

A Sublinear Algorithm for Approximate Shortest Paths in Large Networks

Sabyasachi Basu*
UC Santa Cruz
sbasu3@ucsc.edu

Nadia Kōshima*
MIT
nadianw36@gmail.com

Talya Eden
Bar-Ilan University
talyaa01@gmail.com

Omri Ben-Eliezer
MIT
omrib@mit.edu

C. Seshadhri
UC Santa Cruz
sesh@ucsc.edu

ABSTRACT

Computing distances and finding shortest paths in massive real-world networks is a fundamental algorithmic task in network analysis. There are two main approaches to solving this task. On one hand are traversal-based algorithms like bidirectional breadth-first search (BiBFS), which have no preprocessing step but are slow on individual distance inquiries. On the other hand are indexing-based approaches, which create and maintain a large index. This allows for answering individual inquiries very fast; however, index creation is prohibitively expensive even for moderately large networks. For a graph with 30 million edges, the index created by the state-of-the-art is about 40 gigabytes. We seek to bridge these two extremes: quickly answer distance inquiries without the need for costly preprocessing.

In this work, we propose a new algorithm and data structure, WormHole, for approximate shortest path computations. WormHole leverages structural properties of social networks to build a sublinearly sized index, drawing upon the explicit core-periphery decomposition of Ben-Eliezer et al. [WSDM'22]. Empirically, the preprocessing time of WormHole improves upon index-based solutions by orders of magnitude: for a graph with over a billion edges, indexing takes only a few minutes. Furthermore, individual inquiries are consistently much faster than in BiBFS. The acceleration comes at the cost of a minor accuracy trade-off. Nonetheless, our empirical evidence demonstrates that WormHole accurately answers essentially all inquiries within a maximum additive error of 2 (and an average additive error usually much less than 1). We complement these empirical results with provable theoretical guarantees, showing that WormHole, utilizing a sublinear index, requires $n^{o(1)}$ node queries per distance inquiry in random power-law networks. In contrast, any approach without a preprocessing step (including BiBFS) requires $n^{\Omega(1)}$ queries for the same task.

WormHole offers several additional advantages over existing methods: (i) it does not require reading the whole graph and can thus be used in settings where access to the graph is rate-limited; (ii) unlike the vast majority of index-based algorithms, it returns paths, not just distances; and (iii) for faster inquiry times, it can be combined effectively with other index-based solutions, by running them only on the sublinear core.

1 INTRODUCTION

Scalable computation of distances and shortest paths in a large network is one of the most fundamental algorithmic challenges in graph mining and graph learning tasks, with applications

across science and engineering. Examples of such applications include the identification of important genes or species in biological and ecological networks [18], driving directions in road networks [1–3], redistribution of task processing from mobile devices to cloud [59], computer network design and security [24, 30, 31, 58], and identifying a set of users with the maximum influence in a social network [32, 56], among many others. Thus, a long line of work [5, 21, 25, 28, 62] has developed over the years, constructing scalable algorithms for distance computation for a variety of real-life tasks.

The simplest methods for answering a shortest path inquiry (s, t) use traversals, among which the most basic is a breadth first search (BFS) starting from s until we reach t . However, the inquiry time for BFS is linear in the network size, which is much too slow for real-world networks.¹ A popular modification, Bidirectional BFS (BiBFS), runs BFS from both s and t , alternating between the two, until both ends meet. It has well been observed in the literature that BiBFS performs surprisingly well for shortest path inquiries on a wide range of networks (see, e.g., [8, 12, 62] and the many references within). Because BiBFS does not require any prior knowledge on the network structure, it is suitable when the number of shortest path inquiries being made is relatively small. However, pure traversal-based approaches do not scale well when one is required to answer a large number of shortest pair inquiries. As we show in Figure 1, BiBFS ends up seeing the whole graph within just a few hundred inquiries.

A long line of modern approaches tackles the distance computation problem in a fundamentally different manner, by preprocessing the network and creating an *index*. The index, in turn, supports extremely fast real-time computation of distances. This line of work has been investigated extensively in recent years, with Pruned Landmark Labeling (PLL, Akiba et al. [5]) being perhaps the most influential approach.

In virtually all index-based methods, pre-processing involves choosing a subset \mathcal{L} of nodes, called *landmarks*; computing all shortest paths among them; and keeping an index of the distance of every node in the network to every landmark. Thus, the space requirement for the index is at least of order $n \cdot |\mathcal{L}|$, where n is the total number of nodes. Naively, this memory requirement can be as bad as *quadratic* in n . Despite several improvements to beat the quadratic footprint, existing hub

¹To avoid confusion, throughout this paper we use the term *inquiry* to indicate a request (arriving as an input in real-time to our data structure) to compute a short path $SP(s, t)$ between s and t . The term *query* refers to the act, taken by the algorithm itself, of retrieving information about a specific node. For more details on the query model we consider, see §1.2.)

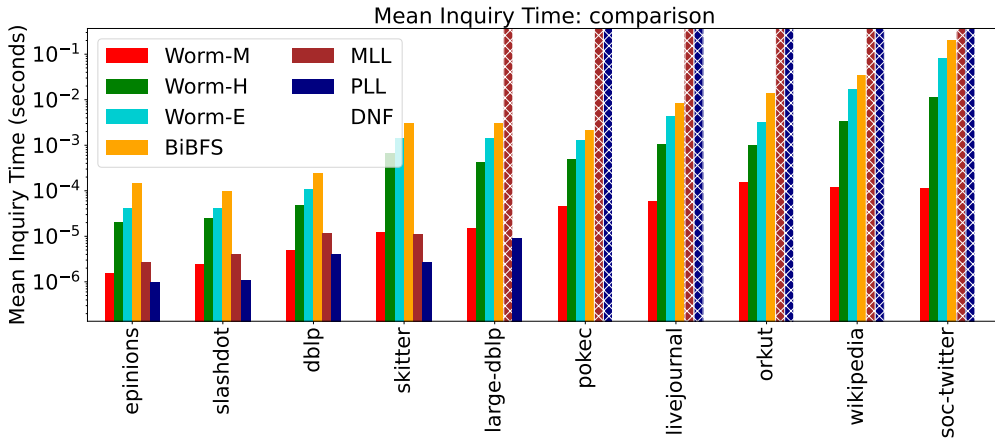


Figure 1: We illustrate the average running time per shortest path inquiry for three variants of WormHole, as compared to index-based (MLL [62] and PLL [5]), and traversal-based (BiBFS) competitors. PLL only finds distances, not paths. DNF marks that the preprocessing (index construction) step did not finish. All three of our variants outperformed BiBFS consistently. Index based solutions, on the other hand, generally failed on medium to large graphs as the index construction phase timed out. We note that even in smaller graphs where the index construction of MLL and PLL completed successfully, our fastest variant WormHole_M has comparable per-inquiry running time.

and landmark-based approaches methods continue to have high cost, and can become infeasible even for moderately-sized graphs [40].

Notably, most index-based approaches only return distances in the graph, and not the paths themselves. The first concrete systematic investigation of solutions outputting shortest paths was made by Zhang et al. [62]. Their work points out that while existing index-based solutions can be adapted to also output shortest paths, these adaptations incur a very high additional space cost on top of that required for distance computations. The authors of [62] then proposed a new approach called monotonic landmark labelling (MLL) for saving on the index construction space cost. While their algorithm is the current state of the art for this problem, it is again index-based, meaning that the preprocessing cost is still rather expensive. Improving the computational complexity of the construction phase remains a fundamental challenge.

Beyond the computational constraints, it is sometimes simply unrealistic to assume access to the whole network; examples of scenarios where access is only given via limited query access include, e.g., social network analysis through external APIs [10], page downloads in web graphs [14], modern lightweight monitoring solutions in enterprise security [26, 60], and state space exploration in software testing, reinforcement learning and robotics [27], among many others. Existing indexing-based approaches are unsuitable for these scenarios since they require reading the whole graph as a prerequisite. Traversal-based methods such as BiBFS are suitable, but as mentioned they do not scale well if one requires multiple distance computations.

The limitations of indexing-based and traversal-based methods give rise to a natural question of whether there is a middle ground solution, with preprocessing that is more efficient than in index-based approaches and inquiry time that is faster than in traversal-based approaches. Namely, we ask:

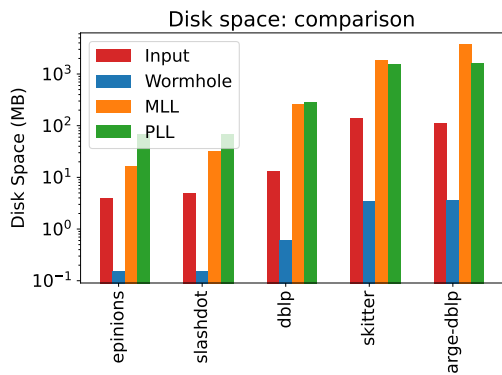
Is it possible to answer shortest-path inquiries in large networks very quickly, without constructing an expensive index, or even seeing the whole graph?

A general solution for any arbitrary graph is perhaps impossible; however, real-world social and information networks admit order and structure that can be exploited. In this work we address this question positively for a slightly relaxed version of the shortest path problem on such networks. Inspired by the core-periphery structure of large networks [43, 50, 52, 63], we provide a solution which constructs a *sublinearly-sized index* and answers inquiries by querying a strictly sublinear subset of vertices. In particular, our solution *does not need to access the whole network*. The algorithm returns *approximate* shortest paths, where the approximation error is additive and very small (almost always zero or one). In practice, the setup time is negligible (a few minutes for billion-edges graphs), and inquiry times improve on those of BiBFS. Moreover, it can be easily combined with existing index-based solutions, to further improve on the inquiry times. We also include theoretical results that shed light on the empirical performance.

