

Network Discovery and Verification

Zuzana Beerliova¹, Felix Eberhard¹, Thomas Erlebach², Alexander Hall¹, Michael Hoffmann²,
Matúš Mihalák², and L. Shankar Ram¹

¹ Department of Computer Science, ETH Zürich
{bzuzana, mhall, lshankar}@inf.ethz.ch, efelix@student.ethz.ch

² Department of Computer Science, University of Leicester
{te17, mh55, mm215}@mcs.le.ac.uk

Abstract. Consider the problem of discovering (or verifying) the edges and non-edges of a network, modelled as a connected undirected graph, using a minimum number of queries. A query at a vertex v discovers (or verifies) all edges and non-edges whose endpoints have different distance from v . In the network discovery problem, the edges and non-edges are initially unknown, and the algorithm must select the next query based only on the results of previous queries. We study the problem using competitive analysis and give a randomized on-line algorithm with competitive ratio $O(\sqrt{n \log n})$ for graphs with n vertices. We also show that no deterministic algorithm can have competitive ratio better than 3. In the network verification problem, the graph is known in advance and the goal is to compute a minimum number of queries that verify all edges and non-edges. This problem has previously been studied as the problem of placing landmarks in graphs or determining the metric dimension of a graph. We show that there is no approximation algorithm for this problem with ratio $o(\log n)$ unless $\mathcal{P} = \mathcal{NP}$. Furthermore, we prove that the optimal number of queries for d -dimensional hypercubes is $\Theta(d/\log d)$.

1 Introduction

In recent years, there has been an increasing interest in the study of networks whose structure has not been imposed by a central authority but arisen from local and distributed processes. Prime examples of such networks are the Internet and unstructured peer-to-peer networks such as Gnutella. For these networks, it is very difficult and costly to obtain a “map” providing an accurate representation of all nodes and the links between them. Such maps would be useful for many purposes, e.g., for studying routing aspects or robustness properties of these networks.

In order to create maps of the Internet, a commonly used technique is to obtain local views of the network from various locations (vantage points) and combine them into a map that is hopefully a good approximation of the real network [2, 10]. More generally, one can view this technique as an approach for discovering the topology of an unknown network by using a certain type of queries—a query corresponds to asking for the local view of the network from one specific vantage point. In this paper, we formalize *network discovery* as a combinatorial optimization problem whose goal is to minimize the number of queries required to discover all edges and non-edges of the network. We study the problem as an on-line problem using competitive analysis. Initially, the network is unknown to the algorithm. To decide the next query to ask, the algorithm can only use the knowledge about the network it has gained from the answers of previously asked queries. In the end, the number of queries asked by the algorithm is compared to the optimal number of queries sufficient to discover the network. In this paper, we consider a query model in which the answer to a query at a vertex v consists of all edges and non-edges whose endpoints have different (graph-theoretic) distance from v .

In the off-line version of the network discovery problem, the network is known to the algorithm from the beginning. The goal is to compute a minimum number of queries that suffice to discover the network. Although an algorithm for this off-line problem would not be useful for network discovery (if the network is known in advance, there is no need to discover it), it could be employed for network verification, i.e., for checking whether a given map is accurate. Therefore, we refer to the off-line version of network discovery as *network verification*. For network verification, we are interested in polynomial-time optimal or approximation algorithms.

Motivation. As mentioned above, the motivation for our research comes from the problem of discovering information about the topology of communication networks such as the Internet or peer-to-peer networks. The query model that we study is motivated by the following considerations. First, notice that our query model can be interpreted in the following way: A query at v yields the shortest-path subgraph rooted at v , i.e., the set of all edges on shortest paths between v and any other vertex. To see that this is equivalent to our definition (where a query yields all edges and non-edges between vertices of different distance from v), note that an edge connects vertices

of different distance from v if and only if it lies on a shortest path between v and one of these two vertices. Furthermore, the shortest-path subgraph rooted at v implicitly confirms the absence of all edges between vertices of different distance from v that are not part of the shortest-path subgraph. Examples of real-life scenarios where the shortest-path subgraph rooted at a node of the network could be determined are the following:

- With traceroute tools, one can determine the path that packets take in the network if they are sent from one’s node to some destination. If each traceroute experiment returns a random shortest path to the destination, one could use repeated traceroute experiments to all destinations to discover all edges of the shortest-path subgraph. Making a query at v would mean getting access to node v and running repeated traceroute experiments from v to all other nodes. If we assume that the cost of getting access to a node is much higher than that of running the traceroute-experiments, minimizing the number of queries (corresponding to the number of nodes one has to get access to) is a meaningful goal.
- If the network routes all packets along arbitrary shortest paths, one could imagine a routing protocol in which each node stores the shortest-path subgraph rooted at that node. In this case, reading out the routing table at a node would correspond to making a query at that node.

It is clear that our model of network discovery is a simplification of reality. In real networks, routing is not necessarily along shortest paths, but may be affected by routing policies, link qualities, or link capacities. Furthermore, routing tables or traceroute experiments will often reveal only a single path (or at most a few different paths) to each destination, but not the whole shortest-path subgraph. Nevertheless, we believe that our model is a good starting point for a theoretical investigation of fundamental issues arising in network discovery.

Related Work. We are not aware of any theoretical work on network discovery problems using query models similar to the one we consider in this paper. Graph discovery problems have been studied in distributed settings where one or several agents move along the edges of the graph (see, e.g., [3]); the problems arising in such settings appear to require very different techniques from the ones in our setting.

It turns out, however, that the network verification problem has previously been considered as the problem of placing landmarks in graphs. Here, the motivation is to place landmarks in as few vertices of the graph as possible in such a way that each vertex of the graph is uniquely identified by the vector of its distances to the landmarks. This problem has been studied by Khuller et al. in [7]. They prove that the problem is \mathcal{NP} -hard in general (by reduction from 3-SAT) and give an $O(\log n)$ -approximation algorithm for graphs with n vertices, based on SETCOVER. For trees, they show that the problem can be solved optimally in polynomial time. For d -dimensional grids, $d \geq 2$, they show that d landmarks suffice and claim that d landmarks are also necessary; we will show that this lower bound is incorrect in the case of d -dimensional hypercubes (grids with side length 2). Furthermore, they prove that one landmark is sufficient if and only if G is a path, and discuss properties of graphs for which 2 landmarks suffice. They also show that if k landmarks suffice for a graph with n vertices and diameter D , we must have $n \leq D^k + k$. The smallest number of landmarks that are required for a given graph G is also called the *metric dimension* of G [6]. For further known results about this concept, we refer to the survey [4]. Results for the problem variant where extra constraints are imposed on the set of landmarks (e.g., connectedness or independence) are surveyed in [8].

