

TensorInference: A Julia package for tensor-based probabilistic inference

Martin Roa-Villescas^{1*} and Jin-Guo Liu^{2*}

¹ Eindhoven University of Technology ² Hong Kong University of Science and Technology (Guangzhou)

* These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Open Journals](#) ↗

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#)).

Summary

TensorInference.jl is a Julia ([Bezanson et al., 2017](#)) package designed for performing probabilistic inference in discrete graphical models. Capitalizing on the recent advances in the field of tensor networks ([Orús, 2014, 2019](#); [Robeva & Seigal, 2019](#)), TensorInference.jl offers high-performance solutions for prevalent inference problems. Specifically, it provides methods to:

1. calculate the partition function (also known as the probability of evidence).
2. compute the marginal probability distribution over each variable given evidence.
3. find the most likely assignment to all variables given evidence.
4. find the most likely assignment to a set of query variables after marginalizing out the remaining variables.
5. draw samples from the posterior distribution given evidence ([Cheng et al., 2019](#); [Han et al., 2018](#)).

The use of a tensor network-based infrastructure ([Fishman et al., 2022](#); [Jutho et al., 2023](#)) offers several advantages when dealing with complex computational tasks. Firstly, it simplifies the process of computing gradients by employing differentiable programming ([Liao et al., 2019](#)), a critical operation for the aforementioned inference tasks. Secondly, it supports generic element types without a significant compromise on performance. The advantage of supporting generic element types lies in the ability to solve a variety of problems using the same tensor network contraction algorithm, simply by varying the element types used. This flexibility has allowed us to seamlessly implement solutions for several of the inference tasks described above ([Jin Guo Liu et al., 2022](#); [Jin-Guo Liu et al., 2021](#)). Thirdly, it allows users to define a hyper-optimized contraction order, which is known to have a significant impact on the computational performance of contracting tensor networks ([Gao et al., 2021](#); [Markov & Shi, 2008](#); [Pan & Zhang, 2022](#)). TensorInference.jl provides a predefined set of state-of-the-art contraction ordering methods. These methods include a *local search based method* (TreeSA) ([Kalachev et al., 2022](#)), two *min-cut based methods* (KaHyParBipartite) ([Gray & Kourtis, 2021](#)) and (SABipartite), and a *greedy method* (GreedyMethod). Finally, tensor networks — and by extension, TensorInference.jl — harness the latest developments in computational technology, including a highly optimized set of BLAS routines ([Blackford et al., 2002](#)) and GPU technology.

TensorInference.jl succeeds JunctionTrees.jl ([Roa-Villescas et al., 2022, 2023](#)), a Julia package implementing the Junction Tree Algorithm (JTA) ([Jensen et al., 1990](#); [Lauritzen & Spiegelhalter, 1988](#)). While the latter employs tensor-based technology to optimize the computation of individual sum-product messages within the JTA context, TensorInference.jl takes a different route. It adopts a holistic tensor network approach, which opens new doors for optimization opportunities, and significantly reduces the algorithm's complexity compared to the JTA.

44 **Statement of need**

45 A major challenge in developing intelligent systems is the ability to reason under uncertainty, a
46 challenge that appears in many real-world problems across various domains, including artificial
47 intelligence, medical diagnosis, computer vision, computational biology, and natural language
48 processing. Reasoning under uncertainty involves calculating the probabilities of relevant
49 variables while taking into account any information that is acquired. This process, which can
50 be thought of as drawing global insights from local observations, is known as *probabilistic*
51 *inference*.

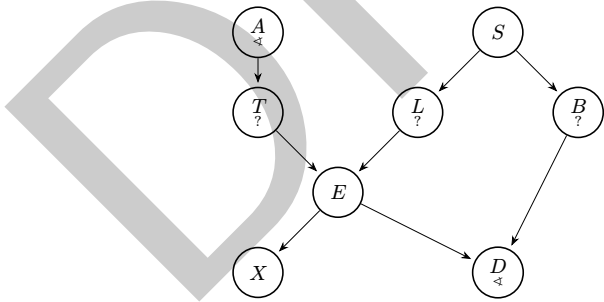
52 *Probabilistic graphical models* (PGMs) provide a unified framework to perform probabilistic
53 inference. These models use graphs to represent the joint probability distribution of complex
54 systems in a concise manner by exploiting the conditional independence between variables in
55 the model. Additionally, they form the foundation for various algorithms that enable efficient
56 probabilistic inference.

