

---

# Bayesian Network Structure Learning using Digital Annealer

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Annealing processors, which efficiently solve a quadratic unconstrained binary  
2 optimization (QUBO), are a potential breakthrough in improving the accuracy  
3 of score-based Bayesian network structure learning. However, currently, the bit  
4 capacity of an annealing processor is very limited. To utilize the power of an-  
5 nealing processors, it is necessary to encode score-based learning problems into  
6 QUBO within the upper bound of bits. In this paper, we propose a novel approach  
7 with direct encoding of candidate parent sets in the form of Cartesian products.  
8 Experimental results on benchmark networks with 27 to 70 variables show that  
9 our approach requires lesser bits than the bit capacity of the second-generation  
10 Fujitsu digital annealer, a fully coupled annealing processor developed by with  
11 semiconductor technology. Moreover, we demonstrate that the digital annealer  
12 with our conversion method consistently outperforms the state-of-the-art heuristic  
13 algorithms on the benchmark networks.

## 14 1 Introduction

15 A Bayesian network is a probabilistic graphical model that represents the structure of a joint probabil-  
16 ity distribution among random variables in a directed acyclic graph (DAG) [Pearl, 1988]. One class  
17 of associated computational problems is learning the structure of a Bayesian network from data. We  
18 focus on score-based Bayesian network structure learning for finding the DAG with a maximal score  
19 that depends on the data [Cooper and Herskovits, 1992, Cowell, 2001].

20 The Bayesian network learning problem is NP-hard [Chickering et al., 2004]; therefore, the standard  
21 methodology is using heuristic approaches. Many algorithms have been proposed to improve the  
22 accuracy and to reduce the running time. A search over the space of orderings [Teyssier and Koller,  
23 2005, Scanagatta et al., 2015] is one of the most successful heuristic approaches.

24 Annealing processors may contribute to finding a high-scoring network structure in a realistic  
25 timeframe. An annealing processor is expected to be an alternative hardware to von Neumann  
26 computers for quadratic unconstrained binary optimization (QUBO) problems. In particular, it is  
27 reported that complementary metal oxide semiconductor (CMOS) annealing processors already  
28 outperform conventional computers on the speed of solving max-cut problems [Gyoten et al., 2018].

29 We note that the bit capacity of an annealing processor is currently limited. Therefore, we need an  
30 efficient conversion method of Bayesian network structure learning into QUBO within the limited  
31 bits. Additionally, it is also important to show the lower bounds of penalty coefficients because the  
32 precision for the biases and variable couplers is limited.

33 Annealing processors are classified into the nearest neighbor type and the fully connected type  
 34 [Yamamoto, 2020]. While the coupling nodes of a nearest neighbor annealing processor is limited to  
 35 only between adjacent nodes, the coupling exists between arbitrary nodes of a fully coupled annealing  
 36 processor. Though the scalability of nearest neighbor annealing processors is high, it is necessary to  
 37 consider the additional bits for minor embedding [Choi, 2008, 2010].

38 O’Gorman et al. 2014 proposed a method to convert score-based Bayesian network structure learning  
 39 into QUBO that requires  $\mathcal{O}(n^2)$  bits for  $n$  random variables and a maximum parent set size  $m = 2$ .  
 40 They also demonstrated the sufficient lower bounds of penalty coefficients. However, when  $m \geq 3$ ,  
 41 the number of necessary auxiliary variables for a quadratization [Boros and Gruber, 2014] is at most  
 42  $\mathcal{O}(n(n-1)^{\frac{m}{2}})$ . This is a significant disadvantage for the current limited bit capacity of annealing  
 43 processors.

44 In this study, we propose an efficient conversion method based on the advanced identification of  
 45 candidate parent sets and their representation in the form of Cartesian products. We also provide a  
 46 greedy algorithm to decompose the candidate parent sets into the form of Cartesian products and  
 47 prove the sufficient lower bounds of penalty coefficients.

48 Experimental results on benchmark networks with 27 to 70 variables show that our conversion method  
 49 reduces the required bits significantly in comparison to the previous work [O’Gorman et al., 2014].  
 50 Our approach allows us to utilize the power of the second generation Fujitsu digital annealer, a fully  
 51 coupled CMOS annealing processor [Aramon et al., 2019]. We demonstrate that the digital annealer  
 52 consistently outperforms the ordering space search algorithms on the benchmark networks.

## 53 2 Background

### 54 2.1 Score-based Bayesian Network Structure Learning

55 The goal of score-based Bayesian network structure learning is to find a DAG with maximal score.  
 56 Given to random variables  $\mathcal{X} = (X_i)_{i=1}^n$  and a complete data set of  $N$  instances  $\mathcal{D} = \{D_1, \dots, D_N\}$ ,  
 57 we optimize the parent set  $\Pi_i$  of each random variable,

$$\Pi_1^*, \dots, \Pi_n^* = \arg \min_{\substack{\Pi_1, \dots, \Pi_n \subset \mathcal{X} \\ \mathcal{G} \in \text{DAG}}} \sum_{i=1}^n -\log S^{(i)}(\Pi_i | \mathcal{D}), \quad (1)$$

58 where  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ ,  $\mathcal{V} = \{1, \dots, n\}$ ,  $\mathcal{E} = \{(j, i) | j, i \in \{1, \dots, n\}, X_j \in \Pi_i\}$ , and  $S_i : \Pi_i \rightarrow \mathbb{R}$  is  
 59 a local score function corresponding to  $X_i$ . The Bayesian Dirichlet equivalent uniform (BDeu) score  
 60 [Buntine, 1991] is one of the commonly used scores,

$$S_{\text{BDeu}}^{(i)}(\Pi_i | \mathcal{D}) \equiv \prod_{j=1}^{\beta_i} \frac{\Gamma(\alpha_{i,j})}{\Gamma(N_{i,j} + \alpha_{i,j})} \prod_{k=1}^{\gamma_i} \frac{\Gamma(N_{i,j,k} + \alpha_{i,j,k})}{\Gamma(\alpha_{i,j,k})}, \quad (2)$$

61 where  $N = \sum_{j=1}^{\beta_i} N_{i,j}$ ,  $N_{i,j} = \sum_{k=1}^{\gamma_i} N_{i,j,k}$ ,  $\alpha_{i,j} = \sum_{k=1}^{\gamma_i} \alpha_{i,j,k}$ ,  $\beta_i$  is the number of joint states of  
 62  $\Pi_i$ ,  $\gamma_i$  is the number of states of  $X_i$ ,  $N_{i,j,k}$  is the number of cases of the parent set  $\Pi_i$  in its  $j$ -th state  
 63 and  $X_i$  in its  $k$ -th state,  $\alpha_{i,j,k} = \frac{\alpha}{\beta_i \gamma_i}$  is the hyperparameter of the Dirichlet function, and  $0 < \alpha \in \mathbb{R}$   
 64 is called equivalent sample size [Heckerman et al., 1995a].

