

Fast Estimation of Diameter and Shortest Paths (without Matrix Multiplication)*

Donald Aingworth[†] Chandra Chekuri[‡] Piotr Indyk[§] Rajeev Motwani[¶]

Abstract

In the recent past, there has been considerable progress in devising algorithms for the all-pairs shortest paths problem running in time significantly smaller than the obvious time bound of $O(n^3)$. Unfortunately, all the new algorithms are based on fast matrix multiplication algorithms that are notoriously impractical. Our work is motivated by the goal of devising purely combinatorial algorithms that match these improved running times. Our results come close to achieving this goal, in that we present algorithms with a small additive error in the length of the paths obtained. Our algorithms are easy to implement, have the desired property of being combinatorial in nature, and the hidden constants in the running time bound are fairly small.

Our main result is an algorithm which solves the all-pairs shortest paths problem in unweighted, undirected graphs with an *additive* error of 2 in time $O(n^{2.5}\sqrt{\log n})$. This algorithm returns actual paths and not just the distances. In addition, we give more efficient algorithms with running time $O(n^{1.5}\sqrt{k \log n} + n^2 \log^2 n)$ for the case where we are only required to determine shortest paths between k *specified* pairs of vertices rather than all pairs of vertices. The starting point for all our results is an $O(m\sqrt{n \log n})$ algorithm for distinguishing between graphs of diameter 2 and 4, and this is later extended to obtaining a ratio $2/3$ approximation to the diameter in time $O(m\sqrt{n \log n} + n^2 \log n)$. Unlike in the case of all-pairs shortest paths, our results for approximate diameter computation can be extended to the case of *directed* graphs with *arbitrary* positive real *weights* on the edges.

*A preliminary version of this paper appeared in the *Proceedings of the Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 547–553, Atlanta, GA, January 1996.

[†]Department of Computer Science, Stanford University. Email: donald@cs.stanford.edu. Supported by an NSF Graduate Fellowship and NSF Grant CCR-9357849.

[‡]Department of Computer Science, Stanford University. Email: chekuri@cs.stanford.edu. Supported by an OTL grant and NSF Grant CCR-9357849.

[§]Department of Computer Science, Stanford University. Email: indyk@cs.stanford.edu. Supported by an OTL grant and NSF Grant CCR-9357849.

[¶]Department of Computer Science, Stanford University. Email: rajeev@cs.stanford.edu. Supported by an Alfred P. Sloan Research Fellowship, an IBM Faculty Development Award, an OTL grant, and NSF Young Investigator Award CCR-9357849, with matching funds from IBM, Schlumberger Foundation, Shell Foundation, and Xerox Corporation.

1 Introduction

Consider the problem of computing all-pairs shortest paths (APSP) in an unweighted, undirected graph G with n vertices and m edges. The recent work of Alon, Galil, and Margalit [AGM91], Alon, Galil, Margalit, and Naor [AGMN92], and Seidel [Sei92] has led to dramatic progress in devising fast algorithms for this problem. These algorithms are based on formulating the problem in terms of matrices with *small integer* entries and using fast matrix multiplications. They achieve a time bound of $\tilde{O}(n^\omega)$ ¹ where ω denotes the exponent in the running time of the matrix multiplication algorithm used. The current best matrix multiplication algorithm is due to Coppersmith and Winograd [CW90] and has $\omega = 2.376$. In contrast, the naive algorithm for APSP performs breadth-first searches from each vertex, and requires time $\Theta(nm)$.

Given the fundamental nature of this problem, it is important to consider the desirability of implementing the algorithms in practice. Unfortunately, fast matrix multiplication algorithms are far from being practical and suffer from large hidden constants in the running time bound. Consequently, we adopt the view of treating these results primarily as indicators of the existence of efficient algorithms and consider the question of devising a purely *combinatorial algorithm* for APSP that runs in time $O(n^{3-\epsilon})$. The (admittedly vague) term “combinatorial algorithm” is intended to contrast with the more algebraic flavor of algorithms based on fast matrix multiplication. To understand this distinction, the reader may find it instructive to try and interpret the “algebraic” algorithms in purely graph-theoretic terms even with the use of the simpler matrix multiplication algorithm of Strassen [Str69]. Previously, the best known combinatorial algorithm is due to Feder and Motwani [FM91] that runs in time $O(n^3/\log n)$, yielding only a marginal improvement over the naive algorithm.

We take a step in the direction of realizing the goals outlined above by presenting an algorithm which solves the APSP problem with an *additive* error of 2 in time $O(n^{2.5}\sqrt{\log n})$. This algorithm returns actual paths and not just the distances. Note that the running time is better than the $\tilde{O}(n^{2.81})$ time bound of the more practical matrix multiplication algorithm of Strassen [Str69]. Further, as explained below, we also give slightly more efficient algorithms (for *sparse* graphs) for approximating the diameter. Our algorithms are easy to implement, have the desired property of being combinatorial in nature, and the hidden constants in the running time bound are fairly small. Our additive approximations are presented only for the case of unweighted, undirected graphs, but they can be easily generalized to the case of undirected graphs with small integer edge weights. In addition, we give more efficient algorithms with running time $O(n^{1.5}\sqrt{k\log n} + n^2 \log^2 n)$ for the case where we are only required to determine shortest paths between k specified pairs of vertices rather than all pairs of vertices.

A crucial step in the development of our result was the shift of focus to the problem of computing the diameter of a graph. This is the maximum over all pairs of vertices of the shortest path distance between the vertices. The diameter can be determined by computing all-pairs shortest path (APSP)

¹The notation $\tilde{O}(f(n))$ denotes $O(f(n) \text{ polylog}(n))$.

distances in the graph, and it appears that this is the only known way to solve the diameter problem. In fact, Fan Chung [Chu87] had earlier posed the question of whether there is an $O(n^{3-\epsilon})$ algorithm for finding the diameter without resorting to fast matrix multiplication. The situation with regard to combinatorial algorithms for diameter is only marginally better than in the case of APSP. Basch, Khanna, and Motwani [BKM95] presented a combinatorial algorithm that verifies whether a graph has diameter 2 in time $O(n^3/\log^2 n)$. A slight adaptation of this algorithm yields a boolean matrix multiplication algorithm which runs in the same time bound, thereby allowing us to verify that the diameter of a graph is d , for any constant d , in $O(n^3/\log^2 n)$ time.

