



HAL
open science

Enumeration of far-apart pairs by decreasing distance for faster hyperbolicity computation

David Coudert, André Nusser, Laurent Viennot

► **To cite this version:**

David Coudert, André Nusser, Laurent Viennot. Enumeration of far-apart pairs by decreasing distance for faster hyperbolicity computation. [Research Report] Inria; I3S, Université Côte d'Azur. 2021. hal-03201405

HAL Id: hal-03201405

<https://inria.hal.science/hal-03201405v1>

Submitted on 18 Apr 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Enumeration of far-apart pairs by decreasing distance for faster hyperbolicity computation*

David Coudert¹, André Nusser², and Laurent Viennot³

¹Université Côte d’Azur, Inria, CNRS, I3S, France

²Max Planck Institute for Informatics and Graduate School of Computer Science, Saarland Informatics Campus, Saarbrücken, Germany

³Inria, Paris University, CNRS, Irif, France

Abstract

Hyperbolicity is a graph parameter which indicates how much the shortest-path distance metric of a graph deviates from a tree metric. It is used in various fields such as networking, security, and bioinformatics for the classification of complex networks, the design of routing schemes, and the analysis of graph algorithms. Despite recent progress, computing the hyperbolicity of a graph remains challenging. Indeed, the best known algorithm has time complexity $O(n^{3.69})$, which is prohibitive for large graphs, and the most efficient algorithms in practice have space complexity $O(n^2)$. Thus, time as well as space are bottlenecks for computing the hyperbolicity.

In this paper, we design a tool for enumerating all far-apart pairs of a graph by decreasing distances. A node pair (u, v) of a graph is far-apart if both v is a leaf of all shortest-path trees rooted at u and u is a leaf of all shortest-path trees rooted at v . This notion was previously used to drastically reduce the computation time for hyperbolicity in practice. However, it required the computation of the distance matrix to sort all pairs of nodes by decreasing distance, which requires an infeasible amount of memory already for medium-sized graphs. We present a new data structure that avoids this memory bottleneck in practice and for the first time enables computing the hyperbolicity of several large graphs that were far out-of-reach using previous algorithms. For some instances, we reduce the memory consumption by at least two orders of magnitude. Furthermore, we show that for many graphs, only a very small fraction of far-apart pairs have to be considered for the hyperbolicity computation, explaining this drastic reduction of memory.

As iterating over far-apart pairs in decreasing order without storing them explicitly is a very general tool, we believe that our approach might also be relevant to other problems.

Keywords: Gromov hyperbolicity; graph algorithms, far-apart pairs iterator.

*This work has been supported by the French government, through the UCA^{JEDI} Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-15-IDEX-01, and the Distancia project with reference number ANR-17-CE40-0015.

1 Introduction

This paper aims at computing the hyperbolicity of graphs whose size ranges from tens of thousands to millions of nodes. The hyperbolicity is a parameter of a metric space generalizing the idea of Riemannian manifolds with negative curvature. When considering the metric of a graph, it measures, to some extent, how much it deviates from a tree metric. This parameter was first introduced by Gromov in the context of automatic groups [43] in relation with their Cayley graphs.

Hyperbolicity has received great attention in computer science in the last decades as it seems to capture important properties of several large practical graphs such as Internet [58], the Web [52] and databases relations [65]. It is also used to classify complex networks [1, 5, 48] and was proposed as a measure of how much a network is “democratic” [4, 11]. Formal relationships between Gromov hyperbolicity and the existence of a core (a subset of vertices intersecting a constant fraction of all the shortest-paths) are investigated in [20]. Reciprocally, the existence of a core is shown to be inherent to any hyperbolic network in [20]. Furthermore, small hyperbolicity has tractability implications and measuring hyperbolicity has applications in routing [9, 19, 49], approximating other graph parameters [18, 32] and bioinformatics [16, 39]. See [1, 37] for recent surveys.

Computing the hyperbolicity is often a prerequisite in the above applications. As hyperbolicity can be defined by a simple 4-point condition, it can be naively computed in $\mathcal{O}(n^4)$ time. As far as we know, the best theoretical algorithm [40] has time complexity $\mathcal{O}(n^{3.69})$. Although its complexity is $o(n^4)$, it is still supercubic and the algorithm appears to be impractical for graphs with a few tens of thousands of nodes. On the lower bounds side it was shown that under the Strong Exponential Time Hypothesis [47] hyperbolicity cannot be computed in subquadratic-time, even for sparse graphs [13, 24, 40].

The only practical algorithms that can manage larger graphs [12, 23] enumerate all pairs of nodes by decreasing distance. For each pair, each 4-tuple obtained with a previous pair is tested with regard to the 4-point condition defining hyperbolicity. Each test of a 4-tuple provides a lower bound of hyperbolicity. This approach allows to stop the enumeration as soon as the distance of the scanned pair equals twice the best hyperbolicity lower bound found so far. As it scans a portion of all 4-tuples, its worst case complexity is $\mathcal{O}(n^4)$ but it appears much faster in practice as first scanning pairs with large distances allows to find good lower bounds early in the enumeration. A main optimization for further reducing the number of pairs scanned, consists in considering only far-apart pairs, that is, pairs such that no neighbor of one node is further apart from the other node, see Section 2 for a formal definition. It can be proven that the 4-point condition defining hyperbolicity holds on all 4-tuples if it holds on 4-tuples made up of two such far-apart pairs [53, 60]. The main bottleneck of this method lies in its inherent quadratic space usage: all far-apart pairs and all-pair distances are stored in the current implementations of this approach. Computing the hyperbolicity of practical graphs with millions of nodes thus remains a great challenge.

1.1 Our approach

To make progress towards this challenge, we propose to enumerate far-apart pairs by decreasing distance without computing all-pair distances. The key of our approach is to first compute all eccentricities, that is, for each node, what is the largest distance from it. Computing all eccentricities is feasible in practice, see Section 1.3. Note that we obtain the diameter D of the graph as a side product, as it simply is the maximum eccentricity. We then scan nodes with eccentricity D and enumerate far-apart nodes at distance D from them, that is, nodes at distance D that form

far-apart pairs with them. We then scan nodes with eccentricity at least $D - 1$ and enumerate far-apart nodes at distance $D - 1$ from them, and so on. We also include various optimizations proposed in [23, 12] to further reduce the number of 4-tuples considered. The main difficulty of our approach is that we have to compute distances on the fly as we do not store all-pair distances. This mainly requires to perform a breadth-first search (BFS) for each node of a far-apart pair considered. Interestingly, we show how to prune these BFS searches based on some of the optimizations proposed in [12]. Storing most recent BFS searches in a cache also increases performance for some instances as we observe some sharing of far-apart nodes in practice.

1.2 Main contributions

Our main contributions are the following.

- We present the first non-naive algorithm for iterating over far-apart pairs that neither computes and stores all distances explicitly, nor sorts the node pairs by recomputing all distances from scratch whenever they are needed.
- This, for the first time, enables enumerating all far-apart pairs of large graphs. Previously this was not possible either due to excessive amounts of time or memory needed for the computation of all far-apart pairs.
- As the prime application of our algorithm for iterating over far-apart pairs, we significantly reduce the memory consumption when computing the graph hyperbolicity. The memory reduction is at least two orders of magnitude for some instances.
- This drastic memory reduction enables us to compute the hyperbolicity of many large graphs for the first time.
- Due to the significance of far-apart pairs in a graph (e.g., they are the defining vertices for radius, eccentricities, diameter, ...), we believe that our contribution of a far-apart pair iterator is also relevant in other settings.

1.3 Other related work

The most advanced practical algorithm for computing hyperbolicity [12] can be seen as a refinement of [23] that further prunes the search space. A complementary approach proposed in [22] consists in splitting the graph further than biconnected components using clique decomposition. Using such a decomposition could also be used in our framework to further reduce memory usage.

Practical algorithms for computing all eccentricities [61, 38, 51] came along in a line of research for improving diameter computation of real world graphs [2, 29, 14, 38].

Several optimizations were proposed for BFS search. Most notably, [45] reduces space usage to $O(n)$ bits, and [3] performs several BFS in parallel from a node and some of its neighbors using bit-parallel word operations. Both methods could be used to further improve our approach.

1.4 Organization

Definitions and notations used in this paper are introduced in Section 2. We then present our far-apart pair iterator in Section 3. Section 4 is devoted to hyperbolicity. We recall its definition

and review the best known algorithmic results on this parameter. We then present our memory efficient algorithm based on the proposed far-apart pair iterator. In Section 5, we report on the experimental evaluation of the algorithms presented in this paper on various graphs. In particular, we compare our algorithm for computing hyperbolicity with the state-of-the-art algorithm [12]. We conclude this paper in Section 6 with some directions for future research.

2 Definitions and notations

We use the graph terminology of [10, 36]. All graphs considered in this paper are finite, undirected, connected, unweighted and simple. The graph $G = (V, E)$ has $n = |V|$ vertices and $m = |E|$ edges. The open neighborhood $N_G(S)$ of a set $S \subseteq V$ consists of all vertices in $V \setminus S$ with at least one neighbor in S .

Given two vertices u and v , a uv -path of length $\ell \geq 0$ is a sequence of vertices $(u = v_0 v_1 \dots v_\ell = v)$, such that $\{v_i, v_{i+1}\}$ is an edge for every i . In particular, a graph G is *connected* if there exists a uv -path for all pairs $u, v \in V$, and in such a case the *distance* $d_G(u, v)$ is defined as the minimum length of a uv -path in G . When G is clear from the context, we write d (resp. N) instead of d_G (resp. N_G). The *eccentricity* $\text{ecc}(u)$ of a vertex u is the maximum distance between u and any other vertex $v \in V$, i.e., $\text{ecc}(u) = \max_{v \in V} d(u, v)$. The maximum eccentricity is the *diameter* $\text{diam}(G)$ and the minimum eccentricity is the *radius* $\text{rad}(G)$.

The notion of *far-apart* pairs of vertices has been introduced in [60, 53] to reduce the number of 4-tuples to consider in the computation of the hyperbolicity (see Section 4). Roughly, we say that two vertices $u, v \in V$ are *far-apart* if for all $w \in V$ neither u lies on a shortest path from w to v , nor v lies on a shortest path from u to w . More formally, we have:

Definition 1. *In a graph $G = (V, E)$, vertex u is far from vertex v , or v -far, if for any neighbor w of u , we have $d(v, w) \leq d(v, u)$. The pair u, v of vertices is far-apart if u is v -far and v is u -far.*

The number of far-apart pairs in a graph can be orders of magnitude smaller than the total number of pairs. For instance, a $p \times q$ grid has only 2 far-apart pairs. On the other hand, all pairs in a clique graph are far-apart.

The set of all far-apart pairs can be determined in time $\mathcal{O}(nm)$ in unweighted graphs through breadth-first search (BFS), and the interested reader will find in Appendix A a discussion on several time and space complexity trade-offs for determining far-apart pairs. We now present some interesting properties of far vertices.

Lemma 1. *For $v \in V$, vertex $u \in V$ is v -far if and only if u is a leaf of all shortest path trees rooted at v .*

Proof. Clearly, a non-leaf vertex u of a shortest path tree rooted at v has a neighbor at distance $d(v, u) + 1$ and cannot be v -far. Reciprocally, if u is not v -far, it has a neighbor w satisfying $d(v, w) = d(v, u) + 1$. Modifying any shortest path tree by setting u as the parent of w yields a valid shortest path tree where u is not a leaf. \square

From Lemma 1, we also get the following immediate results.