Our Results. For network discovery, we give a lower bound showing that no deterministic on-line algorithm can have competitive ratio better than 3, and we present a randomized on-line algorithm with competitive ratio $O(\sqrt{n \log n})$ for networks with n nodes. For the network verification problem, we prove that it cannot be approximated within a factor of $o(\log n)$ unless $\mathcal{P} = \mathcal{NP}$, thus showing that the approximation algorithm from [7] is best possible (up to constant factors). We also give a useful lower bound formula for the optimal number of queries of a given graph, and we show that the optimal number of queries for d -dimensional hypercubes is $\Theta(d/\log d)$.

The remainder of the paper is structured as follows. Section 2 gives preliminaries and defines the problems formally. Sections 3 and 4 give our results for network discovery and network verification, respectively. Section 5 points to open problems and promising directions for future research.

2 Preliminaries and Problem Definitions

Throughout this paper, the term *network* refers to a connected, undirected graph. For a given graph $G = (V, E)$, we denote the number of vertices by $n = |V|$ and the number of edges by $m = |E|$. For two distinct nodes $u, v \in V$, we say that $\{u, v\}$ is an *edge* if $\{u, v\} \in E$ and a *non-edge* if $\{u, v\} \notin E$. The set of non-edges of G

is denoted by \bar{E} . We assume that the set V of nodes is known in advance and that it is the presence or absence of edges that needs to be discovered or verified.

A *query* is specified by a vertex $v \in V$ and called a *query at v* . The query at v is also denoted by v . The answer of a query at v consists of a set E_v of edges and a set \bar{E}_v of non-edges. These sets are determined as follows. Label every vertex $u \in V$ with its distance (number of edges on a shortest path) from v . We refer to sets of vertices with the same distance from v as *layers*. Then E_v is the set of all edges connecting vertices in different layers, and \bar{E}_v is the set of all non-edges whose endpoints are in different layers. Because the query result can be seen as a layered graph, we refer to this query model as the *layered-graph* query model.

A set $Q \subseteq V$ of queries discovers (all edges and non-edges of) a graph $G = (V, E)$ if $\bigcup_{q \in Q} E_q = E$ and $\bigcup_{q \in Q} \bar{E}_q = \bar{E}$. In the off-line case, we also say “verifies” instead of “discovers”. The network verification problem is to compute, for a given network G , a smallest set of queries that verifies G . The network discovery problem is the on-line version of the network verification problem. Its goal is to compute a smallest set of queries that discovers G . Here, the edges and non-edges of G are initially unknown to the algorithm, the queries are made sequentially, and the next query must always be determined based only on the answers of previous queries.

We denote by $OPT(G)$, for a given graph G , the cardinality of an optimal query set for verifying G , and by $A(G)$ the cardinality of the query set produced by an algorithm A . The quality of an algorithm is measured by the worst possible ratio $A(G)/OPT(G)$ over all networks G . In the off-line case, an algorithm is a ρ -approximation algorithm (and achieves approximation ratio ρ) if it runs in polynomial time and satisfies $A(G)/OPT(G) \leq \rho$ for all networks G . In the on-line case, an algorithm is ρ -competitive (and achieves competitive ratio ρ) if $A(G)/OPT(G) \leq \rho$ for all networks G . It is weakly ρ -competitive if $A(G) \leq \rho \cdot OPT(G) + c$ for some constant c . If the on-line algorithm is randomized, $A(G)$ is replaced by $\mathbb{E}[A(G)]$ in these definitions. Note that, as is common in competitive analysis, we do not require on-line algorithms to run in polynomial time.

We refer to the network discovery problem with the layered-graph query model and the goal of discovering all edges and non-edges as LG-ALL-DISCOVERY, and to its off-line version as LG-ALL-VERIFICATION. Related problem formulations with different query models and discovery objectives are discussed in Section 5.

3 Network Discovery

We consider the on-line scenario. It is clear that any algorithm that does not repeat queries has competitive ratio at most $n - 1$, since $n - 1$ queries are always sufficient to discover a network. Furthermore, the inapproximability result that we will derive in Section 4 (Theorem 3) shows that we cannot hope for a polynomial-time on-line algorithm with competitive ratio $o(\log n)$; it may still be possible to obtain such a ratio using exponential-time on-line algorithms, however.

First, we present a lower bound on the competitive ratio of deterministic on-line algorithms.

Theorem 1. *There is no deterministic on-line algorithm for LG-ALL-DISCOVERY with weak competitive ratio $3 - \varepsilon$ for any $\varepsilon > 0$.*

Proof. Let A be any deterministic algorithm for LG-ALL-DISCOVERY. We first give a simpler proof that A cannot be better than 2-competitive. Consider Figure 1(a). We refer to the subgraph induced by the vertices labelled r, x, y , and z as a *2-gadget*. Assume that the given graph G consists of a global root g and $k, k \geq 2$, disjoint copies of the 2-gadget, with the r -vertex of each 2-gadget connected to the global root g . One can easily verify that $OPT(G) = k$ for this graph, and that the set of all x -vertices of the 2-gadgets constitutes an optimal query set. On the other hand, algorithm A can be forced to make the first query at g (as, initially, the vertices are indistinguishable to the algorithm). This will not discover any information about edges or non-edges between vertices x, y and z of each 2-gadget. The only queries that can discover this information are queries at x, y and z . In fact, a query at x or y suffices to discover the edge between x and y and the non-edges between x and z and between y and z . When A makes the first query among the vertices in $\{x, y, z\}$ of a 2-gadget, we can force it to make that query at z , since the three vertices are indistinguishable to the algorithm. The query at z does not discover the edge between x and y . Therefore, the algorithm must make a second query in the 2-gadget to discover that edge. In total, the algorithm must make at least $2k + 1$ queries. As the construction works for arbitrary values of k , this shows that no deterministic on-line algorithm can guarantee weak competitive ratio $2 - \varepsilon$ for any constant $\varepsilon > 0$.

