Graphical Models are All You Need: Per-interaction reconstruction uncertainties in a dark matter detection experiment

 $\begin{array}{cccc} \textbf{Christina Peters}^{1,2,\dagger} & \textbf{Aaron Higuera}^2 & \textbf{Shixiao Liang}^2 & \textbf{Venkat Roy}^3 \\ \textbf{Waheed U. Bajwa}^{3,4} & \textbf{Hagit Shatkay}^{1,*} & \textbf{Christopher D. Tunnell}^{2,5} \end{array}$

Department of Computer and Information Sciences, University of Delaware
Department of Physics and Astronomy, Rice University
Department of Electrical and Computer Engineering, Rutgers, The State University of New Jersey
Department of Statistics, Rutgers, The State University of New Jersey
Department of Computer Science, Rice University
† petersc@udel.edu

Abstract

We demonstrate that Bayesian networks fill a significant methodology gap for uncertainty quantification in particle physics, providing a framework for modeling complex systems with physical constraints. To address the problem of interaction position reconstruction in dark matter direct-detection experiments, we built a Bayesian network that utilizes domain knowledge of the system in both the structure of the graph and the representation of the random variables. This method yielded highly informative per-interaction uncertainties that were previously unattainable using existing methodologies, while also demonstrating comparable precision on reconstructed positions.

1 Introduction

While machine learning research has led to many recent advancements in physical science research, the efficacy of most machine learning methodologies for science is limited by their inability to quantify uncertainties on the parameters they estimate. Within the physical sciences, uncertainties are crucial for hypothesis testing and arise from several sources including limitations in the ability to observe and model systems, as well as inherent nondeterminism. Observations may be uncertain due to only some aspects of the system being observed, measurement error, or measurements having finite resolution. Therefore, it is imperative to have machine learning approaches that take into account both the different possible states of the system and the probabilities of each possible state.

To account for uncertainties, one may construct a probabilistic model representing the joint distribution over the variables in the system. For complex systems with hundreds of variables, constructing the joint distribution is often intractable. Approaches to handling this intractability generally rely on scientists comparing the summary statistics of observed data to those of simulated data [1].

Alternatively, Bayesian networks (BNs) — one of the two broad classes of Probabilistic Graphical Models (PGMs) — use a graph-based representation of the joint distribution as the basis for compactly encoding a high-dimensional distribution. As PGM representations are generalizable and encode domain knowledge of the system, they enable understanding and evaluation of the properties of a complicated distribution, as well as construction of accurate models of a system.

^{*}Deceased.

Both neural networks (NNs) and BNs have been in widespread use since at least the late 1980s. For BNs, the development of a rigorous formalism for probabilistic reasoning by Judea Pearl [2], covering representation and inference, was key to the methodology gaining acceptance. It is important to note that Bayesian networks are not the same as Bayesian NNs — which are NNs that have posterior distribution of weights. While NNs provide a generalizable framework for many classes of problems, they were not designed to provide information about the reliability of their predictions. These methods of uncertainty quantification can lack transparency and the uncertainties predicted by NNs often require re-calibration before being used for analysis [3–5].

Contribution We demonstrate the utility of BNs as a methodology to address questions within the physical sciences where uncertainty quantification is crucial. We specifically explore the problem of position reconstruction — an inverse problem where location is inferred or estimated based on sensor measurements — within astroparticle physics. Robust position reconstruction is paramount for enabling rare-event discoveries by dark matter detection experiments, as it allows for focus on interactions occurring only within the central volume of a detector where there are fewer backgrounds.

Related Work Within astrophysics BNs have been used in the form of Bayesian hierarchical models; recent examples include deriving luminosity-metallicity relations of RR Lyrae stars [6], developing an SED model for type Ia supernovae [7], clarifying the Hubble constant tension [8], and type Ia supernova cosmology fits [9, 10]. Furthermore, they have been used to develop trustworthy estimates of redshift distributions [11]. Outside of astrophysics, Bayesian networks have been used for modeling of nuclear data [12] and to establish the significance of coincident events by gravitational-wave detectors [13].

