



University of Pennsylvania  
**ScholarlyCommons**

---

Departmental Papers (CIS)

Department of Computer & Information Science

---

July 2004

# Approximating Longest Directed Paths and Cycles

Andreas Björklund  
*Lund University*

Thore Husfeldt  
*Lund University*

Sanjeev Khanna  
*University of Pennsylvania, [sanjeev@cis.upenn.edu](mailto:sanjeev@cis.upenn.edu)*

Follow this and additional works at: [http://repository.upenn.edu/cis\\_papers](http://repository.upenn.edu/cis_papers)

---

## Recommended Citation

Andreas Björklund, Thore Husfeldt, and Sanjeev Khanna, "Approximating Longest Directed Paths and Cycles", . July 2004.

Postprint version. Published in *Lecture Notes in Computer Science*, Volume 3142, Automata, Languages and Programming, (ICALP 2004), pages 222-233.

Publisher URL: <http://dx.doi.org/10.1007/b99859>

This paper is posted at ScholarlyCommons. [http://repository.upenn.edu/cis\\_papers/205](http://repository.upenn.edu/cis_papers/205)  
For more information, please contact [libraryrepository@pobox.upenn.edu](mailto:libraryrepository@pobox.upenn.edu).

---

# Approximating Longest Directed Paths and Cycles

## Abstract

We investigate the hardness of approximating the longest path and the longest cycle in directed graphs on  $n$  vertices. We show that neither of these two problems can be polynomial time approximated within  $n^{1-\varepsilon}$  for any  $\varepsilon > 0$  unless  $P = NP$ . In particular, the result holds for digraphs of constant bounded outdegree that contain a Hamiltonian cycle.

Assuming the stronger complexity conjecture that Satisfiability cannot be solved in subexponential time, we show that there is no polynomial time algorithm that finds a directed path of length  $\Omega(f(n) \log^2 n)$ , or a directed cycle of length  $\Omega(f(n) \log n)$ , for any nondecreasing, polynomial time computable function  $f$  in  $\Omega(1)$ . With a recent algorithm for undirected graphs by Gabow, this shows that long paths and cycles are harder to find in directed graphs than in undirected graphs.

We also find a directed path of length  $\Omega(\log^2 n / \log \log n)$  in Hamiltonian digraphs with bounded outdegree. With our hardness results, this shows that long directed cycles are harder to find than a long directed paths. Furthermore, we present a simple polynomial time algorithm that finds paths of length  $\Omega(n)$  in directed expanders of constant bounded outdegree.

## Comments

Postprint version. Published in *Lecture Notes in Computer Science*, Volume 3142, Automata, Languages and Programming, (ICALP 2004), pages 222-233.

Publisher URL: <http://dx.doi.org/10.1007/b99859>

# Approximating Longest Directed Paths and Cycles

Andreas Björklund<sup>1</sup>, Thore Husfeldt<sup>1</sup>, and Sanjeev Khanna<sup>2\*</sup>

<sup>1</sup> Department of Computer Science, Lund University, Box 118, 221 00 Lund, Sweden.  
thore@cs.lu.se

<sup>2</sup> Dept. of CIS, University of Pennsylvania, Philadelphia, PA 19104.  
sanjeev@cis.upenn.edu

**Abstract.** We investigate the hardness of approximating the longest path and the longest cycle in directed graphs on  $n$  vertices. We show that neither of these two problems can be polynomial time approximated within  $n^{1-\epsilon}$  for any  $\epsilon > 0$  unless  $P = NP$ . In particular, the result holds for digraphs of constant bounded outdegree that contain a Hamiltonian cycle.

Assuming the stronger complexity conjecture that Satisfiability cannot be solved in subexponential time, we show that there is no polynomial time algorithm that finds a directed path of length  $\Omega(f(n) \log^2 n)$ , or a directed cycle of length  $\Omega(f(n) \log n)$ , for any nondecreasing, polynomial time computable function  $f$  in  $\omega(1)$ . With a recent algorithm for undirected graphs by Gabow, this shows that long paths and cycles are harder to find in directed graphs than in undirected graphs.

We also find a directed path of length  $\Omega(\log^2 n / \log \log n)$  in Hamiltonian digraphs with bounded outdegree. With our hardness results, this shows that long directed cycles are harder to find than a long directed paths. Furthermore, we present a simple polynomial time algorithm that finds paths of length  $\Omega(n)$  in directed expanders of constant bounded outdegree.

## 1 Introduction

Given an unweighted graph or digraph  $G = (V, A)$  with  $n = |V|$ , the *Longest Path* problem is to find the longest sequence of distinct vertices  $v_1 \cdots v_k$  such that  $v_i v_{i+1} \in A$ . This problem is notorious for the difficulty of understanding its approximation hardness [4]. The present paper establishes a number of upper and lower bounds for the *directed* case.

The best known polynomial time algorithms for directed graphs essentially find such structures of logarithmic length. More precisely, Alon, Yuster, and Zwick find [1] a dipath or dicycle of length exactly  $c \log n$  for any constant  $c$ , provided it exists, and Gabow and Nie [7] find a dicycle of length  $\log n / \log \log n$ , provided it exists (such a cycle may be far longer than logarithmic).