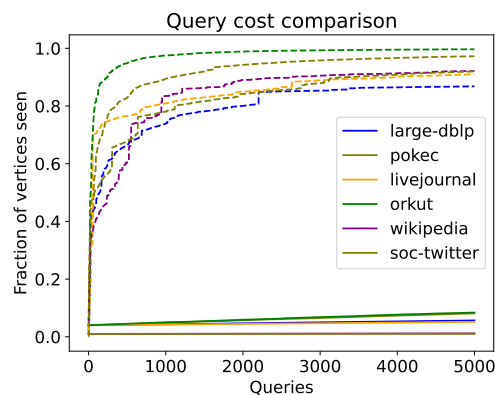
1.1 Our Contribution

We design a new algorithm, WormHole, that creates a data structure allowing us to answer multiple shortest path inquiries by exploiting the typical structure of many social and information networks. WormHole is simple, easy to implement, and theoretically backed. We provide several variants of it, each suitable for a different setting, showing excellent empirical results on a variety of network datasets. Below are some of its key features:

- **Performance-accuracy tradeoff.** To the best of our knowledge, ours is the first approximate *sublinear* shortest paths algorithm in large networks. The fact that we allow small additive error, gives rise to a trade-off between preprocessing time/space and per-inquiry time, and allows us to come



(a) Stored index size on disk for WormHole compared to indexing-based methods.



(b) Fraction of vertices seen by WormHole in large and huge graphs; see Table 1 for the classes. The dotted lines represent BiBFS and solid lines are WormHole.

Figure 2: (a) a comparison of the footprint in terms of disk space for different methods. The indexing based methods did not terminate on graphs larger than these. For WormHole, we consider the sum of C_{in} and C_{out} binary files. Note that PLL here is the *distance* algorithm, solving a weaker problem. The red bar “Input” is the size of the edge list. (b) we look at the number of vertices queried (visited) by BiBFS (dotted lines) and WormHole (solid lines) (the number is the same for all three variants). Observe that while BiBFS ends up seeing between 70% and 100% of the vertices in just a few hundred inquiries, we are well below 20% even after 5000 inquiries.

up with a solution with efficient preprocessing *and* fast per-inquiry time. Notably, our most accurate (but slowest) variant, WormHole_E, has near-perfect accuracy: more than 90% of the inquiries are answered *with no additive error*, and in *all* networks, more than 99% of the inquiries are answered with additive error at most 2. See Table 3 for more details.

- **Extremely rapid setup time.** Our longest index construction time was just two minutes even for billion-edged graphs. For context, PLL and MLL timed out on half of the networks that we tested, and for moderately sized graphs where PLL and MLL did finish their runs, WormHole index construction was $\times 100$ faster. Namely, WormHole finished in seconds while

PLL took hours. See Table 4 and Table 5. This rapid setup time is achieved due to the use of a sublinearly-sized index. For the largest networks we considered, it is sufficient to take an index of about 1% of the nodes to get small mean additive error – see Table 1. For smaller networks, it may be up to 6%.

- **Fast inquiry time.** Compared to BiBFS, the vanilla version WormHole_E (without any index-based optimizations) is $\times 2$ faster for almost all graphs and more than $\times 4$ faster on the three largest graphs that we tested. A simple variant WormHole_H achieves an order of magnitude improvement at some cost to accuracy: consistently $20\times$ faster across almost all graphs, and more than $180\times$ for the largest graph we have. See Table 3 for a full comparison. Indexing based methods typically answer inquiries in microseconds; both of the aforementioned variants are still in the millisecond regime.
- **Combining WormHole and the state of the art.** WormHole works by storing a small subset of vertices on which we compute the exact shortest paths. For arbitrary inquiries, we route our path through this subset, which we call the core. We use this insight to provide a third variant, WormHole_M by implementing the state of the art for shortest paths, MLL, on the core. This achieves inquiry times that are comparable to MLL (with the same accuracy guarantee as WormHole_H) at a fraction of the setup cost, and runs for massive graphs where MLL does not terminate. We explore this combined approach in §5.3, and provide statistics in Table 6.
- **Sublinear query complexity.** The query complexity refers to the number of vertices queried by the algorithm. In a limited query access model where querying a node reveals its list of neighbors (see §1.2), the query complexity of our algorithm scales very well with the number of distance / shortest path inquiries made. To answer 5000 approximate shortest path inquiries, our algorithm only observes between 1% and 20% of the nodes for most networks. In comparison, BiBFS sees more than 90% of the graph to answer a few hundred shortest path inquiries. See Figure 2 and Figure 5 for a comparison.

- **Provable guarantees on error and performance.** In §4 we prove a suite of theoretical results complementing and explaining the empirical performance. The results, stated informally below, are proved for the Chung-Lu model of random graphs with a power-law degree distribution [15–17].

THEOREM 1.1 (INFORMAL). *In a Chung-Lu random graph G with power-law exponent $\beta \in (2, 3)$ on n vertices, WormHole has the following guarantees with high probability:*

- Worst case error: *For all pairs u, v of nodes in G , the path between u and v output by WormHole is at most $O(\log \log n)$ longer than the shortest path.*
- Query complexity: *WormHole has preprocessing query complexity of $o(n)$, and query cost per inquiry of $n^{o(1)}$.*

In contrast, any method that does not preprocess the graph (including BiBFS) must have $n^{\Omega(1)}$ query cost per inquiry.

1.2 Setting

We consider the problem of constructing a data structure for *approximately* answering shortest-path inquiries between pairs

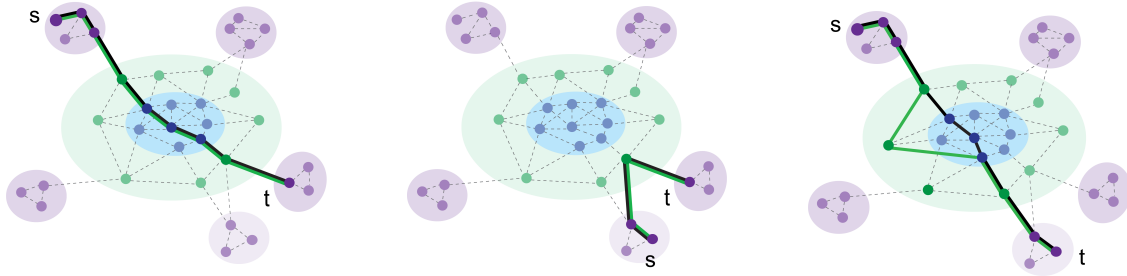


Figure 3: The decomposition and some representative cases where WormHole succeeds or fails. The central blue region is the inner ring C_{in} , the green layer outside is the outer ring C_{out} , and purple regions attached to this are the peripheral components forming \mathcal{P} . Dashed lines are edges. Vertices labelled s and t are respectively the source and destination. The green lines are actual shortest paths, while the black lines are paths output by WormHole. We ignore the case where both source vertices are in the same peripheral component. The first one (a) is the case where the shortest path and the path output by WormHole are identical; no error is incurred in this case. The second (b) is the case where the source and destination are in two different peripheral components, but they encounter a common vertex while traversing to the inner ring. The third (c) is an example of a case where we incur an error: the shortest path(s) interleaves through the outer ring C_{out} , so by restricting the traversal solely to the interior of C_{in} , we incur an error, in this case of 1.

of vertices (s, t) in an undirected graph G , given limited query access to the graph.

Query model. Access to the network is given through the standard node query model [10, 14], where we start with an arbitrary seed vertex as the “access point” to the network, and querying a node v reveals its list of neighbors $\Gamma(v)$. Unlike existing index-based solutions, which perform preprocessing on the whole graph, we aim for a solution that queries and stores only a small fraction of the nodes in the network.

Objective. Following the initialization of the data structure, the task is to answer multiple shortest path inquiries, where each inquiry $SP(s, t)$ needs to be answered with a valid path $p_0 p_1 \dots p_\ell$ between $s = p_0$ and $t = p_\ell$, and the objective is to minimize the mean additive error measured over all inquiries. The additive error for an inquiry $SP(s, t)$ is the difference between the length of the returned $s-t$ path and the actual shortest distance between s and t in G . Depending on the specific application, one would like to minimize (a subset of) the additive error, running time, memory and/or node queries.

Core-periphery structure. The degree distribution in social and information networks often follows a power-law distribution with exponent $2 < \beta < 3$, which results in a *core-periphery* structure [9, 43, 50, 52, 63], where the core is a highly connected component with good expansion properties, consisting of higher degree nodes, while the periphery is a collection of small, poorly connected components of low degree.

Our data structure is designed for networks exhibiting these structural characteristics. It takes advantage of the structure by first performing a preprocessing step to acquire (parts of) the core of the network, and then answering approximate shortest path inquiries by routing through the core. The working hypothesis is that pairs of nodes that are sufficiently far apart will

typically have the shortest path between them (or close to it) routed through the higher degree parts of the network. This is somewhat reminiscent of approaches based on the highway dimension [1–3] for routing in road networks, although the structural characteristics of these network types differ considerably.

1.3 The algorithm

WormHole builds an explicit hierarchical core-periphery type structure with a sublinear inner ring and provides a framework which uses this structure to answer shortest path inquiries. There are two phases:

- A preprocessing step where we decompose the graph into three partitions, storing only the smallest one: a highly dense subgraph of *sublinear size*.
- The phase where we answer inquiries: here the algorithm (approximately) answers shortest path inquiries of the form $SP(s, t)$ for arbitrary vertex pairs (s, t) .

We elaborate on the two phases.

1.3.1 The decomposition. It is well-documented that social networks exhibit a core-periphery structure; see, e.g., [43, 50, 52, 63] and the many references within. The *core* is a highly-connected component with good expansion properties and smaller effective diameter. The *periphery*, denoted \mathcal{P} , consists of smaller isolated communities that connect to the core, but are sparsely connected internally, and whose union is of linear size [16]. Therefore, when answering shortest path inquiries, it is reasonable to first check if the two vertices are in the same peripheral community, and otherwise route through the core.

A recent work of Ben-Eliezer et al. [10] observed that a more fine-grained constructive version of the core-periphery decomposition may be useful for algorithmic purposes in real world networks. In their work, the core is further decomposed into

two layers: a sublinearly-sized *inner ring*, denoted C_{in} , which is very dense and consists of the highest degree vertices in the graph; and its set of neighbors, which we refer to as the *outer ring*, C_{out} , where $C_{out} = \Gamma(C_{in}) \setminus C_{in}$. Their work shows (empirically) that even a sublinearly-sized inner ring is sufficient so that the union of the inner and outer rings effectively contains the core. Our work makes use of this fine-grained core-periphery decomposition. During preprocessing, we acquire the nodes of the inner core, which then constitutes the (sublinear) index of our data structure.

