

# DYNAMIC INFERENCE IN PROBABILISTIC GRAPHICAL MODELS

WEIMING FENG, KUN HE, XIAOMING SUN, AND YITONG YIN

**ABSTRACT.** Probabilistic graphical models, such as Markov random fields (MRFs), are useful for describing high-dimensional distributions in terms of local dependence structures. The probabilistic inference is a fundamental problem related to graphical models, and sampling is a main approach for the problem. In this paper, we study probabilistic inference problems when the graphical model itself is changing dynamically with time. Such dynamic inference problems arise naturally in today’s application, e.g. multivariate time-series data analysis and practical learning procedures.

We give a dynamic algorithm for sampling-based probabilistic inferences in MRFs, where each dynamic update can change the underlying graph and all parameters of the MRF simultaneously, as long as the total amount of changes is bounded. More precisely, suppose that the MRF has  $n$  variables and polylogarithmic-bounded maximum degree, and  $N(n)$  independent samples are sufficient for the inference for a polynomial function  $N(\cdot)$ . Our algorithm dynamically maintains an answer to the inference problem using  $\tilde{O}(nN(n))$  space cost, and  $\tilde{O}(N(n)+n)$  incremental time cost upon each update to the MRF, as long as the Dobrushin-Shlosman condition is satisfied by the MRFs. This well-known condition has long been used for guaranteeing the efficiency of Markov chain Monte Carlo (MCMC) sampling in the traditional static setting. Compared to the static case, which requires  $\Omega(nN(n))$  time cost for redrawing all  $N(n)$  samples whenever the MRF changes, our dynamic algorithm gives a  $\tilde{\Omega}(\min\{n, N(n)\})$ -factor speedup. Our approach relies on a novel dynamic sampling technique, which transforms local Markov chains (a.k.a. single-site dynamics) to dynamic sampling algorithms, and an “algorithmic Lipschitz” condition that we establish for sampling from graphical models, namely, when the MRF changes by a small difference, samples can be modified to reflect the new distribution, with cost proportional to the difference on MRF.

arXiv:1904.11807v2 [cs.DS] 27 Nov 2020

---

(Weiming Feng, Yitong Yin) STATE KEY LABORATORY FOR NOVEL SOFTWARE TECHNOLOGY, NANJING UNIVERSITY. E-mail: [fengwm@smail.nju.edu.cn](mailto:fengwm@smail.nju.edu.cn), [yinyt@nju.edu.cn](mailto:yinyt@nju.edu.cn)

(Kun He) SHENZHEN UNIVERSITY; SHENZHEN INSTITUTE OF COMPUTING SCIENCES. E-mail: [hekun.threebody@foxmail.com](mailto:hekun.threebody@foxmail.com)

(Xiaoming Sun) CAS KEY LAB OF NETWORK DATA SCIENCE AND TECHNOLOGY, INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES. E-mail: [sunxiaoming@ict.ac.cn](mailto:sunxiaoming@ict.ac.cn)

Weiming Feng and Yitong Yin are supported by the National Key R&D Program of China 2018YFB1003202 and the National Natural Science Foundation of China under Grant Nos. 61722207 and 61672275. Kun He and Xiaoming Sun are supported by the National Natural Science Foundation of China Grants No. 61832003, 61433014 and K.C.Wong Education Foundation.

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our results	1
1.2	Related work	2
1.3	Organization of the paper.	2
<b>2</b>	<b>Dynamic inference problem</b>	<b>3</b>
2.1	Markov random fields.	3
2.2	Probabilistic inference and sampling	3
2.3	Dynamic inference problem	4
<b>3</b>	<b>Main results</b>	<b>5</b>
<b>4</b>	<b>Preliminaries</b>	<b>7</b>
<b>5</b>	<b>Outlines of algorithm</b>	<b>8</b>
<b>6</b>	<b>Dynamic Gibbs sampling</b>	<b>9</b>
6.1	Coupling for dynamic instances	10
6.2	Data structure for Gibbs sampling	17
6.3	Single-sample dynamic Gibbs sampling algorithm	18
6.4	Multi-sample dynamic Gibbs sampling algorithm	21
<b>7</b>	<b>Proofs for dynamic Gibbs sampling</b>	<b>24</b>
7.1	Analysis of the couplings	24
7.2	Implementation of the algorithms	29
7.3	Dynamic Gibbs sampling for specific models	31
<b>8</b>	<b>Proofs for dynamic inference</b>	<b>37</b>
8.1	Proof of the main theorem	37
8.2	Dynamic inference on specific models	37
<b>9</b>	<b>Conclusion</b>	<b>38</b>
	<b>References</b>	<b>38</b>

## 1. INTRODUCTION

The probabilistic graphical models provide a rich language for describing high-dimensional distributions in terms of the dependence structures between random variables. The *Markov random field* (MRF) is a basic graphical model that encodes pairwise interactions of complex systems. Given a graph  $G = (V, E)$ , each vertex  $v \in V$  is associated with a function  $\phi_v : Q \rightarrow \mathbb{R}$ , called the *vertex potential*, on a finite domain  $Q = [q]$  of  $q$  *spin states*, and each edge  $e \in E$  is associated with a symmetric function  $\phi_e : Q^2 \rightarrow \mathbb{R}$ , called the *edge potential*, which describes a pairwise interaction. Together, these induce a probability distribution  $\mu$  over all configurations  $\sigma \in Q^V$ :

$$\mu(\sigma) \propto \exp(H(\sigma)) = \exp\left(\sum_{v \in V} \phi_v(\sigma_v) + \sum_{e=\{u,v\} \in E} \phi_e(\sigma_u, \sigma_v)\right).$$

This distribution  $\mu$  is known as the Gibbs distribution and  $H(\sigma)$  is the *Hamiltonian*. It arises naturally from various physical models, statistics or learning problems, and combinatorial problems in computer science [MM09, KFB09].

The *probabilistic inference* is one of the most fundamental computational problems in graphical model. Some basic inference problems ask to calculate the marginal distribution, conditional distribution, or maximum-a-posteriori probabilities of one or several random variables [WJ08]. Sampling is perhaps the most widely used approach for probabilistic inference. Given a graphical model, independent samples are drawn from the Gibbs distribution and certain statistics are computed using the samples to give estimates for the inferred quantity. For most typical inference problems, such statistics are easy to compute once the samples are given, for instance, for estimating the marginal distribution on a variable subset  $S$ , the statistics is the frequency of each configuration in  $Q^S$  among the samples, thus the cost for inference is dominated by the cost for generating random samples [JVV86, ŠVV09].

The classic probabilistic inference assumes a static setting, where the input graphical model is fixed. In today's application, dynamically changing graphical models naturally arise in many scenarios. In various practical algorithms for learning graphical models, e.g. the contrastive divergence algorithm for learning the restricted Boltzmann machine [Hin12] and the iterative proportional fitting algorithm for maximum likelihood estimation of graphical models [WJ08], the optimal model  $\mathcal{I}^*$  is obtained by updating the parameters of the graphical model iteratively (usually by gradient descent), which generates a sequence of graphical models  $\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_M$ , with the goal that  $\mathcal{I}_M$  is a good approximation of  $\mathcal{I}^*$ . Also in the study of the multivariate time-series data, the dynamic Gaussian graphical models [CW07], multiregression dynamic model [QS93], dynamic graphical model [FVY19], and dynamic chain graph models [AQ<sup>+</sup>17], are all dynamically changing graphical models and have been used in a variety of applications. Meanwhile, with the advent of Big Data, scalable machine learning systems need to deal with continuously evolving graphical models (see e.g. [RKD<sup>+</sup>19] and [SWA09]).

The theoretical studies of probabilistic inference in dynamically changing graphical models are lacking. In the aforementioned scenarios in practice, it is common that a sequence of graphical models is presented with time, where any two consecutive graphical models can differ from each other in all potentials but by a small total amount. Recomputing the inference problem from scratch at every time when the graphical model is changed, can give the correct solution, but is very wasteful. A fundamental question is whether probabilistic inference can be solved dynamically and efficiently.

In this paper, we study the problem of probabilistic inference in an MRF when the MRF itself is changing dynamically with time. At each time, the whole graphical model, including all vertices and edges as well as their potentials, are subject to changes. Such *non-local* updates are very general and cover all applications mentioned above. The problem of *dynamic inference* then asks to maintain a correct answer to the inference in a dynamically changing MRF with low incremental cost proportional to the amount of changes made to the graphical model at each time.

**1.1. Our results.** We give a dynamic algorithm for sampling-based probabilistic inferences. Given an MRF instance with  $n$  vertices, suppose that  $N(n)$  independent samples are sufficient to give an approximate solution to the inference problem, where  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  is a polynomial function. We give dynamic algorithms for general inference problems on dynamically changing MRF.

Suppose that the current MRF has  $n$  vertices and polylogarithmic-bounded maximum degree, and each update to the MRF may change the underlying graph and/or all vertex/edge potentials, as long as the total amount of changes is bounded. Our algorithm maintains an approximate solution to the inference with  $\tilde{O}(nN(n))$  space cost, and with  $\tilde{O}(N(n) + n)$  incremental time cost upon each update, assuming that the MRFs satisfy the Dobrushin-Shlosman condition [DS85a, DS85b, DS87]. The condition has been widely used to imply the efficiency of Markov chain Monte Carlo (MCMC) sampling (e.g. see [Hay06, DGJ08]). Compared to the static algorithm, which requires  $\Omega(nN(n))$  time for re-drawing all  $N(n)$  samples each time, our dynamic algorithm significantly improves the time cost with an  $\tilde{\Omega}(\min\{n, N(n)\})$ -factor speedup.

On specific models, the Dobrushin-Shlosman condition has been established in the literature, which directly gives us following efficient dynamic inference algorithms, with  $\tilde{O}(nN(n))$  space cost and  $\tilde{O}(N(n) + n)$  time cost per update, on graphs with  $n$  vertices and maximum degree  $\Delta = O(1)$ :

- for Ising model with temperature  $\beta$  satisfying  $e^{-2|\beta|} > 1 - \frac{2}{\Delta+1}$ , which is close to the uniqueness threshold  $e^{-2|\beta_c|} = 1 - \frac{2}{\Delta}$ , beyond which the static versions of sampling or marginal inference problem for anti-ferromagnetic Ising model is intractable [GŠV16, GŠV15];
- for hardcore model with fugacity  $\lambda < \frac{2}{\Delta-2}$ , which matches the best bound known for sampling algorithm with near-linear running time on general graphs with bounded maximum degree [Vig99, LV99, DG00];
- for proper  $q$ -coloring with  $q > 2\Delta$ , which matches the best bound known for sampling algorithm with near-linear running time on general graphs with bounded maximum degree [Jer95].

Our dynamic inference algorithm is based on a dynamic sampling algorithm, which efficiently maintains  $N(n)$  independent samples for the current MRF while the MRF is subject to changes. More specifically, we give a dynamic version of the *Gibbs sampling* algorithm, a local Markov chain for sampling from the Gibbs distribution that has been studied extensively. Our techniques are based on: (1) couplings for dynamic instances of graphical models; and (2) dynamic data structures for representing single-site Markov chains so that the couplings can be realized algorithmically in sub-linear time. Both these techniques are of independent interest, and can be naturally extended to more general settings with multi-body interactions.

Our results show that on dynamically changing graphical models, sampling-based probabilistic inferences can be solved significantly faster than rerunning the static algorithm at each time. This has practical significance in speeding up the iterative procedures for learning graphical models.

**1.2. Related work.** The problem of dynamic sampling from graphical models was introduced very recently in [FVY19]. There, a dynamic sampling algorithm was given for graphical models with soft constraints, and can only deal with local updates that change a single vertex or edge at each time. The regimes for such dynamic sampling algorithm to be efficient are much more restrictive than the conditions for the rapid mixing of Markov chains. Our algorithm greatly improves the regimes for efficient dynamic sampling for the Ising and hardcore models in [FVY19], and for the first time, can handle non-local updates that change all vertex/edge potentials simultaneously. Besides, the dynamic/online sampling from log-concave distributions was also studied in [NR17, LMV19].

Another related topic is the dynamic graph problems, which ask to maintain a solution (e.g. spanners [FG19, NSWN17, WN17] or shortest paths [BC16, HKN16, HKN14]) while the input graph is dynamically changing. More recently, important progress has been made on dynamically maintaining structures that are related to graph random walks, such as spectral sparsifier [DGGP19, ADK<sup>+</sup>16] or effective resistances [DGGP18, GHP18]. Instead of one particular solution, dynamic inference problems ask to maintain an estimate of a statistics, such statistics comes from an exponential-sized probability space described by a dynamically changing graphical model.

**1.3. Organization of the paper.** In Section 2, we formally introduce the dynamic inference problem. In Section 3, we formally state the main results. Preliminaries are given in Section 4. In Section 5, we outline our dynamic inference algorithm. In Section 6, we present the algorithms for dynamic Gibbs sampling. The analyses of these dynamic sampling algorithms are given in Section 7. The proof of the main theorem on dynamic inference is given in Section 8. The conclusion is given in Section 9.

## 2. DYNAMIC INFERENCE PROBLEM

**2.1. Markov random fields.** An instance of *Markov random field (MRF)* is specified by a tuple  $\mathcal{I} = (V, E, Q, \Phi)$ , where  $G = (V, E)$  is an undirected simple graph;  $Q$  is a domain of  $q = |Q|$  *spin states*, for some finite  $q > 1$ ; and  $\Phi = (\phi_a)_{a \in V \cup E}$  associates each  $v \in V$  a *vertex potential*  $\phi_v : Q \rightarrow \mathbb{R}$  and each  $e \in E$  an *edge potential*  $\phi_e : Q^2 \rightarrow \mathbb{R}$ , where  $\phi_e$  is symmetric.

A *configuration*  $\sigma \in Q^V$  maps each vertex  $v \in V$  to a spin state in  $Q$ , so that each vertex can be interpreted as a variable. And the *Hamiltonian* of a configuration  $\sigma \in Q^V$  is defined as:

$$H(\sigma) \triangleq \sum_{v \in V} \phi_v(\sigma_v) + \sum_{e = \{u, v\} \in E} \phi_e(\sigma_u, \sigma_v).$$

This defines the *Gibbs distribution*  $\mu_{\mathcal{I}}$ , which is a probability distribution over  $Q^V$  such that

$$\forall \sigma \in Q^V, \quad \mu_{\mathcal{I}}(\sigma) = \frac{1}{Z} \exp(H(\sigma)),$$

where the normalizing factor  $Z \triangleq \sum_{\sigma \in Q^V} \exp(H(\sigma))$  is called the *partition function*.

The Gibbs measure  $\mu(\sigma)$  can be 0 as the functions  $\phi_v, \phi_e$  can take the value  $-\infty$ . A configuration  $\sigma$  is called *feasible* if  $\mu(\sigma) > 0$ . To trivialize the problem of constructing a feasible configuration, we further assume the following natural condition for the MRF instances considered in this paper:<sup>1</sup>

$$(1) \quad \forall v \in V, \forall \sigma \in Q^{\Gamma_G(v)} : \sum_{c \in Q} \exp\left(\phi_v(c) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, c)\right) > 0.$$

where  $\Gamma_G(v) \triangleq \{u \in V \mid \{u, v\} \in E\}$  denotes the neighborhood of  $v$  in graph  $G = (V, E)$ .

Some well studied typical MRFs include:

- *Ising model:* The domain of each spin is  $Q = \{-1, +1\}$ . Each edge  $e \in E$  is associated with a *temperature*  $\beta_e \in \mathbb{R}$ ; and each vertex  $v \in V$  is associated with a *local field*  $h_v \in \mathbb{R}$ . For each configuration  $\sigma \in \{-1, +1\}^V$ ,  $\mu_{\mathcal{I}}(\sigma) \propto \exp\left(\sum_{\{u, v\} \in E} \beta_e \sigma_u \sigma_v + \sum_{v \in V} h_v \sigma_v\right)$ .
- *Hardcore model:* The domain is  $Q = \{0, 1\}$ . Each configuration  $\sigma \in Q^V$  indicates an independent set in  $G = (V, E)$ , and  $\mu_{\mathcal{I}}(\sigma) \propto \lambda^{|\sigma|}$ , where  $\lambda > 0$  is the *fugacity* parameter.
- *proper  $q$ -coloring:* uniform distribution over all proper  $q$ -colorings of  $G = (V, E)$ .

**2.2. Probabilistic inference and sampling.** In graphical models, the task of probabilistic inference is to derive the probabilities regarding one or more random variables of the model. Abstractly, this is described by a function  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$  that maps each graphical model  $\mathcal{I} \in \mathfrak{M}$  to a target  $K$ -dimensional probability vector, where  $\mathfrak{M}$  is the class of graphical models containing the random variables we are interested in and the  $K$ -dimensional vector describes the probabilities we want to derive. Given  $\theta(\cdot)$  and an MRF instance  $\mathcal{I} \in \mathfrak{M}$ , the inference problem asks to estimate the probability vector  $\theta(\mathcal{I})$ .

Here are some fundamental inference problems [WJ08] for MRF instances. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance and  $A, B \subseteq V$  two disjoint sets where  $A \uplus B \subseteq V$ .

- *Marginal inference:* estimate the marginal distribution  $\mu_{A, \mathcal{I}}(\cdot)$  of the variables in  $A$ , where

$$\forall \sigma_A \in Q^A, \quad \mu_{A, \mathcal{I}}(\sigma_A) \triangleq \sum_{\tau \in Q^{V \setminus A}} \mu_{\mathcal{I}}(\sigma_A, \tau).$$

- *Posterior inference:* given any  $\tau_B \in Q^B$ , estimate the posterior distribution  $\mu_{A, \mathcal{I}}(\cdot \mid \tau_B)$  for the variables in  $A$ , where

$$\forall \sigma_A \in Q^A, \quad \mu_{A, \mathcal{I}}(\sigma_A \mid \tau_B) \triangleq \frac{\mu_{A \cup B, \mathcal{I}}(\sigma_A, \tau_B)}{\mu_{B, \mathcal{I}}(\tau_B)}.$$

---

<sup>1</sup>This condition guarantees that the marginal probabilities are always well-defined, and the problem of constructing a feasible configuration  $\sigma$ , where  $\mu_{\mathcal{I}}(\sigma) > 0$ , is trivial. The condition holds for all MRFs with soft constraints, or with hard constraints where there is a permissive spin, e.g. the hardcore model. For MRFs with truly repulsive hard constraints such as proper  $q$ -coloring, the condition may translate to the condition  $q \geq \Delta + 1$  where  $\Delta$  is the maximum degree of graph  $G$ , which is necessary for the irreducibility of local Markov chains for  $q$ -colorings.

- *Maximum-a-posteriori (MAP) inference*: find the maximum-a-posteriori (MAP) probabilities  $P_{A,I}^*(\cdot)$  for the configurations over  $A$ , where

$$\forall \sigma_A \in Q^A, \quad P_{A,I}^*(\sigma_A) \triangleq \max_{\tau_B \in Q^B} \mu_{A \cup B, I}(\sigma_A, \tau_B).$$

All these fundamental inference problems can be described abstractly by a function  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$ . For instances, for marginal inference,  $\mathfrak{M}$  contains all MRF instances where  $A$  is a subset of the vertices,  $K = |Q|^{|A|}$ , and  $\theta(I) = (\mu_{A,I}(\sigma_A))_{\sigma_A \in Q^A}$ ; and for posterior or MAP inference,  $\mathfrak{M}$  contains all MRF instances where  $A \uplus B$  is a subset of the vertices,  $K = |Q|^{|A|}$  and  $\theta(I) = (\mu_{A,I}(\sigma_A \mid \tau_B))_{\sigma_A \in Q^A}$  (for posterior inference) or  $\theta(I) = (P_{A,I}^*(\sigma_A))_{\sigma_A \in Q^A}$  (for MAP inference).

One canonical approach for probabilistic inference is by sampling: sufficiently many independent samples are drawn (approximately) from the Gibbs distribution of the MRF instance and an estimate of the target probabilities is calculated from these samples. Given a probabilistic inference problem  $\theta(\cdot)$ , we use  $\mathcal{E}_\theta(\cdot)$  to denote an estimating function that approximates  $\theta(I)$  using independent samples drawn approximately from  $\mu_I$ . For the aforementioned problems of marginal, posterior and MAP inferences, such estimating function  $\mathcal{E}_\theta(\cdot)$  simply counts the frequency of the samples that satisfy certain properties.

The sampling cost of an estimator is captured in two aspects: the number of samples it uses and the accuracy of each individual sample it requires.

**Definition 2.1 (( $N, \epsilon$ )-estimator for  $\theta$ ).** Let  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$  be a probabilistic inference problem and  $\mathcal{E}_\theta(\cdot)$  an estimating function for  $\theta(\cdot)$  that for each instance  $I = (V, E, Q, \Phi) \in \mathfrak{M}$ , maps samples in  $Q^V$  to an estimate of  $\theta(I)$ . Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$ . For any instance  $I = (V, E, Q, \Phi) \in \mathfrak{M}$  where  $n = |V|$ , the random variable  $\mathcal{E}_\theta(X^{(1)}, \dots, X^{(N(n))})$  is said to be an  $(N, \epsilon)$ -estimator for  $\theta(I)$  if  $X^{(1)}, \dots, X^{(N(n))} \in Q^V$  are  $N(n)$  independent samples drawn approximately from  $\mu_I$  such that  $d_{\text{TV}}(X^{(j)}, \mu_I) \leq \epsilon(n)$  for all  $1 \leq j \leq N(n)$ .

In Definition 2.1, an estimator is viewed as a black-box algorithm specified by two functions  $N$  and  $\epsilon$ . Usually, the estimator is more accurate if more independent samples are drawn and each sample provides a higher level of accuracy. Thus, one can choose some large  $N$  and small  $\epsilon$  to achieve a desired quality of estimate.

**2.3. Dynamic inference problem.** We consider the inference problem where the input graphical model is changed dynamically: at each step, the current MRF instance  $I = (V, E, Q, \Phi)$  is updated to a new instance  $I' = (V', E', Q, \Phi')$ . We consider general update operations for MRFs that can change both the **underlying graph** and **all edge/vertex potentials** simultaneously, where the update request is made by a *non-adaptive adversary* independently of the randomness used by the inference algorithm. Such updates are general enough and cover many applications, e.g. analyses of time series network data [CW07, QS93, FVY19, AQ<sup>+</sup>17], and learning algorithms for graphical models [Hin12, WJ08].

The difference between the original and the updated instances is measured as follows.

**Definition 2.2 (difference between MRF instances).** The difference between two MRF instances  $I = (V, E, Q, \Phi)$  and  $I' = (V', E', Q, \Phi')$ , where  $\Phi = (\phi_a)_{a \in V \cup E}$  and  $\Phi' = (\phi'_a)_{a \in V' \cup E'}$ , is defined as

$$(2) \quad d(I, I') \triangleq \sum_{v \in V \cap V'} \|\phi_v - \phi'_v\|_1 + \sum_{e \in E \cap E'} \|\phi_e - \phi'_e\|_1 + |V \oplus V'| + |E \oplus E'|,$$

where  $A \oplus B = (A \setminus B) \cup (B \setminus A)$  stands for the symmetric difference between two sets  $A$  and  $B$ ,  $\|\phi_v - \phi'_v\|_1 \triangleq \sum_{c \in Q} |\phi_v(c) - \phi'_v(c)|$ , and  $\|\phi_e - \phi'_e\|_1 \triangleq \sum_{c, c' \in Q} |\phi_e(c, c') - \phi'_e(c, c')|$ .

Given a probability vector specified by the function  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$ , the *dynamic inference problem* asks to maintain an estimator  $\hat{\theta}(I)$  of  $\theta(I)$  for the current MRF instance  $I = (V, E, Q, \Phi) \in \mathfrak{M}$ , with a data structure, such that when  $I$  is updated to  $I' = (V', E', Q, \Phi') \in \mathfrak{M}$ , the algorithm updates  $\hat{\theta}(I)$  to an estimator  $\hat{\theta}(I')$  of the new vector  $\theta(I')$ , or equivalently, outputs the difference between the estimators  $\hat{\theta}(I)$  and  $\hat{\theta}(I')$ .

It is desirable to have the dynamic inference algorithm which maintains an  $(N, \epsilon)$ -estimator for  $\theta(I)$  for the current instance  $I$ . However, the dynamic algorithm cannot be efficient if  $N(n)$  and  $\epsilon(n)$



change drastically with  $n$  (so that significantly more samples or substantially more accurate samples may be needed when a new vertex is added), or if recalculating the estimating function  $\mathcal{E}_\theta(\cdot)$  itself is expensive. We introduce a notion of *dynamical efficiency* for the estimators that are suitable for dynamic inference.

**Definition 2.3 (dynamical efficiency).** Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$ . Let  $\mathcal{E}(\cdot)$  be an estimating function for some  $K$ -dimensional probability vector of MRF instances. A tuple  $(N, \epsilon, \mathcal{E})$  is said to be *dynamically efficient* if it satisfies:

- **(bounded difference)** there exist constants  $C_1, C_2 > 0$  such that for any  $n \in \mathbb{N}^+$ ,

$$|N(n+1) - N(n)| \leq \frac{C_1 \cdot N(n)}{n} \quad \text{and} \quad |\epsilon(n+1) - \epsilon(n)| \leq \frac{C_2 \cdot \epsilon(n)}{n};$$

- **(small incremental cost)** there is a deterministic algorithm that maintains  $\mathcal{E}(X^{(1)}, \dots, X^{(m)})$  using  $(mn + K) \cdot \text{polylog}(mn)$  bits where  $X^{(1)}, \dots, X^{(m)} \in Q^V$  and  $n = |V|$ , such that when  $X^{(1)}, \dots, X^{(m)} \in Q^V$  are updated to  $Y^{(1)}, \dots, Y^{(m')} \in Q^{V'}$ , where  $n' = |V'|$ , the algorithm updates  $\mathcal{E}(X^{(1)}, \dots, X^{(m)})$  to  $\mathcal{E}(Y^{(1)}, \dots, Y^{(m')})$  within time cost  $\mathcal{D} \cdot \text{polylog}(mm'nn') + O(m + m')$ , where  $\mathcal{D}$  is the size of the difference between two sample sequences defined as:

$$(3) \quad \mathcal{D} \triangleq \sum_{i \leq \max\{m, m'\}} \sum_{v \in V \cup V'} 1 \left[ X^{(i)}(v) \neq Y^{(i)}(v) \right],$$

where an unassigned  $X^{(i)}(v)$  or  $Y^{(i)}(v)$  is not equal to any assigned spin.