57 However, even with the representational aid of PGMs, performing probabilistic inference remains
58 an intractable endeavor on many real-world models. The reason is that performing probabilistic
59 inference involves complex combinatorial optimization problems in very high dimensional spaces.
60 To tackle these challenges, more efficient and scalable inference algorithms are needed.

61 As an attempt to tackle the aforementioned challenges, we present TensorInference.jl, a
62 Julia package for probabilistic inference that combines the representational capabilities of
63 PGMs with the computational power of tensor networks. By harnessing the best of both worlds,
64 TensorInference.jl aims to enhance the performance of probabilistic inference, thereby
65 expanding the tractability spectrum of exact inference for more complex, real-world models.

66 **Usage example**

67 The graph below corresponds to the *ASIA network* (Lauritzen & Spiegelhalter, 1988), a simple
68 Bayesian network (Pearl, 1985) used extensively in educational settings. It describes the
69 probabilistic relationships between different random variables which correspond to possible
70 diseases, symptoms, risk factors and test results.



Random variable	Meaning
A	Recent trip to Asia
T	Patient has tuberculosis
S	Patient is a smoker
L	Patient has lung cancer
B	Patient has bronchitis
E	Patient has T and/or L
X	Chest X-Ray is positive
D	Patient has dyspnoea

Figure 1: The ASIA network: a simplified example of a Bayesian network from the context of medical diagnosis (Lauritzen & Spiegelhalter, 1988).

71 In the example, a patient has recently visited Asia and is now experiencing dyspnea. These
72 conditions serve as the evidence for the observed variables (*A* and *D*). The doctor's task is to
73 assess the likelihood of various diseases — tuberculosis, lung cancer, and bronchitis — which
74 constitute the query variables in this scenario (*T*, *L*, and *B*).

75 We now demonstrate how to use TensorInference.jl for conducting a variety of inference
76 tasks on this toy example. Please note that as the API may evolve, we recommend checking

77 the [examples](#) directory of the official `TensorInference.jl` repository for the most up-to-date
78 version of this example.

```
# Import the TensorInference package, which provides the functionality needed
# for working with tensor networks and probabilistic graphical models.
In [1]: using TensorInference

# Load the ASIA network model from the `asia.uai` file located in the examples
# directory. Refer to the documentation of this package for a description of the
# format of this file.
instance = read_instance(pkgdir(TensorInference, "examples", "asia", "asia.uai"))

# Create a tensor network representation of the loaded model.
tn = TensorNetworkModel(instance)

Out [1]: TensorNetworkModel{Int64, OMEinsum.DynamicNestedEinsum{Int64}, Array{Float64}}
variables: 1, 2, 3, 4, 5, 6, 7, 8
contraction time = 2^6.044, space = 2^2.0, read-write = 2^7.098

# Calculate the log10 partition function.
In [2]: probability(tn) |> first |> log10

Out [2]: 0.0

# Calculate the marginal probabilities of each random variable in the model.
In [3]: marginals(tn)

Out [3]: 8-element Vector{Vector{Float64}}:
 [0.01, 0.99]
 [0.0104, 0.9895999999999999]
 [0.5, 0.49999999999999994]
 [0.0550000000000000014, 0.94500000000000001]
 [0.44999999999999996, 0.5499999999999999]
 [0.064828000000000002, 0.93517200000000001]
 [0.110290040000000002, 0.88970996]
 [0.435970600000000004, 0.5640294]

# Set the evidence: We assume that the "X-ray" result (variable 7) is positive.
In [4]: set_evidence!(instance, 7 => 0)

# Since setting the evidence may affect the contraction order of the tensor
# network, we need to recompute it.
tn = TensorNetworkModel(instance)

# Calculate the maximum log-probability among all configurations.
maximum_logp(tn)

Out [4]: 0-dimensional Array{Float64, 0}:
 -3.6522217920023303

# Generate 10 samples from the probability distribution represented by the model.
In [5]: sample(tn, 10)

Out [5]: 8×10 Matrix{Int64}:
 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 0
 1 0 0 0 1 0 0 0 0 0
 1 1 0 0 1 1 0 0 0 1
```

```

1 0 1 0 0 1 1 0 0 0
1 1 0 0 1 1 0 0 0 0
0 0 0 0 0 0 0 0 0 0
1 0 1 1 0 0 1 0 0 0