To get a stronger lower bound of 3, we create a new gadget, called the *3-gadget*, as shown in Figure 1(b). The 3-gadget is the subgraph induced by all vertices except g in the figure. We claim that A can be forced to make 6 queries in each 3-gadget, whereas the optimum query set consists of only 2 vertices in each 3-gadget (drawn shaded in the figure). If we construct a graph with $k, k \geq 2$, disjoint copies of the 3-gadget, the s -vertex in

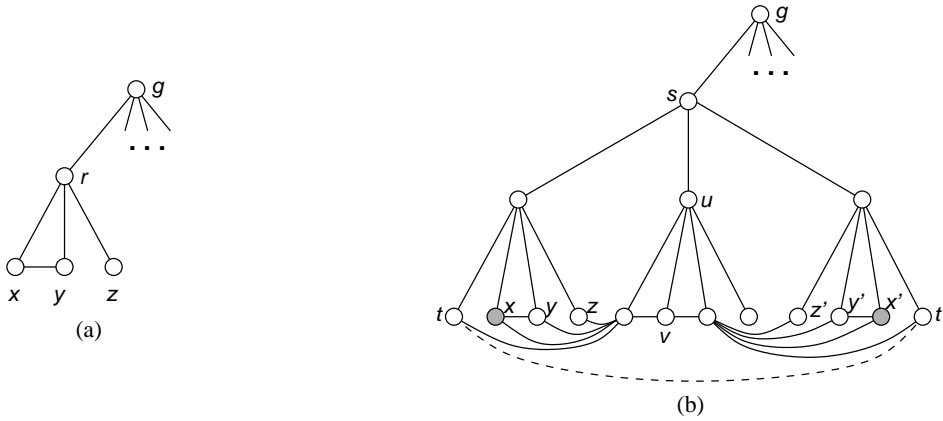


Fig. 1. Lower bound constructions

each of them connected to the global root g as indicated in the figure, we get a graph G for which we claim that $OPT(G) = 2k$ and the algorithm A can be forced to make at least $6k + 1$ queries, showing that no deterministic on-line algorithm can guarantee weak competitive ratio $3 - \varepsilon$ for any constant $\varepsilon > 0$.

To see that $OPT(G) = 2k$, let Q be the set of queries consisting of the two shaded vertices from each copy of the 3-gadget as shown in Figure 1(b). We claim that Q discovers G . This can be verified manually as follows: For every vertex in a 3-gadget Π , consider the 3-tuple whose components are the distances from that vertex to the two query vertices in Π and the distance to an arbitrary query vertex from Q outside Π . One finds that each vertex in Π has a unique 3-tuple, showing that all edges and non-edges of Π are discovered by Q . Each non-edge between two different 3-gadgets is discovered by one of the queries inside these two 3-gadgets. The edges and non-edges between g and each 3-gadget are also discovered by the queries inside that 3-gadget. Hence, $OPT(G) \leq 2k$. It is also easy to see that $OPT(G) \geq 2k$, because each of the edges $\{x, y\}$ and $\{x', y'\}$ (see Figure 1(b)) of each 3-gadget requires a separate query.

To show that $A(G) \geq 6k + 1$, we argue as follows. First, we can force A to make the first query at g . This will not reveal any information about edges within the same layer of any of the 3-gadgets. We view each 3-gadget as consisting of s and a left part, a middle part, and a right part. The left part consists of the left child of s and its four adjacent vertices below (these four vertices are called *bottom vertices*, and the left child of s is called the *root* of that part); the middle and right part are defined analogously. The three parts of a 3-gadget Π are indistinguishable to A until it makes its first query inside Π . A query at s would not discover any new information about Π , so we can ignore queries that A might make at s in the following arguments. When A makes its first query inside Π , we can force this query to be in the middle part, and we can force it to be at u or v . In both cases, the query does not discover any information about the edges and non-edges between the bottom vertices of the left part, nor does it discover any information about the edges and non-edges between the bottom vertices of the right part, nor does it discover the edge drawn dashed. When A chooses its second query in Π , it could be in the left part, in the middle part, or in the right part. Assume that A chooses the left part; since the bottom vertices of the left part are still indistinguishable to A , we can force A to make the query either at the root of the left part or at the bottom vertex t . Similarly, in the right part we can force A to make the query at its root or at t' . In the middle part, A can make the query anywhere. In any case, the second query made by A does not discover any information about edges and non-edges between vertices in the set $\{x, y, z\}$ and in the set $\{x', y', z'\}$. Similarly as in the case of Figure 1(a), for each of these sets we can force A to make the first query at z (at z') and thus require a second query at x or y (at x' or y') to discover everything about these groups. In total, A must make at least 6 queries in each 3-gadget. \square

With the gadget of Figure 1(a) one can prove easily that no randomized on-line algorithm for LG-ALL-DISCOVERY can have weak competitive ratio $4/3 - \varepsilon$ for any $\varepsilon > 0$; just observe that we can force a randomized algorithm to make the first query at z with probability at least $1/3$. Note that all lower bounds on the weak competitive ratio also hold for the (standard) competitive ratio where no additive constant c is allowed.

Next, we present a randomized on-line algorithm that is $O(\sqrt{n \log n})$ -competitive.