The dynamic efficiency basically asks  $N(\cdot)$ ,  $\epsilon(\cdot)$ , and  $\mathcal{E}(\cdot)$  to have some sort of ‘‘Lipschitz’’ properties. To satisfy the bounded difference condition,  $N(n)$  and  $1/\epsilon(n)$  are necessarily polynomially bounded, and they can be any constant, polylogarithmic, or polynomial functions, or multiplications of such functions. The condition with small incremental cost also holds very commonly. In particular, it is satisfied by the estimating functions for all the aforementioned problems for the marginal, posterior and MAP inferences as long as the sets of variables have sizes  $|A|, |B| = O(\log n)$ . We remark that the  $O(\log n)$  upper bound is somehow necessary for the efficiency of inference, because otherwise the dimension of  $\theta(\mathcal{I})$  itself (which is at least  $q^{|A|}$ ) becomes super-polynomial in  $n$ .

### 3. MAIN RESULTS

Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance, where  $G = (V, E)$ . Let  $\Gamma_G(v)$  denote the neighborhood of  $v$  in  $G$ . For any vertex  $v \in V$  and any configuration  $\sigma \in Q^{\Gamma_G(v)}$ , we use  $\mu_{v, \mathcal{I}}^\sigma(\cdot) = \mu_{v, \mathcal{I}}(\cdot \mid \sigma)$  to denote the marginal distribution on  $v$  conditional on  $\sigma$ :

$$\forall c \in Q : \quad \mu_{v, \mathcal{I}}^\sigma(c) = \mu_{v, \mathcal{I}}(c \mid \sigma) \triangleq \frac{\exp(\phi_v(c) + \sum_{u \in \Gamma_G(v)} \phi_{uv}(\sigma_u, c))}{\sum_{a \in Q} \exp(\phi_v(a) + \sum_{u \in \Gamma_G(v)} \phi_{uv}(\sigma_u, a))}.$$

Due to the assumption in (1), the marginal distribution is always well-defined. The following condition is the *Dobrushin-Shlosman condition* [DS85a, DS85b, DS87, Hay06, DGJ08].

**Condition 3.1 (Dobrushin-Shlosman condition).** Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance with Gibbs distribution  $\mu = \mu_{\mathcal{I}}$ . Let  $A_{\mathcal{I}} \in \mathbb{R}_{\geq 0}^{V \times V}$  be the *influence matrix* which is defined as

$$A_{\mathcal{I}}(u, v) \triangleq \begin{cases} \max_{(\sigma, \tau) \in B_{u,v}} d_{\text{TV}}(\mu_v^\sigma, \mu_v^\tau), & \{u, v\} \in E, \\ 0 & \{u, v\} \notin E, \end{cases}$$

where the maximum is taken over the set  $B_{u,v}$  of all  $(\sigma, \tau) \in Q^{\Gamma_G(v)} \times Q^{\Gamma_G(v)}$  that differ only at  $u$ , and  $d_{\text{TV}}(\mu_v^\sigma, \mu_v^\tau) \triangleq \frac{1}{2} \sum_{c \in Q} |\mu_v^\sigma(c) - \mu_v^\tau(c)|$  is the total variation distance between  $\mu_v^\sigma$  and  $\mu_v^\tau$ . An MRF instance  $\mathcal{I}$  is said to satisfy the *Dobrushin-Shlosman condition* if there is a constant  $\delta > 0$  such that

$$\max_{u \in V} \sum_{v \in V} A_{\mathcal{I}}(u, v) \leq 1 - \delta.$$

Our main theorem assumes the following setup: Let  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$  be a probabilistic inference problem that maps each MRF instance in  $\mathfrak{M}$  to a  $K$ -dimensional probability vector, and let  $\mathcal{E}_\theta$  be its estimating function. Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$ . We use  $\mathcal{I} = (V, E, Q, \Phi) \in \mathfrak{M}$ , where  $n = |V|$ , to denote the current instance and  $\mathcal{I}' = (V', E', Q, \Phi') \in \mathfrak{M}$ , where  $n' = |V'|$ , to denote the updated instance.

**Theorem 3.2 (dynamic inference algorithm).** *Assume that  $(N, \epsilon, \mathcal{E}_\theta)$  is dynamically efficient, both  $\mathcal{I}$  and  $\mathcal{I}'$  satisfy the Dobrushin-Shlosman condition, and  $d(\mathcal{I}, \mathcal{I}') \leq L = o(n)$ .*

*There is an algorithm that maintains an  $(N, \epsilon)$ -estimator  $\hat{\theta}(\mathcal{I})$  of the probability vector  $\theta(\mathcal{I})$  for the current MRF instance  $\mathcal{I}$ , using  $\tilde{O}(nN(n) + K)$  bits, such that when  $\mathcal{I}$  is updated to  $\mathcal{I}'$ , the algorithm updates  $\hat{\theta}(\mathcal{I})$  to an  $(N, \epsilon)$ -estimator  $\hat{\theta}(\mathcal{I}')$  of  $\theta(\mathcal{I}')$  for the new instance  $\mathcal{I}'$ , within expected time cost*

$$\tilde{O}\left(\Delta^2 LN(n) + \Delta n\right),$$

where  $\tilde{O}(\cdot)$  hides a  $\text{polylog}(n)$  factor,  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ , where  $\Delta_G$  and  $\Delta_{G'}$  denote the maximum degree of  $G = (V, E)$  and  $G' = (V', E')$  respectively.

Note that the extra  $O(\Delta n)$  cost is necessary for editing the current MRF instance  $\mathcal{I}$  to  $\mathcal{I}'$ .

Typically, the difference between two MRF instances  $\mathcal{I}, \mathcal{I}'$  is small<sup>2</sup>, and the underlying graphs are sparse [DSOR16], that is,  $L, \Delta \leq \text{polylog}(n)$ . In such cases, our algorithm updates the estimator within time cost  $\tilde{O}(N(n)+n)$ , which significantly outperforms static sampling-based inference algorithms that require time cost  $\Omega(n'N(n')) = \Omega(nN(n))$  for redrawing all  $N(n')$  independent samples.

**Dynamic sampling.** The core of our dynamic inference algorithm is a dynamic sampling algorithm: Assuming the Dobrushin-Shlosman condition, the algorithm can maintain a sequence of  $N(n)$  independent samples  $X^{(1)}, \dots, X^{(N(n))} \in Q^V$  that are  $\epsilon(n)$ -close to  $\mu_{\mathcal{I}}$  in total variation distance, and when  $\mathcal{I}$  is updated to  $\mathcal{I}'$  with difference  $d(\mathcal{I}, \mathcal{I}') \leq L = o(n)$ , the algorithm can update the maintained samples to  $N(n')$  independent samples  $Y^{(1)}, \dots, Y^{(N(n'))} \in Q^{V'}$  that are  $\epsilon(n')$ -close to  $\mu_{\mathcal{I}'}$  in total variation distance, using a time cost  $\tilde{O}(\Delta^2 LN(n) + \Delta n)$  in expectation. This shows an “algorithmic Lipschitz” condition holds for sampling from Gibbs distributions: when the MRF changes insignificantly, a population of samples can be modified to reflect the new distribution, with cost proportional to the difference on MRF. We show that such property is guaranteed by the Dobrushin-Shlosman condition. This dynamic sampling algorithm is formally described in Theorem 6.1 and is of independent interest [FVY19].

**Applications on specific models.** On specific models, we have the following results, where  $\delta > 0$  is an arbitrary constant.

model	regime	space cost	time cost for each update
Ising	$e^{-2 \beta } \geq 1 - \frac{2-\delta}{\Delta+1}$	$\tilde{O}(nN(n) + K)$	$\tilde{O}(\Delta^2 LN(n) + \Delta n)$
hardcore	$\lambda \leq \frac{2-\delta}{\Delta-2}$	$\tilde{O}(nN(n) + K)$	$\tilde{O}(\Delta^3 LN(n) + \Delta n)$
$q$ -coloring	$q \geq (2 + \delta)\Delta$	$\tilde{O}(nN(n) + K)$	$\tilde{O}(\Delta^2 LN(n) + \Delta n)$

TABLE 1. Dynamic inference for specific models.

The results for Ising model and  $q$ -coloring are corollaries of Theorem 3.2. The regime for hardcore model is better than the Dobrushin-Shlosman condition (which is  $\lambda \leq \frac{1-\delta}{\Delta-1}$ ), because we use the coupling introduced by Vigoda [Vig99] to analyze the algorithm.

<sup>2</sup>In multivariate time-series data analysis, the MRF instances of two sequential times are similar. In the iterative algorithms for learning graphical models, the difference between two sequential MRF instances generated by gradient descent are bounded to prevent oscillations. Specifically, the difference is very small when the iterative algorithm approaches to the convergence state [Hin12, WJ08].



#### 4. PRELIMINARIES

**Total variation distance and coupling.** Let  $\mu$  and  $\nu$  be two distributions over  $\Omega$ . The *total variation distance* between  $\mu$  and  $\nu$  is defined as

$$d_{\text{TV}}(\mu, \nu) \triangleq \frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)|.$$

A *coupling* of  $\mu$  and  $\nu$  is a joint distribution  $(X, Y) \in \Omega \times \Omega$  such that marginal distribution of  $X$  is  $\mu$  and the marginal distribution of  $Y$  is  $\nu$ . The following coupling lemma is well-known.

**Proposition 4.1 (coupling lemma).** *For any coupling  $(X, Y)$  of  $\mu$  and  $\nu$ , it holds that*

$$\Pr[X \neq Y] \geq d_{\text{TV}}(\mu, \nu).$$

*Furthermore, there is an optimal coupling that achieves equality.*

**Local neighborhood.** Let  $G = (V, E)$  be a graph. For any vertex  $v \in V$ , let  $\Gamma_G(v) \triangleq \{u \in V \mid \{u, v\} \in E\}$  denote the neighborhood of  $v$ , and  $\Gamma_G^+(v) \triangleq \Gamma_G(v) \cup \{v\}$  the inclusive neighborhood of  $v$ . We simply write  $\Gamma_v = \Gamma(v) = \Gamma_G(v)$  and  $\Gamma_v^+ = \Gamma^+(v) = \Gamma_G^+(v)$  for short when  $G$  is clear in the context. We use  $\Delta = \Delta_G \triangleq \max_{v \in V} |\Gamma_v|$  to denote the maximum degree of graph  $G$ .

A notion of **local neighborhood for MRF** is frequently used. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance. For  $v \in V$ , we denote by  $\mathcal{I}_v \triangleq \mathcal{I}[\Gamma_v^+]$  the restriction of  $\mathcal{I}$  on the inclusive neighborhood  $\Gamma_v^+$  of  $v$ , i.e.  $\mathcal{I}_v = (\Gamma_v^+, E_v, Q, \Phi_v)$ , where  $E_v = \{\{u, v\} \in E\}$  and  $\Phi_v = (\phi_a)_{a \in \Gamma_v^+ \cup E_v}$ .

**Gibbs sampling.** The *Gibbs sampling* (a.k.a. *heat-bath*, *Glauber dynamics*), is a classic Markov chain for sampling from Gibbs distributions. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance and  $\mu = \mu_{\mathcal{I}}$  its Gibbs distribution. The chain of Gibbs sampling (Algorithm 1) is on the space  $\Omega \triangleq Q^V$ , and has the stationary distribution  $\mu_{\mathcal{I}}$  [LP17, Chapter 3].

---

#### Algorithm 1: Gibbs sampling

---

**Initialization:** an initial state  $X_0 \in \Omega$  (not necessarily feasible);

- 1 **for**  $t = 1, 2, \dots, T$  **do**
  - 2     pick  $v_t \in V$  uniformly at random;
  - 3     draw a random value  $c \in Q$  from the marginal distribution  $\mu_{v_t}(\cdot \mid X_{t-1}(\Gamma_{v_t}))$ ;
  - 4      $X_t(v_t) \leftarrow c$  and  $X_t(u) \leftarrow X_{t-1}(u)$  for all  $u \in V \setminus \{v_t\}$ ;
- 

**Marginal distributions.** Here  $\mu_v(\cdot \mid \sigma(\Gamma_v)) = \mu_{v, \mathcal{I}}(\cdot \mid \sigma(\Gamma_v))$  denotes the marginal distribution at  $v \in V$  conditioning on  $\sigma(\Gamma_v) \in Q^{\Gamma_v}$ , which is computed as:

$$(4) \quad \forall c \in Q : \quad \mu_v(c \mid \sigma(\Gamma_v)) = \frac{\phi_v(c) \prod_{u \in \Gamma_v} \phi_{uv}(\sigma_u, c)}{\sum_{c' \in Q} \phi_v(c') \prod_{u \in \Gamma_v} \phi_{uv}(\sigma_u, c')}.$$

Due to the assumption (1), this marginal distribution is always well defined, and its computation uses only the information of  $\mathcal{I}_v$ .

**Coupling for mixing time.** Consider a chain  $(X_t)_{t=0}^{\infty}$  on space  $\Omega$  with stationary distribution  $\mu_{\mathcal{I}}$  for MRF instance  $\mathcal{I}$ . The *mixing rate* is defined as: for  $\epsilon > 0$ ,

$$\tau_{\text{mix}}(\mathcal{I}, \epsilon) \triangleq \max_{X_0} \min \{t \mid d_{\text{TV}}(X_t, \mu_{\mathcal{I}}) \leq \epsilon\},$$

where  $d_{\text{TV}}(X_t, \mu_{\mathcal{I}})$  denotes the *total variation distance* between  $\mu_{\mathcal{I}}$  and the distribution of  $X_t$ .

A coupling of a Markov chain is a joint process  $(X_t, Y_t)_{t \geq 0}$  such that  $(X_t)_{t \geq 0}$  and  $(Y_t)_{t \geq 0}$  marginally follow the same transition rule as the original chain. Consider the following type of couplings.

**Definition 4.2 (one-step optimal coupling for Gibbs sampling).** A coupling  $(X_t, Y_t)_{t \geq 0}$  of Gibbs sampling on an MRF instance  $\mathcal{I} = (V, E, Q, \Phi)$  is a *one-step optimal coupling* if it is constructed as follows: For  $t = 1, 2, \dots$ ,

- (1) pick the same random  $v_t \in V$ , and let  $(X_t(u), Y_t(u)) \leftarrow (X_{t-1}(u), Y_{t-1}(u))$  for all  $u \neq v_t$ ;

- (2) sample  $(X_t(v_t), Y_t(v_t))$  from an optimal coupling  $D_{\text{opt}, \mathcal{I}_{v_t}}^{\sigma, \tau}(\cdot, \cdot)$  of the marginal distributions  $\mu_{v_t}(\cdot | \sigma)$  and  $\mu_{v_t}(\cdot | \tau)$  where  $\sigma = X_{t-1}(\Gamma_{v_t})$  and  $\tau = Y_{t-1}(\Gamma_{v_t})$ .

The coupling  $D_{\text{opt}, \mathcal{I}_{v_t}}^{\sigma, \tau}(\cdot, \cdot)$  is an *optimal coupling* of  $\mu_{v_t}(\cdot | \sigma)$  and  $\mu_{v_t}(\cdot | \tau)$  that attains the maximum  $\Pr[\mathbf{x} = \mathbf{y}]$  for all couplings  $(\mathbf{x}, \mathbf{y})$  of  $\mathbf{x} \sim \mu_{v_t}(\cdot | \sigma)$  and  $\mathbf{y} \sim \mu_{v_t}(\cdot | \tau)$ . The coupling  $D_{\text{opt}, \mathcal{I}_{v_t}}^{\sigma, \tau}(\cdot, \cdot)$  is determined by the local information  $\mathcal{I}_v$  and  $\sigma, \tau \in Q^{\text{deg}(v)}$ .

With such a coupling, we can establish the following relation between the Dobrushin-Shlosman condition and the rapid mixing of the Gibbs sampling [DS85a, DS85b, DS87, BD97, Hay06, DGJ08].

**Proposition 4.3** ([BD97, Hay06]). *Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance with  $n = |V|$ , and  $\Omega = Q^V$  the state space. Let  $H(\sigma, \tau) \triangleq |\{v \in V \mid \sigma_v \neq \tau_v\}|$  denote the Hamming distance between  $\sigma \in \Omega$  and  $\tau \in \Omega$ . If  $\mathcal{I}$  satisfies the Dobrushin-Shlosman condition (Condition 3.1) with constant  $\delta > 0$ , then the one-step optimal coupling  $(X_t, Y_t)_{t \geq 0}$  for Gibbs sampling (Definition 4.2) satisfies*

$$\forall \sigma, \tau \in \Omega : \quad \mathbb{E}[H(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \leq \left(1 - \frac{\delta}{n}\right) \cdot H(\sigma, \tau),$$

and hence the mixing rate of Gibbs sampling on  $\mathcal{I}$  is bounded as  $\tau_{\text{mix}}(\mathcal{I}, \epsilon) \leq \left\lceil \frac{n}{\delta} \log \frac{n}{\epsilon} \right\rceil$ .

## 5. OUTLINES OF ALGORITHM

Let  $\theta : \mathfrak{M} \rightarrow \mathbb{R}^K$  be a probabilistic inference problem that maps each MRF instance in  $\mathfrak{M}$  to a  $K$ -dimensional probability vector, and let  $\mathcal{E}_\theta$  be its estimating function. Let  $\mathcal{I} = (V, E, Q, \Phi) \in \mathfrak{M}$  be the current instance, where  $n = |V|$ . Our dynamic inference algorithm maintains a sequence of  $N(n)$  independent samples  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N(n))} \in Q^V$  which are  $\epsilon(n)$ -close to the Gibbs distribution  $\mu_{\mathcal{I}}$  in total variation distance and an  $(N, \epsilon)$ -estimator  $\hat{\theta}(\mathcal{I})$  of  $\theta(\mathcal{I})$  such that

$$\hat{\theta}(\mathcal{I}) = \mathcal{E}_\theta(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(N(n))}).$$

Upon an update request that modifies  $\mathcal{I}$  to a new instance  $\mathcal{I}' = (V', E', Q, \Phi') \in \mathfrak{M}$ , where  $n' = |V'|$ , our algorithm does the followings:

- *Update the sample sequence.* Update  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N(n))}$  to a new sequence of  $N(n')$  independent samples  $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(N(n'))} \in Q^{V'}$  that are  $\epsilon(n')$ -close to  $\mu_{\mathcal{I}'}$  in total variation distance, and output the difference between two sample sequences.
- *Update the estimator.* Given the difference between the two sample sequences, update  $\hat{\theta}(\mathcal{I})$  to  $\hat{\theta}(\mathcal{I}') = \mathcal{E}_\theta(\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(N(n'))})$  by accessing the oracle in Definition 2.3.

Obviously, the updated estimator  $\hat{\theta}(\mathcal{I}')$  is an  $(N, \epsilon)$ -estimator for  $\theta(\mathcal{I}')$ .

Our main technical contribution is to give an algorithm that dynamically maintains a sequence of  $N(n)$  independent samples for  $\mu_{\mathcal{I}}$ , while  $\mathcal{I}$  itself is dynamically changing. The dynamic sampling problem was recently introduced in [FVY19]. The dynamical sampling algorithm given there only handles update of a single vertex or edge and works only for graphical models with soft constraints.

In contrast, our dynamic sampling algorithm maintains a sequence of  $N(n)$  independent samples for  $\mu_{\mathcal{I}}$  within total variation distance  $\epsilon(n)$ , while the entire specification of the graphical model  $\mathcal{I}$  is subject to dynamic update (to a new  $\mathcal{I}'$  with difference  $d(\mathcal{I}, \mathcal{I}') \leq L = o(n)$ ). Specifically, the algorithm updates the sample sequence within expected time  $O(\Delta^2 N(n) L \log^3 n + \Delta n)$ . Note that the extra  $O(\Delta n)$  cost is necessary for just editing the current MRF instance  $\mathcal{I}$  to  $\mathcal{I}'$  because a single update may change all the vertex and edge potentials simultaneously. This incremental time cost dominates the time cost of the dynamic inference algorithm, and is efficient for maintaining  $N(n)$  independent samples, especially when  $N(n)$  is sufficiently large, e.g.  $N(n) = \Omega(n/L)$ , in which case the average incremental cost for updating each sample is  $O(\Delta^2 L \log^3 n + \Delta n / N(n)) = O(\Delta^2 L \log^3 n)$ .

We illustrate the main idea by explaining how to maintain one sample. The idea is to represent the trace of the Markov chain for generating the sample by a dynamic data structure, and when the MRF instance is changed, this trace is modified to that of the new chain for generating the sample for the updated instance. This is achieved by both a set of efficient dynamic data structures and the coupling between the two Markov chains.

Specifically, let  $(X_t)_{t=0}^T$  be the Gibbs sampler chain for distribution  $\mu_I$ . When the chain is rapidly mixing, starting from an arbitrary initial configuration  $X_0 \in Q^V$ , after suitably many steps  $X = X_T$  is an accurate enough sample for  $\mu_I$ . At each step,  $X_{t-1}$  and  $X_t$  may differ only at a vertex  $v_t$  which is picked from  $V$  uniformly and independently at random. The evolution of the chain is fully captured by the initial state  $X_0$  and the sequence of pairs  $\langle v_t, X_t(v_t) \rangle$ , from  $t = 1$  to  $t = T$ , which is called the *execution log* of the chain in the paper.

Now suppose that the current instance  $I$  is updated to  $I'$ . We construct such a coupling between the original chain  $(X_t)_{t=0}^T$  and the new chain  $(Y_t)_{t=0}^T$ , such that  $(Y_t)_{t=0}^T$  is a faithful Gibbs sampling chain for the updated instance  $I'$  given that  $(X_t)_{t=0}^T$  is a faithful chain for  $I$ , and the difference between the two chains is small, in the sense that they have almost the same execution logs except for about  $O(TL/n)$  steps, where  $L$  is the difference between  $I$  and  $I'$ .

To simplify the exposition of such coupling, for now we restrict ourselves to the cases where the update to the instance  $I$  does not change the set of variables. Without loss of generality, we only consider the following two basic update operations that modifies  $I$  to  $I'$ .

- *Graph update.* The update only adds or deletes some edges, while all vertex potentials and the potentials of unaffected edges are not changed.
- *Hamiltonian update.* The update changes (possibly all) potentials of vertices and edges, while the underlying graph remains unchanged.

The general update of graphical model can be obtained by combining these two basic operations.

Then the new chain  $(Y_t)_{t=0}^T$  can be coupled with  $(X_t)_{t=0}^T$  by using the same initial configuration  $Y_0 = X_0$  and the same sequence  $v_1, v_2, \dots, v_T \in V$  of randomly picked vertices. And for  $t = 1, 2, \dots, T$ , the transition  $\langle v_t, Y_t(v_t) \rangle$  of the new chain can be generated using the same vertex  $v_t$  as in the original  $(X_t)_{t=0}^T$  chain, and a random  $Y_t(v_t)$  generated according to a coupling of the marginal distributions of  $X_t(v_t)$  and  $Y_t(v_t)$ , conditioning respectively on the current states of the neighborhood of  $v_t$  in  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$ . Note that these two marginal distributions must be identical unless **(I)**  $X_{t-1}$  and  $Y_{t-1}$  differ from each other over the neighborhood of  $v_t$  or **(II)** the  $v_t$  itself is incident to where the models  $I$  and  $I'$  differ. The event **(II)** occurs rarely due to the following reasons.

- For graph update, the event **(II)** occurs only if  $v_t$  is incident to an updated edge. Since only  $L$  edges are updated, the event occurs in at most  $O(TL/n)$  steps in expectation.
- For Hamiltonian update, all the potentials of vertices and edges can be changed, thus  $I, I'$  may differ everywhere. The key observation is that, as the total difference between the current and updated potentials is bounded by  $L$ , we can apply a filter to first select all candidate steps where the coupling may actually fail due to the difference between  $I$  and  $I'$ , which can be as small as  $O(TL/n)$ , and the actual coupling between  $(X_t)_{t=0}^\infty$  and  $(Y_t)_{t=0}^\infty$  is constructed with such prior.

Finally, when  $I$  and  $I'$  both satisfy the Dobrushin-Shlosman condition, the percolation of disagreements between  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$  is bounded, and we show that the two chains are almost always identically coupled as  $\langle v_t, X_t(v_t) \rangle = \langle v_t, Y_t(v_t) \rangle$ , with exceptions at only  $O(TL/n)$  steps. The original chain  $(X_t)_{t=0}^T$  can then be updated to the new chain  $(Y_t)_{t=0}^T$  by only editing these  $O(TL/n)$  local transitions  $\langle v_t, Y_t(v_t) \rangle$  which are different from  $\langle v_t, X_t(v_t) \rangle$ . This is aided by the dynamic data structure for the execution log of the chain, which is of independent interest.

## 6. DYNAMIC GIBBS SAMPLING

In this section, we give the dynamic sampling algorithm that updates the sample sequences.

In the following theorem, we use  $I = (V, E, Q, \Phi)$ , where  $n = |V|$ , to denote the current MRF instance and  $I' = (V', E', Q, \Phi')$ , where  $n' = |V'|$ , to denote the updated MRF instance. And define

$$d_{\text{graph}}(I, I') \triangleq |V \oplus V'| + |E \oplus E'|$$

$$d_{\text{Hamil}}(I, I') \triangleq \sum_{v \in V \cap V'} \|\phi_v - \phi'_v\|_1 + \sum_{e \in E \cap E'} \|\phi_e - \phi'_e\|_1.$$

Note that  $d(I, I') = d_{\text{graph}}(I, I') + d_{\text{Hamil}}(I, I')$ , where  $d(I, I')$  is defined in (2).

**Theorem 6.1 (dynamic sampling algorithm).** Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  be two functions satisfying the bounded difference condition in Definition 2.3. Assume that  $\mathcal{I}$  and  $\mathcal{I}'$  both satisfy Dobrushin-Shlosman condition,  $d_{\text{graph}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{graph}} = o(n)$  and  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{Hamil}}$ .

There is an algorithm that maintains a sequence of  $N(n)$  independent samples  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N(n))} \in Q^V$  where  $d_{\text{TV}}(\mu_{\mathcal{I}}, \mathbf{X}^{(i)}) \leq \epsilon(n)$  for all  $1 \leq i \leq N(n)$ , using  $O(nN(n) \log n)$  memory words, each of  $O(\log n)$  bits, such that when  $\mathcal{I}$  is updated to  $\mathcal{I}'$ , the algorithm updates the sequence to  $N(n')$  independent samples  $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(N(n'))} \in Q^{V'}$  where  $d_{\text{TV}}(\mu_{\mathcal{I}'}, \mathbf{Y}^{(i)}) \leq \epsilon(n')$  for all  $1 \leq i \leq N(n')$ , within expected time cost

$$(5) \quad O\left(\Delta^2(L_{\text{graph}} + L_{\text{Hamil}})N(n) \log^3 n + \Delta n\right),$$

where  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ , and  $\Delta_G, \Delta_{G'}$  denote the maximum degree of  $G = (V, E)$  and  $G' = (V', E')$ .

Our algorithm is based on the Gibbs sampling algorithm. Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  be two functions in Theorem 6.1. We first give the *single-sample dynamic Gibbs sampling algorithm* (Algorithm 2) that maintains a single sample  $\mathbf{X} \in Q^V$  for the current MRF instance  $\mathcal{I} = (V, E, Q, \Phi)$  where  $n = |V|$  such that  $d_{\text{TV}}(\mathbf{X}, \mu_{\mathcal{I}}) \leq \epsilon(n)$ . We then use this algorithm to obtain the *multi-sample dynamic Gibbs sampling algorithm* that maintains  $N(n)$  independent samples for the current instance.