In the present paper we show that this problem is hard to approximate. Specifically, Theorem 1 states that in directed graphs the length of the longest path cannot be polynomial time approximated within an approximation ratio of  $n^{1-\epsilon}$  for any  $\epsilon > 0$  unless  $P = NP$ .

---

\* Supported in part by an Alfred P. Sloan Research Fellowship and by an NSF Career Award CCR-0093117.

We can claim a stronger bound if we make a stronger assumption called the *Exponential Time Hypothesis* (ETH), namely that Satisfiability has no subexponential time algorithms [8]. Our Theorem 2 states that if we could find a dipath of length  $f(n) \log^2 n$  efficiently (for some polynomial time computable and non-decreasing  $f$  in  $\omega(1)$ ), then there would be a deterministic algorithm for 3-Sat with  $s$  variables with running time  $2^{o(s)}$ , violating ETH. This is relevant to the remaining open question in [1]: “Is there a polynomial time algorithm for deciding if a given graph  $G = (V, E)$  contains a path of length, say,  $\log^2 n$ ?” Even though this question remains open, Alon, Yuster, and Zwick’s choice of time bound was not as capricious as their wording may suggest: any stronger algorithm than  $\log^2 n$  for Longest Dipath would be at variance with the Exponential Time Hypothesis.

*Undirected Graphs versus Directed Graphs.* Our hardness results under ETH are of further interest in the light of a very recent result of Gabow [6] for the undirected case, which shows how to find superpolylogarithmic paths and cycles. More precisely, if the graph contains a cycle of length  $l$  through a given vertex  $v$ , then [6] finds a cycle through  $v$  of length at least  $(\log l)^{c \log \log l}$  for some constant  $c > 0$ . (The same bound for Longest Path follows.)

This shows that paths and cycles in directed graphs are harder to find than in undirected graphs, proving (under ETH) a widely held belief.

*Algorithm for Hamiltonian Digraphs.* Our lower bound holds even if the input digraph is known to be Hamiltonian, which addresses the question to what extent knowledge of the presence of a long path helps in the search for one. We complement this by an algorithm in Theorem 3 to efficiently find paths of length  $\Omega(\log^2 n / \log \log n)$  in Hamiltonian digraphs of constant bounded outdegree; this is close to our own  $f(n) \log^2 n$  lower bound. The best previous upper bound [1] was  $O(\log n)$ .

*Longest Path versus Longest Cycle.* For the related longest *cycle* problem, where we also require  $v_k v_1 \in A$ , we essentially show that one cannot efficiently find a cycle of more than logarithmic length. To be precise, Theorem 2 shows that (under ETH) no polynomial time can find a cycle of length  $\geq f(n) \log n$ , for any nondecreasing, polynomial time computable function  $f$  in  $\omega(1)$ . This is no more than a factor  $\log \log n$  short of the best known approximation algorithm: Recently, Gabow and Nie [7] gave a polynomial time algorithm to find a directed cycle of length  $\geq \log n / \log \log n$  if one exists.

Moreover, together with the longest path guarantee  $\Omega(\log^2 n / \log \log n)$  from Theorem 3, the lower bound separates the complexities of the Longest Path and Longest Cycle problems, at least for the directed, bounded outdegree, Hamiltonian case, and assuming ETH.

*Long Paths in Sparse Expanders.* In contrast to our worst-case inapproximability result, it is well known that almost every digraph contains a path of length  $\Omega(n)$ , and that this path is easy to find [3, Chap. 8]. Thus it would be interesting to understand which natural classes of digraphs admit efficient longest path algorithms.

With Theorem 4 we observe that a very simple algorithm that always finds a path of length  $\Omega(n)$  in a bounded-outdegree directed expander. This provides some insight into the structure of digraphs where long paths are hard to find: The hard instances construct in our lower bound proof have bounded outdegree as well, but can be seen to have very bad expansion properties (for any given size there is a vertex subset of that size with constant size separators).

*Related work.* Among the canonical NP-hard problems, the *undirected* version of this problem has been identified as the one that is least understood [4]. However, a number of recent papers have established increasingly good approximation algorithms [14, 2], culminating in the very recent result by Gabow [6] cited above. Even better bounds exist for restricted classes of graphs; for example, a recent result [4] finds cycles of length  $O(l^\alpha)$  ( $\alpha = \log_3 2$ ) in graphs of maximum degree 3.

However, it remains fair to say that in undirected graphs, the approximation hardness of Longest Path remains open. It has been conjectured [10] that the length of a longest path in undirected graphs cannot be approximated within  $n^\alpha$  for some  $\alpha > 0$  unless  $P = NP$ , a somewhat weaker bound than the one we prove for digraphs, but this is far from being proved: the quoted reference shows that the Longest Path is not in APX, and that no polynomial time algorithm can approximate the length of the longest path within  $2^{\log^{1-\epsilon} n}$  for any  $\epsilon > 0$  unless  $NP \subseteq DTIME(2^{\log^{O(1/\epsilon)} n})$ .

Our lower bound uses a reduction to the  $k$  *Vertex Disjoint Paths* problem in digraphs. Thus there is no direct way to translate our argument to the undirected case, because the problem is known to be polynomially solvable for undirected graphs [12].

## 2 Preliminaries

We write  $uv$  for the arc  $(u, v)$ . The vertex set  $V$  is sometimes identified with  $\{1, 2, \dots, n\}$ . For a subset  $W \subseteq V$  of the vertices of a graph  $G$ , we denote by  $G[W]$  the graph induced by  $W$ .

