
GNN-Guided Block Selection in Gibbs MCMC

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Exact inference in large Bayesian Networks (BNs) is computationally intractable,
2 limiting its practical application. Markov Chain Monte Carlo (MCMC) methods
3 like Gibbs sampling offer a scalable alternative but can be arbitrarily slowed by
4 highly coupled variables— addressable by jointly sampling some variables as a
5 block. We propose an automated block detection method to amortise inference
6 time: training a Graph Neural Network (GNN) to propose blocks directly from the
7 BN structure. We further introduce a novel coupling heuristic based on the Markov
8 chain’s spectral gap, which we show can be more robust than existing heuristics.
9 Our GNN, trained on a dataset of small, randomly generated BNs, generalizes well
10 to larger networks, accelerating MCMC sample efficiency in our experiments.

11 1 Introduction

12 **Bayesian Networks (BNs)** are a class of probabilistic graphical models that represent complex
13 multivariable distributions through local conditional dependencies. They provide an inepretable
14 framework for probabilistic modeling, supporting arbitrary posterior inference queries given observed
15 evidence. BNs can be constructed from human expert knowledge, structure learning algorithms, or
16 even elicited from large language models. However, exact inference is computationally intractable
17 for larger networks. Approximate inference methods like **Markov Chain Monte Carlo (MCMC)**
18 offer a practical alternative for achieving acceptable accuracy within reasonable computation time.

19 **Gibbs sampling** is a common MCMC method for BNs due to its simple computation of local updates
20 and 100% acceptance ratio. However, highly coupled variables can severely slow convergence and
21 mixing. **Blocked Gibbs sampling** addresses this by grouping coupled variables into **blocks** that are
22 jointly sampled, breaking out of sticky Markov chain states.

23 While practitioners can specify blocks a priori using domain knowledge, automatic blocking methods
24 are essential for general applicability. Venugopal and Gogate (2013) present a robust algorithm
25 for dynamically proposing blocks during MCMC based on information revealed by initial samples.
26 However, this approach faces a fundamental catch-22: effective blocks require sufficient posterior
27 samples of the full range of Markov chain modal states to observe coupling, yet collecting this data is
28 slowed by the couplings themselves.

29 Circumventing this issue, we train a **Graph Neural Network (GNN)** Kipf and Welling (2017)
30 to propose effective blocks directly from BN structure before MCMC starts; this can be used in
31 tandem with dynamic block refinement. Our GNN is trained offline on a diverse dataset of randomly
32 generated BNs, amortising the computational cost of block identification. We further introduce an
33 alternative coupling heuristic based on the Markov chain’s **spectral gap**, which demonstrates superior
34 performance on adversarially coupled BNs compared to existing measures.

35 Yoon et al. (2019) trained GNNs to directly do inference, though with fixed approximation error;
36 our method fully integrates into MCMC and can be used with existing techniques. GNN-proposed
37 blocks can still be refined further dynamically refined using MCMC samples (Venugopal and Gogate,

2013). Learned variable sampling distributions allow proposing larger blocks (Wang et al., 2018), non-uniform variable/block selection rates speed up convergence, and of course parallelisation across chains and variables with non-intersecting Markov blankets is still effective (Gonzalez et al., 2011).

2 Guiding Block Selection in Gibbs MCMC with GNNs

In order to propose good blocks, we first must investigate coupling: Section 2.1 derives a good heuristic metric to detect coupling between pairs of variables. We use this in the algorithm of Section 2.2 to propose a block partition $\mathbb{B} = \{B_i\}_{i=1}^k$ to cover the variables of the BN: $\{X_j\}_{X_j \in \mathcal{V}} = \bigsqcup B_i$ (though in some cases overlapping blocks are better). Finally in Section 2.3, we show how we train a GNN to predict this coupling metric; put together we can propose blocks for an arbitrary input BN.

2.1 Deriving the spectral gap coupling heuristic

For a finite discrete-time Markov chain (S, T, μ_0) of state space S , transition matrix $T_{ij} = P(Z_{t+1} = s_j | Z_t = s_i)$ and initial distribution $Z_0 \sim \mu_0$, convergence in **total variation distance (TVD)** to the stationary distribution π is governed by its spectral gap Levin and Peres (2017), per the bound $\|\mu_0^T T^n - \pi^T\|_{TV} = O(\lambda_2^n)$, where $\lambda_2 \neq 1$ is the second-largest eigenvalue. The spectral gap $\gamma = 1 - \lambda_2$ thus provides a natural basis for block selection. Unfortunately, computing the spectral gap for multivariable chains is impractical except in special cases (see e.g. Chimisov et al. (2018)), requiring enumeration of an exponentially large product state space.