Consider the problem of devising a fast algorithm for *approximating* the diameter. It is easy to estimate the diameter within a ratio $1/2$ in $O(m)$ time: perform a breadth-first search (BFS) from any vertex v and let d be the depth of the BFS tree obtained; clearly, the diameter of G lies between d and $2d$. No better approximation algorithm was known for this problem; in fact, it was not even known how to distinguish between graphs of diameter 2 and 4. Our first result is an $O(m\sqrt{n\log n})$ algorithm for distinguishing between graphs of diameter 2 and 4, and this is later extended to obtaining a ratio $2/3$ approximation to the diameter in time $O(m\sqrt{n\log n} + n^2 \log n)$. It should be noted that, unlike in the case of all-pairs shortest paths, our results for approximate diameter computation can be extended to the case of *directed* graphs with *arbitrary* positive real *weights* on the edges.

The problem of computing approximate shortest paths has been considered earlier in the literature, but purely from the point of view of multiplicative errors in the approximation. Awerbuch, Berger, Cowen, and Peleg [ABCP93] and Cohen [Coh93] have presented efficient algorithms for computing t -stretch paths for $t \geq 4$, where a path is said to have stretch t if its length is at most t times the length of the shortest path between its end-points. Cohen [Coh94] gave an algorithm that approximates paths from s sources to all other nodes in a weighted graph in time $O((m+sn)n^\epsilon)$ for any $\epsilon > 0$. This algorithm outputs paths of length $(1 + O(1/\text{polylog } n))d(u, v) + O(w_{max}\text{polylog } n)$, where $d(u, v)$ denotes the distance between vertices u and v , and w_{max} is the largest edge weight in the graph. Her algorithm may be specialized to the unweighted case to compute paths of length $(1 + \delta)d(u, v) + c$ for any $\delta > 0$ within the same time, where the constant c depends on ϵ and δ .

The rest of this paper is organized as follows. We begin in Section 2 by presenting some definitions and an algorithm for a version of the dominating set problem that underlies all our algorithms. In Section 3, we describe the algorithms for distinguishing between graphs of diameter 2 and 4, and the extension to obtaining a ratio $2/3$ approximation to the diameter. As we remarked earlier, these results can be applied to directed, weighted graphs. Then, in Section 4, we apply the ideas developed in estimating the diameter to obtain the promised algorithm for an additive-error approximation for APSP. These ideas are extended in Section 5 to obtain a more efficient algorithm for additive-error approximations to the k -pairs shortest paths problem. Finally, in Section 6 we present an empirical study of the performance of our algorithm for APSP.

2 Preliminaries and a Basic Algorithm

We present some notation and a result concerning dominating sets in graphs that underlies all our algorithms. Initially, all definitions are with respect to some undirected, unweighted, connected graph $G(V, E)$ with n vertices and m edges. Later, we will point out the extension to directed and weighted graphs.

Definition 1 *The distance $d(u, v)$ between two vertices u and v is the length of the shortest path between them.*

Definition 2 *The diameter Δ of a graph G is defined to be $\max_{u, v \in G} d(u, v)$.*

Definition 3 *The k -neighborhood $N_k(v)$ of a vertex v is the set of all vertices other than v that are at distance at most k from v , i.e.,*

$$N_k(v) = \{u \in V \mid 0 < d(v, u) \leq k\}.$$

The degree of a vertex v is denoted by $d_v = |N_1(v)|$. Finally, we will use the notation $N(v) = N_1(v) \cup \{v\}$ to denote the set of vertices at distance at most 1 from v .

It is important to keep in mind that the set $N(v)$ contains not just the neighbors of v , but also includes v itself.

Definition 4 *For any vertex $v \in V$, we denote by $B(v)$ the depth of a breadth-first search (BFS) tree in G rooted at the vertex v .*

Throughout this paper, we will working with a parameter s to be chosen later that will serve as the threshold for classifying vertices as being of *low* degree or *high* degree. This threshold is implicit in the following definition.

Definition 5 *Let $L(V) = \{u \in V \mid d_u < s\}$ and $H(V) = V \setminus L(V) = \{u \in V \mid d_u \geq s\}$.*

The following is a generalization of the standard notion of a dominating set.

Definition 6 *Given a set $A \subseteq V$, a set $D \subseteq V$ is a dominating set for A if and only if for each vertex $v \in A$, $N(v) \cap D \neq \emptyset$. That is, for each vertex in $A \setminus D$, one of its neighbors is in D .*

The following theorem underlies all our algorithms.

Theorem 1 *There exists a dominating set for $H(V)$ of size $O(s^{-1}n \log n)$ and such a dominating set can be found in $O(m + ns)$ time.*

Remark 1 *It is easy to see that choosing a set of $\Theta(s^{-1}n \log n)$ vertices uniformly at random gives the desired dominating set for $H(V)$ with high probability. This construction in the proof of this theorem is in effect a derandomization of this randomized algorithm.*

Proof: Suppose, to begin with, that $H(V) = V$; then, we are interested in the standard dominating set for the graph G . The problem of computing a minimum dominating set for G can be reformulated as a set cover problem, as follows: for every vertex v create a set $S_v = N(v)$. This gives an instance of the set cover problem $\mathcal{S} = \{S_v \mid v \in V\}$, where the goal is to find a minimum cardinality collection of sets whose union is V . Given any set cover solution $\mathcal{C} \subseteq \mathcal{S}$, the set of vertices corresponding to the subsets in \mathcal{C} forms a dominating set for G of the same size as \mathcal{C} . This is because each vertex v occurs in one of the sets $S_w \in \mathcal{C}$, and thus is either in the dominating set itself or has a neighbor therein. Similarly, any dominating set for G corresponds to a set cover for \mathcal{S} of the same cardinality.

The greedy set cover algorithm repeatedly chooses the set that covers the most uncovered elements, and it is known to provide a set cover of size within a factor $\log n$ of the *optimal fractional solution* [Joh74, Lov75]. Since every vertex has degree at least s and therefore the corresponding set S_v has cardinality at least s , assigning a weight of $1/s$ to every set in \mathcal{S} gives a fractional set cover of total weight (fractional size) equal to $s^{-1}n$. Thus, the optimal *fractional* set cover size is $O(n/s)$, and the greedy set cover algorithm must then deliver a solution of size $O(s^{-1}n \log n)$. This gives a dominating set for G of the same size. If we implement the greedy set cover algorithm by keeping the sets in buckets sorted by the number of uncovered vertices, the algorithm can be shown to run in time $O(m)$.