Our proof starts with a reduction from a problem known to be NP-complete for over twenty years. In the  $k$  *Vertex Disjoint Paths* problem we are given a digraph  $G$  of order  $n > 2k$ , and we are asked whether there exists a set of  $k$  vertex disjoint paths in  $G$  such that the  $i$ th path connects vertex  $2i - 1$  to vertex  $2i$ , for  $i = 1, \dots, k$ . This problem is NP-complete [5] even when  $k = 2$ . We need to modify this result slightly to see that it is valid even if we restrict the ‘yes’-instances to be partitionable into two disjoint paths. To be precise, we define the *Two Vertex Disjoint Paths* problem (2VDP): given a digraph  $G$  of order  $n \geq 4$ , decide whether there exists a pair of vertex disjoint paths, one from 1 to 2 and one from 3 to 4. We study the *restricted* version of this problem (R2VDP), where the ‘yes’-instances are guaranteed to contain two such paths that together exhaust all vertices of  $G$ . In other words, the graph  $G$  with the additional arcs 23 and 41 contains a Hamiltonian cycle through these arcs.

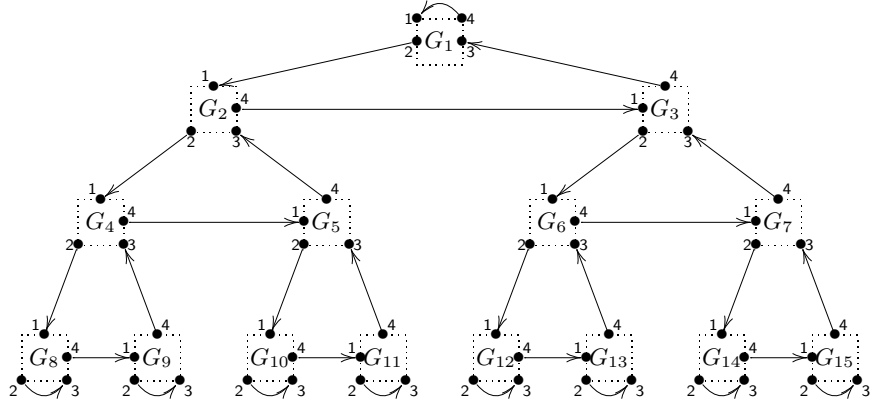


Fig. 1.  $T_4[G]$ .

**Proposition 1** *Restricted Two Vertex Disjoint Paths is NP-complete.*

The proof is an extension of the construction in [5] and can be found in Sec. 7. It replaces a reduction from 3-Sat by a reduction from Monotone 1-in-3-Sat, and uses a more intricate clause gadget to guarantee the existence of two paths that cover all vertices. The modification is necessary to prove the lower bound for Longest Path even for Hamiltonian instances.

### 3 Long Paths Find Vertex Disjoint Paths

We will use instances of R2VDP to build graphs in which long paths must reveal a solution to the original problem. Given an instance  $G = (V, A)$  of R2VDP, define  $T_d[G]$  as a graph made up out of  $m = 2^d - 1$  copies  $G_1 \cdots G_m$  of  $G$  arranged in a balanced binary tree structure. For all  $i < 2^{d-1}$ , we say that the copies  $G_{2i}$  and  $G_{2i+1}$  are the left and right *child* of the copy  $G_i$ . The copy  $G_1$  is the *root* of the tree, and  $G_i$  for  $i \geq 2^{d-1}$  are the *leaves* of the tree. The copies of  $G$  in  $T_d[G]$  are connected by additional arcs as follows. For every copy  $G_i$  having children, three arcs are added (cf. Fig. 1): One arc from 2 in  $G_i$  to 1 in  $G_{2i}$ , one arc from 4 in  $G_{2i}$  to 1 in  $G_{2i+1}$ , and one arc from 4 in  $G_{2i+1}$  to 3 in  $G_i$ . Moreover, in every leaf copy  $G_i$  ( $i \geq 2^{d-1}$ ) we add the arc 23, and in the root  $G_1$  we add the arc 41.

**Lemma 1** *Given an instance  $G = (V, A)$  of R2VDP on  $n = |V|$  vertices, and any integers  $m = 2^d - 1 > 3$ , consider  $T_d[G]$  with  $N = mn$  vertices. Then*

- *If  $G$  has a solution then  $T_d[G]$  contains a path of length  $N - 1$ .*
- *Given any path of length larger than  $(4d - 5)n$  in  $T_d[G]$ , we can in time polynomial in  $N$  construct a solution to  $G$ .*

*Proof.* For the first part of the lemma, consider a solution for  $G$  consisting of two disjoint paths  $P$  and  $Q$  connecting 1 to 2 and 3 to 4, respectively, such that  $P + 23 + Q + 41$  is a Hamiltonian cycle in  $G$ . The copies of  $P$  and  $Q$  in all  $G_i$ s together with the added arcs constitute a Hamiltonian cycle in  $T_d[G]$  of length  $mn$  and thus a path of the claimed length.

For the second part, first consider an internal copy  $G_i$  and observe that if a path traverses all of the four arcs connecting  $G_i$  to the rest of the structure then this path constitutes a solution to R2VDP for  $G$ . Thus we can restrict our attention to paths in  $T_d[G]$  that avoid at least one the four external arcs of each internal  $G_i$ ; we call such paths *avoiding*.

