showing that its running time is asymptotically equivalent to loopy belief propagation. We also describe how the abstraction can be used to implement other variations of belief propagation. We then demonstrate its performance on a variety of tests. We compare loopy belief propagation implemented in Python and in C++ against sparse-matrix belief propagation using `scipy.sparse` and PyTorch on CPUs and PyTorch on GPUs. The results illustrate the advantages of the sparse-matrix abstraction, and represent a first step toward full integration of belief propagation into modern machine learning and deep learning workflows.

1.1 RELATED WORK

Belief propagation is one of the canonical variational inference methods for probabilistic graphical models (Pearl, 2014; Wainwright et al., 2008). Loopy belief propagation is naturally amenable to fine-grained parallelism, as it involves sending messages across edges in a graph in parallel. The algorithm is classical and well studied, but because it involves tight loops and intricate indexing, it cannot be efficiently implemented in high-level or mathematical programming languages. Instead, practitioners rely on implementations in low-level languages such as C++ (Schmidt, 2007; Andres et al., 2012). Specific variations for parallel computing have been proposed (Schwing et al., 2011), and other variations have been implemented in graph-based parallel-computing frameworks (Low et al., 2014; Malewicz et al., 2010). Specialized implementations have also been created for belief propagation on GPUs (Zheng et al., 2012).

While loopy belief propagation is the canonical message-passing inference algorithm, many variations have been created to address some of its shortcomings. Some variations modify the inference objective to make belief propagation a convex optimization, such as tree-reweighted belief propagation (Wainwright et al., 2003) and convexified belief propagation (Meshi et al., 2009). Other variations compute the most likely variable state rather than marginal probabilities, such as max-product belief propagation (Wainwright et al., 2008) and max-product linear programming (Globerson and Jaakkola, 2008).

Linearized belief propagation approximates the message update formulas with linear operations (Gatterbauer et al., 2015), which, like our approach, can benefit from highly optimized linear algebra libraries and specialized hardware. However, our approach aims to retain the exact non-linear formulas of belief propagation (and variants), while linearized belief propagation is an approximation.

Our approach of using sparse matrices as an abstraction layer for implementing belief propagation relies on sparse matrix operations being implemented in efficient, optimized, compiled libraries. Special algorithms have been developed to parallelize these operations on GPUs, enabling sparse-matrix computations to use the thousands of cores typically available in such hardware (Bell and Garland, 2008). Other libraries for sparse-matrix computation, such as those built into MATLAB (Gilbert et al., 1992) and `scipy.sparse` often seamlessly provide multi-core parallelism. Frameworks for large-scale, distributed matrix computation, such as Apache Spark (Bosagh Zadeh et al., 2016), can also be used as backends for our approach. Finally, one of the more important recent advances in computing hardware, field-programmable gate arrays (FPGAs), also support sparse-

matrix operations (Zhuo and Prasanna, 2005).

By formulating belief propagation as a series of matrix and tensor operations, we make it fit into modern deep learning software frameworks, such as PyTorch (Paszke et al., 2017) and TensorFlow (Abadi et al., 2016). Since these frameworks are designed to easily switch between CPU and GPU computation, we use PyTorch for one of our belief propagation implementations in our experiments. The added advantage is that, once these frameworks fully support back-propagation through sparse matrix products, we will be able to immediately back-propagate through belief propagation. Back-propagating through inference has been shown to allow more robust training of graphical model parameters (Domke, 2013).

2 BACKGROUND

In this section, we review belief propagation and some common—but not often documented—simplifications. Define a pairwise Markov random field (MRF) as a factorized probability distribution such that the probability of variable $x \in \mathbb{X}$ is

$$\Pr(X = \mathbf{x}) = \frac{1}{Z} \exp \left(\sum_{i \in \mathcal{V}} \phi_i(x_i) + \sum_{(i,j) \in \mathcal{E}} \phi(x_i, x_j) \right), \quad (1)$$

where the normalizing constant Z is

$$Z = \sum_{\mathbf{x} \in \mathbb{X}} \exp \left(\sum_{i \in \mathcal{V}} \phi_i(x_i) + \sum_{(i,j) \in \mathcal{E}} \phi(x_i, x_j) \right), \quad (2)$$

and \mathbf{x} represents the full state vector of all variables, $x_i \in \mathbb{X}_i$ represents the state of the i th variable, and $G = \{\mathcal{V}, \mathcal{E}\}$ is a graph defining the structure of the MRF.

The goal of marginal inference is to compute or approximate the marginal probabilities of the variables

$$\{\Pr(x_1), \dots, \Pr(x_n)\} \quad (3)$$

and of other subsets of variables. In pairwise MRFs, the marginal inference is often limited to the unary marginals and the pairwise marginals along the edges.

2.1 LOOPY BELIEF PROPAGATION

Belief propagation is a dynamic programming algorithm for computing marginal inference in tree-structured MRFs that is often applied in practice on non-tree, i.e., loopy, graphs. Loopy belief propagation is therefore a heuristic approximation that works well in many practical scenarios. The algorithm operates by passing messages among variables along the edges of the MRF structure.