Consider now the case where $H(V) \neq V$. Construct a graph $G' = (V', E')$, adding a set of dummy vertices $X = \{x_i \mid 1 \leq i \leq s\}$, as follows: define $V' = V \cup X$ and $E' = E \cup \{(x_i, x_j) \mid 1 \leq i < j \leq s\} \cup \{(u, x_i) \mid u \in L(V)\}$. Every vertex in this new graph has degree s or higher, so by the preceding argument we can construct a dominating set for G' of size $O(s^{-1}(n+s) \log(n+s)) = O(s^{-1}n \log n)$. Since none of the new vertices in X are connected to the vertices in $H(V)$, the restriction of this dominating set to V will give a dominating set for $H(V)$ of size $O(s^{-1}n \log n)$. Finally, the running time is increased by the addition of the new vertices and edges, but since the total number of edges added is at most $ns + s^2 = O(ns)$, we get the desired time bound. \square

2.1 Extension to Directed Graphs

We briefly indicate the extension of the preceding definitions, notation, and observations to directed graphs. Given a directed graph $G(V, E)$, we will denote by \overleftarrow{G} the graph obtained from G by *reversing* the direction of all the edges of G . We assume the directed graph G is strongly connected.

We use \rightarrow and \leftarrow to overline quantities defined with respect to G and \overleftarrow{G} respectively. We will use the term *degree* to refer to the *out-degree* of a vertex, and for $v \in V$ we will denote its degree by \overrightarrow{d}_v . The definitions of distance, diameter, neighborhoods, BFS-tree and dominating set given earlier extend naturally to directed graphs as described below. We give the definitions only for G , and definitions for \overleftarrow{G} can be obtained similarly.

Definition 7 For any two vertices $u, v \in V$, we define $d(u, v)$ as the length of the shortest path from u to v . If such a path does not exist, then we assume $d(u, v) = \infty$.

Note that $d(u, v)$ is not symmetric in general.

Definition 8 The diameter Δ of a graph G is defined to be $\max_{u, v \in G} d(u, v)$.

Definition 9 Let $\vec{N}_k(v) = \{u \in V \mid 0 < d(v, u) \leq k\}$. Further, $\vec{N}(v) = \vec{N}_1(v) \cup \{v\}$ denotes the set of vertices at distance at most 1 from v .

Definition 10 \vec{BFS} is a BFS tree in G using only out-going edges. For any vertex $v \in V$, we denote by $\vec{B}(v)$ the depth of a \vec{BFS} tree in G rooted at the vertex v .

We define $\vec{H}(V)$, $\vec{L}(V)$, and dominating set for directed graphs with respect to out-going edges incident at the vertices.

Definition 11 For some s , let $\vec{L}(V) = \{u \in V \mid \vec{d}_u < s\}$ and $\vec{H}(V) = V \setminus \vec{L}(V) = \{u \in V \mid \vec{d}_u \geq s\}$.

Definition 12 Given a set $A \subseteq V$, a set $D \subseteq V$ is an out-dominating set for A if and only if for each vertex $v \in A$, $\vec{N}(v) \cap D \neq \emptyset$.

The following is an easy consequence of Theorem 1.

Corollary 1 Given a directed graph $G(V, E)$, there exists an out-dominating set for $\vec{H}(V)$ of size $O(s^{-1} \log n)$ and such a dominating set can be found in $O(m + ns)$ time.

3 Estimating the Diameter

In this section we will develop an algorithm to find an estimator E such that $\frac{2}{3}\Delta \leq E \leq \Delta$. We first present an algorithm for distinguishing between graphs of diameter 2 and 4. It is then shown that this algorithm generalizes to the promised approximation algorithm.

3.1 Distinguishing Diameter 2 from 4

The basic idea behind the algorithm is rooted in the following lemma whose proof is straightforward.

Lemma 1 Suppose that G has a pair of vertices a and b with $d(a, b) \geq 4$. Then, any \vec{BFS} tree rooted at a vertex $v \in \vec{N}(a)$ and any \overleftarrow{BFS} tree rooted at a vertex $v \in \overleftarrow{N}(b)$ will have depth at least 3.

The algorithm shown in Figure 1, called Algorithm 2-vs-4, computes BFS trees from a small set of vertices that is guaranteed to contain such a vertex, and so one of these BFS trees will certify that the diameter is more than 2.

We are assuming here that the sets $\vec{L}(V)$ and $D(V)$ are provided as a part of the input; otherwise, they can be computed in $O(m + ns)$ time.

Algorithm 2-vs-4

1. **if** $\vec{L}(V) \neq \emptyset$ **then**
 - (a) **choose** $v \in \vec{L}(V)$
 - (b) **compute** a \vec{BFS} tree from v and a \overleftarrow{BFS} tree from each of the vertices in $\vec{N}(v)$
2. **else**
 - (a) **compute** an out-dominating set D for $\vec{H}(V) = V$
 - (b) **compute** a \vec{BFS} tree from each of the vertices in D
3. **endif**
4. **if** all BFS trees have depth 2 **then return** 2
else return 4.

Figure 1: Algorithm 2-vs-4

Theorem 2 *Algorithm 2-vs-4 distinguishes graphs of diameter 2 and 4, and it has running time $O(ms^{-1}n \log n + ms)$.*

Proof: It is clear that the algorithm outputs 2 for graphs of diameter 2 since in such graphs no BFS tree can have depth exceeding 2. Assume then that G has diameter 4 and fix any pair of vertices $a, b \in V$ such that $d(a, b) \geq 4$. We will show that the algorithm performs a properly directed BFS from a vertex $v \in \vec{N}(a) \cup \overleftarrow{N}(b)$. Since, by Lemma 1, the depth of the BFS tree rooted at v is at least 3, the algorithm will output 4.

We consider the two cases that can arise in the algorithm.

Case 1: $[\vec{L}(V) \neq \emptyset]$

If b belongs to $\vec{N}(v)$, then there is nothing to prove. If $\vec{B}(v) > 2$, then again we have nothing to prove. Therefore, the only case that remains is when $\vec{B}(v) = 2$ and $d(v, b) = 2$ (see Figure 2). These assumptions imply that $\vec{N}(v) \cap \overleftarrow{N}(b) \neq \emptyset$, and Lemma 1 completes the proof. The size of $\vec{N}(v)$ is at most s , therefore the time to compute the BFS trees is bounded by $O(ms)$.

Case 2: $[\vec{L}(V) = \emptyset]$

Since D is an out-dominating set for V , it follows immediately that $D \cap \vec{N}(a) \neq \emptyset$, establishing the proof of correctness. From Theorem 1, we have $|D| = O(s^{-1}n \log n)$ and this implies a bound of $O(ms^{-1}n \log n)$ on the cost of computing the BFS trees in this case. \square

Choosing $s = \sqrt{n \log n}$, we obtain the following corollary.

