



Contents lists available at ScienceDirect

Discrete Applied Mathematics

journal homepage: [www.elsevier.com/locate/dam](http://www.elsevier.com/locate/dam)

# Reconstructing Markov processes from independent and anonymous experiments

Silvio Micali, Zeyuan Allen Zhu\*

MIT CSAIL, United States

## ARTICLE INFO

### Article history:

Received 21 June 2014

Received in revised form 7 June 2015

Accepted 21 June 2015

Available online xxxx

### Keywords:

Graph reconstruction

Random walk

Markov process

Local algorithms

## ABSTRACT

We investigate the problem of *exactly* reconstructing, with high confidence and up to isomorphism, the ball of radius  $r$  centered at the starting state of a Markov process from *independent* and *anonymous* experiments. In an anonymous experiment, the states are visited according to the underlying transition probabilities, but no global state names are known: one can only recognize whether two states, *reached within the same experiment*, are the same.

We prove quite tight bounds for such exact reconstruction in terms of both the number of experiments and their lengths.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

The problem of reconstructing a large “object” from partial observations is quite fundamental, and arises in many fields, such as system biology [22,29], social networks [38,27], brain networks [36,19], telecommunication networks [10], and many others.

We investigate a more complex type of reconstruction. In essence, our goal is to reconstruct a Markov process from the records produced by *limited* observers acting *independently*, without coordination, and without even sharing a common “name space”. Let us explain.

### 1.1. Our model

**Our Markov model.** In a Markov process, we denote the underlying transition graph by  $G = (V, E)$  and the starting vertex by  $v$ . In this paper, the graph  $G$  is undirected and has infinitely many vertices, each of finite degree. An infinite sequence of vertices is generated by the following process. The first vertex is  $v$ , and, if the  $i$ th vertex is  $u$ , then the  $(i + 1)$ -st vertex is chosen at random uniformly and independently among the neighbors of  $u$ .

A sequence of vertices so generated is called a *random walk*. If  $(v \Rightarrow) v_0 \rightarrow v_1 \rightarrow \dots$  is a random walk, then  $v_0 \rightarrow \dots \rightarrow v_\ell$  is a *random walk of length  $\ell$* .

**Note.** Assuming that  $G$  is undirected and unweighted allows us to present our results in the cleanest way. We shall discuss how to relax both assumptions in Section 1.4. Assuming that  $G$  has infinitely many vertices is a simple way to force us to consider only “local” algorithms: essentially, algorithms whose performance does not depend on the size of the whole graph, which may be larger than all the parameters we shall care about.

\* Corresponding author.

E-mail addresses: [silvio@csail.mit.edu](mailto:silvio@csail.mit.edu) (S. Micali), [zeyuan@csail.mit.edu](mailto:zeyuan@csail.mit.edu) (Z.A. Zhu).

<http://dx.doi.org/10.1016/j.dam.2015.06.035>

0166-218X/© 2015 Elsevier B.V. All rights reserved.

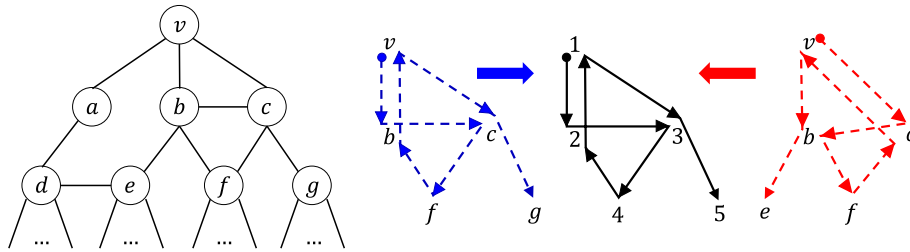


Fig. 1. An example illustrating the definition of anonymous experiment.

**Our anonymous observation model.** If  $v_0 \rightarrow v_1 \rightarrow \dots$  is a random walk, then its corresponding (*anonymous*) experiment is the sequence of integers  $f(v_0) \rightarrow f(v_1) \rightarrow \dots$ , where  $f(v_i) \stackrel{\text{def}}{=} |\{v_0, \dots, v_i\}|$  and  $i^*$  is the smallest integer  $j$  such that  $v_j = v_i$ . Intuitively,  $f(u)$  maps  $u$  to an integer indicating that  $u$  is the  $f(u)$ th distinct vertex in this walk.

**Example.** In the graph of Fig. 1, the length-7 walk  $v \rightarrow b \rightarrow c \rightarrow f \rightarrow b \rightarrow v \rightarrow c \rightarrow g$  corresponds to the anonymous experiment  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 2 \rightarrow 1 \rightarrow 3 \rightarrow 5$ .

Note the walk  $v \rightarrow c \rightarrow b \rightarrow f \rightarrow c \rightarrow v \rightarrow b \rightarrow e$  also corresponds to the same experiment.

**Rationale.** Markov processes naturally model physical systems. In essence, the possible “states” of the system are the vertices of the transition graph  $G$ , and, when put in its “initial state”  $v$ , the system evolves (i.e., new states are generated) according to the transition rules.

When the system is “new” – better said, studied for the first time – no one initially has any idea about the underlying graph  $G$ . However, each individual can, on his own, experiment with the system by putting it into its initial state, and independently observe its “evolution”: that is, a random walk in  $G$ .

In sum, each individual will observe and record the states encountered in a random walk. Since the system is new, no global names exist for the states. Thus each individual may very well use his own name space for the states encountered, and thus his record is an anonymous experiment as defined above.

Of course, an individual observer might consider writing down a full description of every state he sees. However, this may not be possible due to – say – memory limitations [15], or privacy reasons [26]. Also, an observer may not know how many details are sufficient to identify each encountered state. In any case, an anonymous experiment is a most compact and meaningful record.

### 1.2. Our results for the basic reconstruction problem

Whether human or not, a realistic observer has a bounded lifetime, and thus cannot visit more than  $\ell$  nodes in his random walk.<sup>1</sup> Thus, even with an unlimited number of such observers, one can at most learn  $\mathcal{D}_\ell$ , the distribution over the anonymous experiments of length  $\ell$  (that are induced from the random walks of the same length). Since our  $G$  has infinitely many vertices,  $\mathcal{D}_\ell$  cannot suffice to reconstruct the entire graph  $G$ . However, one may be able to use  $\mathcal{D}_\ell$  in order to reconstruct  $B(v, r)$ , the ball of center  $v$  and radius  $r$  (i.e., the subgraph of  $G$  induced by all vertices whose distance from  $v$  is at most  $r$ ). More precisely, one may be able to compute a graph  $G' = (V', E')$  and a distinguished vertex  $v' \in V'$ , such that  $G'$  is isomorphic to  $B(v, r)$  and the isomorphism maps  $v'$  to  $v$ . Thus, our basic reconstruction problem can be formulated as follows:

For every  $r$ , is there a length  $\ell$  such that  $B(v, r)$  is reconstructible (up to isomorphism) from  $\mathcal{D}_\ell$ ?

Notice that, given access to the distribution  $\mathcal{D}_\ell$ , one can also simulate access to the distributions  $\mathcal{D}_1, \dots, \mathcal{D}_{\ell-1}$ . Of course, although for now we are ignoring the complexity of learning these distributions, it would be nice if, given  $(\mathcal{D}_1, \dots, \mathcal{D}_\ell)$  as oracles, the reconstruction algorithm is efficient. Here, we say that  $(\mathcal{D}_1, \dots, \mathcal{D}_\ell)$  are given as oracles, if the algorithm is allowed to ask for the precise probability of  $\mathcal{D}_t(P)$  for any anonymous experiment  $P$  of length  $t \in [\ell]$ .

Notice too that, in principle, our basic reconstruction problem may be impossible. For instance, could there exist two different Markov processes,  $(G_1, v_1)$  and  $(G_2, v_2)$ , having the same distribution  $\mathcal{D}_\ell$  for all  $\ell \geq 0$ ? If this were the case, the two processes would be indistinguishable by any number of anonymous experiments, of any length, which immediately implies a negative answer to the above question. Yet, we provide a constructive proof showing that our basic reconstruction problem is indeed possible, when the underlying graph  $G$  is undirected.

**Theorem 1.** Let  $n$  be the number of vertices in  $B(v, r)$  and  $m$  the number of edges. One can reconstruct  $B(v, r)$  in time  $O(n^2)$  and with  $O(n^2)$  oracle accesses to  $(\mathcal{D}_1, \dots, \mathcal{D}_\ell)$ , where  $\ell = O(m)$ . Moreover, the reconstruction algorithm only makes membership queries to  $\text{supp}(\mathcal{D}_i)$  for  $i \in [\ell]$ .

In contrast, as we shall see in Section 1.4, this reconstruction becomes impossible when the underlying graph is directed but not strongly connected.

Is this algorithm tight? To answer this question we must refine our reconstruction problem.

<sup>1</sup> For concreteness, if he lives for at most 100 years, and each transition from node to node takes 1 s of time, then  $\ell = 100 \times 366 \times 24 \times 60 \times 60$ .



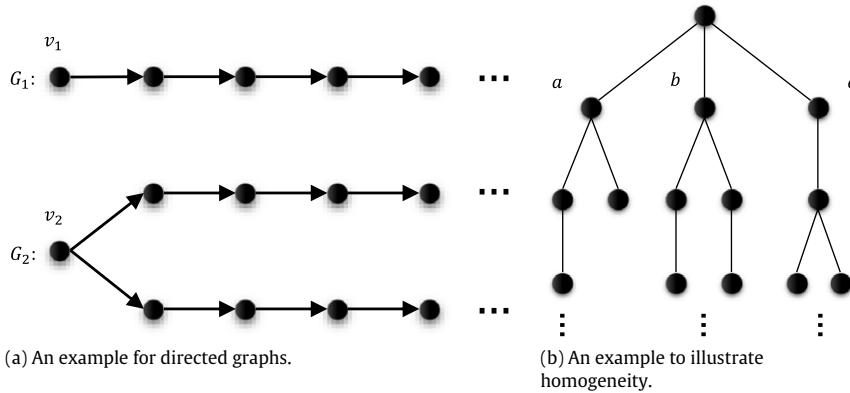


Fig. 2. Examples for extensions and improvements.