### 65 2.2 Hamiltonian

66 The Hamiltonian, which is the objective function of an annealing processor, is a quadratic pseudo-  
 67 Boolean function,

$$H(\boldsymbol{\sigma}) = \sum_{i \in \mathcal{V}_{\text{AP}}} h_i \sigma_i + \sum_{(i,j) \in \mathcal{E}_{\text{AP}}} J_{i,j} \sigma_i \sigma_j, \quad (3)$$

68 where  $\boldsymbol{\sigma} = (\sigma_i)_{i=1}^{|\mathcal{V}_{\text{AP}}|} \in \mathbb{B}^{|\mathcal{V}_{\text{AP}}|}$ , the biases  $h_i \in \mathbb{R}$  for all  $i \in \mathcal{V}_{\text{AP}}$ , the couplers  $J_{i,j} \in \mathbb{R}$  for all  
 69  $(i, j) \in \mathcal{E}_{\text{AP}}$ , and the graph  $\mathcal{G}_{\text{AP}} = (\mathcal{V}_{\text{AP}}, \mathcal{E}_{\text{AP}})$ . Higher degree problems are reformed into quadratic  
 70 ones using auxiliary variables. This reformulation is called quadratization.

71 **Definition 1.** If a quadratic polynomial function  $g(\mathbf{v}, \mathbf{h})$  is a quadratization of a pseudo-Boolean  
 72 function  $f(\mathbf{v})$ , then  $f(\mathbf{v}) = \min_{\mathbf{h} \in \mathbb{B}^J} g(\mathbf{v}, \mathbf{h})$  for all  $\mathbf{v} \in \mathbb{B}^I$ .

73 Anthony et al. 2016 proved that every pseudo-Boolean function of  $I$  variables and of degree  $K$  has  
 74 a quadratization involving at most  $\mathcal{O}(I^{\frac{K}{2}})$  auxiliary variables. In particular, at most  $\mathcal{O}(2^{\frac{I}{2}})$  when  
 75  $K = I$ . It is well known that every pseudo-Boolean function can be uniquely represented as a  
 76 multilinear polynomial in its variables [Boros and Hamme, 2002].

### 77 2.3 Basic Conversion of Score-based Bayesian Network Structure Learning

78 Using  $n(n-1)$  bits to encode the paths into  $\mathbf{d} = ((d_{j,i})_{1 \leq j \leq n, j \neq i})_{i=1}^n \in \mathbb{B}^{n(n-1)}$  ( $d_{j,i} = 1$  if  
 79  $X_j$  is the parent of  $X_i$ ,  $d_{j,i} = 0$  otherwise) and  $\binom{n}{2}$  bits to encode the topological orders into  
 80  $\mathbf{r} = (r_{i,j})_{1 \leq i < j \leq n} \in \mathbb{B}^{\binom{n}{2}}$  ( $r_{i,j} = 0$  if the order of  $X_i$  is higher than  $X_j$ ,  $r_{i,j} = 1$  otherwise), it is  
 81 possible to represent eq. (1) on the Hamiltonian,

$$H_{\text{total}}(\mathbf{d}, \mathbf{r}) \equiv \sum_{i=1}^n H_{\text{score}}^{(i)}(\mathbf{d}, \mathbf{r}, i) + H_{\text{cycle}}(\mathbf{d}, \mathbf{r}). \quad (4)$$

82 The states of  $\mathbf{d}, \mathbf{r}$  are mapped one-to-one to the states of  $\Pi_i$ . Let  $\Pi_i = \pi^{(i)}(\mathbf{d}, \mathbf{r}, i)$  for all  $1 \leq i \leq n$ .  
 83 The local score of the Hamiltonian is

$$H_{\text{score}}^{(i)}(\mathbf{d}, \mathbf{r}, i) \equiv -\log S^{(i)}(\pi^{(i)}(\mathbf{d}, \mathbf{r}, i) \mid \mathcal{D}) + \log S^{(i)}(\phi \mid \mathcal{D}), \quad (5)$$

84 for all  $1 \leq i \leq n$ . The score function has a quadratization involving at most  $\mathcal{O}(n2^{\frac{n-1}{2}})$  auxiliary  
 85 variables. O’Gorman et al. 2014 added the maximum parent set size constraint to the Hamiltonian.  
 86 In this case, the number of auxiliary variables is at most  $\mathcal{O}(n(n-1)^{\frac{m}{2}})$ . The cycle constraint of the  
 87 Hamiltonian consists of the topological order constraint and the consistency constraint,

$$H_{\text{cycle}}(\mathbf{d}, \mathbf{r}) \equiv \sum_{1 \leq i < j < k \leq n} \delta_1 R(r_{i,j}, r_{j,k}, r_{i,k}) + \sum_{1 \leq i < j \leq n} \delta_2 (d_{i,j} r_{i,j} + d_{j,i} (1 - r_{i,j})), \quad (6)$$

88 where  $R(r_1, r_2, r_3) = r_1 r_2 (1 - r_3) + (1 - r_1)(1 - r_2) r_3$  for all  $r_1, r_2, r_3 \in \mathbb{B}$ . When the penalty  
 89 coefficients  $0 < \delta_1, \delta_2 \in \mathbb{R}$  are sufficiently large, the DAG constraint is satisfied indirectly through  
 90 the relationship of the paths  $\mathbf{d}$  and the topological order  $\mathbf{r}$ . If it holds that

$$\max\{0, \max_{\substack{1 \leq j^*, i^* \leq n \\ j^* \neq i^*}} \max_{\substack{\mathbf{d}, \mathbf{r} \in \mathbb{B}^{n(n-1)} \\ d_{j^*, i^*} = 1}} (H_{\text{score}}^{(i^*)}(\mathbf{d}, \mathbf{r}, i^*) - H_{\text{score}}^{(i^*)}(\mathbf{d}, \mathbf{r}, i^*))\} < \delta_1 < \frac{\delta_2}{n-2}, \quad (7)$$

91 then there is no cycle on the paths of the ground state, where  $\mathbf{d}^{(j^*, i^*)} = ((d_{j,i})_{1 \leq j \leq n, j \neq i})_{i=1}^n$ ,  
 92  $d_{j,i}^{(j^*, i^*)} = 0$  if  $(j, i) = (j^*, i^*)$ ,  $d_{j,i}^{(j^*, i^*)} = d_{j,i}$  otherwise. The computational cost to obtain the left  
 93 side of eq. (7) is at most  $\mathcal{O}(n^{m+1})$ . In particular, at most  $\mathcal{O}(n^2 2^{n-2})$  when  $m = n - 1$ .

## 94 3 Candidate Parent Set Decomposition

95 Parent set identification is a major technique to narrow the search space of structure optimization,  
 96 based on the relationship between parent sets and local scores under the DAG constraints [de Campos  
 97 and Ji, 2011, Correia et al., 2020]. The collection of candidate parent sets of a random variable  $X_i$  is  
 98  $\{W \subseteq \mathcal{X} \setminus \{X_i\} \mid W' \subset W \Rightarrow S^{(i)}(W' \mid \mathcal{D}) < S^{(i)}(W \mid \mathcal{D})\}$ . To reduce the required bits of the  
 99 score component of the Hamiltonian, we propose an efficient conversion method with the parent set  
 100 identification. We directly encode the candidate parent sets instead of using the paths  $\mathbf{d}$ .