Within the field of particle physics, there have been significant efforts to apply modern machine learning techniques, demonstrated by the hundreds of papers in Ref. [14]. However, work on uncertainty awareness and quantification when using machine learning techniques has been primarily on estimating uncertainty using deep learning ([15]; e.g., LHC searches [16], and neutrino reconstruction [17]).

2 Brief Review of Bayesian Networks

BNs use a directed acyclic graph to encode a probability distribution [18, 19] by making use of the independencies between variables [20–23]. The directed edges in BNs correspond to direct influence of one variable on another, allowing the networks to be used as interpretable models of physical systems for reasoning about causes and effects within systems [2, 24]. By utilizing the conditional independencies between the variables, the graphical representation of the joint distribution is more compact than the full joint distribution over the variables. The local probability information is a conditional distribution given the immediate parents of the node, i.e. $P(X_i \mid \operatorname{Parents}(X_i))$. Each entry in the joint distribution is defined as the product of the local conditional distributions,

$$P(X_1 = x_1, ..., X_N = x_N) = \prod_{i=1}^{N} P(x_i \mid \text{parents}(X_i)),$$
 (1)

where parents (X_i) are the values of the parent nodes, $Parents(X_i)$, that appear in x_1, \ldots, x_n . Moreover, the posterior distribution over a variable of interest conditioned on an observation can be computed by performing a probability query. This is performed by computing the posterior distribution over the values of the query variables, $\mathbf{Y} = \{Y_1, \ldots, Y_M\}$, conditioned on the observed values, $\{x_1, \ldots, x_N\}$, of the evidence variables, $\mathbf{X} = \{X_1, \ldots, X_N\}$: $P(Y_1, \ldots, Y_M \mid X_1 = x_1, \ldots, X_N = x_N)$. Thus, the BN framework is well-suited for determining the probability of any one of several causes being a contributing factor to an observed event.

3 Building a Bayesian Network for Position Reconstruction

A major and novel component of this work is representing the position reconstruction problem using the BN framework.

Data For dark matter direct detection experiments such as XENONnT [25] or LZ [26], the two-dimensional position of a particle interaction within the cylindrical detector can be inferred from the light detected by the photosensor array on the top of the detector, commonly called a hit pattern [27].

Using Ref. [28] we simulated a data set based on the detector geometry of XENONnT, comprised of the ground-truth position of each physical interaction, the ground-truth number of electrons subsequently produced, and the set of associated photosensor measurements. The spatial location of a particle interaction within the detector is denoted as a two-dimensional position in polar coordinates. The positions were randomly generated from a uniform distribution over the detector area. The number of electrons were generated from a uniform random distribution from 1 to 2000 electrons, which includes the range caused by both dark matter particle interactions as well as common sources of background. See Section 3 of Ref. [29] for a detailed description of the data generation process. We generated training and testing sets of 5×10^6 and 5×10^4 interactions, respectively.

Nodes We include nodes representing the interaction position, the number of electrons, and the photosensor measurements. Which variables to include and the set of values that each of the variables can take on are crucial choices in framing the problem. These choices allow us to encode physical constraints into the BN. As this is a proof-of-concept, we include a limited number of variables to decrease the complexity of the problem.

The interaction position is represented by the random variable L. We represent locations within the detector as discrete values by dividing the detector into $1 \, \mathrm{cm^2}$ pixels, a size chosen to provide the minimum precision necessary for the position reconstruction to be useful in experimental analyses. An advantage to discretizing a continuous random variable is that the conditional probabilities can be specified explicitly for each value, without defining a parametric form. Thus, the two-dimensional position of the interaction is represented by the single discrete multinomial random variable L which maps each particle interaction to the pixel of its ground-truth spatial location. The set of possible values that L can take is $\{0, \ldots, 13845\}$, the indices corresponding to the discrete pixels.