1.3.2 Acquiring the inner core. To acquire the inner core C_{in} we follow the procedure by [10]. The procedure gradually expands C_{in} , starting from an arbitrary seed vertex, and iteratively adding vertices with high connectivity to the current core. See Algorithm 1 for the pseudo code and Figure 4 for an illustration of an expansion step.

1.3.3 Answering shortest-path queries. In the second phase, given a query $SP(s, t)$, WormHole does the following. First, it checks if the two vertices are in the same peripheral component, by performing a truncated BiBFS from both s and t up to depth two. If the two trees collide, it returns the shortest path between s and t . Otherwise, WormHole continues both BFS traversals until it reaches the outer ring (from both s and t). From here, it takes a single step to reach the inner ring, and then performs a restricted BiBFS on the subgraph induced by the inner ring vertices. We note that the choice of BiBFS here is arbitrary, and we can use any shortest-path algorithm (including modern index-based approaches, initialized only on the inner core) as a black-box to find a shortest path in the inner ring.

Figure 3 illustrates a few typical cases encountered by the algorithm; in the first two cases the algorithm returns a true shortest path, and in the third case the returned path is not a shortest path (thus incurring a nonzero additive error).

We stress that a single decomposition is subsequently used to answer *all* shortest path queries. Theorem 1.1 provides a strong theoretical guarantee on the performance of WormHole. It is worth emphasizing that our notion of approximation is inspired by practical relaxations, and is distinct from the one usually considered in theoretical works.

2 RELATED WORK

There is a vast body of research on shortest paths and distances, spanning over many decades, and including hundreds of algorithms and heuristics designed for a variety of settings. Here, we only review several works that are most closely related to ours. For a more comprehensive overview, see the surveys [42, 55] and references therein.

Index-based approaches. As mentioned earlier, a ubiquitous set of algorithms are based on landmark/sketch approaches [37, 61, 62], with Pruned Landmark Labeling (PLL) [5] being perhaps the most influential one. These algorithms follow a two-step procedure: the first step generates an ordering of the vertices according to importance (based on different heuristics), and the second step generates labeling from pruned shortest path trees constructed according to the ordering. Then the shortest distance between an arbitrary pair of vertices s and t

can be answered quickly based on their labels. However, even with pruning, PLL requires significant setup time. Hence, there have been many attempts to parallelize it [29, 37, 39].

Embedding based approaches. Some recent approaches leverage embeddings of graphs to estimate shortest paths. Like in representation learning, they seek to find efficient representations of distances between pairs of nodes [64, 65]. A modern line of work also considers hyperbolic embeddings of the graphs, that are closely related to tree decompositions, to answer shortest path inquiries [11, 33]. Recent work has also looked at accelerating this process by using GPU based deep learning methods [35, 47, 48]. Query-by-Sketch [57] considers the, related but incomparable, task of answering *shortest-path-graph inquiries*, where the goal is to compute a subgraph containing exactly all shortest paths between a given pair of vertices. They propose an alternative labeling scheme to improve the scalability and inquiry times.

BFS-based approaches. Another set of algorithms is based on Breadth First Search (BFS) or Bidirectional Breadth First Search (BiBFS), as they are exact methods with no preprocessing step. There has been substantial work in the last few years proving that BiBFS is sublinear (i.e., proving complexity upper bounds of the form $n^{1-\Omega(1)}$) with high probability for several graph families. These include, e.g., hyperbolic random graphs (Blasius et al., [11]), and graphs with a finite second moment and power-law graphs (Borassi and Natale [12]). Very recently, Alon et al. [8] showed that BiBFS has sublinear query complexity for a broad family of expander graphs. Their work also proves query lower bounds of the form $n^{\Omega(1)}$ for traversal-based algorithms in Erdős-Renyi and d -regular random graphs.

Core-periphery based approaches. Several other works exploit the core-periphery structure of networks [6, 13, 38]. Brady and Cohen [13] use compact routing schemes to design an algorithm with small additive error based on the resemblance of the periphery to a forest. Akiba, Sommer and Kawarabayashi [6] exploit the property of low tree-width outside the core, and in [38], the authors design a Core-Tree based index in order to improve preprocessing times on massive graphs. We note that in all these results, the memory overhead is super-linear.

Worst case graphs. On the theoretical side, there have been many results on exact and approximate shortest paths in worst-case graphs, e.g., [19, 20, 22, 49, 51, 67], with improvements as recently as the past year. Most of these focus on the 2-approximation case, first investigated in the seminal work of Aingworth et al. [4]. We point the readers to Zwick’s survey [66] on exact and approximate shortest-path algorithms, and Sen’s survey [53] on distance oracles for general graphs with an emphasis on lowering pre-processing cost. Notably, because we make a beyond worst case structural assumption that is common in large networks, namely, a core-periphery structure, both our results and algorithm substantially differ from the worst-case theoretical literature.

3 ALGORITHM

WormHole utilizes insights about the structure of real world networks to cleverly decompose the graph and calculate approximate shortest paths. We discuss the algorithm and various steps in detail in this section; our two main components are a sublinear decomposition procedure, adapted from the recent work of Ben-Eliezer et al [10], and a routing algorithm that takes advantage of this decomposition to find highly accurate approximate shortest paths.

3.1 The Structural Decomposition Phase

The former is a simple implementation of a structural decomposition provided in [10] that gives us access to the inner ring; the inner ring is central to our procedure, and the bulk of our computation is done on it. This decomposition stratifies our graph into three sections: the inner ring, C_{in} , a dense component with the highest degree nodes, the outer ring, C_{out} , the set of neighbors of the inner ring, and the rest of the vertices which form the periphery, \mathcal{P} , and typically reside in fragments of size $O(\log n)$ [41]. The construction of the inner ring is an iterative procedure that captures the highest degree vertices in sublinear time; we shall refer to this process as the CoreGen procedure. The procedure gradually expands the inner ring, starting from an arbitrary seed vertex. The outer ring is the set of neighbors of the inner ring vertices (that are not already in the inner ring). At each step, the procedure removes, from C_{out} , the vertex with highest number of neighbors in C_{in} and adds it to C_{in} . It then queries this vertex and adds its neighbors to C_{out} , while keeping track of how many neighbors C_{out} vertices have in the updated C_{in} . See Figure 4 for the first few steps of the inner ring generation procedure. The size of inner ring, in terms of a percentage of the vertex set, is given to the core-generation algorithm as a parameter and the process is iterated till that size is reached. Algorithm 1 has the pseudocode for CoreGen. Here, $\Gamma(v)$ refers to the neighbors of a vertex v .

We emphasize that *the decomposition is performed only once*, and all subsequent operations are done using this precomputed decomposition.

Algorithm 1 CoreGen(v, s (< 1), query access to $G = (V, E)$): starting from a seed vertex v , generates an inner ring of size $s|V|$; returns an inner ring, C_{in} , and an outer ring, C_{out} .

```

1: size = 0
2:  $C_{in} \leftarrow v, C_{out} \leftarrow \Gamma(v)$ 
3: while size <  $s|V|$  do
4:   Pick  $u \in C_{out}$  with maximum number of neighbors in
    $C_{in}$  (break ties randomly)
5:    $C_{in} = C_{in} \cup \{u\}, C_{out} = C_{out} \cup \Gamma(u) \setminus C_{in}$ 
6:   size += 1
7: end while
8: return  $C_{in}, C_{out}$ 

```

3.2 The Routing Phase

The approach is simple: assume the preprocessing phase acquires the inner ring. Upon an inquiry (s, t) , the algorithm

Algorithm 2 WormHole(s, t, C_{in}, C_{out}, G): final algorithm

```

1:  $T(s) = \Gamma(s), T(t) = \Gamma(t)$ 
2:  $C(s), C(t) \leftarrow \emptyset$ 
3: while  $T(s) \cap T(t) = \emptyset$  or if both  $C(s), C(t) = \emptyset$  do
4:   for both  $s$  and  $t$  (denoted by  $i$  in the next block) do
5:     if  $C(i) = \emptyset$  and  $\Gamma(T(i)) \cap C_{in} \neq \emptyset$  then
6:       Expand  $T(i)$  by a single level
7:        $C(i) \leftarrow \Gamma(T(i)) \cap C_{in}$ 
8:       stop further expansion of  $T(i)$ 
9:     else if  $C(i) = \emptyset$  then
10:      Expand  $T(i)$  by a single level
11:     end if
12:   end for
13: end while
14:
15: if  $T(s) \cap T(t) \neq \emptyset$  then
16:   return SP( $s, t$ ) through  $T(s) \cap T(t)$ 
17: else
18:   Expand the BiBFS tree from  $C(s), C(t)$  in  $G[C_{in}]$ 
19:   return SP( $s, t$ ) through  $T(s) \cup G[C_{in}] \cup T(t)$ 
20: end if

```

starts two BFS trees from both s and t . It then expands the tree at each step till one of the followings happens:

- (1) the two trees intersect, or
- (2) the trees reach the outer ring.

If the search trees in the former do not intersect, WormHole mandatorily routes shortest paths through the inner ring. Once in the inner ring, it computes the exact shortest path through it.

3.3 Variants of WormHole

In the default variant of Algorithm 2, WormHole_E, we use the bidirectional breadth-first BiBFS shortest path algorithm as a primitive in order to compute the shortest path between two inner ring vertices; see [46] for a full description of BiBFS. The theoretical analysis in §4 considers this variant.

WormHole_H is a simple variant of Algorithm 2 where we start the expansion of the BiBFS trees in $G[C_{in}]$ only from the vertex $v_s = \arg\max_{v \in C(s)} \text{deg}(v)$ (and similarly v_t for $C(t)$, the highest degree vertices). We note that for this variant, some of the theoretical results in §4 no longer hold: while WormHole_H satisfies all sublinearity properties, the error bounds in Theorem 4.4 no longer hold, as this variant only computes an approximate shortest path inside the inner ring. An empirical comparison of the two is done in §5.1 and Table 3.

In WormHole_M we combine WormHole with an index-based algorithm restricted to the core. While indexing-based algorithms are quite expensive, in terms of both preprocessing and space cost, they lead to much lower times per inquiry. Therefore this variant is suitable when faster inquiry times are preferable to low space requirements. WormHole_M makes the index creation cost substantially lower (compared to generating it for the entire graph) while providing speedups for answering shortest path inquiries compared to the BiBFS implementation. We discuss this briefly in §5.3 but leave a complete systematic exploration of these options for future research.