In the weighted case, our algorithm will reconstruct the “topology” of the underlying ball, that is, all edges in the ball, without their weights. This implies, for instance, if the random walk we studied has laziness – that is, at each vertex it stays at where it is with half probability, and goes to a random neighbor with another half probability – we can still reconstruct  $G$ . In general, reconstructing the weights too will require future work.

**Improvements.** The performance of our algorithm of Theorem 1 can be dramatically improved given reasonable guarantees about the topology of the underlying graph. One such guarantee is “non homogeneity”. Consider the simple graph (indeed a tree) in Fig. 2b.

In this graph, the three children of the root,  $a$ ,  $b$ , and  $c$ , cannot be distinguished until level 3. Indeed, all of them are indistinguishable at level 1, that is, in  $B(v, 1)$ . Vertex  $c$  can be distinguished from the others at level 2: indeed, in  $B(v, 2)$  vertex  $c$  has only one child (equivalently 2 neighbors), while each of  $a$  and  $b$  has 2 children. At this level, however, no way exists to distinguish  $a$  from  $b$ . But one additional level suffices.

Informally, we say that a graph  $G$  has homogeneity  $\omega$ , if for each vertex  $u$  in  $G$ , every two neighbors of  $u$  can be distinguished in a ball centered at  $u$  with at most  $\omega$  edges. Then, if the graph  $G$  is guaranteed to be of homogeneity  $\omega$ , the algorithm for Theorem 1 (without knowing  $\omega$ ) can be extended to reconstruct  $B(v, r)$  with experiments only of length  $\ell = O(r \cdot \omega)$ .

Notice that this specific improvement does not contradict the impossibility result of Theorem 2. Indeed, to prove Theorem 2 we exhibit a ball  $B(v, r)$  whose homogeneity is very large, namely,  $\omega = 2^{2^r}$ . In fact,  $B(v, r)$  is constructed so that  $B(v, r - 1)$  consists of a complete binary tree, and thus the two children of the root cannot be distinguished up to level  $r - 1$ .

1.5. Related work

**Graph reconstruction using queries.** The problem of reconstructing an unknown graph from oracle queries has been studied in many different contexts, and most notably using edge detection queries [16,2,1,5,6], edge counting queries [17,7,25], or distance queries [20,21,32,24].

In an edge detection query model, the oracle, on input a subset  $S$  of the vertices, returns if there exists an edge between any two vertices in  $S$ . Angluin and Chen [6] show that using  $O(\log n)$  adaptive queries per edge is sufficient for reconstructing an arbitrary graph, and this has been generalized to hypergraphs [5].

In an edge counting query model, the oracle, on input a subset  $S$  of the vertices, returns the number of edges between any two vertices in  $S$ . While Grebinski and Kucherov [17] prove tight bounds of  $O(dn)$  and  $O(n^2 / \log n)$  non-adaptive queries for  $d$ -degree-bounded and general graphs, in a more recent work, Mazzawi [25] shows that an information-theoretically tight bound of  $O(m \log(n^2/m) / \log m)$  can be achieved using non-adaptive queries for any graph with  $n$  vertices and  $m$  edges.

In a distance query model, the supported queries are of the form  $dist(u, v)$ , that is, the oracle returns the (possibly approximate) distance between any two given vertices. A lower bound of  $\Omega(n^2)$  queries is shown by Reyzin and Srivastava [32] for general graphs. Mathieu and Zhou [24] generalize this lower bound to allow approximate distance oracles, provide an upper bound of  $\tilde{O}(n^{3/2})$  for constant-degree graphs, and  $\tilde{O}(n)$  for outerplanar graphs.

All the results above are quite different from ours: the “name space” of the vertices are shared between different queries. As a result, if one is satisfied with a polynomial running time – say,  $O(n^2)$  – it is trivial to (even locally) reconstruct any graph using any of the oracles above.

**Learning graphical models.** Much work has been done in the machine learning community on learning the structures of graphical models. While we refer interested readers to Part III of Kollar and Friedman’s book [23], we summarize a few of them below.

A first type of research in this field assumes that the topology of a graphical model (e.g., a Bayesian network) is known, and focuses on estimating the parameters in the model. Two well-known methods are the maximum likelihood estimation and



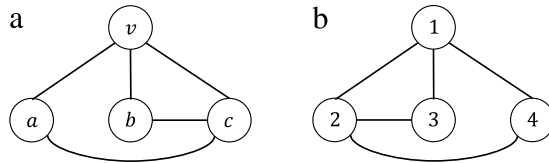


Fig. 3. An example to illustrate the proof of Theorem 1 for  $r = 1$ .

**Supporting graph.** Given an anonymous experiment  $P$  of length  $\ell$  that contains  $n$  distinct integers, one can define its supporting graph  $\text{Graph}(P) \stackrel{\text{def}}{=} (V, E)$ , where  $V = \{1, 2, \dots, n\}$  and  $(a, b) \in E$  if and only if  $a \rightarrow b$  (or  $b \rightarrow a$ ) appears in  $P$ . For instance, letting  $P = 1 \rightarrow 2 \rightarrow 3 \rightarrow 1 \rightarrow 4 \rightarrow 2 \rightarrow 1$ , we have  $\text{Graph}(P)$  equal to the graph in Fig. 3b. As we shall see in detail, a usual property about supporting graphs is that given any  $P \in \text{supp}(\mathcal{D}_{v,\ell})$ , its supporting graph  $\text{Graph}(P)$  is a subgraph of  $G$  (up to renaming of the vertices), where vertex 1 in  $\text{Graph}(P)$  is mapped to  $v$  in  $G$ .

**Path replacement.** Given any experiment  $P$ , we denote by  $\text{Replace}(P, u, P')$  the new experiment after replacing the last occurrence of integer  $u$  in  $P$  by the path  $P'$ . For instance

$$\text{Replace}(1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5 \rightarrow 4 \rightarrow 3, 4, 4 \rightarrow 6 \rightarrow 4) = 1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5 \rightarrow (4 \rightarrow 6 \rightarrow 4) \rightarrow 3$$

where the parentheses are for clarification purpose.

3. Theorem 1: a reconstructability result

In this section we show a positive result on reconstructing  $B(v, r)$  from random anonymous experiments of length  $\ell = O(m)$ , where  $m$  is the number of edges in  $B(v, r)$ .

**Theorem 1 (Restated).** Let  $n$  be the number of vertices in  $B(v, r)$  and  $m$  the number of edges.  $\text{Reconstruct}(v, r)$  (see Fig. 4) reconstructs  $B(v, r)$  with oracle accesses to  $(\mathcal{D}_1, \dots, \mathcal{D}_\ell)$ , where  $\ell = 2(m + 1)$ . More specifically,  $\text{Reconstruct}$  runs in time  $O(n^2)$ , and makes a total of  $O(n^2)$  membership queries to  $\text{supp}(\mathcal{D}_i)$  for  $i \in [\ell]$ .

3.1. An intuitive and non-constructive proof of Theorem 1

In this subsection, we show why Theorem 1 holds in a rather “non-constructive” way, that is, without worrying about the running time of the reconstruction algorithm. In the next subsection we prove Theorem 1, with the claimed running time of its reconstruction algorithm.

**The warm-up case: Reconstruction for  $r = 1$ .** Before proving the theorem, let us build the intuition by studying the special case of  $r = 1$ . Consider the following simple 2-step algorithm for reconstructing  $B(v, 1)$ .

\* (Throughout this section we slightly abuse the notation: for any experiment  $P$  of length no more than  $\ell$ , we use  $P \in \text{supp}(\mathcal{D}_\ell)$  to indicate the fact that  $P \in \text{supp}(\mathcal{D}_i)$  for some  $i \in [\ell]$ .)

1. In the first step we learn the degree of  $v$ . Let  $k \geq 1$  be the maximum integer such that the experiment

$$P = 1 \rightarrow 2 \rightarrow 1 \rightarrow 3 \rightarrow \dots \rightarrow k \rightarrow 1$$

is in  $\text{supp}(\mathcal{D}_\ell)$ . It is easy to show that vertex  $v$  has precisely  $k - 1$  neighbors in  $G$  according to the definition of  $k$ .

2. In the second step we learn the pairwise connections among the  $3 = k - 1$  neighbors of  $v$ . Letting  $P = 1 \rightarrow 2 \rightarrow 1 \rightarrow 3 \rightarrow 1 \rightarrow 4 \rightarrow 1$  be the walk studied in the first step, we proceed as follows.

- We first check if

$$P_1 \stackrel{\text{def}}{=} \text{Replace}(P, 2, 2 \rightarrow 3 \rightarrow 2) = 1 \rightarrow (2 \rightarrow 3 \rightarrow 2) \rightarrow 1 \rightarrow 3 \rightarrow 1 \rightarrow 4 \rightarrow 1$$

is in  $\text{supp}(\mathcal{D}_\ell)$ . If not, it indicates that there is no pairwise connection between any two neighbors of  $v$ , and the algorithm may terminate. Otherwise, there exists at least one pair of neighbors of  $v$  that are connected and the algorithm proceeds. Note that  $P_1$  indeed exists in  $\text{supp}(\mathcal{D}_\ell)$  for the graph of Fig. 3a, because  $v \rightarrow a \rightarrow c \rightarrow a \rightarrow v \rightarrow c \rightarrow v \rightarrow b \rightarrow v$  is such a walk.