Theorem 2. *There is a randomized on-line algorithm that achieves competitive ratio $O(\sqrt{n \log n})$ for LG-ALL-DISCOVERY.*

```

E ← ∅; /* discovered edges */
N ← ∅; /* discovered non-edges */
A ←  $\binom{V}{2}$ ; /* all pairs of distinct nodes */
/* Phase 1 */
for i = 1 to  $3\sqrt{n \ln n}$  do
    v ← randomly chosen node from V;
    (Ev, Nv) ← query(v);
    E ← E ∪ Ev;
    N ← N ∪ Nv;
od;
/* Phase 2 */
while E ∪ N ≠ A do
    {u, v} ← an arbitrary element of A \ (E ∪ N);
    (Eu, Nu) ← query(u);
    (Ev, Nv) ← query(v);
    E ← E ∪ Eu ∪ Ev;
    N ← N ∪ Nu ∪ Nv;
    S ← set of nodes from which the (non-)edge {u, v} is discovered;
    foreach x ∈ S \ {u, v} do
        (Ex, Nx) ← query(x);
        E ← E ∪ Ex;
        N ← N ∪ Nx;
    od;
od;

```

Fig. 2. On-line algorithm for LG-ALL-DISCOVERY

Proof. The on-line algorithm is shown in Figure 2. In the first phase, it makes $3\sqrt{n \ln n}$ queries at nodes chosen uniformly at random. In the second phase, as long as node pairs with unknown status exist, it picks an arbitrary such pair $\{u, v\}$ and proceeds as follows. First, it queries u and v in order to determine the distance of all nodes to u and v . From this it can deduce the set S of nodes from which the edge or non-edge between u and v can be discovered; these are simply the nodes for which the distance to u differs from the distance to v . Then, it queries all remaining nodes in S .

To analyze the algorithm, it is helpful to view LG-ALL-DISCOVERY as a HITTINGSET problem. For every edge or non-edge $\{u, v\}$, let S_{uv} be the set of nodes from which a query discovers $\{u, v\}$. The task of the LG-ALL-DISCOVERY problem translates into the task of computing a subset of V that hits all sets S_{uv} . The goal of the first phase is to hit all sets that have size at least $\sqrt{n \ln n}$ with high probability. If this succeeds, the problem remaining for the second phase is a HITTINGSET problem where all sets have size at most $\sqrt{n \ln n}$. The algorithm of the second phase repeatedly picks an arbitrary set that is not yet hit, and includes all its elements in the solution. As the sets have size at most $\sqrt{n \ln n}$, the number of queries made in the second phase is at most a factor of $\sqrt{n \ln n}$ away from the optimum.

Let us make this analysis precise. Consider a node pair $\{u, v\}$ for which the set S_{uv} has size at least $\sqrt{n \ln n}$. In each query of the first phase, the probability that S_{uv} is not hit is at most $1 - \frac{\sqrt{n \ln n}}{n} = 1 - \frac{\sqrt{\ln n}}{\sqrt{n}}$. Thus, the probability that S_{uv} is not hit throughout the first phase is at most

$$\left(1 - \frac{\sqrt{\ln n}}{\sqrt{n}}\right)^{3\sqrt{n \ln n}} = \left(\left(1 - \frac{\sqrt{\ln n}}{\sqrt{n}}\right)^{\frac{\sqrt{n}}{\sqrt{\ln n}}}\right)^{3 \ln n} \leq e^{-3 \ln n} = \frac{1}{n^3}.$$

There are at most $\binom{n}{2}$ sets S_{uv} of cardinality at least $\sqrt{n \ln n}$. The probability that at least one of them is not hit in the first phase is at most $\binom{n}{2} \cdot \frac{1}{n^3} \leq \frac{1}{n}$.

Now consider the second phase, conditioned on the event that the first phase has hit all sets S_{uv} of size at least $\sqrt{n \ln n}$. In each iteration of the while-loop of the second phase, the algorithm asks at most $\sqrt{n \ln n}$ queries. Let ℓ be the number of iterations. It is clear that the optimum must make at least ℓ queries, because no two unknown pairs $\{u, v\}$ considered in different iterations of the second phase can be resolved by the same query.

Hence, we know $OPT(G) \geq 1$ and $OPT(G) \geq \ell$. Furthermore, the number of queries made by the algorithm is at most $3\sqrt{n \ln n} + \ell\sqrt{n \ln n} = O(\sqrt{n \log n}) \cdot OPT(G)$.

So we have that with probability at least $1 - \frac{1}{n}$, the first phase succeeds and the algorithm makes $O(\sqrt{n \log n}) \cdot OPT(G)$ queries. If the first phase fails, the algorithm makes at most n queries. This case increases the expected number of queries made by the algorithm by at most $\frac{1}{n} \cdot n = 1$. Thus, we have that the expected number of queries is at most

$$O(\sqrt{n \log n}) \cdot OPT(G) + \frac{1}{n} \cdot n = O(\sqrt{n \log n}) \cdot OPT(G).$$

This concludes the proof of the theorem. \square

4 Network Verification

As mentioned in the introduction, the problem LG-ALL-VERIFICATION is known to be \mathcal{NP} -hard, and there is a SETCOVER-based $O(\log n)$ -approximation algorithm for it [7]. In this section, we first show that it is impossible to achieve approximation ratio $o(\log n)$ unless $\mathcal{P} = \mathcal{NP}$.

Theorem 3. *It is \mathcal{NP} -hard to approximate LG-ALL-VERIFICATION within ratio $o(\log n)$.*

Proof. We prove the inapproximability result using an approximation-preserving reduction from the *test collection problem (TCP)*. That problem is defined as follows:

Problem TCP

Input: ground set S and collection \mathcal{C} of subsets of S

Feasible solution: subset $\mathcal{C}' \subseteq \mathcal{C}$ such that for every two distinct elements x and y of S , there exists a set $C \in \mathcal{C}'$ such that exactly one of x and y is in C .

Objective: minimize the cardinality of \mathcal{C}'

In the original application for TCP, S is a set of diseases and \mathcal{C} is a collection of tests. A test $C \in \mathcal{C}$, applied to a patient, will give a positive result if the patient is infected by a disease in C . If a patient is known to be infected by exactly one of the diseases in S , the goal of TCP is to compute a minimum number of tests that together can uniquely identify that disease.

Without loss of generality, we can restrict ourselves to instances of TCP in which any two elements of the ground set can be separated by at least one of the sets in \mathcal{C} ; instances without this property do not have any feasible solutions.