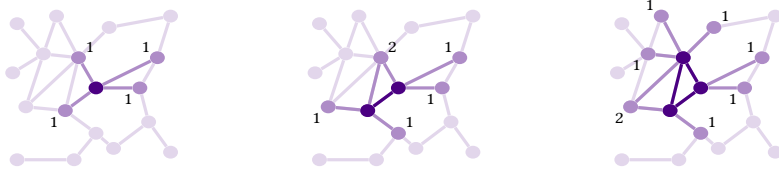


Figure 4: Construction of the inner and outer rings of the core; figure taken from [10]. At any point, the outer ring is the set of vertices adjacent to the inner ring. The algorithm expands the inner ring by adding to it a vertex from the outer ring that has the most neighbors in the inner ring. The image looks at two successive steps: the numbers labelling vertices in the outer ring refer to how many inner ring vertices it is adjacent to. Thus, in the second step, the vertex labelled 2 is added to the inner ring.

4 THEORETICAL ANALYSIS

It is a relatively standard observation that many social networks exhibit, at least approximately, a power-law degree distribution (see, e.g., [9] and the many references within). The Chung-Lu model [41] is a commonly studied random graph model which admits such degree distribution.

In this section we provide a proof-of-concept for the correctness of WormHole on Chung-Lu random graphs, aiming to explain the good performance in practice through the study of a popular theoretical model. We sometimes only include proof sketches instead of full proofs, in the interest of saving space.

4.1 Preliminaries

We start by defining power-law networks and the Chung-Lu model [15–17], followed by a set of useful results from Lu [41]. A network is said to have a power-law degree distribution when for sufficiently large k , the fraction of nodes with degree k , denoted $p(k)$, follows a power-law. That is, $p(k) \propto k^{-\beta}$. Multiple studies have shown that real world graphs indeed follow (or approximately follow) a power-law degree distribution, and the exponent β is typically in the range $2 < \beta < 5$ or so [7, 23, 44, 54]. Theoretically, the regime $2 < \beta < 3$ is the most interesting, and thus we focus our study on this regime.

In a series of works, Chung and Lu suggested a model to generate graphs according to a power-law degree distribution, in order to generate massive graphs that capture properties of real-world graphs [15–17]. In this model, one is given a list of n expected degrees, $w_1 \geq w_2 \geq \dots \geq w_n$, where in our case we assume that the expected degrees follow a power-law degree distribution with $2 < \beta < 3$. Then, for each pair of nodes i, j , the edge between them is instantiated with probability $\frac{w_i \cdot w_j}{W}$, where $W = \sum_{i=1}^n w_i$ (and it is assumed that for any i, j , $w_i \cdot w_j \leq W$). Note that indeed by this definition, it holds that for every vertex v , $E[d(v)] = w(v)$, and therefore that $E[|E(G)|] = W/2$. We denote this distribution on graphs as \mathcal{G}_{CL} , and denote a graph drawn from it by $G \sim \mathcal{G}_{CL}$.

In all that follows we assume that $G \sim \mathcal{G}_{CL}$, and the probabilities are taken over the generation process. We let $\gamma = 1/\log \log n$, and let C_{CL} denote the set of vertices with degree at least $t = n^\gamma$.

THEOREM 4.1 (THEOREM 4 IN [16]). *Suppose a power law random graph with exponent $2 < \beta < 3$, average degree d strictly greater than 1, and maximum degree $d_{max} > \log n / \log \log n$. Then almost surely the diameter is $\Theta(\log n)$, the diameter of*

the C_{CL} core is $O(\log \log n)$ and almost all vertices are within distance $O(\log \log n)$ to C_{CL} .

Claim 4.2 (Fact 2 in [15]). *With probability at least $1 - 1/n$, for all vertices $v \in V$ such that $w(v) \geq 10 \log n$*

$$d(u) \in \left[\frac{1}{2} w(u), 2w(u) \right] \text{ and } |E(G)| \in [W/4, W],$$

and for all vertices with $w(v) \leq 10 \log n$, $d(v) \leq 20 \log n$.

4.2 Sublinearity of Inner Ring

Chung and Lu proved [15–17] that in random graphs with a power-law degree distribution and exponent $2 < \beta < 3$, there exists a core with small diameter, so that most of the vertices in the graph are close to it. Here we extend their result to show that this core is of sublinear size and that CoreGen indeed captures it; i.e., $C_{CL} \subseteq C_{in}$.

THEOREM 4.3 (WormHole INNER RING \supseteq CHUNG-LU CORE). *Consider a graph G generated according to the Chung-Lu process with $\beta > 2$, and suppose we run WormHole on G with an inner ring of size $n^{1-\Theta(1/\log \log n)}$. Then with high constant probability, C_{in} contains all vertices with degree greater than $n^{1/\log \log n}$.*

PROOF SKETCH. Consider a process P which simultaneously constructs G as the core grows as outlined in CoreGen: Initially the graph G has no edges, and C_{in} is an empty set. Every time a vertex v is added to C_{in} , P iterates over all nodes $u \in V \setminus C_{in}$ and draws a coin with probability $\frac{w(v) \cdot w(u)}{2W}$ to determine whether the edge (u, v) is in E or not. Once C_{in} is acquired (i.e., when C_{in} is of size $t = 100 \log n \cdot n^{1-1/\log \log n}$), P iterates over all pairs of vertices that were not yet examined and determines their edges. Clearly this process generates a graph according to the distribution $G \sim \mathcal{G}_{CL}$ (since we only changed the order in which the edges were decided).

Let d denote the expected average degree in the graph, and let C_{in}^ℓ denote the set of vertices added to C_{in} after ℓ steps of the above process. Finally, let $\gamma = 1/\log \log n$. It is easy to show that with high probability, except for a negligible fraction, all vertices added to the core have expected degree $w(V) \geq d$. Conditioned on this event, for every u with $w(u) \geq n^\gamma$, every time a vertex v is added to the core, v creates an edge with u with probability

$$\frac{d(u) \cdot d(v)}{W} \geq \frac{n^\gamma \cdot d}{nd} = \frac{1}{n^{1-\gamma}}.$$

Let χ_i denote the event that the i -th vertex added to C_{in} creates an edge with u . Then by the above, $E[\chi_i] \geq \frac{1}{n^{1-\gamma}}$. Observe

that $d_{C_{in}}(u) = \sum_i \chi_i = \frac{|C_{in}^\ell|}{n^{1-\gamma}}$, and set $\mu := \frac{|C_{in}^\ell|}{n^{1-\gamma}}$. Note that the χ_i variables are independent $\{0, 1\}$ random variables, and that for $\ell = \text{clogn} \cdot n^{1-\gamma}$, $\mu \geq \text{clogn}$. Therefore, by the multiplicative Chernoff bound,

$$\Pr \left[\left| d_{C_{in}^\ell}(u) - \mu \right| > \frac{1}{4} \mu \right] < 2 \exp \left(\frac{-\frac{1}{16} \cdot \mu}{3} \right) \leq \frac{1}{n^2}.$$

That is, we get that with probability $1 - 1/n$, all vertices with expected degree $w(u) \geq n^\gamma$ have $d_{C_{in}}(u) \geq \frac{3}{4} n^\gamma$. A similar analysis can be used to prove that all vertices u with expected degree $w(u) < n^\gamma/8$ will have $d_{C_{in}}(u) \leq \frac{1}{2} n^\gamma$ with high probability. Hence, after performing additional ℓ steps, we will have that with high probability, all vertices with expected degree at least n^γ will be added to $C_{in}^{2\ell}$ before any vertex with expected degree less than $n^\gamma/8$. The proof concludes by noticing that by Lemma 4.2, the degrees of all vertices will either be as expected up to a constant multiplicative factor, or have an expected degree below, say $n^\gamma/8$, and actual degree not higher than $n^\gamma/2$. \square

4.3 Approximation Error

Now that we have a sublinear inner ring that contains the Chung-Lu core, we must show that routing paths through it incurs only a small penalty. Intuitively, the larger the inner ring, the easier this is to satisfy: if the inner ring is the whole graph, the statement holds trivially. Therefore the challenge lies in showing that we can achieve a strong guarantee in terms of accuracy even with a sublinear inner ring. We prove that `WormHoLe` incurs an additive error at most $O(\log \log n)$ for all pairs, which is much smaller than the diameter $\Theta(\log n)$.

THEOREM 4.4 (GOOD ADDITIVE ERROR). *If $C_{CL} \subseteq C_{in}$, then with high probability, for all pairs (s, t) of nodes, the additive error of our algorithm for the inquiry $\text{SP}(s, t)$ is at most $O(\log \log n)$.*

The above result holds with high probability even in the worst case. Namely, for all pairs (s, t) of vertices in the graph, the length of the path returned by `WormHoLe` is at most $O(\log \log n)$ higher than the actual distance between s and t . This trivially implies that the average additive error of `WormHoLe` is, with high probability, bounded by the same amount.

PROOF. Let $d(s, C_{CL})$ and $d(C_{CL}, t)$ denote the distances of s and t to the core. If $d(s, t) \leq d(s, C_{CL}) + d(C_{CL}, t)$, then the BFS trees from both sides will intersect without the need to go through the core C_{CL} .

For pairs (s, t) where there is a shortest path going through C_{CL} , the length of this shortest path is

$$d(s, u) + d(u, v) + d(v, t)$$

where u and v are in the core, whereas the path that `WormHoLe` outputs is of length at most

$$d(s, C_{CL}) + \text{diam}(C_{CL}) + d(C_{CL}, t).$$

By definition, $d(s, C_{CL}) \leq d(s, u)$ and $d(C_{CL}, t) \leq d(v, t)$. The proof follows since the diameter of C_{CL} is bounded by $O(\log \log n)$ with high probability (see Theorem 4.1). \square

4.4 Query Complexity

Recall the node query model in this paper (see §1.2): starting from a single node, we are allowed to iteratively make queries, where each query retrieves the neighbor list of a node v of our choice. We are interested in the query complexity, i.e., the number of queries required to conduct certain operations.