Corollary 1. *For each $u \in V$, any $v \in V$ such that $d(u, v) = \text{ecc}(u)$ is u -far.*

Corollary 2. *For any $u, v \in V$, if $d(u, v) = \text{diam}(G)$, then the pair (u, v) is far-apart.*

In particular, Corollary 1 implies that for any vertex u , the set of u -far vertices is non-empty.

Interestingly, knowing the far vertices of a node u and their distance to u , it is possible to scan distant nodes in a sort of backward BFS as described in Appendix B.

3 Iterator over far-apart pairs

In Appendix A we present several algorithmic choices that result in different time and space complexity trade-offs for determining the set of far-apart pairs. Here, we engineer a data structure and algorithms to determine the set of far-apart pairs and return these pairs sorted by decreasing distances. Our objective is to provide an iterator that determines the next pair to yield on the fly, that postpones computations as much as possible, and with an acceptable memory consumption.

Data structure. To store and organize data, we use an array F indexed by distances in range from 1 to $\text{diam}(G)$, so that cell F^d contains data related to far-apart pairs at distance d . More precisely, F^d is a hash map associating to a vertex $u \in V$ the subset F_u^d of u -far vertices at distance d from u . The subset F_u^d can be implemented using either a set data structure allowing to answer a membership query in $\mathcal{O}(1)$ time, or an array of size $|F_u^d|$ whose elements are sorted according to any total ordering of the vertices to enable membership queries in time $\mathcal{O}(\log_2 |F_u^d|)$.

We assume, as it is the case for most modern programming languages, that it is possible to visit the elements of a hash map in a fixed arbitrary order (e.g., insertion order). We have the same assumption for set data structures.

Initialization. Recall that our aim is to iterate over far-apart pairs by non increasing distances and to postpone computation as much as possible. To this end, we initially store in F , for each vertex u with eccentricity $\text{ecc}(u)$, an empty set of u -far vertices in F^d with $d = \text{ecc}(u)$. This empty set will serve as an indicator to trigger the effective computation of the set of u -far vertices the first time it is needed (recall that by Corollary 1, this set is non-empty).

Clearly, this initialization procedure requires the knowledge of the eccentricities of all vertices. Fortunately, although the determination of the eccentricities has worst-case time complexity in $\mathcal{O}(nm)$, practically efficient algorithms have been proposed to perform this task [38, 61]. These algorithms maintain upper and lower bounds on the eccentricity of each vertex and improve these bounds by computing distances from a few well chosen vertices until all gaps are closed. In practice, these algorithms are orders of magnitude faster than a naive algorithm performing a BFS from each vertex.

Filling. Each time we compute distances from a vertex u that has not been considered before, the corresponding sets F_u^d for all d can be inserted in F . This is done in particular while computing the eccentricities during the initialization of the data structure.

Next. We define a function $\text{next}(F)$ that yields the next far-apart pair in the ordering. Observe that, since F is used to iterate over the far-apart pairs by non increasing distances, when starting to consider far-apart pairs at distance d , we know that all pairs at distance $d' > d$ have already been considered and that at initialization, we have added in F^d an empty set for each vertex with eccentricity d . Hence, all vertices involved in a far-apart pair at distance d are known.

Therefore, we can define a function $\text{next}(F^d)$ that returns either the next item (u, F_u^d) in the fixed arbitrary ordering on the items stored in the hash map F^d , or **Stop** when all items have been considered. Similarly, we define a function $\text{next}(F_u^d)$ that returns either the next vertex in the ordering defined over the vertices in F_u^d , or **Stop** when all vertices have been considered.

For a distance d and a vertex u , we repeat calls to $\text{next}(F_u^d)$ until finding a vertex w such that u is w -far (by testing if u is in F_w^d), in which case the far-apart pair (u, w) at distance d is returned. If no vertex w is found, we go to the next vertex in F^d through a call to $\text{next}(F^d)$. When all vertices in F^d have been considered, we start considering far-apart pairs at distance $d - 1$.

The use of a total ordering on the vertices allows to ensure that a far-apart pair (u, v) is returned only once, i.e., when $u < v$. During these operations, if we encounter an empty set F_x^d , we know from the initialization step that x has never been considered before and that $\text{ecc}(x) = d$. So, we compute distances from x and store the x -far vertices in F via the filling step.

Improvements. Depending on the usage of this data structure, some simple improvements can be done, such as:

1. When iterating over the vertices in F_u^d , each time a vertex $w \in F_u^d$ is found such that u is not w -far, we can remove w from F_u^d . This avoids checking twice if a pair is far-apart, and, at the end of the iterations on F_u^d , this set will contain a vertex w only if u is w -far. Actually, instead of removing elements from F_u^d , it is safer to insert the vertices w such that u is w -far into a temporary array (or set) T , and to replace F_u^d by T when $\text{next}(F_u^d)$ returns **Stop**. We will use this improvement in our algorithm for computing hyperbolicity.
2. When all vertices in F^d have been considered, one can delete F^d to reduce the memory consumption if appropriate with the usage. Similarly, one can avoid storing F^d if only far-apart pairs at distance $d' > d$ are requested, as is the case for instance to enumerate the diameters of a graph or when computing hyperbolicity when a lower-bound has been found.

The time complexity for iterating over all far-apart pairs sorted by non increasing distance is in $\mathcal{O}(nm)$ when using sets for F_u^d (resp., $\mathcal{O}(nm + n^2 \log n)$ when using arrays). Indeed, the algorithm computes BFS distances from each vertex of the graph (some during initialization, and others during queries). Furthermore, the query time to report all far-apart pairs involving a vertex u is in $\mathcal{O}(n)$ since $|\cup_{d=1}^{\text{ecc}(u)} F_u^d| = \mathcal{O}(n)$ and checking if $u \in F_w^d$ requires time $\mathcal{O}(1)$ (resp., $\mathcal{O}(\log n)$). The sorting operation over the far-apart pairs is implicit and thus adds no extra cost. The space complexity of the data structure is in $\mathcal{O}(n^2)$. However, we will see in Section 5 that it is much smaller in practice.

4 Use Case: Gromov Hyperbolicity

We now present an interesting use case for our far-apart pair iterator with the computation of the Gromov hyperbolicity of a graph [43].

4.1 Definitions

A metric space (V, d_G) is a tree metric if there exists a distance-preserving mapping from V to the nodes of an edge-weighted tree. If so, the graph G is said to be 0-hyperbolic because it satisfies the 4-points condition below with parameter $\delta = 0$. Introducing some slack δ allows to measure how much the metric of a graph deviates from that of a tree. This slack δ is called the *hyperbolicity* of the graph:

Definition 2 (4-points Condition, [43]). *Let G be a connected graph. For every 4-tuple u, v, x, y of vertices of G , we define $\delta(u, v, x, y)$ as half of the difference between the two largest sums among $S_1 = d(u, v) + d(x, y)$, $S_2 = d(u, x) + d(v, y)$, and $S_3 = d(u, y) + d(v, x)$.*

The hyperbolicity of G , denoted by $\delta(G)$, is equal to $\max_{u,v,x,y \in V(G)} \delta(u, v, x, y)$. Moreover, we say that G is δ -hyperbolic whenever $\delta \geq \delta(G)$.

Other characterizations exist for 0-hyperbolic graphs and yield other definitions for the hyperbolicity δ of a graph that differ only by a small constant factor [7, 33, 43].

From the 4-points condition above, it is straightforward to compute graph hyperbolicity in $\Theta(n^4)$ -time. In theory, it can be decreased to $\mathcal{O}(n^{3.69})$ by using a clever (max, min)-matrix product [40]; however, in practice, the best-known algorithms still run in $\mathcal{O}(n^4)$ -time [12, 23]. Graphs with small hyperbolicity can be recognized faster. In particular, 0-hyperbolic graphs coincide with *block graphs*, that are graphs whose biconnected components are complete subgraphs [6, 46]. Hence, deciding whether a graph is 0-hyperbolic can be done in time $\mathcal{O}(n + m)$. The latter characterization of 0-hyperbolic graphs follows from a more general result stating that the hyperbolicity of a graph is the maximum hyperbolicity of its biconnected components (see e.g. [23] for a proof). More recently, it has been proven that the recognition of $\frac{1}{2}$ -hyperbolic graphs is computationally equivalent to deciding whether there exists a chordless cycle of length 4 in a graph [24]. The latter problem can be solved in deterministic $\mathcal{O}(n^{3.26})$ -time [50] and in randomized $\mathcal{O}(n^{2.373})$ -time [64] by using fast matrix multiplication.

Several pre-processing methods for reducing the size of the input graph have been proposed. In particular, [60] proved that the hyperbolicity of G is equal to the maximum of the hyperbolicity of the graphs resulting from both a *modular* [41, 44] or a *split* [31, 30] decomposition of G . These decompositions can be computed in linear time [17]. Algorithms with time complexity in $\mathcal{O}(\text{mw}(G)^3 \cdot n + m)$ and $\mathcal{O}(\text{sw}(G)^3 \cdot n + m)$ when parameterized respectively by the *modular width* $\text{mw}(G)$ and the *split width* $\text{sw}(G)$ have been proposed in [26]. Moreover, [22] shows how to use modified versions of the atoms of a decomposition of G by clique-minimal separators [62, 8, 25]. This decomposition can be obtained in time $\mathcal{O}(nm)$.

4.2 Previous algorithm

In this section, we recall the algorithm proposed in [12] for computing hyperbolicity. This algorithm improves upon the algorithm proposed in [23] by adding pruning techniques to further reduce the number of 4-tuples to consider. To the best of our knowledge, the algorithms proposed in [12, 23] are the only algorithms enabling to compute the exact hyperbolicity of graphs with up to 50000 nodes. These algorithms are based on the following lemmas.

Lemma 2 ([60, 23]). *Let G be a connected graph. There exist two far-apart pairs (u, v) and (x, y) satisfying $\delta(u, v, x, y) = \delta(G)$.*

Lemma 3 ([23]). *For every 4-tuple u, v, x, y of vertices of a connected graph G , we have $\delta(u, v, x, y) \leq \min_{a, b \in \{u, v, x, y\}} d(a, b)$. Furthermore, if $S_1 = d(u, v) + d(x, y)$ is the largest of the sums defined in Definition 2 (which can be assumed w.l.o.g.), we have $\delta(u, v, x, y) \leq \frac{1}{2} \min \{d(u, v), d(x, y)\}$.*

For the sake of completeness, we quickly give the proof of Lemma 3. Without loss of generality the second largest sum is $S_2 = d(u, x) + d(v, y)$. By triangle inequality, we have $d(u, v) \leq d(u, x) + d(x, y) + d(y, v) = S_2 + d(x, y)$, yielding $S_1 - S_2 \leq 2d(x, y)$. We can similarly obtain $S_1 - S_2 \leq 2d(u, v)$.

The key idea of the algorithms of [12, 23] is to visit the most promising 4-tuples first, that is, those made of pairs of far-apart vertices at largest distance, and to stop computation as soon as the bounds of Lemma 3 are reached. These algorithms thus need to iterate over far-apart pairs ordered by decreasing distances.