101 Moreover, we decompose the candidate parent sets  $(W_{h,i})_{h=0}^{\lambda_i}$  of each random variable into the form  
 102 of Cartesian products as follows:

- 103 1. Decompose  $(W_{h,i})_{h=0}^{\lambda_i}$  into  $(W_{h,i} \cap Z_i)_{h=0}^{\lambda_i}, (W_{h,i} \cap (\mathcal{X} \setminus Z_i))_{h=0}^{\lambda_i}$ ,
- 104 2. Remove duplicates in the elements of  $(W_{h,i} \cap Z_i)_{h=0}^{\lambda_i}, (W_{h,i} \cap (\mathcal{X} \setminus Z_i))_{h=0}^{\lambda_i}$ ,

105 3. Store  $(W_{h,i} \cap Z_i)_{h=0}^{\lambda_{1,i}}, (W_{h,i} \cap (\mathcal{X} \setminus Z_i))_{h=0}^{\lambda_{2,i}}$  in  $(U_{h,i})_{h=0}^{\lambda_{1,i}}, (V_{h,i})_{h=0}^{\lambda_{2,i}}$ ,

106 where  $Z_i \subseteq \cup_{h=0}^{\lambda_i} W_{h,i}$ ,  $W_{0,i} = U_{0,i} = V_{0,i} = \phi$ ,  $\lambda_i, \lambda_{1,i}, \lambda_{2,i} \in \mathbb{N} \cup \{0\}$  for all  $1 \leq i \leq n$ . There  
107 is a clear relationship,

$$\{W_{0,i}, \dots, W_{\lambda_i,i}\} \subseteq \{U \cup V \mid (U, V) \in \{U_{0,i}, \dots, U_{\lambda_{1,i},i}\} \times \{V_{0,i}, \dots, V_{\lambda_{2,i},i}\}\}, \quad (8)$$

108 for all  $1 \leq i \leq n$ . Here, given that the Hamiltonian is a quadratic pseudo-Boolean function, we can  
109 represent the score against  $U_{h,i} \cup V_{h',i}$  by allocating  $U_{h,i}, V_{h',i}$  to two bits on the Hamiltonian. There-  
110 fore, it is possible to encode the candidate parent sets into the Hamiltonian using  $(U_{h,i})_{h=0}^{\lambda_{1,i}}, (V_{h,i})_{h=0}^{\lambda_{2,i}}$ .  
111 The number of required bits of the score component of the Hamiltonian is  $\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i})$ .

112 **Example 1.** An example of the candidate parent sets in the form of Cartesian products as follows:

$$\begin{aligned} \mathcal{X} &= \{X_1, X_2, X_3, X_4\}, \quad Z_i = \{X_1, X_2\}, \quad \lambda_i = 5, \quad \lambda_{1,i} = 2, \quad \lambda_{2,i} = 1 \\ (W_{h,i})_{h=0}^{\lambda_i} &= (\phi, \{X_1\}, \{X_1, X_2\}, \{X_3, X_4\}, \{X_1, X_3, X_4\}, \{X_1, X_2, X_3, X_4\}), \\ (W_{h,i} \cap Z_i)_{h=0}^{\lambda_i} &= (\phi, \{X_1\}, \{X_1, X_2\}, \phi, \{X_1\}, \{X_1, X_2\}), \\ (W_{h,i} \cap (\mathcal{X} \setminus Z_i))_{h=0}^{\lambda_i} &= (\phi, \phi, \phi, \{X_3, X_4\}, \{X_3, X_4\}, \{X_3, X_4\}), \\ (U_{h,i})_{h=0}^{\lambda_{1,i}} &= (\phi, \{X_1\}, \{X_1, X_2\}), \quad (V_{h,i})_{h=0}^{\lambda_{2,i}} = (\phi, \{X_3, X_4\}). \end{aligned}$$

113 We optimize  $Z_i \subseteq \cup_{h=0}^{\lambda_i} W_{h,i}$  to minimize  $\lambda_{1,i} + \lambda_{2,i}$ . However, it is often infeasible to search all  
114 elements of the power set  $\mathcal{P}(\cup_{h=0}^{\lambda_i} W_{h,i})$ . Therefore, we heuristically search  $Z_i$  adding elements one  
by one, as algorithm 1. The computational cost is at most  $\mathcal{O}(\lambda_i^3)$  for all  $1 \leq i \leq n$ .

---

**Algorithm 1** Greedy Candidate Parent Set Decomposition

---

1: **Input:**  $(W_{h,i})_{h=0}^{\lambda_i}$  **Output:**  $Z$  **Initialize:**  $\lambda \leftarrow \lambda_i, Z' \leftarrow \phi, Z \leftarrow \phi$ .  
2: **for**  $d = 1$  to  $|\cup_{h=0}^{\lambda_i} W_{h,i}| - 1$  **do**  
3:   **for**  $X$  in  $\cup_{h=0}^{\lambda_i} W_{h,i} \setminus Z$  **do**  
4:     **if**  $\lambda_{1,i} + \lambda_{2,i} < \lambda$  for  $Z_i = Z \cup \{X\}$  **then**  $\lambda \leftarrow \lambda_{1,i} + \lambda_{2,i}, Z' \leftarrow Z \cup \{X\}$ .  
5:   **if**  $Z \neq Z'$  **then**  $Z \leftarrow Z'$  **else break**

---

115

116 **Example 2.** An example of the bit reduction flow of algorithm 1 is as follows:

$$\begin{aligned} Z_i &= \phi, \lambda_{1,i} = 0, \lambda_{2,i} = 5 : (\phi) \times (\phi, \{X_1\}, \{X_1, X_2\}, \{X_3, X_4\}, \{X_1, X_3, X_4\}, \{X_1, X_2, X_3, X_4\}), \\ Z_i &= \{X_1\}, \lambda_{1,i} = 1, \lambda_{2,i} = 3 : (\phi, \{X_1\}) \times (\phi, \{X_2\}, \{X_3, X_4\}, \{X_2, X_3, X_4\}), \\ Z_i &= \{X_1, X_2\}, \lambda_{1,i} = 2, \lambda_{2,i} = 1 : (\phi, \{X_1\}, \{X_1, X_2\}) \times (\phi, \{X_3, X_4\}). \end{aligned}$$

## 117 4 Efficient Conversion of Score-based Bayesian Network Structure Learning

118 We make  $(U_{h,i})_{h=0}^{\lambda_{1,i}}, (V_{h,i})_{h=0}^{\lambda_{2,i}}$  correspond to  $(p_{h,i})_{h=0}^{\lambda_{1,i}}, (q_{h,i})_{h=0}^{\lambda_{2,i}}$  one-to-one, where  $p_{h,i}, q_{h',i} \in \mathbb{B}$   
119 for all  $0 \leq h \leq \lambda_{1,i}, 0 \leq h' \leq \lambda_{2,i}, 1 \leq i \leq n$ . To identify the parent sets, we use the one-to-one  
120 correspondence constraint that  $\sum_{h=0}^{\lambda_{1,i}} p_{h,i} = \sum_{h=0}^{\lambda_{2,i}} q_{h,i} = 1$  for all  $1 \leq i \leq n$ . The Hamiltonian  
121 consists of the score component, the one-to-one correspondence constraint, and the cycle constraint,

$$H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) \equiv \sum_{i=1}^n (H_{\text{score}}^{*(i)}(\mathbf{p}, \mathbf{i}, \mathbf{q}, \mathbf{i}) + H_{\text{one}}^{*(i)}(\mathbf{p}, \mathbf{i}, \mathbf{q}, \mathbf{i})) + H_{\text{cycle}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}), \quad (9)$$

122 where  $\mathbf{p} = ((p_{h,i})_{h=0}^{\lambda_{1,i}})_{i=1}^n$ ,  $\mathbf{q} = ((q_{h,i})_{h=0}^{\lambda_{2,i}})_{i=1}^n$ . Under the one-to-one correspondence constraint,  
123 we can represent the paths among random variables indirectly using  $\mathbf{p}, \mathbf{q}$  without additional auxiliary  
124 variables,

$$d_{j,i}^* \equiv \sum_{\substack{1 \leq h \leq \lambda_{1,i} \\ X_j \in U_{h,i}}} p_{h,i} + \sum_{\substack{1 \leq h \leq \lambda_{2,i} \\ X_j \in V_{h,i}}} q_{h,i}, \quad (10)$$