Given  $T_d[G]$  define  $e_d[G]$  as the length of the longest avoiding path in  $T_d[G]$  ending in vertex 4 of its root copy, and  $s_d[G]$  as the length of the longest avoiding path starting in vertex 1 of the root copy. Consider a path  $P$  ending in vertex 4 of the root copy, for  $d > 1$ . At most  $n$  vertices of  $P$  are in  $G_1$ . The path  $P$  has entered  $G_1$  via vertex 3 from  $G_3$ 's vertex 4. There are two possibilities. Either the first part of  $P$  is entirely in the subtree rooted at  $G_3$ , in which case  $P$  has length at most  $n + e_{d-1}[G]$ . Or it entered  $G_3$  via 1 from the subtree rooted at  $G_2$ , in which case it may pass through at most  $n$  vertices in  $G_3$ , amounting to length at most  $2n + e_{d-1}[G]$ . (Especially,  $P$  cannot leave via  $G_3$ 's vertex 2, because then it wouldn't be avoiding). A symmetric argument for  $s_d[G]$  for  $d > 1$  shows an equivalent relation. Thus we have that

$$\begin{aligned} e_1[G] &\leq n, & e_{d+1}[G] &\leq 2n + e_d[G], \\ s_1[G] &\leq n, & s_{d+1}[G] &\leq 2n + s_d[G]. \end{aligned}$$

Furthermore, note that a longest avoiding path in  $T_d[G]$  connects a path amounting to  $e_{d-1}[G]$  in the right subtree, through a bridge consisting of as many vertices as possible in the root, with a path amounting to  $s_{d-1}[G]$  in the left subtree. Consequently, a typical longest avoiding path starts in a leaf copy of the right subtree, travels to its sister copy, goes up a level and over to the sister of that copy, continues straight up in this zigzag manner to the root copy, and down in the same fashion on the other side. Formally, the length of a longest avoiding path in  $T_d[G]$  for  $d > 1$  is bounded from above by  $e_{d-1}[G] + n + s_{d-1}[G] \leq (4d - 5)n$ .  $\square$

**Theorem 1** *There can be no deterministic, polynomial time approximation algorithm for Longest Path or Longest Cycle in a Hamiltonian directed graph on  $n$  vertices with performance ratio  $n^{1-\epsilon}$  for any fixed  $\epsilon > 0$ , unless  $P = NP$ .*

*Proof.* First consider the path case. Given an instance  $G = (V, A)$  of R2VDP with  $n = |V|$ , fix  $k = 1/\epsilon$  and construct  $T_d[G]$  for the smallest integers  $m = 2^d - 1 \geq (4dn)^k$ . Note that the graph  $T_d[G]$  has order  $N = n^{O(k)}$ . Assume there is a deterministic algorithm finding a long path of length  $l_{\text{apx}}$  in time polynomial in  $N$ , and let  $l_{\text{opt}}$  denote the length of a longest path. Return ‘yes’ if and only if  $l_{\text{apx}} > (4d - 5)n$ . To see that this works note that if  $G$  is a ‘yes’-instance and if indeed  $l_{\text{opt}}/l_{\text{apx}} \leq N^{1-\epsilon}$  then  $l_{\text{apx}} > (4d - 5)n$ , so Lem. 1 gives a solution to  $G$ .

If on the other hand  $G$  is a ‘no’-instance then the longest path must be *avoiding* as defined in the proof of Lem. 1, so its length is at most  $(4d - 5)n$ . Thus we can solve the R2VDP problem in polynomial time, which by Prop. 1 requires  $P = NP$ .

For the cycle case, we may use a simpler construction. Simply connect  $m$  copies  $G_1, \dots, G_m$  of  $G$  on a string, by adding arcs from vertex 2 in  $G_i$  to vertex 1 in  $G_{i+1}$ , and arcs from vertex 4 in  $G_i$  to vertex 3 in  $G_{i-1}$ . Finally, add the arc 41 in  $G_1$  and the arc 23 in  $G_m$ . The resulting graph has a cycle of length  $mn$  whenever  $G$  is a ‘yes’-instance, but any cycle of size at least  $2n + 1$  must reveal a solution to  $G$ .  $\square$

## 4 Subexponential Algorithms for Satisfiability

In this section we show that good dipath and dicycle algorithms imply subexponential time algorithms for Satisfiability.

We need the well-known reduction from Monotone 1-in-3-Sat to 3-Sat. It can be verified that the number of variables in the construction (see also [11, Exerc. 9.5.3]) is not too large:

**Lemma 2 ([13])** *Given a 3-Sat instance  $\varphi$  with  $s$  variables and  $r$  clauses we can construct an instance of Monotone 1-in-3-Sat with  $O(r)$  clauses and variables that is satisfiable if and only if  $\varphi$  is.*

**Lemma 3** *There is a deterministic algorithm for Monotone 1-in-3-Sat on  $r$  variables running in time  $2^{o(r)}$ , if there is*

1. *a polynomial time deterministic approximation algorithm  $A_{LP}$  for Longest Path in  $N$ -node Hamiltonian digraphs with guarantee  $f(N) \log^2 N$ , or*
2. *a polynomial time deterministic approximation algorithm  $A_{LC}$  for Longest Cycle in  $N$ -node Hamiltonian digraphs with guarantee  $f(N) \log N$ ,*

*where  $f$  is any polynomial time computable, nondecreasing function in  $\omega(1)$ .*

*Proof.* We need to verify that our constructions obey the necessary size bounds. The R2VDP instance  $G$  build from the instance to Monotone 1-in-3-Sat described in Sec. 7 has size  $n = O(r)$ .

