

# Expander Flows, Geometric Embeddings and Graph Partitioning

Sanjeev Arora\*      Satish Rao†      Umesh Vazirani‡

April 12, 2004

## Abstract

We give a  $O(\sqrt{\log n})$ -approximation algorithm for SPARSEST CUT, BALANCED SEPARATOR, and GRAPH CONDUCTANCE problems. This improves the  $O(\log n)$ -approximation of Leighton and Rao (1988). We use a well-known semidefinite relaxation with triangle inequality constraints. Central to our analysis is a geometric theorem about projections of point sets in  $\mathfrak{R}^d$ , whose proof makes essential use of a phenomenon called measure concentration.

We also describe an interesting and natural “certificate” for a graph’s expansion, by embedding an  $n$ -node expander in it with appropriate dilation and congestion. We call this an expander flow.

## 1 Introduction

Partitioning a graph into two (or more) large pieces while minimizing the size of the “interface” between them is a fundamental combinatorial problem. Graph partitions or separators are central objects of study in the theory of Markov chains, geometric embeddings and are a natural algorithmic primitive in numerous settings, including clustering, divide and conquer approaches, PRAM emulation, VLSI layout, and packet routing in distributed networks. Since finding optimal separators is NP-hard, one is forced to settle for approximation algorithms (see [31]).

Here we give new approximation algorithms for some of the important problems in this class. In a graph  $G = (V, E)$ , for any cut  $(S, \bar{S})$  where  $|S| \leq |V|/2$ , the *edge expansion* of the cut is  $|E(S, \bar{S})| / |S|$ . In the SPARSEST CUT problem we wish to determine the cut with the smallest edge expansion:

$$\alpha(G) = \min_{S \subseteq V, |S| \leq |V|/2} \frac{|E(S, \bar{S})|}{|S|}. \quad (1)$$

A cut  $(S, \bar{S})$  is *c-balanced* if both  $S, \bar{S}$  have at least  $c|V|$  vertices. In the  $c$ -BALANCED-SEPARATOR problem we wish to determine  $\alpha_c(G)$ , the minimum expansion of  $c$ -balanced cuts. In the GRAPH CONDUCTANCE problem we wish to determine

$$\Phi(G) = \min_{S \subseteq V, |E(S)| \leq |E|/2} \frac{|E(S, \bar{S})|}{|E(S)|}, \quad (2)$$

---

\*Computer Science Department, Princeton University. [www.cs.princeton.edu/~arora](http://www.cs.princeton.edu/~arora). Work done partly while visiting UC Berkeley (2001-2002) and the Institute for Advanced Study (2002-2003). Supported by David and Lucille Packard Fellowship, and NSF Grants CCR-0098180 and CCR-0205594.

†Computer Science Division, UC Berkeley. [www.cs.berkeley.edu/~satishr](http://www.cs.berkeley.edu/~satishr). Partially supported by NSF award CCR-0105533

‡Computer Science Division, UC Berkeley. [www.cs.berkeley.edu/~vazirani](http://www.cs.berkeley.edu/~vazirani) Partially supported by NSF ITR Grant CCR-0121555

where  $E(S)$  denotes the set of edges incident to nodes in  $S$ . We can reduce each of these problems to constant degree graphs and moreover for this class, edge expansion and conductance are related by a constant factor.

A weak approximation for GRAPH CONDUCTANCE follows from the connection —first discovered in context of Riemannian manifolds [8]—between conductance and the eigenvalue gap of the Laplacian:  $2\Phi(G) \geq \lambda \geq \Phi(G)^2/2$  [3, 2, 21]. The approximation factor is  $1/\Phi(G)$ , and hence  $\Omega(n)$  in the worst case, however for constant  $\phi(G)$  it is an excellent bound. This connection between eigenvalues and expansion has had enormous influence in a variety of fields (see e.g. [9]).

Leighton and Rao [22] designed the first true approximation by giving  $O(\log n)$ -approximations for SPARSEST CUT and GRAPH CONDUCTANCE and  $O(\log n)$ -pseudo-approximations for  $c$ -BALANCED SEPARATOR. They used a linear programming relaxation of the problem based on multicommodity flow proposed in [30]. This led to approximation algorithms for numerous NP-hard problems, see [31]. However, the integrality gap of the LP is  $\Omega(\log n)$ , and crossing this  $\log n$  barrier therefore calls for new techniques.

In this paper we give  $O(\sqrt{\log n})$ -approximations for SPARSEST CUT and GRAPH CONDUCTANCE and  $O(\sqrt{\log n})$ -pseudo-approximation to  $c$ -BALANCED SEPARATOR.

Now we give a quick overview of our results. Our algorithm uses *semidefinite programming* (SDP). These are mathematical programs in which each vertex  $i$  is assigned some point  $v_i$  on the unit sphere in  $\mathbb{R}^n$ . In our case the goal is to find an assignment such that the average distance between all pairs of points is “large” whereas the average distance between endpoints of edges is minimized.