More precisely, Algorithm 1 below (see also [12]) iterates first over the far-apart pairs sorted by non increasing distances. Then, given the i -th far-apart pair (x_i, y_i) , it iterates over the previous far-apart pairs (v, w) , such that $d(v, w) \geq d(x_i, y_i)$ in order to consider quadruples (v, w, x_i, y_i) such that $S_1 = d(v, w) + d(x_i, y_i)$ is the largest sum. Note that such pairs (v, w) have been considered previously and satisfy $(v, w) = (x_j, y_j)$ for some $j < i$. This ensures by Lemma 3 that as soon as $d(x_i, y_i) \leq 2\delta_L$, where δ_L is the current best solution, no further improvements can be done. So, the hyperbolicity δ_L of the graph is then returned in Line 5. To iterate over the pairs (v, w) such that $d(v, w) \geq d(x_i, y_i)$, the algorithm maintains the *mates* of each vertex. Vertex w is a mate of v if (v, w) is a far-apart pair satisfying $d(v, w) \geq d(x_i, y_i)$. In other words, (v, w) is a far-apart pair previously considered for some $j < i$ such that $(x_j, y_j) = (v, w)$ or $(x_j, y_j) = (w, v)$.

To further prune the search space, [12] introduces the notions of *skippable*, *acceptable* and *valuable* vertices that are computed by `computeAccVal` according to the definitions given below. Algorithm 1 can be read before considering the details of this optimization. Its time complexity is in $\mathcal{O}(n^4)$ and its space complexity is in $\Theta(n^2)$. Indeed, it not only needs to store the list of far-apart pairs, but also the distance matrix, and the lists of mates.

Skippable, acceptable and valuable. Given a pair x, y of nodes and a lower bound δ_L on hyperbolicity, [12] proposes a classification of the nodes to prune as those that cannot lead to any improvement of the lower-bound δ_L known so far. For instance, a node v such that $\min\{d(x, v), d(y, v)\} \leq \delta_L$ can be skipped, since by Lemma 3 we then have $\delta(x, y, v, w) \leq \delta_L$ for any w . In Lemma 4 we summarize the conditions defining (x, y, δ_L) -skippable nodes.

Lemma 4 ([12]). *A node v is (x, y, δ_L) -skippable if it satisfies any of the following conditions:*

1. v does not belong to any far-apart pair considered before (x, y) (Lemma 5 in [12]);
2. $\min\{d(x, v), d(y, v)\} \leq \delta_L$ (Lemma 3);
3. $2 \text{ecc}(v) - d(x, v) - d(y, v) < 4\delta_L + 2 - d(x, y)$ (Lemma 8 in [12]);
4. $\text{ecc}(v) + d(x, y) - 3\delta_L - \frac{3}{2} < \max\{d(x, v), d(y, v)\}$ (Lemma 9 in [12]).

A node that does not satisfy any condition of Lemma 4 is defined as (x, y, δ_L) -acceptable, and so it must be considered. This class is further refined in [12] with the subset of *c-valuable* vertices, where c is any fixed node (a good choice is a node with small eccentricity or centrality) as specified in the following Lemma.

Algorithm 1: Algorithm for computing the hyperbolicity proposed in [12]

Input: $\mathcal{F} = (\{x_1, y_1\}, \dots, \{x_N, y_N\})$, an ordered list of far-apart pairs.
Input: d , the distance matrix.

```

1  $\delta_L \leftarrow 0$ 
2  $\text{mates}[v] \leftarrow \emptyset$  for each  $v$ 
3 for  $i \in [1, N]$  do
4   if  $d(x_i, y_i) \leq 2\delta_L$  then
5     return  $\delta_L$ 
6    $(\text{acceptable}, \text{valuable}) \leftarrow \text{computeAccVal}()$ 
7   for  $v \in \text{valuable}$  do
8     for  $w \in \text{mates}[v]$  do
9       if  $w \in \text{acceptable}$  then
10         $\delta_L \leftarrow \max\{\delta_L, \delta(x_i, y_i, v, w)\}$ 
11   add  $y_i$  to  $\text{mates}[x_i]$ 
12   add  $x_i$  to  $\text{mates}[y_i]$ 
13 return  $\delta_L$ 

```

Lemma 5 ([12]). *Let c be any fixed node. A (x, y, δ_L) -acceptable node v is c -valuable if $2d(c, v) - 2\delta_L > d(x, v) + d(y, v) - d(x, y)$.*

In Algorithm 1, the 4-tuples considered with far-apart pair (x_i, y_i) are such that v is c -valuable, w is (x, y, δ_L) -acceptable and (v, w) is a far-apart pair seen previously. Overall, the classification of the nodes is done in overall time $\mathcal{O}(n)$ for a given pair (x_i, y_i) and lower bound δ_L . The experiments reported in [12] show that this classification leads to a significant reduction of the number of considered 4-tuples as well as computation time.

4.3 Hub labeling

A main bottleneck of Algorithm 1 comes from the $\Theta(n^2)$ memory usage. This can be alleviated by using hub labeling [34], a technique that allows to encode distances in a graph. The technique is also called two-hop labeling [21]. It appears to give a very efficient space-time tradeoff in practice. We tried to use it as a replacement of the distance matrix in Algorithm 1. It perfectly fits in memory with all practical graphs we could test. However, computing all distances from a given vertex is orders of magnitude slower than performing a BFS from that vertex. As the technique appeared as an inefficient way to extend Algorithm 1, we abandoned it.

4.4 Our algorithm

To improve upon Algorithm 1, and in particular to reduce the memory usage in practice, we do the following.

- We use the far-apart pair iterator presented in Section 3 to avoid the pre-computation and storage of the list of far-apart pairs. The use of Improvement 1 of Section 3 ensures that at the end of the visit of F_u^d , it contains only the vertices w such that u is w -far, and so the pair (v, w) is far-apart.

Algorithm 2: New algorithm for computing the hyperbolicity

```
Input:  $G = (V, E)$ 
1 Initialize the far-apart pair iterator  $F$ 
2  $\delta_L \leftarrow \text{lowerBoundHeuristic}(G)$ 
3 while  $\text{has\_next}(F)$  do
4    $(x, y) \leftarrow \text{next}(F)$  // provides  $d(x, y)$ 
5   if  $d(x, y) \leq 2\delta_L$  then
6      $\perp$  return  $\delta_L$ 
7    $d_x \leftarrow c_{x,y,\delta_L}\text{-prunedBFS}(x)$ 
8    $d_y \leftarrow c_{y,x,\delta_L}\text{-prunedBFS}(y)$ 
9    $(\text{acceptable}, \text{valuable}) \leftarrow \text{computeAccVal}()$ 
10  for  $v \in \text{valuable}$  do
11    for  $w \in \text{mates}(F, v, d)$  do // provides  $d(v, w)$ 
12      if  $w \in \text{acceptable}$  then
13         $\perp$   $\delta_L \leftarrow \max\{\delta_L, \delta(x, y, v, w)\}$ 
14 return  $\delta_L$ 
```

- We design a function $\text{mates}(F, v, d)$ to iterate over the far-apart pairs at distance $d \geq d(x, y)$ that involve v and that have previously been reported by $\text{next}(F)$.

When $d(x, y) < d \leq \text{diam}(G)$, this function simply yields vertices from F_v^d since the order of operations of the algorithm when using Improvement 1 ensures that F_v^d contains only the vertices forming far-apart pairs with v .

When $d = d(x, y)$, we have to ensure that a v -far vertex w is yielded if and only if the pair (v, w) has previously been reported by a call to $\text{next}(F)$ (so the pair (v, w) is far-apart). To do so, we modify the far-apart pairs iterator and its $\text{next}(F)$ function as follows. We use an extra hash map T (initially empty) associating to a vertex u the subset of u -far vertices at distance d from u that have previously been reported by $\text{next}(F)$. Then, when $\text{next}(F)$ yields a far-apart pair (x, y) , we store y in T_x and x in T_y . This way, when $d = d(x, y)$, function $\text{mates}(F, v, d)$ simply has to yield vertices from T_v . Finally, as soon as function $\text{next}(F)$ starts reporting pairs at distance $d - 1$, we exchange hash maps T and F^d (alternative implementation of Improvement 1), and proceed with a cleared hash map T .

- Instead of giving the distance matrix as input to the algorithm, we compute for each pair (x, y) the BFS distances from x and y before the call to $\text{computeAccVal}()$. Since the distance $d(v, w)$ is obtained while extracting the mates of v from the far-apart pair iterator, we get all the needed distances to compute $\delta(x, y, v, w)$.

Although these repeated computations of BFS distances have no impact on the overall time complexity of the algorithm, which remains in $\mathcal{O}(n^4)$, they represent a significant computation time in practice. To reduce the impact on the overall computation time, we propose Optimisations 2 (Section 4.4.2) and 3 (Section 4.4.3) below.

See Algorithm 2 for the overall presentation of our algorithm. In the following, we present some optimisations aiming at reducing the computation time.

4.4.1 Optimisation 1: Lower bound initialisation

The technique of acceptable and valuable nodes becomes more efficient when δ_L is larger as Inequalities 3 and 4 of Lemma 4 and the inequality of Lemma 5 become stricter. For that reason, we first use the heuristic described in [23] to set an initial value to δ_L . It is referenced as lowerBoundHeuristic in Algorithm 2.

4.4.2 Optimisation 2: Cache of BFSs

We use a cache of BFSs to avoid recomputing a BFS that has recently been computed. This cache has a bounded capacity (e.g., 1000 BFSs). Observe that even a cache of 2 BFSs is beneficial as function `next(F)` reports successively all far-apart pairs at distance d involving a vertex x and such that $x < y$.

4.4.3 Optimisation 3: Pruned BFS for searching acceptable nodes

When considering a pair (x, y) , we perform BFS searches from both x and y to obtain distances from x and y and detecting both acceptable and valuable nodes. Lemma 4 allows to restrict both searches as follows.

First, we observe that if a node v satisfies $\text{ecc}(v) - d(x, v) < 3\delta_L - \frac{3}{2} - d(x, y)$, then Lemma 4 applies and v is (x, y, δ_L) -skippable. A (x, y, δ_L) -acceptable node v must thus satisfy:

$$\text{ecc}(v) - d(x, v) \geq c_{x,y,\delta_L}, \quad \text{where} \quad c_{x,y,\delta_L} = 3\delta_L - \frac{3}{2} - d(x, y). \quad (1)$$

We then define the c_{x,y,δ_L} -pruned BFS search from x as a BFS search from x that visits only nodes satisfying Equation (1). More precisely, when visiting a node u , we enqueue only neighbors v of u that satisfy Equation (1). Note that c_{x,y,δ_L} is constant given x, y and δ_L . We can then safely replace the regular BFS from x by a pruned BFS as stated by the following lemma.

Lemma 6. *A c_{x,y,δ_L} -pruned BFS search from x visits all (x, y, δ_L) -acceptable nodes.*

Proof. The proof follows from two facts. First, any (x, y, δ_L) -acceptable node must satisfy Equation (1). Second, any node $v \neq x$ satisfying Equation (1) must have some neighbor closer to x that satisfies Equation (1). This second condition proven below allows to easily prove by induction that all nodes satisfying Equation (1) are visited by a c_{x,y,δ_L} -pruned BFS search from x .