Given the error function  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$ , suppose that  $T(\mathcal{I})$  is an easy-to-compute integer-valued function that upper bounds the mixing time on instance  $\mathcal{I}$ , such that

$$(6) \quad T(\mathcal{I}) \geq \tau_{\text{mix}}(\mathcal{I}, \epsilon(n)),$$

where  $\tau_{\text{mix}}(\mathcal{I}, \epsilon(n))$  denotes the mixing rate for the Gibbs sampling chain  $(\mathbf{X}_t)_{t \geq 0}$  on instance  $\mathcal{I}$ . By Proposition 4.3, if the Dobrushin-Shlosman condition is satisfied, we can set

$$(7) \quad T(\mathcal{I}) = \left\lceil \frac{n}{\delta} \log \frac{n}{\epsilon(n)} \right\rceil.$$

Our algorithm for single-sample dynamic Gibbs sampling maintains a random process  $(\mathbf{X}_t)_{t=0}^T$ , which is a Gibbs sampling chain on instance  $\mathcal{I}$  of length  $T = T(\mathcal{I})$ , where  $T(\mathcal{I})$  satisfies (6). Clearly  $\mathbf{X}_T$  is a sample for  $\mu_{\mathcal{I}}$  with  $d_{\text{TV}}(\mathbf{X}_T, \mu_{\mathcal{I}}) \leq \epsilon(n)$ .

When the current instance  $\mathcal{I}$  is updated to a new instance  $\mathcal{I}' = (V', E', Q, \Phi')$  where  $n' = |V'|$ , the original process  $(\mathbf{X}_t)_{t=0}^T$  is transformed to a new process  $(\mathbf{Y}_t)_{t=0}^{T'}$  such that the following holds as an invariant:  $(\mathbf{Y}_t)_{t=0}^{T'}$  is a Gibbs sampling chain on  $\mathcal{I}'$  with  $T' = T(\mathcal{I}')$ . Hence  $\mathbf{Y}_T$  is a sample for the new instance  $\mathcal{I}'$  with  $d_{\text{TV}}(\mathbf{Y}_T, \mu_{\mathcal{I}'}) \leq \epsilon(n')$ . This is achieved through the following two steps:

- (1) We construct couplings between  $(\mathbf{X}_t)_{t=0}^T$  and  $(\mathbf{Y}_t)_{t=0}^{T'}$ , so that the new process  $(\mathbf{Y}_t)_{t=0}^{T'}$  for  $\mathcal{I}'$  can be obtained by making small changes to the original process  $(\mathbf{X}_t)_{t=0}^T$  for  $\mathcal{I}$ .
- (2) We give a data structure which represents  $(\mathbf{X}_t)_{t=0}^T$  incrementally and supports various updates and queries to  $(\mathbf{X}_t)_{t=0}^T$  so that the above coupling can be generated efficiently.

**6.1. Coupling for dynamic instances.** The Gibbs sampling chain  $(\mathbf{X}_t)_{t=0}^T$  can be uniquely and fully recovered from: the initial state  $\mathbf{X}_0 \in Q^V$ , and the pairs  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  that record the transitions. We call  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  the *execution-log* for the chain  $(\mathbf{X}_t)_{t=0}^T$ , and denote it with

$$\text{Exe-Log}(\mathcal{I}, T) \triangleq \langle v_t, X_t(v_t) \rangle_{t=1}^T.$$

The following invariants are assumed for the random execution-log with an initial state.

**Condition 6.2 (invariants for Exe-Log).** Fixed an initial state  $\mathbf{X}_0 \in Q^V$ , the followings hold for the random execution-log  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for the Gibbs sampling chain  $(\mathbf{X}_t)_{t=0}^T$  on instance  $\mathcal{I} = (V, E, Q, \Phi)$ :

- $T = T(\mathcal{I})$  where  $T(\mathcal{I})$  satisfies (6);
- the random process  $(\mathbf{X}_t)_{t=0}^T$  uniquely recovered from the transitions  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  and the initial state  $\mathbf{X}_0$ , is identically distributed as the Gibbs sampling (Algorithm 1) on instance  $\mathcal{I}$  starting from initial state  $\mathbf{X}_0$  with  $v_t$  as the vertex picked at the  $t$ -th step.

Such invariants guarantee that  $X_T$  provides a sample for  $\mu_I$  with  $d_{TV}(X_T, \mu_I) \leq \epsilon(|V|)$ .

Suppose the current instance  $\mathcal{I}$  is updated to a new instance  $\mathcal{I}'$ . We construct couplings between the execution-log  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  with initial state  $X_0 \in Q^V$  for  $\mathcal{I}$  and the execution-log  $\text{Exe-Log}(\mathcal{I}', T') = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}$  with initial state  $Y_0 \in Q^{V'}$  for  $\mathcal{I}'$ . Our goal is as follows: assuming Condition 6.2 for  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T)$ , the same condition should hold invariantly for  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T')$ .

Unlike traditional coupling of Markov chains for the analysis of mixing time, where the two chains start from arbitrarily distinct initial states but proceed by the same transition rule, here the two chains  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^{T'}$  start from similar states but have to obey different transition rules due to differences between instances  $\mathcal{I}$  and  $\mathcal{I}'$ .

Due to the technical reason, we divide the update from  $\mathcal{I} = (V, E, Q, \Phi)$  to  $\mathcal{I}' = (V', E', Q, \Phi')$  into two steps: we first update  $\mathcal{I} = (V, E, Q, \Phi)$  to

$$(8) \quad \mathcal{I}_{\text{mid}} = (V, E, Q, \Phi^{\text{mid}}),$$

where the potentials  $\Phi^{\text{mid}} = (\phi_a^{\text{mid}})_{a \in V \cup E}$  in the middle instance  $\mathcal{I}_{\text{mid}}$  are defined as

$$\forall a \in V \cup E, \quad \phi_a^{\text{mid}} \triangleq \begin{cases} \phi'_a & \text{if } a \in V' \cup E' \\ \phi_a & \text{if } a \notin V' \cup E'; \end{cases}$$

then we update  $\mathcal{I}_{\text{mid}} = (V, E, Q, \Phi^{\text{mid}})$  to  $\mathcal{I}' = (V', E', Q, \Phi')$ . In other words, the update from  $\mathcal{I}$  to  $\mathcal{I}_{\text{mid}}$  is only caused by updating the potentials of vertices and edges, while the underlying graph remains unchanged; and the update from  $\mathcal{I}_{\text{mid}}$  to  $\mathcal{I}'$  is only caused by updating the underlying graph, i.e. adding vertices, deleting vertices, adding edges and deleting edges.

The dynamic Gibbs sampling algorithm can be outlined as follows.

- **UpdateHamiltonian:** update  $X_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  to a new initial state  $Z_0$  and a new execution log  $\text{Exe-Log}(\mathcal{I}_{\text{mid}}, T) = \langle u_t, Z_t(u_t) \rangle_{t=1}^T$  such that the random process  $(Z_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}_{\text{mid}}$ .
- **UpdateGraph:** update  $Z_0$  and  $\langle u_t, Z_t(u_t) \rangle_{t=1}^T$  to a new initial state  $Y_0$  and a new execution log  $\text{Exe-Log}(\mathcal{I}', T) = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}$  such that the random process  $(Y_t)_{t=0}^{T'}$  is the Gibbs sampling on instance  $\mathcal{I}'$ .
- **LengthFix:** change the length of the execution log  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  from  $T$  to  $T'$ , where  $T' = T(\mathcal{I}')$  and  $T(\mathcal{I}')$  satisfies (6).

The dynamic Gibbs sampling algorithm is given in Algorithm 2.

---

**Algorithm 2:** Dynamic Gibbs sampling

---

**Data** :  $X_0 \in Q^V$  and  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for current  $\mathcal{I} = (V, E, Q, \Phi)$ .

**Update:** an update that modifies  $\mathcal{I}$  to  $\mathcal{I}' = (V', E', Q, \Phi')$ .

- 1 compute  $T' = T(\mathcal{I}')$  satisfying (6) and construct  $\mathcal{I}_{\text{mid}} = (V', E', Q, \Phi^{\text{mid}})$  as in (8);
  - 2  $(Z_0, \langle u_t, Z_t(u_t) \rangle_{t=1}^T) \leftarrow \text{UpdateHamiltonian}(\mathcal{I}, \mathcal{I}_{\text{mid}}, X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$ ;  
// update the potentials:  $\mathcal{I} \rightarrow \mathcal{I}_{\text{mid}}$
  - 3  $(Y_0, \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}) \leftarrow \text{UpdateGraph}(\mathcal{I}_{\text{mid}}, \mathcal{I}', Z_0, \langle u_t, Z_t(u_t) \rangle_{t=1}^T)$ ;  
// update the underlying graph:  $\mathcal{I}_{\text{mid}} \rightarrow \mathcal{I}'$
  - 4  $(Y_0, \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}) \leftarrow \text{LengthFix}(\mathcal{I}', Y_0, \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T, T')$ , where  $T' = T(\mathcal{I}')$ ;  
// change the length of the execution log from  $T$  to  $T' = T(\mathcal{I}')$
  - 5 update the data to  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T') = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}$ ;
- 

The subroutine LengthFix is given in Algorithm 3. We then describe UpdateHamiltonian (Section 6.1.1) and UpdateGraph (Section 6.1.2).



---

**Algorithm 3:** LengthFix  $(\mathcal{I}, X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T, T')$ 


---

**Data** :  $X_0 \in Q^V$  and Exe-Log  $(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for current  $\mathcal{I} = (V, E, Q, \Phi)$ .

**Input** : the new length  $T' > 0$ .

1 **if**  $T' < T$  **then**

2   | truncate  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  to  $\langle v_t, X_t(v_t) \rangle_{t=1}^{T'}$ ;

3 **else**

4   | extend  $\langle v_t, X_t(v_t) \rangle_{t=1}^{T'}$  to  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  by simulating the Gibbs sampling chain on  $\mathcal{I}$  for  $T - T'$  more steps;

5 update the data to  $X_0$  and Exe-Log  $(\mathcal{I}, T') = \langle v_t, X_t(v_t) \rangle_{t=1}^{T'}$

---

6.1.1. *Coupling for Hamiltonian update.* We consider the update of changing potentials of vertices and edges. The update do not change the underlying graph. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance. Let  $X_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  be the current initial state and execution log such that the random process  $(X_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}$ . Upon such an update, the new instance becomes  $\mathcal{I}' = (V, E, Q, \Phi')$ . The algorithm UpdateHamiltonian  $(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$  updates the data to  $Y_0$  and  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

We transform the pair of  $X_0 \in Q^V$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  to a new pair of  $Y_0 \in Q^V$  and  $\langle v_t, Y_t(v_t) \rangle_{t=1}^T$  for  $\mathcal{I}'$ . This is achieved as follows: the vertex sequence  $(v_t)_{t=1}^T$  is identically coupled and the chain  $(X_t)_{t=0}^T$  is transformed to  $(Y_t)_{t=0}^T$  by the following one-step local coupling between  $X$  and  $Y$ .

**Definition 6.3 (one-step local coupling for Hamiltonian update).** The two chains  $(X_t)_{t=0}^\infty$  on instance  $\mathcal{I} = (V, E, Q, \Phi)$  and  $(Y_t)_{t=0}^\infty$  on instance  $\mathcal{I}' = (V, E, Q, \Phi')$  are coupled as:

- Initially  $X_0 = Y_0 \in Q^V$ ;
- for  $t = 1, 2, \dots$ , the two chains  $X$  and  $Y$  jointly do:
  - (1) pick the same  $v_t \in V$ , and let  $(X_t(u), Y_t(u)) \leftarrow (X_{t-1}(u), Y_{t-1}(u))$  for all  $u \in V \setminus \{v_t\}$ ;
  - (2) sample  $(X_t(v_t), Y_t(v_t))$  from a coupling  $D_{I_{v_t}, I'_{v_t}}^{\sigma, \tau}(\cdot, \cdot)$  of the marginal distributions  $\mu_{v_t, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v_t, \mathcal{I}'}(\cdot | \tau)$  with  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_G(v_t))$ , where  $G = (V, E)$ .

The local coupling  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  for Hamiltonian update is specified as follows.

**Definition 6.4 (local coupling  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  for Hamiltonian update).** Let  $v \in V$  be vertex and  $\sigma, \tau \in Q^{\Gamma_G(v)}$  two configurations, where  $G = (V, E)$ . We say a random pair  $(c, c') \in Q^2$  is drawn from the coupling  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  if  $(c, c')$  is generated by the following two steps:

- **sampling step:** sample  $(c, c') \in Q^2$  jointly from an optimal coupling  $D_{\text{opt}, I_v}^{\sigma, \tau}$  of the marginal distributions  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$ , such that  $c \sim \mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $c' \sim \mu_{v, \mathcal{I}'}(\cdot | \tau)$ ;
- **resampling step:** flip a coin independently with the probability of HEADS being

$$(9) \quad p_{I_v, I'_v}^\tau(c') \triangleq \begin{cases} 0 & \text{if } \mu_{v, \mathcal{I}}(c' | \tau) \leq \mu_{v, \mathcal{I}'}(c' | \tau), \\ \frac{\mu_{v, \mathcal{I}}(c' | \tau) - \mu_{v, \mathcal{I}'}(c' | \tau)}{\mu_{v, \mathcal{I}}(c' | \tau)} & \text{otherwise;} \end{cases}$$

if the outcome of coin flipping is HEADS, resample  $c'$  from the distribution  $v_{I_v, I'_v}^\tau$  independently, where the distribution  $v_{I_v, I'_v}^\tau$  is defined as

$$(10) \quad \forall b \in Q : \quad v_{I_v, I'_v}^\tau(b) \triangleq \frac{\max\{0, \mu_{v, \mathcal{I}'}(b | \tau) - \mu_{v, \mathcal{I}}(b | \tau)\}}{\sum_{x \in Q} \max\{0, \mu_{v, \mathcal{I}}(x | \tau) - \mu_{v, \mathcal{I}'}(x | \tau)\}}.$$

**Lemma 6.5.**  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  in Definition 6.4 is a valid coupling between  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$ .

By Lemma 6.5, the resulting  $(Y_t)_{t=0}^T$  is a faithful copy of the Gibbs sampling on instance  $\mathcal{I}'$ , assuming that  $(X_t)_{t=0}^T$  is such a chain on instance  $\mathcal{I}$ .

Next we give an upper bound for the probability  $p_{I_v, I'_v}^\tau(\cdot)$  defined in (9).



**Lemma 6.6.** For any two instances  $\mathcal{I} = (V, E, Q, \Phi)$  and  $\mathcal{I}' = (V, E, Q, \Phi')$  of MRF model, and any  $v \in V, c \in Q$  and  $\sigma \in Q^{\Gamma_G(v)}$ , it holds that

$$(11) \quad p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c) \leq 2 \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right),$$

where  $\|\phi_v - \phi'_v\|_1 = \sum_{c \in Q} |\phi_v(c) - \phi'_v(c)|$  and  $\|\phi_e - \phi'_e\|_1 = \sum_{c, c' \in Q} |\phi_e(c, c') - \phi'_e(c, c')|$ .

By Lemma 6.6, for each vertex  $v \in V$ , we define an upper bound of the probability  $p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(\cdot)$  as

$$(12) \quad p_v^{\text{up}} \triangleq \min \left\{ 2 \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right), 1 \right\}.$$

With  $p_v^{\text{up}}$ , we can implement the one-step local coupling in Definition 6.3 as follows. We first sample each  $v_i \in V$  for  $1 \leq i \leq T$  uniformly and independently. For each vertex  $v \in V$ , let  $T_v \triangleq \{1 \leq t \leq T \mid v_t = v\}$  be the set of all the steps that pick the vertex  $v$ . We select each  $t \in T_v$  independently with probability  $p_v^{\text{up}}$  to construct a random subset  $\mathcal{P}_v \subseteq T_v$ , and let

$$(13) \quad \mathcal{P} \triangleq \bigcup_{v \in V} \mathcal{P}_v.$$

We then couple the two chains  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$ . First set  $X_0 = Y_0$ . For each  $1 \leq t \leq T$ , we set  $(X_t(u), Y_t(u)) \leftarrow (X_{t-1}(u), Y_{t-1}(u))$  for all  $u \in V \setminus \{v_t\}$ ; then generate the random pair  $(X_t(v_t), Y_t(v_t))$  by the following procedure.

- **sampling step:** Let  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_G(v_t))$ . We draw a random pair  $(c, c') \in Q^2$  from the optimal coupling  $D_{\text{opt}, \mathcal{I}_v}^{\sigma, \tau}$  of the marginal distributions  $\mu_{v, \mathcal{I}}(\cdot \mid \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot \mid \tau)$  such that  $c \sim \mu_{v, \mathcal{I}}(\cdot \mid \sigma)$  and  $c' \sim \mu_{v, \mathcal{I}'}(\cdot \mid \tau)$ ;
- **resampling step:** If  $t \notin \mathcal{P}$ , set  $X_t(v_t) = c$  and  $Y_t(v_t) = c'$ . Otherwise, set  $X_t(v_t) = c$  and

$$(14) \quad Y_t(v_t) = \begin{cases} b \sim \nu_{\mathcal{I}_v, \mathcal{I}'_v}^\tau & \text{with probability } p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c')/p_{v_t}^{\text{up}} \\ c' & \text{with probability } 1 - p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c')/p_{v_t}^{\text{up}}. \end{cases}$$

Note that  $p_{v_t}^{\text{up}} > 0$  if  $t \in \mathcal{P}$ . By Lemma 6.6, it must hold that  $p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c') \leq p_{v_t}^{\text{up}}$ . Hence, the probability  $p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c')/p_{v_t}^{\text{up}}$  is valid. Note that the probability that  $Y_t(v_t)$  is set as  $b$  is

$$\Pr[Y_t(v_t) \text{ is set as } b] = \Pr[t \in \mathcal{P}] \cdot \frac{p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c')}{p_{v_t}^{\text{up}}} = p_{v_t}^{\text{up}} \cdot \frac{p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c')}{p_{v_t}^{\text{up}}} = p_{\mathcal{I}_v, \mathcal{I}'_v}^\tau(c').$$

Hence, our implementation perfectly simulates the coupling in Definition 6.3.

Let  $\mathcal{D}_t$  denote the set of disagreements between  $X_t$  and  $Y_t$ . Formally,

$$(15) \quad \mathcal{D}_t \triangleq \{v \in V \mid X_t(v) \neq Y_t(v)\}.$$

Note that if  $v_t \notin \Gamma_G(\mathcal{D}_{t-1})$ , the random pair  $(c, c')$  drawn from the coupling  $D_{\text{opt}, \mathcal{I}_v}^{\sigma, \tau}$  must satisfy  $c = c'$ . Thus it is easy to make the following observation for the  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$  coupled as above.

**Observation 6.7.** For any integer  $t \in [1, T]$ , if  $v_t \notin \Gamma_G^+(\mathcal{D}_{t-1})$  and  $t \notin \mathcal{P}$ , then  $X_t(v_t) = Y_t(v_t)$  and  $\mathcal{D}_t = \mathcal{D}_{t-1}$ .

With this observation, the new  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T) = \langle v_t, Y_t(v_t) \rangle_{t=1}^T$  can be generated from  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  as Algorithm 4.

Observation 6.7 says that the nontrivial coupling between  $X_t(v_t)$  and  $Y_t(v_t)$  is only needed when  $v_t \in \Gamma_G^+(\mathcal{D}_{t-1})$  or  $t \in \mathcal{P}$ , which occurs rarely as long as  $\mathcal{D}_{t-1}$  and  $\mathcal{P}$  are small. This is a key to ensure the small incremental time cost of Algorithm 4. For the  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$  coupled as above and any  $1 \leq t \leq T$ , let  $\gamma_t$  indicate whether the event  $t \in \mathcal{P} \vee v_t \in \Gamma_G^+(\mathcal{D}_{t-1})$  occurs:

$$(16) \quad \gamma_t \triangleq 1 \left[ t \in \mathcal{P} \vee v_t \in \Gamma_G^+(\mathcal{D}_{t-1}) \right],$$

---

**Algorithm 4:** UpdateHamiltonian  $\left(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T\right)$

---

**Data** :  $X_0 \in Q^V$  and Exe-Log  $(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for  $\mathcal{I} = (V, E, Q, \Phi)$ .

**Update:** an update that modifies  $\mathcal{I}$  to  $\mathcal{I}' = (V, E, Q, \Phi')$ .

- 1  $t_0 \leftarrow 0$ ,  $\mathcal{D} \leftarrow \emptyset$ , and construct a  $Y_0 \leftarrow X_0$ ;
  - 2 for each  $v \in V$ , construct a random subset  $\mathcal{P}_v \subseteq T_v \triangleq \{1 \leq t \leq T \mid v_t = v\}$  such that each element in  $T_v$  is selected independently with probability  $p_v^{\text{up}}$  defined in (12);
  - 3 construct the set  $\mathcal{P} \leftarrow \bigcup_{v \in V} \mathcal{P}_v$ ;
  - 4 **while**  $\exists t_0 < t \leq T$  such that  $v_t \in \Gamma_G^+(\mathcal{D})$  or  $t \in \mathcal{P}$  **do**
  - 5     find the smallest  $t > t_0$  such that  $v_t \in \Gamma_G^+(\mathcal{D})$  or  $t \in \mathcal{P}$ ;
  - 6     for all  $t_0 < i < t$ , let  $Y_i(v_i) = X_i(v_i)$ ;
  - 7     sample  $Y_t(v_t) \in Q$  conditioning on  $X_t(v_t)$  according to the optimal coupling between  $\mu_{v_t, \mathcal{I}}(\cdot \mid X_{t-1}(\Gamma_G(v_t)))$  and  $\mu_{v_t, \mathcal{I}'}(\cdot \mid Y_{t-1}(\Gamma_G(v_t)))$ ;
  - 8     **if**  $t \in \mathcal{P}$  **then**
  - 9         **with probability**  $p_{\mathcal{I}_{v_t}, \mathcal{I}'_{v_t}}^\tau(Y_t(v_t)) / p_{v_t}^{\text{up}}$  where  $\tau = Y_{t-1}(\Gamma_G(v_t))$  **do**
  - 10             resample  $Y_t(v_t) \sim v_{\mathcal{I}_{v_t}, \mathcal{I}'_{v_t}}^\tau$ , where  $v_{\mathcal{I}_{v_t}, \mathcal{I}'_{v_t}}^\tau$  is defined in (10) ;
  - 11     **if**  $X_t(v_t) \neq Y_t(v_t)$  **then**  $\mathcal{D} \leftarrow \mathcal{D} \cup \{v_t\}$  **else**  $\mathcal{D} \leftarrow \mathcal{D} \setminus \{v_t\}$ ;
  - 12      $t_0 \leftarrow t$ ;
  - 13 for all remaining  $t_0 < i \leq T$ : let  $Y_i(v_i) = X_i(v_i)$ ;
  - 14 update the data to  $Y_0$  and Exe-Log  $(\mathcal{I}', T) = \langle v_t, Y_t(v_t) \rangle_{t=1}^T$ ;
- 

and  $R_{\text{Hamil}}$  denote the number of occurrences of such bad events:

$$(17) \quad R_{\text{Hamil}} \triangleq \sum_{t=1}^T \gamma_t.$$

The following lemma bounds the expectation of  $R_{\text{Hamil}}$ .

**Lemma 6.8 (cost of the coupling for UpdateHamiltonian).** *Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance and  $\mathcal{I}' = (V, E, Q, \Phi')$  the updated instance. Assume that  $\mathcal{I}$  satisfies Dobrushin-Shlosman condition (Condition 3.1) with constant  $\delta > 0$ , and  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') = \sum_{v \in V} \|\phi_v - \phi'_v\|_1 + \sum_{e \in E} \|\phi_e - \phi'_e\|_1 \leq L$ . It holds that  $\mathbb{E}[R_{\text{Hamil}}] = O\left(\frac{\Delta T L}{n \delta}\right)$ , where  $n = |V|$ ,  $\Delta$  is the maximum degree of graph  $G = (V, E)$ .*

6.1.2. *Coupling for graph update.* Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance, where  $\Phi = (\phi_a)_{a \in V \cup E}$ . Let  $X_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  be the current initial state and execution log such that the random process  $(X_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}$ . Let  $\mathcal{I}' = (V', E', Q, \Phi')$  be the new instance obtained by updating the underlying graph, where  $\Phi' = (\phi_a)_{a \in V' \cup E'}$  satisfies

$$\forall a \in (V \cap V') \cap (E \cap E'), \quad \phi_a = \phi'_a.$$

Given the update from  $\mathcal{I}$  to  $\mathcal{I}'$ , the subroutine UpdateGraph  $\left(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T\right)$  updates the data to a new initial state  $Y_0$  and a new execution-log  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

The subroutine UpdateGraph does as the following three steps.

- AddVertex: add isolated vertices in  $V' \setminus V$  with potentials  $(\phi_v)_{v \in V' \setminus V}$ , and update the instance  $\mathcal{I} = (V, E, Q, \Phi)$  to a new instance

$$(18) \quad \mathcal{I}_1 = \mathcal{I}_1(\mathcal{I}, \mathcal{I}') \triangleq (V \cup V', E, Q, \Phi \cup (\phi_v)_{v \in V' \setminus V});$$

then update  $X_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  to  $Z_0$  and Exe-Log  $(\mathcal{I}_1, T) = \langle u_t, Z_t(u_t) \rangle_{t=1}^T$  such that the random process  $(Z_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}_1$ .

- **UpdateEdge**: add new edges in  $E' \setminus E$  with potentials  $(\phi_e)_{e \in E' \setminus E}$ , delete edges in  $E \setminus E'$ , and update the instance  $\mathcal{I}_1$  to a new instance

$$(19) \quad \begin{aligned} \mathcal{I}_2 &= \mathcal{I}_2(\mathcal{I}, \mathcal{I}') \triangleq (V \cup V', E', Q, \Phi \cup (\phi_v)_{v \in V' \setminus V} \cup (\phi_e)_{e \in E' \setminus E} \setminus (\phi_e)_{e \in E \setminus E'}) \\ &= (V \cup V', E', Q, \Phi' \cup (\phi_v)_{v \in V' \setminus V}); \end{aligned}$$

then update  $Z_0$  and  $\langle u_t, Z_t(u_t) \rangle_{t=1}^T$  to  $Z'_0$  and  $\text{Exe-Log}(\mathcal{I}_2, T) = \langle w_t, Z'_t(w_t) \rangle_{t=1}^T$  such that the random process  $(Z'_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}_2$ .