The number of electrons produced by the interaction is represented by the random variable E, which maps each particle interaction to its ground-truth number of electrons. The set of possible values that E can take is all positive integers. In addition, we define a set of 253 random variables, $\mathbf{S} = \{S_0, \ldots, S_{252}\}$, which maps each particle interaction to the ground-truth photosensor measurement, with all non-negative integers as the set of possible values.

Graph Structure The graph structure we choose, shown in Fig. 1, provides a compact representation of a high-dimensional joint distribution using the strong assumption that each S node is conditionally independent of all other S nodes, given the L and E nodes. This simple graph structure has been proven to be quite effective for classification in practice, even in cases where the independence assumptions are violated [30]. Its advantages are that it is easier to interpret, faster to learn and query, and smaller to store in memory than a more complex structure.

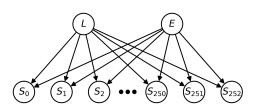


Figure 1: Graph structure used in this work.

Local Probability Distributions Both L and E have no parent nodes, thus the local probability distributions are priors — they are not dependent on other random variables. We define the prior over L to be the fraction of training set interactions within each pixel and the prior over E to be a uniform distribution from 1 to 2000 electrons. The S nodes each have both L and E as parent nodes. To make the network both faster to query and smaller to store in memory, we represent the local probability distribution for the S nodes with a Poisson distribution. This choice also allows us to encode domain knowledge of the system, as the number of photons detected by a sensor over a given time period is well represented as a Poisson distribution. The Poisson distribution has one parameter λ , which is both the expected value and the variance, and is a distribution over nonnegative integers. We define the local probability distribution at each S node, given any assignment to the L and E nodes, to be

$$P(S_j = s_j \mid L = l, E = e) = \frac{(\lambda_{l,j,e})^{s_j} \exp(-\lambda_{l,j,e})}{s_j!},$$
 (2)

where $\lambda_{l,j,e}$ is learned from the training data.

Probability Queries The posterior distribution over the L and E nodes, given the observed values of the S nodes, s_0, \ldots, s_{252} , is

$$P(L, E \mid S_0 = s_0, ..., S_{252} = s_{252}) = P(L, E) \prod_j P(S_j = s_j, L, E).$$
 (3)

The distribution in (3) is a categorical distribution which assigns the probabilities from the joint distribution over L and E. The marginal probability distribution over L can be calculated by summing over the set of possible values of E.

4 Results

We built a Bayesian network with the structure shown in Fig. 1 and local probability distributions learned from the simulation. We then performed probability queries to attain the posterior distribution over position and number of electrons for the interactions in the test set.

Fig. 2 shows the marginal distribution over L given the photosensor measurements (left) visualized as a heatmap out to the limits of the 5- σ confidence interval (center) for an example interaction from the test set. The root mean square (RMS) of the differences between the ground-truth position and expectation value of the position for the test set is 0.753 cm. This is dominated by the 1 cm² pixel area and reduced by increasing the number of pixels. A feature of the method which resulted from our decision to discretize the value of the position node is that it naturally constrains the reconstructed position to be within the physical volume of the detector, in contrast to NNs which often produce nonphysical results [27] or require a customized layer to place physical constraints on outputs [29]. Fig. 2 also presents a comparison of two sets of photosensor measurements: the first is a simulated interaction and the second is the same interaction with the intensities randomly reassigned to different sensors. One major shortcoming of other existing reconstruction methods is that they only provide a single point estimate for the interaction position, which can be particularly uninformative when the photosensor measurements have no corresponding particle interactions (i.e., noise). Fig. 2 demonstrates that BNs overcome this shortcoming by providing informative posteriors for position reconstruction even in the case of noise, which could make the posteriors useful for anomaly detection.