- We then check if

$$P_2 \stackrel{\text{def}}{=} \text{Replace}(P_1, 2, 2 \rightarrow 4 \rightarrow 2) = 1 \rightarrow (2 \rightarrow 3 \rightarrow (2 \rightarrow 4 \rightarrow 2)) \rightarrow 1 \rightarrow 3 \rightarrow 1 \rightarrow 4 \rightarrow 1$$

is in  $\text{supp}(\mathcal{D}_\ell)$ . If not, it indicates that there does not exist a neighbor of  $v$  that is connected to two other neighbors, and the algorithm may terminate (in the case of  $k = 4$ ). Otherwise, like in Fig. 3a where  $v \rightarrow (c \rightarrow a \rightarrow (c \rightarrow b \rightarrow c)) \rightarrow v \rightarrow a \rightarrow v \rightarrow b \rightarrow v$  is such a walk, there exists a neighbor of  $v$  connected to two other neighbors, and the algorithm proceeds.

- We finally check if  $P_3 \stackrel{\text{def}}{=} \text{Replace}(P_2, 3, 3 \rightarrow 4 \rightarrow 3)$  is in  $\text{supp}(\mathcal{D}_\ell)$ . If not, like in Fig. 3a, we know the other two neighbors of  $v$  are not connected; otherwise they are connected. In both cases the algorithm may terminate here (in the case of  $k = 4$ ).

**Input:** Membership access to  $\text{supp}(\mathcal{D}_\ell)$ , a starting vertex  $v$  and a radius  $r$ .  
**Output:** A graph  $G'$  that is isomorphic to  $B(v, r)$ , and the isomorphism maps vertex 1 to  $v$ .

```

1:  $P \leftarrow 1$ .
2: for  $r_0 \leftarrow 1$  to  $r$  do
3:    $P_0 \leftarrow P; G_0 \leftarrow \text{Graph}(P)$ . ▷  $G_0$  is a reconstruction of  $B(v, r_0 - 1)$ 
4:    $n_0 \leftarrow$  the number of vertices in  $G_0$ .
5:    $u_1, \dots, u_k \leftarrow$  the vertices in  $G_0$  of distance precisely  $r_0 - 1$  from vertex 1. ▷  $u_i \in [n_0]$ 
6:   for  $i \leftarrow 1$  to  $k$  do
7:      $P' \leftarrow P_{i-1}$ .
8:      $x \leftarrow$  the smallest integer not appearing in  $P'$ .
9:     while  $\text{Replace}(P', u_i, u_i \rightarrow x \rightarrow u_i) \in \mathcal{D}_\ell$  do
10:       $P' \leftarrow \text{Replace}(P', u_i, u_i \rightarrow x \rightarrow u_i)$ .
11:      for all  $u' \in \{u_{i+1}, \dots, u_k\} \cup \{n_0 + 1, \dots, x - 1\}$  do
12:        if  $\text{Replace}(P', x, x \rightarrow u' \rightarrow x) \in \mathcal{D}_\ell$  then
13:           $P' \leftarrow \text{Replace}(P', x, x \rightarrow u' \rightarrow x)$ .
14:        end if
15:      end for
16:       $x \leftarrow$  the smallest integer not appearing in  $P'$ .
17:    end while
18:     $P_i \leftarrow P'$ .
19:  end for
20:   $P \leftarrow P_k$ .
21: end for
22: return  $\text{Graph}(P)$ .

```

Fig. 4.  $\text{Reconstruct}^{\mathcal{D}_\ell}(v, r)$ .

In the end of the algorithm, we output the supporting graph of the last experiment seen in  $\text{supp}(\mathcal{D}_\ell)$  by the above steps. In our example, this is  $\text{Graph}(P_2)$ , shown in Fig. 3b. Note that Fig. 3b is isomorphic to Fig. 3a and the isomorphism maps vertex 1 to vertex  $v$ , so is indeed a reconstruction of  $B(v, 1)$ . In this example, the longest experiment ever queried is  $P_3$ , of length  $12 = 2(m + 1) = \ell$ .

**The general case: Reconstruction for  $r > 1$ .** One can learn from the above warm-up case that, for any experiment  $P$  of length no more than  $\ell$ ,

- if  $P \in \text{supp}(\mathcal{D}_\ell)$ , then  $\text{Graph}(P)$  is a subgraph of  $G$  (up to renaming with 1 being mapped to  $v$  in  $G$ ), and conversely
- if  $\text{Graph}(P)$  is a subgraph of  $G$  (up to renaming with 1 being mapped to  $v$  in  $G$ ), then  $P \in \text{supp}(\mathcal{D}_\ell)$ .

We summarize this as

$$P \in \text{supp}(\mathcal{D}_\ell) \iff \text{Graph}(P) \text{ is a subgraph of } G \text{ (up to renaming with 1 mapped to } v \text{)}. \tag{3.1}$$

Therefore, one would hope to enumerate over all possible experiments  $P$  and use the information of whether  $P$  is in  $\text{supp}(\mathcal{D}_\ell)$  to reconstruct  $B(v, r)$ . Let us formalize this.

We call an experiment  $P$  *economical* if for any two integers  $a, b$  in the path, the segment  $a \rightarrow b$  appears at most once in  $P$ . All paths studied in the warm-up case are economical.

One can now study the following algorithm  $\text{NaiveReconstruct}$ . It enumerates over all valid experiments by the increasing order of their lengths, in order to find the longest experiment  $P^* \in \text{supp}(\mathcal{D}_\ell)$  such that

$$\text{both } P^* \text{ is economical and } \text{Graph}(P^*) \text{ is of radius } r \text{ from vertex } 1.$$

Owing to (3.1), this  $P^*$  satisfies that  $\text{Graph}(P^*)$  is isomorphic to  $B(v, r)$  and the isomorphism maps vertex 1 to vertex  $v$ . Since any economical experiment  $P$  of length  $2(m + 1)$  has at least  $m + 1$  edges in its supporting graph,  $\text{Graph}(P)$  cannot be a subgraph of  $B(v, r)$  and thus  $P \notin \text{supp}(\mathcal{D}_\ell)$ . This implies that  $\text{NaiveReconstruct}$  only needs oracle access to  $\text{supp}(\mathcal{D}_\ell)$  for  $\ell \leq 2(m + 1)$  in order to determine that  $P^*$  is the longest such experiment.

### 3.2. A constructive proof of Theorem 1

Although being sufficient for reconstructing  $B(v, r)$  given oracle access to  $\text{supp}(\mathcal{D}_\ell)$ ,  $\text{NaiveReconstruct}$  is still unsatisfactory because (1) the enumeration procedure is too slow and (2) the algorithm is not generalizable to the improvement

case studied in Section 1.4. We thus propose a more constructive algorithm *Reconstruct* that only makes  $O(n^2)$  membership queries to  $\text{supp}(\mathcal{D}_\ell)$ .

At a high level, *Reconstruct* builds  $B(v, r)$  by learning  $B(v, 1), \dots, B(v, r)$  layer by layer, and for each layer, by learning the vertices one by one. At any time of the algorithm, we maintain an economical experiment  $P$  whose supporting graph  $\text{Graph}(P)$  is a subgraph of  $B(v, r)$ . We incrementally “add” new vertices or edges to  $\text{Graph}(P)$ , verify if the new graph is still a subgraph of  $B(v, r)$  using (3.1), and if so, we update the current experiment  $P$  and continue. The details are as follows.

We describe *Reconstruct* in Fig. 4 and show its correctness by an induction on  $r$ . Suppose that we have reconstructed  $B(v, r_0 - 1)$  for some value  $r_0 - 1 \geq 0$ , and we now want to reconstruct  $B(v, r_0)$  using  $\mathcal{D}_\ell$  where  $\ell = 2(m + 1)$ .

Let  $n_0$  be the number of vertices in  $B(v, r_0 - 1)$ , and  $P_0$  an arbitrary experiment such that  $G_0 \stackrel{\text{def}}{=} \text{Graph}(P_0)$  is a reconstruction of  $B(v, r_0 - 1)$ .<sup>5</sup> We also denote by  $u_1, \dots, u_k \in [n_0]$  the vertices in  $G_0$  that have distance precisely  $r_0 - 1$  from vertex 1. We iterate over all  $i = 1, 2, \dots, k$ , and for each  $i$  we first let  $P' = P_{i-1}$  and repeatedly do the following (see Line 7 through 8 in Fig. 4).

Whenever  $\text{Replace}(P', u_i, u_i \rightarrow \star \rightarrow u_i)$  exists in  $\text{supp}(\mathcal{D}_\ell)$ , where  $\star$  is the smallest integer not appearing in  $P'$ , we know that there is at least one more vertex neighboring to  $u_i$  that is not explored so far, and we add it to  $P'$  by updating  $P' \leftarrow \text{Replace}(P', u_i, u_i \rightarrow \star \rightarrow u_i)$ . Equivalently, this update on  $P'$  can be understood as we are introducing a new vertex  $x$  along with a new edge  $(x, u_i)$  to  $\text{Graph}(P')$ .

As soon as a new vertex  $\star$  is added to  $P'$ , we add the edges connecting  $\star$  to other vertices in  $\text{Graph}(P')$  as follows. In principle,  $\star$  may be connected to any vertex in  $u' \in \{u_{i+1}, \dots, u_k\} \cup \{n_0 + 1, \dots, \star - 1\}$ , and we check them one by one. For each such a candidate neighbor  $u'$ , we check if  $\text{Replace}(P', \star, \star \rightarrow u' \rightarrow \star)$  exists in  $\text{supp}(\mathcal{D}_\ell)$ , and if so, we update  $P' \leftarrow \text{Replace}(P', \star, \star \rightarrow u' \rightarrow \star)$  and continue to the next  $u'$ . Equivalently, this update can be understood as we are adding an extra edge between  $x$  and  $u'$  into  $\text{Graph}(P')$ .