The complexity of finding such embeddings depends crucially on the notion of distance. Under the standard Euclidean norm (or even  $\ell_1$  norm) the problem is NP hard: the optimum cut can be efficiently recovered from the optimum vectors. The notion of distance that is more tractable (and used in SDPs) is the *square* of the Euclidean norm, the so-called  $\ell_2^2$  norm. With this distance function, the embedding problem is related to finding eigenvectors of the adjacency matrix of the graph and thus yields only weak approximations. To tighten the relaxation, one can ask that  $\ell_2^2$  distances between the  $v_i$ 's form a metric: every triple  $i, j, k$  satisfies the triangle inequality, i.e.  $|v_i - v_j|^2 + |v_j - v_k|^2 \geq |v_i - v_k|^2$ . The general SDP framework allows such constraints. Furthermore, these constraints correspond exactly to the linear constraints in the Leighton-Rao LP relaxation and therefore this  $\ell_2^2$  embedding subsumes both the eigenvalue as well as the  $O(\log n)$  linear programming bound. The conjectured integrality gap of the resulting relaxation to the cut problem is  $O(1)$  [15], and is known to be at least  $10/9$  [34].

Our  $O(\sqrt{\log n})$ -approximation relies on a new result about the geometric structure of such embeddings: they contain  $\Omega(n)$  sized sets  $S$  and  $T$  that are well-separated, in the sense that every pair of points  $v_i \in S$  and  $v_j \in T$  must be at least  $\Delta = \Omega(1/\sqrt{\log n})$  apart in  $\ell_2^2$  distance. (We also present a randomized algorithm to find such sets.) This result is tight for an  $n$ -vertex hypercube —whose natural embedding into  $\mathbb{R}^{\log n}$  defines an  $\ell_2^2$  metric— where any two large sets are within  $O(1/\sqrt{\log n})$  distance.

Finding such a well-separated subset pair suffices for a good approximation. Since the sum of the  $\ell_2^2$  distances between endpoints of edges is small in the embedding, the sets  $S$  and  $T$  cannot have many edges between them, and this is the basis of finding a small cut. Formally, finding a good separator involves shrinking  $S$  to a point and performing a breadth first search from it and outputting the level with fewest edges (Section 2.1.1).

The algorithm for finding the above-mentioned well-separated pair  $S, T$  is complicated (Section 5) but we also describe a simpler algorithm (Section 3) that works for a somewhat smaller separation  $\Delta = \Omega(1/\log^{2/3} n)$ . In that case the idea is to partition the vectors with a randomly oriented hyperplane slice of prescribed thickness. Points that fall inside the slice are discarded. The sets of points on the two sides of the slice are our first candidates for  $S, T$ . However, they can contain a few pairs of points  $v_i \in S, v_j \in T$  whose squared distance is less than  $\Delta$ , which we discard. The technically hard part in the analysis is to prove that not too many points get discarded. This makes essential use of a phenomenon called *measure concentration*, a cornerstone of modern convex geometry [5].

**Graph embeddings and expander flows:** Our ideas also imply a new structural result in graph theory:

an embedding of expander graphs in any arbitrary graph that is more efficient (in terms of maximum edge congestion, which is the number of expander edges routed through a single graph edge) than any known before. This result is proved using techniques similar to the ones used to prove the existence of the  $\Delta$ -separated sets (though it is not an immediate corollary and requires some work). To understand the connection to approximation algorithms, note that any algorithm that approximates edge expansion  $\alpha = \alpha(G)$  must implicitly certify that every cut has large expansion. One way to do this is to embed a complete graph into the given graph with minimum congestion,  $\mu$ . Clearly, every cut must have expansion at least  $1/n\mu$ . (See Section 7.) This is exactly the certificate used in the Leighton-Rao paper, where it is shown that congestion  $O(\alpha/n \log n)$  suffices (and this amount of congestion is required on some graphs.)

This paper considers a generalization of this approach, where we embed not the complete graph but some flow that is an expander. We show how this idea can be used to derive a certificate to the effect that the expansion is  $\Omega(\alpha/\sqrt{\log n})$  (see Section 7). The conjectures presented in the full version imply that this approach can be improved to certify that the expansion is  $\Omega(\alpha)$ .

Expander flows also provide a different and possibly more efficient (though the current writeup ignores efficiency issues besides polynomiality)  $O(\sqrt{\log n})$ -approximation algorithm for graph separators that uses multicommodity flows combined with eigenvalue computations. A polynomial bound follows by observing that embedding a particular graph (the expander flow) with minimum congestion is a multicommodity flow problem. The condition that the embedded graph is an expander can be imposed by exponentially many linear constraints, one for each cut. A violated (within a constant factor) constraint can be efficiently found by an eigenvalue computation, and thus the linear program can be solved by the Ellipsoid method. More details appear in the full version. In fact, the algorithms of this paper (including the SDP rounding) were discovered in this setting.

## Related Work.

**Semidefinite programming and approximation algorithms:** Semidefinite programs (SDPs) have numerous applications in optimization. They are solvable in polynomial time via the ellipsoid method [18], and more efficient interior point methods are now known [1, 27]. In a seminal paper, Goemans and Williamson [17] used SDPs to design good approximation algorithms for MAX-CUT and MAX- $k$ -SAT. Researchers soon extended their techniques to other problems [20, 19, 15], but lately progress in this direction has stalled. Especially in the context of minimization problems, the GW approach of analysing “random hyperplane” rounding in an edge-by-edge fashion runs into well-known problems. By contrast, our main theorem about  $\ell_2^2$  spaces (and the “rounding” technique that follows from it) takes a more global view of the metric space. The ideas may prove useful for other problems where triangle inequality constraints are conjectured to tighten SDP relaxations.