5 Conclusion

We found the BN framework to be well suited to uncertainty quantification of position reconstruction in an astroparticle experiment. The reconstructed positions inferred by our model have precision comparable to existing reconstruction methods for dark matter detection experiments [27, 29, 31], and can be improved by using a smaller pixel size. More importantly, the posterior distributions over position are informative and enable a variety of experimental analyses ranging from anomaly detection to determining the boundaries of the fiducial volume, the central region of a detector where there are fewer backgrounds. We found that our choice of representation for the values of the nodes (as a categorical distribution in the case of L) and the local probability distributions (as a Poisson distribution in the case of S) greatly impacted both the computational efficiency of the inference and the accuracy of the inferred positions. This method can be extended to energy and three-dimensional position reconstruction, as well as signal classification. Improvements to the reconstructed positions can be made by utilizing a more complex graph structure. For example, adding edges between sensor nodes to account for correlations between sensors and adding nodes for experimentally relevant variables. Alternatively, the structure can be learned directly from data.

The BN framework is based on a formalism for probabilistic reasoning and can incorporate scientists' knowledge about the system to build a physically interpretable model. Based on this proof-of-concept, we conclude that the BN framework, while not applicable to all analyses, is a well suited uncertainty quantification method for a variety of other reconstruction tasks where per-event uncertainties are crucial.

Broader Impact

There are numerous applications of the BN framework within particle physics, as well as more generally throughout the physical sciences. As uncertainty quantification using BNs and other methods

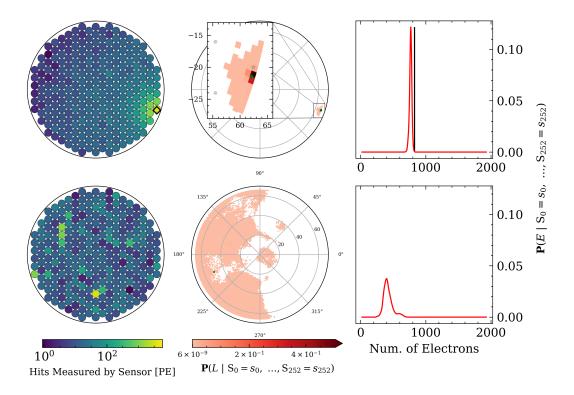


Figure 2: Set of photosensor measurements (left), posterior distribution over position (center), and posterior distribution over number of electrons (right) for a simulated interaction (top row) and scrambled photosensor measurements where the intensities measured by the sensors are randomly reassigned (bottom row). This demonstrates that, unlike reconstruction methods that only provide a point estimate, the Bayesian network method can provide an informative position reconstruction even in the bottom case where there is no interaction signal in the photosensor measurements. True position and true number of electrons are shown as a black diamonds and vertical black line in the top row. The bottom row has no ground-truth position or ground-truth number of electrons. The expectation value of position is shown as a green triangle in both center panels. The top center panel inset shows the sensor positions as filled grey circles.

is further developed specifically for applications within the physical sciences, we anticipate that scientists will increasingly account for uncertainties in their analyses and hypothesis testing, which will in turn enhance scientific research as a whole. We do not envision that this methodology will result in any negative ethical or societal impacts in the future.

Acknowledgments and Disclosure of Funding

This work is supported by the National Science Foundation through awards 1940074, 1940209, and 1940080. The authors would like to acknowledge the following open source projects which were used in this work: *numpy* [32], *scipy* [33], and *matplotlib* [34].

References

- [1] K. Cranmer, J. Brehmer, and G. Louppe, Proc Natl Acad Sci 117, 30055 (2020).
- [2] J. Pearl, Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference (Morgan-Kaufmann, 1988).
- [3] L. Hirschfeld, K. Swanson, K. Yang, R. Barzilay, and C. W. Coley, J. Chem. Inf. Model. **60**, 3770–3780 (2020).
- [4] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, arXiv e-prints (2017), arXiv:1706.04599.