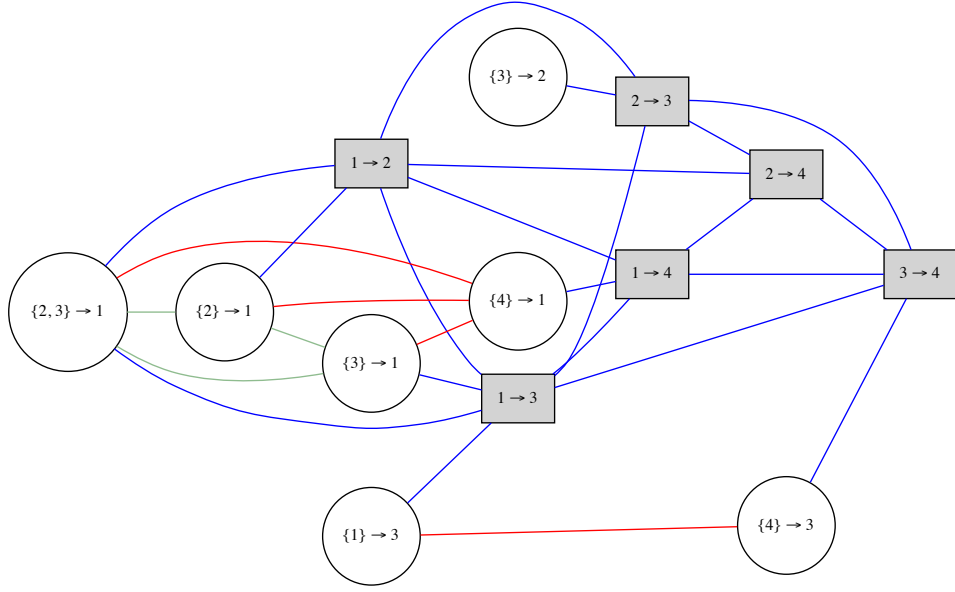


Figure 1: An example of bit allocation for our conversion method.  $n = 4$ ,  $\lambda_{1,1} = 3$ ,  $\lambda_{2,1} = 1$ ,  $\lambda_{1,2} = 1$ ,  $\lambda_{2,2} = 0$ ,  $\lambda_{1,3} = 1$ ,  $\lambda_{2,3} = 1$ ,  $\lambda_{1,4} = 0$ ,  $\lambda_{2,4} = 0$ ,  $U_{1,1} = \{2, 3\}$ ,  $U_{2,1} = \{3\}$ ,  $U_{3,1} = \{2\}$ ,  $V_{1,1} = \{4\}$ ,  $U_{1,2} = \{3\}$ ,  $U_{1,3} = \{1\}$ ,  $V_{1,3} = \{4\}$ . Circle :  $\mathbf{p}, \mathbf{q}$ . Square :  $\mathbf{r}$ . Red lines include in the score component of the Hamiltonian, a green line in the one-to-one correspondence constraint, and blue lines in the cycle constraint.

125 for all  $1 \leq j, i \leq n$ . Figure 1 is an example of bit allocation using our conversion method. The  
 126 number of bits required in our conversion method is  $\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i}) + \binom{n}{2}$ . Note that we do not  
 127 directly encode  $p_{0,i}, q_{0,i}$  on the Hamiltonian.

128 **Score Component.** The local score component of the Hamiltonian is

$$H_{\text{score}}^{*(i)}(\mathbf{p}, \mathbf{i}, \mathbf{q}, \mathbf{i}) \equiv \sum_{h=1}^{\lambda_{1,i}} s_{1,h,i} p_{h,i} + \sum_{h=1}^{\lambda_{2,i}} s_{2,h,i} q_{h,i} + \sum_{h=1}^{\lambda_{1,i}} \sum_{h'=1}^{\lambda_{2,i}} t_{h,h',i} p_{h,i} q_{h',i}, \quad (11)$$

129 for all  $1 \leq i \leq n$ . We can get these coefficients by solving simultaneous equations under the  
 130 one-to-one correspondence constraint,  $s_{1,h,i} = -\log S^{(i)}(U_{h,i} | \mathcal{D}) + \log S^{(i)}(\phi | \mathcal{D})$ ,  $s_{2,h,i} =$   
 131  $-\log S^{(i)}(V_{h,i} | \mathcal{D}) + \log S^{(i)}(\phi | \mathcal{D})$ ,  $t_{h,h',i} = -\log S^{(i)}(U_{h,i} \cup V_{h',i} | \mathcal{D}) + \log S^{(i)}(U_{h,i} | \mathcal{D}) +$   
 132  $\log S^{(i)}(V_{h',i} | \mathcal{D}) - \log S^{(i)}(\phi | \mathcal{D})$ .

133 **One-to-One Correspondence Constraint.** We penalize the connection among bits to select each  
 134 element from  $(U_{h,i})_{h=0}^{\lambda_{1,i}}, (V_{h,i})_{h=0}^{\lambda_{2,i}}$ ,

$$H_{\text{one}}^{*(i)}(\mathbf{p}, \mathbf{i}, \mathbf{q}, \mathbf{i}) \equiv \sum_{1 \leq h < h' \leq \lambda_{1,i}} \xi_{1,i} p_{h,i} p_{h',i} + \sum_{1 \leq h < h' \leq \lambda_{2,i}} \xi_{2,i} q_{h,i} q_{h',i}, \quad (12)$$

135 for all  $1 \leq i \leq n$ , where the penalty coefficient  $0 < \xi_{1,i}, \xi_{2,i} \in \mathbb{R}$ . If  $\xi_{1,i}, \xi_{2,i}$  is sufficient large,  
 136  $\sum_{h=0}^{\lambda_{1,i}} p_{h,i} = \sum_{h=0}^{\lambda_{2,i}} q_{h,i} = 1$  is induced indirectly.

137 **Cycle Constraint.** Compared to eq. (6), the cycle constraint of the Hamiltonian is

$$H_{\text{cycle}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) \equiv \sum_{1 \leq i < j < k \leq n} \delta_1^* R(r_{i,j}, r_{j,k}, r_{i,k}) + \sum_{1 \leq i < j \leq n} \delta_2^* (d_{i,j}^* r_{i,j} + d_{j,i}^* (1 - r_{i,j})), \quad (13)$$

138 where the penalty coefficients  $0 < \delta_1^*, \delta_2^* \in \mathbb{R}$ . By setting  $\delta_1^*, \delta_2^*$  appropriately, we can prevent the  
 139 cycle from occurring.

## 140 5 Sufficient Lower Bounds of Penalty Coefficients

141 We demonstrate the sufficient lower bounds of penalty coefficients. The basic idea is that we find the  
 142 range of penalty coefficients so that the change in return value of the Hamiltonian is negative when  
 143 the input state changes to the state we desire to induce.