The message sent from variable x_i to variable x_j is

$$m_{i \rightarrow j}[x_j] = \log \left(\sum_{x_i} \exp(a_{x_i}) \right), \quad (4)$$

where (as shorthand to fit the page)

$$a_{x_i} = \phi_{ij}(x_i, x_j) + \phi_i(x_i) + \sum_{k \in N_i \setminus j} m_{k \rightarrow i}[x_i] - d_{ij}; \quad (5)$$

N_i is the set of neighbors of variable i , i.e., $N_i = \{k | (i, k) \in \mathcal{E}\}$; and d_{ij} is any constant value that is eventually cancelled out by normalization (see Eq. (7)). The message itself is a function of the receiving variable's state, which in a discrete setting can be represented by a vector (hence the bracket notation $m_{i \rightarrow j}[x_j]$ for indexing by the state of the variable).

The estimated unary marginals, i.e., the *beliefs*, are computed by fusing the incoming messages with the potential function and normalizing:

$$\Pr(x_j) \approx \exp(b_j[x_j]) = \exp \left(\phi_j(x_j) + \sum_{i \in N_j} m_{i \rightarrow j}[x_j] - z_j \right), \quad (6)$$

where z is a normalizing constant

$$z_j = \log \left(\sum_{x_j} \exp \left(\phi_j(x_j) + \sum_{i \in N_j} m_{i \rightarrow j}[x_j] \right) \right). \quad (7)$$

For convenience and numerical stability, we will consider the log-beliefs

$$b_j[x_j] = \phi_j(x_j) + \sum_{i \in N_j} m_{i \rightarrow j}[x_j] - z_j. \quad (8)$$

We again use bracket notation for indexing into the beliefs because, for discrete variables, the beliefs can most naturally be stored in a lookup table, which can be viewed as a vector.

2.2 SIMPLIFICATIONS

The message update in Eq. (4) contains an expression that is nearly identical to the belief definition in Eq. (8). In the message update, the exponent is

$$a_{x_i} = \phi_{ij}(x_i, x_j) + \phi_i(x_i) + \sum_{k \in N_i \setminus j} m_{k \rightarrow i}[x_i] - d_{ij}. \quad (9)$$

The only difference between this expression and the belief update is that the belief uses the constant z_j to ensure normalization and the message update omits the message

from the receiving variable (j). For finite message values, we can use the equivalence

$$\begin{aligned} \phi_i(x_i) + \sum_{k \in N_i \setminus j} m_{k \rightarrow i}[x_i] - d_{ij} &= \\ \phi_i(x_i) + \sum_{k \in N_i} m_{j \rightarrow i}[x_i] - m_{j \rightarrow i}[x_i] - d_{ij} &= \\ b_i[x_i] - m_{j \rightarrow i}[x_i], \end{aligned} \quad (10)$$

where the last equality sets $d_{ij} = z_i$. The message and belief updates can then be written as

$$\begin{aligned} b_j[x_j] &= \phi_j(x_j) + \sum_{i \in N_j} m_{i \rightarrow j}[x_j] - z_j, \\ m_{i \rightarrow j}[x_j] &= \\ \log \left(\sum_{x_i} \exp(\phi_{ij}(x_i, x_j) + b_i[x_i] - m_{j \rightarrow i}[x_i]) \right). \end{aligned} \quad (11)$$

These two update equations repeat for all variables and edges, and when implemented in low-level languages optimized by pipelining compilers, yield the fastest computer programs that can run belief propagation. However, the indexing requires reasoning about the graph structure, making any code that implements such updates cumbersome to maintain, prone to errors, and difficult to adapt to new computing environments such as parallel computing settings.

3 SPARSE-MATRIX BELIEF PROPAGATION

In this section, we describe sparse-matrix belief propagation and analyze its complexity. Instead of directly implementing the updates, sparse-matrix belief propagation uses an abstraction layer of matrices and tensors. This abstraction allows belief propagation to be written in high-level languages such as Python and MATLAB, delegating the majority of computation into highly optimized linear algebra libraries.

3.1 TENSOR REPRESENTATIONS OF THE MRF, BELIEFS, AND MESSAGES

We store a c by n belief matrix \mathbf{B} , where n is the number of variables and c is the maximum number of states of any variable. For simplicity, assume all variables have the cardinality c , i.e., $|\mathbb{X}_i| = c$. This belief matrix is simply the stacked belief vectors, such that $B_{ij} = b_j[i]$. Similarly, we rearrange this matrix into an analogously shaped and ordered unary potential function matrix Φ , where $\Phi_{ij} = \phi_j(x_j = i)$.

We store the pairwise potentials in a three-dimensional tensor Γ of size $c \times c \times |E|$. Each k th slice of the tensor

is a matrix representing the k th edge potential function as a table, i.e., $\Gamma_{ijk} = \phi_{s_k t_k}(i, j)$.

Consider a view E of the edge set \mathcal{E} in which each edge appears in forward and reverse order, i.e., $(i, j) \in E$ if and only if $(j, i) \in E$. Let the edges also be indexed in an arbitrary but fixed order

$$E = \{(s_1, t_1), (s_2, t_2), \dots, (s_{|E|}, t_{|E|})\}, \quad (12)$$

and define vectors $\mathbf{s} = [s_1, \dots, s_{|E|}]^\top$ and $\mathbf{t} = [t_1, \dots, t_{|E|}]^\top$. These variables encode a fixed ordering for the messages. The particular order does not matter, but to convert message passing into matrix operations, we need the order to be fixed.