To prove the above second condition, consider a neighbor u of v at distance $d(x, v) - 1$ from x . By triangle inequality, we have $\text{ecc}(u) \geq \text{ecc}(v) - 1$, implying $\text{ecc}(u) - d(x, u) \geq \text{ecc}(v) - d(x, v)$. As v satisfies Equation (1), so does u . \square

Notice that the set V_{x,y,δ_L} of nodes visited by a c_{x,y,δ_L} -pruned BFS search from x is larger than the set of (x, y, δ_L) -acceptable nodes. Indeed, Conditions 1 to 3 of Lemma 4 cannot be used for the search. Furthermore, remark that the set V_{x,y,δ_L} depends on the distance $d(x, y)$ and not on the precise node y . That is,

Lemma 7. *For any z such that $d(x, y) = d(x, z)$, we have $V_{x,z,\delta_L} = V_{x,y,\delta_L}$.*

The following results are direct consequences of Equation (1) and Lemma 7.

Corollary 3. *For any z such that $d(x, y) \geq d(x, z)$, we have $V_{x,z,\delta_L} \subseteq V_{x,y,\delta_L}$.*

Corollary 4. *For any $\delta > \delta_L$, we have $V_{x,y,\delta} \subseteq V_{x,y,\delta_L}$.*

Lemma 7 and corollaries 3 and 4 enable the use of our pruned BFS in combination with a cache of BFSs. Indeed, during the execution of the algorithm both the considered distance $d(x, y)$ decreases and the lower bound δ_L increases. Hence, a cached c_{x,y,δ_L} -pruned BFS remains valid for future use. Hence, for Line 7 of Algorithm 2, we first check if a BFS from x is in the cache, and if so we retrieve it. Otherwise, we perform a c_{x,y,δ_L} -pruned BFS search from x and add it to the cache. We proceed similarly for y .

Observe also that to determine the sets of (x, y, δ_L) -acceptable and c -valuable vertices, it suffices to consider vertices that have been visited by the c_{x,y,δ_L} -pruned BFS search, or by the c_{y,x,δ_L} -pruned BFS search if this set is smaller.

5 Experimental results

In this section we conduct experiments to test the performance of our new algorithm in comparison to the previous state-of-the-art one. We additionally test the impact that each proposed optimization has on the overall performance. To gain further insights into our main tool, the algorithm for efficiently enumerating far-apart pairs, we also conduct experiments to analyze the performance and the number of far-apart pairs that graphs exhibit.

5.1 Implementation notes

We have implemented all the algorithms in C++ and our code is available [28]. In this section, we discuss some implementation choices.

We have implemented a cache of BFSs with bounded capacity κ that additionally maintains the age information of the data it stores. We use a counter τ , initialized to 0, that is increased by one each time a BFS that is already in the cache is accessed, or a BFS that is not in the cache is added. We associate to a cached BFS from a vertex x an age information a_x , initialized to the value of τ at insertion time. Then, we set a_x to the current value of τ each time the BFS from x is accessed. Hence, the last accessed BFS is such that $a_x = \tau$. The use of a hash map associating to a vertex the corresponding BFS and age information enables to decide in time $\mathcal{O}(1)$ if a BFS from x is in the cache, and if so to return a pointer on the corresponding data. Updating the age information is also done in time $\mathcal{O}(1)$. The insertion of a BFS in the cache takes time $\mathcal{O}(1)$ if the cache is not full, and time $\mathcal{O}(\kappa)$ as soon as it has reached its maximum capacity. Indeed, the insertion of a BFS when the cache is full requires to remove first the BFS with largest age information, that is, the one of the vertex x maximizing $\tau - a_x$, and so with smallest a_x . Note that for κ much smaller than n , the time required for managing the cache is negligible with respect to the time required for a BFS. Observe that this cache will be accessed $\mathcal{O}(n^2)$ times by Algorithm 2, and more precisely at most twice the number of far-apart pairs at distance $2\delta(G)$ or more.


5.2 Data & Hardware

We test graphs from the BioGRID interaction database (BG-*) [54]; a protein interactions network (dip20170205) [56]; and graphs of the autonomous systems from the Internet (CAIDA_as-* and DIMES-*) [63, 57]. We also test social networks (Epinions, Hollywood, Slashdot, Twitter), co-author graphs (ca-*, dblp), computer networks (Gnutella, Skitter), web graphs (NotreDame),

road networks (oregon2, FLA-t), a 3D triangular mesh (buddha), and grid-like graphs from VLSI applications (alue7065) and from computer games (FrozenSea). The data is available from snap.stanford.edu, webgraph.di.unimi.it, www.dis.uniroma1.it/challenge9, graphics.stanford.edu, steinlib.zib.de, and movingai.com. Furthermore, we test synthetic inputs: grid300-10 and grid500-10 are square grids with respective sizes 301×301 , and 501×501 where 10% of the edges were randomly deleted. Each graph is taken as an undirected unweighted graph and we consider only its largest bi-connected component (available from [28]). See Table 1 for the characteristics of these graphs.

We used a computer equipped with two Intel Xeon Gold 6240 CPUs operating at 2.6GHz and 192G RAM to run our experiments. Note that our code uses a single thread.

5.3 Parameter choice

As there are instances which do not terminate in reasonable time or need unreasonable amounts of memory, we cap both resources at a fixed value to obtain a clear picture. More precisely, for each graph and each experiment, we kill the process as soon as it takes more than 6 hours or uses more than 192 GB of memory. We use the  symbol to indicate a killed process and, if applicable, put this symbol in the column (i.e., time or memory) that caused the process to be killed. Furthermore, for the cache described in Section 5.1, we use a size of 1000, unless mentioned otherwise.

5.4 Comparison to previous work

To evaluate the performance improvement over previous approaches, we compare to the practical state-of-the-art algorithm, which was given in [12]. To this end, we measure the single-threaded computation time, the memory, and the best upper and lower bound found for our new algorithm as well as the one by Borassi et al. [12]. The results of this experiment are shown in Table 2.

There are several interesting observations that one can derive from these experiments. First, the new approach that we present in this paper needs significantly less memory. More precisely, the memory consumption is up to a factor of 28 times lower than for [12] (see the slashdot0902-d graph), only considering the graphs where both approaches stay within the limits. If we also consider the graphs where [12] runs out of memory, we use at least a factor 177 less memory (see grid300-10). Mainly due to this significant reduction of memory consumption, we are able to compute the hyperbolicity for graphs which were previously out of reach due to excessive amounts of memory which would have been necessary. However, there are also graphs for which we can compute the hyperbolicity, which hit the timeout limit of 6 hours when using [12]. Note that in these two cases, the computation roughly takes between 4 and 5 hours, so [12] takes at least 1 to 2 hours more time to finish on these instances. In general, we note that all graphs for which [12] computes the hyperbolicity within the time and memory limits, also stay within the limits using our approach. Over all 40 benchmark instances, our new approach computes the hyperbolicity for 8 more graphs than [12]. If the hyperbolicity cannot be computed within the limits with our new approach, then this is due to running out of time instead of running out of memory. This again shows that we achieve a drastic reduction in memory consumption for computing the hyperbolicity.

While we are sometimes faster than [12], one could also assume that sometimes it is the other way around. This is indeed the case, and while for almost all graphs we are slightly faster or slightly slower, there is one instance (namely, DIMES_201012) where we are a factor of 13 slower than [12].

Table 1: Some graph parameters of all the graphs that we use in our experiments. Note that for all graphs, we extracted the largest biconnected component and restrict to this subgraph in our experiments.

Graph	#nodes	#edges	radius	diameter	mean ecc.
BG-MV-Physical	9851	45558	11	22	14.45
BG-S-Affinity_Capture-MS	17793	174210	6	12	9.02
BG-S-Affinity_Capture-RNA	3339	10408	3	5	4.00
BG-S-Affinity_Capture-Western	9971	44331	10	20	12.60
BG-S-Biochemical_Activity	2944	10444	7	13	9.20
BG-S-Dosage_Rescue	1521	4143	9	17	12.21
BG-S-Synthetic_Growth_Defect	3013	21341	4	8	5.72
BG-S-Synthetic_Lethality	2258	12187	4	7	5.45
dip20170205	13969	60621	10	17	12.95
CAIDA_as_20000102	4009	10101	4	8	5.62
CAIDA_as_20040105	10424	27061	4	8	5.98
CAIDA_as_20050905	12957	33541	5	8	6.11
CAIDA_as_20110116	23214	89783	4	8	6.04
CAIDA_as_20120101	25614	109180	4	8	6.10
CAIDA_as_20130101	27454	124672	5	10	7.21
CAIDA_as_20131101	29432	143000	5	9	6.46
DIMES_201012	18764	84851	4	7	5.16
DIMES_201204	16907	66489	4	7	5.07
p2p-Gnutella09	5606	23510	5	8	6.31
gnutella31-d	33812	119127	6	9	7.41
notreDame-d	134958	833732	18	36	20.99
ca-CondMat	17234	84595	6	12	8.44
ca-HepPh	9025	114046	6	11	7.83
ca-HepTh	5898	20983	7	11	8.63
com-dblp.ungraph	211409	883570	8	15	10.92
dblp-2010	140610	572873	10	17	11.99
email-Enron	20416	163257	5	9	6.54
epinions1-d	36111	365253	5	9	6.36
facebook_combined	3698	85963	4	6	5.26
loc-brightkite	33187	188577	6	11	7.95
loc-gowalla_edges	137519	887929	6	11	7.91
slashdot0902-d	51528	473218	5	8	6.04
oregon2_010331	7602	27870	4	8	5.89
t.FLA-w	691175	941893	890	1780	1378.52
buddha-w	543652	1631574	244	487	360.44
froz-w	749520	2895228	812	1451	1130.38
grid300-10	90211	162152	300	600	450.50
grid500-10	250041	449831	500	1000	750.49
z-alue7065	34040	54835	213	426	319.43

Table 2: Comparing the memory and time consumption for our algorithm and [12]. If there is a timeout, we state the best lower and upper bounds on the hyperbolicity that we obtained.