An important simplification that Venugopal and Gogate (2013) make is to consider only variable pairs. For two variables X, Y in a BN, the joint posterior $\pi(X, Y)$ can be feasibly estimated in a short MCMC run as the $|\mathcal{X}| \times |\mathcal{Y}|$ matrix P . The hellinger distance of P vs its independent assumption rank one approximation, $Q = p_X p_Y^T$, where e.g. $(p_X)_x = \sum_{y \in \mathcal{Y}} P_{xy}$, is a good metric of coupling as if the independence approximation is exact, X and Y would seem to be not at all coupled. Their proposed heuristic is thus $\text{HD}(P, Q) := \frac{1}{\sqrt{2}} \sqrt{\sum_{i,j} (\sqrt{P_{ij}} - \sqrt{Q_{ij}})^2}$.

We develop a spectral gap based pairwise score by constructing an averaged subset transition matrix $T^{(X,Y)} : |\mathcal{X}| \times |\mathcal{Y}| \rightarrow |\mathcal{X}| \times |\mathcal{Y}|$ that models marginal dynamics for variables X, Y under standard Gibbs sampling. This decomposes as $T^{(X,Y)} = \frac{1}{2} T_X^{(X,Y)} + \frac{1}{2} T_Y^{(X,Y)}$, where $T_Y^{(X,Y)}$ represents sampling Y with X fixed. Denoting their joint Markov Blanket as $\text{MB}(X, Y) = \text{MB}$, the transition probabilities $T_Y^{(X,Y)} : (x_0, y_0) \rightarrow (x_0, y')$ are derived by averaging over the distribution of MB, which we approximate as $\text{MB} \sim \pi(\text{MB}|x_0, y_0)$ (as we perform Gibbs sampling over the whole BN)

$$P_{\text{Gibbs}}(Y|X = x_0, Y = y_0, \text{MB}) \propto P(Y|Pa(Y)) \prod_{Z: Y \in Pa(Z)} P(Z|Pa(Z)) \quad (1)$$

$$P_{\text{Gibbs}}(Y|X = x_0, Y = y_0) \approx \int \pi(\text{MB}|x_0, y_0) P_{\text{Gibbs}}(y|\text{MB}, x_0) d\text{MB} \quad (2)$$

Since $\pi(\text{MB}|x_0, y_0)$ is intractable, we approximate by averaging over $y_0 \sim \pi(Y = y|X = x) = \int \pi(\text{MB}|x) P(y|\text{MB}, x) d\text{MB}$ to get uniform probabilities $\tilde{T}_Y^{(X,Y)} : (x_0, *) \rightarrow (x_0, y')$. This represents a best-case mixing scenario where sampling y' jumps straight to the true posterior $\pi(Y|X = x_0)$, so is a lower bound for mixing time of the actual Gibbs chain.

$$E_{Y_0 \sim \pi(Y|x_0)} [P_{\text{Gibbs}}(Y|X = x_0, Y = y_0)] \approx \quad (3)$$

$$\int \left[\int \pi(\text{MB}|x_0, y_0) P_{\text{Gibbs}}(y|\text{MB}, x_0) d\text{MB} \right] \pi(y_0|x_0) dy_0 \quad (4)$$

$$\int \left[\int \pi(\text{MB}|x_0, y_0) \pi(y_0|x_0) dy_0 \right] P_{\text{Gibbs}}(y|\text{MB}, x_0) d\text{MB} = \quad (5)$$

$$\int \pi(\text{MB}|x_0) P_{\text{Gibbs}}(y|\text{MB}, x_0) d\text{MB} = \pi(Y|X = x_0) \quad (6)$$

This allows us to derive an approximate matrix $\tilde{T}^{(X,Y)}$ and compute its spectral gap using only the posterior $\pi(X, Y)$: the same information required for the Hellinger distance heuristic. While these heuristics perform similarly on aggregate across random BNs, in certain adversarially constructed BNs (e.g. see A.2) spectral gap substantially outperforms Hellinger distance for block proposals.