Using the inner ring as our index, and routing the shortest paths through it, we get an algorithm with small additive error and sublinear setup query cost: we prove that our query cost per inquiry is also subpolynomial (i.e., of the form $n^{o(1)}$, where the $o(1)$ term tends to zero as $n \rightarrow \infty$). Moreover, we prove that preprocessing is necessary to achieve such query complexity per inquiry; anything else *must* incur a cost of $n^{\Omega(1)}$.

The first result is the upper bound on our performance.

THEOREM 4.5 (SUBPOLYNOMIAL QUERY COMPLEXITY FOR SHORTEST PATHS). *Suppose that the preprocessing phase of our algorithm acquires C_{in} . The average query complexity of computing a path between a pair of vertices is bounded by $n^{\Theta(1/\log \log n)}$.*

PROOF SKETCH. For a given inquiry $\text{SP}(u, v)$, we give an upper bound on the query complexity of the BFS that starts at u , and similarly for v ; the total query complexity is the sum of these two quantities.

Let C denote the subset of vertices obtained by the algorithm during preprocessing. If $u \in C_{in}$, then the algorithm only performs a BiBFS restricted on $G[C_{in}]$ to another vertex $v \in C_{in}$. Since C_{in} is queried during the preprocessing phase, this step requires no additional queries. The non-trivial case is when the source u is not in C_{in} . In this case, the algorithm performs a BFS until it either collides with the simultaneous BFS that is taking place from the other vertex in the inquiry, or until it reaches C_{in} , at which point it performs a BiBFS inside C_{in} .

Fix some vertex $u \notin C_{in}$. By the definition of the Chung-Lu model, when performing a walk in the graph from some vertex v , for every step in the walk, the probability it reaches C_{in} is at least

$$p = \sum_{w \in C_{in}} \frac{d(w)}{2m} = \frac{\text{Vol}(C_{in})}{2m} \geq \left(\frac{|C_{in}| \cdot d}{2n \cdot d} \right) = \frac{1}{n^{\Theta(1/\log \log n)}}. \quad (1)$$

Therefore, after expanding the BFS tree from u for $O(1/p)$ many times, we reach C_{in} with high constant probability. Hence, the BFS from u requires $\Theta(1/p) = n^{\Theta(1/\log \log n)}$ many queries until reaching C_{in} . This concludes the proof in the case that all shortest paths between u and v contain at least one edge in C_{in} .

It remains to consider the case in our algorithm where the searches from u and v collide outside C_{in} , or reach the same vertex in C_{in} . This requires us, at most, to query all vertices encountered during the BFS until it reaches C_{in} for the first time. By Theorem 4.3, the degree of all vertices outside C_{in} is bounded by $O(n^{1/\log \log n})$. Therefore the overall query complexity of the walk outside the core is $n^{1/\log \log n} \cdot n^{\Theta(1/\log \log n)} = n^{\Theta(1/\log \log n)}$. \square

Finally, this brings us to the lower bound. We prove that any method for finding a path between two nodes u and v in a Chung-Lu random graph that does not employ preprocessing, requires $n^{\Omega(1)}$ node queries to succeed with good probability. Due to space considerations we only provide here a proof

sketch; for a more complete proof of similar results for Erdős-Renyi and other random graphs, see Alon et al. [8].

THEOREM 4.6 (POLYNOMIAL QUERY COMPLEXITY WITHOUT PREPROCESSING). *Fix $2 < \beta < 3$. Let G be a graph generated by the Chung-Lu process with power law parameter β and average expected degree d (possibly depending on n , but satisfying $d = n^{o(1)}$). Let $u \neq v$ be a random pair of nodes in G . Any algorithm that receives node query access to G starting at u and v must make, with high probability, $n^{\Omega(1)}$ queries to G in order to find a path between u and v .*

PROOF SKETCH. Consider a Chung-Lu graph with parameters β, d , and suppose that n is large enough. The highest expected degree of a vertex in the graph is bounded by n^{1-c} for a constant c that may depend on β and d (but not on n).

For a given time step of the algorithm, let A_u denote the component consisting of all nodes queried by virtue of being reached from u in the algorithm, and all edges in G intersecting at least one such node. Define A_v similarly for v . Consider one step of the algorithm where some vertex w is added to A_u , together with all new (previously unseen) edges incident to w . Let e be any such new edge. The probability that e intersects A_v (and thereby closes the desired path between u and v) is bounded by

$$\sum_{x \in A_u} \frac{w(x)}{2W} \leq |A_v| \cdot \frac{n^{1-c}}{n} = \frac{|A_v|}{n^c},$$

where $|A_v|$ is the number of edges in the component A_v . Union bounding over all edges e added in this process, we get that the probability of finding a path is bounded by $|A_u| \cdot |A_v| \cdot n^{-c}$. It follows that to find a path with good probability, at least one of A_u and A_v needs to have at least $n^{c/3}$ edges. The rest of the proof sketch is devoted to proving this last claim.

We use the following fact on Chung-Lu graphs with the relevant parameters. Fix $\alpha > 0$. The probability that a new edge (emanating from a newly queried vertex w) will hit, at its other endpoint, a vertex of expected degree at least n^α is at most $n^{-\gamma}$, where γ is a constant that depends only on β, d, α , and $\gamma \in (0, \alpha)$.

Pick $\alpha = c/6$ and the corresponding $\gamma < \alpha$, and suppose we make an arbitrary traversal involving $n^{Y/2}$ queries starting at u . The probability that any specific query hits a node of degree at least n^α is bounded by $n^{-\gamma}$. Union bounding over all $n^{Y/2}$ queries, we conclude that the event that an n^α -degree node is queried throughout the whole traversal is bounded by $n^{Y/2} \cdot n^{-\gamma} = n^{-Y/2}$. Conditioning on this event not happening, the size of A_u at the end of the traversal is bounded by $n^{Y/2} \cdot n^\alpha < n^{c/3}$, as desired.

The reasoning in the last paragraph implicitly assumes that outgoing edges from A_u will never intersect A_u itself. This is true with high probability by essentially the same argument, since the probability of each edge e to intersect A_u is bounded by $|A_u|/n^c < n^{-2c/3}$. \square

5 EXPERIMENTAL RESULTS

In this section, we experimentally evaluate the performance of our algorithm. We look at several metrics to evaluate performance in different aspects. We compare with BiBFS, a traversal-based approach, and with the indexing algorithms

PLL and MLL. We test several aspects, summarized next. Detailed results are provided in the rest of this section.

- (1) **Query cost:** By query cost, we refer to the number of vertices queried by WormHole, consistent with our access model (see §1.2). We show that WormHole actually does remarkably well in terms of query cost, seeing a small fraction of the whole graph even for several thousands of shortest path inquiries. See Figures 2(b) and 5.
- (2) **Inquiry time:** We demonstrate that WormHole_E achieves consistent speedups over traditional BiBFS, even while using it as the sole primitive in the procedure. More complex methods such as PLL and MLL time out for the majority of large graphs. We also provide variants that achieve substantially higher speedups. Finally, in §5.3, we show how using the existing indexing-based state or the art methods on the core lets us achieve indexing-level inquiry times. See Figure 1.
- (3) **Accuracy:** We show that our estimated shortest paths are accurate up to an additive error 2 on 99% of the inquiries for the default version WormHole_E; a faster heuristic, WormHole_H, shows lower accuracy, but still over 90% of inquiries satisfy this condition. See §5.1 and Table 3 for details.
- (4) **Setup:** We look at the setup time and disk space with each associated method. Perhaps as expected, WormHole_E beats the indexing based algorithms by a wide margin in terms of both space and time: see Figure 7. In §5.3 we further show that using these methods restricted to C_{in} results in a variant WormHole_M with much lower setup cost (Table 6).

Datasets. The experiments have been carried out on a series of datasets of varying sizes, as detailed in Table 2. The datasets have been taken either from the SNAP large networks database [36] or the KONECT project [34]. We organize the results into two broad sections: we first introduce two variants

Class	$ C_{in} $	Networks
small	6%	dblp, epinions, slashdot, skitter
med	4%	large-dblp, pokec, livejournal, orkut
large	1%	wikipedia, soc-twitter

Table 1: Classification of networks used in experiments by size of C_{in} used in experiments.

Network	$ V $	$ E $	BiBFS	PLL	MLL
epinions	$7.6 \cdot 10^4$	$5.1 \cdot 10^5$	✓	✓	✓
slashdot	$7.9 \cdot 10^4$	$5.2 \cdot 10^5$	✓	✓	✓
dblp	$3.2 \cdot 10^5$	$1.0 \cdot 10^6$	✓	✓	✓
skitter	$1.7 \cdot 10^6$	$1.1 \cdot 10^7$	✓	✓	✓
large-dblp	$1.8 \cdot 10^6$	$2.9 \cdot 10^7$	✓	✓	✗
soc-pokec	$1.6 \cdot 10^6$	$3.1 \cdot 10^7$	✓	✗	✗
soc-live	$4.8 \cdot 10^6$	$6.8 \cdot 10^7$	✓	✗	✗
soc-orkut	$3.1 \cdot 10^6$	$1.2 \cdot 10^8$	✓	✗	✗
wikipedia	$1.4 \cdot 10^7$	$4.4 \cdot 10^8$	✓	✗	✗
soc-twitter	$4.2 \cdot 10^7$	$1.5 \cdot 10^9$	✓	✗	✗

Table 2: Network datasets used for experimental evaluation with their corresponding sizes. We observe that BiBFS finishes on all the datasets, but the indexing based methods do not on the medium and large networks. We were able to set up MLL on large-dblp in reasonable time, but the subsequent shortest path inquiries were met with consistent segmentation faults that we were unable to debug.