- [5] J. Gawlikowski *et al.*, arXiv e-prints (2021), arXiv:2107.03342.
- [6] A. Garofalo, H. E. Delgado, L. M. Sarro, et al., Monthly Notices of the Royal Astronomical Society 513, 788 (2022).
- [7] K. S. Mandel, S. Thorp, G. Narayan, *et al.*, Monthly Notices of the Royal Astronomical Society **510**, 3939 (2022).
- [8] S. M. Feeney, D. J. Mortlock, and N. Dalmasso, Monthly Notices of the Royal Astronomical Society **476**, 3861 (2018).
- [9] S. R. Hinton et al., Astrophysical Journal 876, 15 (2019).
- [10] D. Rubin et al. (The Supernova Cosmology Project), Astrophysical Journal 813, 137 (2015).
- [11] A. I. Malz and D. W. Hogg, Astrophysical Journal 928, 127 (2022).
- [12] G. Schnabel et al., arXiv e-prints (2021), arXiv:2110.10322.
- [13] R. Essick, G. Mo, and E. Katsavounidis, Physical Review D 103, 042003 (2021).
- [14] HEP ML Community, "A Living Review of Machine Learning for Particle Physics," imlwg.github.io/HEPML-LivingReview/ (2022).
- [15] T. Y. Chen, B. Dey, A. Ghosh, M. Kagan, B. Nord, and N. Ramachandra, Submitted to the Proceedings of the US Community Study on the Future of Particle Physics (2022), arXiv:2208.03284.
- [16] B. Nachman, SciPost Phys. 8, 090 (2020).
- [17] D. H. Koh, A. Mishra, and K. Terao, in *Bayesian Deep Learning Workshop, NeurIPS 2021* (2021).
- [18] R. A. Howard and J. E. Matheson, Readings on the Principles and Applications of Decision Analysis (Strategic Decisions Group, 1989).
- [19] J. Q. Smith, The Annals of Statistics, 654 (1989).
- [20] T. Verma and J. Pearl, in 4th Workshop on Uncertainty in Artificial Intelligence (1988) p. 352.
- [21] D. Geiger and J. Pearl, in *Proceedings of the Fourth Conference on Uncertainty in Artificial Intelligence* (1988).
- [22] D. Geiger, T. Verma, and J. Pearl, in *Machine Intelligence and Pattern Recognition*, Vol. 10 (Elsevier, 1990) pp. 139–148.
- [23] D. Geiger, T. Verma, and J. Pearl, Networks 20, 507 (1990).
- [24] M. P. Wellman, Artificial intelligence 44, 257 (1990).
- [25] E. Aprile et al. (XENON Collaboration), JCAP 11, 031 (2020).
- [26] D. S. Akerib et al. (The LZ Collaboration), Nucl.Instrum.Meth.A 953, 163047 (2020).
- [27] E. Aprile et al. (XENON Collaboration), Physical Review D 100, 052014 (2019).
- [28] S. Liang and C. Tunnell, "Domain-informed Neural Networks," 10.5281/zenodo.5771941 (2021).
- [29] S. Liang, A. Higuera, C. Peters, V. Roy, W. U. Bajwa, H. Shatkay, and C. D. Tunnell, Front. Artif. Intell. (2022).
- [30] P. Domingos and M. Pazzani, Machine Learning **29**, 103–130 (1997).
- [31] D. Zhang et al. (PandaX-4T Collaboration), Journal of Instrumentation 16 (2021).
- [32] C. R. Harris et al., Nature 585, 357 (2020).
- [33] P. Virtanen et al. (SciPy 1.0 Contributors), Nature Methods 17, 261 (2020).
- [34] J. D. Hunter, Computing in Science & Engineering 9, 90 (2007).

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default [TODO] to [Yes], [No], or [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? [Yes] See Section X.
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

- 1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] Limitations are discussed in Sec. 3 and Sec. 4.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
- 3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No] These will be available as a public GitHub repository following the publication of the full work.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [No] . The general approach to training is discussed in Sec. 3 and full details will be provided with the code.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [No]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [No] This is listed on the Zenodo page of the assests.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
- 5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]