In addition to the fixed order, we also define a *message reversal* order vector \mathbf{r} where

$$E[r_i] = (t, s) \text{ if } E[i] = (s, t). \quad (13)$$

We can represent this reversal vector as a sparse, $|E| \times |E|$ permutation matrix \mathbf{R} where $R_{ij} = 1$ iff $r_i = j$.

Define a c by $|E|$ message matrix \mathbf{M} whose i th column is the i th message. In other words, $M_{ij} = m_{s_j \rightarrow t_j}[i]$, or equivalently, \mathbf{M} is a horizontal stack of all messages:

$$\mathbf{M} = [m_{s_1 \rightarrow t_1}, m_{s_2 \rightarrow t_2}, \dots, m_{s_{|E|} \rightarrow t_{|E|}}]. \quad (14)$$

Finally, we define a binary *sparse* matrix \mathbf{T} (for *to*) with $|E|$ rows and n columns whose nonzero entries are as follows:

$$T_{ij} = \begin{cases} 1.0 & \text{if } t_i = j \\ 0 & \text{otherwise.} \end{cases} \quad (15)$$

This matrix \mathbf{T} maps the ordered messages to the variables that receive the messages.

We define an analogous matrix \mathbf{F} (for *from*), also binary and sparse with $|E|$ rows and n columns, whose nonzero entries are as follows:

$$F_{ij} = \begin{cases} 1.0 & \text{if } s_i = j \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

This matrix \mathbf{T} maps the ordered messages to the variables that *send* the messages.

3.1.1 The logsumexp operation

A common operation that occurs in various steps of computing log-space belief propagation is the logsumexp operation, which is defined as follows:

$$\text{logsumexp}(\mathbf{A}) = \log(\mathbf{1}^\top \exp(\mathbf{A})), \quad (17)$$

where the log and exp operations are performed element-wise, and we use the matrix-vector product with the ones

vector ($\mathbf{1}$) as a compact notation for summing across the rows of the exponentiated input matrix.¹ The resulting output is a column vector with the same number of rows as the input matrix \mathbf{A} .

3.1.2 Slice Indexing

To describe the message updates in a compact notation, we use slice indexing, which is common in matrix and tensor software because it can be implemented efficiently in linear time and easy to parallelize. We borrow syntax from `numpy` that is also reminiscent of the famous MATLAB syntax, where

$$\mathbf{A}[:, \mathbf{i}] = [\mathbf{A}[:, i_1], \mathbf{A}[:, i_2], \dots]. \quad (18)$$

This slice indexing allows the reordering, selection, or repetition of the rows or columns of matrices or tensors.

3.2 BELIEF PROPAGATION AS TENSOR OPERATIONS

Using these constructed matrices, the belief matrix is updated with the operations

$$\begin{aligned} \tilde{\mathbf{B}} &\leftarrow \Phi + \mathbf{M}\mathbf{T} \\ \mathbf{B} &\leftarrow \tilde{\mathbf{B}} - \mathbf{1} \text{logsumexp}(\tilde{\mathbf{B}}), \end{aligned} \quad (19)$$

where the last operation uses the logsumexp operation to compute the normalizing constants of each belief column vector and multiplies by the ones vector ($\mathbf{1}$) to broadcast it across all rows.²

The message matrix \mathbf{M} is updated with the following formula:

$$\mathbf{M} \leftarrow \text{logsumexp}(\mathbf{\Gamma} + \mathbf{B}[:, \mathbf{s}] - \mathbf{M}[:, \mathbf{r}]). \quad (20)$$

This expression uses two forms of shorthand that require further explanation. First, the addition of the tensor $\mathbf{\Gamma}$ and the matrix $(\mathbf{B}[:, \mathbf{s}] - \mathbf{M}[:, \mathbf{r}])$ requires broadcasting. The tensor $\mathbf{\Gamma}$ is of size $c \times c \times |E|$, and the matrix $(\mathbf{B}[:, \mathbf{s}] - \mathbf{M}[:, \mathbf{r}])$ is of size $c \times |E|$. The broadcasting copies the matrix c times and stacks them as rows to form the same shape as $\mathbf{\Gamma}$. Second, the logsumexp operation sums across the columns of the summed tensor, outputting

¹It is especially useful to form this abstraction because this operation is notoriously unstable in real floating-point arithmetic. Numerically stable implementations that adjust the exponent by subtracting a constant and adding the constant back to the output of the logarithm are possible. These stable implementations add a linear-time overhead to the otherwise linear-time operation, so they maintain the same asymptotic running time of the original logsumexp operation.

²In `numpy`, this broadcasting is automatically inferred from the size of the matrices being subtracted.

a tensor of shape $c \times 1 \times |E|$, which is then squeezed into a matrix of size $c \times |E|$.

The message matrix can equivalently be updated with this formula:

$$\mathbf{M} \leftarrow \text{logsumexp}(\mathbf{\Gamma} + \mathbf{B}\mathbf{F}^\top - \mathbf{M}\mathbf{R}). \quad (21)$$

Here $\mathbf{B}\mathbf{F}^\top - \mathbf{M}\mathbf{R}$ is once again $c \times |E|$, and it is equivalent to the slice-notation form above.

Belief propagation is then implemented by iteratively running Eq. (19) and then either of the equivalent Eqs. (20) and (21).