Corollary 2 *Graphs of diameter 2 and 4 can be distinguished in $O(m\sqrt{n \log n})$ time.*

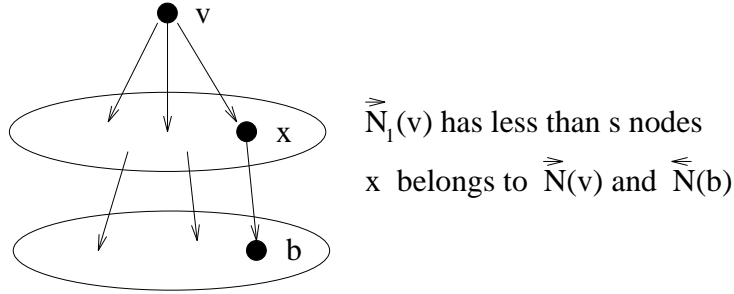


Figure 2: Case 1 in Algorithm 2-vs-4.

3.2 Approximating the Diameter

The ideas used in Algorithm 2-vs-4 can be generalized to estimate the diameter in the directed case: fix any two vertices a and b for which $d(a, b) = \Delta$, where Δ is the diameter of the graph. Suppose we can find a vertex v in $\vec{N}_{\Delta/3}(a)$ or v' in $\overleftarrow{N}_{\Delta/3}(b)$, then it is clear that $\vec{B}(v) \geq \frac{2}{3}\Delta$ or $\overleftarrow{B}(v') \geq \frac{2}{3}\Delta$, and we can use $\vec{B}(v)$ or $\overleftarrow{B}(v')$ as our estimator. As before, we will find a small set of vertices which is guaranteed to have a vertex in $\vec{N}_{\Delta/3}(a) \cup \overleftarrow{N}_{\Delta/3}(b)$. Then, we can compute the BFS tree from each of these vertices and use the maximum of the depths of these trees as our estimator E . The reason for choosing the fraction $1/3$ will become apparent in the analysis of the algorithm. In what follows, it will simplify notation to assume that $\Delta/3$ is an integer; in general though, our analysis needs to be modified to use $\lfloor \Delta/3 \rfloor$. Also, we assume that $\Delta \geq 3$, and it is easy to see that the case $\Delta \leq 2$ is easy to handle separately.

A key tool in the rest of our algorithms will be the notion of a *partial-BFS* defined in terms of a parameter k .

Definition 13 A k -partial-BFS tree is obtained by performing a BFS up to the point where exactly k vertices (not including the root) have been visited.

Lemma 2 A k -partial-BFS tree can be computed in time $O(k^2)$.

Proof: The number of edges examined for each vertex visited is bounded by k since the k -partial-BFS process is terminated when k distinct vertices have been examined. This implies that the total number of edges examined is $O(k^2)$, and that dominates the running time. \square

Note that a k -partial-BFS tree contains the k vertices closest to the root, but that this set is not uniquely defined due to the need to break ties, which is done arbitrarily. Typically, k will be clear from the context and we will not specify it explicitly.

Definition 14 Let $\vec{PBFS}_k(v)$ be the set of vertices visited by a k -partial-BFS from v . Denote by $\vec{PB}(v)$ the depth of the tree constructed in this fashion. Further, \overleftarrow{PBFS} and $\overleftarrow{PB}(v)$ are defined similarly.

Consider now the formal description of the approximation algorithm for diameter, Algorithm Approx-Diameter, as shown in Figure 3.

Algorithm Approx-Diameter

1. **compute** an s -partial- \overrightarrow{BFS} tree from each vertex in V
2. **let** w be the vertex with the maximum depth ($\overrightarrow{PB}(w)$) partial- \overrightarrow{BFS} tree
3. **compute** a \overrightarrow{BFS} tree from w and a \overleftarrow{BFS} tree from each vertex in $P\overrightarrow{BFS}_s(w)$
4. **compute** a new graph \hat{G} from G by adding all edges of the form (v, u) where $u \in P\overrightarrow{BFS}_s(v)$
5. **compute** an out-dominating set D in \hat{G}
6. **compute** a \overrightarrow{BFS} tree from each vertex in D
7. **return** estimator E equal to the maximum depth of all BFS trees from Steps 3 and 6.

Figure 3: Algorithm Approx-Diameter

The following lemmas constitute the analysis of this algorithm.

Lemma 3 *The dominating set D found in Step 5 is of size $O(s^{-1}n \log n)$.*

Proof: In \hat{G} , each vertex $v \in V$ is adjacent to all vertices in $P\overrightarrow{BFS}_s(v)$ with respect to the graph G . Since $|P\overrightarrow{BFS}_s(v)| = s$ for every vertex v , the out-degree of each vertex in \hat{G} is at least s . From Theorem 1, it follows that we can find a dominating set of size $O(s^{-1}n \log n)$. \square

Lemma 4 *If $|\overrightarrow{N}_{\Delta/3}(v)| \geq s$ for all $v \in V$, then $D \cap (\overrightarrow{N}_{\Delta/3}(v) \cup \{v\}) \neq \emptyset$ for each vertex $v \in V$.*

Proof: Consider any particular vertex $v \in V$. If v is in D , then there is nothing to prove. Otherwise, since D is a dominating set in \hat{G} , there is a vertex $u \in D$ such that (v, u) is an edge in \hat{G} . If (v, u) is in G , then again we are done since $u \in \overrightarrow{N}(v) \subset \overrightarrow{N}_{\Delta/3}(v)$. The other possibility is that u is not a neighbor of v in G , but then it must be the case that $u \in P\overrightarrow{BFS}_s(v)$. The condition $|\overrightarrow{N}_{\Delta/3}(v)| \geq s$ implies that $P\overrightarrow{BFS}_s(v) \subset \overrightarrow{N}_{\Delta/3}(v)$, which in turn implies that $u \in \overrightarrow{N}_{\Delta/3}(v)$, and hence $u \in D \cap \overrightarrow{N}_{\Delta/3}(v)$. \square

The reader should notice the similarity between the preceding lemma and Case 2 in Theorem 2. Lemma 4 follows from the more general set cover ideas used in the proof of Theorem 1 and as such it holds even if we replace $\Delta/3$ by some other fraction of Δ . The more crucial lemma is given below.