- **DeleteVertex**: delete isolated vertices in  $V \setminus V'$ , and update the instance  $\mathcal{I}_2$  to  $\mathcal{I}' = (V', E', Q, \Phi')$ ; then update  $Z'_0$  and  $\langle w_t, Z'_t(w_t) \rangle_{t=1}^T$  to  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T) = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

The algorithm **UpdateGraph** is given in Algorithm 5.

---

**Algorithm 5:** UpdateGraph  $(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$

---

**Data** :  $X_0 \in Q^V$  and  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for current  $\mathcal{I} = (V, E, Q, \Phi)$ .

**Update**: an update of the underlying graph that modifies  $\mathcal{I}$  to  $\mathcal{I}' = (V', E', Q, \Phi')$ .

- 1 construct instances  $\mathcal{I}_1$  and  $\mathcal{I}_2$  as in (18) and (19);
  - 2  $(Z_0, \langle u_t, Z_t(u_t) \rangle_{t=1}^T) \leftarrow \text{AddVertex}(\mathcal{I}, \mathcal{I}_1, X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$ ;  
// add isolated vertices to update  $\mathcal{I}$  to  $\mathcal{I}_1$
  - 3  $(Z'_0, \langle w_t, Z'_t(w_t) \rangle_{t=1}^T) \leftarrow \text{UpdateEdge}(\mathcal{I}_1, \mathcal{I}_2, Z_0, \langle u_t, Z_t(u_t) \rangle_{t=1}^T)$ ;  
// add and delete edges to update  $\mathcal{I}_1$  to  $\mathcal{I}_2$
  - 4  $(Y_0, \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T) \leftarrow \text{DeleteVertex}(\mathcal{I}_2, \mathcal{I}', Z'_0, \langle w_t, Z'_t(w_t) \rangle_{t=1}^T)$ ;  
// delete isolated vertices to update  $\mathcal{I}_2$  to  $\mathcal{I}'$
  - 5 update the data to  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}') = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$ ;
- 

The subroutines **AddVertex** and **DeleteVertex** are simple, because they only deal with isolated variables. We first describe the main subroutine **UpdateEdge**, then describe **AddVertex** and **DeleteVertex**.

**The coupling for UpdateEdge.** We first consider the update of adding and deleting edges. The update does not change the set of variables. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance. Let  $X_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  be the current initial state and execution log such that the random process  $(X_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}$ . Upon such an update, the new instance becomes  $\mathcal{I}' = (V, E', Q, \Phi')$ , where  $\phi'_a = \phi_a$  for all  $a \in V \cup (E \cap E')$ . The subroutine  $\text{UpdateEdge}(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$  updates the data to  $Y_0$  and  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

We use  $\mathcal{S} \subseteq V$  to denote the set of vertices affected by the update from  $\mathcal{I}$  to  $\mathcal{I}'$ :

$$(20) \quad \mathcal{S} \triangleq \bigcup_{(u,v) \in E \oplus E'} \{u, v\},$$

where  $E \oplus E'$  is the symmetric difference between  $E$  and  $E'$ .

We transform this pair of  $X_0 \in Q^V$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  to a new pair of  $Y_0 \in Q^V$  and  $\langle v_t, Y_t(v_t) \rangle_{t=1}^T$  for  $\mathcal{I}'$ . This is achieved as follows: the vertex sequence  $(v_t)_{t=1}^T$  is identically coupled and the chain  $(X_t)_{t=0}^T$  is transformed to  $(Y_t)_{t=0}^T$  by the following one-step local coupling between  $X$  and  $Y$ .

**Definition 6.9 (one-step local coupling for UpdateEdge).** The two chains  $(X_t)_{t=0}^\infty$  on instance  $\mathcal{I} = (V, E, Q, \Phi)$  and  $(Y_t)_{t=0}^\infty$  on instance  $\mathcal{I}' = (V, E', Q, \Phi')$  are coupled as:

- Initially  $X_0 = Y_0 \in Q^V$ ;
- for  $t = 1, 2, \dots$ , the two chains  $X$  and  $Y$  jointly do:
  - (1) pick the same  $v_t \in V$ , and let  $(X_t(u), Y_t(u)) \leftarrow (X_{t-1}(u), Y_{t-1}(u))$  for all  $u \in V \setminus \{v_t\}$ ;

- (2) sample  $(X_t(v_t), Y_t(v_t))$  from a coupling  $D_{\mathcal{I}_v, \mathcal{I}'_v}^{\sigma, \tau}(\cdot, \cdot)$  of the marginal distributions  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$  with  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_{G'}(v_t))$ , where  $G = (V, E)$  and  $G' = (V, E')$ .

The local coupling  $D_{\mathcal{I}_v, \mathcal{I}'_v}^{\sigma, \tau}(\cdot, \cdot)$  for UpdateEdge is specified as follows.

$$(21) \quad \forall \sigma \in Q^{\Gamma_G(v)}, \tau \in Q^{\Gamma_{G'}(v)} : \quad D_{\mathcal{I}_v, \mathcal{I}'_v}^{\sigma, \tau}(\cdot, \cdot) = \begin{cases} D_{\text{opt}, \mathcal{I}_v}^{\sigma, \tau}(\cdot, \cdot) & \text{if } v \notin \mathcal{S}, \\ \mu_{v, \mathcal{I}}(\cdot | \sigma) \times \mu_{v, \mathcal{I}'}(\cdot | \tau) & \text{if } v \in \mathcal{S}, \end{cases}$$

where  $D_{\text{opt}, \mathcal{I}_v}^{\sigma, \tau}$  is an optimal coupling of marginal distributions  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$ . Recall  $\mathcal{I}_v = (\Gamma_v^+, E_v, Q, \Phi_v)$  where  $E_v = \{\{u, v\} \in E\}$  and  $\Phi_v = (\phi_a)_{a \in \Gamma_v^+ \cup E_v}$ . Obviously,  $D_{\mathcal{I}_v, \mathcal{I}'_v}^{\sigma, \tau}$  is a valid coupling of  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$ . Because for any  $v \notin \mathcal{S}$ , we have  $\mathcal{I}_v = \mathcal{I}_{v'}$  and hence  $\mu_{v, \mathcal{I}}(\cdot | \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot | \tau)$  are the same, both defined by (4) on  $\mathcal{I}_v$ . Thus they can be coupled by  $D_{\text{opt}, \mathcal{I}_v}^{\sigma, \tau}$ .

Obviously the resulting  $(Y_t)_{t=0}^T$  is a faithful copy of the Gibbs sampling on instance  $\mathcal{I}'$ , assuming that  $(X_t)_{t=0}^T$  is such a chain on instance  $\mathcal{I}$ .

Recall  $\mathcal{D}_t \triangleq \{v \in V \mid X_t(v) \neq Y_t(v)\}$  is set of disagreements between  $X_t$  and  $Y_t$ . The following observation is easy to make for the  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$  coupled as above.

**Observation 6.10.** For any  $t \in [1, T]$ , if  $v_t \notin \mathcal{S} \cup \Gamma_G^+(\mathcal{D}_{t-1})$  then  $X_t(v_t) = Y_t(v_t)$  and  $\mathcal{D}_t = \mathcal{D}_{t-1}$ .

With this observation, the new  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T) = \langle v_t, Y_t(v_t) \rangle_{t=1}^T$  can be generated from  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  as in Algorithm 6.

---

**Algorithm 6:** UpdateEdge( $\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T$ )

---

**Data** :  $X_0 \in Q^V$  and  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for current  $\mathcal{I} = (V, E, Q, \Phi)$ .

**Update:** an update of adding and deleting edges that modifies  $\mathcal{I}$  to  $\mathcal{I}' = (V, E', Q, \Phi')$ .

- 1  $t_0 \leftarrow 0, \mathcal{D} \leftarrow \emptyset, Y_0 \leftarrow X_0$  and construct  $\mathcal{S} \leftarrow \bigcup_{(u,v) \in E \oplus E'} \{u, v\}$ ;
  - 2 **while**  $\exists t_0 < t \leq T$  such that  $v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D})$  **do**
  - 3     find the smallest  $t > t_0$  such that  $v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D})$ ;
  - 4     for all  $t_0 < i < t$ , let  $Y_i(v_i) = X_i(v_i)$ ;
  - 5     sample  $Y_t(v_t)$  conditioning on  $X_t(v_t)$  according to the coupling  $D_{v_t}^{\sigma, \tau}(\cdot, \cdot)$  (constructed in (21)), where  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_{G'}(v_t))$ ;
  - 6     **if**  $X_t(v_t) \neq Y_t(v_t)$  **then**  $\mathcal{D} \leftarrow \mathcal{D} \cup \{v_t\}$  **else**  $\mathcal{D} \leftarrow \mathcal{D} \setminus \{v_t\}$ ;
  - 7      $t_0 \leftarrow t$ ;
  - 8 for all remaining  $t_0 < i \leq T$ : let  $Y_i(v_i) = X_i(v_i)$ ;
  - 9 update the data to  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T) = \langle v_t, Y_t(v_t) \rangle_{t=1}^T$ ;
- 

Observation 6.10 says that the nontrivial coupling between  $X_t(v_t)$  and  $Y_t(v_t)$  is only needed when  $v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D}_{t-1})$ , which occurs rarely as long as  $\mathcal{D}_{t-1}$  remains small. This is a key to ensure the small incremental time cost of Algorithm 6. Formally, for the  $(X_t)_{t=0}^T$  and  $(Y_t)_{t=0}^T$  coupled as above, for any  $1 \leq t \leq T$ , let  $\gamma_t$  indicate whether this bad event occurs:

$$(22) \quad \gamma_t \triangleq 1 \left[ v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D}_{t-1}) \right],$$

and let  $R_{\text{graph}}$  denote the number of occurrences of such bad events:

$$(23) \quad R_{\text{graph}} \triangleq \sum_{t=1}^T \gamma_t.$$

We will see that  $R_{\text{graph}}$  dominates the cost of Algorithm 6, once a data structure is given to encode the execution-log and resolve the updates in Line 9 and various queries (in Lines 2, 3 and 5) to the data.

**Lemma 6.11 (cost of the coupling for UpdateEdge).** Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance and  $\mathcal{I}' = (V, E', Q, \Phi')$  the updated instance. Assume that  $\mathcal{I}'$  satisfies Dobrushin-Shlosman condition

(Condition 3.1) with constant  $\delta > 0$ , and  $|E \oplus E'| \leq L$ . It holds that  $\mathbb{E} [R_{\text{graph}}] = O\left(\frac{\Delta T L}{n\delta}\right)$ , where  $n = |V|$ ,  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ , and  $\Delta_G, \Delta_{G'}$  denote the maximum degree of  $G = (V, E)$  and  $G' = (V, E')$ .

**Coupling for AddVertex.** Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance. Let  $\mathbf{X}_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  be the current initial state and execution log such that the random process  $(\mathbf{X}_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}$ . The update adds a set of *isolated* vertices  $S$  with potentials  $(\phi_a)_{a \in S}$ . Upon such an update, the new instance becomes

$$\mathcal{I}' = (V', E, Q, \Phi') = (V \cup S, E, Q, \Phi \cup (\phi_a)_{a \in S}).$$

The subroutine  $\text{AddVertex}(\mathcal{I}, \mathcal{I}', \mathbf{X}_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$  updates the data to  $Y_0$  and  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

Since the new instance  $\mathcal{I}'$  is the same as  $\mathcal{I}$  except the isolated vertices in  $S$ , we can construct  $Y_0(V) = \mathbf{X}_0$  and  $Y_0(S) \in Q^S$  is arbitrary, and  $\text{Exe-Log}(\mathcal{I}', T) = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  can be constructed by inserting random appearances of vertices in  $S$  into  $(v_t)_{t=1}^T$ , while for any  $v \in S$ , the  $Y_t(v)$  at the inserted steps  $t$  are sampled i.i.d. from the marginal distribution  $\mu_{v, \mathcal{I}'}(\cdot)$ , which is just a distribution over  $Q$  proportional to  $\exp(\phi_v(\cdot))$  in the case of Gibbs sampling, since  $v$  is an isolated vertex. Let  $[T] \triangleq \{1, 2, \dots, T\}$ . Formally:

- (1) Let  $P \subseteq [T]$  be a random subset such that each  $t \in [T]$  is selected into  $P$  independently with probability  $\frac{|S|}{|S \cup V|}$ . Let  $h = |P|$  and enumerate all elements in  $P$  as  $r_1 < r_2 < \dots < r_h$ . Let  $m = T - h$  and enumerate all elements in  $[T] \setminus P$  as  $\ell_1 < \ell_2 < \dots < \ell_m$ .
- (2) For each  $1 \leq i \leq h$ , sample  $u_i \in S$  uniformly and independently.
- (3) Let  $\langle v_t, X_t(v_t) \rangle_{t=1}^m \leftarrow \text{LengthFix}(\mathcal{I}, \mathbf{X}_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T, m)$ .
- (4) Construct  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  as follows:

$$\begin{aligned} \forall t = r_k \in P: \quad & v'_t = u_k \quad \text{and} \quad Y_t(v'_t) \sim \mu_{u_k, \mathcal{I}'}(\cdot), \text{ where } \mu_{u_k, \mathcal{I}'}(c) \propto \exp(\phi_{u_k}(c)); \\ \forall t = \ell_k \in [T'] \setminus P: \quad & v'_t = v_k \quad \text{and} \quad Y_t(v'_t) = X_k(v'_t) = X_k(v_k). \end{aligned}$$

It is easy to see that  $(Y_t)_{t=0}^T$  is a faithful copy of the Gibbs sampling on instance  $\mathcal{I}'$ .

**Coupling for DeleteVertex.** Let  $\mathcal{I} = (V, E, Q, \Phi)$  be the current MRF instance. The update deletes a set of *isolated* variables  $S \subseteq V$ . Let  $\mathbf{X}_0$  and  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  be the current initial state and execution log such that the random process  $(\mathbf{X}_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}$ . Upon such update, the instance is updated to  $\mathcal{I}' = (V', E, Q, \Phi')$ , where  $V' = V \setminus S$  and  $\Phi' = \Phi \setminus (\phi_v)_{v \in S}$ . The subroutine  $\text{DeleteVertex}(\mathcal{I}, \mathcal{I}', \mathbf{X}_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$  updates the data to  $Y_0$  and  $\langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  such that the random process  $(Y_t)_{t=0}^T$  is the Gibbs sampling on instance  $\mathcal{I}'$ .

We can simply construct  $Y_0 = \mathbf{X}_0(V')$ . The new execution-log  $\text{Exe-Log}(\mathcal{I}', \epsilon) = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^T$  can be constructed from the original  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  by simply deleting all appearances of vertices  $v \in S$  in  $(v_t)_{t=1}^T$  and the corresponding trivial transitions  $X_t(v)$ , followed by calling  $\text{LengthFix}$  on instance  $\mathcal{I}'$  to properly append the chain to the length  $T$ .

It is easy to see that  $(Y_t)_{t=0}^T$  is a faithful copy of the Gibbs sampling on instance  $\mathcal{I}'$ .

**6.2. Data structure for Gibbs sampling.** We now describe an efficient data structure for Gibbs sampling  $(\mathbf{X}_t)_{t=0}^T$ . Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance. The data structure should provide the following functionalities.

- **Data:** an initial state  $\mathbf{X}_0 \in Q^V$  and an execution-log  $\langle v_t, X_t(v_t) \rangle_{t=1}^T \in (V \times Q)^T$  that records the  $T$  transitions of the Gibbs sampling  $(\mathbf{X}_t)_{t=0}^T$ ;
- **updates:**
  - $\text{Insert}(t, v, c)$ , which inserts a transition  $\langle v, c \rangle$  after the  $(t-1)$ -th transition  $\langle v_{t-1}, X_{t-1}(v_{t-1}) \rangle$ ;
  - $\text{Remove}(t)$ , which deletes the  $t$ -th transition  $\langle v_t, X_t(v_t) \rangle$ ;
  - $\text{Change}(t, c)$ , which changes the  $t$ -th transition  $\langle v_t, X_t(v_t) \rangle$  to  $\langle v_t, c \rangle$ ;

Note that the updates  $\text{Insert}(t, v, c)$  and  $\text{Remove}(t)$  change the length  $T$  of the chain, as well as the order-numbers of all transitions after the inserted/deleted transition.



- **queries:**

- $\text{Eval}(t, v)$ , which returns the value of  $X_t(v)$  for arbitrary  $t$  and  $v$  (not necessarily  $= v_t$ );
- $\text{Succ}(t, v)$ , which returns  $i$  for the smallest  $i > t$  such that  $v_i = v$  if such  $i$  exists, or returns  $\perp$  if otherwise.

It is not difficult to realize that the query  $\text{Eval}(t, v)$  can actually be solved by a predecessor search defined symmetrically to  $\text{Succ}(t, v)$ . This data structure problem for Gibbs sampling is quite natural and is of independent interest.

**Theorem 6.12 (data structure for Gibbs sampling).** *There exists a deterministic dynamic data structure which stores an arbitrary initial state  $X_0 \in Q^V$  and an execution-log  $\langle v_t, X_t(v_t) \rangle_{t=1}^T \in (V \times Q)^T$  for Gibbs sampling using  $O(T + |V|)$  memory words, each of  $O(\log T + \log |V| + \log |Q|)$  bits, such that each operation among Insert, Remove, Change, Eval and Succ can be resolved in time  $O(\log^2 T + \log |V|)$ .*

*Proof.* The initial state and execution-log are stored by separate data structures.

The initial state  $X_0 \in Q^V$  is maintained by a deterministic dynamic dictionary, with  $(v, X_0(v))$  for vertices  $v \in V$  as the key-value pairs. Such a deterministic data structure answers queries of  $X_0(v)$  given any  $v \in V$  while  $V$  is dynamically changing.

The execution-log  $\langle v_t, X_t(v_t) \rangle_{t=1}^T \in (V \times Q)^T$  is stored by  $|V|$  balanced search trees  $(\mathcal{T}_v)_{v \in V}$  (e.g., red-black trees). In each tree  $\mathcal{T}_v$ , each node in  $\mathcal{T}_v$  stores a distinct transition  $\langle v_t, X_t(v_t) \rangle$  with  $v_t = v$ , such that the in-order tree walk of  $\mathcal{T}_v$  prints all  $\langle v_t, X_t(v_t) \rangle$  with  $v_t = v$  in the order they appear in the execution-log  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$ . Altogether these trees  $(\mathcal{T}_v)_{v \in V}$  have  $T$  nodes in total. Besides, these trees  $(\mathcal{T}_v)_{v \in V}$  are indexed by another deterministic dynamic dictionary, with  $(v, p_v)$  for vertices  $v \in V$  as key-value pairs, where each  $p_v$  is the pointer to the root of tree  $\mathcal{T}_v$ . This dictionary provides random accesses to the trees  $\mathcal{T}_v$  for all  $v \in V$ , while  $V$  is dynamically changing.

Given any  $t$ , we want to answer predecessor (or successor) search for the largest  $i \leq t$  (or smallest  $i > t$ ) such that  $v_i = v$ . This is achieved with assistance from another data structure, an *order-statistic tree* (or *OS-tree*)  $\widehat{\mathcal{T}}$  [CLRS09, Section 14]. In  $\widehat{\mathcal{T}}$ , each node stores the “identity” of an individual transition  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  (which is actually a pointer to the node storing the transition  $\langle v_t, X_t(v_t) \rangle$  in the tree  $\mathcal{T}_v$  with  $v_t = v$ ). In particular, the in-order tree walk of  $\widehat{\mathcal{T}}$  prints all  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$  in that order. Such a data structure supports two query functions: (1) Select: given any  $t$ , returns the identity of the  $t$ -th transition  $\langle v_t, X_t(v_t) \rangle$ ; and (2) Rank: given the identity of any transition  $\langle v_t, X_t(v_t) \rangle$ , returns its rank  $t$  in the sequence  $\langle v_t, X_t(v_t) \rangle_{t=1}^T$ . Besides, the OS-tree  $\widehat{\mathcal{T}}$  also supports standard insertion (of a new transition  $\langle v, c \rangle$  to a given rank  $t$ ) and deletion (of the transition  $\langle v_t, X_t(v_t) \rangle$  at a given rank  $t$ ). As a balanced tree, all these queries and updates for the OS-tree  $\widehat{\mathcal{T}}$  can be resolved in  $O(\log T)$  time.

The successor and predecessor searches mentioned above for any  $v \in T$  and  $t$ , can then be resolved by binary searches in the balanced search tree  $\mathcal{T}_v$  while querying the OS-tree  $\widehat{\mathcal{T}}$  as an oracle for ordering, which takes time at most  $O(\log^2 T + \log |V|)$  in total, where the  $\log |V|$  cost is used for accessing the root of  $\mathcal{T}_v$  via the dynamic dictionary that indexes the trees  $(\mathcal{T}_v)_{v \in V}$ .

This solves the successor query  $\text{Succ}(t, v)$  as well as the evaluation query  $\text{Eval}(t, v)$  for Gibbs sampling, both within time cost  $O(\log^2 T + \log |V|)$ , where the latter is actually solved by the predecessor search for the largest  $i \leq t$  such that  $v_i = v$  and returning the value of  $X_i(v_i)$  recorded in the  $i$ -th transition  $\langle v_i, X_i(v_i) \rangle$  or returning the value of  $X_0(v)$  if no such  $i$  exists.

It is also easy to verify that with the above dynamic data structures, all updates, including:  $\text{Insert}(t, v, c)$ ,  $\text{Remove}(t)$  and  $\text{Change}(t, c)$ , can be implemented with cost at most  $O(\log^2 T + \log |V|)$ , and the data structures together use  $O(T + |V|)$  words in total, where each word consists of  $O(\log T + \log |V| + \log |Q|)$  bits.  $\square$

**6.3. Single-sample dynamic Gibbs sampling algorithm.** With the data structure for Gibbs sampling stated in Theorem 6.12, the couplings constructed in Section 6.1 can be implemented as the algorithm for dynamic Gibbs sampling. Recall  $d_{\text{graph}}(\cdot, \cdot)$  and  $d_{\text{Hamil}}(\cdot, \cdot)$  are defined in (2).

**Lemma 6.13 (single-sample dynamic Gibbs sampling algorithm).** *Let  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  be an error function. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance with  $n = |V|$  and  $\mathcal{I}' = (V', E', Q, \Phi')$  the updated instance with  $n' = |V'|$ . Denote  $T = T(\mathcal{I})$ ,  $T' = T(\mathcal{I}')$  and  $T_{\max} = \max\{T, T'\}$ . Assume*



$d_{\text{graph}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{graph}} = o(n)$ ,  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{Hamil}}$ , and  $T, T' \in \Omega(n \log n)$ . The single-sample dynamic Gibbs sampling algorithm (Algorithm 2) does the followings:

- **(space cost)** The algorithm maintains an explicit copy of a sample  $X \in Q^V$  for the current instance  $\mathcal{I}$ , and also a data structure using  $O(T)$  memory words, each of  $O(\log T)$  bits, for representing an initial state  $X_0 \in Q^V$  and an execution-log  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  for the Gibbs sampling  $(X_t)_{t=0}^T$  on  $\mathcal{I}$  generating sample  $X = X_T$ .
- **(correctness)** Assuming that Condition 6.2 holds for  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T)$  for the Gibbs sampling on  $\mathcal{I}$ , upon each update that modifies  $\mathcal{I}$  to  $\mathcal{I}'$ , the algorithm updates  $X$  to an explicit copy of a sample  $Y \in Q^{V'}$  for the new instance  $\mathcal{I}'$ , and correspondingly updates the  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T)$  represented by the data structure to a  $Y_0 \in Q^{V'}$  and  $\text{Exe-Log}(\mathcal{I}', T') = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}$  for the Gibbs sampling  $(Y_t)_{t=0}^{T'}$  on  $\mathcal{I}'$  generating the new sample  $Y = Y_{T'}$ , where  $Y_0$  and  $\text{Exe-Log}(\mathcal{I}', T')$  satisfy Condition 6.2 for the Gibbs sampling on  $\mathcal{I}'$ , therefore,

$$d_{\text{TV}}(Y, \mu_{\mathcal{I}'}) \leq \epsilon(n').$$

- **(time cost)** Assuming Condition 6.2 for  $X_0$  and  $\text{Exe-Log}(\mathcal{I}, T)$  for the Gibbs sampling on  $\mathcal{I}$ , the expected time complexity for resolving an update is

$$O\left(\Delta n + \Delta \left(|T - T'| + \frac{T_{\max}(L_{\text{Hamil}} + L_{\text{graph}})}{n} + \mathbb{E}[R_{\text{Hamil}}] + \mathbb{E}[R_{\text{graph}}]\right) \log^2 T_{\max}\right),$$

where  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ ,  $\Delta_G, \Delta_{G'}$  denote the maximum degrees of  $G = (V, E)$  and  $G' = (V', E')$ ,  $R_{\text{Hamil}}$  is defined in (17) for the subroutine UpdateHamiltonian in Algorithm 2, and  $R_{\text{graph}}$  is defined in (23) for the subroutine UpdateEdge in Algorithm 2.

We remark that the  $O(\Delta n)$  in time cost is necessary because the update from  $\mathcal{I}$  to  $\mathcal{I}'$  may change all the potentials of vertices and edges. One can reduce the  $O(\Delta n)$  from the time cost if we further restrict that one update can only change constant number of vertices, edges, and potentials.

The following result is a corollary from Lemma 6.13.

**Corollary 6.14.** Assume  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  in Lemma 6.13 satisfies the bounded difference condition in Definition 2.3. Assume  $\mathcal{I}$  and  $\mathcal{I}'$  in Lemma 6.13 both satisfy Dobrushin-Shlosman condition (Condition 3.1) with constant  $\delta > 0$ . The single-sample dynamic Gibbs sampling algorithm (Algorithm 2) uses  $O(n \log n)$  memory words, each of  $O(\log n)$  bits to maintain the sample for current instance  $\mathcal{I}$ , and resolves the update from  $\mathcal{I}$  to  $\mathcal{I}'$  with expected time cost  $O\left(\Delta n + \Delta^2(L_{\text{graph}} + L_{\text{Hamil}}) \log^3 n\right)$ .

*Proof of Lemma 6.13.* The dynamic Gibbs sampling algorithm is implemented as follows. The algorithm uses the dynamic data structure in Theorem 6.12 to maintain the initial state  $X_0$  and execution-log  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$ . Besides, the algorithm maintains the explicit copy of the sample  $X \in Q^V$  by a deterministic dynamic dictionary, with  $(v, X(v))$  for vertices  $v \in V$  as the key-value pairs. The lemma is proved as follows.

**Space cost:** Note that  $T = \Omega(n \log n)$ ,  $|V| = n$  and  $|Q| = O(1)$ . We have  $O(n) = O(T)$  and  $O(\log T + \log |V| + \log |Q|) = O(\log T)$ . The dynamic dictionary for sample  $X$  uses  $O(n)$  memory words, each of  $O(\log |V| + \log |Q|)$  bits. Combining with Theorem 6.12, we have the algorithm uses  $O(T)$  memory words to maintain the initial state, execution-log and the random sample, each word is of  $O(\log T + \log |V| + \log |Q|) = O(\log T)$  bits.