144 **One-to-One Correspondence Constraint.** We consider to decrease the value of  $\sum_{h=1}^{\lambda_{1,i^*}} p_{h,i^*}$  by  
 145 one until it reaches 1. In the case of  $p_{h^*,i^*} = 1, \sum_{h=1}^{\lambda_{1,i^*}} p_{h,i^*} > 1$ , it holds that  $H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) -$   
 146  $H_{\text{total}}^*(\mathbf{p}^{(h^*,i^*)}, \mathbf{q}, \mathbf{r}) \geq \xi_{1,i^*} + s_{1,h^*,i^*} + \sum_{h=1}^{\lambda_{2,i}} t_{h^*,h,i^*} q_{h,i^*}$ , where  $\mathbf{p}^{(h^*,i^*)} = ((p_{h,i}^{(h^*,i^*)})_{h=0}^{\lambda_{1,i}})_{i=1}^n$   
 147 and  $p_{h,i}^{(h^*,i^*)} = 0$  if  $(h,i) = (h^*,i^*), p_{h,i}^{(h^*,i^*)} = p_{h,i}$  otherwise. Considering the case where  $\mathbf{p}$  and  $\mathbf{q}$   
 148 are swapped in the above, if  $\xi_{1,i}, \xi_{2,i}$  satisfy that

$$\max_{0 \leq h \leq \lambda_{1,i}} (-s_{1,h,i} - \sum_{h'=1}^{\lambda_{2,i}} \min\{0, t_{h,h',i}\}) < \xi_{1,i}, \quad (14)$$

$$\max_{0 \leq h \leq \lambda_{2,i}} (-s_{2,h,i} - \sum_{h'=1}^{\lambda_{1,i}} \min\{0, t_{h',h,i}\}) < \xi_{2,i}, \quad (15)$$

149 for all  $1 \leq i \leq n$ , then the grand state does not violate the one-to-one correspondence constraint.  
 150 The computational cost to obtain the left side of eq. (14) and eq. (15) is at most  $\mathcal{O}(\lambda_{1,i}\lambda_{2,i})$  for all  
 151  $1 \leq i \leq n$ .

152 **Cycle Constraint.** We consider four patterns of  $(r_{i^*,j^*}, d_{j^*,i^*}^*, d_{i^*,j^*}^*)$  violating the consistency  
 153 constraint. It is assumed that  $X_{j^*} \in U_{h^*,i^*}, X_{j^*} \notin U_{h^{**},i^*} \subset U_{h^*,i^*}, p_{h^*,i^*} = 1, p_{h^{**},i^*} = 0$ .  
 154 In the case of  $(0, 1, 0)$ , it holds that  $H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) - H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}^{(i^*,j^*)}) \geq \delta_2^* - (n-2)\delta_1^*$ ,  
 155 where  $\mathbf{r}^{(i^*,j^*)} = (r_{i,j}^{(i^*,j^*)})_{1 \leq i < j \leq n}$ , and  $r_{i,j}^{(i^*,j^*)} = 1 - r_{i,j}$  if  $(i,j) = (i^*,j^*), r_{i,j}^{(i^*,j^*)} = r_{i,j}$   
 156 otherwise. Similarly, it is possible to consider the case of  $(1, 0, 1)$ . In the case of  $(0, 1, 1)$ ,  
 157 it holds that  $H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) - H_{\text{total}}^*(\mathbf{p}^{(h^*,h^{**},i^*)}, \mathbf{q}, \mathbf{r}) \geq \delta_2^* + s_{1,h^*,i^*} + \sum_{h=1}^{\lambda_{2,i}} t_{h^*,h,i^*} q_{h,i^*} -$   
 158  $s_{1,h^{**},i^*} - \sum_{h=1}^{\lambda_{2,i}} t_{h^{**},h,i^*} q_{h,i^*}$ , where  $\mathbf{p}^{(h^*,h^{**},i^*)} = ((p_{h,i}^{(h^*,h^{**},i^*)})_{h=0}^{\lambda_{1,i}})_{i=1}^n$ , and  $p_{h,i}^{(h^*,h^{**},i^*)} = 0$   
 159 if  $(h,i) = (h^*,i^*), p_{h,i}^{(h^*,h^{**},i^*)} = 1$  if  $(h,i) = (h^{**},i^*), p_{h,i}^{(h^*,h^{**},i^*)} = p_{h,i}$  otherwise. Sim-  
 160 ilarly, it is possible to consider the case of  $(1, 1, 1)$ . These results suggest the relationship of  
 161  $\delta_1^*, \delta_2^*$  to induce the consistency constraint. Here, based on theorem 1, we consider a strategy  
 162 to repeat picking up one element from  $\mathbf{r}$  and switching its value until  $H_{\text{trans}}(\mathbf{r}) = 0$ . It is as-  
 163 sumed that  $H_{\text{trans}}(\mathbf{r}) > H_{\text{trans}}(\mathbf{r}^{(i^*,j^*)})$ . In the case of  $(1, 1, 0)$ , it holds that  $H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) -$   
 164  $H_{\text{total}}^*(\mathbf{p}^{(h^*,h^{**},i^*)}, \mathbf{q}, \mathbf{r}^{(i^*,j^*)}) \geq \delta_1^* + s_{1,h^*,i^*} + \sum_{h=1}^{\lambda_{2,i}} t_{h^*,h,i^*} q_{h,i^*} - s_{1,h^{**},i^*} - \sum_{h=1}^{\lambda_{2,i}} t_{h^{**},h,i^*} q_{h,i^*}$ .  
 165 Similarly, it is possible to consider the case of  $(0, 0, 1)$ . In the case of  $(1, 0, 0)$  or  $(0, 0, 0)$ , it holds  
 166 that  $H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}) - H_{\text{total}}^*(\mathbf{p}, \mathbf{q}, \mathbf{r}^{(i^*,j^*)}) \geq \delta_1^*$ . These results suggest the lower bound of  $\delta_1^*$   
 167 to induce the topological order constraint. Considering the case where  $\mathbf{p}$  and  $\mathbf{q}$  are swapped in the  
 168 above, if  $\delta_1^*, \delta_2^*$  satisfy that

$$\max_{1 \leq i \leq n} \max\{\eta_{1,i}, \eta_{2,i}\} < \delta_1^* < \frac{\delta_2^*}{n-2}, \quad (16)$$

$$\eta_{1,i} \equiv \max_{1 \leq j \leq n} \max_{0 \leq h \leq \lambda_{1,i}} \max_{\substack{X_j \in U_{h,i} \\ X_j \notin U_{h',i} \subset U_{h,i}}} \max_{0 \leq h'' \leq \lambda_{2,i}} (-s_{1,h,i} - t_{h,h'',i} + s_{1,h',i} + t_{h',h'',i}),$$

$$\eta_{2,i} \equiv \max_{1 \leq j \leq n} \max_{0 \leq h \leq \lambda_{2,i}} \max_{\substack{X_j \in V_{h,i} \\ X_j \notin V_{h',i} \subset V_{h,i}}} \max_{0 \leq h'' \leq \lambda_{1,i}} (-s_{2,h,i} - t_{h'',h,i} + s_{2,h',i} + t_{h',h'',i}),$$

169 for  $n \geq 3$ , then the grand state does not violate the cycle constraint under the one-to-one cor-  
 170 respondence constraint. The computational cost to obtain the left side of eq. (16) is at most  
 171  $\mathcal{O}(\sum_{i=1}^n n\lambda_{1,i}\lambda_{2,i}(\lambda_{1,i} + \lambda_{2,i}))$ .