Lemma 5 *Let S be the set of vertices v such that $|\overrightarrow{N}_{\Delta/3}(v)| < s$. If $S \neq \emptyset$ then the vertex w found in Step 2 belongs to S . In addition if $\overrightarrow{PB}(w) < \frac{2}{3}\Delta$, then for every vertex v ,*

$$P\overrightarrow{BFS}_s(w) \cap \overleftarrow{N}_{\Delta/3}(v) \neq \emptyset.$$

Proof: It can be verified that for any vertex $u \in S$, $\overrightarrow{\text{PB}}(u) > \Delta/3$; conversely, for any vertex v in $V \setminus S$, $\overrightarrow{\text{PB}}(v) \leq \Delta/3$. From this we can conclude that if S is nonempty, then the vertex of largest depth belongs to S .

Also, for each vertex $u \in S$, we must have $\overrightarrow{N}_{\Delta/3}(u) \subset P\overrightarrow{BFS}_s(u)$. If $\overrightarrow{\text{B}}(w) < \frac{2}{3}\Delta$ then every vertex is within a distance $\frac{2}{3}\Delta$ from w . From this and the fact that $\overrightarrow{N}_{\Delta/3}(w) \subset P\overrightarrow{BFS}_s(w)$, it follows that $P\overrightarrow{BFS}_s(w) \cap \overleftarrow{N}_{\Delta/3}(v) \neq \emptyset$. \square

The proof of the above lemma makes clear the reason why our estimate is only within $\frac{2}{3}$ of the diameter. Essentially, we need to ensure that the Δ/k neighborhood of w intersects the Δ/k neighborhood of every other vertex. This can happen only if $\overrightarrow{\text{B}}(w)$ is sufficiently small. If it is not small enough, we want $\overrightarrow{\text{B}}(w)$ itself to be a good estimator. Balancing these conditions gives us $k = 3$ and the ratio $2/3$.

Theorem 3 *Algorithm Approx-Diameter gives an estimate E such that $\frac{2}{3}\Delta \leq E \leq \Delta$ in time $O(ms + ms^{-1}n \log n + ns^2)$. Choosing $s = \sqrt{n \log n}$ gives a running time of $O(m\sqrt{n \log n} + n^2 \log n)$.*

Proof: The analysis is partitioned into two cases. Let a and b be two vertices such that $d(a, b) = \Delta$.

Case 1: [For all vertices v , $|\overrightarrow{N}_{\Delta/3}(v)| \geq s$.]

If either a or b is in D , we are done. Otherwise from the proof of Lemma 4, the set D has a vertex $v \in \overrightarrow{N}_{\Delta/3}(a)$. Since in Step 6 we compute \overrightarrow{BFS} trees from each vertex in D , one of these is v and $\overrightarrow{\text{B}}(v)$ is the desired estimator.

Case 2: [There exists a vertex $v \in V$ such that $|\overrightarrow{N}_{\Delta/3}(v)| < s$.]

Let w be the vertex in Step 2. If $\overrightarrow{\text{B}}(w) \geq \frac{2}{3}\Delta$, $\overrightarrow{\text{B}}(w)$ is our estimator and we are done. Otherwise from Lemma 5, $P\overrightarrow{BFS}_s(w)$ has a vertex $v \in \overleftarrow{N}_{\Delta/3}(b)$. Since in Step 3 we compute \overleftarrow{BFS} trees from each vertex in $P\overrightarrow{BFS}_s(w)$, one of these is v and $\overleftarrow{\text{B}}(v)$ is the desired estimator.

The running time is easy to analyze. Each partial-BFS in Step 1 takes at most $O(s^2)$ time by Lemma 2; thus, the total time spent on Step 1 is $O(ns^2)$. Step 2 can be implemented in $O(n)$ time. In Step 3, we compute BFS trees from s vertices, which requires a total of $O(ms)$ time. The time required in Step 4 is dominated by the time required to compute the partial-BFS trees in Step 1. Theorem 1 implies that Step 5 requires only $O(n^2 + ns)$ time (note that the graph \widehat{G} could have many more edges than m). By Lemma 3, Step 6 takes $O(ms^{-1}n \log n)$ time. Finally, the cost of Step 7 is dominated by the cost of computing the various BFS trees in Steps 3 and 6. The running time is dominated by the cost of Steps 1, 3, and 6, and adding the bounds for these gives the desired result. \square

3.3 Extension to Weighted Graphs

The algorithm for estimating the diameter extends to the case of weighted graphs as well, provided all edge weights are positive. This requires some minor modifications to Algorithm Approx-Diameter that are listed below.

- The BFS is replaced by Dijkstra’s algorithm [CLR90] for shortest paths, and the depth of the tree now refers to the distance to the farthest vertex found so far.
- In forming the new graph \hat{G} in Step 4 we need to remove all the original edges of G before we add the new edges. Note that \hat{G} is an unweighted graph.

The last modification is necessary because in a weighted graph it is not necessarily the case that a neighbor of a vertex v belongs to $N_{\Delta/3}(v)$. The running time remains the same because the time required by Dijkstra’s algorithm (implemented with Fibonacci heaps [CLR90]) is $O(m)$ when $m = \Omega(n \log n)$.

We obtain the following theorem.

Theorem 4 *Given a directed graph with positive edge weights, there is an algorithm that gives an estimate E such that $\frac{2}{3}\Delta \leq E \leq \Delta$ in time $O(m\sqrt{n \log n} + n^2 \log n)$.*

4 Estimating All-Pairs Shortest Paths

We now turn to the problem of approximate APSP computations. We restrict ourselves to undirected and unweighted graphs for the rest of the paper, although it should be noted that there is an obvious extension of the results below to the case of undirected graphs with edge weights that are small integers.

It is possible to determine not only the diameter, but the all-pairs shortest path distances to within an additive error of 2. The basic idea is that a dominating set, since it contains a neighbor of every vertex in the graph, must contain a vertex that is within distance 1 of any shortest path. Since we can only find a small dominating set for vertices in $H(V)$, we have to treat $L(V)$ vertices differently, but their low degree allows us to manage with only a partial-BFS, which we can combine with the information we have gleaned from the dominating set.

We give a detailed description of the approximate APSP algorithm, Algorithm Approx-APSP, in Figure 4. In Figure 5 we illustrate the main ideas behind this algorithm.

Theorem 5 *In Algorithm Approx-APSP, for all vertices $u, v \in V$, the distances returned in \hat{d} satisfy the inequality*

$$0 \leq \hat{d}(u, v) - d(u, v) \leq 2.$$

Further, the algorithm can be modified to produce paths of length \hat{d} rather than merely returning the approximate distances. This algorithm runs in time $O(n^2 s + n^3 s^{-1} \log n)$; choosing $s = \sqrt{n \log n}$ gives a running time of $O(n^{2.5} \sqrt{\log n})$.