Halldorsson et al. [5] prove that TCP cannot be approximated with ratio $o(\log |S|)$ unless $\mathcal{P} = \mathcal{NP}$. Their proof uses an approximation-preserving reduction from SETCOVER; the latter problem was shown \mathcal{NP} -hard to approximate within $o(\log n)$, where n is the cardinality of the ground set, by Arora and Sudan [1]. The proof by Arora and Sudan establishes the inapproximability result for SETCOVER even for instances in which the size of the ground set and the number of sets are polynomially related. The reduction from SETCOVER to TCP maintains this property. Hence, we know that it is \mathcal{NP} -hard to approximate TCP with ratio $o(\log |S|)$ even for instances satisfying $|\mathcal{C}| \leq |S|^g$ for some positive constant g .

Let an instance (S, \mathcal{C}) of TCP be given. Let $n_{\text{TCP}} = |S|$ and $m_{\text{TCP}} = |\mathcal{C}|$. By the remark above, we can assume that $m_{\text{TCP}} = n_{\text{TCP}}^{O(1)}$. We construct an instance $G = (V, E)$ of LG-ALL-VERIFICATION as follows. First, we add $n_{\text{TCP}} + m_{\text{TCP}}$ vertices to V : an *element vertex* v_s for every element $s \in S$ and a *test vertex* u_C for every $C \in \mathcal{C}$. We add the following edges to E :

- Any two element vertices are joined by an edge.
- Every test vertex u_C is joined to all element vertices v_s with $s \in C$.

The idea behind this construction is that queries at test vertices verify all edges in the clique of element vertices if and only if the corresponding tests form a test cover. We have to extend the construction slightly since, in LG-ALL-VERIFICATION, the edges and non-edges incident to the test vertices need to be verified as well. We add a small number of auxiliary vertices that take care of this. We add $h = 2(\lceil \log m_{\text{TCP}} \rceil + 2)$ auxiliary vertices w_1, \dots, w_h to G . For each i , $1 \leq i \leq h/2$, the auxiliary vertices w_{2i-1} and w_{2i} are said to form a *pair*. In addition, we add one extra node z . We add the following edges:

- The two auxiliary vertices in each pair are joined by an edge.
- Number the m_{TCP} test vertices arbitrarily from 0 to $m_{\text{TCP}} - 1$. Both auxiliary vertices in the i -th pair, $1 \leq i \leq h/2 - 2$, are joined to those of the m_{TCP} test vertices whose number has a 1 in the i -th position of its binary representation.

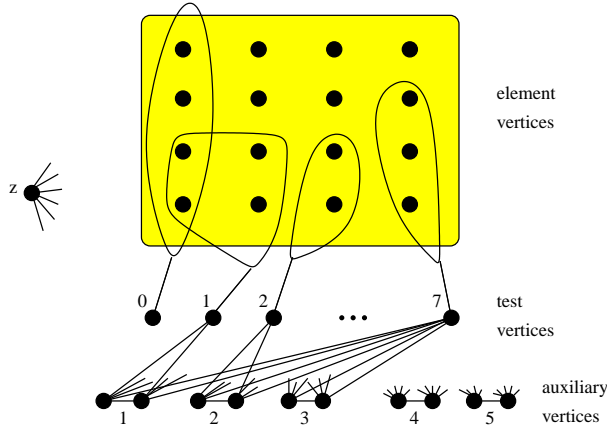


Fig. 3. Illustration of the construction of the graph $G = (V, E)$ that is an instance of LG-ALL-VERIFICATION. The auxiliary vertices in pairs 4 and 5 are adjacent to all test vertices. The auxiliary vertices in pair i , $1 \leq i \leq 3$, are adjacent to the test vertices whose number has a 1 in position i of the binary representation. For example, the auxiliary vertices in pair 2 are adjacent to test vertices 2, 3, 6 and 7

- Both auxiliary vertices in the last two pairs are joined to all test vertices.
- The extra node z is joined to all other vertices of the graph.

The graph constructed in this way is denoted by $G = (V, E)$. See Figure 3 for an illustration. We prove two claims:

Claim 1. Given a solution \mathcal{C}' to the TCP instance (S, \mathcal{C}) , there is a solution Q of the constructed instance $G = (V, E)$ of LG-ALL-VERIFICATION satisfying $|Q| = |\mathcal{C}'| + \lceil \log m_{\text{TCP}} \rceil + 2$.

Proof (of Claim 1). Let a solution \mathcal{C}' to the TCP instance (S, \mathcal{C}) be given. Let Q contain all test vertices corresponding to sets $C \in \mathcal{C}'$ as well as the first vertex of every pair of auxiliary vertices. Obviously, we have $|Q| = |\mathcal{C}'| + \lceil \log m_{\text{TCP}} \rceil + 2$. We claim that Q discovers all edges and non-edges of G . First, it is clear that Q discovers all edges between element vertices, as \mathcal{C}' was a solution to TCP. We see that Q discovers all edges and non-edges between test vertices and element vertices by the query at the first vertex in the last pair of auxiliary vertices. The non-edges between the test vertices are discovered by the queries at the auxiliary vertices in the first $h/2 - 2$ pairs: for every two test vertices, there exists a position i in which their binary representations differ, and the auxiliary vertex from the i -th pair discovers the non-edge between these two test vertices. The edges connecting each pair of auxiliary vertices and the non-edges between different pairs of auxiliary vertices are discovered by the queries at the auxiliary vertices. The edges and non-edges between test vertices and auxiliary vertices are discovered by the queries in the last two pairs of auxiliary vertices. Finally, the edges incident to z are discovered by queries at two different auxiliary vertices. \square

Claim 2. Given a solution Q to the constructed instance $G = (V, E)$ of LG-ALL-VERIFICATION, one can construct in polynomial time a solution \mathcal{C}' of the original TCP instance (S, \mathcal{C}) satisfying $|\mathcal{C}'| \leq |Q| - \lceil \log m_{\text{TCP}} \rceil - 2$.

Proof (of Claim 2). Observe that Q must contain at least one vertex from each pair of auxiliary vertices; otherwise, the edge joining this pair would not be discovered. The queries at these vertices do not discover any edges between element vertices (note that all element vertices are at distance 2 from any auxiliary vertex because of the extra vertex z). Let Q' be the vertices in Q that are not auxiliary vertices. We have $|Q'| \leq |Q| - \lceil \log m_{\text{TCP}} \rceil - 2$. Now, Q' is a set of element vertices and test vertices that, in particular, discovers all edges between element vertices.