Let  $P_i$  be the final experiment  $P'$  after exploring all the vertices neighboring to  $u_i$ , and  $G_i = \text{Graph}(P_i)$ . We have, according to (3.1), that  $G_i$  is a subgraph of  $G$ . In fact, the last such subgraph  $G_k$  reconstructs  $B(v, r_0)$ :

**Claim 3.1.**  $G_k$  is isomorphic to  $B(v, r_0)$  and the isomorphism maps vertex 1 to vertex  $v$ .

**Proof.** First of all,  $G_k$  must be a subgraph of  $B(v, r_0)$  because  $P_k \in \text{supp}(\mathcal{D}_\ell)$  and, by construction, all vertices of  $G_k$  are within distance  $r_0$  from vertex 1. Therefore, we only need to verify if there is any vertex or edge in  $B(v, r_0)$  missing from  $G_k$ .

Let  $\sigma$  be an arbitrary embedding of  $G_k$  into  $B(v, r_0)$ , i.e., a mapping from the vertex set of  $G_k$  to that of  $B(v, r_0)$ , preserving edges, and mapping vertex 1 to vertex  $v$ .

For the missing vertex case, we prove by way of contradiction and suppose there is a vertex  $w$  in  $B(v, r_0) \setminus B(v, r_0 - 1)$  missing from  $G_k$  under this embedding  $\sigma$ . Because  $w$  is at distance  $r_0$  from  $v$ , it must be connected to some vertex at distance  $r_0 - 1$  from  $v$ . Let this vertex be  $\sigma(u_i)$  for some  $i \in [k]$ . (There must exist such a  $u_i$  because  $G_0$  reconstructs  $B(v, r_0 - 1)$  from the inductive step.)

Next, since  $w$  is missing from  $G_k$ , vertex  $u_i$  must have fewer neighbors in  $G_k$  than vertex  $\sigma(u_i)$  does in  $B(v, r_0)$ . At the time we finish constructing  $P_i$  (so the while loop in Line 9 from Fig. 4 terminates),  $G_i = \text{Graph}(P_i)$  can be embedded into  $G$  under the same  $\sigma$ . Letting  $\hat{P} = \text{Replace}(P_i, u_i, u_i \rightarrow x \rightarrow u_i)$ , the same embedding  $\sigma$ , while appended with  $\sigma(x) \mapsto w$ , should provide a valid embedding of  $\text{Graph}(\hat{P})$  into  $G$ , and according to (3.1) this implies  $\hat{P} \in \text{supp}(\mathcal{D}_\ell)$ . This contradicts the termination condition of the while loop in Line 9 that says  $\hat{P} \notin \text{supp}(\mathcal{D}_\ell)$ . Therefore there is no missing vertex.

One can perform a similar argument for the missing edge case.  $\square$

In sum, we have shown that  $B(v, r_0)$  can be constructed by the algorithm above, and by induction, *Reconstruct* outputs a reconstruction of  $B(v, r)$ . Notice that the experiment  $P$ , at the end of the algorithm, has a total length of  $2m$  because each edge in  $B(v, r)$  is traversed precisely once in each direction. Therefore the longest experiment *Reconstruct* has ever queried is of length  $2(m + 1)$ , and choosing  $\ell = 2(m + 1)$  is sufficient for our purpose. In addition, *Reconstruct* makes no more than  $O(n^2)$  membership queries to  $\text{supp}(\mathcal{D}_\ell)$ . ■

#### 4. Theorem 2: a lower bound on experiment length

In this section, for any integer  $h \geq 1$ , we construct two (infinite) binary trees  $T_1 = T_1^{(h)}$  and  $T_2 = T_2^{(h)}$  with the starting vertex being the root for both cases. We show, quite surprisingly, although  $T_1$  and  $T_2$  are different at depth  $r = 2h + 3$ , any anonymous experiment of length no longer than  $\ell = O(2^h)$  has the same probability to be generated from  $T_1$  and  $T_2$ . Formally,

**Lemma 4.1.** *There exists a constant  $c$  such that, given two binary trees  $T_1 = T_1^{(h)}$  and  $T_2 = T_2^{(h)}$  (as constructed in Fig. 5), and letting the starting vertex  $v_1$  and  $v_2$  be their roots, we have:*

- $B^{T_1}(v_1, 2h + 3)$  and  $B^{T_2}(v_2, 2h + 3)$  are different (i.e., non-isomorphic), but
- the distributions over random experiments of length  $\ell \leq c \cdot 2^h$  in  $T_1$  and  $T_2$  are the same.

<sup>5</sup> I.e.,  $P_0$  satisfies that  $\text{Graph}(P_0)$  is isomorphic to  $B(v, r_0 - 1)$  and the isomorphism maps vertex 1 to  $v$ . In fact,  $P_0$  is inherited from the inductive step of the algorithm, and corresponds to an arbitrary walk that starts from  $v$  and traverses each edge in  $B(v, r_0 - 1)$  exactly once in each direction.



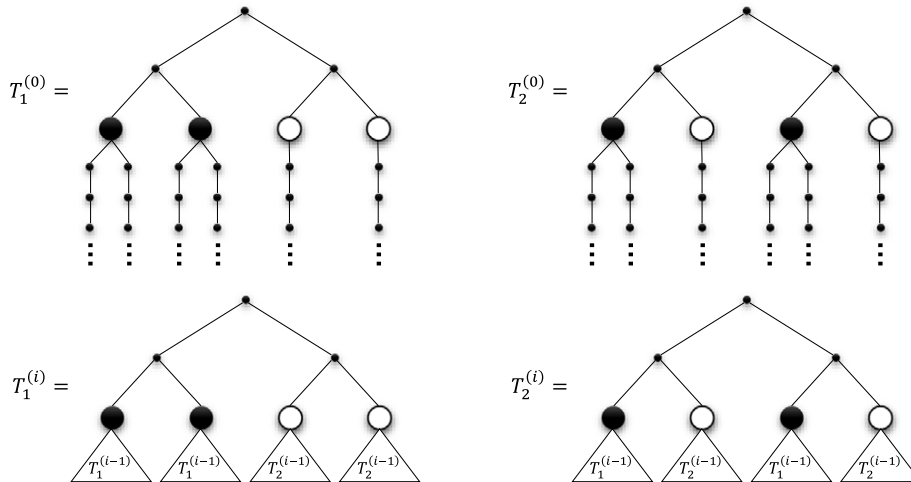


Fig. 5. The recursive definition of the hard instance for Theorem 2.

Theorem 2 is immediately implied by the above lemma, because it rules out the possibility of reconstructing  $B(v, r)$ , even for binary trees, with oracle access to  $(\mathcal{D}_1, \dots, \mathcal{D}_\ell)$  for any  $\ell = 2^{o(r)}$ .

4.1. Our hard instance

We define  $T_1 = T_1^{(h)}$  and  $T_2 = T_2^{(h)}$  recursively.

Let  $T_1^{(0)}$  and  $T_2^{(0)}$  be defined as follows (see Fig. 5): the roots of both trees have two children and each child in turns has two children; among the four grandchildren of the root, two of them are “black”, having two infinite chains of descendants, and two of them are “white”, having one infinite chain of descendants.

$T_1^{(i)}$  and  $T_2^{(i)}$  are defined similarly (see Fig. 5): the roots of both trees have two children and each child in turns has two children; among the four grandchildren of the root, two of them are “black”, having  $T_1^{(i-1)}$  as subtrees, and two of them are “white”, having  $T_2^{(i-1)}$  as subtrees.

4.2. A warm-up property

For  $j \in \{1, 2\}$ , let  $\mathcal{D}_{j,\ell}$  be the distribution over random experiments of length  $\ell$  generated from the Markov process starting from the root of  $T_j$ . Given an experiment  $P$  of length  $\ell$ , we denote by  $\Pr[P \mid T_j]$  the probability that  $P$  is generated from  $\mathcal{D}_{j,\ell}$ .

Recall that one can associate  $P$  with its supporting graph  $G_P = \text{Graph}(P)$ . Since  $T_1$  and  $T_2$  are binary trees, if the supporting graph  $G_P$  has cycles or is non-binary,  $P$  cannot exist in  $\mathcal{D}_{j,\ell}$ . We thus focus only on the experiments  $P$  for which  $G_P$  is a binary tree. We make the following claim:

**Claim 4.2.** *If the root (i.e., vertex 1) of  $G_P$  has at most one grandchild, then  $\Pr[P \mid T_1] = \Pr[P \mid T_2]$ .*

Before proving Claim 4.2, we summarize the high level intuition as follows.

Any experiment  $P$  is consistent with a set of walks  $\mathcal{Q}_1$  on  $T_1$ , and a set of walks  $\mathcal{Q}_2$  on  $T_2$ . The probability  $\Pr[P \mid T_j]$  is equal to  $\sum_{Q \in \mathcal{Q}_j} \Pr[Q \mid T_j]$ , the sum of probabilities over the walks in  $\mathcal{Q}_j$ , i.e., those walks consistent with  $P$ . We show that, under the condition  $P$  visits only one grandchild of the root, there is a one-to-one mapping  $\tau$  between  $\mathcal{Q}_1$  and  $\mathcal{Q}_2$  that preserves probabilities. This immediately implies that  $\Pr[P \mid T_1] = \Pr[P \mid T_2]$ . The one-to-one mapping  $\tau$  is illustrated in Fig. 6, and note that if  $P$  visits two grandchildren such a mapping may not exist.

**Proof of Claim 4.2.** We prove the claim when the root has only one grandchild in  $G_P$ . The other case – when the root has no grandchild – is only simpler. We denote by  $u \in \mathbb{Z}_+$  this unique grandchild, and focus on the case of  $h = 0$ ; the case of  $h > 0$  is similar.