Proof: We first show that the algorithm can be easily modified to return actual paths rather than only the distances. To achieve this, in Steps 3 and 4 we can associate with each updated entry in the matrix the path from the BFS tree used for the update. In Step 5, we merely concatenate the two paths from Step 3 the sum of whose lengths determine the minimum value of \hat{d} .

Algorithm Approx-APSP

Comment: Define $G[L(V)]$ to be the subgraph of G induced by $L(V)$.

1. **initialize** all entries in the distance matrix \hat{d} to ∞
2. **compute** a dominating set D for $H(V)$ of size $s^{-1}n \log n$
3. **compute** a BFS tree from each vertex $v \in D$, and update \hat{d} with the shortest path lengths for v so obtained
4. **compute** a BFS tree in $G[L(V)]$ for each vertex $v \in L(V)$, and update \hat{d} with the shortest path lengths for v so obtained
5. **for all** $u, v \in V \setminus D$ **do**

$$\hat{d}(u, v) \leftarrow \min\{\hat{d}(u, v), \min_{w \in D}\{\hat{d}(w, u) + \hat{d}(w, v)\}\}$$

6. **return** \hat{d} as the APSP matrix, and its largest entry as the diameter.

Figure 4: Algorithm Approx-APSP

For a vertex u , it is clear that the shortest path distance to any vertex $v \in V$ that is returned cannot be smaller than the correct values, since they correspond to actual paths. To see that they differ by no more than 2, we need to consider three cases:

Case 1: [$u \in D$]

In this case, the BFS tree from v is computed in Step 3 and so clearly the distances returned are correct.

Case 2: [$u \in H(V) \setminus D$]

By the definition of D , it must be the case that u has a neighbor w in D . Clearly, the distances from u and w to any other vertex cannot differ by more than 1, and the distances from w are always correct as per Case 1. The assignment in Step 6 guarantees $\hat{d}(u, v) \leq \hat{d}(w, v) + \hat{d}(w, u) = d(w, v) + 1 \leq d(u, v) + 2$.

Case 3: [$u \in L(V)$]

Fix any shortest path from u to v . Suppose that the path from u to v is entirely contained in $L(V)$; then, $\hat{d}(u, v)$ is set correctly in Step 4. Otherwise, the path must contain a vertex $w \in H(V)$. If w is contained in D , then the correct distance is computed as per Case 1. Finally, if $w \in H(V) \setminus D$, then D contains a neighbor x of w . Clearly, in Step 6, one of the possibilities considered will involve a path from u to x and a path from x to v . Since the distances involving x are correctly computed in Step 3, this means that $\hat{d}(u, v) \leq d(x, u) + d(x, v) \leq d(w, u) + d(w, v) + 2 = d(u, v) + 2$.

Finally, we analyze the running time of this algorithm. Step 1 requires only $O(n^2)$ time, and

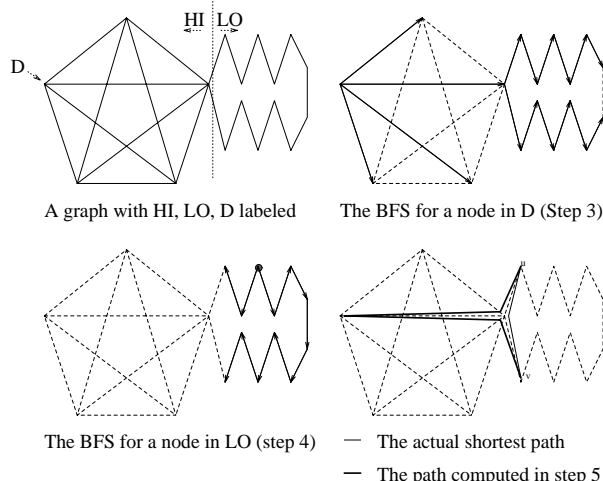


Figure 5: Illustration of Algorithm Approx-APSP.

Theorem 1 implies that we can perform Step 2 in the stated time bound. Step 3 requires $ms^{-1}n \log n$ for computing the BFS trees. Step 4 may compute as many as $\Omega(n)$ BFS trees, but $G[L(V)]$ only has $O(ns)$ edges and so this requires only $O(n^2s)$ time. Finally, Step 5 takes all n^2 vertex pairs, and compares them against the $s^{-1}n \log n$ vertices in D . This implies the desired time bound. \square

Although the error in this algorithm is 2, it can be improved for the special case of distinguishing diameter 2 from 4 based on the following two observations.

Fact 1 *If $u \in H(V)$ is at distance Δ from some vertex v , then $\widehat{d}(u, v) \leq \Delta + 1$.*

Proof: Consider w , the vertex that dominates u . If the algorithm were to have set $\widehat{d}(u, v) > \Delta + 1$ then Step 5 of the algorithm would imply $\widehat{d}(w, v) > \Delta$. Since \widehat{d} is exact for vertices in D , this is not possible. \square

Fact 2 *If the algorithm reports for some $u \in L(V)$ that $B(u) > 2$, we can verify this in time $O(ns)$ per vertex.*

Thus, by performing a verification for each of the $L(V)$ vertices that report distance over 2, we can improve Algorithm Approx-APSP so that it always performs as well as the diameter approximation algorithms of the previous section. The first fact also appears to be useful in bringing the diameter error down to 1, but unfortunately, the vertices in $L(V)$ cannot be handled as easily for larger diameters.

5 Estimating k -Pairs Shortest Paths

In this section we consider the problem where we only seek to determine the distances between a given set of distinguished pairs of vertices denoted by P . We show that the Algorithm Approx-APSP

can be generalized to handle this problem with the same error bounds. The generalized algorithm, called Approx- k PSP, works in time $O(n^{1.5}\sqrt{k \log n} + n^2 \log^2 n)$, where k is the cardinality of the set P . When P contains all pairs of vertices, the behavior of Approx- k PSP is identical to that of Approx-APSP. However, for small k , the algorithm is significantly faster than Approx-APSP.

The main idea behind the speed-up is the observation that the choice of $s = \sqrt{n \log n}$ is not optimal when we do not need to find the distances between all pairs. Step 5 of Approx-APSP (the concatenation of paths) requires at most $O(kns^{-1} \log n)$ time (instead of $O(n^3s^{-1} \log n)$), so the total time taken is $O(mns^{-1} \log n + n^2s + kns^{-1} \log n)$. Now, the first term of the sum is not necessarily dominated by the last one. We have two cases: if $k \geq m$, the last term dominates the first but we get the desired running time that depends on k when we balance the second and third terms; conversely, when $k < m$, we observe that the first term dominates the last. Note that the first term is the cost of performing BFS from all the dominating set vertices.