**Correctness:** The invariants for execution-log (Condition 6.2) are preserved by the coupling simulated by the algorithm. The correctness holds as a consequence.

**Time cost:** Consider the update that modifies  $\mathcal{I}$  to  $\mathcal{I}'$ . We divide the algorithm into two stages.

- **Preparation stage:** construct the updated instances  $\mathcal{I}'$  and other middle instances  $\mathcal{I}_{\text{mid}}, \mathcal{I}_1, \mathcal{I}_2$  in (8), (18), (19); compute  $p_v^{\text{up}}$  in (12) for all  $v \in V$  and construct the random set  $\mathcal{P} \subseteq [T] = \{1, 2, \dots, T\}$  in (13).
- **Update stage:** given  $\mathcal{P}$  and  $p_v^{\text{up}}$  for all  $v \in V$ , update the initial state  $X_0$  to  $Y_0$ , the execution-log  $\text{Exe-Log}(\mathcal{I}, T) = \langle v_t, X_t(v_t) \rangle_{t=1}^T$  to  $\text{Exe-Log}(\mathcal{I}', T') = \langle v'_t, Y_t(v'_t) \rangle_{t=1}^{T'}$ , and the sample  $X$  to  $Y$ .

We make the following two claims.

**Claim 6.15.** *The expected running time of the preparation stage is*

$$\mathbb{E} \left[ T_{\text{preparation}}^{\text{single}} \right] = O \left( \Delta n + \mathbb{E} [|\mathcal{P}|] \log^2 T_{\text{max}} \right),$$

and the expected size of  $\mathcal{P}$  is at most  $\frac{4T_{\text{max}}L_{\text{Hamil}}}{n}$ .

**Claim 6.16.** *The expected running time of the update stage is*

$$(24) \quad \mathbb{E} \left[ T_{\text{update}}^{\text{single}} \right] = O \left( \Delta \left( |T - T'| + \frac{T_{\text{max}}L_{\text{graph}}}{n} + \mathbb{E} [R_{\text{Hamil}}] + \mathbb{E} [R_{\text{graph}}] \right) \log^2 T_{\text{max}} \right),$$

$R_{\text{Hamil}}$  is defined in (17) for the subroutine *UpdateHamiltonian* in Algorithm 2, and  $R_{\text{graph}}$  is defined in (23) for the subroutine *UpdateEdge* in Algorithm 2.

By the linearity of expectation, the expected time cost of the algorithm is  $\mathbb{E} \left[ T_{\text{preparation}}^{\text{single}} \right] + \mathbb{E} \left[ T_{\text{update}}^{\text{single}} \right]$ . This proves the time cost.  $\square$

We introduce the following technique lemma to prove Corollary 6.14.

**Lemma 6.17.** *Let  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  be a function such that there exists a constant  $C > 0$  such that*

$$\forall n \in \mathbb{N}^+ : \quad |\epsilon(n+1) - \epsilon(n)| \leq \frac{C}{n} \epsilon(n).$$

Then the function  $N$  has the following properties

- for any  $n \in \mathbb{N}^+$ , it holds that  $\epsilon(n) \geq \frac{1}{\text{poly}(n)}$ ;
- let  $\alpha \geq 1$  be a constant, given any  $n, n' \in \mathbb{N}^+$  such that  $\frac{1}{\alpha} \leq \frac{n'}{n} \leq \alpha$ ,

$$\left| n \log \frac{n}{\epsilon(n)} - n' \log \frac{n'}{\epsilon(n')} \right| = C' |n' - n| \log n.$$

where  $C'$  is a constant that depends only on  $\alpha, C$  and  $\epsilon(3\lceil C \rceil)$ .

*Proof.* By the condition, we have  $\epsilon(t) \leq (1 + \frac{C}{t-C}) \epsilon(t+1)$  for all  $t > \lceil C+1 \rceil$ . Thus for all  $n > l = 3\lceil C \rceil$ ,

$$(25) \quad \epsilon(l) \leq \prod_{i=l}^{n-1} \left( 1 + \frac{C}{i-C} \right) \epsilon(n) \leq \epsilon(n) \exp \left( C \sum_{i=2}^{n-1} \frac{1}{i} \right) \leq \epsilon(n) \exp(C \ln n) = \epsilon(n) n^C.$$

Thus, we have  $\epsilon(n) \geq \frac{1}{\text{poly}(n)}$ .

We then prove the second property. It is lossless to assume that  $\min\{n, n'\} \geq l$ , since otherwise we can choose  $C'$  sufficiently large so that the second property holds. Firstly, we prove for the case  $n > n'$ . We have  $\left| \log \frac{n}{n'} \right| \leq \frac{n-n'}{n'}$ . By  $\epsilon(t) \leq (1 + \frac{C}{t-C}) \epsilon(t+1)$  for all  $t > \lceil C+1 \rceil$ , we also have

$$\epsilon(n') \leq \prod_{i=n'}^{n-1} \left( 1 + \frac{C}{i-C} \right) \epsilon(n) \leq \epsilon(n) \exp \left( \frac{C(n-n')}{n'-C} \right).$$

Thus,

$$(26) \quad \left| \log \frac{n}{\epsilon(n)} - \log \frac{n'}{\epsilon(n')} \right| \leq \left| \log \frac{n}{n'} \right| + \left| \log \frac{\epsilon(n)}{\epsilon(n')} \right| \leq \frac{n-n'}{n'} + \frac{C(n-n')}{n'-C} \leq \frac{(2C+1)(n-n')}{n'}.$$

The last equality is due to  $2(n'-C) \geq n'+l-2C \geq n'$ . Let  $C' = 2 + \lceil \log \epsilon(l) \rceil + 3C$ . We have

$$\left| n \log \frac{n}{\epsilon(n)} - n' \log \frac{n'}{\epsilon(n')} \right| \leq \left| (n'-n) \log \frac{n}{\epsilon(n)} \right| + \left| n' \left( \log \frac{n}{\epsilon(n)} - \log \frac{n'}{\epsilon(n')} \right) \right| \leq C' |n' - n| \log n.$$

The last inequality is due to (25) and (26). Similarly, we can also prove the lemma if  $n < n'$ .  $\square$

*Proof of Corollary 6.14.* By  $L_{\text{graph}} = o(n)$ , we have  $n' = \Theta(n)$ . Since  $\mathcal{I}$  and  $\mathcal{I}'$  both satisfy Dobrushin-Shlosman condition (Condition 3.1) with constant  $\delta > 0$ , we can set  $T, T'$  as in (7) such that

$$T = \left\lceil \frac{n}{\delta} \log \frac{n}{\epsilon(n)} \right\rceil = \Theta(n \log n)$$

$$T' = \left\lceil \frac{n'}{\delta} \log \frac{n'}{\epsilon(n')} \right\rceil = \Theta(n \log n).$$

The equations hold because  $n' = \Theta(n)$  and the error function  $\epsilon$  satisfies  $\epsilon(\ell) \geq \frac{1}{\text{poly}(\ell)}$  by Lemma 6.17. Thus, we have

$$(27) \quad T_{\max} = \max\{T, T'\} = O(n \log n).$$

By Lemma 6.17 and  $|n' - n| \leq L_{\text{graph}} = o(n)$ , we have

$$(28) \quad |T - T'| = O(L_{\text{graph}} \log n).$$

Let  $\mathcal{I}_{\text{mid}} = (V, E, Q, \Phi^{\text{mid}})$  be the middle instance constructed as in (8). In Algorithm 2, we call the subroutine UpdateHamiltonian for instances  $\mathcal{I}$  and  $\mathcal{I}_{\text{mid}}$ . Since  $\mathcal{I}$  satisfies the Dobrushin-Shlosman condition, by Lemma 6.8 and  $d(\mathcal{I}, \mathcal{I}_{\text{mid}}) \leq d(\mathcal{I}, \mathcal{I}') \leq L_{\text{Hamil}}$ , we have

$$(29) \quad \mathbb{E}[R_{\text{Hamil}}] = O\left(\frac{\Delta T L_{\text{Hamil}}}{\delta n}\right) = O(\Delta L_{\text{Hamil}} \log n),$$

where  $R_{\text{Hamil}}$  is defined in (17) for the subroutine UpdateHamiltonian.

We also call the subroutine UpdateGraph for instances  $\mathcal{I}_{\text{mid}}$  and  $\mathcal{I}'$  in Algorithm 2. The subroutine is shown in Algorithm 5. We first add isolated vertices to update  $\mathcal{I}_{\text{mid}}$  to  $\mathcal{I}_1$ , then update edges to update  $\mathcal{I}_1$  to  $\mathcal{I}_2$ , finally delete isolated vertices to update  $\mathcal{I}_2$  to  $\mathcal{I}'$ . Since  $\mathcal{I}'$  satisfies Dobrushin-Shlosman condition and the only difference between  $\mathcal{I}_2$  and  $\mathcal{I}'$  is that  $\mathcal{I}_2$  contains extra isolated vertices, it is easy to verify that  $\mathcal{I}_2$  also satisfies Dobrushin-Shlosman condition. In Algorithm 5, the subroutine UpdateEdge is called for  $\mathcal{I}_1$  and  $\mathcal{I}_2$ . By Lemma 6.11, we have

$$(30) \quad \mathbb{E}[R_{\text{graph}}] = O\left(\frac{\Delta T L_{\text{graph}}}{\Delta n}\right) = O(\Delta L_{\text{graph}} \log n).$$

where  $R_{\text{graph}}$  is defined in (23) for the subroutine UpdateEdge.

Combining (27), (28), (29), (30) with Lemma 6.13, we have the expected time cost is

$$\begin{aligned} \mathbb{E}[T_{\text{cost}}] &= O\left(\Delta n + \Delta \left(|T - T'| + \frac{T_{\max}(L_{\text{Hamil}} + L_{\text{graph}})}{n} + \mathbb{E}[R_{\text{Hamil}}] + \mathbb{E}[R_{\text{graph}}]\right) \log^2 T_{\max}\right) \\ &= O\left(\Delta n + \Delta^2(L_{\text{graph}} + L_{\text{Hamil}}) \log^3 n\right). \quad \square \end{aligned}$$

**6.4. Multi-sample dynamic Gibbs sampling algorithm.** In this section, we give an *Multi-sample dynamic Gibbs sampling algorithm* that maintains multiple independent random samples for the current MRF instance. Theorem 6.1 follows immediately from the following lemma.

**Lemma 6.18 (multi-sample dynamic Gibbs sampling algorithm).** *Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  and  $\epsilon : \mathbb{N}^+ \rightarrow (0, 1)$  be two functions satisfying the bounded difference condition in Definition 2.3. Let  $\mathcal{I} = (V, E, Q, \Phi)$  be an MRF instance with  $n = |V|$  and  $\mathcal{I}' = (V', E', Q, \Phi')$  the updated instance with  $n' = |V'|$ . Assume that  $\mathcal{I}$  and  $\mathcal{I}'$  both satisfy Dobrushin-Shlosman condition with constant  $\delta > 0$ ,  $d_{\text{graph}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{graph}} = o(n)$  and  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{Hamil}}$ . Denote  $T = \lceil \frac{n}{\delta} \log \frac{n}{\epsilon(n)} \rceil$ ,  $T' = \lceil \frac{n'}{\delta} \log \frac{n'}{\epsilon(n')} \rceil$ .*

*There is an algorithm which does the followings:*

- (**space cost**) *The algorithm maintains  $N(n)$  explicit copies of independent samples  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N(n))}$ , where  $\mathbf{X}^{(i)} \in Q^V$  for all  $1 \leq i \leq N(n)$ , for the current instance  $\mathcal{I}$ , and also a data structure using  $O(nN(n) \log n)$  memory words, each of  $O(\log n)$  bits, for representing the initial state  $\mathbf{X}_0^{(i)} \in Q^V$  and the execution-log  $\text{Exe-Log}^{(i)}(\mathcal{I}, T) = \left\langle v_t^{(i)}, X_t^{(i)}(v_t^{(i)}) \right\rangle_{t=1}^T$  for  $1 \leq i \leq N(n)$  such that each Gibbs sampling  $(\mathbf{X}_t^{(i)})_{t=0}^T$  on  $\mathcal{I}$  generating an independent sample  $\mathbf{X}^{(i)} = \mathbf{X}_T^{(i)}$ .*

- **(correctness)** Assuming that Condition 6.2 holds for each  $X_0^{(i)}$  and  $\text{Exe-Log}^{(i)}(I, T)$  for the Gibbs sampling on  $I$ , upon each update that modifies  $I$  to  $I'$ , the algorithm updates  $X^{(1)}, X^{(2)}, \dots, X^{(N(n))}$  to  $N(n')$  explicit copies of independent samples  $Y^{(1)}, Y^{(2)}, \dots, Y^{(N(n'))} \in Q^{V'}$  for the new instance  $I'$ , and correspondingly updates the data represented by the data structure to  $Y_0^{(i)} \in Q^{V'}$  and  $\text{Exe-Log}^{(i)}(I', T') = \left\langle u_t^{(i)}, Y_t^{(i)}(u_t^{(i)}) \right\rangle_{t=1}^{T'}$  for  $1 \leq i \leq N(n')$  such that each Gibbs sampling chain  $(Y_t^{(i)})_{t=0}^{T'}$  on  $I'$  generating a new sample  $Y^{(i)} = Y_{T'}^{(i)}$ , where each  $Y_0^{(i)}$  and  $\text{Exe-Log}^{(i)}(I', T')$  satisfy Condition 6.2 for the Gibbs sampling on  $I'$ , therefore,

$$d_{\text{TV}}(Y^{(i)}, \mu_{I'}) \leq \epsilon(n').$$

- **(time cost)** Assuming Condition 6.2 for each  $X_0^{(i)}$  and  $\text{Exe-Log}^{(i)}(I, T)$  for the Gibbs sampling on  $I$ , the time complexity for resolving an update is

$$O\left(\Delta^2(L_{\text{Hamil}} + L_{\text{graph}})N(n) \cdot \log^3 n + \Delta n\right),$$

where  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ , and  $\Delta_G, \Delta_{G'}$  denote the maximum degree of  $G = (V, E)$  and  $G' = (V', E')$ .

The following technique lemma will be used to prove Lemma 6.18.

**Lemma 6.19.** Let  $N : \mathbb{N}^+ \rightarrow \mathbb{N}^+$  be a function such that there exists a constant  $C > 0$  such that

$$\forall n \in \mathbb{N}^+ : |N(n+1) - N(n)| \leq \frac{C}{n}N(n).$$

Then the function  $N$  has the following properties

- for any  $n \in \mathbb{N}^+$ , it holds that  $N(n) \leq \text{poly}(n)$ ;
- let  $\alpha \geq 1$  be a constant, given any  $n, n' \in \mathbb{N}^+$  such that  $\frac{1}{\alpha} \leq \frac{n'}{n} \leq \alpha$ ,

$$|N(n) - N(n')| = C'(\alpha, C) \cdot \frac{|n - n'|}{n}N(n),$$

where  $C'(\alpha, C)$  is a constant that depends only on  $\alpha$  and  $C$ .

*Proof.* By the condition, we have  $N(n+1) \leq (1 + \frac{C}{n})N(n)$ . Thus for all  $n \in \mathbb{N}^+$ ,

$$N(n) \leq N(1) \prod_{i=1}^{n-1} \left(1 + \frac{C}{i}\right) \leq N(1) \exp\left(C \sum_{i=1}^{n-1} \frac{1}{i}\right) = N(1) \exp(\Theta(\ln n)) = \text{poly}(n).$$

We then prove the second property. Note that  $\frac{|n-n'|}{n} \leq \alpha$ , it suffices to prove

$$(31) \quad \left| \frac{N(n')}{N(n)} - 1 \right| \leq C'(\alpha, C) \cdot \frac{|n - n'|}{n}.$$

Assume that  $\min\{n, n'\} \leq 2C\alpha$ . Then, we have  $\max\{n, n'\} \leq 2C\alpha^2$ . We can choose  $C'(\alpha, C)$  sufficiently large so that (31) holds. Assume  $n' > n > 2C\alpha$ . Note that  $\frac{|n-n'|}{n} \leq \alpha$ . We have

$$1 - \frac{C|n-n'|}{n} \leq \left(1 - \frac{C}{n}\right)^{|n-n'|} \leq \frac{N(n')}{N(n)} \leq \left(1 + \frac{C}{n}\right)^{|n-n'|} \leq 1 + \frac{C \exp(\alpha C) |n-n'|}{n},$$

which implies (31) holds if  $C'(\alpha, C) \geq C \exp(\alpha C)$ . Assume  $n > n' > 2C\alpha$ . Note that  $\frac{|n-n'|}{n} \leq \alpha$  and  $n' \geq \frac{n}{\alpha}$ . We have

$$1 - \frac{\alpha C |n-n'|}{n} \leq \left(1 - \frac{\alpha C}{n}\right)^{|n-n'|} \leq \frac{N(n')}{N(n)} \leq \left(1 + \frac{\alpha C}{n}\right)^{|n-n'|} \leq 1 + \frac{C\alpha \exp(\alpha^2 C) |n-n'|}{n}.$$

which implies (31) holds if  $C'(\alpha, C) \geq C\alpha \exp(\alpha^2 C)$ .  $\square$

*Proof.* The main idea of the multi-sample dynamic Gibbs sampling algorithm is to use single-sample dynamic Gibbs sampling algorithm (Algorithm 2) to maintain each sample  $X^{(i)} \in Q^V$  for  $1 \leq i \leq N(n)$ . We need a careful implementation of the algorithm to guarantee the time cost in Lemma 6.18.

**Space cost:** Note that  $T = \left\lceil \frac{n}{\delta} \log \frac{n}{\epsilon(n)} \right\rceil = \Theta(n \log n)$  due to Lemma 6.17 and  $N(n) \leq \text{poly}(n)$  due to Lemma 6.19. The dynamic dictionary for each sample  $X^{(i)}$  uses  $O(n)$  memory words, each of  $O(\log n)$  bits. Hence, the algorithm uses  $O(T \cdot N(n)) = O(nN(n) \log n)$  memory words to maintain all the initial states, execution-logs and the random samples due to Theorem 6.12.

**Correctness:** The invariants for execution-log (Condition 6.2) are preserved by the coupling simulated by the algorithm. The correctness holds as a consequence.

**Time cost:** Define  $N_{\min} \triangleq \min\{N(n), N(n')\}$ . Fix  $1 \leq k \leq N_{\min}$ . We use the Algorithm 2 to update the sample  $X^{(k)}$  to  $Y^{(k)}$ . Let  $\mathcal{P}_k \subseteq [T]$  denote the set defined in (13) for the subroutine UpdateHamiltonian in Algorithm 2. The multi-sample dynamic Gibbs sampling has the following three stages.

- **Preparation stage:** construct the updated instances  $\mathcal{I}'$  and other middle instances  $\mathcal{I}_{\text{mid}}, \mathcal{I}_1, \mathcal{I}_2$  in (8), (18), (19); compute  $p_v^{\text{up}}$  in (12) for all  $v \in V$ ; and construct the random sets  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_{N_{\min}}$ .
- **Update stage:** given the  $(p_v^{\text{up}})_{v \in V}$  and  $(\mathcal{P}_i)_{1 \leq i \leq N_{\min}}$ , for each  $1 \leq i \leq N_{\min}$ , use Algorithm 2 to update the initial state  $X_0^{(i)}$  to  $Y_0^{(i)}$ , the execution-log  $\text{Exe-Log}^{(i)}(\mathcal{I}, T) = \left\langle v_t^{(i)}, X_t^{(i)}(v_t^{(i)}) \right\rangle_{t=1}^T$  to  $\text{Exe-Log}^{(i)}(\mathcal{I}', T') = \left\langle u_t^{(i)}, Y_t^{(i)}(u_t^{(i)}) \right\rangle_{t=1}^{T'}$ , and the sample  $X^{(i)}$  to  $Y^{(i)}$ .
- **Completion stage:** If  $N(n') < N(n)$ , for each  $N(n') < i \leq N(n)$ , remove the sample  $X^{(i)}$ , the initial state  $X_0^{(i)}$  and the execution-log  $\text{Exe-Log}^{(i)}(\mathcal{I}, T) = \left\langle v_t^{(i)}, X_t^{(i)}(v_t^{(i)}) \right\rangle_{t=1}^T$  from the data; if  $N(n') > N(n)$ , for each  $N(n) < i \leq N(n')$ , construct an independent Gibbs sampling chain  $(Y_t^{(i)})_{t=0}^{T'}$  on instance  $\mathcal{I}'$ , write the sample  $Y^{(i)} = Y_{T'}^{(i)}$ , the initial state  $Y_0^{(i)}$  and the execution-log  $\text{Exe-Log}^{(i)}(\mathcal{I}', T') = \left\langle u_t^{(i)}, Y_t^{(i)}(u_t^{(i)}) \right\rangle_{t=1}^{T'}$  into the data.

Let  $T_{\text{preparation}}^{\text{multi}}$ ,  $T_{\text{update}}^{\text{multi}}$  and  $T_{\text{completion}}^{\text{multi}}$  denote the running time of the corresponding stages. Note that the update stage of the multi-sample dynamic sampling algorithm repeats the update stage of the single-sample algorithm for  $N_{\min}$  times. Also note that both  $\mathcal{I}$  and  $\mathcal{I}'$  satisfies Dobrushin-Shlosman condition. Combining (24), (27), (28), (29), and (30), we have

$$(32) \quad \begin{aligned} \mathbb{E} \left[ T_{\text{update}}^{\text{multi}} \right] &= \sum_{i=1}^{N_{\min}} \mathbb{E} \left[ T_{\text{update}}^{\text{single},(i)} \right] = O(N_{\min} \Delta^2 (L_{\text{graph}} + L_{\text{Hamil}}) \log^3 n) \\ &\text{(by } N_{\min} \leq N(n)) \quad = O(N(n) \Delta^2 (L_{\text{graph}} + L_{\text{Hamil}}) \log^3 n) \end{aligned}$$

where  $T_{\text{update}}^{\text{single},(i)}$  is the running time of the update stage of the Algorithm 2 that updates the  $i$ -th sample.

In completion stage, we either remove the chains from the data structure, or generate the new chains and write them into data structure. It is easy to see the running time of the completion stage satisfies

$$\begin{aligned} \mathbb{E} \left[ T_{\text{completion}}^{\text{multi}} \right] &= O(|N(n) - N(n')| T_{\text{max}} \log T_{\text{max}}) = O(n |N(n) - N(n')| \log^2 n) \\ &\text{(by Lemma 6.19)} \quad = O(|n - n'| N(n) \log^2 n) = O(L_{\text{graph}} N(n) \log^2 n), \end{aligned}$$

where  $T_{\text{max}} = \max\{T, T'\} = O(n \log n)$  since  $n' = \Theta(n)$  and  $\epsilon(n') \geq \frac{1}{\text{poly}(n')}$  (by  $L_{\text{graph}} = o(n)$  and Lemma 6.17).

We make the following claim about the preparation stage.

**Claim 6.20.** *The expected running time of the preparation stage is*

$$\mathbb{E} \left[ T_{\text{preparation}}^{\text{multi}} \right] = O \left( \Delta n + \log^2 n \sum_{i=1}^{N_{\min}} \mathbb{E} [|\mathcal{P}_i|] \right),$$

and the expected size of  $\mathcal{P}_i$  is at most  $\frac{4T_{\text{max}}L_{\text{Hamil}}}{n}$  for each  $1 \leq i \leq N_{\min}$ .

By Claim 6.20, we have

$$\mathbb{E} \left[ T_{\text{preparation}}^{\text{multi}} \right] = O \left( \Delta n + N(n) L_{\text{Hamil}} \log^3 n \right).$$

By the linearity of expectation, the expected time cost of the algorithm is  $\mathbb{E} \left[ T_{\text{preparation}}^{\text{multi}} \right] + \mathbb{E} \left[ T_{\text{update}}^{\text{multi}} \right] + \mathbb{E} \left[ T_{\text{completion}}^{\text{multi}} \right]$ . This proves the time cost.  $\square$

## 7. PROOFS FOR DYNAMIC GIBBS SAMPLING

**7.1. Analysis of the couplings.** We analysis the couplings in dynamic Gibbs sampling algorithm. In Section 7.1.1, we analysis the coupling for Hamiltonian update. In Section 7.1.2, we analysis the coupling for graph update.

**7.1.1. Proofs for the coupling for Hamiltonian update.** In this section, we prove Lemma 6.5, Lemma 6.6, and Lemma 6.8.

**The validity of the coupling (proof of Lemma 6.5).** We first prove that the distribution  $v_{I_v, I'_v}^\tau(\cdot)$  in (10) is valid. We draw samples from  $v_{I_v, I'_v}^\tau(\cdot)$  only if the result of coin flipping is HEADS, which implies  $\mu_{v, I}(x | \tau) > \mu_{v, I'}(x | \tau)$  for some  $x \in Q$ . Thus, the two distributions  $\mu_{v, I}(\cdot | \tau)$  and  $\mu_{v, I'}(\cdot | \tau)$  are not identical, and

$$\sum_{x \in Q} \max \{0, \mu_{v, I}(x | \tau) - \mu_{v, I'}(x | \tau)\} > 0.$$

Hence, the denominator of  $v_{I_v, I'_v}^\tau(\cdot)$  is positive. Besides, since both  $\mu_{v, I}(\cdot | \tau)$  and  $\mu_{v, I'}(\cdot | \tau)$  are distributions over  $Q$ , we have

$$\sum_{x \in Q} \max \{0, \mu_{v, I'}(x | \tau) - \mu_{v, I}(x | \tau)\} = \sum_{x \in Q} \max \{0, \mu_{v, I}(x | \tau) - \mu_{v, I'}(x | \tau)\}.$$

Thus we have  $\sum_{x \in Q} v_{I_v, I'_v}^\tau(x) = 1$ . Hence,  $v_{I_v, I'_v}^\tau(\cdot)$  a valid distribution.

We next prove the coupling  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  in Definition 6.4 is a valid coupling between  $\mu_{v, I}(\cdot | \tau)$  and  $\mu_{v, I'}(\cdot | \tau)$ . If  $\mu_{v, I}(\cdot | \tau)$  and  $\mu_{v, I'}(\cdot | \tau)$  are identical, the result holds trivially. We may assume  $\mu_{v, I}(\cdot | \tau)$  and  $\mu_{v, I'}(\cdot | \tau)$  are not identical, thus the distribution  $v_{I_v, I'_v}^\tau(\cdot)$  is well-defined.