Graph	Borassi et al. [12]			Our Algorithm		
	time (s)	memory	hyperb.	time (s)	memory	hyperb.
BG-MV-Physical	10429	1.03 GB	4.5	6725	166.33 MB	4.5
BG-S-Affinity_Capture-MS	⌚	–	[2.0, 4.0]	⌚	–	[2.0, 4.0]
BG-S-Affinity_Capture-RNA	0.77	157.18 MB	2.0	0.67	57.29 MB	2.0
BG-S-Affinity_Capture-Western	7827	1.05 GB	4.0	5255	154.66 MB	4.0
BG-S-Biochemical_Activity	8.45	104.90 MB	3.0	16.32	46.16 MB	3.0
BG-S-Dosage_Rescue	11.97	30.54 MB	4.0	14.31	24.43 MB	4.0
BG-S-Synthetic_Growth_Defect	1.73	108.76 MB	2.0	2.93	43.82 MB	2.0
BG-S-Synthetic_Lethality	1.11	77.09 MB	2.0	1.75	33.40 MB	2.0
dip20170205	⌚	–	[4.5, 5.0]	18318	339.30 MB	4.5
CAIDA_as_20000102	1.28	260.87 MB	2.5	1.25	56.40 MB	2.5
CAIDA_as_20040105	18.43	1.90 GB	2.5	25.87	166.40 MB	2.5
CAIDA_as_20050905	16.9	2.37 GB	3.0	21.77	203.72 MB	3.0
CAIDA_as_20110116	13806	8.41 GB	2.0	13491	556.72 MB	2.0
CAIDA_as_20120101	⌚	–	[2.0, 2.5]	⌚	–	[2.0, 2.5]
CAIDA_as_20130101	2961.95	10.09 GB	2.5	3535.82	1.47 GB	2.5
CAIDA_as_20131101	⌚	–	[2.0, 2.5]	16108	593.03 MB	2.5
DIMES_201012	465.76	4.86 GB	2.0	6043	482.41 MB	2.0
DIMES_201204	66.35	4.34 GB	2.0	56.17	304.48 MB	2.0
p2p-Gnutella09	2.9	381.64 MB	3.0	2.35	98.38 MB	3.0
gnutella31-d	139.54	13.13 GB	3.5	109.49	1.51 GB	3.5
notreDame-d	–	⌚	–	4514	53.02 GB	8.0
ca-CondMat	112.76	3.38 GB	3.5	172.94	281.18 MB	3.5
ca-HepPh	36.97	1.17 GB	3.0	86.35	152.86 MB	3.0
ca-HepTh	4.02	407.78 MB	4.0	2.99	91.54 MB	4.0
com-dblp.ungraph	–	⌚	–	5449	14.20 GB	5.0
dblp-2010	–	⌚	–	6904	7.51 GB	5.5
email-Enron	754.13	5.36 GB	2.5	897.72	404.38 MB	2.5
epinions1-d	696.73	18.60 GB	2.5	2331.33	759.20 MB	2.5
facebook_combined	6919	248.00 MB	1.5	4551	136.65 MB	1.5
loc-brightkite	910.39	12.81 GB	3.0	910.12	812.98 MB	3.0
loc-gowalla_edges	–	⌚	–	14478	4.04 GB	3.5
slashdot0902-d	9560	37.54 GB	2.5	4718	1.34 GB	2.5
oregon2_010331	72.99	980.82 MB	2.0	65.13	133.76 MB	2.0
t.FLA-w	–	⌚	–	⌚	–	[81.0, 835.5]
buddha-w	–	⌚	–	⌚	–	[93.0, 221.0]
froz-w	–	⌚	–	⌚	–	[367.5, 633.5]
grid300-10	–	⌚	–	10.13	1.08 GB	280.0
grid500-10	–	⌚	–	99.64	2.99 GB	463.0
z-alue7065	35.01	9.07 GB	138.0	33.47	431.18 MB	138.0

We suspect that for this graph, there are many nodes from which we repeatedly perform BFSs to obtain distances.

We specifically want to highlight two grid graphs, grid300-10 and grid500-10, for which our approach only takes 10 seconds and 90 seconds, respectively, while [12] fails to compute the hyperbolicity due to exceeding the 192 GB memory limit. Note that for the former grid graph, only 1.08 GB of memory is needed for our approach and thus the memory usage is lower by at least a factor of 177. This drastic reduction of memory comes from the highly structured input: a perfect grid only has 2 far-apart pairs, the two pairs of opposing corners. Thus, the hyperbolicity computation only needs to consider these two pairs and computing the shortest paths between all pairs is hugely wasteful. Our approach includes this natural intuition to lazily compute the distances between the pairs of nodes to avoid huge running times and memory consumption in cases like these.

For a better overview, we plotted a comparison of the time and memory usage of both algorithms, see Figures 1 and 2. In Figure 1 we can see that the times indeed are very similar for both approaches except for the single outlier mentioned above. In Figure 2 we can see the drastic reduction in memory usage. Note again that the graphs included in this figure are only the graphs on which both algorithms terminate within the limits. The graphs where [12] exceeds the limits while our approach stays within the limits are not shown. In particular, the memory reduction by at least a factor of 177 is not shown in this plot. Furthermore, we can see that with increasing memory consumption, also the advantage that our algorithm has over [12] with respect to memory usage becomes more and more pronounced.

5.5 Impact of different optimizations

Additional to the comparison with the previous state-of-the-art algorithm, we also conduct experiments to obtain insights into the impact of the different optimizations that we used in our algorithm. In particular, we want to know the impact of the lower bound initialization (see Section 4.4.1), the cache size (see Section 4.4.2), and the pruning (see Section 4.4.3). We conduct experiments with all our graphs with the normal cache size of 1000 with no heuristic and no pruning, heuristic but no pruning, pruning but no heuristic, and heuristic and pruning. To gain more insight into the effect of the BFS cache, we also conduct the last experiment with a cache size of 2, which means that we only store the current two BFSs in memory. See Table 3 for the results of these experiments.

The overall results of these experiments are that there is no significant benefit for the running time when using the heuristically computed initial lower bound. However, the pruning has significant impact depending on the graph: sometimes it does not help much, but in several cases it does decrease the running time by a larger factor – up to a factor of 5 (see the DIMES_201204). The cache size also has very different impact on different graphs. It can reduce the time by a factor of 1.4, but it also sometimes increases the running time. However, the positive effects outweigh the negative effects and we thus consider it a worthwhile optimization.

5.6 BFS cache size experiments

To gain further insights into how the BFS cache we use affects the running time behavior, we run experiments with different cache sizes on selected graphs, see Figure 3. We again see different behavior on different instances. While the times monotonously decrease with increasing cache size for two of the graphs, for the other two we have increase-decrease and decrease-increase patterns. One reason for this behavior might be the cache size of the processor. In particular, if the graph

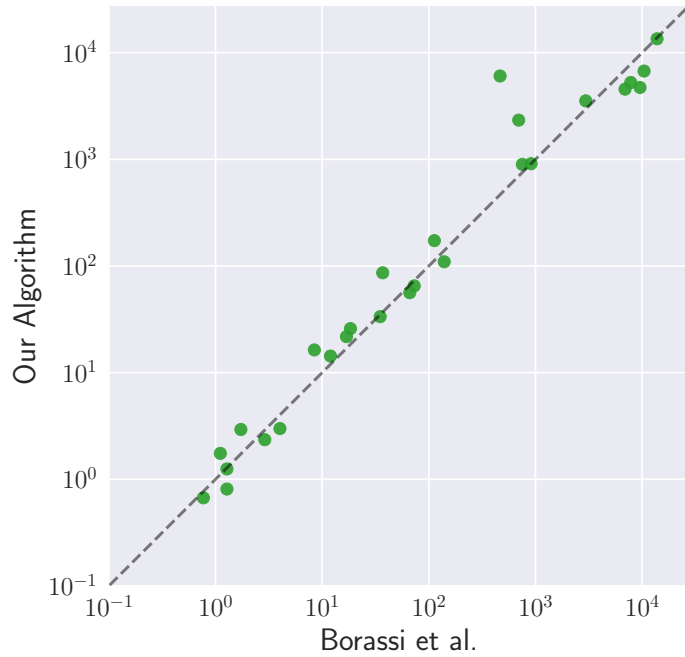


Figure 1: Comparing the *time* in seconds for computing the hyperbolicity on all graphs that finish using both algorithms. The dashed line is the identity, i.e., where both algorithms would take the same time.

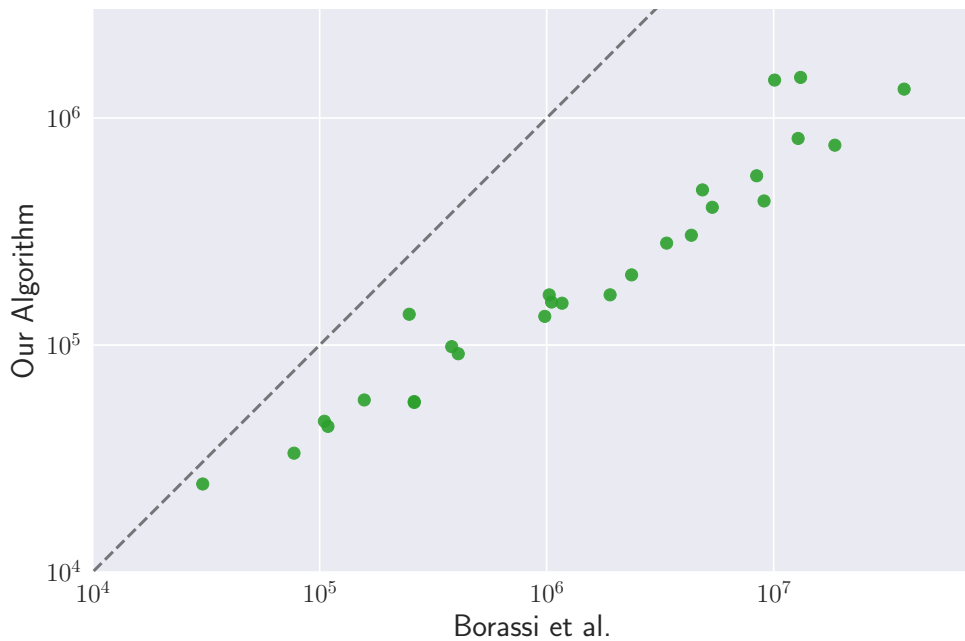


Figure 2: Comparing the *memory* in MB for computing the hyperbolicity on all graphs that finish using both algorithms. The dashed line is the identity, i.e., where both algorithms would use the same amount of memory.

Table 3: Times for computing the hyperbolicity with different optimizations enabled. All entries are in seconds. The columns with “heur” use the lower bound initialization presented in Section 4.4.1, and the columns with “prune” use the pruning of Section 4.4.3. The value of c in the second row gives the size of the BFS cache presented in Section 4.4.2.

Graph	–	heur	prune	heur & prune	
	$c = 1000$	$c = 1000$	$c = 1000$	$c = 2$	$c = 1000$
BG-MV-Physical	7040	7047	6745	6867	6725
BG-S-Affinity_Capture-MS	✗	✗	✗	✗	✗
BG-S-Affinity_Capture-RNA	0.65	0.69	0.67	0.72	0.67
BG-S-Affinity_Capture-Western	5464	5450	5236	5297	5255
BG-S-Biochemical_Activity	16.52	16.51	16.25	20.17	16.32
BG-S-Dosage_Rescue	14.91	14.74	14.58	16.47	14.31
BG-S-Synthetic_Growth_Defect	3.25	3.3	2.99	4.61	2.93
BG-S-Synthetic_Lethality	2.14	2.17	1.75	2.49	1.75
dip20170205	19088	19144	18492	18473	18318
CAIDA_as_20000102	1.24	1.26	1.23	1.72	1.25
CAIDA_as_20040105	39.89	40.03	25.69	34.97	25.87
CAIDA_as_20050905	35.63	36.91	21.66	23.76	21.77
CAIDA_as_20110116	16079	16132	13700	13412	13491
CAIDA_as_20120101	✗	✗	✗	✗	✗
CAIDA_as_20130101	3292.8	3313.78	3594.91	3848	3535.82
CAIDA_as_20131101	15785	15766	16036	16683	16108
DIMES_201012	6144	6139	6032	6574	6043
DIMES_201204	294.08	292.47	56.01	68.74	56.17
p2p-Gnutella09	2.0	2.41	1.87	2.27	2.35
gnutella31-d	172.25	191.22	91.08	104.68	109.49
notreDame-d	4317	4315	4472	4795	4514
ca-CondMat	385.97	387.39	177.16	197.9	172.94
ca-HepPh	144.67	145.0	87.82	103.41	86.35
ca-HepTh	3.84	4.06	2.75	4.71	2.99
com-dblp.ungraph	✗	✗	5352	5713	5449
dblp-2010	14677	14675	6610	6560	6904
email-Enron	1121.16	1108.76	906.45	1184.8	897.72
epinions1-d	2626.57	2607.03	2331.39	3563.06	2331.33
facebook_combined	4677	4623	4552	4609	4551
loc-brightkite	1696.58	1688.57	928.46	1108.55	910.12
loc-gowalla_edges	✗	✗	14259	13169	14478
slashdot0902-d	4492	4463	4775	6766	4718
oregon2.010331	124.18	124.45	65.62	79.84	65.13
t.FLA-w	✗	✗	✗	✗	✗
buddha-w	✗	✗	✗	✗	✗
froz-w	✗	✗	✗	✗	✗
grid300-10	10.2	10.36	10.13	8.7	10.13
grid500-10	102.21	102.31	99.84	83.08	99.64
z-alue7065	36.52	37.04	33.76	33.63	33.47