3.3 VARIATIONS OF BELIEF PROPAGATION

Many variations of belief propagation can similarly be converted into a sparse-matrix format. We describe some of these variations here.

3.3.1 Tree-Reweighted Belief Propagation

The tree-reweighted variation of belief propagation (TRBP) computes messages corresponding to a convex combination of spanning trees over the input graph. The result is a procedure that optimizes a convex inference objective (Wainwright et al., 2003; Wainwright et al., 2008). The belief and message updates for TRBP are adjusted according to edge-appearance probabilities in a distribution of spanning trees over the MRF graph. These updates can be implemented in matrix form by using a length $|E|$ vector $\boldsymbol{\rho}$ containing the appearance probabilities ordered according to edge set E . The matrix-form updates for the beliefs and messages are then

$$\begin{aligned} \tilde{\mathbf{B}} &\leftarrow \mathbf{\Phi} + (\boldsymbol{\rho} \circ \mathbf{M})\mathbf{T} \\ \mathbf{B} &\leftarrow \tilde{\mathbf{B}} - \mathbf{1} \text{logsumexp}(\tilde{\mathbf{B}}), \\ \mathbf{M} &\leftarrow \text{logsumexp}(\mathbf{\Gamma}/\boldsymbol{\rho} + \mathbf{B}\mathbf{F}^\top - \mathbf{M}\mathbf{R}), \end{aligned} \quad (22)$$

where element-wise product \circ and element-wise division $/$ are applied with appropriate broadcasting.

3.3.2 Convexified Belief Propagation

Another important variation of loopy belief propagation uses counting numbers to adjust the weighting of terms in the factorized entropy approximation. The resulting message update formulas weight each marginal by these counting numbers. Under certain conditions, such as when all counting numbers are non-negative, the inference objective can be shown to be concave, so this method is often referred to as *convexified Bethe belief propagation* (Meshi et al., 2009). We can exactly mimic the message and belief update formulas for convexified belief propagation by instantiating a vector \mathbf{c} containing the counting

numbers of each edge factor, resulting in the updates

$$\begin{aligned} \tilde{\mathbf{B}} &\leftarrow \mathbf{\Phi} + \mathbf{M}\mathbf{T} \\ \mathbf{B} &\leftarrow \tilde{\mathbf{B}} - \mathbf{1} \text{logsumexp}(\tilde{\mathbf{B}}), \\ \mathbf{M} &\leftarrow \text{logsumexp}(\mathbf{\Gamma} + (\mathbf{B}\mathbf{F}^\top - \mathbf{M}\mathbf{R})/\mathbf{c}) \circ \mathbf{c}. \end{aligned} \quad (23)$$

The counting numbers for unary factors can be used to compute the inference objective, but they do not appear in the message-passing updates.

3.3.3 Max-Product Belief Propagation

Finally, we illustrate that sparse tensor operations can be used to conduct approximate *maximum a posteriori* (MAP) inference. The max-product belief propagation algorithm (Wainwright et al., 2008) is one method for approximating MAP inference, and it can be implemented with the following updates:

$$\begin{aligned} \mathbf{B} &\leftarrow \text{onehotmax}(\mathbf{\Phi} + \mathbf{M}\mathbf{T}) \\ \mathbf{M} &\leftarrow \text{logsumexp}(\mathbf{\Gamma} + \mathbf{B}\mathbf{F}^\top - \mathbf{M}\mathbf{R}), \end{aligned} \quad (24)$$

where onehotmax is a function that returns an indicator vector with 1 for entries that are the maximum of each column and zeros everywhere else, e.g., the “one-hot” encoding. Similar conversions are also possible for variations of max-product, such as max-product linear programming (Globerson and Jaakkola, 2008).

3.4 TIME-COMPLEXITY ANALYSIS

To analyze our sparse-matrix formulation of loopy belief propagation, and to show that it requires an asymptotically equivalent running time to normal loopy belief propagation, we first revisit the sizes of all matrices involved in the update equations. The first step is the belief update operation, Eq. (19), which updates the c by n belief matrix \mathbf{B} . The potential matrix $\mathbf{\Phi}$ is also c by n ; the message matrix \mathbf{M}^\top is c by $|E|$; and the *sparse* message-recipient indexing matrix \mathbf{T} is $|E|$ by n .

The second step in Eq. (19) normalizes the beliefs. It subtracts from \mathbf{B} the product of $\mathbf{1}$, which is a c by 1 vector, and $\text{logsumexp}(\tilde{\mathbf{B}})$, which is 1 by n . Explicitly writing the numerically stable logsumexp operation, the right side of this line can be expanded to

$$\tilde{\mathbf{B}} - \mathbf{1} \log(\text{sum}(\exp(\mathbf{B} - \max(\mathbf{B}))) + \max(\mathbf{B})). \quad (25)$$

We next examine the message update Eq. (21), which updates \mathbf{M} . The three-dimensional tensor $\mathbf{\Gamma}$ is of size $c \times c \times |E|$; the *sparse* message-sender indexing matrix \mathbf{F}^\top is c by $|E|$; and the *sparse* reversal permutation matrix \mathbf{R} is $|E| \times |E|$. The message matrix \mathbf{M} is c by $|E|$.