Let Q_S be the set of element vertices in Q' and let Q_C be the set of test vertices in Q' . If Q_S is empty, the queries at the vertices in Q_C discover all edges of the clique of element vertices. In particular, this means that for any two distinct element vertices v_s and v_t in V , there must be a query at a vertex adjacent to one of v_s, v_t but not to the other. This shows that the set $\mathcal{C}' = \{C \in \mathcal{C} \mid u_C \in Q'\}$ is a solution of the original TCP instance of the required size.

Now assume that Q_S is nonempty. The set of edges between element vertices that are not discovered by the queries in Q_C is a disjoint union of cliques. The queries in Q_S must discover all edges in these cliques. As the

only edges between element vertices that a query at an element vertex discovers are the edges incident to that vertex, a clique of size k requires $k - 1$ queries. Assume that there are p cliques and denote the number of vertices in these cliques by k_1, k_2, \dots, k_p . Then Q_S must contain at least $\sum_{i=1}^p (k_i - 1)$ vertices. On the other hand, all edges in a clique of size k can always be discovered by $k - 1$ queries at test vertices: we can simply select these queries greedily by choosing, as long as there is an edge $\{u, v\}$ in the clique that has not yet been discovered, any test vertex that is adjacent to one of u, v but not the other. Hence, we can replace the queries in Q_S by at most $\sum_{i=1}^p (k_i - 1)$ queries at test vertices and add these to Q_C , obtaining a set of queries at test vertices that discovers all edges between element vertices. As in the previous paragraph, this set of test vertices gives a solution to the original TCP instance of cardinality at most $|Q'|$. \square

The claims imply that an approximation algorithm with ratio $o(\log n)$ for LG-ALL-VERIFICATION would allow approximating TCP with ratio $o(\log n_{\text{TCP}})$. To show this, we can argue as follows. Assume that there is an approximation algorithm A for LG-ALL-VERIFICATION that achieves approximation ratio $o(\log n)$. Consider the algorithm B for TCP that, given an instance of TCP, constructs an instance of LG-ALL-VERIFICATION as described above, applies A to this instance, and transforms the result into a solution to the TCP instance following Claim 2. Recall that we can assume that the TCP instance satisfies $m_{\text{TCP}} = n_{\text{TCP}}^{O(1)}$. We claim that B achieves approximation ratio $o(\log n_{\text{TCP}})$ for TCP. Let OPT_{TCP} be the optimum objective value for the given TCP instance and OPT_{LG} be the optimum objective value for the constructed instance of LG-ALL-VERIFICATION. Let B_{TCP} and A_{LG} denote the objective values of the solutions computed by B and A , respectively. Note that $OPT_{\text{TCP}} \geq \log n_{\text{TCP}}$ always holds, since n_{TCP} elements cannot be separated by fewer than $\log n_{\text{TCP}}$ test sets.

Claims 1 and 2 imply that $OPT_{\text{TCP}} = OPT_{\text{LG}} - \lceil \log m_{\text{TCP}} \rceil - 2$. We have:

$$OPT_{\text{LG}} = OPT_{\text{TCP}} + \lceil \log m_{\text{TCP}} \rceil + 2 \leq OPT_{\text{TCP}} + O(\log n_{\text{TCP}}) = O(OPT_{\text{TCP}})$$

Thus we can calculate:

$$\begin{aligned} B_{\text{TCP}} &\leq A_{\text{LG}} && \text{(by Claim 2)} \\ &\leq o(\log n) \cdot OPT_{\text{LG}} \\ &= o(\log n) \cdot O(OPT_{\text{TCP}}) \\ &= o(\log n_{\text{TCP}}) \cdot O(OPT_{\text{TCP}}) && \text{(as } n = n_{\text{TCP}} + m_{\text{TCP}} + 2(\lceil \log m_{\text{TCP}} \rceil + 2) + 1 = n_{\text{TCP}}^{O(1)}) \\ &= o(\log n_{\text{TCP}}) \cdot OPT_{\text{TCP}} \end{aligned}$$

This completes the proof of Theorem 3. \square

We have seen that no approximation algorithm can have approximation ratio $o(\log n)$ for the problem LG-ALL-VERIFICATION unless $\mathcal{P} = \mathcal{NP}$. Now we investigate the optimal number of queries for specific graphs. The following is a useful lower bound on $OPT(G)$.

Theorem 4. *If a graph $G = (V, E)$ contains a subgraph H of diameter D_H with n_H vertices, then $OPT(G) \geq \log_{D_H+1} n_H$.*

Proof. Imagine the queries being performed sequentially. At any instant, the unknown edges and non-edges induce disjoint cliques, which we call *unknown groups*. Two vertices are in the same unknown group if and only if they were in the same layer of all queries made so far. Consider the n_H vertices of subgraph H . Initially, all vertices form an unknown group. For each query, the n_H vertices of H will be in at most $D_H + 1$ consecutive layers of the layered graph returned by the query. Therefore, after the first query, at least $n_H / (D_H + 1)$ vertices of H will still be in the same unknown group. Similarly, after k queries, at least $n_H / (D_H + 1)^k$ vertices of H will be in an unknown group together. If k queries suffice to verify all edges and non-edges, the unknown groups must be singletons in the end. So we must have $n_H / (D_H + 1)^k \leq 1$. This implies the statement of the theorem. \square

Note that this theorem implies that a graph containing a clique on k vertices requires at least $\log_2 k$ queries, and a graph with maximum degree Δ requires at least $\log_3(\Delta + 1)$ queries. For the former, take H to be the clique on k vertices, and for the latter, take H to be the subgraph induced by a vertex of degree Δ and its neighbors.

Now, we consider a specific class of graphs, namely d -dimensional hypercubes. The d -dimensional hypercube, denoted H_d , is a graph on 2^d vertices, represented as d -dimensional binary vectors, in which two vertices are adjacent if and only if their binary vectors have Hamming distance 1. It was claimed in [7] that the metric dimension of d -dimensional grids is exactly d . While their upper bound of d is correct, our following result for hypercubes shows that the lower bound of d does not hold for all d -dimensional grids, since d -dimensional hypercubes (a special case of d -dimensional grids) have metric dimension $\Theta(d / \log d)$.