Network	BiBFS	WormHole _E					WormHole _H				
	MIT	MIT	SU/I	+0(%)	≤+1 (%)	≤+2 (%)	MIT	SU/I	+0(%)	≤+1 (%)	≤+2 (%)
epinions	144	41	4.5	98.06	99.99	100.00	20	24	66.97	99.54	100.00
slashdot	99	46	2.8	73.43	95.37	99.28	24	14	63.09	98.78	99.98
dblp	247	110	2.4	97.02	99.96	100.00	48	11	44.72	82.42	96.53
skitter	3004	1439	2.3	94.71	99.89	100.00	660	24	58.99	96.78	99.98
large-dblp	3041	1447	2.3	85.37	99.10	99.95	417	21	47.61	89.74	99.04
pokec	2142	1317	1.8	51.37	92.15	99.63	506	11	14.52	59.51	90.71
livejournal	8565	4318	2.1	71.98	97.95	99.86	1054	29	28.86	77.93	97.83
orkut	14k	3213	4.4	58.50	94.56	99.64	1030	35	20.66	68.11	93.93
wikipedia	35k	17k	2.4	94.94	99.92	100.00	3394	36	44.65	98.74	100.00
soc-twitter	204k	81k	3.4	93.30	99.98	100.00	12k	181	35.13	99.30	99.99

Table 3: Summary of WormHole with the two cases: WormHole_E, with the exact shortest path through the inner ring, and WormHole_H that picks only the shortest path between the highest degree vertices – refer to §5.1. We note the mean inquiry times per inquiry (MIT) in microseconds, and average speed up *per inquiry* (SU/I) compared to BiBFS for each method. We also note the percentiles of inquiries by absolute error: for WormHole_E, we get absolute error under 2 for over 99% of the inquiries. This drops for WormHole_H, but it is still above 99% for six of the ten datasets, and over 90% in all of them. Accuracy numbers are highlighted in green, where darker is better. Similarly, we have a gradient of violet for speedups; darker is faster. For WormHole_E, speedup over BiBFS per inquiry on average is usually between 2× and 3×, but this increases to consistently between 20–30× in WormHole_H, and reaches a max of 181× in our largest dataset, soc-twitter.

of our algorithm. We then compare it with BiBFS as well as indexing based methods – PLL and MLL. The latter two did not terminate in 12 hours for most of the graphs, while BiBFS completed on even our largest networks.

We classify the examined graphs into three different classes and use a fixed percentage as the ‘optimal’ inner ring size for graphs of comparable size (where the inner core size as % of the total size decreases for larger networks, an indication for the sublinearity of our approach). This takes into account the tradeoff between accuracy and the query/memory costs incurred by a larger inner ring. The classification is summarized in Table 1. For the experimental section, we default to these sizes unless mentioned otherwise.

Implementation details. We run our experiments on an AWS ec2 instance with 32 AMD EPYC™ 7R32 vCPUs and 64GB of RAM. The code is written in C++ and is available in the supplementary material as a zipped folder, with links to the datasets. The backbone of the graph algorithms is a subgraph counting library that uses compressed sparse representations [45].

5.1 WormHole_E, WormHole_H and BiBFS

We run two separate versions of WormHole: the one, as described in Algorithm 2, is what we refer to as WormHole_E (exact). Another variant we consider is where we pick just the highest degree vertices from $C(s)$ and $C(t)$ respectively, and do a BFS from those. This cuts down on the inquiry time, but also reduces the accuracy. We call this variant WormHole_H. We take a deeper look into this tradeoff in the following sections. Note that both of these methods have the same query cost per shortest path inquiry, and the only difference is in running time (as all vertices in C_{in} have already been queried during CoreGen). For all runs, we do approximately 10,000 inquiries of uniformly chosen source and destination pairs (discarding disconnected pairs and other invalid inquiries). Our program runs on a compressed sparse representation (CSR) graph, and we store binary arrays for both C_{in} and C_{out} for practical purposes.

Query Cost. To examine query cost, we look at the number of vertices seen by the algorithm over the first 5000 inquiries and compare it to BiBFS. We direct the reader to Figure 2(b) for a summary of the query cost for our larger graphs, and Figure 5 for the same in the smaller ones. Consistently across all graphs, BiBFS quickly views between 70% and 100% of the vertices in just a few hundred inquiries. In comparison, the query cost of WormHole is quite small for all networks: in the smaller networks, we see less than 30% of the vertices even after 5000 inquiries, and in the larger ones this number is less than 10% (in the largest ones, wikipedia and soc-twitter, it is <2%).

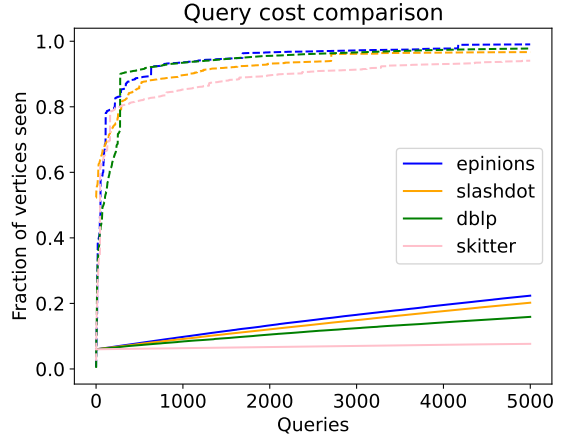


Figure 5: Query cost of WormHole and that of BiBFS different datasets: fraction of the graph seen by WormHole vs BiBFS over the first 10k inquiries in small and medium graphs. The dotted lines refer to the query cost by BiBFS while the solid lines are due to WormHole. The results for large and huge graphs appear in item (b) of Figure 1.

Accuracy. We consider two error measures, absolute (additive) and relative (multiplicative). Clearly, an additive

error of, say, 2, is less preferable for shorter paths than for longer ones. The *relative error* $re(s, t)$ is more refined, as it measures the error with respect to the actual distance of the pair at question. Formally,

$$re(s, t) = \frac{\bar{d}(s, t) - d(s, t)}{d(s, t)}, \quad (2)$$

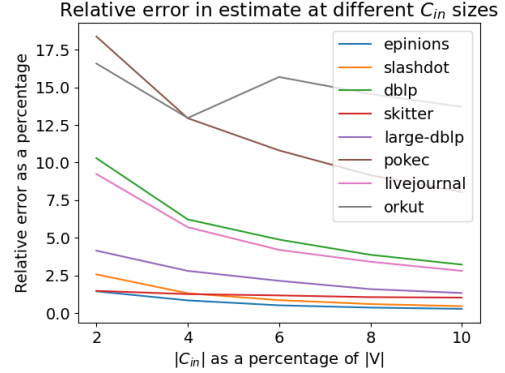
where for a vertex pair (s, t) , $d(s, t)$ is the true distance and $\bar{d}(s, t)$ is the approximate distance estimated by `WormHole`. We investigate the relative error as a function of the core size – see Figure 6. In general, the accuracy drops as we decrease the size of the core. Moreover, we observe that larger graphs give comparably good results at much smaller inner ring sizes, keeping in line with our hypothesis of a sublinear inner ring.

The other key accuracy statistic is additive error: this is summarized in Table 3. For `WormHoleE` across almost all networks, the vast majority of the pair inquiries are *estimated perfectly*. Our worst performance is on soc-pokec and soc-live. Even there, we have perfect estimates for 60% of vertices, and over 94% of vertices have an additive error of less than 1. In *all networks*, more than 99% of the pairs are estimated with absolute error lesser or equal to 2 in `WormHoleE`. The accuracy is poorer in `WormHoleH` (recall that in this variant we do not compute all-pairs shortest paths in the core, resorting to an approximate heuristic instead). However, we note that even then, in most graphs, we have an additive error of at most 2 in over 99% of the queries, and over 90% for all graphs.

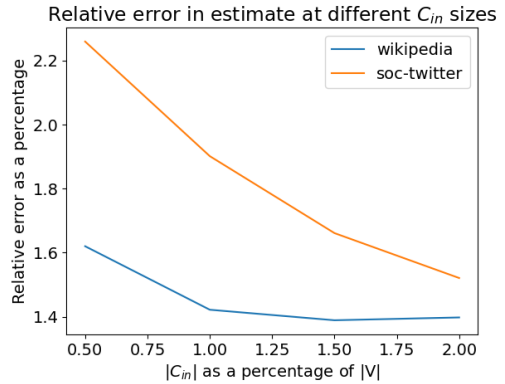
Speedups over BiBFS in inquiry time. The main utility of `WormHoleH` is in exhibiting how much faster our algorithm becomes if we sacrifice some accuracy. This is also documented in Table 3. `WormHoleE` already achieves speedups per inquiry over BiBFS: typically at least 2×, but up to over 4× in some networks. The variant `WormHoleH` further speeds up each inquiry by another order of magnitude, up to a massive 181× in case of our largest network, soc-twitter (while utilizing the same decomposition and data structure). We note that BiBFS does not have any setup cost, while `WormHole` does (for both variants). However, we show that our setup costs are typically very low: our highest setup time is about two minutes, and the highest setup space requirement is under 100MB. The complete statistics are provided in Table 4. (In our implementation, setup time is the time needed for us to capture C_{in} and C_{out} in Algorithm CoreGen, and setup space is the stored binary files for C_{in} and C_{out} .)

5.2 Comparison with index-based methods

As discussed, index construction allows for much faster inquiry times, but setup times that can take up to several hours even for relatively small-sized graphs. We attempt to benchmark our algorithm against two state of the art methods, PLL and MLL. PLL solves the easier task of finding distances, while MLL does explicit shortest path construction (the authors note that PLL may be extended to output paths, but no code is publicly available). However, once the graphs hit a few million vertices, these methods either take too long to run the setup, or even if they succeed in index construction, they may be too large to load into memory. We summarize the results in Table 5, and expand on the discussion in the following paragraphs.



(a) Mean relative error with varying inner ring sizes across different networks



(b) Mean relative error with varying inner ring sizes for soc-twitter

Figure 6: Accuracy of `WormHoleE` in different settings across different datasets. We plot relative error at different inner ring sizes for the datasets.

Mean inquiry time. In the cases where index based methods do succeed (limited to graphs with fewer than 30 million edges), they have a clear advantage in per inquiry cost. Their typical inquiry time is in the microsecond range, where PLL is faster since it only computes distances, and MLL is about 3 times slower than PLL.