CPU computation These three update steps are the entirety of each belief propagation iteration. From the first line of Eq. (19), the main operation to analyze is the dense-sparse matrix multiplication **MT**. Considering an $n \times m$ dense matrix A and a sparse matrix B of size $m \times p$ with s nonzero entries (i.e., $\|B\|_0 = s$), the sparse dot product has time complexity $O(ms)$ in sequential computation, as on a CPU. The time complexity of the sparse dot product depends upon the number of rows m and the number of sparse elements in the sparse matrix B . Every other computation in Eq. (19) involves element-wise operations on c by n matrices. Thus, Eq. (19) requires $O(nc + \|\mathbf{T}\|_0 c)$ time. Since the sparse indexing matrix \mathbf{T} is defined to have a single nonzero entry per column, corresponding to edges, the time complexity of this step is $O(nc + |E|c)$.

In the message-update step, Eq. (21), the outer logsumexp operation and the additions involve element-wise operations over $c \times c \times |E|$ tensors. The matrix multiplications are all dense-sparse dot products, so the total cost of Eq. (21) is $O(|E|c^2 + \|\mathbf{F}\|_0 c + \|\mathbf{R}\|_0 c)$ (Gilbert et al., 1992). Since both indexing matrices \mathbf{F} and \mathbf{R} have one entry per edge, the total time complexity of the message update is $O(|E|c^2 + |E|c)$.

The combined computational cost of both steps is $O(nc + |E|c + |E|c^2)$. This iteration cost is the same as traditional belief propagation.

GPU computation Since the matrix operations in our method are simple, they are easily parallelized on GPUs. Ignoring the overhead of transferring data to GPU memory, we focus on the time complexity of the message passing. First consider the dense-sparse matrix multiplication mentioned previously, with a dense n by m matrix A and sparse m by p matrix B with s nonzero entries. GPU algorithms for sparse dot products use all available cores to run threads of matrix operations (Bell and Garland, 2008). In this case, each thread can run the multiplication operation of a single column in the sparse matrix B .

Given k cores/threads, we assume that there will be two cases: (1) when the number of sparse columns m is *less than or equal to* the number of cores k and (2) when the number of sparse columns m is *more than* the number of cores k . For either case, let s_i be the number of nonzero entries in column i . The time complexity of case (1) is $O(\max_i s_i)$, which is the time needed to process whichever column requires the most multiplications. For case (2), the complexity is $O(\lceil \frac{m}{k} \rceil \max_i s_i)$, which is the time for each set of cores to process k columns. In case (2), we are limited by our largest columns, as the rest of our smaller columns will be processed much sooner. Overall, the GPU time complex-

ity of this operation is $O\left(\max\left(\max_i s_i, \lceil \frac{m}{k} \rceil \max_i s_i\right)\right)$. In our sparse indexing matrices, each column has at most one element.

For the element-wise dense matrix operations, which have time complexity $O(nm)$ on the CPU, we can again multi-thread each entry over the number of cores k in our GPU such that the time complexity is $O(\lceil \frac{nm}{k} \rceil)$.

The running time of the belief propagation steps is then $O\left(\lceil \frac{n}{k} \rceil c + \lceil \frac{|E|}{k} \rceil c + \lceil \frac{|E|c^2}{k} \rceil\right)$.

Based on our parallelism analysis, we expect significantly faster running times when running on the GPU, especially in cases where we have a large number of cores k . While we expect some time loss due to data-transfer overhead to the GPU, this overhead may be negligible when considering the overall time cost for every iteration of the message-passing matrix operations.

Finally, this analysis also applies to other shared-memory, multi-core, multithreaded environments (e.g., on CPUs), since in both CPU and GPU settings, matrix rows can be independently processed.

4 EMPIRICAL EVALUATION

In this section, we describe our empirical evaluation of sparse-matrix belief propagation. We measure the running time for belief propagation on a variety of MRFs using different software implementations and hardware, including optimized and compiled code for CPU-based computation and sparse-matrix belief propagation in a high-level language for both CPU- and GPU-based computation.

4.1 EXPERIMENTAL SETUP

We generate grid MRFs of varying sizes. We randomly generate potential functions for each MRF such that the log potentials are independently normally distributed with variance 1.0. We use MRFs with different variable cardinalities c from the set $\{8, 16, 32, 64\}$. We run experiments with MRFs structured as square, two-dimensional grids, where the number of rows and columns in the grids are $\{8, 16, 32, 64, 128, 256, 512\}$. In other words, the number of variables in models with these grid sizes are, respectively, 64, 256, 1024, 4,096, 16,384, 64,536, and 262,144. We run all implementations of belief propagation until the total absolute change in the messages is less than 10^{-8} .

We run our experiments on different hardware setups. We use two different multi-core CPUs: a 2.4 Ghz Intel i7 with 4 cores and a 4 Ghz Intel i7 with 4 cores.