2.2 Adapting the greedy pairwise heuristic blocking algorithm

We adapt the greedy block merging algorithm from Venugopal and Gogate (2013). Given a pairwise distance function $\rho(X, Y)$ between variables in the BN, where larger values indicate stronger coupling, we construct a block partition \mathbb{B} of a BN of n variables. Starting with n singleton blocks, we iteratively merge the pair of blocks with maximum inter-block score $\sum_{X \in B_i, Y \in B_j} \rho(X, Y)$ (or mean/max). We compare replacing hellinger distance based ρ_{HD} with spectral gap. We explore replacing the hellinger distance based ρ_{HD} with one from spectral gap ρ_{spectral} . For simplicity we constrain merged block size to a uniform cap rather than via more precise computational cost of Gibbs sampling each block’s joint distribution — otherwise larger blocks are generally better.

2.3 Producing block proposals with GNNs

GNNs provide a natural architecture to analyse an input BN so as to produce block proposals. Rather than learning to propose blocks directly, we supervise train the GNN to produce either ρ_{HD} or ρ_{spectral} with a final bilinear layer between node embeddings. For ease of computing pairwise posteriors for these ground truth heuristics, we limit our training dataset to randomly generated BNs of 10-40 nodes and average degree 1.7. Despite this constraint, our GNNs generalise effectively to larger networks, with spectral gap showing slight performance advantages over Hellinger distance: see Figure 1.

Loopy Belief Propagation (LBP) Koller and Friedman (2009) inspires our architecture through its effective graph-based message passing for computing single variable marginal posteriors. LBP operates on **factor graphs**, where factors generalize **Conditional Probability Tables (CPTs)** $P(X|Pa(X))$ by removing normalisation constraints. We convert BNs to factor graphs by replacing each CPT hyperedge $P(X|Pa(X) = \{Y_i\}_{i=1}^{\deg(X)})$ with a **factor node** F_X . Variable nodes $X, Y_1, \dots, Y_{\deg(X)}$ connect bidirectionally to F_X , enabling distinct message passing phases.

We employ 2 relational graph convolution layers Schlichtkrull et al. (2018) with GRU units Cho et al. (2014) across 2 message passing rounds. While deeper networks with additional rounds would likely improve performance, scaling proved challenging; future work could explore alternative architectures like graph attention networks Veličković et al. (2018).

We trained on 22,732 discrete BNs using an 80-20 train-validation split. CPT encoding constrains our architecture: we flatten each CPT as an initial node feature vector of size $N_d^{N_p+1}$ (max domain size $N_d = 5$, max number of parents $N_p = 6$) for factor nodes as well as a factor indicating one-hot flag, while variable nodes are zero initialised. Edge types encode direction (to/from factor) and position (0 for child, 1, ..., N_p for parents), determining relational graph convolution parameters. Following LBP’s un-normalised factors, we handle evidence by disconnecting observed variables from their factors and contracting those CPTs to conforming entries, which become unnoramlised. To encourage GNN comprehension across all factors, we further add multiplicative noise to all CPTs making them all unnormalised - this improves performance.

Randomly generated BNs were filtered to fit the max number of parents constraint, with CPTs deliberately extremised with higher occurrence of entries in $[0.99, 1.0]$ to encourage highly coupled variables and hence the presence of some higher spectral gap / hellinger distance scores. Additional training details appear in Appendix A.1.

3 Experiments

We evaluate our method’s practical utility, particularly generalisation performance on larger BNs, where exact inference is intractable.

We test our GNN’s block proposal performance on 100 BNs of 85-115 nodes each, measuring mean TVD between MCMC 200-sample predicted and “ground truth” (7,000-sample) single variable marginals. To account for stochasticity in highly coupled networks, we average results across 25