Setup cost. In terms of setup times, both `WormHoleE` and `WormHoleH` have a massive advantage over the index-based methods. We direct the reader to Table 4 and Table 5 for a comparison of setup times: we let all methods run for 12 hours, and terminate if they do not finish in that time. Both PLL and MLL failed to complete setup for any graph with more than 30 million vertices. In comparison, even for our largest graph, soc-twitter, of over a 1.5 billion edges, the setup time for `WormHoleE` is just minutes. Even in the cases where these methods do terminate, the storage footprint is massive. We observed that if allowed to run, MLL completes index construction on soc-pokec in a little under 24 hours, but the constructed files are almost a combined 45 gigabytes in size; in comparison, the input COO file is only 250 megabytes! A detailed comparison

Metric	epinions	slashdot	dblp	skitter	large-dblp	pokec	livejournal	orkut	wikipedia	soc-twitter
Setup time (s)	0.02	0.02	0.03	0.42	0.40	0.46	1.68	1.49	11.80	123.57
$C_{in}+C_{out}(MB)$	0.15	0.15	0.62	3.31	3.56	3.19	9.47	6.00	26.56	81.35

Table 4: Setup cost for WormHole: this holds for both WormHole_E and WormHole_H. Setup time is the time needed to capture C_{in} and C_{out} in Algorithm 1. The last row, space, is the footprint on disk of our binary arrays.

Network	Setup (sec)		Inq. time (μs)		Breakeven	
	PLL	MLL	PLL	MLL	PLL	MLL
epinions	4.1	1.5	0.96	2.66	101k	39k
slashdot	6.8	3.6	1.08	3.98	151k	85k
dblp	218	52.4	4.05	11.99	2.1M	535k
skitter	1.1k	466	2.72	11.06	769k	326k
large-dblp	9.2k	1.6k	9.25	N/A	6.4M	inf

Table 5: Comparisons with PLL and MLL. We look at setup time, mean inquiry time, and breakeven compared to WormHole_E. The indexing based methods do not terminate on graphs larger than this. For large-dblp, setup completes for MLL but we are unable to make inquiries.

Network	Setup for MLL on C_{in}		MIT (μs)
	Time (sec)	Space (MB)	
epinions	0.41	2	1.57
slashdot	0.54	5	2.45
dblp	2.99	20	5.00
skitter	28.24	106	12.04
large-dblp	55.68	182	15.22
pokec	144.34	328	45.52
livejournal	803.95	1303	57.89
orkut	1156.28	1476	157.67
wikipedia	551.19	452	120.65
soc-twitter	20949.82	6215	115.44

Table 6: WormHole_M: running MLL on C_{in} .

of the space footprint is given in Figure 2 in §1, and the time comparisons can be found in Figure 7, and in Table 5.

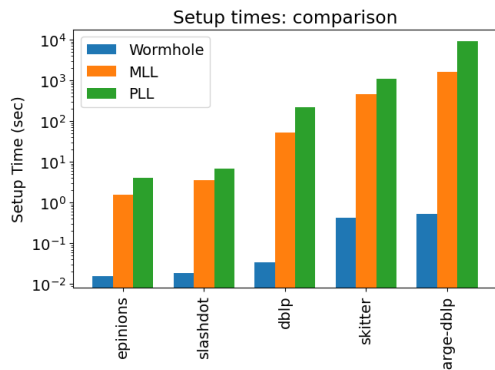


Figure 7: A comparison of setup time between different methods. The index-based methods did not terminate on graphs larger than these.

When is indexing better? A running time comparison: Given the very high setup time of index-based methods, WormHole_E has a head start. However, index-based solutions do catch up and outdo WormHole after sufficiently many shortest path inquiries. We quantify this threshold to give the reader a sense of when to favor each approach. We refer to the threshold as *Breakeven* (BE_X for a competing method X) and provide its values for different graphs in Table 5. Since the index-based methods have setup times between 10^6 and 10^9 orders of magnitude higher than the inquiry times, we look at how many inquiries are needed for them to have an average gain over WormHole_E in the net time taken. Formally, we define breakeven for PLL (likewise MLL) as:

$$BE_{PLL} = \frac{\text{setup}_{PLL} - \text{setup}_{\text{WormHole}_E}}{MIT_{\text{WormHole}_E} - MIT_{PLL}}$$

We can similarly compute the breakeven with respect to WormHole_H, but we note that it is already quite high even against the much slower version WormHole_E. This implies

that even with the low time per inquiry, the setup cost is so prohibitively high that the gains take hundreds of thousands to even millions of inquiries to set in, even on the small networks.

5.3 WormHole as a primitive: WormHole_M

WormHole_E and WormHole_H perform well in terms of query cost and accuracy. The inquiry times, especially in the latter variant, are also huge improvements over BiBFS, but lag behind indexing based methods that perform lookups to find shortest paths. However, as evident by our experiments, landmark based index creation is often prohibitively expensive, both in terms of the time taken to create the index and the space required to store it, to the extent that it may even be impossible for large networks. The success of WormHole_E comes from exact shortest paths computed solely on a small core, which is as low as 1% of the graph. In practice, it takes very little time for a traversal based algorithm to reach C_{in} , and the bulk of the cost comes from the exact path computation inside C_{in} . We thus ask, how much faster can we make our algorithm if we speed this process up, perhaps by doing indexing solely on the core?

To this end, we propose as a third alternative, WormHole_M: it functions almost identical to WormHole_H, but instead of running BiBFS to find the shortest paths in C_{in} , it sets up MLL on all of C_{in} and then uses the MLL index to compute shortest paths inside C_{in} . This is an illustrative example to show how our decomposition can be used in combination with existing techniques. The accuracy guarantees of this variant as presented will be identical to WormHole_H; however, we do not analyze the index size theoretically and suspect it will not be sublinear.

We conduct similar experiments for WormHole_M. Remarkably, while MLL fails to complete setup on most of the graphs, WormHole_M successfully runs it on the core in *all* cases. Moreover, as noted in Table 6, the cost in both time and space is orders of magnitude smaller than for the full graph, though still significantly larger than the default WormHole_H. We note about

two orders of magnitude of improvement in time per inquiry over WormHole_H , but at the same time, the setup cost is also about two orders of magnitude higher in both time and space. Notably, even in cases where MLL does complete on the full graph, we are able to answer inquiries in roughly the same time (see Figure 1) at a fraction of the setup cost (Figure 7). We leave a more systematic investigation of this approach of combining WormHole with existing methods on the core to future work.