**Analysis of random walks:** The mixing time of a random walk on a graph is related to the first nonzero eigenvalue of the Laplacian, and hence to the conductance. Of various techniques known for upper-bounding the mixing time, most rely on lowerbounding the conductance. Diaconis and Saloff-Coste [11] describe a very general idea called the *comparison technique*, whereby the conductance of a graph is lowerbounded by embedding a *known* graph with *known* conductance into it. (The embedding need not be constructive; existence suffices.) Sinclair [32] suggested a similar technique and also noted that the Leighton-Rao multicommodity flow can be viewed as a generalization of the Jerrum-Sinclair [21] canonical path argument. Our results on expander flows imply that the comparison technique can be used to always get to within  $O(\sqrt{\log n})$  of the proper bound for conductance.

**Metric spaces and relaxations of the cut cone:** The *cut cone* is the cone of all cut semi-metrics, and is equivalent to the cone of all  $\ell_1$  semi-metrics. Graph separation problems can often be viewed as the optimization of a linear function over the cut cone (possibly with some additional constraints imposed). Thus optimization over the cut cone is NP-hard. However, one could relax the problem and optimize over some other metric space, embed this metric space in  $\ell_1$  (hopefully, with low distortion), and then derive an approximation algorithm. This approach was pioneered in [23] and [4]; see Shmoys [31]

for a survey. A major open problem in this area is to show that  $\ell_2^2$  metrics (i.e., solutions to the SDP with triangle inequality constraints) embed into  $\ell_1$  with  $O(1)$  distortion. Showing this would prove an integrality gap of  $O(1)$  not only for SPARSEST CUT but also for a more general version of the problem involving nonuniform demands between vertex pairs. The current paper does not address this conjecture. However, James Lee pointed out to us that our results (specifically, Theorem 1 which was earlier implicit in our paper and now explicit thanks to his observation) do represent partial progress: they give an embedding of  $\ell_2^2$  metrics into  $\ell_1$  in which the *average* edge distorts by at most  $\sqrt{\log n}$  factor. Furthermore, we do note in the full version of the paper that for the version of SPARSEST CUT considered here, the embedding conjecture is overkill. Instead we present weaker conjectures that are sufficient to prove an  $O(1)$  integrality gap. We mention a related conjecture about  $\ell_1$  spaces that is also of interest.

## 2 Definitions and Results

Throughout the paper we will assume that we are dealing with constant degree unweighted graphs, since the general case can be reduced to this case, as is wellknown. Furthermore, GRAPH CONDUCTANCE also reduces to SPARSEST CUT on constant degree graphs.

DEFINITION 1 ( $\ell_2^2$  REPRESENTATION) *A vector representation of a graph is an assignment of a vector to each node, say  $v_i$  assigned to node  $i$ . It is called an  $\ell_2^2$ -representation if for all  $i, j, k$ :*

$$|v_i - v_j|^2 + |v_j - v_k|^2 \geq |v_i - v_k|^2 \quad (\Delta\text{-inequality}) \quad (3)$$

An  $\ell_2^2$ -representation is called a unit- $\ell_2^2$  representation if all its vectors have unit length.

REMARK 1 Equivalently, one can say that the unit- $\ell_2^2$ -representation associates a positive semidefinite  $n \times n$  matrix  $M$  with the graph with diagonal entries 1 and  $\forall i, j, k, M_{ij} + M_{jk} - M_{ik} \leq 1$ . The vector representation  $v_1, v_2, \dots, v_n$  is the Cholesky factorization of  $M$ , namely  $M_{ij} = \langle v_i, v_j \rangle$ .

Every cut  $(S, \bar{S})$  gives rise to a natural unit- $\ell_2^2$  representation, namely, one that assigns some unit vector  $v_0$  to every vertex in  $S$  and  $-v_0$  to every vertex in  $\bar{S}$ . Thus the following SDP is a relaxation for  $\alpha_c(G)$  (scaled by  $cn$ .)

$$\min \quad \frac{1}{4} \sum_{\{i,j\} \in E} |v_i - v_j|^2 \quad (4)$$

$$\forall i \quad |v_i|^2 = 1 \quad (5)$$

$$\forall i, j, k \quad |v_i - v_j|^2 + |v_j - v_k|^2 \geq |v_i - v_k|^2 \quad (6)$$

$$\sum_{i < j} |v_i - v_j|^2 \geq 4c(1-c)n^2 \quad (7)$$

This SDP motivates the following definition.

DEFINITION 2 *An  $\ell_2^2$ -representation is  $c$ -spread if equation (7) holds.*

Similarly the following is a relaxation for sparsest cut (up to scaling by  $n$ ; see Section 6).

$$\min \quad \sum_{\{i,j\} \in E} |v_i - v_j|^2 \quad (8)$$

$$\forall i, j, k \quad |v_i - v_j|^2 + |v_j - v_k|^2 \geq |v_i - v_k|^2 \quad (9)$$

$$\sum_{i < j} |v_i - v_j|^2 = 1 \quad (10)$$

As we mentioned before the SDPs subsume both the eigenvalue approach and the Leighton-Rao approach [15]. We show that the optimum value of the SPARSEST CUT SDP is  $\Omega(\alpha(G)n/\sqrt{\log n})$ , which shows that the integrality gap is  $O(\sqrt{\log n})$ .