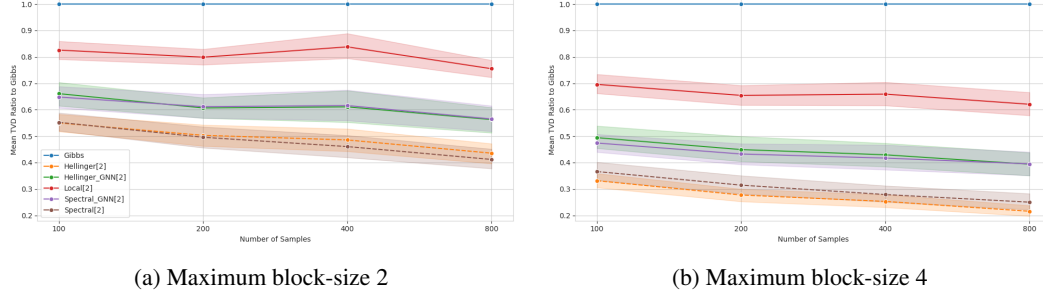


Figure 1: GNN block proposals on 100 small BNs (10-40 nodes) test-set achieve lower mean TVDs than Gibbs and the control of random local blocking, across a range of MCMC sample sizes. Mean and 95% confidence intervals shown.

independent runs per BN. Figure 2 shows that both GNN-predicted spectral gap and Hellinger distance blocking substantially outperform random local blocking across maximum block sizes (2 and 4 shown). Spectral gap shows only very modest improvements over Hellinger distance.

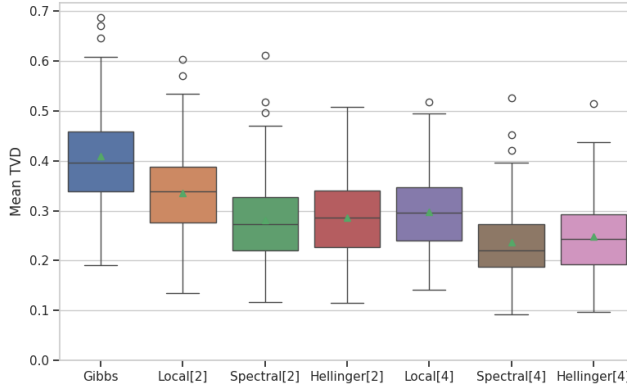


Figure 2: 200 sample MCMC runs of various maximum block-sizes, on 100 large BNs (85-115 nodes). GNN guided block proposals generalise well to larger BNs where exact inference is infeasible, outperforming the random local blocking control; spectral slightly outperforms Hellinger. Means shown as green triangles.

4 Conclusion

We presented a method to accelerate MCMC inference in Bayesian Networks using Graph Neural Networks to propose variable blocks for joint sampling. This amortised approach circumvents the catch-22 of dynamic blocking methods that require samples to identify the couplings that slow sampling. Our experiments demonstrate that GNN-proposed blocks, trained on spectral gap or Hellinger distance heuristics, substantially accelerate posterior convergence on large networks.

Future work could explore more expressive architectures like Graph Attention Networks with additional layers, to better capture global dependencies. A more ambitious direction would be to replace heuristic supervision with direct block proposal mechanisms trained via reinforcement learning on convergence speed, perhaps even with active exploration of sampling moves. Further analysis is warranted to characterize the topological or parametric properties of BNs for which our proposed spectral gap heuristic offers the greatest advantage over hellinger distance.

References

- Chimisov, C., Latuszynski, K., and Roberts, G. (2018). Adapting the gibbs sampler. *arXiv preprint arXiv:1801.09299*.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation.

- 141 Gonzalez, J., Low, Y., Gretton, A., and Guestrin, C. (2011). Parallel gibbs sampling: From colored
142 fields to thin junction trees. In *Proceedings of the Fourteenth International Conference on Artificial
143 Intelligence and Statistics*, pages 324–332. JMLR Workshop and Conference Proceedings.
- 144 Kipf, T. N. and Welling, M. (2017). Semi-supervised classification with graph convolutional networks.
- 145 Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT
146 press.
- 147 Levin, D. A. and Peres, Y. (2017). *Markov chains and mixing times*, volume 107. American
148 Mathematical Soc.
- 149 Schlichtkrull, M., Kipf, T. N., Bloem, P., Van Den Berg, R., Titov, I., and Welling, M. (2018).
150 Modeling relational data with graph convolutional networks. In *European semantic web conference*,
151 pages 593–607. Springer.
- 152 Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. (2018). Graph
153 attention networks.
- 154 Venugopal, D. and Gogate, V. (2013). Dynamic blocking and collapsing for gibbs sampling. *arXiv
155 preprint arXiv:1309.6870*.
- 156 Wang, T., Wu, Y., Moore, D., and Russell, S. J. (2018). Meta-learning mcmc proposals. *Advances in
157 neural information processing systems*, 31.
- 158 Yoon, K., Liao, R., Xiong, Y., Zhang, L., Fetaya, E., Urtasun, R., Zemel, R., and Pitkow, X. (2019).
159 Inference in probabilistic graphical models by graph neural networks. In *2019 53rd Asilomar
160 Conference on Signals, Systems, and Computers*, pages 868–875. IEEE.