Let the four grandchildren of the root in  $T_1^{(0)}$  be denoted by  $a_1, a_2, a_3, a_4$  respectively, and the four grandchildren of the root in  $T_2^{(0)}$  be denoted by  $b_1, b_2, b_3, b_4$ . We order them according to Fig. 6 so  $a_1, a_2, b_1, b_3$  are black, and  $a_3, a_4, b_2, b_4$  are white.

We now construct a one-to-one mapping  $\tau$  from the walks on  $T_1^{(0)}$  that are consistent with  $P$  to the walks on  $T_2^{(0)}$  that are consistent with  $P$ . Our  $\tau$  is defined “by picture”, with four representative examples given in Fig. 6.

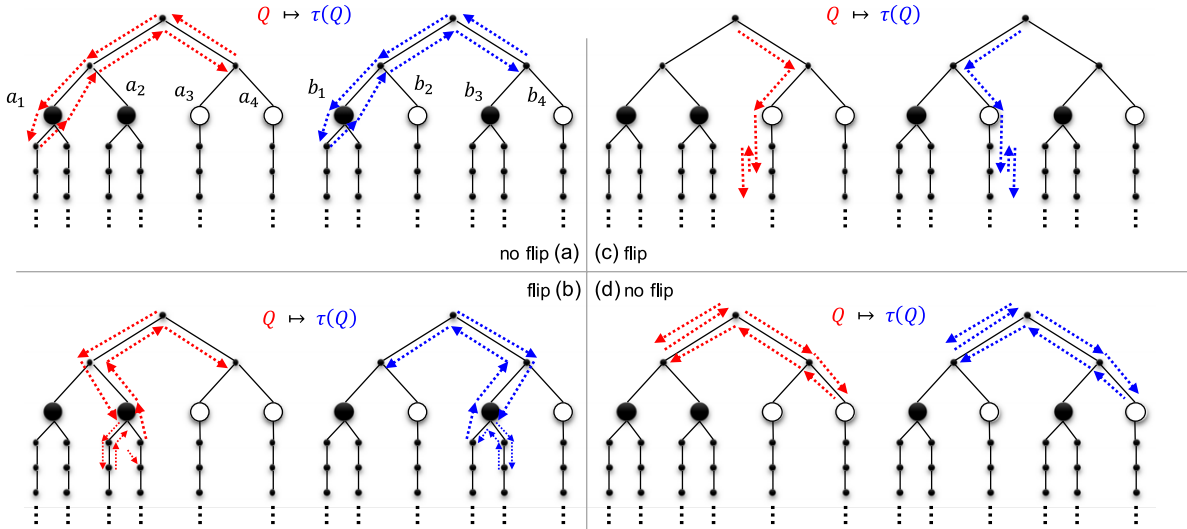


Fig. 6. The illustration of our mapping  $\tau$  used in Claim 4.2.

More precisely, to define  $\tau$ , we first draw  $T_1^{(0)}$  and  $T_2^{(0)}$  on the plane with four grandchildren of the root sorted as  $a_1, a_2, a_3, a_4$  and  $b_1, b_2, b_3, b_4$  from left to right. Then, given a walk  $Q$  on  $T_1^{(0)}$  (starting from the root) that is consistent with  $P$ , denoted as  $Q \triangleleft (T_1^{(0)}, P)$ , vertex  $u$  in  $P$  must be mapped to one of  $\{a_1, a_2, a_3, a_4\}$ .

- If  $u$  is mapped to  $a_2$  or  $a_3$  in  $Q$  (see Fig. 6(b) and (c)), we let  $Q' = \tau(Q)$  be the walk in  $T_2^{(0)}$  that is *flipped* left and right (with respect to the plane), and thus  $u$  is mapped to  $b_3$  or  $b_2$  respectively in  $Q'$ .
- If  $u$  is mapped to  $a_1$  or  $a_4$  in  $Q$  (see Fig. 6(a) and (d)), we let  $Q' = \tau(Q)$  be the “same” walk  $Q$  under translation on the plane, and thus  $u$  is mapped to  $b_1$  or  $b_4$  respectively in  $Q'$ .

It is not hard to verify that  $\tau$  is a one-to-one mapping. In addition, the  $i$ th vertex in  $Q$  has the same degree as the  $i$ th vertex in  $Q' = \tau(Q)$  for any  $i$  and any  $Q$  satisfying  $Q \triangleleft (T_1^{(0)}, P)$ . Therefore we have  $\Pr[Q \mid T_1^{(0)}] = \Pr[Q' \mid T_2^{(0)}]$ , i.e.,  $Q$  and  $Q'$  have the same probability to be generated in the random walk from  $T_1^{(0)}$  and  $T_2^{(0)}$  respectively. This implies

$$\begin{aligned} \Pr[P \mid T_1^{(0)}] &= \sum_{Q \triangleleft (T_1^{(0)}, P)} \Pr[Q \mid T_1^{(0)}] = \sum_{Q \triangleleft (T_1^{(0)}, P)} \Pr[\tau(Q) \mid T_2^{(0)}] \\ &= \sum_{Q' \triangleleft (T_2^{(0)}, P)} \Pr[Q' \mid T_2^{(0)}] = \Pr[P \mid T_2^{(0)}], \end{aligned}$$

that is,  $P$  has the same chance to be generated as an experiment in  $T_1^{(0)}$  and  $T_2^{(0)}$ .  $\square$

### 4.3. A general property

For any  $i \in \{0, 1, \dots, 2h\}$ , we denote by  $L_i$  the set of vertices (in the form of integer numbers) in  $G_P = \text{Graph}(P)$  at depth  $i$  from the root (where the root itself is in  $L_0$ ). We prove the following property about a shortest experiment in which  $\Pr[P \mid T_1] \neq \Pr[P \mid T_2]$ .

**Lemma 4.3.** *Given a shortest experiment  $P$  in which  $\Pr[P \mid T_1] \neq \Pr[P \mid T_2]$ , any  $i \in \{0, 1, \dots, h\}$ , and any  $u \in L_{2i}$ , vertex  $u$  has at least two grandchildren in  $G_P$ .*

Notice that the case of  $i = 0$  is a direct consequence of Claim 4.2, but the proof for the  $i \geq 1$  case is more involved. Before proving it formally, we summarize the basic idea as follows.

If  $P$  is a shortest such experiment, and if there exists some  $u$  in  $P$  with only one grandchild, we shorten  $P$  to a new experiment  $P'$  by essentially removing all occurrences of  $u$  and the descendants of  $u$ . In a rough sense,  $\Pr[P \mid T_j]$  equals to  $\Pr[P' \mid T_j] \times \Pr[P \setminus P' \mid T_j]$  where  $P \setminus P'$  is an experiment corresponding to the removed segment of vertices. Because  $u$  has only one grandchild in  $G_P$ , this removed subsegment  $P \setminus P'$  has the same probability to be generated in  $T_1$  and  $T_2$  (owing to Claim 4.2). We therefore conclude that  $\Pr[P' \mid T_1] \neq \Pr[P' \mid T_2]$ , contradicting to the fact that  $P$  is the shortest such experiment.

**Proof of Lemma 4.3.** The case of  $i = 0$  is inherited from Claim 4.2, so the rest of this section is devoted to proving Lemma 4.3 for  $i \geq 1$ .

For  $j \in \{1, 2\}$ , let  $\mathcal{D}_{j,\ell}$  be the distribution over random experiments of length  $\ell$  in tree  $T_j$ , and  $\mathcal{D}_{j,\ell}^{\text{walk}}$  the distribution over random walks in tree  $T_j$ . We make a quick observation first.

Given an experiment  $P$ , the probability  $\Pr[P \mid T_j]$  is the sum of the probabilities  $\Pr[\sigma(P) \mid T_j]$  over all choices of embeddings  $\sigma : G_P \rightarrow T_j$ :

$$\Pr[P \mid T_j] = \sum_{\text{embedding } \sigma : G_P \rightarrow T_j} \Pr[\sigma(P) \mid T_j]. \tag{4.1}$$

Here an embedding  $\sigma$  is a mapping from the vertices in  $G_P$  to the vertices in  $T_j$ , while preserving edges and mapping vertex 1 to vertex  $v$ . Accordingly,  $\sigma$  maps an experiment  $P$  to an actual walk  $\sigma(P)$  on  $T_j$ , and  $\Pr[\sigma(P) \mid T_j]$  is the probability for  $\sigma(P)$  to be generated from  $\mathcal{D}_{j,\ell}^{\text{walk}}$ . We also recall a useful fact by the definition of random walk:

$$\Pr[\sigma(P) \mid T_j] = \prod_{i=1}^{\ell} \frac{1}{\deg(\sigma(P^{(i)}))}, \tag{4.2}$$

where  $P^{(i)}$  is the  $i$ th integer in the experiment  $P$ , and thus  $\deg(\sigma(P^{(i)}))$  is the degree of the  $i$ th vertex in the length- $\ell$  walk  $\sigma(P)$ .

We are now ready to prove Lemma 4.3. Suppose that Lemma 4.3 does not hold for some  $i \in \{1, \dots, h\}$ , and vertex  $u \in L_{2i}$  has only one grandchild in  $G_P$ , we will show that one can shorten  $P$  to construct a new experiment  $P'$  where it also satisfies  $\Pr[P' \mid T_1] \neq \Pr[P' \mid T_2]$ , contradicting the fact that  $P$  is the shortest such experiment. In order to shorten  $P$ , we first discover that  $P$  must be of some special structure, described as follows.

We note that  $P$  can be viewed as a “walk” on its supporting graph  $G_P = \text{Graph}(P)$ , and let the  $w$  be parent of  $u$  in  $G_P$ . Clearly,  $P$  must visit  $w$  before it visits  $u$  in this walk, but we claim that  $P$  can only be one of the two forms:

- either it enters the subtree rooted at  $u$ , then comes back to  $w$  and never visits  $u$  again;
- or it enters the subtree rooted at  $u$  and never comes back to  $w$ .