## 2.1 Main theorem about $\ell_2^2$ representations

In general,  $\ell_2^2$ -representations are not well-understood<sup>1</sup>.

This is not surprising since in  $\mathfrak{R}^d$  the representation can have at most  $2^d$  distinct vectors [10], so our three-dimensional intuition is of limited use for graphs with more than  $2^3$  vertices. The technical core of our paper is a new theorem about unit  $\ell_2^2$  representations. Note that we assume the dimension  $d \gg \log n$ ; this is without loss of generality since we could always embed the vectors in a higher dimensional space.

**DEFINITION 3 ( $\Delta$ -SEPARATED)** *If  $v_1, v_2, \dots, v_n \in \mathfrak{R}^d$ , and  $\Delta \geq 0$ , two disjoint sets of vectors  $S, T$  are  $\Delta$ -separated if for every  $v_i \in S, v_j \in T$ ,  $|v_i - v_j|^2 \geq \Delta$ .*

**THEOREM 1 (MAIN)**

*For every  $c > 0$ , any  $c$ -spread unit- $\ell_2^2$  representation with  $n$  points contains  $\Delta$ -separated subsets  $S, T$  of size  $\Omega(n)$ , where  $\Delta = \Omega(1/\sqrt{\log n})$ . Furthermore, there is a randomized polynomial-time algorithm for finding these subsets  $S, T$ .*

**REMARK 2** The natural embedding of the boolean hypercube  $\{-1, 1\}^d$  (appropriately scaled) shows that this theorem is tight to within a constant factor. This follows from the isoperimetric inequality for hypercubes.

### 2.1.1 Immediate corollary: $\sqrt{\log n}$ -approximation

Let  $c'$  be the constant in the  $\Omega(n)$  bound on the sizes of sets  $S$  and  $T$  in theorem 1. Let  $W = \sum_{\{i,j\} \in E} |v_i - v_j|^2$  be the optimum value for the SDP defined by equations (4)-(7). (Specifically, it is the objective scaled by 4.) Since the vectors  $v_i$ 's obtained from solving the SDP satisfy the hypothesis of Theorem 1, as an immediate corollary to the theorem we show how to produce a  $c'$ -balanced cut whose expansion is  $O(\sqrt{\log n}W/n)$ .

**COROLLARY 2**

*There is a randomized polynomial-time algorithm that finds with high probability a cut  $(S_{obs}, \overline{S_{obs}})$  that is  $c'$ -balanced, and has expansion  $\alpha_{obs} = O(W\sqrt{\log n}/n)$ .*

**PROOF:** We use the algorithm of Theorem 1 to produce  $\Delta$ -separated subsets  $S, T$  for  $\Delta = g/\sqrt{\log n}$ . Let  $V_0$  denote the vertices whose vectors are in  $S$ . Associate with each edge  $e = \{i, j\}$  a length  $w_e = |v_i - v_j|^2$ . (Thus  $W = \sum_{e \in E} w_e$ .) In the rest of the proof “distance” in the graph is measured with respect to this length function.

Denote by  $V_s$  the set of vertices within distance  $d$  of  $V_0$  and by  $E_s$  the set of edges leaving  $V_s$ . We do breadth-first search and find  $s \leq \Delta/2$  that minimizes  $|E_s|/|V_s|$ . We output the cut  $(V_s, \overline{V_s})$ ; let  $\alpha_{obs}$  denote the expansion of this cut. (We can assume without loss of generality that  $|V_{\Delta/2}| \leq n/2$ , since we could switch  $S$  and  $T$  otherwise.)

For any  $s \leq \Delta/2$ , we have

$$|E_s| \geq \alpha_{obs}|V_s| \geq \alpha_{obs}c'n,$$

since  $|V_0| = |S| \geq c'n$ .

<sup>1</sup>A well-known—but alas, wide-open—conjecture says that they are closely related to the better-understood  $\ell_1$  metrics.

The total length of the edges  $W = \sum_e w_e$ , is thus at least  $\Delta/2$  times the minimum number of edges crossing at any point along this length  $\Delta/2$  interval. More formally, the total length of the edges

$$W = \sum_e w_e \geq \int_{s=0}^{\Delta/2} |E_s| ds \geq \frac{\Delta}{2} \cdot \alpha_{obs} c' n.$$

The corollary follows by solving for  $\alpha_{obs}$ .  $\square$

### 3 $\Delta = \Omega(\log^{-2/3} n)$ -separated sets

We now give an algorithm that given a  $c$ -spread  $\ell_2^2$  representation finds  $\Delta$ -separated sets of size  $\Omega(n)$  for  $\Delta = \Theta(1/\log^{2/3} n)$ . Our correctness proof assumes a key theorem (Theorem 5) whose proof appears in Section 4. The algorithm will be improved in Section 5 to allow  $\Delta = \Theta(1/\sqrt{\log n})$ .

The algorithm is given a  $c$ -spread  $\ell_2^2$ -representation. We select constants  $c', \sigma > 0$  depending on  $c$ .

SET-FIND:

Input: A  $c$ -spread unit vector representation  $v_1, v_2, \dots, v_n \in \mathfrak{R}^d$ .

Parameters: Desired separation  $\Delta$ , desired balance  $c'$ , and projection gap,  $\sigma$ .