For the path case, set  $d = 4n/f^{1/2}(n)$  and construct  $T_d[G]$  as in Sec. 3. Observe that the entire construction will have  $(2^d - 1)n = 2^{o(n)} = 2^{o(r)}$  nodes. Running  $A_{LP}$  on a ‘yes’ instance instance will reveal a cycle of length  $f(n(2^d - 1)) \log^2(n(2^d - 1)) \geq f(n) \log^2(2^{d/2}) \geq 4n^2 > (4d - 5)n$ , so Lem. 1 tells us how to use  $A_{LP}$  to solve the R2VDP instance, and hence the 1-in-3-Sat instance.

For the cycle case, choose the number of copies  $m = 2^{3n/f(n)}$ . Observe that the entire construction has size  $mn = 2^{o(n)} = 2^{o(r)}$ . Running  $A_{LC}$  on this graph will reveal a cycle of length  $f(mn) \log(mn) \geq f(n) \log m = f(n) \cdot 3n/f(n) = 3n > 2n + 1$ , and the conclusion follows similarly to the proof of Theorem 1.  $\square$



**Theorem 2** *There is a deterministic algorithm for 3-Sat on  $s$  variables running in time  $2^{o(s)}$  if there is*

1. *a polynomial time deterministic approximation algorithm for Longest Path in  $N$ -node Hamiltonian digraphs with guarantee  $f(N) \log^2 N$ , or*
2. *a polynomial time deterministic approximation algorithm for Longest Cycle in  $N$ -node Hamiltonian digraphs with guarantee  $f(N) \log N$ ,*

where  $f$  is any polynomial time computable, nondecreasing function in  $\omega(1)$ .

*Proof.* The previous two lemmas give an algorithm that runs in time  $2^{o(r)}$ , where  $r$  is the number of clauses in the input instance. This implies a  $2^{o(s)}$ -algorithm by the Sparsification Lemma of [9].  $\square$

## 5 Finding Long Paths in Hamiltonian Digraphs

Vishwanathan [14] presents a polynomial time algorithm that finds a path of length  $\Omega(\log^2 n / \log \log n)$  in *undirected* Hamiltonian graphs of constant bounded degree. We show in this section that with some modifications the algorithm and its analysis apply to the directed case as well.

**Theorem 3** *There is a polynomial time algorithm always finding a path of length  $\Omega(\log^2 n / \log \log n)$  in any Hamiltonian digraph of constant bounded outdegree on  $n$  vertices.*

To prove the theorem, we need some additional notation. Let  $G = (V, A)$  be a digraph. We say that a vertex  $v \in V$  *spans* the subgraph  $G_v = G[V_v]$  where  $V_v \subseteq V$  is the set of vertices reachable from  $v$  in  $G$ . Consider the algorithm below. It takes a digraph  $G = (V, A)$  on  $n = |V|$  vertices and a specified vertex  $v \in V$  as input, and returns a long path starting in  $v$ .

1. Enumerate all paths in  $G$  starting in  $v$  of length  $\log n$ , if none return the longest found.
2. For each such path  $P = (v, \dots, w)$ , let  $V_w$  be the set of vertices reachable from  $w$  in  $G[V - P + \{w\}]$ .
3. Compute a depth first search tree rooted at  $w$  in  $G[V_w]$ .
4. If the deepest path in the tree is longer than  $\log^2 n$ , return this path.
5. Otherwise, select the enumerated path  $P$  whose end vertex  $w$  spans as large a subgraph as possible after removal of  $P - \{w\}$  from the vertex set, i.e the path maximising  $|V_w|$ .
6. Search recursively for a long path  $R$  starting from  $w$  in  $G[V_w]$ , and return  $(P - \{w\}) + R$ .

First note that the algorithm indeed runs in polynomial time. The enumeration of all paths of length  $\log n$  takes no more than polynomial time since the outdegree is bounded by a constant  $k$ , and thus there cannot be more than  $k^{\log n}$  paths. Computing a depth first search tree is also a polynomial time task, and it

is seen to be performed a polynomial number of times, since the recursion does not branch at all.

To prove that the length of the resulting path is indeed  $\Omega(\log^2 n / \log \log n)$ , we need to show that at each recursive call of the algorithm, there is still a long enough path starting at the current root vertex.