161 A Technical Appendices and Supplementary Material

162 Technical appendices with additional results, figures, graphs and proofs may be submitted with
163 the paper submission before the full submission deadline (see above), or as a separate PDF in the
164 ZIP file below before the supplementary material deadline. There is no page limit for the technical
165 appendices.

166 A.1 GNN training and examples

167 We trained both the spectral gap and hellinger distance predicting GNNs on a single A600 for 20
168 epochs each. We used pyAgrum’s random bayesian network generator, which produces diverse
169 graphs, of varying degree (of given average, 1.7), varying CPTs (which we further extremised to
170 promote coupling), varying domain size up to a given cap (we constrained to $N_d = 5$). We rejection
171 sample out BNs satisfying the max number of parents constraint ($N_p = 6$), and finally randomly
172 choose a number of evidence variables and their observed assignments, of between 1 and 20% of the
173 BN’s variables.

174 For spectral gap, we trained the GNN to predict the second largest eigenvalue of the approximate
175 Markov chain subset transition kernel of a pair of variables in the BN. For our randomly generated
176 BNs, this is quite an unbalanced distribution, with high concentration of values close to 0, same for
177 hellinger distance - however a few pairs of nodes have higher values corresponding to a higher degree
178 of coupling. In order to encourage the GNN to learn not just to predict 0 distance uniformly, we
179 actually employ mean cubed loss rather than mean squared, to more sharply penalise large errors.

180 A.2 Example Adversarial BN requiring spectral gap analysis

181 We constructed by hand a small BN of three nodes X, Y, Z and conditional structure $Y \rightarrow X; Y \rightarrow$
182 Z where both X, Y and Y, Z are desirable blocking candidates. However assuming a max block-
183 size of 2 where we only do a partition style blocking, the spectral gap and hellinger distance
184 metrics disagree on which is preferential to block. Y is a root node with uniform distribution
185 $P(Y = i) = 0.25$ for $i \in \{0, 1, 2, 3\}$, and the CPTs of $P(X|Y)$ and $P(Z|Y)$, shown in Figures 5b
186 and 5c, both ensure that the marginal probabilities $P(X)$ and $P(Z)$ are likewise uniform. The far more

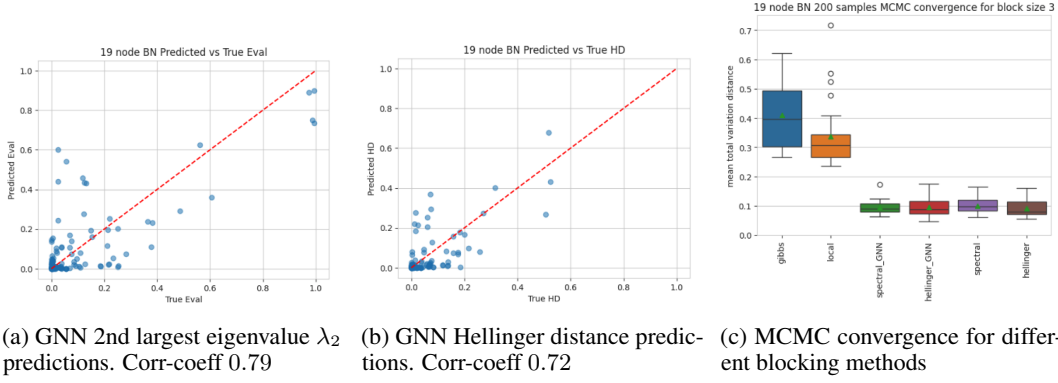


Figure 3: GNN predictive accuracy of pairwise variable distance metrics, and MCMC convergence performance, shown for a single validation set BN of 19 nodes.