Pick a random line  $u$  passing through the origin, and let

$$S_u = \{v_i : \langle v_i, u \rangle \geq \frac{\sigma}{\sqrt{d}}\},$$

$$T_u = \{v_i : \langle v_i, u \rangle \leq -\frac{\sigma}{\sqrt{d}}\}.$$

If  $|S_u| < 2c'n$  or  $|T_u| < 2c'n$ , HALT, else proceed as follows. Pick any  $v_i \in S_u, v_j \in T_u$  such that  $\|v_i - v_j\|^2 \leq \Delta$ , and delete  $i$  from  $S_u$  and  $j$  from  $T_u$ . Repeat until no such  $v_i, v_j$  can be found and output the remaining sets  $S, T$ .

REMARK 3 The procedure SET-FIND can be seen as a rounding procedure of sorts. It starts with a “fat” random hyperplane cut (cf. Goemans-Williamson [17]) to identify the sets  $S_u, T_u$  of vertices that project far apart. It then prunes these sets to find sets  $S, T$ .

Notice that if SET-FIND does not HALT prematurely, it returns a  $\Delta$ -separated pair of sets. Thus, we need to show that in SET-FIND often both  $S_u$  and  $T_u$  are larger than  $2c'n$  and that no more than  $c'n$  points are deleted from  $S_u$  and  $T_u$ . The first claim is relatively easy, and we show this in the next subsection. Analysing the deletion process is much harder and forms the bulk of the paper. We state the formal claims about the process in the following subsection, and proved it in Section 4.

#### 3.1 Projection and $S_u, T_u$

We first remind the reader that in  $\mathfrak{R}^d$ , the projection of any unit vector on a random direction is distributed essentially like a Gaussian with expectation 0 and standard deviation  $1/\sqrt{d}$ .

LEMMA 3 (GAUSSIAN BEHAVIOR OF PROJECTIONS)

If  $v$  is a vector of length  $\ell$  in  $\mathfrak{R}^d$  and  $u$  is a randomly chosen unit vector then (i) for  $x \leq 1$ ,  $\Pr[|\langle v, u \rangle| \leq \frac{x\ell}{\sqrt{d}}] \leq 3x$ . (ii) for  $x \leq \sqrt{d}/4$ ,  $\Pr[|\langle v, u \rangle| \geq \frac{x\ell}{\sqrt{d}}] \leq e^{-x^2/4}$ .

If the projection length in a particular direction  $u$  is  $t\ell/\sqrt{d}$ , we say that  $t$  is the *stretch* of  $v$  in direction  $u$ . (This definition is motivated by the fact that  $\ell/\sqrt{d}$  is the root mean square of the projection length of  $v$  in a random direction.) Lemma 3-(ii) implies that a vector  $v$  has stretch  $t$  in a random

direction  $u$  with probability at most  $e^{-t^2/4}$ . We will use this notion of stretch and this fact extensively in subsequent sections.

Now using part (i) of Lemma 3 and Goemans-Williamson, it is easy to prove that if the  $v_i$ 's are  $c$ -spread then with constant probability,  $S_u$  and  $T_u$  are large.

LEMMA 4

*For every positive  $c < 1/3$ , there are  $c', \sigma > 0$  such that the probability (over the choice of  $u$ ) is at least  $c/8$  that the sets  $S_u, T_u$  defined in SET-FIND- $(c', \sigma)$  have size at least  $2c'n$ .*

PROOF: Goemans-Williamson show that for any two points  $x, y$  on a unit sphere,

$$\Pr[\text{a random hyperplane separates } x, y] \geq .878 \frac{|x-y|^2}{4}.$$

By definition of  $c$ -spread, the sum of the distances between the points is at least  $c(1-c)n^2$ . Therefore the expected number of pairs that are separated by a random hyperplane is at least  $an^2$ , where  $a = .878c(1-c)$ . By Markov's bound the probability that the number of separated pairs is less than  $an^2/2$  is at most  $(1-a)/(1-a/2) \leq 1-a/2$ .

Since these  $an^2/2$  pairs of nodes are split by the hyperplane, there must be at least  $an/2$  nodes on the smaller side.

By Lemma 3-(i), the probability that the projection of a point on the unit sphere falls within  $\sigma/\sqrt{d}$  of the origin is at most  $3\sigma$ . By choosing  $\sigma$  appropriately, and applying the Markov bound, we can ensure that the probability that more than  $an/4$  points fall within  $\sigma/\sqrt{d}$  is at most  $a/4$ . Now by the union bound both  $S_u$  and  $T_u$  have at least  $a/4$  points with probability at least  $a/4$ .

The lemma follows by noting that  $a/4 > c/8$ , and choosing  $2c' = a/4$ .  $\square$

## 3.2 Number of deletions

To analyse the number of deletions, we note that any deleted pair  $v_i \in S_u, v_j \in T_u$  is such that the vector  $v_i - v_j$  has stretch  $t = 2\sigma/\sqrt{\Delta}$ , since its length was at most  $\sqrt{\Delta}$  and its projection length was at least  $2\sigma/\sqrt{d}$ . (Note: From now on, we will often say the pair  $v_i, v_j$  has a certain stretch when we mean  $v_i - v_j$ .) If  $\Delta$  were  $(16\sigma^2/\log n)$  then the analysis would be trivial since any deleted pair has stretch at least  $4\sqrt{\log n}$ ; this event occurs with probability less than  $e^{-4\log n} \ll 1/n^2$  by Lemma 3-(ii). Thus, we expect *no* pairs to be deleted. (Aside: This is an alternative version of Leighton-Rao.)