Formally,

**Claim 4.4.**  $P$  must be of the form:

$$P = P_1 \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow w \rightarrow P_3 \quad \text{or} \quad P = P_1 \rightarrow w \rightarrow u \rightarrow P_2$$

where  $P_2$  consists of only vertices that are  $u$  or descendants of  $u$  (in  $G_P$ ), while  $P_1$  and  $P_3$  consist of only vertices that are neither  $u$  nor descendants of  $u$  (in  $G_P$ ).

**Proof.** Suppose that  $P$  is of neither of the two forms above, then  $P$  must visit some descendants of  $u$  first, then non-descendants, and then descendants again. For instance, such a walk could be

$$P = P_1 \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow w \rightarrow P_3 \rightarrow w \rightarrow u \rightarrow P'_2$$

where  $P_2$  and  $P'_2$  consist of only  $u$  or descendants of  $u$ , while  $P_1$  and  $P_3$  consist of only vertices that are neither  $u$  nor descendants of  $u$ . We only prove the claim for this case above, and other cases are similar.

We first swap the order of the vertices in  $P$  and construct the following experiment  $P'$ :

$$P' = P_1 \rightarrow w \rightarrow P_3 \rightarrow w \rightarrow u \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow P'_2.$$

Since for any two integers  $a$  and  $b$  the directed edge  $a \rightarrow b$  appears exactly the same number of times in  $P$  and  $P'$ , we have that  $\Pr[P \mid T_j] = \Pr[P' \mid T_j]$  according to (4.1) and (4.2), for both  $j = 1$  and  $2$ .

We next observe that the subsequence  $w \rightarrow u \rightarrow w \rightarrow u$  is redundant: since  $u$  and  $w$  are of depth  $2i$  and  $2i - 1$  respectively, they will always be mapped to vertices with degree 3 in  $T_1$  or  $T_2$ . As a result, if we define

$$P'' = P_1 \rightarrow w \rightarrow P_3 \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow P'_2$$

we must have  $\Pr[P' \mid T_j] = (\frac{1}{3})^2 \Pr[P'' \mid T_j]$  for both  $j = 1$  and  $j = 2$ , according to (4.1) and (4.2) again. This indicates  $\Pr[P' \mid T_1] \neq \Pr[P' \mid T_2]$ , contradicting the choice of  $P$  which is the shortest experiment that makes  $\Pr[P \mid T_1] \neq \Pr[P \mid T_2]$ .  $\square$

Now we focus only on the case of  $P = P_1 \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow w \rightarrow P_3$  because the other one is only simpler. We want to shorten it to  $P_1 \rightarrow w \rightarrow P_3$ .

**Claim 4.5.** If an experiment  $P = P_1 \rightarrow w \rightarrow u \rightarrow P_2 \rightarrow u \rightarrow w \rightarrow P_3$  satisfies  $\Pr[P \mid T_1] \neq \Pr[P \mid T_2]$  (where the definitions of  $P_1, P_2$  and  $P_3$  are the same as Claim 4.4), and  $u$  has at most one grandchild in  $\text{Graph}(P)$ , then we have

$$\Pr[P_1 \rightarrow w \rightarrow P_3 \mid T_1] \neq \Pr[P_1 \rightarrow w \rightarrow P_3 \mid T_2].$$

**Proof.** To see this, we consult (4.1) again. For each  $j \in \{1, 2\}$ , and given an embedding  $\sigma : G_P \rightarrow T_j$ , we write it as a pair  $\sigma = (\sigma_1, \sigma_2)$  where

- $\sigma_1$  maps all the vertices excluding the descendants of  $u$  (so including  $u$ ) to  $T_j$ ,
- $\sigma_2$  maps all the descendants of  $u$  (including  $u$ ) to  $T_j$ , and
- $\sigma_1$  and  $\sigma_2$  map  $u$  to the same vertex in  $T_j$ . (When this is the case, we denote by  $\sigma_1 \sim \sigma_2$ .)

We therefore rewrite:

$$\Pr[P \mid T_j] = \sum_{\sigma} \Pr[\sigma(P) \mid T_j] = \sum_{\sigma_1} \sum_{\sigma_2: \sigma_1 \sim \sigma_2} \Pr[(\sigma_1, \sigma_2)(P) \mid T_j]. \tag{4.3}$$

Next, recall that  $Q = \sigma(P) = (\sigma_1, \sigma_2)(P)$  is a walk on the tree  $T_j$ , and  $\Pr[Q \mid T_j]$  can be written as a product of the reciprocal of degrees, i.e.,  $\Pr[Q \mid T_j] = \prod_{i=1}^{\ell} \frac{1}{\deg(Q^{(i)})}$  in which  $\deg(Q^{(i)})$  is the degree of the  $i$ th vertex in the walk  $Q$ . This allows us to break  $Q$  into five segments:  $\sigma_1(P_1 \rightarrow w)$ ,  $\sigma_1(w \rightarrow u)$ ,  $\sigma_2(u \rightarrow P_2 \rightarrow u)$ ,  $\sigma_1(u \rightarrow w)$ , and  $\sigma_1(w \rightarrow P_3)$ , and compute

$$\begin{aligned} \Pr[(\sigma_1, \sigma_2)(P) \mid T_j] &= \Pr[\sigma_1(P_1 \rightarrow w) \mid T_j] \cdot \Pr[\sigma_1(w \rightarrow u) \mid T_j] \\ &\quad \cdot \Pr[\sigma_2(u \rightarrow P_2 \rightarrow u) \mid T_j] \cdot \Pr[\sigma_1(u \rightarrow w) \mid T_j] \cdot \Pr[\sigma_1(w \rightarrow P_3) \mid T_j]. \end{aligned}$$

We reorder them into four segments  $\sigma_1(P_1 \rightarrow w \rightarrow P_3)$ ,  $\sigma_1(w \rightarrow u)$ ,  $\sigma_1(u \rightarrow w)$ ,  $\sigma_2(u \rightarrow P_2 \rightarrow u)$ , and conclude that

$$\Pr[(\sigma_1, \sigma_2)(P) \mid T_j] = \Pr[\sigma_1(P_1 \rightarrow w \rightarrow P_3) \mid T_j] \cdot \Pr[\sigma_1(w \rightarrow u) \mid T_j] \cdot \Pr[\sigma_1(u \rightarrow w) \mid T_j] \cdot \Pr[\sigma_2(u \rightarrow P_2 \rightarrow u) \mid T_j].$$

However, we must have  $\Pr[\sigma_1(w \rightarrow u) \mid T_j] = \Pr[\sigma_1(u \rightarrow w) \mid T_j] = 1/3$  because any embedding  $\sigma = (\sigma_1, \sigma_2)$  maps  $u$  and  $w$  to vertices with degree 3. This, combined with (4.3) gives us

$$\begin{aligned} \Pr[P \mid T_j] &= \sum_{\sigma_1} \sum_{\sigma_2: \sigma_1 \sim \sigma_2} \Pr[\sigma_1(P_1 \rightarrow w \rightarrow P_3) \mid T_j] \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \Pr[\sigma_2(u \rightarrow P_2 \rightarrow u) \mid T_j] \\ &= \frac{1}{9} \sum_{\sigma_1} \Pr[\sigma_1(P_1 \rightarrow w \rightarrow P_3) \mid T_j] \cdot \sum_{\sigma_2: \sigma_1 \sim \sigma_2} \Pr[\sigma_2(u \rightarrow P_2 \rightarrow u) \mid T_j]. \end{aligned}$$

Now, fixing any  $\sigma_1$ , we know that  $u$  is mapped to vertex  $\sigma_1(u)$  in  $T_j$ , and  $\sigma_1(u)$  must be the root of some  $T_k^{(h-i)}$  tree for  $k \in \{1, 2\}$ . Here the value of  $k$  depends on the choice of  $\sigma_1$ . We observe that the summation

$$\sum_{\sigma_2: \sigma_1 \sim \sigma_2} \Pr[\sigma_2(u \rightarrow P_2 \rightarrow u) \mid T_j]$$

is precisely the probability for the experiment  $u \rightarrow P_2 \rightarrow u$  (after renaming so that the integers are 1-based) to be generated in  $T_k^{(h-i)}$ , and this value does not depend on the choice of  $k$  owing to Claim 4.2 and the fact that  $u$  has at most one grandchild in  $P_2$ . Let this value be  $p \in [0, 1]$ , and we conclude that

$$\Pr[P \mid T_j] = \frac{1}{9} \sum_{\sigma_1} \Pr[\sigma_1(P_1 \rightarrow w \rightarrow P_3) \mid T_j] \cdot p = \frac{p}{9} \cdot \Pr[P_1 \rightarrow w \rightarrow P_3 \mid T_j],$$

that is, the value of  $\Pr[P \mid T_j]$  is a fixed constant  $\frac{p}{9}$  multiplied by that of a shorter experiment  $P_1 \rightarrow w \rightarrow P_3$  on the same tree  $T_j$ . Since this is true for both  $j \in \{1, 2\}$ , we conclude that  $\Pr[P_1 \rightarrow w \rightarrow P_3 \mid T_1] \neq \Pr[P_1 \rightarrow w \rightarrow P_3 \mid T_2]$ .  $\square$

Since the above claim contradicts the choice of  $P$  which is the shortest such sequence that makes  $\Pr[P \mid T_1] \neq \Pr[P \mid T_2]$ , we finish the proof of Lemma 4.3.  $\blacksquare$