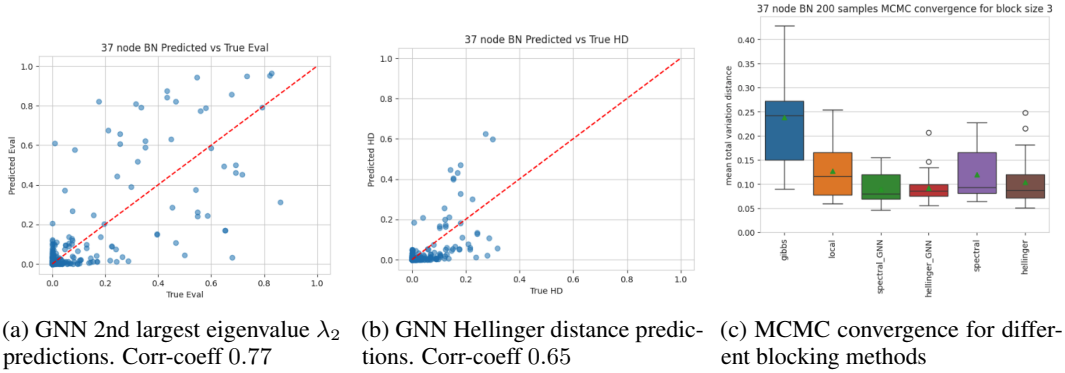
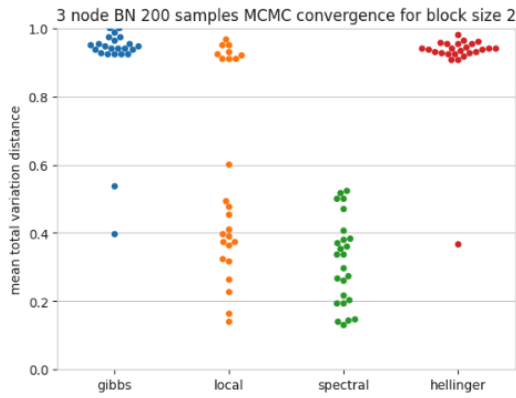


Figure 4: GNN predictive accuracy of pairwise variable distance metrics, and MCMC convergence performance for a real life BN of 37 nodes: <https://www.bnlearn.com/bnrepository/discrete-medium.html#alarm>, with evidence VENTALV: 0, HYPOVOLEMIA: 1, INSUF-FANESTH: 0, HRBP: 1

187 extreme, 2x2-modal concentration of probabilities in $P(X|Y)$ actually lead to far slower convergence
 188 of the Markov chain, as evidenced in Figure 5a, unlike the more weakly 4-modally concentrated
 189 coupling in $P(Z|Y)$. Hellinger distance scores are however insensitive to this large discrepancy,
 190 indeed giving preference to blocking Z, Y over X, Y : $\text{HD}(X, Y) = 0.532$ and $\text{HD}(Y, Z) = 0.543$.
 191 The eigenvalues of the approximate transition matrices $\tilde{T}^{(Y,Z)}$ has 2nd largest eigenvalue 0.8652 so
 192 spectral gap 0.1348 as opposed to $\tilde{T}^{(X,Y)}$ which has $\lambda_2 = 0.9994$ and so a far smaller spectral gap
 193 of 0.0006, so spectral gap based blocking correctly blocks X, Y . In this case random local blocking
 194 blocks correctly half of the time.



(a) MCMC convergence for max block-size 2

	X			
Y	0	1	2	3
0	0.0001	0.4999	0.0001	0.4999
1	0.4999	0.0001	0.4999	0.0001
2	0.0001	0.4999	0.0001	0.4999
3	0.4999	0.0001	0.4999	0.0001

(b) X Given Y CPT

	Z			
Y	0	1	2	3
0	0.0455	0.9082	0.0455	0.0009
1	0.0009	0.0455	0.9082	0.0455
2	0.0455	0.0009	0.0455	0.9082
3	0.9082	0.0455	0.0009	0.0455

(c) Z Given Y CPT

Figure 5: Adversarial 3 node example BN $Y \rightarrow X$; $Y \rightarrow Z$ has high degrees of coupling, but magnitudes more so between X, Y than Z, Y . The spectral gap heuristic correctly identifies this unlike hellinger distance, and so exhibits far better MCMC convergence. Shown here is the