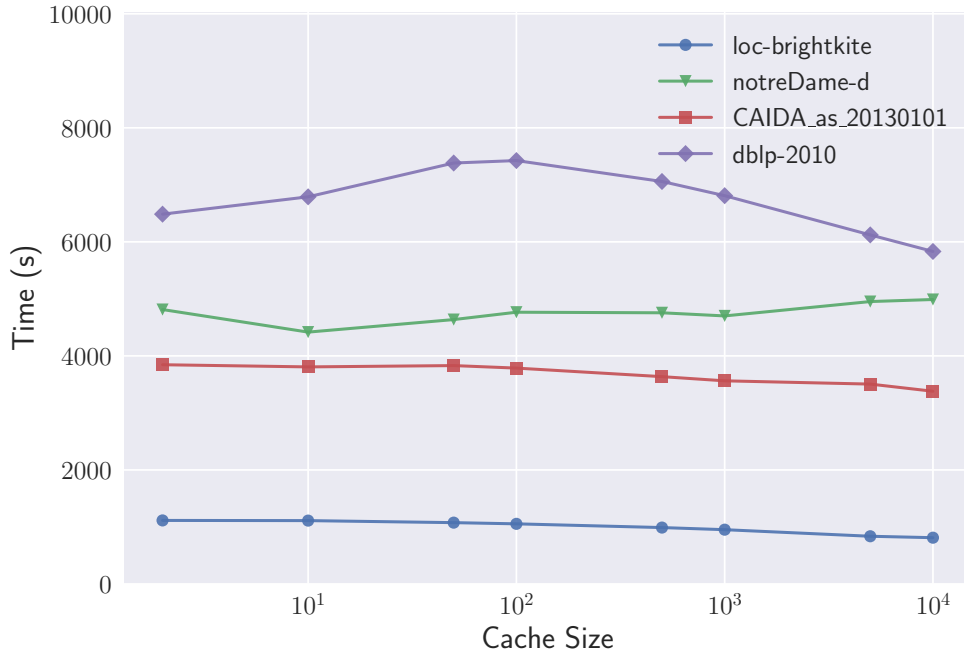


Figure 3: Plot for four graphs showing the running time development depending on the BFS cache size.

already fits in the CPU cache (e.g., L1), then computing a BFS is quite fast. Especially for large BFS cache sizes which might push the graph out of the CPU cache and also might reside in a higher level CPU cache, the computation can become slower.

5.7 Far-apart pairs iterator experiments

Finally, we perform experiments to only analyze the behavior of the far-apart pairs iterator. To this end, we let the far-apart pairs iterator run on all our benchmark graphs and measure the time as well as memory consumption. Additionally, we are interested in the number of far-apart pairs that the different instances have as well as how many pairs of an instance are necessary for the hyperbolicity calculation of our new algorithm. As the graphs are of different sizes, we put the number of pairs in the last two columns in relation to the total number of node pairs in the graph. The results of these experiments are shown in Table 4. For all except four graphs, the far-apart pairs iterator runs through in the given memory and time limits. This is one more tractable instance compared to hyperbolicity computation. Note, however, that there are graphs for which we can compute the hyperbolicity but the far-apart pairs iterator does not run through within the limits (e.g., dblp-2010). This is explained by the fact that to compute the hyperbolicity, we can stop at some point of iterating through the far-apart pairs and do not have to compute them all. More specifically, we show what percentage of all pairs are necessary to compute the hyperbolicity with our algorithm on this specific instance. Observe that for the dblp-2010 graph, only around 0.27% of pairs are relevant for the hyperbolicity computation. In general, most instances only have a single-digit percentage or less of relevant pairs for the hyperbolicity computation. The grid graphs, grid300-10 and grid500-10, even have such a small number of relevant pairs, that they are rounded

to 0 in the precision that we choose for the numbers in the table. These low numbers of relevant pairs explain the drastic memory reduction that we achieve with our algorithm.

We can also see that for many graphs, the far-apart pairs iterator is very fast, while the hyperbolicity computation takes a long time, which is explained by the fact that we might spend quadratic time per pair to compute the hyperbolicity. Considering the percentage of far-apart pairs, we see that most graphs have roughly 30% to 70% far-apart pairs. The graphs with grid structure are extreme outliers, exhibiting a very low percentage of far-apart pairs. This can be explained by the grid structure, for which – considering a grid without missing edges – only the two pairs of opposing corners of the grid are far-apart. Furthermore, the facebook_combined graph has a very large percentage of far-apart pairs. This is explained by the fact that it has, by far, the highest average degree of all the graphs we consider. For such a well-connected graph, BFS trees have a very large number of leaves as shortest paths are not extended further from most nodes, which is necessary to produce such a large number of far-apart pairs.

To further gain insights into the distribution of far-apart pairs, i.e., how they are distributed over various distances, we put all the distances between far-apart pairs into a histogram for four selected graphs, see Figure 4. While the distribution looks somewhat reasonable and expected for three of the four graphs, the notreDame-d graph shows an interesting distribution with two modes. Note that anomalies with respect to coreness were already found in this graph [59] where a special structure around a propeller-shaped subgraph was identified. We suspect that this structure connects a large fraction of pairs and could be related to this bi-modal observation. Finally, we also marked the part of the histogram that contains the far-apart pairs that are relevant for the hyperbolicity calculation. We observe that this region occurs after the peak of the histogram in all four instances that we show. In the notreDame-d instance, it even completely excludes the first, much larger mode.

6 Conclusion

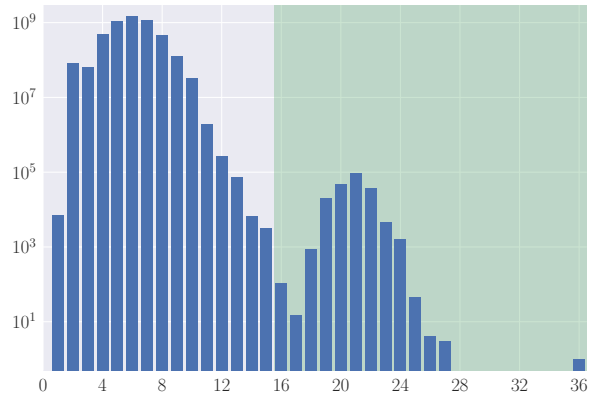
In this work, we developed a fundamental algorithm to iterate over all far-apart pairs in non increasing distance. As primary application we consider the computation of graph hyperbolicity. Our new algorithm enables us to compute, for the first time, the hyperbolicity of some graphs with more than a hundred thousand nodes with non trivial structure (e.g., notreDame-d, loc-gowalla_edges, com-dblp_undgraph). We reduce the memory usage significantly, while not compromising on performance. Non-trivial graphs with more than five hundred thousands nodes unfortunately still remain out of reach with our method. We thus plan to investigate alternative approaches in future work in order to get closer to the million nodes barrier. Furthermore, we believe that iterating over far-apart pairs in non increasing distance is such a fundamental task, that our work will enable faster algorithms also in other settings.

References

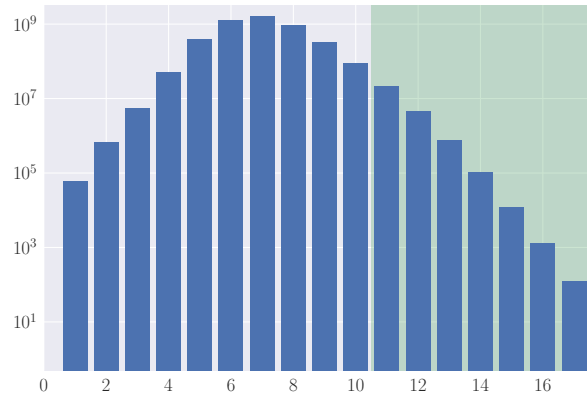
- [1] Muad Abu-Ata and Feodor F. Dragan. Metric tree-like structures in real-life networks: an empirical study. *Networks*, 67(1):49–69, 2016.

Table 4: Time and memory consumption of the far-apart pairs iterator only. We additionally show the percentage of far-apart pairs (“far pairs”) and also the percentage of pairs which have to be considered during the hyperbolicity computation of our algorithm (“hyp pairs”).

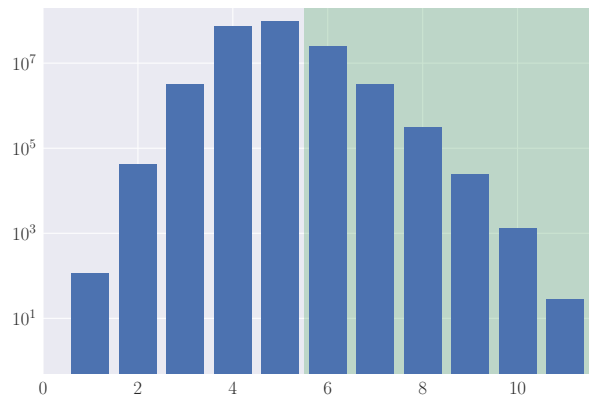
Graph	time (s)	memory	far pairs (%)	hyp pairs (%)
BG-MV-Physical	21.85	391.23 MB	30.91	4.043
BG-S-Affinity_Capture-MS	94.32	1.33 GB	36.92	–
BG-S-Affinity_Capture-RNA	3.29	124.32 MB	70.86	0.602
BG-S-Affinity_Capture-Western	23.27	416.42 MB	32.42	5.355
BG-S-Biochemical_Activity	1.82	75.47 MB	42.95	8.196
BG-S-Dosage_Rescue	0.35	29.13 MB	24.67	5.223
BG-S-Synthetic_Growth_Defect	2.13	87.80 MB	46.12	17.973
BG-S-Synthetic_Lethality	1.08	54.61 MB	42.75	24.658
dip20170205	45.85	686.76 MB	30.12	11.245
CAIDA_as_20000102	4.4	150.73 MB	63.94	5.456
CAIDA_as_20040105	40.36	835.38 MB	67.35	5.806
CAIDA_as_20050905	65.66	1.25 GB	69.01	0.502
CAIDA_as_20110116	261.31	3.64 GB	69.56	40.066
CAIDA_as_20120101	361.92	4.39 GB	68.86	–
CAIDA_as_20130101	438.98	5.02 GB	68.89	5.292
CAIDA_as_20131101	490.41	5.81 GB	69.26	5.075
DIMES_201012	174.08	3.19 GB	74.89	13.431
DIMES_201204	144.55	2.40 GB	72.51	21.459
p2p-Gnutella09	6.47	190.00 MB	28.42	4.036
gnutella31-d	407.55	4.67 GB	29.49	3.787
notreDame-d	9817	86.14 GB	54.38	0.002
ca-CondMat	105.14	1.49 GB	44.0	4.785
ca-HepPh	25.7	453.93 MB	42.31	8.63
ca-HepTh	7.82	196.91 MB	33.72	1.961
com-dblp.ungraph	–	⚠	–	–
dblp-2010	14008	87.83 GB	48.38	0.268
email-Enron	172.75	2.67 GB	55.87	10.95
epinions1-d	630.49	7.38 GB	45.46	6.683
facebook_combined	4.5	158.98 MB	89.08	71.779
loc-brightkite	425.11	4.45 GB	36.35	5.019
loc-gowalla_edges	10509	63.43 GB	33.19	0.683
slashdot0902-d	1442.04	15.35 GB	45.72	5.238
oregon2_010331	19.54	453.85 MB	66.43	31.757
t.FLA-w	–	⚠	–	–
buddha-w	⚠	–	–	–
froz-w	⚠	–	–	–
grid300-10	345.36	4.24 GB	0.04	0.0
grid500-10	2723.73	19.24 GB	0.04	0.0
z-alue7065	46.95	940.21 MB	0.06	0.002