The coupling  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  in Definition 6.4 returns a pair  $(c, c') \in Q^2$ . It is easy to see  $c$  follows the law  $\mu_{v, I}(\cdot | \sigma)$ . We prove that  $c'$  follows the law  $\mu_{v, I'}(\cdot | \sigma)$ . By the definition of  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$ ,  $c' \in Q$  is generated by the following procedure:

- sample  $a \in Q$  from the distribution  $\mu_{v, I}(\cdot | \tau)$ ;
- sample  $b \in Q$  from the distribution  $v_{I_v, I'_v}^\tau$  defined in (10), set

$$c' = \begin{cases} b & \text{with probability } p_{I_v, I'_v}^\tau(a) \\ a & \text{with probability } 1 - p_{I_v, I'_v}^\tau(a). \end{cases}$$

Note that  $a$  follows the law  $\mu_{v, I}(\cdot | \tau)$ . We have for each  $x \in Q$ ,

$$\begin{aligned} \Pr[c' = x] &= \Pr[a = x] \cdot (1 - p_{I_v, I'_v}^\tau(x)) + \sum_{y \in Q} \Pr[a = y] \cdot p_{I_v, I'_v}^\tau(y) \cdot v_{I_v, I'_v}^\tau(x) \\ &= \mu_{v, I}(x | \tau) \cdot (1 - p_{I_v, I'_v}^\tau(x)) + v_{I_v, I'_v}^\tau(x) \sum_{y \in Q} \mu_{v, I}(y | \tau) \cdot p_{I_v, I'_v}^\tau(y). \end{aligned}$$

By the definition of  $p_{I_v, I'_v}^\tau(y)$  in (9), we have

$$\forall y \in Q, \quad \mu_{v, I}(y | \tau) \cdot p_{I_v, I'_v}^\tau(y) = \begin{cases} 0 & \text{if } \mu_{v, I}(y | \tau) \leq \mu_{v, I'}(y | \tau) \\ \mu_{v, I}(y | \tau) - \mu_{v, I'}(y | \tau) & \text{otherwise.} \end{cases}$$



This implies  $\mu_{v,I}(y | \tau) \cdot p_{I_v, I'_v}^\tau(y) = \max \{0, \mu_{v,I}(y | \tau) - \mu_{v,I'}(y | \tau)\}$ . We have

$$\begin{aligned} & v_{I_v, I'_v}^\tau(x) \sum_{y \in Q} \mu_{v,I}(y | \tau) \cdot p_{I_v, I'_v}^\tau(y) \\ &= \frac{\max \{0, \mu_{v,I'}(x | \tau) - \mu_{v,I}(x | \tau)\}}{\sum_{y \in Q} \max \{0, \mu_{v,I}(y | \tau) - \mu_{v,I'}(y | \tau)\}} \sum_{y \in Q} \max \{0, \mu_{v,I}(y | \tau) - \mu_{v,I'}(y | \tau)\} \\ &= \max \{0, \mu_{v,I'}(x | \tau) - \mu_{v,I}(x | \tau)\}. \end{aligned}$$

Hence, we have

$$\Pr[c' = x] = \mu_{v,I}(x | \tau) \cdot (1 - p_{I_v, I'_v}^\tau(x)) + \max \{0, \mu_{v,I'}(x | \tau) - \mu_{v,I}(x | \tau)\}.$$

Suppose  $\mu_{v,I}(x | \tau) \leq \mu_{v,I'}(x | \tau)$ , then we have  $p_{I_v, I'_v}^\tau(x) = 0$ . In this case, we have

$$\Pr[c' = x] = \mu_{v,I}(x | \tau) + \mu_{v,I'}(x | \tau) - \mu_{v,I}(x | \tau) = \mu_{v,I'}(x | \tau).$$

Suppose  $\mu_{v,I}(x | \tau) > \mu_{v,I'}(x | \tau)$ , then we have

$$\Pr[c' = x] = \mu_{v,I}(x | \tau) \cdot (1 - p_{I_v, I'_v}^\tau(x)) = \mu_{v,I'}(x | \tau).$$

Combining these two cases proves that  $c'$  follows the law  $\mu_{v,I'}(\cdot | \tau)$ .  $\square$

**The upper bound of the probability  $p_{I_v, I'_v}^\tau(\cdot)$  (proof of Lemma 6.6).** It suffices to prove that for any two instances  $I = (V, E, Q, \Phi)$  and  $I' = (V, E, Q, \Phi')$  of MRF model, and any  $v \in V, c \in Q$  and  $\sigma \in Q^{\Gamma_G(v)}$ ,

$$(33) \quad \mu_{v,I}(c | \sigma) - \mu_{v,I'}(c | \sigma) \leq 2\mu_{v,I}(c | \sigma) \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right).$$

Note that if  $\mu_{v,I}(c | \sigma) = 0$ , then  $p_{I_v, I'_v}^\tau(c) = 0$ ; otherwise  $p_{I_v, I'_v}^\tau(c) = \max \left\{ 0, \frac{\mu_{v,I}(c | \sigma) - \mu_{v,I'}(c | \sigma)}{\mu_{v,I}(c | \sigma)} \right\}$ . Hence, inequality (33) proves the lemma.

We now prove (33). Suppose  $\mu_{v,I}(c | \sigma) = 0$ . Then the LHS of (33)  $\leq 0$ . Since the RHS  $\geq 0$ , the inequality holds.

We next assume  $\mu_{v,I}(c | \sigma) > 0$ . Then it suffices to prove

$$\frac{\mu_{v,I}(c | \sigma) - \mu_{v,I'}(c | \sigma)}{\mu_{v,I}(c | \sigma)} = 1 - \frac{\mu_{v,I'}(c | \sigma)}{\mu_{v,I}(c | \sigma)} \leq 2 \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right).$$

By the definitions of  $\phi_v, \phi'_v, \phi_e, \phi'_e$ , we can write the ratio as

$$\frac{\mu_{v,I'}(c | \sigma)}{\mu_{v,I}(c | \sigma)} = \frac{\exp(\phi'_v(c) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, c)) \sum_{a \in Q} \exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))}{\exp(\phi_v(c) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, c)) \sum_{a \in Q} \exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))},$$

where  $\Gamma_v$  denotes the neighborhood of  $v$  in  $G$ . Next, we assume that

$$(34) \quad \begin{aligned} \forall c \in Q: \quad \phi_v(c) = -\infty &\iff \phi'_v(c) = -\infty \\ \forall u \in \Gamma_v, c, c' \in Q: \quad \phi_{uv}(c, c') = -\infty &\iff \phi'_{uv}(c, c') = -\infty. \end{aligned}$$

Otherwise, it must hold that the RHS of (33) is  $\infty$ , then (33) holds trivially. Thus we can define the set

$$Q' \triangleq \left\{ a \in Q \mid \phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a) \neq -\infty \right\} = \left\{ a \in Q \mid \phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a) \neq -\infty \right\}.$$

Since  $\exp(-\infty) = 0$ , we have

$$\frac{\mu_{v,I'}(c | \sigma)}{\mu_{v,I}(c | \sigma)} = \frac{\exp(\phi'_v(c) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, c)) \sum_{a \in Q'} \exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))}{\exp(\phi_v(c) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, c)) \sum_{a \in Q'} \exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))}.$$

We then show that

$$(35) \quad \begin{aligned} \forall a \in Q' : \frac{\exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))}{\exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))} &\geq \exp\left(-\|\phi_v - \phi'_v\|_1 - \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1\right) \\ \forall a \in Q' : \frac{\exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))}{\exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))} &\geq \exp\left(-\|\phi_v - \phi'_v\|_1 - \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1\right) \end{aligned}$$

We first use (35) to prove the (33). Since  $\mu_{v,I}(c \mid \sigma) > 0$ , then we have  $c \in Q'$ . By (35), we have

$$\begin{aligned} 1 - \frac{\mu_{v,I'}(c \mid \sigma)}{\mu_{v,I}(c \mid \sigma)} &\leq 1 - \exp\left(-2\|\phi_v - \phi'_v\|_1 - 2 \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1\right) \\ &\leq 2 \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right). \end{aligned}$$

This proves the lemma.

We now prove (35). For any  $a \in Q'$ , it holds that

$$\frac{\exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))}{\exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))} = \exp\left(\phi_v(a) - \phi'_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a) - \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a)\right).$$

Then (35) holds because

$$\begin{aligned} \phi_v(a) - \phi'_v(a) &\geq - \sum_{c \in Q} |\phi_v(c) - \phi'_v(c)| = -\|\phi_v - \phi'_v\|_1; \\ \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a) - \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a) &\geq - \sum_{e=\{u,v\} \in E} \sum_{c, c' \in Q} |\phi_e(c, c') - \phi'_e(c, c')| = - \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1. \end{aligned}$$

The lower bound of  $\frac{\exp(\phi'_v(a) + \sum_{u \in \Gamma_v} \phi'_{uv}(\sigma_u, a))}{\exp(\phi_v(a) + \sum_{u \in \Gamma_v} \phi_{uv}(\sigma_u, a))}$  can be proved in a similar way.  $\square$

**The cost of the coupling for UpdateHamiltonian (proof of Lemma 6.8).** By the definition of the indicator random variable  $\gamma_t$  in (17), we have

$$\begin{aligned} \Pr[\gamma_t = 1 \mid \mathcal{D}_{t-1}] &\leq \Pr[t \in \mathcal{P} \mid \mathcal{D}_{t-1}] + \Pr[v_t \in \Gamma_G^+(\mathcal{D}_{t-1}) \mid \mathcal{D}_{t-1}] \\ &\leq \frac{(\Delta + 1)|\mathcal{D}_{t-1}|}{n} + \sum_{v \in V} \frac{p_v^{\text{up}}}{n}. \end{aligned}$$

By the definition of  $p_v^{\text{up}}$  in (12) and  $d_{\text{Hamil}}(I, I') = \sum_{v \in V} \|\phi_v - \phi'_v\|_1 + \sum_{e \in E} \|\phi_e - \phi'_e\|_1 \leq L$ , we have

$$\Pr[\gamma_t = 1 \mid \mathcal{D}_{t-1}] \leq \frac{(\Delta + 1)|\mathcal{D}_{t-1}|}{n} + \frac{4L}{n}.$$

By the definition of  $R_{\text{Hamil}} \triangleq \sum_{t=1}^T \gamma_t$ , we have

$$(36) \quad \mathbb{E}[R_{\text{Hamil}}] = \sum_{t=1}^T \mathbb{E}[\gamma_t] = \sum_{t=1}^T \mathbb{E}[\mathbb{E}[\gamma_t \mid \mathcal{D}_{t-1}]] \leq \sum_{t=1}^T \left( \frac{(\Delta + 1)\mathbb{E}[|\mathcal{D}_{t-1}|]}{n} + \frac{4L}{n} \right).$$

Next, we bound the expectation  $\mathbb{E}[|\mathcal{D}_t|]$ . Recall that the one-step local coupling for Hamiltonian update (Definition 6.3) is implemented as follows. We first construct the random set  $\mathcal{P} \subseteq V$  in (13). In the  $t$ -th step, where  $1 \leq t \leq T$ , given any  $X_{t-1}$  and  $Y_{t-1}$ , the  $X_t$  and  $Y_t$  is generated as follows.

- Let  $X'(u) = X_{t-1}(u)$  and  $Y'(u) = Y_{t-1}(u)$  for all  $u \in V \setminus \{v_t\}$ , sample  $(X'(v_t), Y'(v_t)) \in Q^2$  jointly from the optimal coupling  $D_{\text{opt}, I_{v_t}}^{\sigma, \tau}$  of the marginal distributions  $\mu_{v_t, I}(\cdot \mid \sigma)$  and  $\mu_{v_t, I}(\cdot \mid \tau)$ , where  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_G(v_t))$ .
- Let  $X_t = X'$  and  $Y_t = Y'$ . If  $t \in \mathcal{P}$ , update the value of  $Y_t(v_t)$  using (14).

Hence, for any vertex  $v \in V$ ,  $X_t(v) \neq Y_t(v)$  only if one of the following two events occurs (1)  $X'(v) \neq Y'(v)$ ; (2)  $v_t = v$  and  $t \in \mathcal{P}$ . Then for any  $v \in V$ , we have

$$(37) \quad \begin{aligned} \Pr[X_t(v) \neq Y_t(v) \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] &\leq \Pr[X'(v) \neq Y'(v) \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] + \Pr[v = v_t \wedge t \in \mathcal{P} \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] \\ &= \Pr[X'(v) \neq Y'(v) \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] + \Pr[v = v_t \wedge t \in \mathcal{P}], \end{aligned}$$

where the equation holds because  $v = v_t \wedge t \in \mathcal{P}$  is independent of  $\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}$ . Given  $\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}$ , the random pair  $X', Y'$  are obtained by the one-step optimal coupling for Gibbs sampling on instance  $\mathcal{I}$  (Definition 4.2). Since  $\mathcal{I}$  satisfies the Dobrushin-Shlosman condition with constant  $0 < \delta < 1$ , by Proposition 4.3, we have

$$(38) \quad \mathbb{E}[H(X', Y') \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] \leq \left(1 - \frac{\delta}{n}\right) H(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) = \left(1 - \frac{\delta}{n}\right) |\mathcal{D}_{t-1}|.$$

where  $H(X, Y) = |\{v \in V \mid X(v) \neq Y(v)\}|$  denote the Hamming distance. Combining (37) and (38),

$$\begin{aligned} \mathbb{E}[|\mathcal{D}_t| \mid \mathcal{D}_{t-1}] &\leq \sum_{v \in V} \Pr[X'(v) \neq Y'(v) \mid \mathcal{D}_{t-1}] + \sum_{v \in V} \Pr[t \in \mathcal{P} \wedge v = v_t \mid \mathcal{D}_{t-1}] \\ &\leq \left(1 - \frac{\delta}{n}\right) |\mathcal{D}_{t-1}| + \sum_{v \in V} \frac{p_v^{\text{up}}}{n} \\ \text{(by (12))} \quad &\leq \left(1 - \frac{\delta}{n}\right) |\mathcal{D}_{t-1}| + \frac{2}{n} \sum_{v \in V} \left( \|\phi_v - \phi'_v\|_1 + \sum_{e=\{u,v\} \in E} \|\phi_e - \phi'_e\|_1 \right) \\ \text{(by } d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L) \quad &\leq \left(1 - \frac{\delta}{n}\right) |\mathcal{D}_{t-1}| + \frac{4L}{n}. \end{aligned}$$

Thus, we have

$$\mathbb{E}[|\mathcal{D}_t|] \leq \left(1 - \frac{\delta}{n}\right) \mathbb{E}[|\mathcal{D}_{t-1}|] + \frac{4L}{n}.$$

Note that  $|\mathcal{D}_0| = 0$ . This implies

$$(39) \quad \mathbb{E}[|\mathcal{D}_t|] \leq \frac{8L}{\delta}.$$

Thus, by (36), we have

$$\mathbb{E}[R_{\text{Hamil}}] \leq \frac{20\Delta TL}{\delta n} = O\left(\frac{\Delta TL}{\delta n}\right).$$

□

7.1.2. *Proofs for the coupling for graph update.* In this section, we prove Lemma 6.11.

**Cost of the coupling for UpdateEdge (Proof of Lemma 6.11).** By the definition of  $R_{\text{graph}}$  in (23) and the linearity of the expectation, we have

$$\mathbb{E}[R_{\text{graph}}] = \sum_{t=1}^T \mathbb{E}[\gamma_t] = \sum_{t=1}^T \mathbb{E}[\mathbb{E}[\gamma_t \mid \mathcal{D}_{t-1}]].$$

Recall  $\gamma_t = 1[v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D}_{t-1})]$  and  $v_t \in V$  is uniformly at random given  $\mathcal{D}_{t-1}$ . Note that  $|\Gamma_G^+(\mathcal{D}_{t-1})| \leq (\Delta + 1)|\mathcal{D}_{t-1}|$  and  $|\mathcal{S}| \leq 2|E \oplus E'| \leq 2L$ . We have

$$(40) \quad \mathbb{E}[R_{\text{graph}}] \leq \sum_{t=1}^T \mathbb{E}\left[\frac{(\Delta + 1)|\mathcal{D}_{t-1}| + 2L}{n}\right] = \frac{(\Delta + 1)}{n} \sum_{t=1}^T \mathbb{E}[|\mathcal{D}_{t-1}|] + \frac{2LT}{n}.$$

Suppose  $\mathcal{I}'$  satisfies Dobrushin-Shlosman condition (Condition 3.1) with the constant  $\delta > 0$ , we claim

$$(41) \quad \forall 0 \leq t \leq T : \quad \mathbb{E}[|\mathcal{D}_t|] \leq \frac{8L}{\delta}.$$

Combining (40) and (41), we have

$$\mathbb{E} [R_{\text{graph}}] \leq \frac{18\Delta L T}{\delta n} = O\left(\frac{\Delta L T}{n}\right).$$

This proves the lemma.

We now prove (41). Let  $(X_t, Y_t)_{t \geq 0}$  be the one-step local coupling for updating edges (Definition 6.9). We claim the following result

$$(42) \quad \forall \sigma, \tau \in \Omega : \quad \mathbb{E} [H(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \leq \left(1 - \frac{\delta}{n}\right) \cdot H(\sigma, \tau) + \frac{4L}{n},$$

where  $H(\sigma, \tau) = |\{v \in V \mid \sigma(v) \neq \tau(v)\}|$  denotes the Hamming distance. Assume (42) holds. Taking expectation over  $X_{t-1}$  and  $Y_{t-1}$ , we have

$$(43) \quad \mathbb{E} [H(X_t, Y_t)] \leq \left(1 - \frac{\delta}{n}\right) \mathbb{E} [H(X_{t-1}, Y_{t-1})] + \frac{4L}{n}.$$

Note that  $X_0 = Y_0$ , we have

$$(44) \quad H(X_0, Y_0) = 0.$$

Combining (43) with (44) implies

$$(45) \quad \forall 0 \leq t \leq T : \quad \mathbb{E} [|\mathcal{D}_t|] = \mathbb{E} [H(X_t, Y_t)] \leq \frac{8L}{\delta}.$$

This proves the claim in (41).

We finish the proof by proving the claim in (42). The main idea is to compare the one-step local coupling for updating edges (Definition 6.9) with the one-step optimal coupling for Gibbs sampling on instance  $\mathcal{I}'$  (Definition 4.2). Let  $(X'_t, Y'_t)_{t \geq 0}$  be the coupling for Gibbs sampling on  $\mathcal{I}'$ . Since  $\mathcal{I}'$  satisfies Dobrushin-Shlosman condition, by Proposition 4.3, we have

$$(46) \quad \forall \sigma, \tau \in \Omega = Q^V : \quad \mathbb{E} [H(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] \leq \left(1 - \frac{\delta}{n}\right) \cdot H(\sigma, \tau).$$

According to the coupling, we can rewrite the expectation in (46) as follows:

$$(47) \quad \mathbb{E} [H(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] = \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ H \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right],$$

where  $(C_v^X, C_v^Y) \sim D_{\text{opt}, I'_v}^{\sigma, \tau}$ ,  $D_{\text{opt}, I'_v}^{\sigma, \tau}$  is the optimal coupling between  $\mu_{v, \mathcal{I}'}(\cdot \mid \sigma)$  and  $\mu_{v, \mathcal{I}'}(\cdot \mid \tau)$ , and the configuration  $\sigma^{v \leftarrow C_v^{X'}} \in Q^V$  is defined as

$$\sigma^{v \leftarrow C_v^{X'}}(u) \triangleq \begin{cases} C_v^{X'} & \text{if } u = v \\ \sigma(u) & \text{if } u \neq v \end{cases}$$

and the configuration  $\tau^{v \leftarrow C_v^{Y'}} \in Q^V$  is defined in a similar way.

Similarly, we can rewrite the expectation in (42) as follows:

$$(48) \quad \mathbb{E} [H(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] = \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ H \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right],$$

where  $(C_v^X, C_v^Y) \sim D_{I_v, I'_v}^{\sigma, \tau}$ , where  $D_{I_v, I'_v}^{\sigma, \tau}$  is the local coupling defined in (21).

The following two properties hold for (47) and (48).

- If  $v \notin \mathcal{S}$ , by the definition of  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  in (21), it holds that  $D_{I_v, I'_v}^{\sigma, \tau} = D_{\text{opt}, I'_v}^{\sigma, \tau}$ . Hence

$$\forall v \notin \mathcal{S} : \quad \mathbb{E} \left[ H \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right] = \mathbb{E} \left[ H \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right].$$

- If  $v \in \mathcal{S}$ , then it holds that  $H(\sigma^{v \leftarrow C_v^X}, \sigma^{v \leftarrow C_v^{X'}}) \leq 1$  and  $H(\tau^{v \leftarrow C_v^{Y'}}, \tau^{v \leftarrow C_v^Y}) \leq 1$ . By the triangle inequality of the Hamming distance, we have

$$\begin{aligned} H(\sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y}) &\leq H(\sigma^{v \leftarrow C_v^X}, \sigma^{v \leftarrow C_v^{X'}}) + H(\sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}}) + H(\tau^{v \leftarrow C_v^{Y'}}, \tau^{v \leftarrow C_v^Y}) \\ &\leq H(\sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}}) + 2. \end{aligned}$$

This implies

$$\forall v \in \mathcal{S} : \quad \mathbb{E} \left[ H(\sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y}) \right] \leq \mathbb{E} \left[ H(\sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}}) \right] + 2.$$

Combining above two properties with (47) and (48), we have for any  $\sigma, \tau \in \Omega$ ,

$$\begin{aligned} &\mathbb{E} [H(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \\ &= \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ H(\sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y}) \right] \\ &\leq \frac{1}{n} \sum_{v \notin \mathcal{S}} \mathbb{E} \left[ H(\sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}}) \right] + \frac{1}{n} \sum_{v \in \mathcal{S}} \left( \mathbb{E} \left[ H(\sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}}) \right] + 2 \right) \\ (*) &\leq \mathbb{E} [H(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] + \frac{4L}{n} \\ &\leq \left( 1 - \frac{\delta}{n} \right) \cdot H(\sigma, \tau) + \frac{4L}{n}, \end{aligned}$$

where (\*) holds due to  $|\mathcal{S}| \leq 2L$ . This proves the claim in (42).  $\square$

**7.2. Implementation of the algorithms.** In this section, we prove the Claim 6.15, Claim 6.16 and Claim 6.20 by giving the implementation of the algorithms.

**7.2.1. Proofs of Claim 6.15 and Claim 6.20.** We prove Claim 6.20, then Claim 6.15 can be proved in a similar way.

It is easy to verify the updated sample  $\mathcal{I}'$ , all the probabilities  $(p_v^{\text{up}})_{v \in V}$  in (12), all middle instances  $\mathcal{I}_{\text{mid}}, \mathcal{I}_1, \mathcal{I}_2$  in (8), (18), (19) can be computed with time cost  $O(\Delta n)$ . We focus on constructing  $\mathcal{P}_i$  for  $1 \leq i \leq N_{\min}$ .

The multi-sample dynamic Gibbs sampling algorithm use the data structure in Theorem 6.12 to maintain  $N(n)$  independent Gibbs sampling chain on instance  $\mathcal{I}$  represented by  $X_0^{(i)}$  and  $\text{Exe-Log}(\mathcal{I}, T) = \left\langle v_t^{(i)}, X_t^{(i)}(v_t^{(i)}) \right\rangle_{t=1}^T$ . To construct the random sets  $\mathcal{P}_i$  for  $1 \leq i \leq N_{\min}$ , we need an additional data structure to maintain the following data. Define the set  $H_v$  as

$$H_v \triangleq \{(i, t) \in [N(n)] \times [T] \mid v_t^{(i)} = v\}.$$

$H_v$  contains all the transition steps in  $N(n)$  independent chains that picks the vertex  $v$ . The algorithm uses an extra data structure  $\mathcal{H}$  to maintain all  $(H_v)_{v \in V}$ . The data structure  $\mathcal{H}$  contains  $n$  balanced binary search trees  $(\mathcal{H}_v)_{v \in V}$ , where each  $\mathcal{H}_v$  maintains the set  $H_v$  in a similar way as in the main data structure in Theorem 6.12. Since  $T = O(n \log n)$ ,  $N(n) \leq \text{poly}(n)$ , the space cost of  $\mathcal{H}$  is  $O(nN(n) \log n)$  memory words, each of  $O(\log n)$  bits, which is dominated by the space cost in Lemma 6.18. And the time cost of adding, deleting, and searching a transition step in  $\mathcal{H}$  is  $O(\log^2 n)$ . We need to update  $\mathcal{H}$  when  $\mathcal{I}$  is updated to  $\mathcal{I}'$ . One can verify that such time cost is dominated by the time cost in Lemma 6.18.

Then for each  $v \in V$ , we pick each element in  $H_v$  with probability  $p_v^{\text{up}}$  to construct the set

$$\mathcal{B}_v \subseteq H_v.$$

This is the standard Bernoulli process. With the data structure  $\mathcal{H}_v$ , the time complexity of constructing the set  $\mathcal{B}_v$  is  $O(|\mathcal{B}_v| \log^2 n)$ . Given all the sets  $\mathcal{B}_v$ , it is easy to construct all the sets  $\mathcal{P}_i$ . Hence,

$$T_{\text{preparation}}^{\text{multi}} = O \left( \Delta n + \sum_{v \in V} |\mathcal{B}_v| \log^2 n \right) = O \left( \Delta n + \sum_{i=1}^{N_{\min}} |\mathcal{P}_i| \log^2 n \right).$$

In the preparation stage of multi-sample dynamic Gibbs sampling algorithm, we first construct the  $\mathcal{I}_{\text{mid}} = (V, E, Q, \Phi^{\text{mid}})$  as in (8), and each  $\mathcal{P}_i$  ( $1 \leq i \leq N_{\text{min}}$ ) is constructed with respect to  $\mathcal{I}$  and  $\mathcal{I}_{\text{mid}}$ . Note that  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}_{\text{mid}}) \leq d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}')$ . By (12), we have for each  $1 \leq i \leq N_{\text{min}}$ ,

$$\mathbb{E}[|\mathcal{P}_i|] \leq \sum_{t=1}^T \sum_{v \in V} \frac{p_v^{\text{up}}}{n} \leq \frac{4TL_{\text{Hamil}}}{n}.$$

This proves the claim.  $\square$

7.2.2. *Proof of Claim 6.16.* We give the implementation of the update stage of the single-sample dynamic Gibbs sampling algorithm (Algorithm 2). The algorithm updates the MRF instance from  $\mathcal{I}$  to  $\mathcal{I}'$  as follows,

$$\mathcal{I} \rightarrow \mathcal{I}_{\text{mid}} \rightarrow \mathcal{I}_1 \rightarrow \mathcal{I}_2 \rightarrow \mathcal{I}',$$

where  $\mathcal{I}_{\text{mid}}$  is defined in (8),  $\mathcal{I}_1 = \mathcal{I}_1(\mathcal{I}_{\text{mid}}, \mathcal{I}')$  is defined in (18), and  $\mathcal{I}_2 = \mathcal{I}_2(\mathcal{I}_{\text{mid}}, \mathcal{I}')$  is defined in (19). Then the algorithm calls LengthFix to modify the length of the execution log from  $T$  to  $T'$ .