#### 4.4. Proof of Lemma 4.1

**Proof.** It is immediate that Lemma 4.3 implies Lemma 4.1: the shortest experiment  $P$  that distinguishes  $\Pr[P \mid T_1]$  and  $\Pr[P \mid T_2]$  must branch out at least once for every two levels, and therefore  $|L_{2i}| \geq 2^i$  and in particular  $L_{2(h+1)} \geq 2^{h+1}$ . This shows that the length of  $P$  must be at least  $\Omega(2^h)$  (in order to visit  $2^{h+1}$  distinct vertices at depth  $2(h+1)$ ). In other words, there exists some constant  $c$  where  $\Pr[P \mid T_1] = \Pr[P \mid T_2]$  for any experiment of length  $\ell \leq c \cdot 2^h$ .  $\blacksquare$

### 5. Theorem 3: a lower bound on the number of experiments

#### 5.1. Our new hard instance

We slightly modify our hard instance in Fig. 5, by replacing the definitions of  $T_1^{(0)}$  and  $T_2^{(0)}$  with Fig. 7: instead of having a black vertex to be the root of two infinite chains and the white vertex to be the root of one (recall Fig. 5), we let a black vertex be the parent of three infinite complete binary trees, and the white one be the parent of two. The new trees  $T_1 = T_1^{(h)}$

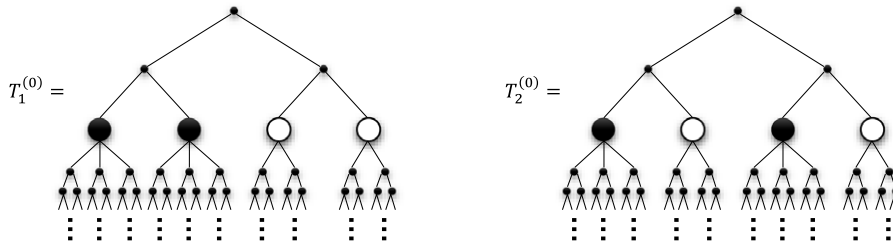


Fig. 7. The new choices of  $T_1^{(0)}$  and  $T_2^{(0)}$  for Theorem 3.

and  $T_2 = T_2^{(h)}$  so constructed are ternary. (In fact, one can verify that the binary variants of  $T_1$  and  $T_2$  also suffice for the purpose of Theorem 3, but they will make the analysis more involved.)

It is a simple exercise to verify that all the proofs in Section 4 remain true for this new hard instance pair  $(T_1, T_2)$ , and therefore Lemma 4.1 still applies: that is, there exists a constant  $c$  such that, letting the starting vertex  $v_1$  and  $v_2$  be the corresponding roots, we have:

- $B^{T_1}(v_1, 2h + 3)$  and  $B^{T_2}(v_2, 2h + 3)$  are different (i.e., non-isomorphic), but
- the distributions over random experiments of length  $\ell \leq c \cdot 2^h$  in  $T_1$  and  $T_2$  are the same.

This new construction satisfies an additional property:

- in a random walk on either  $T_1$  or  $T_2$ , if the current vertex is at depth  $d$  for some  $d \in \mathbb{Z}_{\geq 0}$ , then with probability at least  $2/3$ , the next vertex is going to be at depth  $d + 1$ .

### 5.2. A structural lemma

We say that an experiment  $P$  is *bad*, if it has less than  $2^{h+1}$  vertices at depth  $2(h + 1)$  in its supporting graph  $G_P = \text{Graph}(P)$ . We denote by  $\mathcal{BAD}$  the set of bad experiments. According to the proof of Lemma 4.1, any bad experiment  $P \in \mathcal{BAD}$  has the same chance to be seen in  $T_1$  and  $T_2$ , that is,  $\Pr[P \mid T_1] = \Pr[P \mid T_2]$ .

We now compute a lower bound on the chance of a random experiment to be bad.

**Lemma 5.1.** *For any  $j \in \{1, 2\}$ , and any value of  $\ell$ , with probability at least  $1 - e^{-\Omega(2^h)}$ , the random experiment of length  $\ell$  generated from  $T_j$  is bad.*

The proof of Lemma 5.1 mostly consists of careful applications of Chernoff and union bounds, and we summarize its intuition as follows.

By our construction of the trees, any random walk (on either  $T_1$  or  $T_2$ ) of length  $t$  is likely to arrive at a vertex at depth  $\Omega(t)$ . This is because, in each step of the random walk, the depth increases by 1 with probability at least  $2/3$ , and decreases by 1 with probability at most  $1/3$ . More precisely, by Chernoff bound, the random walk will land at a vertex of depth  $\Omega(2^h)$  after  $t = 2^h$  steps, with probability at least  $1 - e^{-\Omega(2^h)}$ . Since  $2^h < 2^{h+1}$ , in order for this random walk to correspond to a good experiment, it has to come back from depth  $\Omega(2^h)$  to depth  $2(h + 1)$  in order to visit  $2^{h+1}$  vertices at that depth. This, again using Chernoff bound, is a very unlikely event, because going back from depth  $\Omega(2^h)$  to depth  $2(h + 1)$  has a probability at most  $e^{-\Omega(2^h)}$ , no matter how long the random walk is. It is crucial here that the probability does not depend on  $\ell$ .

**Proof of Lemma 5.1.** It suffices to prove the lemma for  $\ell \geq 2^{h+1}$ , because otherwise any experiment  $P$  of length  $\ell$  cannot visit  $2^{h+1}$  vertices at depth  $2(h + 1)$  and is by definition bad.

Let  $\text{dep}_i \in \mathbb{Z}_{\geq 0}$  be the random variable indicating the depth of the  $i$ th vertex in the random walk on  $T_j$ , where  $i \in \{0, 1, \dots, \ell\}$ . We have  $\text{dep}_0 = 0$ . Let the random variable  $x_i$  be defined as  $\frac{\text{dep}_i - \text{dep}_{i-1} + 1}{2} \in \{0, 1\}$ , so that  $\text{dep}_i = \text{dep}_{i-1} + (-1 + 2x_i) = \text{dep}_{i-1} \pm 1$ . By the construction of our graph (either  $T_1$  or  $T_2$ ), we always have  $\mathbb{E}[x_i] \geq \frac{2}{3}$ , that is, with probability at least  $\frac{2}{3}$  the depth increases by 1 in a step.

Let us consider a special timestamp in the random walk: time  $t = 2^h$ . Using Chernoff bound, we deduce below that with very high probability (i.e.,  $1 - e^{-\Omega(t)}$ ), we have that  $\text{dep}_t \geq \frac{1}{6}t = \frac{2^h}{6}$ :

$$\begin{aligned} \Pr \left[ \text{dep}_t < \frac{1}{6}t \right] &= \Pr \left[ -t + 2(x_1 + \dots + x_t) < \frac{1}{6}t \right] = \Pr \left[ x_1 + \dots + x_t < \frac{7}{12}t \right] \\ &\leq \Pr \left[ x_1 + \dots + x_t < \frac{7}{8}\mathbb{E}[x_1 + \dots + x_t] \right] \leq e^{-\Omega(t)}. \end{aligned}$$

Recall that within  $t = 2^h$  steps the random walk cannot visit enough (i.e., at least  $2^{h+1}$ ) vertices at depth  $2(h + 1)$ , so for a random walk of length  $\ell$  to correspond to a good experiment, it must come back from depth  $\text{dep}_t$  to depth  $2(h + 1)$  in the remaining  $\ell - t$  steps.

Conditioning on that  $\text{dep}_t \geq \frac{1}{6}t = \frac{2^h}{6}$ , we compute the chance of the random experiment to reach back to a vertex at depth  $\leq 2(h+1)$  at time  $t+t'$  where  $t' \in \{1, \dots, \ell-t\}$ .

$$\begin{aligned} \Pr[\text{dep}_{t+t'} \leq 2(h+1)] &\leq \Pr\left[\text{dep}_{t+t'} - \text{dep}_t \leq 2(h+1) - \frac{2^h}{6}\right] \\ &= \Pr\left[-t' + 2(x_{t+1} + \dots + x_{t+t'}) \leq 2(h+1) - \frac{2^h}{6}\right] \\ &= \Pr\left[x_{t+1} + \dots + x_{t+t'} \leq \frac{t'}{2} + (h+1) - \frac{2^h}{12}\right]. \end{aligned} \tag{5.1}$$

We assume that  $h$  is sufficiently large (e.g.  $h \geq 8$ ) so that  $\frac{t'}{2} + (h+1) - \frac{2^h}{12} \leq \frac{t'}{2} - \frac{2^h}{24}$ . Then obviously  $t'$  has to be at least  $\frac{2^h}{12}$  before this probability in (5.1) becomes non-zero. Therefore, we only focus on the choices of  $t' \geq \frac{2^h}{12}$  and continue our calculation using Chernoff bound:

$$\begin{aligned} \Pr[\text{dep}_{t+t'} \leq 2(h+1)] &\leq \Pr\left[x_{t+1} + \dots + x_{t+t'} \leq \frac{t'}{2} - \frac{2^h}{24}\right] \\ &\leq \Pr\left[x_{t+1} + \dots + x_{t+t'} \leq \frac{t'}{2}\right] \\ &\leq \Pr\left[x_{t+1} + \dots + x_{t+t'} \leq \frac{3}{4}\mathbb{E}[x_{t+1} + \dots + x_{t+t'}]\right] \leq e^{-\Omega(t')}. \end{aligned}$$

Finally, since we only need to focus on  $t' \geq \frac{2^h}{12}$  due to the discussed reason, by a union bound over all integers  $t' \in [\frac{2^h}{12}, \ell-t]$ , we have that the chance for a random experiment to visit back to depth  $2(h+1)$  is at most  $\sum_{t'=\frac{2^h}{12}}^{\ell-t} e^{-\Omega(t')} = e^{-\Omega(t)}$ .