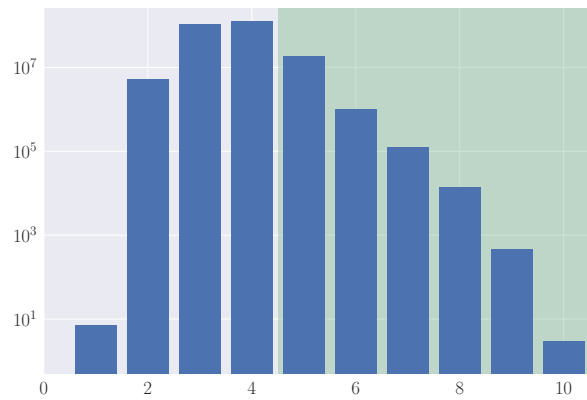
(a) notreDame-d



(b) dblp-2010



(c) loc-brightkite



(d) CAIDA_as_20130101

Figure 4: Histograms of distribution of far-apart pairs for selected graphs. On the x-axis we plot the distance and on the y-axis the number of far-apart node pairs that have this distance. The green area shows the distance range for pairs that have to be evaluated in our hyperbolicity algorithm. We can see that this area always excludes the peak of the histogram in the shown examples.

- [2] Takuya Akiba, Yoichi Iwata, and Yuki Kawata. An exact algorithm for diameters of large real directed graphs. In *International Symposium on Experimental Algorithms - SEA*, volume 9125 of *Lecture Notes in Computer Science*, pages 56–67. Springer, 2015.
- [3] Takuya Akiba, Yoichi Iwata, and Yuichi Yoshida. Fast exact shortest-path distance queries on large networks by pruned landmark labeling. In *ACM SIGMOD International Conference on Management of Data - SIGMOD*, SIGMOD, pages 349–360, New York, NY, USA, 2013. Association for Computing Machinery.
- [4] Réka Albert, Bhaskar DasGupta, and Nasim Mobasher. Topological implications of negative curvature for biological and social networks. *Physical Review E*, 89(3):032811, 2014.
- [5] Hend Alrasheed and Feodor F. Dragan. Core-periphery models for graphs based on their δ -hyperbolicity: An example using biological networks. In *6th Workshop on Complex Networks - CompleNet*, volume 597 of *Studies in Computational Intelligence*, pages 65–77. Springer, 2015.
- [6] Hans-Jürgen Bandelt and Henry Martyn Mulder. Distance-hereditary graphs. *Journal of Combinatorial Theory, Series B*, 41(2):182–208, 1986.
- [7] Sergio Bermudo, José M. Rodríguez, José M. Sigarreta, and Jean-Marie Vilaire. Gromov hyperbolic graphs. *Discrete Mathematics*, 313(15):1575–1585, 2013.
- [8] Anne Berry, Romain Pogorelnik, and Geneviève Simonet. An introduction to clique minimal separator decomposition. *Algorithms*, 3(2):197–215, 2010.
- [9] Marián Boguñá, Fragkiskos Papadopoulos, and Dmitri V. Krioukov. Sustaining the Internet with hyperbolic mapping. *Nature Communications*, 1(62):1–18, October 2010.
- [10] John Adrian Bondy and Uppaluri Siva Ramachandra Murty. *Graph theory with applications*, volume 290. Macmillan London, 1976.
- [11] Michele Borassi, Alessandro Chessa, and Guido Caldarelli. Hyperbolicity measures democracy in real-world networks. *Physical Review E*, 92(3):032812, 2015.
- [12] Michele Borassi, David Coudert, Pierluigi Crescenzi, and Andrea Marino. On computing the hyperbolicity of real-world graphs. In *European Symposium on Algorithms - ESA*, volume 9294 of *Lecture Notes in Computer Science*, pages 215–226, Patras, Greece, September 2015. Springer.
- [13] Michele Borassi, Pierluigi Crescenzi, and Michel Habib. Into the square: On the complexity of some quadratic-time solvable problems. *Electronic Notes in Theoretical Computer Science*, 322:51–67, 2016.
- [14] Michele Borassi, Pierluigi Crescenzi, Michel Habib, Walter A. Kosters, Andrea Marino, and Frank W. Takes. Fast diameter and radius BFS-based computation in (weakly connected) real-world graphs: With an application to the six degrees of separation games. *Theoretical Computer Science*, 586:59–80, 2015.
- [15] Norberto Castillo-García and Paula Hernández Hernández. *A New Heuristic Algorithm for the Vertex Separation Problem*, pages 487–500. Springer International Publishing, 2018.

- [16] John Chakerian and Susan Holmes. Computational tools for evaluating phylogenetic and hierarchical clustering trees. *Journal of Computational and Graphical Statistics*, 21(3):581–599, 2012.
- [17] Pierre Charbit, Fabien de Montgolfier, and Mathieu Raffinot. Linear time split decomposition revisited. *SIAM Journal on Discrete Mathematics*, 26(2):499–514, 2012.
- [18] Victor Chepoi, Feodor F. Dragan, Bertrand Estellon, Michel Habib, and Yann Vaxès. Diameters, centers, and approximating trees of delta-hyperbolic geodesic spaces and graphs. In *Annual Symposium on Computational Geometry - SoCG*, pages 59–68. ACM, 2008.
- [19] Victor Chepoi, Feodor F. Dragan, Bertrand Estellon, Michel Habib, Yann Vaxès, and Yang Xiang. Additive spanners and distance and routing labeling schemes for hyperbolic graphs. *Algorithmica*, 62(3-4):713–732, 2012.
- [20] Victor Chepoi, Feodor F. Dragan, and Yann Vaxès. Core congestion is inherent in hyperbolic networks. In Philip N. Klein, editor, *ACM-SIAM Symposium on Discrete Algorithms - SODA*, pages 2264–2279. SIAM, 2017.
- [21] Edith Cohen, Eran Halperin, Haim Kaplan, and Uri Zwick. Reachability and distance queries via 2-hop labels. In David Eppstein, editor, *ACM-SIAM Symposium on Discrete Algorithms - SODA*, pages 937–946. ACM/SIAM, 2002.
- [22] Nathann Cohen, David Coudert, Guillaume Ducoffe, and Aurélien Lancin. Applying clique-decomposition for computing gromov hyperbolicity. *Theoretical Computer Science*, 690:114–139, 2017.
- [23] Nathann Cohen, David Coudert, and Aurélien Lancin. On computing the gromov hyperbolicity. *ACM Journal of Experimental Algorithmics*, 20:1.6:1–1.6:18, 2015.
- [24] David Coudert and Guillaume Ducoffe. Recognition of C_4 -free and $1/2$ -hyperbolic graphs. *SIAM Journal on Discrete Mathematics*, 28(3):1601–1617, September 2014.
- [25] David Coudert and Guillaume Ducoffe. Revisiting Decomposition by Clique Separators. *SIAM Journal on Discrete Mathematics*, 32(1):682 – 694, January 2018.
- [26] David Coudert, Guillaume Ducoffe, and Alexandru Popa. Fully polynomial FPT algorithms for some classes of bounded clique-width graphs. *ACM Transactions on Algorithms*, 15(3):1–57, June 2019.
- [27] David Coudert, Dorian Mazaauric, and Nicolas Nisse. Experimental evaluation of a branch-and-bound algorithm for computing pathwidth and directed pathwidth. *ACM Journal of Experimental Algorithmics*, 21(1):1.3:1–1.3:23, 2016.
- [28] David Coudert, André Nusser, and Laurent Viennot. Hyperbolicity (version 1.0). <https://gitlab.inria.fr/dcoudert/hyperbolicity/>, 2021.
- [29] Pierluigi Crescenzi, Roberto Grossi, Michel Habib, Leonardo LANZI, and Andrea Marino. On computing the diameter of real-world undirected graphs. *Theoretical Computer Science*, 514:84–95, 2013.

- [30] William H. Cunningham. Decomposition of directed graphs. *SIAM Journal on Algebraic Discrete Methods*, 3(2):214–228, 1982.
- [31] William H. Cunningham and Jack Edmonds. A combinatorial decomposition theory. *Canadian Journal of Mathematics*, 32(3):734–765, 1980.
- [32] Bhaskar DasGupta, Marek Karpinski, Nasim Mobasher, and Farzane Yahyanejad. Effect of Gromov-hyperbolicity parameter on cuts and expansions in graphs and some algorithmic implications. *Algorithmica*, 80(2):772–800, 2018.
- [33] Pierre de La Harpe and Etienne Ghys. *Sur les groupes hyperboliques d’après Mikhael Gromov*, volume 83. Progress in Mathematics, 1990.
- [34] Daniel Delling, Andrew V. Goldberg, Thomas Pajor, and Renato F. Werneck. Robust distance queries on massive networks. In Andreas S. Schulz and Dorothea Wagner, editors, *European Symposium on Algorithms - ESA*, volume 8737 of *Lecture Notes in Computer Science*, pages 321–333. Springer, 2014.
- [35] Josep Díaz, Jordi Petit, and Maria Serna. A survey of graph layout problems. *ACM Computing Surveys*, 34(3):313–356, September 2002.
- [36] Reinhard Diestel. *Graph Theory, 5th edition*, volume 173 of *Graduate Texts in Mathematics*. Springer, Heidelberg, 2017.
- [37] Feodor F. Dragan. Tree-like structures in graphs: A metric point of view. In *39th International Workshop on Graph-Theoretic Concepts in Computer Science - WG*, volume 8165 of *Lecture Notes in Computer Science*, pages 1–4. Springer, 2013.
- [38] Feodor F. Dragan, Michel Habib, and Laurent Viennot. Revisiting radius, diameter, and all eccentricity computation in graphs through certificates. *CoRR*, abs/1803.04660, 2018.
- [39] Andreas Dress, Katharina Huber, Jacobus Koolen, Vincent Moulton, and Andreas Spillner. *Basic Phylogenetic Combinatorics*. Cambridge University Press, Cambridge, UK, December 2011.
- [40] Hervé Fournier, Anas Ismail, and Antoine Vigneron. Computing the gromov hyperbolicity of a discrete metric space. *Information Processing Letters*, 115(6):576–579, 2015.
- [41] Tibor Gallai. Transitiv orientierbare graphen. *Acta Mathematica Hungarica*, 18(1):25–66, 1967.
- [42] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., USA, 1979.
- [43] Micha Gromov. Hyperbolic groups. In S.M. Gersten, editor, *Essays in Group Theory*, volume 8 of *Mathematical Sciences Research Institute Publications*, pages 75–263. Springer, New York, 1987.
- [44] Michel Habib and Christophe Paul. A survey of the algorithmic aspects of modular decomposition. *Computer Science Review*, 4(1):41–59, 2010.