172 **Theorem 1.** If it holds that  $H_{\text{trans}}(\mathbf{r}) \equiv \sum_{1 \leq i < j < k \leq n} R(r_{i,j}, r_{j,k}, r_{i,k}) > 0$ , then there exists  
 173 at least one index pair  $1 \leq i^* < j^* \leq n$  which satisfy  $H_{\text{trans}}(\mathbf{r}) > H_{\text{trans}}(\mathbf{r}^{(i^*,j^*)})$ , where  
 174  $R(r_1, r_2, r_3) = r_1 r_2 (1 - r_3) + (1 - r_1)(1 - r_2)r_3$  for all  $r_1, r_2, r_3 \in \mathbb{B}$ ,  $\mathbf{r} = (r_{i,j})_{1 \leq i < j \leq n} \in \mathbb{B}^{\binom{n}{2}}$ ,  
 175  $\mathbf{r}^{(i^*,j^*)} = (r_{i,j}^{(i^*,j^*)})_{1 \leq i < j \leq n}$ , and  $r_{i,j}^{(i^*,j^*)} = 1 - r_{i,j}$  if  $(i,j) = (i^*,j^*), r_{i,j}^{(i^*,j^*)} = r_{i,j}$  otherwise.

Table 1: The benchmark networks from Bayesian network repository.

Name	$n$	$m$	$\sum_{i=1}^n  \Pi_i $	$\sum_{i=1}^n \beta_i(\gamma_i - 1)$	$\sum_{i=1}^n \lambda_i^*$		
					$N = 100$	$N = 1000$	$N = 10000$
insurance	27	3	52	984	353	883	4036
water	32	5	66	10083	165	216	735
alarm	37	4	46	509	1829	2272	9081
barley	48	4	84	114005	181	310	1552
hailfinder	56	4	66	2656	144	692	4277
hepar2	70	6	123	1453	4837	665	4782

\* The average for 10 simulated datasets.

176 *Proof.* It does not lose the generality by considering the case of  $(r_{1,2}, r_{2,3}, r_{1,3}) = (1, 1, 0)$ . Here, it  
177 holds that  $R(r_{1,2}, r_{2,3}, r_{1,3}) - R(1 - r_{1,2}, r_{2,3}, r_{1,3}) + R(r_{1,2}, r_{2,3}, r_{1,3}) - R(r_{1,2}, 1 - r_{2,3}, r_{1,3}) +$   
178  $R(r_{1,2}, r_{2,3}, r_{1,3}) - R(r_{1,2}, r_{2,3}, 1 - r_{1,3}) = 3$ . Additionally, it holds that  $R(r_{1,2}, r_{2,i}, r_{1,i}) - R(1 -$   
179  $r_{1,2}, r_{2,i}, r_{1,i}) + R(r_{2,3}, r_{3,i}, r_{2,i}) - R(1 - r_{2,3}, r_{3,i}, r_{2,i}) + R(r_{1,3}, r_{3,i}, r_{1,i}) - R(1 - r_{1,3}, r_{3,i}, r_{1,i}) =$   
180  $0$  for all  $3 < i$ . Therefore, it holds that  $H_{\text{trans}}(\mathbf{r}) - H_{\text{trans}}(\mathbf{r}^{(1,2)}) + H_{\text{trans}}(\mathbf{r}) - H_{\text{trans}}(\mathbf{r}^{(2,3)}) +$   
181  $H_{\text{trans}}(\mathbf{r}) - H_{\text{trans}}(\mathbf{r}^{(1,3)}) = 3$ . From this result, it holds that  $H_{\text{trans}}(\mathbf{r}) - H_{\text{trans}}(\mathbf{r}^{(i^*, j^*)}) > 0$  for  
182 at least one index pair  $(i^*, j^*) \in \{(1, 2), (2, 3), (1, 3)\}$ .  $\square$

## 183 6 Experimental Results

184 To validate the performance of our approach, we use 10 simulated datasets for each instance size  
185  $N = 100, 1000, 10000$  and each benchmark network. The benchmark networks are discrete networks  
186 from Bayesian network repository<sup>1</sup>. The score function is the BDeu score with  $\alpha = 1$ . It is often  
187 infeasible to identify exact candidate parent sets by searching the power set  $\mathcal{P}(\mathcal{X} \setminus \{X_i\})$  in a realistic  
188 timeframe. We use the candidate parent sets from algorithm 2. Note that the candidate parent sets  
189 depend on the heuristic search algorithms, but we do not focus on their performance in this study.  
190 Table 1 displays the information of benchmark networks. The code to replicate each experiment in  
191 this paper is available<sup>2</sup>.

---

### Algorithm 2 Greedy Candidate Parent Set Identification

---

```

1: Input:  $\mathcal{D}, i, m$  Output:  $\mathcal{L}$  Initialize:  $\mathcal{L} \leftarrow \{\phi\}, \mathcal{L}' \leftarrow \{\phi\}, \mathcal{L}'' \leftarrow \phi$ 
2: for  $d = 1$  to  $m$  do
3:   for  $W$  in  $\mathcal{L}'$  do
4:     for  $X$  in  $\mathcal{X} \setminus \{X_i\} \setminus W$  do
5:       if  $S_i(W' \mid \mathcal{D}) < S_i(W \cup \{X\} \mid \mathcal{D})$  for all  $W' \subset W \cup \{X\}, W' \in \mathcal{L}$  then
6:          $\mathcal{L}'' \leftarrow \mathcal{L}'' \cup \{W \cup \{X\}\}$ .
7:       if  $\mathcal{L}'' \neq \phi$  then  $\mathcal{L} \leftarrow \mathcal{L} \cup \mathcal{L}'', \mathcal{L}' \leftarrow \mathcal{L}'', \mathcal{L}'' \leftarrow \phi$  else break
8:   for  $W$  in  $\mathcal{L}$  do
9:     if there exist  $W' \subset W$  that satisfies  $S_i(W \mid \mathcal{D}) \leq S_i(W' \mid \mathcal{D})$  then  $\mathcal{L} \leftarrow \mathcal{L} \setminus \{W\}$ .

```

---

### 192 6.1 Number of Required Bits for Score Component

193 In comparison to the existing method [O’Gorman et al., 2014], we reduce the number of required bits  
194 for the score component by encoding the candidate parent sets directly. While  $\sum_{i=1}^n \lambda_i$  candidate  
195 parent sets is encoded in our approach,  $n(n-1)$  paths plus at most  $\mathcal{O}(n(n-1)^{\frac{m}{2}})$  auxiliary variables  
196 for  $m > 2$  in the existing method. The left side of table 2 shows the reduction rate of the number of  
197 required bits for the score component. Moreover, we reduce the number of required bits for the score  
198 component to  $\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i})$  by decomposing the candidate parent sets in the form of Cartesian

<sup>1</sup><https://www.bnlearn.com/bnrepository/>

<sup>2</sup>See supplemental material.

Table 2: The reduction rate of the number of required bits for score component.

Name	$\sum_{i=1}^n \lambda_i / n(n-1)^{\frac{m}{2}}$ *			$\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i}) / \sum_{i=1}^n \lambda_i$ *		
	$N = 100$	$N = 1000$	$N = 10000$	$N = 100$	$N = 1000$	$N = 10000$
insurance	0.09873	0.24677	1.12742	0.61367	0.47285	0.32476
water	0.00097	0.00126	0.00429	0.72680	0.70588	0.44014
alarm	0.03814	0.04738	0.18938	0.45332	0.35537	0.21617
barley	0.00171	0.00292	0.01464	0.76717	0.75538	0.54149
hailfinder	0.00085	0.00409	0.02525	0.82773	0.60178	0.33365
hepar2	0.00021	0.00003	0.00021	0.49694	0.63346	0.31284

\* The average ratio for 10 simulated datasets.