When  $\Delta = \Omega(\log^{-2/3} n)$ , it may be quite likely that many pairs are deleted. However, we observe that for a direction  $u$  the deleted pairs form a matching  $M_u$ . Moreover, if the procedure fails for a direction  $u$  the matching  $M_u$  is of size at least  $c'n$ . Thus if the procedure does not succeed with constant probability, we have large matchings  $M_u$  for most directions  $u$  where each matching edge has stretch  $2\sigma/\sqrt{\Delta}$ . We will show (Theorem 5) that this is impossible. Now we formalize the property of the matchings when SET-FIND often fails to produce a  $\Delta = 1/t^2$ -separated pair.

DEFINITION 4 (( $t, \gamma, \beta$ )-STRETCHED) *An  $\ell_2^2$  set of points  $v_1, v_2, \dots, v_n \in \mathfrak{R}^d$  are  $(t, \gamma, \beta)$ -stretched at scale  $l$  if for at least  $\gamma$  fraction of directions  $u$ , there is a (partial) matching  $M_u$  with  $\beta n$  disjoint pairs  $(i_1, j_1), (i_2, j_2), \dots$ , such that each  $i_m, j_m$  satisfies  $|v_{i_m} - v_{j_m}|^2 \leq l^2$  and  $\langle u, (v_{i_m} - v_{j_m}) \rangle \geq tl/\sqrt{d}$ . (In particular, pair  $v_{i_m}, v_{j_m}$  has stretch at least  $t$  in direction  $u$ .)*

THEOREM 5

*For any  $\gamma, \beta > 0$  there is a  $C = C(\gamma, \beta)$  such that if  $t > C(\log n)^{1/3}$  then a unit- $\ell_2^2$  representation cannot be  $(t, \gamma, \beta)$ -stretched for any scale  $l$ .*

Applying Theorem 5 with  $l = \sqrt{\Delta}$  and  $t = 2\sigma/\sqrt{\Delta}$  shows that there is some  $\Delta = O(\log^{-2/3} n)$ , such that the probability that SET-FIND removes a matching of size  $c'n$  is  $o(1)$ . We conclude that SET-FIND outputs  $S, T$  of size  $\geq c'n$  with probability  $\Omega(1)$ . This completes our analysis of SET-FIND.

## 4 Proof of Theorem 5

The main idea in the proof is to show that if the large matchings  $M_u$  mentioned in Definition 4 exist for most directions  $u$ , then for  $\Omega(1)$  fraction of directions we can string together  $r = \Omega(t)$  pairs from these matchings to produce a vector whose projection is  $\Omega(rtl/\sqrt{d})$ . Triangle inequality implies that any such vector has squared length at most  $rl^2$ , which means that the stretch is  $\sqrt{rt}$ .

Recall that for almost all directions, no pair of vectors has stretch more than  $4\sqrt{\log n}$ . Since the stringing together referred to above is possible for  $\Omega(1)$  directions, we conclude that the stretch  $O(\sqrt{rt})$  cannot exceed  $4\sqrt{\log n}$ , which proves that  $t$  can only be  $O(\log^{1/3} n)$ .

### 4.1 Matching covers

The definition of  $(t, \gamma, \beta)$ -stretched pointsets suggests that for many direction there are many disjoint pairs of points which are stretched. We will work with a related notion.

**DEFINITION 5** ( $(\epsilon, \delta)$ -MATCHING-COVERED POINT SET) *A set of points  $V \subseteq \mathfrak{R}^d$  is  $(\epsilon, \delta)$ -matching-covered at scale  $l$  if for every unit vector  $u \in \mathfrak{R}^d$ , there is a (partial) matching  $M_u$  of  $V$  such that every  $(v_i, v_j) \in M_u$  satisfies  $|v_i - v_j|^2 \leq l^2$  and  $|\langle u, v_i - v_j \rangle| \geq \epsilon$ , and for every  $i$ ,  $\mu(u : v_i \text{ matched in } M_u) \geq \delta$ . We refer to the set of matchings  $M_u$  to be the matching cover of  $V$ .*

**REMARK 4** The main difference from Definition 4 is that every point participates in  $M_u$  with constant probability for a random direction  $u$ .

**LEMMA 6**

*If a set of  $n$  vectors is  $(t, \gamma, \beta)$ -stretched at some scale  $l$ , then they contain a subset  $X$  of  $\Omega(n\gamma\beta)$  vectors that are  $(\epsilon, \delta)$ -matching covered at scale  $l$ , where  $\delta = \Omega(\beta\gamma)$ ,  $\epsilon \geq tl/\sqrt{d}$ , and for every pair  $v_i, v_j$  in the matching cover,  $|v_i - v_j|^2 \leq l^2$ .*

**PROOF:** Consider the multigraph consisting of the union of all partial matchings  $M_u$ 's as described in Definition 4. The average node is in  $M_u$  for  $\gamma\beta$  measure of directions. Remove all nodes that are matched on fewer than  $\gamma\beta/2$  measure of directions (and remove the corresponding matched edges from the  $M_u$ 's). Repeat. The aggregate measure of directions removed is  $\gamma\beta n/2$ . Thus at least  $\gamma\beta n/2$  aggregate measure on directions remains. This implies that there are at least  $\gamma\beta n/4$  nodes left, each matched in at least  $\gamma\beta/4$  measure of directions. This is the desired subset  $X$ .  $\square$

**NOTATION** From now on we restrict attention to the subset  $X$  mentioned in Lemma 6. Let  $H$  denote the multigraph on  $X$  formed by taking the union of all matchings  $M_u$  in the matching cover. For each  $v_i \in X$ , let  $\text{Ball}(v_i, r)$  denote the set of  $v_j$ 's whose distance from  $v_i$  in  $H$  is at most  $r$ .