- [45] Torben Hagerup. Fast breadth-first search in still less space. In Ignasi Sau and Dimitrios M. Thilikos, editors, *International Workshop on Graph-Theoretic Concepts in Computer Science - WG*, volume 11789 of *Lecture Notes in Computer Science*, pages 93–105. Springer, 2019.
- [46] Edward Howorka. On metric properties of certain clique graphs. *Journal of Combinatorial Theory, Series B*, 27(1):67–74, 1979.
- [47] Russell Impagliazzo, Ramamohan Paturi, and Francis Zane. Which problems have strongly exponential complexity? *Journal of Computer and System Sciences*, 63(4):512–530, December 2001.
- [48] W. Sean Kennedy, Iraj Saniee, and Onuttom Narayan. On the hyperbolicity of large-scale networks and its estimation. In *International Conference on Big Data*, pages 3344–3351. IEEE, 2016.
- [49] Robert Krauthgamer and James R. Lee. Algorithms on negatively curved spaces. In *IEEE Symposium on Foundations of Computer Science - FOCS*, pages 119–132. IEEE, 2006.
- [50] François Le Gall. Faster algorithms for rectangular matrix multiplication. In *IEEE Symposium on Foundations of Computer Science - FOCS*, pages 514–523, New Brunswick, NJ, USA, 2012. IEEE.
- [51] Wentao Li, Miao Qiao, Lu Qin, Ying Zhang, Lijun Chang, and Xuemin Lin. Exacting eccentricity for small-world networks. In *IEEE International Conference on Data Engineering - ICDE*, pages 785–796, April 2018.
- [52] Tamara Munzner and Paul Burchard. Visualizing the structure of the world wide web in 3d hyperbolic space. In David R. Nadeau and John L. Moreland, editors, *Symposium on Virtual Reality Modeling Language - VRML*, pages 33–38. ACM, 1995.
- [53] Damien Noguès. δ -hyperbolicité et graphes. Master’s thesis, MPRI, Université Paris 7, 2009.
- [54] Rose Oughtred, Chris Stark, Bobby-Joe Breitkreutz, Jennifer Rust, Lorrie Boucher, Christie Chang, Nadine Kolas, Lara O’Donnell, Genie Leung, Rochelle McAdam, et al. The biogrid interaction database: 2019 update. *Nucleic acids research*, 47(D1):D529–D541, 2019.
- [55] Jordi Petit. Addenda to the survey of layout problems. *Bulletin of the EATCS*, 105:177–201, 2011.
- [56] Lukasz Salwinski, Christopher S. Miller, Adam J. Smith, Frank K. Pettit, James U. Bowie, and David Eisenberg. The database of interacting proteins: 2004 update. *Nucleic acids research*, 32(suppl_1):D449–D451, 2004.
- [57] Yuval Shavitt and Eran Shir. DIMES: Let the internet measure itself. *ACM SIGCOMM Computer Communication Review*, 35(5):71–74, October 2005.
- [58] Yuval Shavitt and Tomer Tankel. On the curvature of the internet and its usage for overlay construction and distance estimation. In *Annual Joint Conference of the IEEE Computer and Communications Societies - INFOCOM*. IEEE, 2004.

- [59] Kijung Shin, Tina Eliassi-Rad, and Christos Faloutsos. Patterns and anomalies in k-cores of real-world graphs with applications. *Knowl. Inf. Syst.*, 54(3):677–710, 2018.
- [60] Mauricio A. Soto Gómez. *Quelques propriétés topologiques des graphes et applications à internet et aux réseaux*. PhD thesis, Univ. Paris Diderot (Paris 7), 2011.
- [61] Frank W. Takes and Walter A. Kosters. Computing the eccentricity distribution of large graphs. *Algorithms*, 6(1):100–118, 2013.
- [62] Robert Endre Tarjan. Decomposition by clique separators. *Discrete Mathematics*, 55(2):221–232, 1985.
- [63] The Cooperative Association for Internet Data Analysis (CAIDA). The CAIDA AS relationships dataset. <http://www.caida.org/data/active/as-relationships/>, 2013.
- [64] Virginia Vassilevska Williams, Joshua R. Wang, Richard Ryan Williams, and Huacheng Yu. Finding four-node subgraphs in triangle time. In *ACM-SIAM Symposium on Discrete Algorithms - SODA*, pages 1671–1680. SIAM, 2015.
- [65] Jörg A. Walter and Helge J. Ritter. On interactive visualization of high-dimensional data using the hyperbolic plane. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD*, pages 123–132. ACM, 2002.

A Time and space trade-offs for determining all far-apart pairs

In this section, we present algorithms offering different time and space trade-offs for the problem of computing all far-apart pairs. The space complexity considered here is the working memory, hence excluding the space needed to store the result, unless needed during computations.

A first algorithm to determine the set of far-apart pairs is to: 1) determine for each vertex $u \in V$ the set F_u of u -far vertices using BFS, and then 2) check for each vertex $v \in F_u$ if u is v -far (i.e., if $u \in F_v$). This can be done in time $\mathcal{O}(nm)$ and space $\mathcal{O}(n^2)$ since $|F_u| \in \mathcal{O}(n)$.

Another algorithm is to execute two steps for each vertex $u \in V$: 1) determine the set F_u of u -far vertices using BFS, and then 2) for each vertex $v \in F_u$, check if there is $w \in N(u)$ such that $d(u, v) < d(w, v)$. The second step requires to compute distances from w and so the time complexity of this algorithm is $\mathcal{O}((n + m) \sum_{u \in V} (1 + |N(u)|)) = \mathcal{O}(m^2)$. Observe that during the processing of each vertex u , this algorithm stores the set F_u , the distances from u , and the distances from one neighbor w of u . Hence, the space complexity is $\mathcal{O}(n)$. The second method can be improved using the bit-parallel BFS proposed in [3]. Indeed, this algorithm computes simultaneously distances from u and b of its neighbors in time $\mathcal{O}(n + m)$ and space $\mathcal{O}(n)$, assuming that b is a constant and using bitwise operations on bit vectors of size b (typically 32 or 64). Hence, the number of BFSs to perform is reduced to $\sum_{u \in V} \lceil |N(u)|/b \rceil$.

Let us now show how to modify the above algorithm to obtain an algorithm with time complexity in $\mathcal{O}(nm)$ and space complexity in $\mathcal{O}(n \text{ pw}(G))$, where $\text{pw}(G)$ denotes the pathwidth of G [35, 55]. The main idea is to compute the distances from u only once and to store them during as few iterations as possible. For that, let $\pi : V \rightarrow [n]$ be a linear ordering of the vertices (i.e., a bijective mapping) related to the pathwidth as described later and let $\Pi(V)$ be the set of all such orderings. The algorithm iterates over the vertices in the order $\pi^{-1}(1), \pi^{-1}(2), \dots, \pi^{-1}(n)$. The distances

from u are used for the processing of vertex u , and for the processing of each neighbor $w \in N(u)$. Let $w_L := \arg \min_{w \in N(u)} \pi(w)$ and $w_R := \arg \max_{w \in N(u)} \pi(w)$. Hence, distances from u must be computed at iteration $\pi(w_L)$ and stored until iteration $\pi(w_R)$. Consequently, at iteration i , we have stored distances from all the vertices in

$$\{u \in V \mid \exists uv \in E \text{ s.t. } \pi(u) < i \leq \pi(v)\} \cup \{\pi^{-1}(i)\} \cup \{v \in V \mid \exists uv \in E \text{ s.t. } \pi(u) \leq i \leq \pi(v)\}.$$

Now observe that the *pathwidth* [35, 55] of a graph G is defined as $\text{pw}(G) = \min_{\pi \in \Pi(V)} p(G, \pi)$, where $p(G, \pi) = \max_{i=1}^n |\{u \in V \mid \exists uv \in E \text{ such that } \pi(u) < i \leq \pi(v)\}|$. Hence, at iteration i we store distances from at most $2p(G, \pi) + 1$ vertices and there is an ordering ensuring to store distances from at most $2\text{pw}(G) + 1$ vertices. Consequently, the space complexity of this algorithm is in $\mathcal{O}(n \text{pw}(G))$ and its time complexity is $\mathcal{O}(nm)$ assuming that the ordering π such that $p(G, \pi) = \text{pw}(G)$ is given. However, the problem of determining an ordering π such that $p(G, \pi) = \text{pw}(G)$ is NP-complete [35, 55]. Nonetheless, efficient heuristic algorithms have been proposed [27, 15].

Finally, recall that the *bandwidth* of a graph G is defined as $\text{bw}(G) = \min_{\pi \in \Pi(V)} b(G, \pi)$, where $b(G, \pi) = \max_{uv \in E} |\pi(u) - \pi(v)|$ [35, 55]. Therefore, distances from u are stored for at most $2b(G, \pi) + 1$ iterations, and there is an ordering ensuring that this number of iterations is at most $2\text{bw}(G) + 1$. However, the problem of determining such an ordering is NP-hard [42].

B Retrieving distances from far vertices

First note the following corollary which is a direct consequence of Lemma 1:

Corollary 5. *For any $v \in V$, let F_v be the set of v -far vertices. The number of leaves of any shortest path tree rooted at v is at least $|F_v|$.*

Observe however that the lower bound of Corollary 5 does not imply the existence of a shortest path tree with $|F_v|$ leaves. For instance, consider the 4-cycle (u_1, u_2, u_3, u_4) . We have $|F_{u_1}| = |\{u_3\}| = 1$, but all shortest path trees rooted at u_1 have 2 leaves (either $\{u_2, u_3\}$ or $\{u_3, u_4\}$).

We now show that, given the v -far vertices and a constant c , one can determine the vertices at distance at least c from v .

Lemma 8. *Let $v \in V$ and let S be the set of all vertices u for which $d(u, v) \geq c$. Given the set F_v of v -far vertices and the distances from v to each of these vertices, one can compute the set S in time $\mathcal{O}(|F_v| + \sum_{u \in S} |N(u)|)$.*

Proof. First, observe that if u is v -far, then for each $w \in N(u)$ it holds that

$$d(v, u) - 1 \leq d(v, w) \leq d(v, u).$$

Furthermore, a neighbor $w \in N(u)$ with $d(v, w) = d(v, u) - 1$ is not v -far, and if a neighbor $w' \in N(u)$ is v -far then $d(v, u) = d(v, w')$.

To determine all the nodes that are at distance at least c from v , it suffices to perform a reverse breadth-first search, starting from far vertices. More precisely, let $\{L_d\}_{d \in \{c, \dots, \text{ecc}(v)+1\}}$ be a set family where L_d is initialized with the v -far vertices at distance d from v , and let $L_{\text{ecc}(v)+1} = \emptyset$. Then, consider these sets in decreasing distance d from v , starting with $d = \text{ecc}(v)$, and stopping when $d = c - 1$. For each vertex $u \in L_d$, add each $w \in N(u)$ to L_{d-1} for which $w \notin L_d \cup L_{d+1}$. At the end of this procedure, we have that if $d(v, u) = d$, then $u \in L_d$ and all the sets are disjoint, i.e. $L_d \cap L_{d'} = \emptyset$ for any $c \leq d, d' \leq \text{ecc}(v)$ with $d \neq d'$. The claimed time bound follows immediately from the description of the algorithm. \square