Table 3: The number of required bits for fully coupled and nearest neighbor annealing processors.

Name	$\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i}) + \binom{n}{2}$ *			$\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i})(\lambda_{1,i} + \lambda_{2,i} + 1) + \binom{n}{2}$ *		
	$N = 100$	$N = 1000$	$N = 10000$	$N = 100$	$N = 1000$	$N = 10000$
insurance	566	767	1661	3375	9023	85482
water	613	648	820	1434	1720	5881
alarm	1489	1472	2628	36761	27985	169004
barley	1247	1362	1968	6758	3446	24796
hailfinder	1659	1957	2967	2212	7084	80578
hepar2	4777	2836	3910	449916	9164	136939

\* The average ratio for 10 simulated datasets.

199 products. The right side of table 2 shows that algorithm 1 reduces the number of required bits for the  
 200 score component although there is some variation among the networks.

## 201 6.2 Selection of Annealing Processor

202 From the following discussion, the Fujitsu digital annealer is suitable for our approach from the  
 203 viewpoint of bit capacity.

204 **Fully Connected Type.** To the best of our knowledge, the bit capacity of the Fujitsu digital annealer  
 205 is the largest in fully coupled annealing processors. The second generation Fujitsu digital annealer  
 206 can deal with problems on a scale of 8192 bits [Matsubara et al., 2020]. The left side of table 3 shows  
 207 that it is possible to encode all the logical conversion results for benchmark networks to the circuit of  
 208 the digital annealer within bit capacity.

209 **Nearest Neighbor Type.** The number of additional bits required for minor embedding depends on  
 210 the design of the hardware graphs. Oku et al. 2019 proposed a heuristic minor embedding algorithm  
 211 for the Hitachi CMOS annealing machine [Masanao et al., 2010]. Using this algorithm, the number  
 212 of required physical spins when embedding a fully connected graph is  $I^2 + I$  for  $I$  variables. The  
 213 conversion method proposed in this study has  $n$  local fully connected graphs on  $p, q$ . Therefore, the  
 214 number of required physical spins must be at least  $\sum_{i=1}^n (\lambda_{1,i} + \lambda_{2,i})(\lambda_{1,i} + \lambda_{2,i} + 1) + \binom{n}{2}$ . From  
 215 the right side of table 3, it is currently infeasible to encode logical conversion results for at least some  
 216 networks to the circuit of CMOS annealing machine within its 102400 nodes [Sugie et al., 2021]. As  
 217 far as we know, the bit capacity of the Hitachi CMOS annealing machine is the largest in nearest  
 218 neighbor annealing processors.

## 219 6.3 Score Maximization

220 We demonstrate the performance of Fujitsu digital annealer for score-based Bayesian network  
 221 structure learning using the conversion results of  $N = 10000$  simulated datasets. The running time  
 222 for each simulated dataset is 6000 [s].



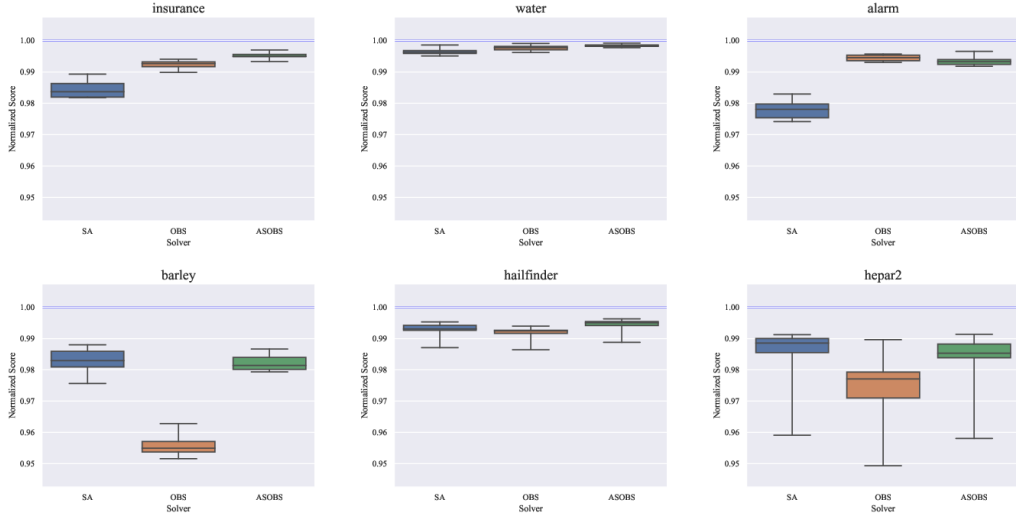


Figure 2: Results of score maximization by the baseline algorithms. For each simulated dataset and each baseline algorithm, we normalized  $\sum_{i=1}^n (\log S^{(i)}(\Pi_i | \mathcal{D}) - \log S^{(i)}(\phi | \mathcal{D}))$  by dividing it by the corresponding value of the Fujitsu digital annealer. In this experiment, we used the second-generation Fujitsu digital annealer. SA : simulated annealing, OBS : ordering-based search, ASOBS : acyclic selection ordering-based search.

223 **Baselines.** We compare the results obtained by the digital annealer with those of three heuristic  
 224 algorithms. One algorithm is the simulated annealing algorithm [Heckerman et al., 1995b] with a  
 225 QUBO same as the one encoded into the digital annealer. Other algorithms are the ordering space  
 226 search algorithms, i.e., ordering-based search and acyclic selection ordering-based search. For a fair  
 227 comparison, the running time of the simulated annealing algorithm for each simulated dataset is 6000  
 228 [s] and that of the ordering space search algorithms is 6000 [s] plus the running time of algorithm 1.  
 229 The computing environment is Microsoft Windows 10 Pro, 3.6 GHz Intel Core i9 processor, and 64  
 230 GB memory.

231 **Result.** Figure 2 shows that the digital annealer is better than all the baselines for all the simulated  
 232 datasets from all the benchmark networks.

## 233 7 Conclusion

234 We proposed a novel approach of converting a score-based Bayesian network structure learning  
 235 into QUBO. The essence of this approach lies in reducing the number of required bits through the  
 236 advanced identification of candidate parent sets and their representation as Cartesian products. The  
 237 Fujitsu digital annealer with our conversion method improved the BDeu score for 27 to 70 variables  
 238 benchmark networks over existing methods. The bit capacity limitation of annealing processor is  
 239 being relaxed rapidly<sup>3</sup>. Though our approach is still a disadvantage for larger-scale networks, we  
 240 expect that our proposed algorithms will be effectively applied to larger-scale score-based Bayesian  
 241 network structure learning in the near future.

242 **Potential Negative Societal Impacts.** The development of annealing processor technology could  
 243 have an impact on various industry fields. However, the number of companies that have commer-  
 244 cialized the API usage of annealing processors is still small. Therefore, there is a concern that the  
 245 market of annealing processors will not work well and the disparities among stakeholders will be  
 246 widen. Researchers are required to properly evaluate the value of technology and communicate it to  
 247 the business side.