We sometimes say that " $v_j$  is  $r$  matching hops from  $v_i$ ." Note that since the matching cover consists of edges of length  $\leq l$ , the triangle inequality for  $\ell_2^2$  representations implies that any such  $v_j, v_i$  satisfy  $|v_i - v_j|^2 \leq rl^2$ .

Now we define a related object.

**DEFINITION 6** ( $(\epsilon, \delta)$ -COVER) *A set  $\{w_1, w_2, \dots\}$  of vectors in  $\mathfrak{R}^d$  is an  $(\epsilon, \delta)$ -cover if every  $|w_j| \leq 1$  and for at least  $\delta$  fraction of unit vectors  $u \in \mathfrak{R}^d$ , there exists an  $j$  such that  $\langle u, w_j \rangle \geq \epsilon$ .*

**REMARK 5** Since  $-u$  is also a random unit vector, the probability is also  $\geq \delta$  that there is a  $w_j$  that  $\langle u, w_j \rangle \leq -\epsilon$ . This will be important later in Lemma 10.

**REMARK 6** Whenever we study these covers, we have a fixed  $v_i \in X$  in mind and the vectors in the cover are of the form  $v_j - v_i$ . In such a case, we say that  $v_i$  is "centrally  $(\epsilon, \delta)$ -covered" by the  $v_j$ 's in question.

Note that if  $v_i \in X$ , then the vectors  $v_j - v_i$  for  $v_j \in \text{Ball}(v_i, 1)$  form an  $(\epsilon, \delta)$ -cover. (The converse is not true: taking the union of such  $(\epsilon, \delta)$ -covers may not always give a matching cover.)

NOTATION: Let  $S_r \subseteq X$  consist of all  $v_i \in X$  such that the vectors  $\{v_j - v_i : v_j \in \text{Ball}(v_i, r)\}$  form an  $(\epsilon r/2, 1 - \delta/2)$ -cover.

Note that thus far there is no reason to believe even that  $S_1$  is nonempty, since we only know that for each  $v_i \in X$ , the set  $\{v_j - v_i : v_j \in \text{Ball}(v_i, 1)\}$  is an  $(\epsilon, \delta)$ -cover, whereas in order for  $v_i$  to be in  $S_1$  these vectors must form an  $(\epsilon/2, 1 - \delta/2)$ -cover.

Thus the main technical step is the following. It assumes that stretch of the edges in the matching cover, namely,  $t = \epsilon\sqrt{d}/l$ , is larger than some fixed constant.

LEMMA 7 (MAIN)

(i)  $S_1 = X$ . (ii) There are constants  $\eta = \eta(\delta), \rho = \rho(\delta)$  such that for  $r \leq \eta t$ , we have  $|S_{r+1}| \geq \rho |S_r|$ .

Theorem 5 follows immediately from Lemma 7.

PROOF:(Theorem 5) If the hypothesis of Theorem 5 is true, then Lemma 6 implies the existence of a set  $X$  of  $\Omega(n)$  vectors  $v_i$ 's that form an  $(\epsilon, \delta)$ -matching covered point set using edges of squared length at most  $l^2$ . Here  $\epsilon = tl/\sqrt{d}$ . Then Lemma 7 and a simple induction implies that for  $r = \eta t$ ,

$$|S_r| \geq \rho^{r-1} |S| = \Omega(\rho^{r-1} n) \gg 1,$$

where we're using the fact that  $r = o(\log n)$ . Thus  $S_r$  is nonempty.

Let  $v_i \in S_r$ . Then for at least  $1 - \delta/2$  fraction of directions  $u$ , some  $v_j \in \text{Ball}(v_i, r)$  satisfies  $|\langle v_i - v_j, u \rangle| \geq r\epsilon/2$ . However,  $|v_j - v_i| \leq \sqrt{r}l$ , so we conclude that the stretch of  $v_j - v_i$  is

$$\frac{r\epsilon/2 \times \sqrt{d}}{\sqrt{r}l} = \sqrt{r}t/2 = \sqrt{\eta}t^{3/2}/2 = \Omega(t^{3/2}).$$

But recall that for any set of  $n$  vectors, at most  $1/n$  of the directions  $u$  are such that one of the  $\binom{n}{2}$  pairs of vectors has stretch  $> 4\sqrt{\log n}$ . But since  $S_r$  is nonempty, we know that the probability is at least  $1 - \delta/2$  that some stretch exceeds  $\Omega(t^{3/2})$ . We conclude that  $t = O(\log^{1/3} n)$ .  $\square$

## 4.2 Proving Lemma 7

We prove Lemma 7 by induction. Recall that it was unclear even that  $S_1$  is nonempty. In fact a phenomenon called *measure concentration* implies that  $S_1 = X$ . We first introduce this idea.