**Lemma 4** *Let  $G = (V, A)$  be a Hamiltonian digraph. Let  $S \subseteq V, v \in V \setminus S$ . Suppose that on removal of the vertices of  $S$ ,  $v$  spans the subgraph  $G_v = (V_v, A_v)$  of size  $t$ . If each vertex  $w \in V_v$  is reachable from  $v$  on a path of length less than  $d$ , then there is a path of length  $t/d|S|$  in  $G_v$  starting in  $v$ .*

*Proof.* Consider a Hamiltonian cycle  $C$  in  $G$ . The removal of  $S$  cuts  $C$  into at most  $|S|$  paths  $P_1 \cdots P_{|S|}$ . Since each vertex in  $V$  lies on  $C$ , the subgraph  $G_v$  must contain at least  $t/|S|$  vertices  $W$  from one of the paths, say  $P_j$ . In fact,  $G_v$  must contain a path of length  $t/|S|$ , since the vertex in  $W$  first encountered along  $P_j$  implies the presence in  $G_v$  of all the subsequent vertices on  $P_j$ , and these are at least  $|W|$ . Denote one such path by  $P = p_0 \cdots p_{|W|-1}$ , and let  $R = r_0 \cdots r_{l-1}$  be a path from  $r_0 = v$  to  $r_{l-1} = p_0$ , of length  $l \leq d$ . Set  $s = |P \cap R|$  and enumerate the vertices on  $P$  from 0 to  $|W| - 1$  and let  $i_1 \cdots i_s$  denote the indices of vertices in  $P \cap R$ , in particular  $i_1 = 0$ . Let  $i_{s+1} = |W|$ . An averaging argument shows that there exists  $j$ , such that  $i_{j+1} - i_j \geq |W|/s$ . Let  $q$  be the index for which  $r_q = p_{i_j}$ . The path along  $R$  from  $r_0$  to  $r_q$  and continuing along  $P$  from  $p_{i_{j+1}}$  to  $p_{i_{j+1}-1}$  has the claimed length.  $\square$

Observe that the algorithm removes no more than  $\log n$  vertices from the graph at each recursive call. Thus, at call  $i$  we have removed at most  $i \log n$  vertices from the original graph; the very same vertices constituting the beginning of our long path. Lemma 4 tells us that we still are in a position were it is possible to extend the path, as long as we can argue that the current end vertex of the path we are building spans large enough a subgraph. Note that whenever we stand at a vertex  $v$  starting a long path  $P$  of length  $> \log n$  in step 1 of the algorithm, the path consisting of the first  $\log n$  vertices of  $P$  is one of the paths of length  $\log n$  being enumerated. This is our guarantee that the subgraph investigated at the next recursive call is not all that smaller than the graph considered during the previous one. It must consist of at least  $|P| - \log n$  vertices. Of course, we cannot be sure that exactly this path is chosen at step 5, but this is of no concern, since it is sufficient for our purposes to assure that there are still enough vertices reachable.

Formally, let  $V_i$  denote the vertex set of the subgraph considered at the recursive call  $i$ . In the beginning, we know that regardless of the choice of start vertex  $v$ , we span the whole graph and thus  $V_0 = V$ , and furthermore, that a path of length  $n$  starts in  $v$ . Combining the preceding discussion with Lem. 4, we establish the following inequality for the only non-trivial case that no path of length  $\log^2 n$  is ever found during step 4 of the algorithm:

$$|V_{i+1}| > \frac{|V_i|}{i \log^3 n} - \log n$$

It is readily verified that  $|V_i| > 0$  for all  $i < c \log n / \log \log n$  for some constant  $c$ , which completes the proof of Theorem. 3.

## 6 Finding Long Paths in Sparse Expanders

In this section we show that in a sparse *expander* graph, a relatively long path is easily found.

A digraph  $G = (V, A)$  on  $n$  vertices is a  $c$ -expander if  $|\delta U| \geq c(1 - \frac{|U|}{n})|U|$  for every subset  $U \subset V$  where  $\delta U = \{v \notin U \mid \exists u \in U : uv \in A\}$ . A standard probabilistic argument shows that with high probability a random digraphs with outdegree  $k$  ( $k > 2$ ), are  $c_k$ -expanders for some constant  $c_k$ , for large enough  $n > n_k$ .

We propose the following algorithm for finding a long path  $p_0 \cdots p_l$  in a sparse expander.

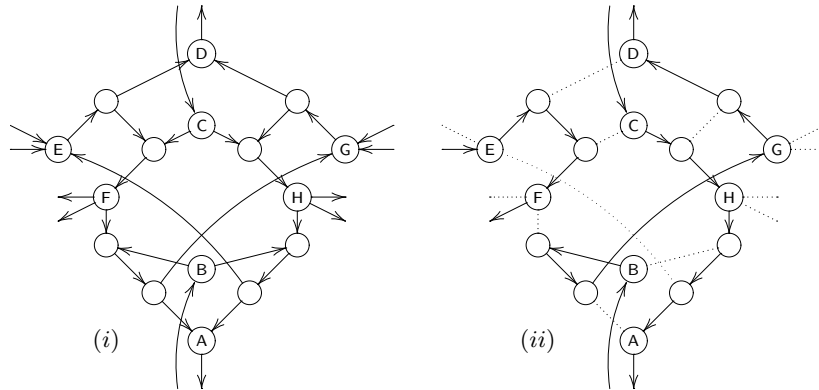
1. Pick an arbitrary start vertex  $p_0$ , and set  $i = 0$ .
2. Let  $G_i = (V_i, A_i)$  be the subgraph spanned by  $p_i$  in  $G[V \setminus (\bigcup_{j=0}^{i-1} p_j)]$ .
3. If  $G_i$  consists only of  $p_i$ , exit.
4. For each neighbour  $v$  of  $p_i$  in  $G_i$ , evaluate the size of the subgraph spanned by  $v$  in  $G_i[V_i \setminus p_i]$ .
5. Choose the neighbour who has the largest spanned subgraph as  $p_{i+1}$ .
6. Set  $i = i + 1$  and goto 2.