<sup>3</sup>Fujitsu announced that they achieved a megabit-class performance with digital annealer

## 248 References

- 249 Judea Pearl, editor. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan  
250 Kaufmann, 1988.
- 251 Gregory F. Cooper and Edward Herskovits. A bayesian method for the induction of probabilistic networks from  
252 data. *Journal of Machine Learning*, 9(4), 1992.
- 253 Robert G. Cowell. Conditions under which conditional independence and scoring methods lead to identical  
254 selection of bayesian network models. In Jack Breese and Daphne Koller, editors, *Proceedings of the 17th*  
255 *Conference on Uncertainty in Artificial Intelligence*, pages 91–97. Morgan Kaufmann Publishers, 2001.
- 256 David M. Chickering, David Heckerman, and Christopher Meek. Large-sample learning of bayesian networks is  
257 np-hard. *Journal of Machine Learning Research*, 20:1287–1330, 2004.
- 258 Marc Teysier and Daphne Koller. Ordering-based search: A simple and effective algorithm for learning bayesian  
259 networks. In *Proceedings of the 21th Conference on Uncertainty in Artificial Intelligence*, pages 584–590,  
260 2005.
- 261 Mauro Scanagatta, Cassio Polpo de Campos, Giorgio Corani, and Marco Zaffalon. Learning bayesian networks  
262 with thousands of variables. In *Proceedings of the 28th International Conference on Neural Information*  
263 *Processing Systems*, pages 1864–1872, 2015.
- 264 Hidenori Gyoten, Masayuki Hiramoto, and Takashi Sato. Area efficient annealing processor for ising model  
265 without random number generator. *IEICE Transactions on Information and Systems*, E101.D(2):314–323,  
266 2018.
- 267 Kasho Yamamoto. *Research on Annealing Processors for Large-Scale Combinatorial Optimization Problems*.  
268 PhD thesis, Graduate School of Information Science and Technology Hokkaido University, 2020.
- 269 Vicky Choi. Minor-embedding in adiabatic quantum computation: I. the parameter setting problem. *Quantum*  
270 *Information Processing*, 7:193–209, 2008.
- 271 Vicky Choi. Minor-embedding in adiabatic quantum computation: Ii. minor-universal graph design. *Quantum*  
272 *Information Processing*, 10:343–352, 2010.
- 273 Bryan A. O’Gorman, Alejandro Perdomo-Ortiz, Ryan Babbush, Alan Aspuru-Guzik, and Vadim Smelyanskiy.  
274 Bayesian network structure learning using quantum annealing. *The European Physical Journal Special Topics*,  
275 225(1), 2014.
- 276 Endre Boros and Aritanan Gruber. On quadratization of pseudo-boolean functions. *CoRR*, abs/1404.6538, 2014.
- 277 Malihah Aramon, Gili Rosenberg, Elisabetta Valiante, Toshiyuki Miyazawa, Hirotaka Tamura, and Helmut G.  
278 Katzgraber. Physics-inspired optimization for quadratic unconstrained problems using a digital annealer.  
279 *Frontiers in Physics*, 7(48), 2019.
- 280 Wray Buntine. Theory refinement of bayesian networks. In *Proceedings of the 7th Conference on Uncertainty*  
281 *in Artificial Intelligence*, pages 52–60, 1991.
- 282 David Heckerman, Dan Geiger, and David M. Chickering. Learning bayesian networks: The combination of  
283 knowledge and statistical data. *Journal of Machine Learning*, 20(3):197–243, 1995a.
- 284 Martin Anthony, Endre Boros, Yves Crama, and Aritanan Gruber. Quadratic reformulations of nonlinear binary  
285 optimization problems. *Mathematical Programming*, 162(1):115–144, 2016.
- 286 Endre Boros and Peter L. Hamme. Pseudo-boolean optimization. *Journal of Discrete Applied Mathematics*, 123:  
287 155–225, 2002.
- 288 Cassio P. de Campos and Qiang Ji. Efficient structure learning of bayesian networks using constraints. *Journal*  
289 *of Machine Learning Research*, 12:663–689, 2011.
- 290 Alvaro H. C. Correia, James Cussens, and Cassio de Campos. On pruning for score-based bayesian network  
291 structure learning. *Journal of Machine Learning Research*, 108:2709–2718, 2020.
- 292 Satoshi Matsubara, Motomu Takatsu, Toshiyuki Miyazawa, Takayuki Shibasaki, Yasuhiro Watanabe, Kazuya  
293 Takemoto, and Hirotaka Tamura. Digital annealer for high-speed solving of combinatorial optimization  
294 problems and its applications. In *Proceedings of the 25th Asia and South Pacific Design Automation*  
295 *Conference*, pages 667–672, 2020.
- 296 David Eppstein. Finding large clique minors is hard. *Graph Algorithms and Applications*, 13(2):197–204, 2009.

- 297 Daisuke Oku, Kotaro Terada, Masato Hayashi, Masanao Yamaoka, Shu Tanaka, and Nozomu Togawa. A  
 298 fully-connected ising model embedding method and its evaluation for cmos annealing machines. *IEICE*  
 299 *Transactions on Information and Systems*, E102-D(9):1696–1706, 2019.
- 300 Yamaoka Masanao, Yoshimura Chihiro, Hayashi Masato, Okuyama Takuya, Aoki Hidetaka, and Mizuno  
 301 Hiroyuki. A 20k-spin ising chip to solve combinatorial optimization problems with cmos annealing. *IEEE*  
 302 *Journal of Solid-State Circuits*, 51(1), 2010.
- 303 Yuya Sugie, Yuki Yoshida, Normann Mertig, Takashi Takemoto, Hiroshi Teramoto, Atsuyoshi Nakamura,  
 304 Ichigaku Takigawa, Shin ichi Minato, Masanao Yamaoka, and Tamiki Komatsuzaki. Minor-embedding  
 305 heuristics for large-scale annealing processors with sparse hardware graphs of up to 102,400 nodes. *Soft*  
 306 *Computing*, 2021.
- 307 Daphne Koller and Nir Friedman, editors. *Probabilistic Graphical Models: Principles and Techniques*. The  
 308 MIT Press, Cambridge, Massachusetts, 2009.
- 309 David Heckerman, Dan Geiger, and David M. Chickering. Learning bayesian networks: Search methods and  
 310 experimental results. In *Preliminary Papers of the 5th International Workshop on Artificial*, pages 112–128,  
 311 1995b.

## 312 Checklist

- 313 1. For all authors...
- 314 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-  
 315 tions and scope? [Yes]
- 316 (b) Did you describe the limitations of your work? [Yes] See section 7.
- 317 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See section 7.
- 318 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 319 2. If you are including theoretical results...
- 320 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See section 4.
- 321 (b) Did you include complete proofs of all theoretical results? [Yes] See section 5.
- 322 3. If you ran experiments...
- 323 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
 324 results (either in the supplemental material or as a URL)? [Yes] in the supplemental material
- 325 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?  
 326 [Yes] See section 6.
- 327 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
 328 multiple times)? [Yes] See fig. 2.
- 329 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,  
 330 internal cluster, or cloud provider)? [Yes] See section 6.3.
- 331 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 332 (a) If your work uses existing assets, did you cite the creators? [Yes] in the supplemental material
- 333 (b) Did you mention the license of the assets? [N/A]
- 334 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 335 (d) Did you discuss whether and how consent was obtained from people whose data you’re us-  
 336 ing/curating? [N/A]
- 337 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-  
 338 tion or offensive content? [N/A]
- 339 5. If you used crowdsourcing or conducted research with human subjects...
- 340 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?  
 341 [N/A]
- 342 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)  
 343 approvals, if applicable? [N/A]
- 344 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on  
 345 participant compensation? [N/A]