```

```

# Retrieve both the maximum log-probability and the most probable configuration.
# In this configuration, the most likely outcomes are that the
# patient smokes (variable 3) and has lung cancer (variable 4).
In [6]: logp, cfg = most_probable_config(tn)

```

```

Out [6]: (-3.6522217920023303, [1, 1, 0, 0, 0, 0, 0, 0])

```

```

# Compute the most probable values of certain variables (e.g., 4 and 7) while
# marginalizing over others. This is known as Maximum a Posteriori (MAP)
# estimation.

```

```

In [7]: set_query!(instance, [4, 7])
mmap = MMAPModel(instance)

```

```

Out [7]: MMAPModel{Int64, Array{Float64}}
variables: 4, 7 (evidence → 0)
query variables: [[1, 2, 6, 5, 3, 8]]
contraction time = 2^6.022, space = 2^2.0, read-write = 2^7.033

```

```

# Get the most probable configurations for variables 4 and 7.
In [8]: most_probable_config(mmap)

```

```

Out [8]: (-2.8754627318176693, [1, 0])

```

```

# Compute the total log-probability of having lung cancer. The results suggest
# that the probability is roughly half.

```

```

In [9]: log_probability(mmap, [1, 0]), log_probability(mmap, [0, 0])

```

```

Out [9]: (-2.8754627318176693, -2.9206248010671856)

```

Acknowledgments

This work is partially funded by the Netherlands Organization for Scientific Research and the Guangzhou Municipal Science and Technology Project (No. 2023A03J0003). We extend our gratitude to Madelyn Cain and Patrick Wijnings for their insightful discussions on the intersection of tensor networks and probabilistic graphical models.

References

- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. (2017). Julia: A Fresh Approach to Numerical Computing. *SIAM Review*, 59(1), 65–98. <https://doi.org/10.1137/141000671>
- Blackford, L. S., Petitet, A., Pozo, R., Remington, K., Whaley, R. C., Demmel, J., Dongarra, J., Duff, I., Hammarling, S., Henry, G., & others. (2002). An updated set of basic linear algebra subprograms (BLAS). *ACM Transactions on Mathematical Software*, 28(2), 135–151. <https://doi.org/10.1145/567806.567807>
- Cheng, S., Wang, L., Xiang, T., & Zhang, P. (2019). Tree tensor networks for generative modeling. *Physical Review B*, 99(15), 155131. <https://doi.org/10.1103/PhysRevB.99.155131>
- Fishman, M., White, S., & Stoudenmire, E. (2022). The ITensor software library for tensor network calculations. *SciPost Physics Codebases*, 004. <https://doi.org/10.21468/SciPostPhysCodeb.4>