In sum, we know that with probability at least  $1 - e^{-\Omega(t)} = 1 - e^{-\Omega(2^h)}$ , the random walk generated (from either  $T_1$  or  $T_2$ ) will: (1) have  $\text{dep}_t \geq \frac{2^h}{12}$  and (2) never visit back to depth  $2(h+1)$ . The experiment corresponding to this walk has to be bad, and therefore we finish the proof.  $\square$

### 5.3. Proof of Theorem 3

We argue that in order to distinguish  $T_1 = T_1^{(h)}$  from  $T_2 = T_2^{(h)}$  with probability at least  $\frac{1}{2}$ , one needs at least  $e^{\Omega(2^h)}$  samples of random experiments of arbitrary lengths.

Indeed, let  $\mathcal{D}_{1,\ell}$  be the distribution over random experiments of length  $\ell$  for tree  $T_1$ , and  $\mathcal{D}_{2,\ell}$  that for  $T_2$ . By definition,  $\mathcal{D}_{1,\ell}$  and  $\mathcal{D}_{2,\ell}$  are identical on the support of  $\mathcal{BAD}$ , the set of bad experiments. Therefore, owing to Lemma 5.1, the total variation distance (i.e., half of the 1-norm distance) between them  $\|\mathcal{D}_{1,\ell} - \mathcal{D}_{2,\ell}\|_{TV}$  is at most  $e^{-\Omega(2^h)}$  for any  $\ell$ ; that is, any algorithm that samples an experiments from  $\mathcal{D}_{1,\ell}$  or  $\mathcal{D}_{2,\ell}$ , can only tell the difference with probability at most  $e^{-\Omega(2^h)}$ .

Using union bound, given any algorithm that takes samples from distributions  $(\mathcal{D}_{j,1}, \mathcal{D}_{j,2}, \dots)$ , unless it takes  $e^{\Omega(2^h)}$  samples, it cannot distinguish  $T_1$  from  $T_2$  with any constant probability.

Let  $A$  be an algorithm that reconstructs  $B(v, r)$  – even only for the case when the underlying graph is a ternary tree – with probability  $1/2$  using  $N$  random experiments. If  $N = 2^{2^{\Omega(r)}}$ , then using  $A$  one can reconstruct  $B(v, 2h+3)$  for  $T_1$  and  $T_2$  respectively, and thus distinguish  $T_1$  from  $T_2$ . This leads to a contradiction because no algorithm can distinguish  $T_1$  from  $T_2$  using only  $e^{\Omega(2^h)}$  samples; therefore we must have  $N = 2^{2^{\Omega(r)}}$ .  $\blacksquare$

## References

- [1] Noga Alon, Vera Asodi, Learning a hidden subgraph, *SIAM J. Discret. Math.* 18 (4) (2005) 697–712.
- [2] Noga Alon, Richard Beigel, Simon Kasif, Steven Rudich, Benny Sudakov, Learning a hidden matching, in: Proceedings of the 43rd Symposium on Foundations of Computer Science, FOCS'02, Washington, DC, USA, 2002, pp. 197–206.
- [3] Reid Andersen, Fan Chung, Kevin Lang, Using pagerank to locally partition a graph, *Internet Math.* 4 (1) (2007) 35–64. An extended abstract appeared in FOCS'2006.
- [4] Reid Andersen, Yuval Peres, Finding sparse cuts locally using evolving sets, in: Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing, STOC'09, 2009, pp. 235–244.
- [5] Dana Angluin, Jiang Chen, Learning a hidden hypergraph, *J. Mach. Learn. Res.* 7 (2006) 2215–2236.
- [6] Dana Angluin, Jiang Chen, Learning a hidden graph using  $O(\log n)$  queries per edge, *J. Comput. System Sci.* 74 (4) (2008) 546–556. Appeared in COLT 2004.
- [7] Mathilde Bouvel, Vladimir Grebinski, Gregory Kucherov, Combinatorial search on graphs motivated by bioinformatics applications: a brief survey, in: Proceedings of the 31st International Conference on Graph-Theoretic Concepts in Computer Science, WG'05, 2005, pp. 16–27.
- [8] Wray Buntine, Theory refinement on bayesian networks, in: Proceedings of the Seventh International Conference on Uncertainty in Artificial Intelligence, 1991, pp. 52–60.

- [9] C.N. Chow, C.K. Liu, Approximating discrete probability distributions with dependence trees, *IEEE Trans. Inf. Theory* 14 (3) (1968) 462–467.
- [10] Mark Coates, Alfred O. Hero III, Robert Nowak, Bin Yu, Internet tomography, *IEEE Signal Process. Mag.* 19 (3) (2002) 47–65.
- [11] Arthur P. Dempster, Nan M. Laird, Donald B. Rubin, Maximum likelihood from incomplete data via the em algorithm, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 39 (1) (1977) 1–38.
- [12] Richard O. Duda, Peter E. Hart, *Pattern Classification and Scene Analysis*, Vol. 3, Wiley, New York, 1973.
- [13] Nir Friedman, Learning belief networks in the presence of missing values and hidden variables, in: *ICML*, vol. 97, 1997, pp. 125–133.
- [14] Nir Friedman, The Bayesian structural EM algorithm, in: *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, 1998, pp. 129–138.
- [15] Leszek Gasieniec, Tomasz Radzik, Graph-theoretic concepts in computer science, in: *Chapter Memory Efficient Anonymous Graph Exploration*, Springer-Verlag, Berlin, Heidelberg, 2008, pp. 14–29.
- [16] Vladimir Grebinski, Gregory Kucherov, Reconstructing a hamiltonian cycle by querying the graph: application to dna physical mapping, *Discrete Appl. Math.* 88 (1–3) (1998) 147–165.
- [17] Vladimir Grebinski, Gregory Kucherov, Optimal reconstruction of graphs under the additive model, *Algorithmica* 28 (1) (2000) 104–124.
- [18] Vincent Gripon, Michael Rabbat, Reconstructing a graph from path traces, in: *Proceedings of International Symposium on Information Theory*, July 2013.
- [19] Patric Hagmann, Leila Cammoun, Xavier Gigandet, Reto Meuli, Christopher J. Honey, Van J. Wedeen, Olaf Sporns, Mapping the structural core of human cerebral cortex, *PLoS biology* 6 (7) (2008) e159.
- [20] Jotun J Hein, An optimal algorithm to reconstruct trees from additive distance data, *Bull. Math. Biol.* 51 (5) (1989) 597–603.
- [21] Valerie King, Li Zhang, Yunhong Zhou, On the complexity of distance-based evolutionary tree reconstruction, in: *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA'03*, Philadelphia, PA, USA, 2003, pp. 444–453.
- [22] Edda Klipp, Ralf Herwig, Axel Kowald, Christoph Wierling, Hans Lehrach, *Systems Biology in Practice: Concepts, Implementation and Application*, John Wiley & Sons, 2008.
- [23] Daphne Kollar, Nir Friedman, *Probabilistic Graphical Models: Principles and Techniques*, The MIT Press, 2009.
- [24] Claire Mathieu, Hang Zhou, Graph reconstruction via distance oracles, in: *ICALP (1)*, in: *Lecture Notes in Computer Science*, vol. 7965, Springer, 2013, pp. 733–744.
- [25] Hanna Mazzawi, Optimally reconstructing weighted graphs using queries, in: *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA'10*, Philadelphia, PA, USA, 2010, pp. 608–615.
- [26] Arvind Narayanan, Vitaly Shmatikov, De-anonymizing social networks, in: *Security and Privacy, 2009 30th IEEE Symposium on*, IEEE, 2009, pp. 173–187.
- [27] Mark Newman, Albert-László Barabási, Duncan J Watts, *The Structure and Dynamics of Networks*, Princeton University Press, 2006.
- [28] Lorenzo Orecchia, Zeyuan Allen Zhu, Flow-based algorithms for local graph clustering, in: *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA'14*, 2014.
- [29] Bernhard O. Palsson, *Properties of reconstructed networks*, in: *Cambridge: Systems Biology*, 2006.
- [30] Judea Pearl, Thomas S. Verma, A theory of inferred causation, *Stud. Logic Found. Math.* 134 (1995) 789–811. Appeared in *Proc. Conference on Knowledge Representation and Reasoning (KR)* 1991.
- [31] Michael G. Rabbat, Mário A.T. Figueiredo, Robert D. Nowak, Network inference from co-occurrences, *IEEE Trans. Inf. Theory* 54 (9) (2008) 4053–4068.
- [32] Lev Reyzin, Nikhil Srivastava, Learning and verifying graphs using queries with a focus on edge counting, in: *Proceedings of the 18th international conference on Algorithmic Learning Theory, ALT'07*, Springer-Verlag, Berlin, Heidelberg, 2007, pp. 285–297.
- [33] Ronitt Rubinfeld, Gil Tamir, Shai Vardi, Ning Xie, Fast local computation algorithms. in *ICS*, 2011, pp. 223–238.
- [34] David J. Spiegelhalter, Steffen L. Lauritzen, Sequential updating of conditional probabilities on directed graphical structures, *Networks* 20 (5) (1990) 579–605.
- [35] Daniel A. Spielman, Shang-Hua Teng, Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems, in: *Proceedings of the Thirty-sixth Annual ACM Symposium on Theory of Computing, STOC'04*, 2004, pp. 81–90.
- [36] Olaf Sporns, Dante R. Chialvo, Marcus Kaiser, Claus C. Hilgetag, Organization, development and function of complex brain networks, *Trends Cogn. Sci.* 8 (9) (2004) 418–425.
- [37] Thomas Verma, Judea Pearl, An algorithm for deciding if a set of observed independencies has a causal explanation, in: *Proceedings of the Eighth International Conference on Uncertainty in Artificial Intelligence*, 1992, pp. 323–330.
- [38] Stanley Wasserman, *Social Network Analysis: Methods and Applications*, Vol. 8, Cambridge university press, 1994.
- [39] Zeyuan Allen Zhu, Silvio Lattanzi, Vahab Mirrokni, A local algorithm for finding well-connected clusters, in: *Proceedings of the 30th International Conference on Machine Learning, ICML'13*, 2013.