The preparation stage computes all probabilities  $(p_v^{\text{up}})_{v \in V}$  in (12), the set  $\mathcal{P}$  in (13), and all instances  $\mathcal{I}_{\text{mid}}, \mathcal{I}_1, \mathcal{I}_2$ . Consider the time cost of the update stage. In the update from  $\mathcal{I}_{\text{mid}}$  to  $\mathcal{I}_1$ , we only add isolated vertices in  $V' \setminus V$ , using the data structure in Theorem 6.12, the expected time cost is

$$\mathbb{E}[T_{\mathcal{I}_{\text{mid}} \rightarrow \mathcal{I}_1}] = O\left(\frac{|V' \setminus V|}{|V|} T_{\text{max}} \log^2 T_{\text{max}}\right) = O\left(\frac{L_{\text{graph}}}{n} T_{\text{max}} \log^2 T_{\text{max}}\right).$$

In the update from  $\mathcal{I}_2$  to  $\mathcal{I}'$ , we only delete isolated vertices in  $V \setminus V'$ , thus

$$\mathbb{E}[T_{\mathcal{I}_2 \rightarrow \mathcal{I}'}] = O\left(\frac{|V \setminus V'|}{|V \cup V'|} T_{\text{max}} \log^2 T_{\text{max}}\right) = O\left(\frac{L_{\text{graph}}}{n} T_{\text{max}} \log^2 T_{\text{max}}\right).$$

It is also easy to observe that the expected time cost of LengthFix is

$$\mathbb{E}[T_{\text{LengthFix}}] = O\left(\Delta |T - T'| \log^2 T_{\text{max}}\right).$$

We then prove that

$$(49) \quad \mathbb{E}[T_{\mathcal{I} \rightarrow \mathcal{I}_{\text{mid}}}] = O\left(\Delta \mathbb{E}[R_{\text{Hamil}}] \log^2 T_{\text{max}}\right)$$

$$(50) \quad \mathbb{E}[T_{\mathcal{I}_1 \rightarrow \mathcal{I}_2}] = O\left(\Delta \mathbb{E}[R_{\text{graph}}] \log^2 T_{\text{max}}\right).$$

Combining all the running time together proves Claim 6.16.

We give the implementation of Algorithm 4 to prove (49). The Algorithm 6 can be implemented in a similar way to prove (50). Since  $(p_v^{\text{up}})_{v \in V}$  and  $\mathcal{P}$  are given, the running time of Algorithm 4 is dominated by the while-loop. We implement Algorithm 4 such that after each execution of the while-loop, the first  $t_0$  transition steps of the Gibbs sampling on instance  $\mathcal{I}$  is updated to the first  $t_0$  transition steps of the Gibbs sampling on instance  $\mathcal{I}'$ , namely,  $(X_t)_{t=0}^{t_0}$  is updated to  $(Y_t)_{t=0}^{t_0}$ , where  $t_0$  is the variable in Algorithm 4. Recall the sets  $\mathcal{D}$  and  $\mathcal{P}$  in Algorithm 4. We need some temporary data structures:

- a balanced binary search tree  $\mathcal{T}$  to maintain the set  $\mathcal{D}$  and the configuration  $X_{t_0-1}(\mathcal{D})$ ;
- a heap  $\mathcal{H}_1$  to maintain the set  $\mathcal{P}$ ;
- a heap  $\mathcal{H}_2$  such that once a vertex  $v$  is added into  $\mathcal{D}$ , the update times  $\text{Succ}(t_0, u)$  for all  $u \in \Gamma_G(v) \cup \{v\}$  are added into  $\mathcal{H}_2$ , where Succ is the operation of the data structure in Theorem 6.12.

Line 5 can be implemented using  $\mathcal{H}_1, \mathcal{H}_2, \mathcal{T}$ . And Line 7 and Line 10 can be implemented using  $\mathcal{T}$  and the main data structure in Theorem 6.12. Note that the time cost of each operation of  $\mathcal{T}$  is  $O(\log n) = O(\log T_{\text{max}})$ . Also note that at most  $\Delta R_{\text{Hamil}}$  elements can be added into  $\mathcal{H}_2$ . Hence, all the time cost contributed by  $\mathcal{H}_2$  is  $O(\Delta R_{\text{Hamil}} \log(\Delta R_{\text{Hamil}})) = O(\Delta R_{\text{Hamil}} \log T_{\text{max}})$ . One can verify that the total running time is

$$T_{\mathcal{I} \rightarrow \mathcal{I}_{\text{mid}}} = O\left(\Delta R_{\text{Hamil}} \log^2 T_{\text{max}}\right).$$

This proves (49).  $\square$



**7.3. Dynamic Gibbs sampling for specific models.** In this section, we apply our algorithm on Ising model, graph  $q$ -coloring, and hardcore model. We prove the following theorem.

**Theorem 7.1.** *There exist dynamic sampling algorithms as stated in Theorem 6.1 with the same space cost  $O(nN(n) \log n)$ , and expected time cost  $O\left(\Delta^2(L_{\text{graph}} + L_{\text{Hamil}})N(n) \log^3 n + \Delta n\right)$  for each update, if the input instance  $\mathcal{I}$  with  $n$  vertices and the updated instance  $\mathcal{I}'$  satisfying  $d_{\text{graph}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{graph}} = o(n)$ ,  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L_{\text{Hamil}}$  both are:*

- Ising models with temperature  $\beta$  and arbitrary local fields where  $\exp(-2|\beta|) \geq 1 - \frac{2-\delta}{\Delta+1}$ ;
- proper  $q$ -colorings with  $q \geq (2 + \delta)\Delta$ ;
- hardcore models with fugacity  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , but with an alternative time cost for each update

$$(51) \quad O\left(\Delta^3(L_{\text{graph}} + L_{\text{Hamil}})N(n) \log^3 n + \Delta n\right),$$

where  $\delta > 0$  is a constant,  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ ,  $\Delta_G$  denotes the maximum degree of the input graph, and  $\Delta_{G'}$  denotes the maximum degree of the updated graph.

In Theorem 7.1, the regime for Ising model and  $q$ -coloring match the Dobrushin-Shlosman condition, thus the results are corollaries of Theorem 6.1. The regime for hardcore model is better than the Dobrushin-Shlosman condition. We give the proof for hardcore model.

We use  $\mathcal{I} = (V, E, \lambda)$  to specify the hardcore model on graph  $G = (V, E)$  with fugacity  $\lambda$ . A configuration of hardcore model is  $\sigma \in \{0, 1\}^V$ , where  $\sigma_v = 1$  indicates  $v$  is occupied,  $\sigma_v = 0$  indicates  $v$  is unoccupied. If  $\sigma$  forms an independent set, then  $\mu_{\mathcal{I}}(\sigma) \propto \lambda^{|\sigma|}$ ; otherwise,  $\mu_{\mathcal{I}}(\sigma) = 0$ . We need the following lemma proved by Vigoda's coupling technique [Vig99].

**Lemma 7.2.** *Let  $\delta > 0$  be a constant. Let  $\mathcal{I} = (V, E, \lambda)$  be a hardcore instance, where  $n = |V|$ , and  $\Omega_{\mathcal{I}} \triangleq \{\sigma \in \{0, 1\}^V \mid \mu_{\mathcal{I}}(\sigma) > 0\}$ . Assume  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , where  $\Delta$  is the maximum degree of  $G = (V, E)$ . There exist a potential function  $\rho_{\mathcal{I}} : \Omega_{\mathcal{I}} \times \Omega_{\mathcal{I}} \rightarrow \mathbb{R}_{\geq 0}$ , where  $\forall \sigma, \tau \in \Omega_{\mathcal{I}}$ ,  $\rho_{\mathcal{I}}(\sigma, \tau) = 0$  if  $\sigma = \tau$  and  $\rho_{\mathcal{I}}(\sigma, \tau) \geq 1$  if  $\sigma \neq \tau$ , and  $\text{Diam}_{\mathcal{I}} \triangleq \max_{\sigma, \tau \in \Omega_{\mathcal{I}}} \rho_{\mathcal{I}}(\sigma, \tau) \leq \Delta n$ , such that the one-step optimal coupling (Definition 4.2)  $(X_t, Y_t)_{t \geq 0}$  of Gibbs sampling on  $\mathcal{I}$  satisfies*

(1) **(step-wise decay)** for the coupling  $(X_t, Y_t)_{t \geq 0}$  of Gibbs sampling, it holds that

$$(52) \quad \forall \sigma, \tau \in \Omega_{\mathcal{I}} : \quad \mathbb{E}[\rho_{\mathcal{I}}(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \leq \left(1 - \frac{\beta}{n}\right) \cdot \rho_{\mathcal{I}}(\sigma, \tau),$$

where  $\beta = \frac{1}{96\delta}$ , which implies  $\tau_{\text{mix}}(\mathcal{I}, \epsilon) \leq \lceil \frac{n}{\beta} \log \frac{\text{Diam}_{\mathcal{I}}}{\epsilon} \rceil = O(n \log \frac{n}{\epsilon})$ .

(2) **(up-bound to Hamming)** for all  $\sigma, \tau \in \Omega_{\mathcal{I}}$ ,  $H(\sigma, \tau) \leq \rho_{\mathcal{I}}(\sigma, \tau)$ , where  $H(\sigma, \tau)$  denotes the Hamming distance between  $\sigma$  and  $\tau$ .

(3) **(Lipschitz)** function  $\rho_{\mathcal{I}}(\cdot, \cdot)$ , seen as a function of  $2n$  variables, is  $K$ -Lipschitz, that is,

$$\max_{\sigma, \sigma', \tau, \tau' \in \Omega_{\mathcal{I}}} |\rho_{\mathcal{I}}(\sigma, \tau) - \rho_{\mathcal{I}}(\sigma', \tau')| \leq K \cdot H(\sigma\tau, \sigma'\tau'),$$

where  $K = 12\Delta$ .

Compared with Proposition 4.3, the step-wise decay property in (52) holds only for feasible configurations  $\sigma$  and  $\tau$ , and the decay property is established on the potential function  $\rho_{\mathcal{I}}$  rather than the Hamming distance  $H$ . We first use Lemma 7.2 to prove Theorem 7.1, then we prove Lemma 7.2 in the end of this section.

Recall that the error function  $\epsilon$  satisfies  $\epsilon(\ell) \geq \frac{1}{\text{poly}(\ell)}$  by Lemma 6.17. Recall  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ . By Lemma 7.2 and  $n' = \Theta(n)$  (since  $L_{\text{graph}} = o(n)$ ), we can set

$$T = T(\mathcal{I}) = \left\lceil \frac{96n}{\delta} \log \frac{n\Delta}{\epsilon(n)} \right\rceil = O(n \log n)$$

$$T' = T(\mathcal{I}') = \left\lceil \frac{96n'}{\delta} \log \frac{n'\Delta}{\epsilon(n')} \right\rceil = O(n \log n).$$

We modify Algorithm 2 for the hardcore model as follows. Suppose the current instance is  $\mathcal{I} = (V, E, \lambda)$ , we set the initial configuration  $X_0$  as

$$\forall v \in V, \quad X_0(v) = 0.$$

Thus  $X_0$  is feasible. Suppose the instance  $\mathcal{I} = (V, E, \lambda)$  is updated to  $\mathcal{I}' = (V', E', \lambda')$ . We divide the update into the following steps

$$\mathcal{I} \rightarrow \mathcal{I}_{\text{mid}} \rightarrow \mathcal{I}_1 \rightarrow \mathcal{I}_2 \rightarrow \mathcal{I}_3 \rightarrow \mathcal{I}'$$

- change fugacity to update  $\mathcal{I} = (V, E, \lambda)$  to  $\mathcal{I}_{\text{mid}} = (V, E, \lambda')$  using UpdateHamiltonian;
- add isolated vertices in  $V' \setminus V$  to update  $\mathcal{I}_{\text{mid}} = (V, E, \lambda')$  to  $\mathcal{I}_1 = (V \cup V', E, \lambda')$  using AddVertex;
- delete edges in  $E \setminus E'$  to update  $\mathcal{I}_1 = (V \cup V', E, \lambda')$  to  $\mathcal{I}_2 = (V \cup V', E \cap E', \lambda')$  using UpdateEdge;
- add edges in  $E' \setminus E$  to update  $\mathcal{I}_2 = (V \cup V', E \cap E', \lambda')$  to  $\mathcal{I}_3 = (V \cup V', E', \lambda')$  using UpdateEdge;
- delete isolated vertices in  $V' \setminus V$  to update  $\mathcal{I}_3 = (V \cup V', E', \lambda')$  to  $\mathcal{I}' = (V', E', \lambda')$ ;
- fix the length of the execution log from  $T$  to  $T'$ .

Compared to Algorithm 2, we further divide the update of edges into two steps: at first delete edges, then add edges. Thus, we have the following observation.

**Observation 7.3.** *The following results holds:*

- $\Omega_{\mathcal{I}} = \Omega_{\mathcal{I}_{\text{mid}}}$ ,  $\Omega_{\mathcal{I}_1} \subseteq \Omega_{\mathcal{I}_2}$  and  $\Omega_{\mathcal{I}_3} \subseteq \Omega_{\mathcal{I}_2}$ , where  $\Omega_{\mathcal{I}}$  is the set of feasible configurations for any instance  $\mathcal{I}$ .
- the instances  $\mathcal{I}, \mathcal{I}_2, \mathcal{I}_3, \mathcal{I}'$  all satisfy  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , where  $\lambda$  and  $\Delta$  are the fugacity and maximum degree of the corresponding instance.

By the observation, we know that  $\Omega_{\mathcal{I}} = \Omega_{\mathcal{I}_{\text{mid}}}$ ,  $\Omega_{\mathcal{I}_1} \subseteq \Omega_{\mathcal{I}_2}$  and  $\Omega_{\mathcal{I}_3} \subseteq \Omega_{\mathcal{I}_2}$ , thus we can use Lemma 7.2, because the step-wise decay property (52) is established only on feasible configurations.

We need to analyze  $R_{\text{Hamil}}$  and  $R_{\text{graph}}$  defined in (17) and (23) for the hardcore model. We prove the following two lemmas for hardcore model.

**Lemma 7.4.** *Consider UpdateHamiltonian  $(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$ . Let  $\mathcal{I} = (V, E, \lambda)$  be the current instance and  $\mathcal{I}' = (V, E, \lambda')$  the updated instance. Assume  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , where  $\delta > 0$  is a constant and  $\Delta$  is the maximum degree of  $G = (V, E)$ . Also assume  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') = n |\ln \lambda - \ln \lambda'| \leq L$ . Then  $\mathbb{E}[R_{\text{Hamil}}] = O\left(\frac{\Delta^2 TL}{n\delta}\right)$ , where  $n = V$ ,  $\Delta$  is the maximum degree of graph  $G = (V, E)$ .*

**Lemma 7.5.** *Consider UpdateEdge  $(\mathcal{I}, \mathcal{I}', X_0, \langle v_t, X_t(v_t) \rangle_{t=1}^T)$ . Let  $\mathcal{I} = (V, E, \lambda)$  be the current instance and  $\mathcal{I}' = (V, E', \lambda)$  the updated instance. Assume  $|E \oplus E'| \leq L$ . Also assume one of the following two conditions holds for some constant  $\delta > 0$ :*

- $\lambda \leq \frac{2-\delta}{\Delta_G-2}$  and  $\Omega_{\mathcal{I}'} \subseteq \Omega_{\mathcal{I}}$ , where  $\Delta_G$  is the maximum degree of  $G = (V, E)$ ;
- $\lambda \leq \frac{2-\delta}{\Delta_{G'}-2}$  and  $\Omega_{\mathcal{I}} \subseteq \Omega_{\mathcal{I}'}$ , where  $\Delta_{G'}$  is the maximum degree of  $G' = (V, E')$ .

Then  $\mathbb{E}[R_{\text{graph}}] = O\left(\frac{\Delta^2 TL}{n\delta}\right)$ , where  $n = V$ ,  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ .

Note that we call the subroutine UpdateHamiltonian for the update modifying  $\mathcal{I}$  to  $\mathcal{I}_{\text{mid}}$ . By Observation 7.3, the condition in Lemma 7.4 holds. We call the subroutine UpdateEdge for the update modifying  $\mathcal{I}_1$  to  $\mathcal{I}_2$  and the update modifying  $\mathcal{I}_2$  to  $\mathcal{I}_3$ . By Observation 7.3, in both two calls of UpdateEdge, the condition in Lemma 7.5 holds. Then Theorem 7.1 for hardcore can be proved by going through the proof in Section 6. Compared to Lemma 6.8 and Lemma 6.11,  $\mathbb{E}[R_{\text{Hamil}}], \mathbb{E}[R_{\text{graph}}]$  in Lemma 7.4 and Lemma 7.5 are bounded by  $O\left(\frac{\Delta^2 TL}{n\delta}\right)$  rather than  $O\left(\frac{\Delta TL}{n\delta}\right)$ . This is why the hardcore model has an alternative running time in (51).

The proofs of Lemma 7.4 and Lemma 7.5 are similar to the proofs of Lemma 6.8 and Lemma 6.11. We give the proofs here for the completeness.

*Proof of Lemma 7.4.* By the definition of the indicator  $\gamma_t$  in (17), we have

$$\Pr[\gamma_t = 1 \mid \mathcal{D}_{t-1}] \leq \Pr[t \in \mathcal{P}] + \Pr[v_t \in \Gamma_G^+(\mathcal{D}_{t-1})] = \frac{(\Delta+1)|\mathcal{D}_{t-1}|}{n} + \sum_{v \in V} \frac{p_v^{\text{up}}}{n}.$$

By the definition of  $p_v^{\text{up}}$  in (12) and  $d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') = n |\ln \lambda - \ln \lambda'| \leq L$ , we have

$$\Pr[\gamma_t = 1 \mid \mathcal{D}_{t-1}] \leq \frac{(\Delta + 1)|\mathcal{D}_{t-1}|}{n} + \frac{2L}{n}.$$

By the definition of  $R_{\text{Hamil}} \triangleq \sum_{t=1}^T \gamma_t$ , we have

$$(53) \quad \mathbb{E}[R_{\text{Hamil}}] = \sum_{t=1}^T \mathbb{E}[\gamma_t] = \sum_{t=1}^T \mathbb{E}[\mathbb{E}[\gamma_t \mid \mathcal{D}_{t-1}]] \leq \sum_{t=1}^T \left( \frac{(\Delta + 1)\mathbb{E}[|\mathcal{D}_{t-1}|]}{n} + \frac{2L}{n} \right).$$

Next, we bound the expectation  $\mathbb{E}[|\mathcal{D}_t|]$ . In our implementation of the one-step local coupling for Hamiltonian update (Definition 6.3), we first construct the random set  $\mathcal{P} \subseteq V$  in (13). In the  $t$ -th step, where  $1 \leq t \leq T$ , given any  $\mathbf{X}_{t-1}$  and  $\mathbf{Y}_{t-1}$ , the  $\mathbf{X}_t$  and  $\mathbf{Y}_t$  is generated as follows.

- Let  $X'(u) = X_{t-1}(u)$  and  $Y'(u) = Y_{t-1}(u)$  for all  $u \in V \setminus \{v_t\}$ , sample  $(X'(v_t), Y'(v_t)) \in \{0, 1\}^2$  jointly from the optimal coupling  $D_{\text{opt}, \mathcal{I}, v_t}^{\sigma, \tau}$  of the marginal distributions  $\mu_{v_t, \mathcal{I}}(\cdot \mid \sigma)$  and  $\mu_{v_t, \mathcal{I}}(\cdot \mid \tau)$ , where  $\sigma = X_{t-1}(\Gamma_G(v_t))$  and  $\tau = Y_{t-1}(\Gamma_G(v_t))$ .
- Let  $\mathbf{X}_t = \mathbf{X}'$  and  $\mathbf{Y}_t = \mathbf{Y}'$ . If  $t \in \mathcal{P}$ , update the value of  $Y_t(v_t)$  using (14).

Note that  $\Omega_{\mathcal{I}} = \Omega_{\mathcal{I}'}$ . Since  $\mathcal{I}$  satisfies  $\lambda \leq \frac{2-\delta}{\Delta-2}$  with constant  $\delta > 0$ , by Lemma 7.2, for any feasible  $\mathbf{X}_{t-1}, \mathbf{Y}_{t-1} \in \Omega_{\mathcal{I}} = \Omega_{\mathcal{I}'}$ , we have

$$(54) \quad \mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}', \mathbf{Y}') \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] \leq \left(1 - \frac{\delta}{96n}\right) \rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}).$$

By Lemma 7.2, function  $\rho_{\mathcal{I}}(\cdot, \cdot)$ , seen as a function of  $2n$  variables, is  $12\Delta$ -Lipschitz. Let  $\mathcal{F}$  indicates whether  $t \in \mathcal{P}$ . We flip the value of  $Y_t(v_t)$  only if  $\mathcal{F}$  occurs. By (54), we have

$$\begin{aligned} \mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}_t, \mathbf{Y}_t) \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] &\leq \mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}', \mathbf{Y}') + 12\Delta\mathcal{F} \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] \\ &= \mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}', \mathbf{Y}') \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] + \mathbb{E}[12\Delta\mathcal{F} \mid \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}] \\ (\mathcal{F} \text{ is independent with } \mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) &\leq \left(1 - \frac{\delta}{96n}\right) \rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) + 12\Delta\mathbb{E}[\mathcal{F}] \\ &\leq \left(1 - \frac{\delta}{96n}\right) \rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) + 12\Delta \sum_{v \in V} \frac{p_v^{\text{up}}}{n} \\ (\text{by (12)}) &\leq \left(1 - \frac{\delta}{96n}\right) \rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) + \frac{24\Delta}{n} \sum_{v \in V} |\ln \lambda - \ln \lambda'| \\ (\text{by } d_{\text{Hamil}}(\mathcal{I}, \mathcal{I}') \leq L) &\leq \left(1 - \frac{\delta}{96n}\right) \rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1}) + \frac{24L\Delta}{n}. \end{aligned}$$

Note that  $\rho_{\mathcal{I}}(\mathbf{X}_0, \mathbf{Y}_0) = 0$  and  $\mathbf{X}_0(v) = \mathbf{Y}_0(v) = 0$  for all  $v \in V$ , the configurations  $\mathbf{X}_t, \mathbf{Y}_t$  are feasible for all  $t \geq 0$ . Thus, we have

$$\mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}_t, \mathbf{Y}_t)] \leq \left(1 - \frac{\delta}{96n}\right) \mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}_{t-1}, \mathbf{Y}_{t-1})] + \frac{24L\Delta}{n}.$$

Thus  $\mathbb{E}[\rho_{\mathcal{I}}(\mathbf{X}_t, \mathbf{Y}_t)] \leq \frac{5000L\Delta}{\delta}$ . By the up-bound to Hamming in Lemma 7.2, we have

$$\mathbb{E}[|\mathcal{D}_t|] \leq \frac{5000L\Delta}{\delta}.$$

Thus, by (53), we have

$$\mathbb{E}[R_{\text{Hamil}}] \leq \frac{50000\Delta^2 TL}{\delta n} = O\left(\frac{\Delta^2 TL}{\delta n}\right).$$

□

*Proof of Lemma 7.5.* By the definition of  $R_{\text{graph}}$  in (23) and the linearity of the expectation, we have

$$\mathbb{E}[R_{\text{graph}}] = \sum_{t=1}^T \mathbb{E}[\gamma_t] = \sum_{t=1}^T \mathbb{E}[\mathbb{E}[\gamma_t \mid \mathcal{D}_{t-1}]].$$

Recall  $Y_t = 1 [v_t \in \mathcal{S} \cup \Gamma_G^+(\mathcal{D}_{t-1})]$  and  $v_t \in V$  is uniformly at random given  $\mathcal{D}_{t-1}$ . Note that  $|\Gamma_G^+(\mathcal{D}_{t-1})| \leq (\Delta + 1)|\mathcal{D}_{t-1}|$  and  $|\mathcal{S}| \leq 2|E \oplus E'| \leq 2L$ . We have

$$(55) \quad \mathbb{E} [R_{\text{graph}}] \leq \sum_{t=1}^T \mathbb{E} \left[ \frac{(\Delta + 1)|\mathcal{D}_{t-1}| + 2L}{n} \right] = \frac{(\Delta + 1)}{n} \sum_{t=1}^T \mathbb{E} [|\mathcal{D}_{t-1}|] + \frac{2LT}{n}.$$

Suppose  $\lambda \leq \frac{2-\delta}{\Delta_G-2}$  and  $\Omega_{I'} \subseteq \Omega_I$ . The other condition follows from symmetry. We claim that

$$(56) \quad \forall 0 \leq t \leq T : \quad \mathbb{E} [|\mathcal{D}_t|] \leq \frac{10000\Delta L}{\delta}.$$

Combining (55) and (56), we have

$$\mathbb{E} [R_{\text{graph}}] \leq \frac{100000\Delta LT}{n\delta} = O\left(\frac{\Delta^2 LT}{n\delta}\right).$$

This proves the lemma.

We now prove (56). Let  $(X_t, Y_t)_{t \geq 0}$  be the one-step local coupling for updating edges (Definition 6.9). We claim the following result

$$(57) \quad \forall \sigma \in \Omega_I, \tau \in \Omega_{I'} \subseteq \Omega_I, \mathbb{E} [\rho_I(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \leq \left(1 - \frac{\delta}{96n}\right) \cdot \rho_I(\sigma, \tau) + \frac{48\Delta L}{n},$$

where  $\rho_I$  is the potential function in Lemma 7.2. Assume (57) holds. Since  $X_0 = Y_0 = \{0\}^V$  and  $\Omega_{I'} \subseteq \Omega_I$ , we must have  $X_{t-1}, Y_{t-1} \in \Omega_I$ . Taking expectation over  $X_{t-1}$  and  $Y_{t-1}$ , we have

$$(58) \quad \mathbb{E} [\rho_I(X_t, Y_t)] \leq \left(1 - \frac{\delta}{96n}\right) \mathbb{E} [\rho_I(X_{t-1}, Y_{t-1})] + \frac{48\Delta L}{n}.$$

Note that  $X_0 = Y_0$ , we have

$$(59) \quad \rho_I(X_0, Y_0) = 0.$$

Combining (58), (59) and upper-bound Hamming property in Lemma 7.2 implies

$$\forall 0 \leq t \leq T : \quad \mathbb{E} [|\mathcal{D}_t|] \leq \mathbb{E} [\rho_I(X_t, Y_t)] \leq \frac{10000\Delta L}{\delta}.$$

This proves the claim in (56).