Theorem 5. For the d -dimensional hypercube $H_d = (V, E)$, we have $OPT(H_d) = \Theta(d/\log d)$.

Proof. First, the lower bound of Theorem 4, applied to $H = H_d$, shows that $OPT(H_d) \geq \log_{d+1} 2^d = \Omega(d/\log d)$.

It remains to prove an upper bound of $O(d/\log d)$. Consider a query set Q consisting of the vertex $(0, 0, \dots, 0)$ and $k = 4\alpha \frac{d}{\log d}$ vertices chosen uniformly at random among all vertices of the hypercube (where α is a sufficiently large constant). We show that the probability that Q does not verify the graph, denoted by $P[Q \text{ fails}]$, is smaller than one (actually, tends to zero for increasing d); this proves that there exists a set of $k+1 = O(d/\log d)$ queries that verifies H_d .

$P[Q \text{ fails}]$ is the probability that there exist two vertices $u, v \in V$ with the property “ $d(u, q) = d(v, q)$ for every query vertex $q \in Q$ ”. We now show that Q fails if and only if there exist disjoint, non-empty sets $D_1, D_2 \subseteq \{1, 2, \dots, d\}$ of the same size such that $\forall q \in Q : \sum_{i \in D_1} q_i = \sum_{i \in D_2} q_i$, where $q = (q_1, q_2, \dots, q_d)$. If there are disjoint, non-empty sets D_1 and D_2 of the same size with this property, let u and v be the two vertices whose binary vectors have zeros on positions not in $D_1 \cup D_2$, with u having zeros on positions in D_1 and ones on positions in D_2 , and v having ones on positions in D_1 and zeros on positions in D_2 . It is easy to see that u and v have the same distance from every query vertex $q \in Q$.

For the other direction, assume that Q fails. Then there must be distinct vertices u and v for which $d(u, q) = d(v, q)$ for all $q \in Q$. Let $u = (u_1, u_2, \dots, u_d)$ and $v = (v_1, v_2, \dots, v_d)$. Take D_1 and D_2 as the sets of indices where u and v differ, in such a way that $u_i = 0$ and $v_i = 1$ for all $i \in D_1$, $u_i = 1$ and $v_i = 0$ for all $i \in D_2$, and $u_i = v_i$ for all $i \notin D_1 \cup D_2$. Let $d_1 = |D_1|$ and $d_2 = |D_2|$. For $q \in Q$, let k_1 denote the number of ones in $\{q_i \mid i \in D_1\}$ and k_2 the number of ones in $\{q_i \mid i \in D_2\}$. Then $d(u, q) = k_1 + (d_2 - k_2) + x$ and $d(v, q) = (d_1 - k_1) + k_2 + x$, where x is the contribution of coordinates where u and v agree to the distance of u and v from q . From $d(u, q) = d(v, q)$ we get $(k_1 - k_2) = \frac{1}{2}(d_1 - d_2)$. For the query vertex $q = (0, 0, \dots, 0)$, we get $k_1 = k_2 = 0$ and thus $k_1 - k_2 = 0$, implying $d_1 = d_2$. Thus, D_1 and D_2 must have the same size. This shows that for all $q \in Q$ we must have $k_1 - k_2 = \frac{1}{2}(0 - 0) = 0$. We see that if Q fails then there exist disjoint, non-empty sets $D_1, D_2 \subseteq \{1, 2, \dots, d\}$ of the same size such that $\forall q \in Q : \sum_{i \in D_1} q_i = \sum_{i \in D_2} q_i$.

We can see $Q \setminus \{(0, 0, \dots, 0)\}$ as a random $k \times d$ 0-1 matrix (q_{ij}) where $q_{ij} = 1$ with probability 0.5, $1 \leq i \leq k$ and $1 \leq j \leq d$. Then the probability that Q fails is at most

$$\sum_{\substack{D_1, D_2 \\ D_1 \cap D_2 = \emptyset \\ |D_1| = |D_2| > 0}} P \left[\forall i : \sum_{j \in D_1} q_{ij} = \sum_{j \in D_2} q_{ij} \right] = \sum_{s=1}^{\lfloor d/2 \rfloor} \underbrace{\binom{d}{s} \binom{d-s}{s}}_{\substack{\# \text{ pairs of disjoint} \\ \text{sets of size } s}} P_1^{(s)} \cdot P_2^{(s)} \cdot \dots \cdot P_k^{(s)},$$

where $P_r^{(s)}$ is the probability that two fixed disjoint sets of positions D_1 and D_2 of size s have the same number of ones in row r of the matrix. We have $P_r^{(s)} = P \left[\sum_{j \in D_1} q_{rj} = \sum_{j \in D_2} q_{rj} \right] = \sum_{j=0}^s \binom{s}{j} \frac{1}{2^s} \binom{s}{j} \frac{1}{2^s} = \frac{1}{4^s} \binom{2s}{s}$. Therefore

$$P[Q \text{ fails}] \leq \sum_{s=1}^{\lfloor d/2 \rfloor} \frac{1}{2} \binom{d}{s} \binom{d-s}{s} \left(\frac{1}{4^s} \binom{2s}{s} \right)^k. \quad (1)$$

We show now that the right-hand side of (1) tends to zero for increasingly large values of d by showing that each term of the sum is less than $\frac{1}{2} d^2 2^{-d/\log d}$, i.e. that

$$X(s) := \binom{d}{s} \binom{d-s}{s} \left(\frac{1}{4^s} \binom{2s}{s} \right)^k \leq d^2 \left(\frac{1}{2} \right)^{d/\log d}. \quad (2)$$

This inequality is satisfied for $s = 1$ and sufficiently large α . Hence we consider now the case $1 < s \leq d/2$. We start by bounding the first two factors of $X(s)$ using $s! \geq \left(\frac{s}{e}\right)^s$ (which holds by Stirling's inequality [9, p. 433]):

$$\binom{d}{s} \binom{d-s}{s} \leq \binom{d}{s}^2 \leq \left(\frac{d!}{s!(d-s)!} \right)^2 \leq \left(\frac{d^s}{s!} \right)^2 \leq \left(\frac{e \cdot d}{s} \right)^{2s}.$$