The intuition for the improvement comes from the following example: suppose that all vertices have degree d . Then we could take $s = d$, so $m = O(nd)$ and $|D| = O(\frac{n}{d} \log n)$, and the first term would be equal to $m|D| = O(n^2 \log n)$. This example shows that performing BFS from all the dominating set vertices is not expensive if the degrees are more or less uniform. Of course, in general, such an assumption is not true. However, we can exploit this observation by partitioning the vertex set into $O(\log n)$ classes, such that the i th class consists of vertices of degree between $\frac{n}{2^i}$ and $\frac{n}{2^{i-1}}$, and computing the dominating sets for each class separately. This effectively reduces the the first term to $O(n^2 \text{polylog}(n))$, and now we can balance the second and third terms as in the other case.

This algorithm, called Approx- k PSP, is described in Figure 6. The algorithm is recursive and at the top level it is invoked with parameter value $i = 1$, assuming $G_1 = G$. The algorithm uses a parameter t that is chosen later so as to minimize the running time.

We begin the analysis by identifying the optimal choice of t .

Lemma 6 *The parameter t can be chosen such that the running time of Algorithm Approx- k PSP is $O(n^{1.5}\sqrt{k \log n} + n^2 \log^2 n)$.*

Proof: Observe that for each i and $v \in V_i$, the degree of v in G_i is less than s_{i-1} (assume $s_0 = n$). This implies that $|E_i| = O(ns_{i-1})$. Let C_i denote the running time of the invocation of Approx- k PSP with argument i . It is easy to see that $C_t = n|E_t|$ and that for $i < t$, C_i is dominated by the time required for Steps 3(c) and 3(d) which require $|D_i||E_i|$ and $k|D_i|$ time, respectively. The total time can now be estimated as follows:

$$\begin{aligned} \sum_{i=1}^t C_i &= n|E_t| + \sum_{i=1}^{t-1} (|D_i||E_i| + k|D_i|) \leq \sum_{i=1}^{t-1} \frac{n}{s_i} ns_{i-1} \log n + k \sum_{i=1}^{t-1} \frac{n}{s_i} \log n + n^2 s_t \\ &\leq 2n^2 \log^2 n + \frac{n}{s_t} k \log n + n^2 s_t \end{aligned}$$

Choosing t such that $s_t = \Theta\left(\sqrt{\frac{k \log n}{n}}\right)$ gives the desired bound on the running time of the algorithm. \square

Algorithm Approx- k PSP(i)

Comment: Define $s_i = \frac{n}{2^i}$

1. **initialize** all $\hat{d}(u, v)$ to ∞

2. **if** $i = t$ **then** compute a BFS tree from each vertex $v \in V_i$ in G_i and update \hat{d} with the shortest path lengths obtained

3. **else let**

$U_i = \{v \in V_i \mid \text{degree of } v \text{ in } G_i \text{ is at least } s_i\}$,

$V_{i+1} = V_i - U_i$, and

G_{i+1} be the subgraph of G_i induced by V_{i+1}

(a) **call** Approx- k PSP($i+1$)

(b) **compute** a dominating set D_i for U_i in G_i

(c) **compute** a BFS tree from each $v \in D_i$ and update \hat{d} with the shortest path lengths obtained

(d) **for each** $\{u, v\} \in P$ **do**

$$\hat{d}(u, v) \leftarrow \min\{\hat{d}(u, v), \min_{w \in D_i} \hat{d}(w, u) + \hat{d}(w, v)\}$$

4. **return**

Figure 6: Algorithm Approx- k PSP

We can now complete the analysis of the algorithm.

Theorem 6 *For all pairs $(u, v) \in P$, the distances returned in \hat{d} by Approx- k PSP satisfy the inequalities:*

$$0 \leq \hat{d}(u, v) - d(u, v) \leq 2.$$

The algorithm runs in time $O(n^2 \log^2 n + n^{1.5} \sqrt{k \log n})$.

Proof: Let $d_i(u, v)$ denote the distance between u and v in G_i . Clearly, it is sufficient to show by induction on i that $0 \leq \hat{d}(u, v) - d_i(u, v) \leq 2$ after finishing Approx- k PSP(i). The base case (for $i = t$) holds trivially since we compute the exact shortest paths. The proof of the inductive step is similar to the proof Theorem 5, hence we omit the details. The time bound follows from Lemma 6. \square

5.1 Application: Randomized Approximation Scheme for Diameter

Algorithm Approx- k PSP can be used to obtain a randomized approximation scheme for the diameter of a graph. Let $u, v \in V$ be such that $d(u, v) = \Delta$. If we choose a vertex w uniformly at random

	Approx-APSP speedup	Fast Approx-APSP speedup	Approx-APSP accuracy	Fast Approx-APSP accuracy
GB Median	0.59	3.95	0.69	0.53
GB Average	2.44	10.18	0.72	0.47
GB Std Dev	0.24	1.73	0.16	0.13
RG Median	0.52	5.30	0.39	0.51
RG Average	0.63	4.75	0.39	0.55
RG Std Dev	0.23	1.70	0.14	0.12

Table 1: Summary of Experimental Results

from V , the probability of $d(u, w) \leq \frac{\epsilon}{2}\Delta$ is $\Theta(\frac{\epsilon\Delta}{n})$. This guarantees that a set P of $O(\frac{n^2}{\epsilon^2\Delta^2} \log n)$ vertices chosen uniformly at random contains vertices x, y such that $d(u, x) \leq \frac{\epsilon}{2}\Delta$ and $d(v, y) \leq \frac{\epsilon}{2}\Delta$, hence $d(x, y) \geq (1-\epsilon)\Delta$ with high probability. We can use Algorithm Approx- k PSP to approximate distances between all pairs of vertices in P in $O(n^{1.5}\sqrt{|P|}\log n + n^2 \log^2 n) = O(\frac{n^{2.5}}{\epsilon\Delta} \log n + n^2 \log^2 n)$ time. For large Δ , say $\Delta = \Omega(n^\delta)$ for some $\delta > 0$, the improvement is significant. We obtain the following theorem.

Theorem 7 *For any $0 < \epsilon < 1$ there exists a Monte Carlo algorithm which finds an estimator E such that $(1-\epsilon)\Delta \leq E \leq \Delta + 2$, in time $O(\frac{n^{2.5}}{\epsilon\Delta} \log n + n^2 \log^2 n)$.*

Note that while this randomized algorithm assumes knowledge of Δ , it is sufficient to provide it with a constant-factor approximation to Δ . The depth of a BFS tree rooted at an arbitrary vertex of G is a 2-approximation for Δ and can be used for this purpose.