We also run sparse-matrix belief propagation on various

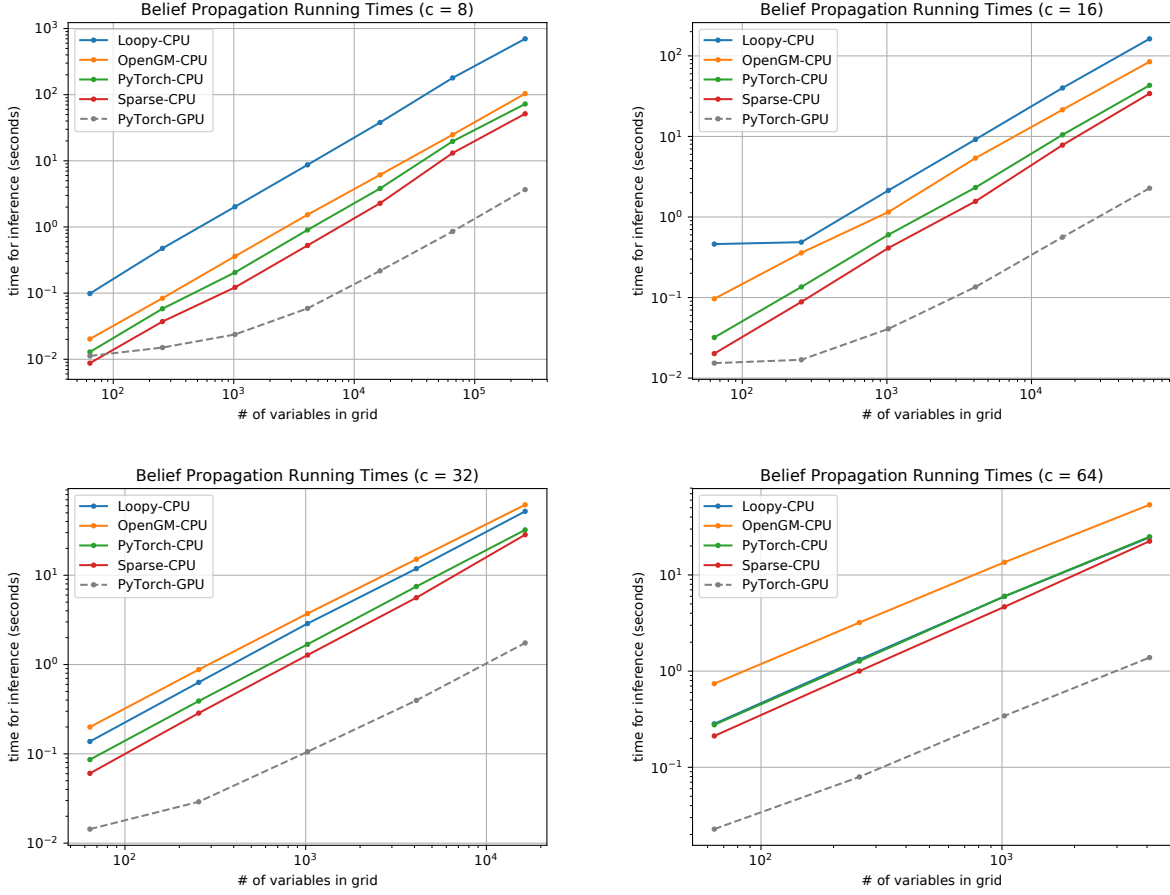


Figure 1: Log-log plots of belief propagation running times for four different implementations on the CPU. Each plot shows the results for different variable cardinalities c . OpenGM refers to the compiled C++ library, Loopy refers to the direct Python implementation. Sparse uses implements sparse-matrix belief propagation with `scipy.sparse`. And PyTorch implements it with the PyTorch library. We also plot the running time using PyTorch on the least powerful GPU we tested (Nvidia GTX780M) for comparison. The CPU runs plotted here use a 2.4 Ghz Intel i7.

GPUs. We run on an Nvidia GTX 780M (1,536 cores, 4 GB memory), an Nvidia GTX 1080 (2,560 cores, 8 GB), an Nvidia GTX 1080Ti (3,584 cores, 11 GB), and an Nvidia Tesla P40 (3840 cores, 24 GB).

Additional experiments will be available in this paper’s supplemental material.

4.2 IMPLEMENTATION DETAILS

We compare four different implementations of belief propagation. First, we use the compiled and optimized C++ implementation of belief propagation in the OpenGM library (Andres et al., 2012). This software represents the low-level implementation. Second, we use a direct implementation of simplified belief propagation (see Section 2.2) in Python and `numpy` using Python loops and dictionaries (hash maps) to manage indexing over graph structure.

Third, we use an implementation of sparse-matrix belief propagation in Python using `scipy.sparse`. Fourth, we use an implementation of sparse-matrix belief propagation in Python using the deep learning library PyTorch, which enables easily switching between CPU and GPU computation.

4.3 RESULTS AND DISCUSSION

Considering the results for CPU-based belief propagation in Fig. 1, the sparse-matrix belief propagation is faster than any other CPU-based belief propagation for all MRF sizes and all variable cardinalities. Similarly, the curves show a clear linearity, with a slope suggesting that all the CPU-based belief propagation algorithms increase in time complexity at a linear rate. It is also evident that the PyTorch implementation is consistently slower