**Theorem 4** *The algorithm finds a path of length  $\frac{c}{2(k+1)}n$  in every  $c$ -expander digraph  $G = (V, A)$  with maximum outdegree  $k$ .*

*Proof.* Consider step  $i$ . Enumerate the neighbours of  $p_i$  in  $G_i$  as  $r_1 \cdots r_{k'}$ . Let  $V_i[r_j]$  be the vertices reachable from  $r_j$  in  $G_i[V_i - \{p_i\}]$ . Now observe that the  $V_i[r_j]$  either are very small or really large for small  $i$ , since the set of vertices outside  $V_i[r_j]$  in  $G$  which are directly connected by an arc from a vertex in  $V_i[r_j]$  must lie on the prefix path  $p_0 \cdots p_i$  by definition, and there must be a lot of them because of the expander criterion. Specifically, when  $i$  is small, there must be a  $j$  for which  $V_i[r_j]$  is large, since  $k' \leq k$  and  $\bigcup V_i[r_j] = V_i - \{p_i\}$ . Observe that  $V_{i+1}$  is the largest  $V_i[r_j]$ , to obtain  $|V_{i+1}| \geq n - \frac{2(i+1)}{c}$  whenever at least one  $V_i[r_j]$  is too large to be a small subgraph, i.e. as long as  $\frac{c(|V_i| - 1)}{2k} \geq i + 1$ , where we for the sake of simplicity have used the expansion factor  $c/2$  which holds for all set sizes. Observing that  $V_0 = n$ , we may solve for the smallest  $i$ , when the inequality above fails to hold. This will not happen unless  $i \geq \frac{c}{2(k+1)}n$ , as promised.  $\square$

## 7 Proof of Proposition 1

We review the construction in [5], in which the *switch* gadget from Fig. 2 plays a central role. Its key property is captured in the following statement.



**Fig. 2.** (i) A switch. Only the labelled vertices are connected to the rest of the graph, as indicated by the arrows. (ii) Three vertex-disjoint paths through a switch.

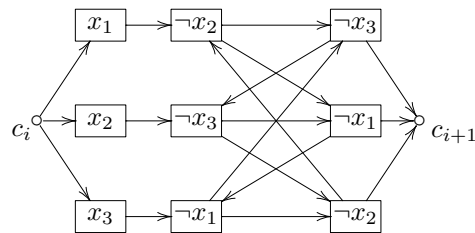
**Lemma 5 ([5])** Consider the subgraph in Fig. 2. Suppose that there are two vertex-disjoint paths passing through the subgraph—one leaving at A and the other entering at B. Then the path leaving A must have entered at C and the path entering at B must leave at D. Furthermore, there is exactly one additional path through the subgraph and it connects either E to F or G to H, depending on the actual routing of the path leaving at A.

Also, if one of these additional paths is present, all vertices are traversed.

To prove Prop. 1 we reduce from Monotone 1-in-3-Satisfiability, rather than 3-Satisfiability as used in [5]. An instance of 1-in-3-Sat is a Boolean expression in conjunctive normal form in which every clause has three literals. The question is if there is a truth assignment such that in every clause, exactly one literal is true. It is known that even when all literals are positive (*Monotone 1-in-3-Sat*) the problem is NP-complete [13].

Given such an instance  $\varphi$  with clauses  $t_1, \dots, t_m$  on variables  $x_1, \dots, x_n$ , we construct an instance  $G_\varphi$  of R2VDP as follows.

*Clause gadgets.* Every clause  $t_i$  is represented by a gadget consisting of a vertex  $c_i$  and nine switches, three for every literal in  $t_i$ . Consider the clause  $t_i = (x_1 \vee x_2 \vee x_3)$ . The vertices  $c_i, c_{i+1}$  and the E and F vertices in the nine switches are connected as shown in Fig. 3. Thus all clause gadgets are connected on a string ending in a dummy vertex  $c_{m+1}$ .



**Fig. 3.** A clause gadget consisting of 9 switches. Every incoming arc to a switch enters the switch's vertex E; every outgoing arc leaves the switch's vertex F.

The clause gadget has the following desirable properties: Call a path from  $c_i$  to  $c_{i+1}$  *valid* if it is consistent with a truth assignment to  $\{x_1, x_2, x_3\}$  in the sense that if it passes through a switch labelled with a literal (like  $\neg x_2$ ) then it cannot pass through its negation (like  $x_2$ ). The following claims are easily verified:

**Lemma 6** *Consider the construction in Fig. 3.*

1. *Every valid path from  $c_i$  to  $c_{i+1}$  corresponds to a truth assignment to  $\{x_1, x_2, x_3\}$  that sets exactly one variable to true.*
2. *If there is a truth assignment to  $\{x_1, x_2, x_3\}$  that sets exactly one variable to true then there is a valid path from  $c_i$  to  $c_{i+1}$  corresponding to the assignment. Moreover, there is such a valid path passing through all five switches whose labels are consistent with the assignment.*

*Variable gadgets.* Every variable  $x_i$  is represented by a vertex  $v_i$ . (Again, vertex  $v_{n+1}$  is a dummy vertex.) All switches in the clause gadgets representing the positive literal of the variable  $v_i$  are connected in series (the ordering of the switches on this string is not important): the vertex H in a switch is connected to vertex G of the next switch with the same label. Furthermore, there is an arc from  $v_i$  to vertex G in the first switch on its literal path, and an arc from vertex H in the last switch on the path to vertex  $v_{i+1}$ .