6 Experimental Results

To evaluate the usefulness of our algorithm, we ran it on two families of graphs and compared the results against a carefully coded algorithm based on breadth-first searches. The algorithm Approx-APSP was tweaked with the following heuristic improvement to Step 5 that avoids many needless iterations: when a node has a neighbor in D , then we copy the distances of its neighbor (since they can differ by at most 1). This algorithm (called Fast Approx-APSP) occasionally has a higher fraction of incorrect entries, but seems to be the fastest way to solve the all-pairs shortest path problem.

The first family of graphs were random graphs from the $G_{n,m}$ model [Bol85], which are graphs chosen uniformly at random from those with n vertices and m edges. In our experiments, we chose random graphs with n ranging from 10 to 1000, and $2m/n^2$ ranging from 0.03 to 0.90. On these graphs, Fast Approx-APSP runs about 5 times faster than the BFS implementation, and about half of the distances are off by one.

The second family of graphs come from the Stanford GraphBase [Knu93]. We tested all of the connected, undirected graphs from Appendix C in Knuth [Knu93] (ignoring edge weights). This is a very heterogeneous family of graphs, including graphs representing highway connections for American cities, athletic schedules, 5-letter English words, and expander graphs, as well as more combinatorial graphs. Thus the results here are quite indicative of practical performance. Although the BFS-based algorithm runs fastest for certain subfamilies of the GraphBase, Fast Approx-APSP outperformed all other algorithms overall.

The results are summarized in Table 1. In the table, GB and RG refer to GraphBase and random graphs, respectively. The speedup numbers indicate the inverse of the ratio of the execution time of the algorithms to that of the carefully coded BFS algorithm. The accuracy refers to the ratio of the total number of *exact* entries in the distance matrix to the total number of entries in the matrix. In both of these families, the accuracy of Approx-APSP could be improved by subtracting 1 in Step 5. This did not seem necessary given that the BFS approach performed about as fast as Approx-APSP, and that Fast- Approx APSP performed faster with roughly 50% accuracy. The numbers indicate that for general graphs where an additive factor error is acceptable, Fast Approx-APSP is the algorithm of choice, and for more specific families of graphs, the parameters can be adjusted for even better performance.

7 Conclusions and Further Work

Our work suggests several interesting directions for future work, the most elementary being: Is there a combinatorial algorithm running in time $O(n^{3-\epsilon})$ for distinguishing between graphs of diameter 2 and 3? It is our belief that the problem of efficiently computing the diameter can be solved given such a decision algorithm, and our work provides some evidence in support of this belief. In fact, it is our view that the bottleneck in obtaining a faster combinatorial APSP algorithm is precisely the problem of distinguishing graphs of diameter 2 and 3. This also raises the question of whether there is some strong equivalence between the diameter and APSP problems, e.g., that their complexity is the same within poly-logarithmic factors. Finally, of course, removing the additive error from our results remains a major open problem.

Acknowledgements

We are grateful to Noga Alon for his comments and suggestions, and to Nati Linial for helpful discussions. We are also indebted to Edith Cohen for comments that helped us extend some of our results. Thanks also to Michael Goldwasser, David Karger, Sanjeev Khanna, and Eric Torng for their comments.

References

- [ABCP93] B. AWERBUCH, B. BERGER, L. COWEN, AND D. PELEG. *Near-linear cost sequential and distributed constructions of sparse neighborhood covers*. in Proceedings of the 34th Annual IEEE Symposium on Foundations of Computer Science, pages 638-647, 1993.
- [AGM91] N. ALON, Z. GALIL, AND O. MARGALIT. *On the exponent of the all pairs shortest path problem*. in Proceedings of the 32nd Annual IEEE Symposium on Foundations of Computer Science, pages 569–575, 1991.
- [AGMN92] N. ALON, Z. GALIL, O. MARGALIT, AND M. NAOR. *Witnesses for Boolean Matrix Multiplication and for Shortest Paths*. in Proceedings of the 33rd Annual IEEE Symposium on Foundations of Computer Science, pages 417–426, 1992.
- [BKM95] J. BASCH, S. KHANNA, AND R. MOTWANI. *On Diameter Verification and Boolean Matrix Multiplication*. Report No. STAN-CS-95-1544, Department of Computer Science, Stanford University (1995).
- [Bol85] B. BOLLOBÁS. *Random Graphs*. Academic Press, 1985.
- [Chu87] FAN R.K. CHUNG. *Diameters of Graphs: Old Problems and New Results*. *Congressus Numerantium*, 60:295–317, 1987.
- [Coh93] E. COHEN. *Fast algorithms for t -spanners and stretch- t paths*. in Proceedings of the 34th Annual IEEE Symposium on Foundations of Computer Science, pages 648-658, 1993.
- [Coh94] E. COHEN. *Polylog-time and near-linear work approximation scheme for undirected shortest paths*. in Proceedings of 26th Annual ACM Symposium on Theory of Computing, pages 16–26, 1994.
- [CW90] D. COPPERSMITH AND S. WINOGRAD. *Matrix multiplication via arithmetic progressions*. *Journal of Symbolic Computation*, 9:251–280, 1990.
- [CLR90] T. CORMEN, C.E. LEISERSON, AND R.L. RIVEST. *Introduction to Algorithms*. MIT Press and Mc-Graw Hill, New York, 1990.
- [FM91] T. FEDER AND R. MOTWANI. *Clique Partitions, Graph Compression and Speeding-up Algorithms*. *Journal of Computer and System Sciences*, 51:261–272, 1995.
- [Joh74] D.S. JOHNSON. *Approximation algorithms for combinatorial problems*. *Journal of Computer and System Sciences*, 9:256–278, 1974.
- [Knu93] D.E. KNUTH. *The Stanford GraphBase: A platform for combinatorial computing*. Addison-Wesley, 1993.

- [Lov75] L. LOVÁSZ. *On the ratio of optimal integral and fractional covers*. Discrete Mathematics, 13:383–390, 1975
- [Sei92] R. G. SEIDEL. *On the all-pairs-shortest-path problem*. in Proceedings of the 24th Annual ACM Symposium on Theory of Computing, pages 745–749, 1992.
- [Str69] V. STRASSEN. *Gaussian elimination is not optimal*. Numerische Mathematik, 13:354–356, 1969.