REFERENCES

- [1] Ittai Abraham, Daniel Delling, Amos Fiat, Andrew V. Goldberg, and Renato F. Werneck. 2011. VC-Dimension and Shortest Path Algorithms. In *Automata, Languages and Programming (ICALP '11)*.
- [2] Ittai Abraham, Daniel Delling, Amos Fiat, Andrew V. Goldberg, and Renato F. Werneck. 2016. Highway Dimension and Provably Efficient Shortest Path Algorithms. *J. ACM* 63, 5, Article 41 (2016).
- [3] Ittai Abraham, Daniel Delling, Andrew V. Goldberg, and Renato F. Werneck. 2011. A Hub-Based Labeling Algorithm for Shortest Paths in Road Networks. In *Experimental Algorithms*, Panos M. Pardalos and Steffen Rebennack (Eds.).
- [4] D. Aingworth, C. Chekuri, P. Indyk, and R. Motwani. 1999. Fast Estimation of Diameter and Shortest Paths (Without Matrix Multiplication). *SIAM J. Comput.* 28, 4 (1999).
- [5] Takuya Akiba, Yoichi Iwata, and Yuichi Yoshida. 2013. Fast Exact Shortest-Path Distance Queries on Large Networks by Pruned Landmark Labeling (*SIGMOD '13*).
- [6] Takuya Akiba, Christian Sommer, and Ken-ichi Kawarabayashi. 2012. Shortest-Path Queries for Complex Networks: Exploiting Low Tree-Width Outside the Core (*EDBT '12*).
- [7] Réka Albert, Hawoong Jeong, and Albert-László Barabási. 1999. Diameter of the World-Wide Web. *Nature* 401, 6749 (1999).
- [8] Noga Alon, Allan Grönlund, Søren Fuglede Jørgensen, and Kasper Green Larsen. 2023. Sublinear Time Shortest Path in Expander Graphs. arXiv:2307.06113 [cs.DS]
- [9] Igor Artico, I Smolyarenko, Veronica Vinciotti, and Ernst C Wit. 2020. How rare are power-law networks really? *Proceedings of the Royal Society A* 476, 2241 (2020).
- [10] Omri Ben-Eliezer, Talya Eden, Joel Oren, and Dimitris Fotakis. 2022. Sampling Multiple Nodes in Large Networks: Beyond Random Walks (*WSDM '22*). 37–47.
- [11] Thomas Blásius, Cedric Freiberger, Tobias Friedrich, Maximilian Katzmann, Felix Montenegro-Retana, and Marianne Thieffry. 2022. Efficient Shortest Paths in Scale-Free Networks with Underlying Hyperbolic Geometry. *ACM Trans. Algorithms* 18, 2, Article 19 (mar 2022).
- [12] Michele Borassi and Emanuele Natale. 2019. KADABRA is an Adaptive Algorithm for Betweenness via Random Approximation. *ACM J. Exp. Algorithmics* 24, Article 1.2 (feb 2019).
- [13] Arthur Brady and Lenore Cowen. 2006. Compact Routing on Power Law Graphs with Additive Stretch. In *Proceedings of the Workshop on Algorithm Engineering and Experiments (ALENEX '06)*.
- [14] Flavio Chiericetti, Anirban Dasgupta, Ravi Kumar, Silvio Lattanzi, and Tamás Sarlós. 2016. On Sampling Nodes in a Network. In *Proceedings of the 25th International Conference on World Wide Web (WWW '16)*.
- [15] Fan Chung and Linyuan Lu. 2002. Connected components in random graphs with given expected degree sequences. *Annals of combinatorics* 6, 2 (2002).
- [16] Fan Chung and Linyuan Lu. 2004. The average distance in a random graph with given expected degrees. *Internet Mathematics* 1, 1 (2004).
- [17] Fan Chung and Linyuan Lu. 2006. The volume of the giant component of a random graph with given expected degrees. *SIAM Journal on Discrete Mathematics* 20, 2 (2006).
- [18] Andrea Costa, Ana M. Martín González, Katell Guizien, Andrea M. Doglioli, José María Gómez, Anne A. Petrenko, and Stefano Allesina. 2019. Ecological networks: Pursuing the shortest path, however narrow and crooked. *Scientific Reports* 9, 1 (Nov. 2019).
- [19] Mingyang Deng, Yael Kirkpatrick, Victor Rong, Virginia Vassilevska Williams, and Qiyan Zhong. 2022. New Additive Approximations for Shortest Paths and Cycles (*ICALP '22, Vol. 229*).
- [20] Jörg Derungs, Riko Jacob, and Peter Widmayer. 2007. Approximate Shortest Paths Guided by a Small Index. In *Algorithms and Data Structures, 10th International Workshop, WADS (Lecture Notes in Computer Science, Vol. 4619)*.
- [21] E. W. Dijkstra. 1959. A Note on Two Problems in Connexion with Graphs. *Numer. Math.* 1, 1 (dec 1959).
- [22] Michal Dory, Sebastian Forster, Yael Kirkpatrick, Yasamin Nazari, Virginia Vassilevska Williams, and Tijn de Vos. 2023. Fast 2-Approximate All-Pairs Shortest Paths. *CoRR abs/2307.09258* (2023).
- [23] Michalis Faloutsos, Petros Faloutsos, and Christos Faloutsos. 1999. On Power-Law Relationships of the Internet Topology (*SIGCOMM '99*).
- [24] John A Fitch III and Lance J Hoffman. 1993. A shortest path network security model. *Computers & Security* 12, 2 (1993), 169–189.
- [25] Michael L. Fredman and Robert Endre Tarjan. 1987. Fibonacci Heaps and Their Uses in Improved Network Optimization Algorithms. *J. ACM* 34, 3 (jul 1987).
- [26] Peng Gao, Xusheng Xiao, Zhichun Li, Kangkook Jee, Fengyuan Xu, Sanjeev R. Kulkarni, and Prateek Mittal. 2019. A Query System for Efficiently Investigating Complex Attack Behaviors for Enterprise Security. *Proc. VLDB Endow.* 12, 12 (2019).
- [27] Patrice Godefroid, Nils Klarlund, and Koushik Sen. 2005. DART: Directed Automated Random Testing. In *Proceedings of the 2005 ACM SIGPLAN conference on Programming language design and implementation (PLD '05)*.
- [28] Andrew V. Goldberg and Chris Harrelson. 2005. Computing the Shortest Path: A Search Meets Graph Theory (*SODA '05*).
- [29] Ruoming Jin, Zhen Peng, Wendell Wu, Feodor Dragan, Gagan Agrawal, and Bin Ren. 2020. Parallelizing Pruned Landmark Labeling: Dealing with Dependencies in Graph Algorithms (*ICS '20*). Article 11.
- [30] Xin Jin, Charalampos Katsis, Fan Sang, Jiahao Sun, Elisa Bertino, Ramana Rao Kompella, and Ashish Kundu. 2023. Prometheus: Infrastructure Security Posture Analysis with AI-generated Attack Graphs. *CoRR* (2023). arXiv:2312.13119
- [31] Siddharth Kaza, Jennifer Xu, Byron Marshall, and Hsinchun Chen. 2009. Topological Analysis of Criminal Activity Networks: Enhancing Transportation Security. *IEEE Transactions on Intelligent Transportation Systems* 10, 1 (2009).
- [32] Masahiro Kimura and Kazumi Saito. 2006. Tractable Models for Information Diffusion in Social Networks. In *Knowledge Discovery in Databases: PKDD 2006*.
- [33] Maksim Kitsak, Alexander Ganin, Ahmed Elmokashfi, Hongzhu Cui, Daniel A. Eisenberg, David L. Alderson, Dmitry Korkin, and Igor Linkov. 2023. Finding shortest and nearly shortest path nodes in large substantially incomplete networks by hyperbolic mapping. *Nature Communications* 14, 1 (2023).
- [34] Jérôme Kunegis. 2013. KONECT: The Koblenz Network Collection (*WWW '13 Companion*).
- [35] Andrey Kutuzov, Mohammad Dorgham, Oleksiy Oliynyk, Chris Biemann, and Alexander Panchenko. 2019. Learning Graph Embeddings from WordNet-based Similarity Measures. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*.
- [36] Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>.
- [37] Wentao Li, Miao Qiao, Lu Qin, Ying Zhang, Lijun Chang, and Xuemin Lin. 2020. Scaling up distance labeling on graphs with core-periphery properties (*SIGMOD '20*). 1367–1381.
- [38] Wentao Li, Miao Qiao, Lu Qin, Ying Zhang, Lijun Chang, and Xuemin Lin. 2020. Scaling Up Distance Labeling on Graphs with Core-Periphery Properties (*SIGMOD '20*).
- [39] Wentao Li, Miao Qiao, Lu Qin, Ying Zhang, Lijun Chang, and Xuemin Lin. 2022. Distance labeling: on parallelism, compression, and ordering. *The VLDB Journal* 31 (01 2022).
- [40] Ye Li, Leong Hou U, Man Lung Yiu, and Ngai Meng Kou. 2017. An Experimental Study on Hub Labeling Based Shortest Path Algorithms. *Proc. VLDB Endow.* 11, 4 (dec 2017).
- [41] Linyuan Lincoln Lu. 2002. *Probabilistic methods in massive graphs and Internet computing*. Ph.D. Dissertation. University of California, San Diego.
- [42] Amgad Madkour, Walid G Aref, Faizan Ur Rehman, Mohamed Abdur Rahman, and Saleh Basalamah. 2017. A survey of shortest-path algorithms. *arXiv preprint arXiv:1705.02044* (2017).
- [43] Fragkiskos D. Malliaros, Christos Giatsidis, Apostolos N. Papadopoulos, and Michalis Vazirgiannis. 2019. The Core Decomposition of Networks: Theory, Algorithms and Applications. *The VLDB Journal* 29, 1 (nov 2019).
- [44] M. E. J. Newman. 2003. The Structure and Function of Complex Networks. *SIAM Rev.* 45, 2 (2003).
- [45] Ali Pinar, C. Seshadhri, and Vaidyanathan Vishal. 2017. ESCAPE: Efficiently Counting All 5-Vertex Subgraphs (*WWW '17*).
- [46] Ira Pohl. 1969. *Bi-directional and heuristic search in path problems*. Ph.D. Dissertation. Stanford University, USA.
- [47] Jianzhong Qi, Wei Wang, Rui Zhang, and Zhuowei Zhao. 2020. A Learning Based Approach to Predict Shortest-Path Distances (*EDBT'20*).
- [48] Fatemeh Salehi Rizi, Joerg Schloetterer, and Michael Granitzer. 2020. Shortest Path Distance Approximation Using Deep Learning Techniques (*ASONAM '18*).
- [49] Liam Roditty. 2023. New Algorithms for All Pairs Approximate Shortest Paths (*STOC '23*).
- [50] M. Puck Rombach, Mason A. Porter, James H. Fowler, and Peter J. Mucha. 2014. Core-Periphery Structure in Networks. *SIAM J. Appl. Math.* 74, 1 (2014).

- [51] Barna Saha and Christopher Ye. 2023. Faster Approximate All Pairs Shortest Paths. *CoRR* abs/2309.13225 (2023).
- [52] Stephen B. Seidman. 1983. Network structure and minimum degree. *Social Networks* 5, 3 (1983).
- [53] Sandeep Sen. 2009. Approximating shortest paths in graphs. In *International Workshop on Algorithms and Computation*. 32–43.
- [54] Matteo Serafino, Giulio Cimini, Amos Maritan, Andrea Rinaldo, Samir Suweis, Jayanth R. Banavar, and Guido Caldarelli. 2021. True scale-free networks hidden by finite size effects. *Proceedings of the National Academy of Sciences* 118, 2 (2021).
- [55] Christian Sommer. 2014. Shortest-path queries in static networks. *ACM Computing Surveys (CSUR)* 46, 4 (2014).
- [56] Chi Wang, Wei Chen, and Yajun Wang. 2012. Scalable Influence Maximization for Independent Cascade model in Large-scale Social Networks. *Data Mining and Knowledge Discovery Journal* 25, 3 (2012).
- [57] Ye Wang, Qing Wang, Henning Koehler, and Yu Lin. 2021. Query-by-Sketch: Scaling Shortest Path Graph Queries on Very Large Networks (*SIGMOD '21*).
- [58] Jennifer J Xu and Hsinchun Chen. 2004. Fighting organized crimes: using shortest-path algorithms to identify associations in criminal networks. *Decision Support Systems* 38, 3 (2004).
- [59] Xiaolong Xu, Yuancheng Li, Tao Huang, Yuan Xue, Kai Peng, Lianyong Qi, and Wanchun Dou. 2019. An energy-aware computation offloading method for smart edge computing in wireless metropolitan area networks. *Journal of Network and Computer Applications* (2019).
- [60] Zhiqiang Xu, Pengcheng Fang, Changlin Liu, Xusheng Xiao, Yu Wen, and Dan Meng. 2022. DEPCOMM: Graph Summarization on System Audit Logs for Attack Investigation. In *IEEE Symposium on Security and Privacy (SP)*.
- [61] Hongyang Zhang, Huacheng Yu, and Ashish Goel. 2019. Pruning Based Distance Sketches with Provable Guarantees on Random Graphs (*WWW '19*).
- [62] Junhua Zhang, Wentao Li, Long Yuan, Lu Qin, Ying Zhang, and Lijun Chang. 2022. Shortest-Path Queries on Complex Networks: Experiments, Analyses, and Improvement. *Proc. VLDB Endow.* 15, 11 (2022).
- [63] Xiao Zhang, Travis Martin, and M. E. J. Newman. 2015. Identification of core-periphery structure in networks. *Physical Review E* 91, 3 (mar 2015).
- [64] Xiaohan Zhao, Alessandra Sala, Christo Wilson, Haitao Zheng, and Ben Y. Zhao. 2010. Orion: Shortest Path Estimation for Large Social Graphs (*WOSN'10*).
- [65] Xiaohan Zhao, Alessandra Sala, Haitao Zheng, and Ben Y. Zhao. 2011. Efficient shortest paths on massive social graphs. In *7th International Conference on Collaborative Computing: Networking, Applications and Worksharing (CollaborateCom)*.
- [66] Uri Zwick. 2001. Exact and Approximate Distances in Graphs – A Survey. In *Algorithms – ESA 2001*.
- [67] Uri Zwick. 2002. All pairs shortest paths using bridging sets and rectangular matrix multiplication. *Journal of the ACM (JACM)* 49, 3 (2002), 289–317.