We finish the proof by proving the claim in (57). Let  $(X'_t, Y'_t)_{t \geq 0}$  be the one-step optimal coupling for Gibbs sampling on instance  $\mathcal{I}$  (Definition 4.2). Since  $\mathcal{I}$  satisfies  $\lambda \leq \frac{2-\delta}{\Delta_G-2}$ , by Lemma 7.2, we have

$$(60) \quad \forall \sigma, \tau \in \Omega_I : \quad \mathbb{E} [\rho_I(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] \leq \left(1 - \frac{\delta}{96n}\right) \cdot \rho_I(\sigma, \tau).$$

According to the coupling, we can rewrite the expectation in (60) as follows:

$$(61) \quad \mathbb{E} [\rho_I(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] = \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right],$$

where  $(C_v^{X'}, C_v^{Y'}) \sim D_{\text{opt}, I_v}^{\sigma, \tau}, D_{\text{opt}, I_v}^{\sigma, \tau}$  is the optimal coupling between  $\mu_{v, \mathcal{I}}(\cdot \mid \sigma)$  and  $\mu_{v, \mathcal{I}}(\cdot \mid \tau)$ , and the configuration  $\sigma^{v \leftarrow C_v^{X'}} \in Q^V$  is defined as

$$\sigma^{v \leftarrow C_v^{X'}}(u) \triangleq \begin{cases} C_v^{X'} & \text{if } u = v \\ \sigma(u) & \text{if } u \neq v \end{cases}$$

and the configuration  $\tau^{v \leftarrow C_v^{Y'}} \in Q^V$  is defined in a similar way.

Similarly, we can rewrite the expectation in (57) as follows:

$$(62) \quad \mathbb{E} [\rho_I(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] = \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right],$$

where  $(C_v^X, C_v^Y) \sim D_{I_v, I'_v}^{\sigma, \tau}$ , where  $D_{I_v, I'_v}^{\sigma, \tau}$  is the local coupling defined in (21).

The following two properties hold for (61) and (62).

- If  $v \notin \mathcal{S}$ , by the definition of  $D_{I_v, I'_v}^{\sigma, \tau}(\cdot, \cdot)$  in (21), it holds that  $D_{I_v, I'_v}^{\sigma, \tau} = D_{\text{opt}, I_v}^{\sigma, \tau}$ . Hence

$$\forall v \notin \mathcal{S} : \quad \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right] = \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right].$$

- If  $v \in \mathcal{S}$ , then it holds that  $H(\sigma^{v \leftarrow C_v^X}, \sigma^{v \leftarrow C_v^{X'}}) \leq 1$  and  $H(\tau^{v \leftarrow C_v^Y}, \tau^{v \leftarrow C_v^{Y'}}) \leq 1$ , where  $H$  is the Hamming distance. Since  $\Omega_{I'} \subseteq \Omega_I$ , it holds that  $\sigma^{v \leftarrow C_v^{X'}}, \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y}, \tau^{v \leftarrow C_v^{Y'}} \in \Omega_I$ . Since the function  $\rho_I$  is  $12\Delta$ -Lipschitz, we have

$$\forall v \in \mathcal{S} : \quad \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right] \leq \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right] + 24\Delta.$$

Combining above two properties with (60), (61) and (62), we have for any  $\sigma \in, \tau \in \Omega$ ,

$$\begin{aligned} & \mathbb{E} [\rho_I(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \\ &= \frac{1}{n} \sum_{v \in V} \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^X}, \tau^{v \leftarrow C_v^Y} \right) \right] \\ &\leq \frac{1}{n} \sum_{v \notin \mathcal{S}} \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right] + \frac{1}{n} \sum_{v \in \mathcal{S}} \left( \mathbb{E} \left[ \rho_I \left( \sigma^{v \leftarrow C_v^{X'}}, \tau^{v \leftarrow C_v^{Y'}} \right) \right] + 24\Delta \right) \\ (*) &\leq \mathbb{E} [\rho_I(X'_t, Y'_t) \mid X'_{t-1} = \sigma \wedge Y'_{t-1} = \tau] + \frac{48L\Delta}{n} \\ &\leq \left( 1 - \frac{\delta}{96n} \right) \cdot \rho_I(\sigma, \tau) + \frac{48L\Delta}{n}, \end{aligned}$$

where (\*) holds due to  $|\mathcal{S}| \leq 2L$ . This proves the claim in (57).  $\square$

Finally, we prove Lemma 7.2. This proof is based on the coupling technique in [Vig99].

*Proof of Lemma 7.2.* We give a potential function  $\rho_I$  for the hard core instance  $\mathcal{I}$ . We mainly use Vigoda's potential function in [Vig99]. However, we need to slightly modify Vigoda's potential function to handle the isolated vertices.

Recall that for hard core model,  $Q = \{0, 1\}$ . For any  $\sigma \in Q^V$ ,  $\sigma(v) = 1$  represents  $v$  is occupied and  $\sigma(v) = 0$  represents  $v$  is unoccupied. For each vertex  $v \in V$ , we use  $\deg(v)$  to denote the degree of  $v$  in graph  $G = (V, E)$ . We divide the graph  $G = (V, E)$  into two graphs  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$  such that

$$\begin{aligned} V_1 &= \{v \in V \mid \deg(v) = 0\}, & E_1 &= \emptyset, \\ V_2 &= V \setminus V_1, & E_2 &= E. \end{aligned}$$

Thus  $G_1$  is an empty graph and  $G_2$  contains no isolated vertex. The potential function  $\rho_I$  is defined as

$$\forall \sigma, \tau \in \Omega_I : \quad \rho_I(\sigma, \tau) \triangleq 4\rho_1(\sigma(V_1), \tau(V_1)) + 4\rho_2(\sigma(V_2), \tau(V_2)).$$

Here,  $\rho_1$  is the potential function on  $G_1$ , which is the Hamming distance:

$$\rho_1(\sigma(V_1), \tau(V_1)) = \sum_{v \in V_1} 1 [\sigma(v) \neq \tau(v)].$$

And  $\rho_2(\sigma(V_2), \tau(V_2))$  is the Vigoda's potential function [Vig99] on the graph  $G_2$ . Formally, let  $D = \{v \in V_2 \mid \sigma(v) \neq \tau(v)\}$ . For each  $v \in V_2$ , let  $d_v = |D \cap \Gamma_{G_2}(v)|$ . Let  $c = \frac{\Delta\lambda}{\Delta\lambda+2}$ , where  $\Delta$  is the maximum degree of graph  $G$ . Note that the maximum degree of graph  $G_2$  is also  $\Delta$ . The potential function  $\rho_2(\sigma(V_2), \tau(V_2))$  is defined as

$$\alpha_v = \begin{cases} \deg(v) & \text{if } v \in D \\ 0 & \text{otherwise;} \end{cases} \quad \beta_v = \begin{cases} -cd_v & \text{if } \exists w \in \Gamma_{G_2}(v) \text{ such that } \sigma(w) = \tau(w) = 1 \\ -c(d_v - 1) & \text{if there is no such } w \text{ and } d_v > 1 \\ 0 & \text{otherwise;} \end{cases}$$

$$\rho_2(\sigma(V_2), \tau(V_2)) = \sum_{v \in V_2} (\alpha_v + \beta_v).$$

It is easy to see  $\rho_I(\sigma, \sigma) = 0$  and  $\max_{\sigma, \tau \in \Omega_I} \rho_I(\sigma, \tau) = \Delta n$ . We then verify other properties for  $\rho_I$ .

At first, we prove the upper-bound to Hamming property. For function  $\rho_1$ , it holds that

$$\rho_1(\sigma(V_1), \tau(V_1)) = H(\sigma(V_1), \tau(V_1)).$$

For function  $\rho_2$ , it holds that

$$\rho_2(\sigma(V_2), \tau(V_2)) = \sum_{v \in V_2} (\alpha_v + \beta_v) = \sum_{v \in D} \alpha_v + \sum_{v \in V_2} \beta_v \geq \sum_{v \in D} \sum_{w \in \Gamma_{G_2}(v)} (1 - c),$$

where the last inequality holds due to  $\sum_{v \in V_2} \beta_v \geq -\sum_{v \in V_2} cd_v = -c \sum_{v \in D} \deg(v)$ . Since graph  $G_2$  contains no isolated vertex, then  $|\Gamma_{G_2}(v)| = \deg(v) \geq 1$  for all  $v \in D$ . Note  $c < 1$ . Thus

$$\rho_2(\sigma(V_2), \tau(V_2)) \geq |D|(1 - c) = |D| \frac{2}{\Delta\lambda + 2} \geq \frac{|D|}{4} = \frac{1}{4}H(\sigma(V_2), \tau(V_2)),$$

where  $\frac{2}{\lambda\Delta + 2} \geq \frac{1}{4}$  is because  $\lambda < \frac{2}{\Delta - 2}$  and  $\Delta \geq 3$ . Combining together we have

$$\rho_I(\sigma, \tau) = 4\rho_1(\sigma(V_1), \tau(V_1)) + 4\rho_2(\sigma(V_2), \tau(V_2)) \geq H(\sigma, \tau).$$

This also implies  $\rho_I(\sigma, \tau) \geq 1$  [ $\sigma \neq \tau$ ].

Next, we show the function  $\rho_I$  is  $12\Delta$ -Lipschitz. Recall  $V_1 \cap V_2 = \emptyset$ ,  $V_1 \cup V_2 = V$  and

$$\rho_I(\sigma, \tau) = 4\rho_1(\sigma(V_1), \tau(V_1)) + 4\rho_2(\sigma(V_2), \tau(V_2)).$$

Since  $\rho_1$  is the Hamming distance, it is easy to see  $\rho_1$  is 1-Lipschitz. To give the Lipschitz constant for  $\rho_2$ . We extend the function  $\rho_2$  as follows. Suppose the function  $\rho_2$  is defined over  $Q^{V_2} \times Q^{V_2}$ , where  $Q = \{0, 1\}$ . For any  $x, y, x', y' \in Q^{V_2}$  such that  $H(xy, x'y') = 1$ , it is easy to verify the extended function  $\rho_2$  satisfies

$$|\rho_2(x, y) - \rho_2(x', y')| \leq 3\Delta.$$

This implies the original function  $\rho_2$  is  $3\Delta$ -Lipschitz. Hence, the function  $\rho_I$  is  $12\Delta$ -Lipschitz.

Finally, we prove the step-wise decay property. Let  $(X_t^{(1)})_{t \geq 0}, (Y_t^{(1)})_{t \geq 0}$  be the Gibbs sampling chains for hard core model on graph  $G_1$ . Since  $G_1$  is a graph consisting of isolated vertices, then the one step optimal coupling  $(X_t^{(1)}, Y_t^{(1)})_{t \geq 0}$  satisfies

$$\forall \sigma, \tau \in \Omega_I : \mathbb{E} \left[ \rho_1 \left( X_t^{(1)}, Y_t^{(1)} \mid X_{t-1}^{(1)} = \sigma(V_1) \wedge Y_{t-1}^{(1)} = \tau(V_1) \right) \right] \leq \left( 1 - \frac{1}{|V_1|} \right) \rho_1(\sigma(V_1), \tau(V_1)).$$

Let  $(X_t^{(2)})_{t \geq 0}, (Y_t^{(2)})_{t \geq 0}$  be the Gibbs sampling chains for hard core model on graph  $G_2$ . If  $\lambda \leq \frac{2-\delta}{\Delta-2} = \frac{2(1-\delta/2)}{\Delta-2}$ , then due to Vigoda's proof<sup>3</sup>, the one step optimal coupling  $(X_t^{(2)}, Y_t^{(2)})_{t \geq 0}$  satisfies:

$$\forall \sigma, \tau \in \Omega_I : \mathbb{E} \left[ \rho_2 \left( X_t^{(2)}, Y_t^{(2)} \mid X_{t-1}^{(2)} = \sigma(V_2) \wedge Y_{t-1}^{(2)} = \tau(V_2) \right) \right] \leq \left( 1 - \frac{\delta}{96|V_2|} \right) \rho_2(\sigma(V_2), \tau(V_2)).$$

Let  $(X_t)_{t \geq 0}, (Y_t)_{t \geq 0}$  be the Gibbs sampling chains for hard core model on graph  $G$ . If  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , then the one step optimal coupling  $(X_t, Y_t)_{t \geq 0}$  satisfies:

$$\begin{aligned} \forall \sigma, \tau \in \Omega_I : \mathbb{E} [\rho_I(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] & \\ &= \frac{|V_1|}{n} \left( \left( 1 - \frac{1}{|V_1|} \right) 4\rho_1(\sigma(V_1), \tau(V_1)) + 4\rho_2(\sigma(V_2), \tau(V_2)) \right) \\ &\quad + \frac{|V_2|}{n} \left( 4\rho_1(\sigma(V_1), \tau(V_1)) + \left( 1 - \frac{\delta}{96|V_2|} \right) 4\rho_2(\sigma(V_2), \tau(V_2)) \right) \\ &\leq \left( 1 - \frac{\min\{\delta/96, 1\}}{n} \right) \rho_I(\sigma, \tau). \end{aligned}$$

Thus, the potential function  $\rho_I$  satisfies the step-wise decay property.

$$\forall \sigma, \tau \in \Omega_I : \mathbb{E} [\rho_I(X_t, Y_t) \mid X_{t-1} = \sigma \wedge Y_{t-1} = \tau] \leq \left( 1 - \frac{\delta/96}{n} \right) \rho_I(\sigma, \tau).$$

<sup>3</sup>It can be verified that in Vigoda's proof [Vig99], the Markov chain for sampling hard core is indeed the Gibbs sampling and the coupling for analysis is indeed the one step-optimal coupling for Gibbs sampling.



This proves the lemma.  $\square$

## 8. PROOFS FOR DYNAMIC INFERENCE

**8.1. Proof of the main theorem.** Our dynamic inference algorithm is given as follows. For each MRF instance  $\mathcal{I} = (V, E, Q, \Phi)$ , where  $n = |V|$ , our dynamic inference algorithm maintains  $N(n)$  independent samples  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(N(n))} \in Q^V$  satisfying each  $d_{\text{TV}}(\mu_{\mathcal{I}}, \mathbf{X}^{(i)}) \leq \epsilon(n)$  and the estimator  $\hat{\theta}(\mathcal{I}) = \mathcal{E}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(N(n))})$  for  $\theta(\mathcal{I})$ . Given an update that modifies  $\mathcal{I}$  to  $\mathcal{I}' = (V', E', Q, \Phi')$  where  $n' = |V'|$ , our algorithm does as follows.

- *Update the sample sequence.* Update  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(N(n))}$  to  $N(n')$  independent random samples  $\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}, \dots, \mathbf{Y}^{(N(n'))} \in Q^{V'}$  such that each  $d_{\text{TV}}(\mu_{\mathcal{I}'}, \mathbf{Y}^{(i)}) \leq \epsilon(n')$  and output the difference between two sample sequences.
- *Update the estimator.* Given the difference between two sample sequences  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(N(n))}$  and  $\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}, \dots, \mathbf{Y}^{(N(n'))}$ , update  $\hat{\theta}(\mathcal{I})$  to  $\hat{\theta}(\mathcal{I}') = \mathcal{E}_{\theta}(\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}, \dots, \mathbf{Y}^{(N(n'))})$  using the black-box algorithm in Definition 2.3.

Obviously,  $\hat{\theta}(\mathcal{I}')$  is an  $(N, \epsilon)$ -estimator for  $\theta(\mathcal{I}')$ .

The sample sequence is maintained and updated by the dynamic sampling algorithm in Theorem 6.1. By Theorem 6.1, we have the space cost for maintaining the sample sequence is  $O(nN(n) \log n)$  memory words, each of  $O(\log n)$  bits. By following the proof of Theorem 6.1, it is easy to verify that the expected time cost for each update is  $O(\Delta^2 LN(n) \log^3 n + \Delta n)$ .

The estimator is maintained and updated by the black-box algorithm in Definition 2.3. By Lemma 6.19, we have  $N(n) \leq \text{poly}(n)$ . Combining with Definition 2.3, we have the space cost for maintaining the estimator is  $(n \cdot N(n) + K) \text{polylog}(n)$  bits. Let  $\mathcal{D}$  be the size of the difference between two sample sequences as defined in (3). We can follow the proof of Theorem 6.1 to bound the expectation of  $\mathcal{D}$ . Let  $T = \left\lceil \frac{n}{\delta} \log \frac{n}{\epsilon(n)} \right\rceil$  and  $T' = \left\lceil \frac{n'}{\delta} \log \frac{n'}{\epsilon(n')} \right\rceil$ . Since  $|n - n'| \leq L = o(n)$ , we have  $|T - T'| = O(L \log n)$  (due to Lemma 6.17). Combining (39), (45) and (7) yields

$$\mathbb{E}[\mathcal{D}] = |N(n) - N(n')| \cdot \max\{n, n'\} + O(L + |T - T'|) \cdot N(n) = O(LN(n) \log n),$$

where the last equation holds because  $N(n) - N(n') = O(\frac{N(n)}{n})$  (due to Lemma 6.19). Combining with Definition 2.3, we have the expected time cost for updating the estimator is  $LN(n) \text{polylog}(n)$ .

In summary, our dynamic inference algorithm maintains an estimator for the current MRF instance  $\mathcal{I}$ , using extra  $\tilde{O}(nN(n) + K)$  memory words, each of  $O(\log n)$  bits, such that when  $\mathcal{I}$  is updated to  $\mathcal{I}'$ , the algorithm updates the estimator within expected time cost

$$\begin{aligned} \mathbb{E}[T_{\text{cost}}] &= \mathbb{E}[T_{\text{sample}}] + \mathbb{E}[T_{\text{estimator}}] \\ &= O(\Delta^2 LN(n) \log^3 n + \Delta n) + LN(n) \text{polylog}(n) \\ &= \tilde{O}(\Delta^2 LN(n) + \Delta n). \end{aligned}$$

**8.2. Dynamic inference on specific models.** Applying our dynamic inference algorithm on Ising model,  $q$ -coloring and hardcore model yields the following result.

**Theorem 8.1.** *There exist dynamic inference algorithms as stated in Theorem 3.2 with the same space cost  $\tilde{O}(nN(n) + K)$ , and expected time cost  $\tilde{O}(\Delta^2 LN(n) + \Delta n)$  for each update, if the input instance  $\mathcal{I}$  with  $n$  vertices and the updated instance  $\mathcal{I}'$  with  $d(\mathcal{I}, \mathcal{I}') \leq L = o(n)$  both are:*

- *Ising models with temperature  $\beta$  and arbitrary local fields where  $\exp(-2|\beta|) \geq 1 - \frac{2-\delta}{\Delta+1}$ ;*
- *proper  $q$ -colorings with  $q \geq (2 + \delta)\Delta$ ;*
- *hardcore models with fugacity  $\lambda \leq \frac{2-\delta}{\Delta-2}$ , but with an alternative time cost for each update*

$$\tilde{O}(\Delta^3 LN(n) + \Delta n),$$

where  $\delta > 0$  is a constant,  $\Delta = \max\{\Delta_G, \Delta_{G'}\}$ ,  $\Delta_G$  and  $\Delta_{G'}$  denote the maximum degree of the input graph and updated graph respectively.

With the dynamic sampling algorithm in Theorem 7.1, Theorem 8.1 can be proved by going through the same proof in Section 8.1.

## 9. CONCLUSION

In this paper we study probabilistic inference problem in a graphical model when the model itself is changing dynamically with time. We study the non-local updates so that two consecutive graphical models may differ everywhere as long as the total amount of their difference is bounded. This general setting covers many typical applications. We give a sampling-based dynamic inference algorithm that maintains an inference solution efficiently against the dynamic inputs. The algorithm significantly improves the time cost compared to the static sampling-based inference algorithm.

Our algorithm generically reduces the dynamic inference to dynamic sampling problem. Our main technical contribution is a dynamic Gibbs sampling algorithm that maintains random samples for graphical models dynamically changed by non-local updates. Such technique is extendable to all single-site dynamics. This gives us a systematic approach for transforming classic MCMC samplers on static inputs to the sampling and inference algorithms in a dynamic setting. Our dynamic algorithms are efficient as long as the one-step optimal coupling exhibits a step-wise decay, a key property that has been widely used in supporting efficient MCMC sampling in the classic static setting and captured by the Dobrushin-Shlosman condition.

Our result is the first one that shows the possibility of efficient probabilistic inference in dynamically changing graphical models (especially when the graphical models are changed by non-local updates). Our dynamic inference algorithm has potentials in speeding up the iterative algorithms for learning graphical models, which deserves more theoretical and experimental research. In this paper, we focus on discrete graphical models and sampling-based inference algorithms. Important future directions include considering more general distributions and the dynamic algorithms based on other inference techniques.

## REFERENCES

- [ADK<sup>+</sup>16] Ittai Abraham, David Durfee, Ioannis Koutis, Sebastian Krinninger, and Richard Peng. On fully dynamic graph sparsifiers. In *FOCS*, 2016.
- [AQ<sup>+</sup>17] Osvaldo Anacleto, Catriona Queen, et al. Dynamic chain graph models for time series network data. *Bayesian Anal.*, 12(2):491–509, 2017.
- [BC16] Aaron Bernstein and Shiri Chechik. Deterministic decremental single source shortest paths: beyond the  $o(mn)$  bound. In *STOC*, 2016.
- [BD97] Russ Bubley and Martin Dyer. Path coupling: A technique for proving rapid mixing in Markov chains. In *FOCS*, 1997.
- [CLRS09] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2009.
- [CW07] Carlos M. Carvalho and Mike West. Dynamic matrix-variate graphical models. *Bayesian Anal.*, 2(1):69–97, 2007.
- [DG00] Martin Dyer and Catherine Greenhill. On Markov chains for independent sets. *J. Algorithms*, 35(1):17–49, 2000.
- [DGGP18] David Durfee, Yu Gao, Gramoz Goranci, and Richard Peng. Fully dynamic effective resistances. *arXiv preprint arXiv:1804.04038*, 2018.
- [DGGP19] David Durfee, Yu Gao, Gramoz Goranci, and Richard Peng. Fully dynamic spectral vertex sparsifiers and applications. In *STOC*, 2019.
- [DGJ08] Martin Dyer, Leslie Ann Goldberg, and Mark Jerrum. Dobrushin conditions and systematic scan. *Combin. Probab. Comput.*, 17(6):761–779, 2008.
- [DS85a] Roland L Dobrushin and Senya B Shlosman. Completely analytical Gibbs fields. In *Statistical Physics and Dynamical Systems*, pages 371–403. Springer, 1985.
- [DS85b] Roland Lvovich Dobrushin and Senya B Shlosman. Constructive criterion for the uniqueness of Gibbs field. In *Statistical Physics and Dynamical Systems*, pages 347–370. Springer, 1985.

- [DS87] RL Dobrushin and SB Shlosman. Completely analytical interactions: constructive description. *J. Statist. Phys.*, 46(5-6):983–1014, 1987.
- [DSOR16] Christopher De Sa, Kunle Olukotun, and Christopher Ré. Ensuring rapid mixing and low bias for asynchronous Gibbs sampling. In *ICML*, 2016.
- [FG19] Sebastian Forster and Gramoz Goranci. Dynamic low-stretch trees via dynamic low-diameter decompositions. In *STOC*, pages 377–388, 2019.
- [FVY19] Weiming Feng, Nisheeth K Vishnoi, and Yitong Yin. Dynamic sampling from graphical models. In *STOC*, 2019.
- [GHP18] Gramoz Goranci, Monika Henzinger, and Pan Peng. Dynamic Effective Resistances and Approximate Schur Complement on Separable Graphs. In *ESA*, volume 112, 2018.
- [GŠV15] Andreas Galanis, Daniel Štefankovič, and Eric Vigoda. Inapproximability for antiferromagnetic spin systems in the tree nonuniqueness region. *J. ACM*, 62(6):50, 2015.
- [GŠV16] Andreas Galanis, Daniel Štefankovič, and Eric Vigoda. Inapproximability of the partition function for the antiferromagnetic Ising and hard-core models. *Combin. Probab. Comput.*, 25(04):500–559, 2016.
- [Hay06] Thomas P Hayes. A simple condition implying rapid mixing of single-site dynamics on spin systems. In *FOCS*, 2006.
- [Hin12] Geoffrey E Hinton. A practical guide to training restricted boltzmann machines. In *Neural Networks: Tricks of the Trade*, pages 599–619. Springer, 2012.
- [HKN14] Monika Henzinger, Sebastian Krinninger, and Danupon Nanongkai. Decremental single-source shortest paths on undirected graphs in near-linear total update time. In *FOCS*, 2014.
- [HKN16] Monika Henzinger, Sebastian Krinninger, and Danupon Nanongkai. Dynamic approximate all-pairs shortest paths: Breaking the  $O(mn)$  barrier and derandomization. *SIAM J. Comput.*, 45(3):947–1006, 2016.
- [Jer95] Mark Jerrum. A very simple algorithm for estimating the number of  $k$ -colorings of a low-degree graph. *Random Structures & Algorithms*, 7(2):157–165, 1995.
- [JVV86] Mark Jerrum, Leslie G. Valiant, and Vijay V. Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoret. Comput. Sci.*, 43:169–188, 1986.
- [KFB09] Daphne Koller, Nir Friedman, and Francis Bach. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [LMV19] Holden Lee, Oren Mangoubi, and Nisheeth Vishnoi. Online sampling from log-concave distributions. In *NIPS*, 2019.
- [LP17] David A Levin and Yuval Peres. *Markov chains and mixing times*. American Mathematical Soc., 2017.
- [LV99] Michael Luby and Eric Vigoda. Fast convergence of the Glauber dynamics for sampling independent sets. *Random Structures & Algorithms*, 15(3-4):229–241, 1999.
- [MM09] Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
- [NR17] Hariharan Narayanan and Alexander Rakhlin. Efficient sampling from time-varying log-concave distributions. *J. Mach. Learn. Res.*, 18(1):4017–4045, 2017.
- [NSWN17] Danupon Nanongkai, Thatchaphol Saranurak, and Christian Wulff-Nilsen. Dynamic minimum spanning forest with subpolynomial worst-case update time. In *FOCS*, 2017.
- [QS93] Catriona M. Queen and Jim Q. Smith. Multiregression dynamic models. *J. Roy. Statist. Soc. Ser. B*, 55(4):849–870, 1993.
- [RKD<sup>+</sup>19] Cedric Renggli, Bojan Karlaš, Bolin Ding, Feng Liu, Kevin Schawinski, Wentao Wu, and Ce Zhang. Continuous integration of machine learning models: A rigorous yet practical treatment. In *SysML*, 2019.
- [ŠVV09] Daniel Štefankovič, Santosh Vempala, and Eric Vigoda. Adaptive simulated annealing: A near-optimal connection between sampling and counting. *J. ACM*, 56(3):18, 2009.
- [SWA09] Padhraic Smyth, Max Welling, and Arthur U Asuncion. Asynchronous distributed learning of topic models. In *NIPS*, 2009.
- [Vig99] Eric Vigoda. Fast convergence of the Glauber dynamics for sampling independent sets: Part II. Technical Report TR-99-003, International Computer Science Institute, 1999.

- [WJ08] Martin J. Wainwright and Michael I. Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- [WN17] Christian Wulff-Nilsen. Fully-dynamic minimum spanning forest with improved worst-case update time. In *STOC*, 2017.