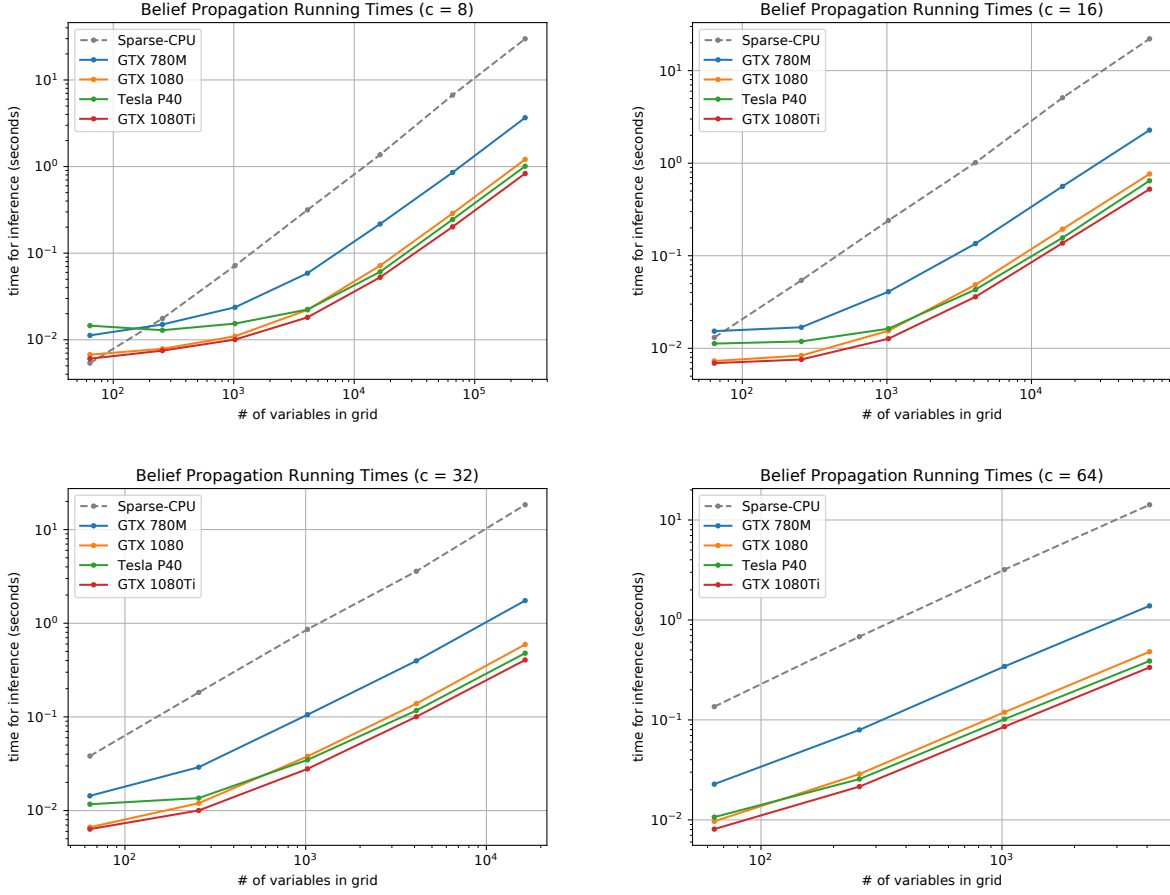


Figure 2: Log-log plots of PyTorch-GPU belief propagation running times for four different GPUs (780M, 1080, Tesla P40, 1080Ti) and the fastest CPU method (Sparse-CPU) with different variable cardinalities c . The CPU is the 4Ghz Intel i7.

than the `scipy.sparse` implementation, which is to be expected because PyTorch operations incur an additional overhead for their ability to integrate with deep learning procedures (e.g., back-propagation and related tasks).

Notably, we can see that the direct Python loop-based implementation is by far the slowest of these options. However, when the variable cardinality increases to large values, the Python implementation nearly matches the speed of sparse-matrix belief propagation with PyTorch on the CPU. While OpenGM’s belief propagation does offer some benefits compared to Python initially at a lower values of c , it actually results in the slowest running times at $c \geq 32$. We can conclude that despite the compiled and optimized C++ code, the number of variables and states can eventually overshadow any speedups initially seen at lower values.

In Fig. 1, we also include the running times for sparse-matrix belief propagation with PyTorch on the GPU—

shown with the dotted gray line. A direct comparison is not exactly fair, since we are comparing across different computing hardware, though these curves were measured on the same physical computer. The comparison makes clear that the GPU offers significant speed increases over any CPU-based belief propagation, even that of the faster `scipy.sparse` implementation. Interestingly, there is a trend with the GPU runtimes that are sub-linear in the log-log plots, representing the cases where the number of sparse columns is not yet more than the number of cores of the GPU.

Examining the results for GPU-based belief propagation in Fig. 2, a majority of the GPUs are fairly close in running time between the three powerful Nvidia units: the 1080, the 1080Ti, and the P40. The 780M understandably lags behind. As seen previously, until the cores are saturated with operations, there appears to be a pseudo-constant time cost. And once the cores are saturated, the running times grow linearly. This trend is best seen at

$c = 16$ for the first two or three points. We also include the fastest CPU running time, using `scipy.sparse` on an Intel 4 Ghz i7, to illustrate the drastic difference in time between sparse-matrix belief propagation on the CPU and GPU.

These results demonstrate that sparse-matrix belief propagation enables the fastest running times for inference in these grid MRFs on the CPU. And using different software backends (`scipy.sparse` or PyTorch) for the sparse-matrix operations leads to different behavior, with PyTorch incurring some overhead resulting in slower computation. Once ported to the GPU, the speedups are even more drastic, resulting in running times that are many factors faster than those seen on the CPU, easily outweighing the overhead cost of using software backends like PyTorch that support seamless switching from CPU to GPU computation.