Likewise, all switches labelled with negated literals of this variable are connected. Thus there are two strings of switches leaving  $v_i$ : one contains all the positive literals, and one contains all the negated literals. Both end in  $v_{i+1}$ .

Also, *all* the switches are arranged on a path and connected by added arcs from vertex A in a switch to vertex C in the next one, and arcs back from vertex D in a switch to vertex B of the preceding switch. The ordering of the switches on this switch path is not important.

Finally, there is an arc from  $v_{n+1}$  to  $c_1$  and an arc from vertex D in the first switch on the switch path to  $v_1$ .

To finish the construction of an instance of R2VDP it remains to identify the first four vertices. Vertex 1 is vertex B of the last switch on the switch path, vertex 2 is  $c_{m+1}$ , vertex 3 is vertex C of the first switch on the switch path, and vertex 4 is vertex A of the last switch on the switch path.

**Lemma 7**  *$G_\varphi$  has two vertex disjoint paths from 1 to 2 and from 3 to 4 if and only if  $\varphi$  has a solution. Moreover, if  $G_\varphi$  contains such paths then it contains two such paths that together exhaust all its vertices.*

*Proof.* Assume  $\varphi$  can be satisfied so that exactly one variable in every clause is true. Walk through  $G_\varphi$  starting in vertex 1. This path is forced to traverse all switches until it reaches  $v_1$ . In general, assume that we reached  $v_i$ . To continue to  $v_{i+1}$  traverse the G–H paths of the string of negative literal switches if  $x_i$  is true; otherwise take the string of positive literal switches. Note that this forces us to avoid the E–F paths in these switches later.

Arriving at  $v_{n+1}$  continue to  $c_1$ . To travel from  $c_i$  to  $c_{i+1}$  we are forced to traverse the clause gadget of Fig. 3. Note that the truth assignment has set

exactly one of the variables to true, blocking the E–F path in the two switches labelled by its negative literal. Likewise, two of the variables are false, blocking the (two) switches labelled by their positive literal. The remaining five switches are labelled by the positive literal of the true variable or negative literals of the falsified variables. The valid path ensured by Lem. 6 passes through exactly these five switches.

Finally, the path arrives at  $v_{m+1} = 2$ . The path travelling from 3 to 4 is now unique. Observe that the two paths exhaust all the vertices and thus form a Hamiltonian cycle if we add 23 and 41.

Conversely, assume there are two paths from 1 to 2 and from 3 to 4. The subpaths connecting  $v_i$  to  $v_{i+1}$  ensure that all literal switches are consistent in the sense that if the E–F path in a switch labelled  $x_i$  is blocked then it is blocked in *all* such switches, and not blocked in any switch labelled  $\neg x_i$ . This forces the subpaths from  $c_i$  to  $c_{i+1}$  to be valid. Lem. 6 ensures that the corresponding truth assignment is satisfying and sets exactly one variable in each clause.  $\square$

**Acknowledgements.** The third author would like to express his thanks to Chandra Chekuri for many useful discussions on this problem. Hal Gabow suggested the formulation of the bound in Thm. 2.

## References

1. N. Alon, R. Yuster, and U. Zwick. Color-coding. *Journal of the ACM*, 42(4):844–856, 1995.
2. A. Björklund and T. Husfeldt. Finding a path of superlogarithmic length. *SIAM Journal on Computing*, 32(6):1395–1402, 2003.
3. Béla Bollobás. *Random graphs*. Cambridge University Press, 2nd edition, 2001.
4. T. Feder, R. Motwani, and C. Subi. Approximating the longest cycle problem in sparse graphs. *SIAM Journal on Computing*, 31(5):1596–1607, 2002.
5. S. Fortune, J. Hopcroft, and J. Wyllie. The directed subgraph homeomorphism problem. *Theoretical Computer Science*, 10:111–121, 1980.
6. H. N. Gabow. Finding paths and cycles of superlogarithmic length. In *Proc. 36th STOC*, 2004.
7. H. N. Gabow and S. Nie. Finding a long directed cycle. In *Proc. 15th SODA*, 2004.
8. R. Impagliazzo and R. Paturi. On the complexity of  $k$ -SAT. *Journal of Computer and Systems Sciences*, 62(2):367–375, 2001.
9. R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? In *Proc. 39th FOCS*, pages 653–663, 1998.
10. D. Karger, R. Motwani, and G.D.S. Ramkumar. On approximating the longest path in a graph. *Algorithmica*, 18(1):82–98, 1997.
11. C. H. Papadimitriou. *Computational Complexity*. Addison-Wesley, 1994.
12. N. Robertson and P. D. Seymour. Graph minors XIII: The disjoint paths problem. *J. Combinatorial Theory Ser. B*, 35, 1983.
13. T. J. Schaefer. The complexity of satisfiability problems. In *Proc. 10th STOC*, pages 216–226, 1978.
14. S. Vishwanathan. An approximation algorithm for finding a long path in Hamiltonian graphs. In *Proc. 11th SODA*, pages 680–685, 2000.