Next we consider the third factor of $X(s)$ and again apply Stirling's approximation formula, using the form $n! \leq \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(\frac{12n}{12n-1}\right)$ [9, p. 433] to bound $(2s)!$, yielding:

$$\left(\frac{1}{4^s} \binom{2s}{s} \right)^k = \left(\frac{(2s)!}{4^s s! s!} \right)^k \leq \left(\frac{\left(\frac{24s}{24s-1}\right)}{\sqrt{\pi s}} \right)^k \leq \frac{1}{s^{k/2}},$$

where the last inequality holds since $\frac{24s}{\sqrt{\pi(24s-1)}} < 1$. We obtain $X(s) \leq (\frac{e \cdot d}{s})^{2s} \cdot s^{-k/2} = 2^{\log(\frac{e \cdot d}{s}) \cdot 2s - \log s \cdot k/2}$. Considering only the exponent and using the definition of k , we proceed to show that $Y(s) := \log(\frac{e \cdot d}{s}) \cdot 2s - 2\alpha \cdot d \cdot \frac{\log s}{\log d} \leq -d/\log d$, implying $X(s) \leq 2^{-d/\log d}$.

We consider two cases for s : (i) $2 \leq s \leq d^{1/2}$, (ii) $d^{1/2} \leq s \leq d/2$.

For (i) we get $Y(s) \leq \log(\frac{e \cdot d}{s}) \cdot 2d^{1/2} - 2\alpha \cdot d \cdot \frac{\log s}{\log d} \leq \log(e \cdot d) \cdot 2d^{1/2} - 2\alpha \cdot \frac{d}{\log d} \leq -d/\log d$, for large enough α .

For (ii) we obtain $Y(s) \leq \log(\frac{e \cdot d}{s}) \cdot 2s - 2\alpha \cdot d \cdot \frac{\log d^{1/2}}{\log d} = \log e \cdot 2s + \log((\frac{d}{s})^{2s}) - \alpha \cdot d$. The term $(\frac{d}{s})^{2s}$ has its maximum in point $s = d/e$ (the first derivative is $(\frac{d}{s})^{2s}(2 \ln(\frac{d}{s}) - 2)$ and the second derivative is $(\frac{d}{s})^{2s}(2 \ln(\frac{d}{s}) - 2)^2 - 2(\frac{d}{s})^{2s}s^{-1}$) and therefore $Y(s) \leq \log e \cdot (2s + 2d/e) - \alpha \cdot d \leq -d/\log d$, for large enough α .

We have shown that $P[Q \text{ fails}] < 1$ and, therefore, H_d can be verified with $O(\frac{d}{\log d})$ queries. \square

We remark that since the probability that Q fails is close to zero, one can choose $O(d/\log d)$ query vertices uniformly at random to discover or verify the hypercube with probability close to 1. This shows that the simple and practical strategy of choosing query vertices at random works well for hypercubes.

5 Directions for Future Work

In this paper, we have considered network discovery and network verification problems in the layered-graph query model. The goal was to discover or verify all edges and non-edges of a network. For the network discovery problem, a significant gap remains between our randomized upper bound of $O(\sqrt{n \log n})$ and the small constant lower bounds. Thus, the major problem left open by our work is to close this gap. Another interesting question is finding a deterministic construction of query sets of size $O(d/\log d)$ for hypercubes.

The subject of our study can be seen as one example of a family of problem settings in which the goal is to discover or verify information about a graph using certain kinds of queries. Different problems are obtained if the query model is varied, or if the objective is changed. Other natural query models are, for example, that a query at v returns only the distances from v to all other vertices of the graph; that a query is specified by two vertices u and v , and returns the set of all edges on shortest paths (or just one shortest path) between u and v ; or that a query returns an arbitrary shortest-path tree rooted at v . Concerning the objective, the goal could be to discover or verify a certain graph parameter such as the diameter, the average path length, or the independence number. One could also relax the requirement and only ask for an approximate answer. For example, one could pose the problem of minimizing the number of queries required to approximate the average path length of the graph within a factor of $1 + \epsilon$.

We believe that the study of such problems could be a fruitful area of research with applications in the monitoring and analysis of communication networks such as the Internet or peer-to-peer networks.

References

1. S. Arora and M. Sudan. Improved low-degree testing and its applications. In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing (STOC'97)*, pages 485–495, 1997.
2. P. Barford, A. Bestavros, J. Byers, and M. Crovella. On the marginal utility of deploying measurement infrastructure. In *Proceedings of the ACM SIGCOMM Internet Measurement Workshop 2001*, November 2001.
3. M. A. Bender and D. K. Slonim. The power of team exploration: Two robots can learn unlabeled directed graphs. In *Proceedings of the 35th Annual IEEE Symposium on Foundations of Computer Science (FOCS'94)*, pages 75–85, 1994.
4. G. Chartrand and P. Zhang. The theory and applications of resolvability in graphs: A survey. *Congressus Numerantium*, 160:47–68, 2003.
5. B. V. Halldórsson, M. M. Halldórsson, and R. Ravi. On the approximability of the minimum test collection problem. In *Proceedings of the 9th Annual European Symposium on Algorithms (ESA'01)*, LNCS 2161, pages 158–169. Springer-Verlag, 2001.
6. F. Harary and R. Melter. The metric dimension of a graph. *Ars Combinatorica*, pages 191–195, 1976.
7. S. Khuller, B. Raghavachari, and A. Rosenfeld. Landmarks in graphs. *Discrete Applied Mathematics*, 70:217–229, 1996.
8. V. Saenpholphat and P. Zhang. Conditional resolvability in graphs: A survey. *International Journal of Mathematics and Mathematical Sciences*, 38:1997–2017, 2004.
9. S. M. Selby and B. Girling, editors. *Standard Mathematical Tables*. The Chemical Rubber Company, Cleveland, OH, 14th edition, 1965.
10. L. Subramanian, S. Agarwal, J. Rexford, and R. Katz. Characterizing the Internet hierarchy from multiple vantage points. In *Proceedings of INFOCOM'02*, June 2002.