5 CONCLUSION

We presented sparse-matrix belief propagation, which exactly reformulates loopy belief propagation (and its variants) as a series of matrix and tensor operations. This reformulation creates an abstract interface between belief propagation and a variety of highly optimized libraries for sparse-matrix and tensor operations. We demonstrated how sparse-matrix belief propagation scales as efficiently as low-level, compiled and optimized implementations, yet it can be implemented in high-level mathematical programming languages. We also demonstrated how the abstraction layer allows easy portability to advanced computing hardware by running sparse-matrix belief propagation on GPUs. The immediately resulting parallelization benefits required little effort once the sparse-matrix abstraction was in place. Our software library with these implementations is available at <https://bitbucket.org/berthuang/mrftools/>.

Open Questions and Next Steps There are still a number of research directions that we would like to pursue. We mainly focused on analyzing the benefits of sparse-matrix belief propagation on grid-based MRFs, but there are many different structures of MRFs used in important applications (chain models, random graphs, graphs based on structures of real networks). Similarly, applying our approach to real-world examples would help confirm the utility of it in current problems within machine learning. Likewise, we focused on fairly sparse graphs that did not have many non-zero entries per column. It would be interesting to explore the difference between how sparse-matrix belief propagation behaves on the dense and sparse matrices on different hardware and whether fully dense matrices would still result in notable speed improvements,

or if overhead from the sparse-matrix format would become a bottleneck.

Our sparse-matrix formulation of belief propagation is derived for pairwise MRFs, so it remains an open question what modifications are necessary for higher-order MRFs which may have arbitrarily large factors. Benefits similar to those we measured on GPUs can arise from the sparse-matrix abstraction for other computing backends, such as the use of compute-cluster parallelism through libraries such as Apache Spark, or the computation of belief propagation on FPGAs. Finally, the integration of belief propagation into deep learning models is straightforward with the matrix abstraction. Though many popular frameworks do not yet support back-propagation through sparse dot products, such support is forthcoming according to their respective developer communities.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. (2016). Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283.
- Andres, B., Beier, T., and Kappes, J. H. (2012). Opengm: A C++ library for discrete graphical models. *CoRR*, abs/1206.0111.
- Bell, N. and Garland, M. (2008). Efficient sparse matrix-vector multiplication on CUDA. Technical report, Nvidia Technical Report NVR-2008-004, Nvidia Corporation.
- Bosagh Zadeh, R., Meng, X., Ulanov, A., Yavuz, B., Pu, L., Venkataraman, S., Sparks, E., Staple, A., and Zaharia, M. (2016). Matrix computations and optimization in Apache Spark. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 31–38. ACM.
- Domke, J. (2013). Learning graphical model parameters with approximate marginal inference. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(10):2454–2467.
- Gatterbauer, W., Günnemann, S., Koutra, D., and Faloutsos, C. (2015). Linearized and single-pass belief propagation. In *Proceedings of the VLDB Endowment*, volume 8, pages 581–592. VLDB Endowment.
- Gilbert, J. R., Moler, C., and Schreiber, R. (1992). Sparse matrices in MATLAB: Design and implementation. *SIAM Journal on Matrix Analysis and Applications*, 13(1):333–356.

- Globerson, A. and Jaakkola, T. S. (2008). Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. In *Advances in Neural Information Processing Systems*, pages 553–560.
- Low, Y., Gonzalez, J. E., Kyrola, A., Bickson, D., Guestrin, C. E., and Hellerstein, J. (2014). Graphlab: A new framework for parallel machine learning. *arXiv preprint arXiv:1408.2041*.
- Malewicz, G., Austern, M. H., Bik, A. J., Dehnert, J. C., Horn, I., Leiser, N., and Czajkowski, G. (2010). Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 135–146. ACM.
- Meshi, O., Jaimovich, A., Globerson, A., and Friedman, N. (2009). Convexifying the Bethe free energy. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 402–410. AUAI Press.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. (2017). Automatic differentiation in PyTorch. In *NIPS Workshop on Autodiff*.
- Pearl, J. (2014). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann.
- Schmidt, M. (2007). UGM: A Matlab toolbox for probabilistic undirected graphical models.
- Schwing, A., Hazan, T., Pollefeys, M., and Urtasun, R. (2011). Distributed message passing for large scale graphical models. In *Computer Vision and Pattern Recognition*.
- Wainwright, M. J., Jaakkola, T. S., and Willsky, A. S. (2003). Tree-reweighted belief propagation algorithms and approximate ml estimation by pseudo-moment matching. In *International Conference in Artificial Intelligence and Statistics*.
- Wainwright, M. J., Jordan, M. I., et al. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305.
- Zheng, L., Mengshoel, O., and Chong, J. (2012). Belief propagation by message passing in junction trees: Computing each message faster using GPU parallelization. *arXiv preprint arXiv:1202.3777*.
- Zhuo, L. and Prasanna, V. K. (2005). Sparse matrix-vector multiplication on fpgas. In *Proceedings of the 2005 ACM/SIGDA 13th International Symposium on Field-Programmable Gate Arrays*, pages 63–74. ACM.