- 97 Gao, X., Kalinowski, M., Chou, C.-N., Lukin, M. D., Barak, B., & Choi, S. (2021). *Limitations*
98 *of linear cross-entropy as a measure for quantum advantage*. [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2112.01657)
99 [arXiv.2112.01657](https://doi.org/10.48550/arXiv.2112.01657)
- 100 Gray, J., & Kourtis, S. (2021). Hyper-optimized tensor network contraction. *Quantum*, 5, 410.
101 <https://doi.org/10.22331/q-2021-03-15-410>
- 102 Han, Z.-Y., Wang, J., Fan, H., Wang, L., & Zhang, P. (2018). Unsupervised generative
103 modeling using matrix product states. *Physical Review X*, 8(3), 031012. [https://doi.org/](https://doi.org/10.1103/PhysRevX.8.031012)
104 [10.1103/PhysRevX.8.031012](https://doi.org/10.1103/PhysRevX.8.031012)
- 105 Jensen, F. V., Lauritzen, S. L., & Olesen, K. G. (1990). Bayesian updating in causal probabilistic
106 networks by local computations. *Computational Statistics Quarterly*, 4, 269–282.
- 107 Jutho, Lukas, ho-oto, maartenvd, getzdan, Liu, J., Aluthge, D., Florian, Lyon, S., Morley, A.,
108 Privett, A., Brann, D., louchtchenko, D., Saba, E., Otto, F., Garrison, J., Bhattacharya, J.,
109 Feist, J., TagBot, J., ... jemiryguo. (2023). *Jutho/TensorOperations.jl: v4.0.0* (Version
110 v4.0.0). Zenodo. <https://doi.org/10.5281/zenodo.8166121>
- 111 Kalachev, G., Panteleev, P., & Yung, M.-H. (2022). *Multi-tensor contraction for XEB*
112 *verification of quantum circuits*. <https://doi.org/10.48550/arXiv.2108.05665>
- 113 Lauritzen, S. L., & Spiegelhalter, D. J. (1988). Local computations with probabilities on
114 graphical structures and their application to expert systems. *Journal of the Royal Statistical*
115 *Society: Series B (Methodological)*, 50(2), 157–194.
- 116 Liao, H.-J., Liu, J.-G., Wang, L., & Xiang, T. (2019). Differentiable programming tensor
117 networks. *Physical Review X*, 9(3), 031041. <https://doi.org/10.1103/PhysRevX.9.031041>
- 118 Liu, Jin Guo, Gao, X., Cain, M., Lukin, M. D., & Wang, S.-T. (2022). *Computing solution*
119 *space properties of combinatorial optimization problems via generic tensor networks*. arXiv.
120 <https://doi.org/10.48550/ARXIV.2205.03718>
- 121 Liu, Jin-Guo, Wang, L., & Zhang, P. (2021). Tropical tensor network for ground states of spin
122 glasses. *Physical Review Letters*, 126(9). <https://doi.org/10.1103/physrevlett.126.090506>
- 123 Markov, I. L., & Shi, Y. (2008). Simulating quantum computation by contracting tensor net-
124 works. *SIAM Journal on Computing*, 38(3), 963–981. <https://doi.org/10.1137/050644756>
- 125 Orús, R. (2014). A practical introduction to tensor networks: Matrix product states and
126 projected entangled pair states. *Annals of Physics*, 349, 117–158. [https://doi.org/10.](https://doi.org/10.1016/j.aop.2014.06.013)
127 [1016/j.aop.2014.06.013](https://doi.org/10.1016/j.aop.2014.06.013)
- 128 Orús, R. (2019). Tensor networks for complex quantum systems. *Nature Reviews Physics*,
129 1(9), 538–550. <https://doi.org/10.1038/s42254-019-0086-7>
- 130 Pan, F., & Zhang, P. (2022). Simulation of quantum circuits using the big-batch tensor network
131 method. *Phys. Rev. Lett.*, 128, 030501. <https://doi.org/10.1103/PhysRevLett.128.030501>
- 132 Pearl, J. (1985). Bayesian networks: A model of self-activated memory for evidential reasoning.
133 *Proc. Of Cognitive Science Society (CSS-7)*.
- 134 Roa-Villescas, M., Liu, J. G., Wijnings, P. W. A., Stuijk, S., & Corporaal, H. (2023). Scal-
135 ing probabilistic inference through message contraction optimization. *2023 Congress in*
136 *Computer Science, Computer Engineering, & Applied Computing (CSCE)*.
- 137 Roa-Villescas, M., Wijnings, P. W. A., Stuijk, S., & Corporaal, H. (2022). Partial evaluation in
138 junction trees. *2022 25th Euromicro Conference on Digital System Design (DSD)*, 429–437.
139 <https://doi.org/10.1109/DSD57027.2022.00064>
- 140 Robeva, E., & Seigal, A. (2019). Duality of graphical models and tensor networks. *Information*
141 *and Inference: A Journal of the IMA*, 8(2), 273–288. [https://doi.org/10.1093/imaia/](https://doi.org/10.1093/imaia/iaiy009)
142 [iaiy009](https://doi.org/10.1093/imaia/iaiy009)