### 4.2.1 Measure concentration.

Let  $S^{d-1}$  denote the surface of the unit ball in  $\mathfrak{R}^d$  and let  $\mu(\cdot)$  denote the standard measure on it. For any set of points  $A$ , we denote by  $A_\gamma$  the  $\gamma$ -neighborhood of  $A$ , namely, the set of all points that have distance at most  $\gamma$  to some point in  $A$ .

LEMMA 8 (CONCENTRATION OF MEASURE)

If  $A \subseteq S^{d-1}$  is measurable and  $\gamma > \frac{2\sqrt{\log(1/\mu(A))+t}}{\sqrt{d}}$ , where  $t > 0$ , then  $\mu(A_\gamma) \geq 1 - \exp(-t^2/2)$ .

PROOF: P. Levy's isoperimetric inequality ([5]) states that  $\mu(A_\gamma)/\mu(A)$  is minimized for spherical caps<sup>2</sup>. The lemma now follows by a simple calculation using the standard formula for (d-1)-dimensional volume of spherical caps, which says that the cap of points whose distance is at least  $s/\sqrt{d}$  from an equatorial plane is  $\exp(-s^2/2)$ .  $\square$

The following Lemma is an immediate corollary.

<sup>2</sup>Levy's isoperimetric inequality is not trivial; see [29] for a sketch. However, results qualitatively the same—but with worse constants—as Lemma 8 can be derived from the more elementary Brunn-Minkowski inequality; this “approximate isoperimetric inequality” of Ball, de Arias and Villa also appears in [29].

LEMMA 9

Let  $\{v_1, v_2, \dots\}$  be a finite set of vectors that is an  $(\epsilon, \delta)$ -cover, and  $|v_i| \leq \ell$ . Then, for any  $\gamma > \frac{\sqrt{2 \log(2/\delta) + t}}{\sqrt{d}}$ , the vectors are also a  $(\epsilon - 2\ell\gamma, \delta')$ -cover, where  $\delta' = 1 - \exp(-t^2/2)$ .

PROOF: Let  $A$  denote the set of directions  $u$  for which there is an  $i$  such that  $\langle u, v_i \rangle \geq \epsilon$ . Since  $|v_i - u|^2 = 1 + |v_i|^2 - 2\langle u, v_i \rangle$  we also have:

$$A = S^{d-1} \cap \bigcup_i \text{Ball}\left(v_i, \sqrt{1 + |v_i|^2 - 2\epsilon}\right),$$

which also shows that  $A$  is measurable. Thus by Lemma 8,  $\mu(A_\gamma) \geq 1 - \exp(-t^2/2)$ .

We argue that for each direction  $u$  in  $A_\gamma$ , there is a vector  $v_i$  in the  $(\epsilon, \delta)$  cover with  $\langle v_i, u \rangle \geq \epsilon - 2\ell\gamma$  as follows. Let  $u \in A$ ,  $u' \in A_\gamma \cap S^{d-1}$  be such that  $|u - u'| \leq \gamma$ .

The projection length of  $v_i$  on  $u$  is  $|v_i| \cos \theta$  where  $\theta$  is the angle between  $v_i$  and  $u$ . The projection length of  $v_i$  on  $u'$  is  $|v_i| \cos \theta'$  where  $\theta'$  is the angle between  $v_i$  and  $u'$ . The angle formed by  $u$  and  $u'$  is at most  $2\gamma$  since  $\alpha \leq 2 \sin \alpha$ , so  $\theta - 2\gamma \leq \theta' \leq \theta + 2\gamma$ . Since the absolute value of the slope of the cosine function is at most 1, we can conclude that the differences in projection is at most  $2\gamma|v_i| \leq 2\gamma\ell$ . That is,  $\langle v_i, u' \rangle \geq \epsilon - 2\gamma\ell$ .

Combined with the lower bound on the  $\mu(A_\gamma)$ , we conclude that the set of directions  $u'$  such that there is an  $i$  such that  $\langle u', v_i \rangle \geq \epsilon - 2\ell\gamma$  has measure at least  $1 - \exp(-t^2/2)$ .  $\square$

We can thus use Lemma 9 to boost  $\delta$  to almost 1. However, we choose to do it only sometimes, not always. This explains the slightly strange hypothesis in the next Lemma. The proof of Lemma 7 will not use this lemma *per se* but will use the argument in it.

LEMMA 10

If  $\{w_1, w_2, \dots, w_k\} \subseteq \mathfrak{X}^d$  is an  $(\epsilon_1, 1 - \delta_1)$ -cover and  $\{w'_1, w'_2, \dots, w'_l\}$  is an  $(\epsilon_2, \delta_2)$ -cover then the set  $\{w_e - w'_f : 1 \leq e \leq k, 1 \leq f \leq l\}$  is a  $(\epsilon_1 + \epsilon_2, \delta_2 - \delta_1)$ -cover.

PROOF: Let  $u \in \mathfrak{X}^d$  be a random unit vector. The probability is at least  $1 - \delta_1$  that there is a  $w_e$  such that  $\langle u, w_e \rangle \geq \epsilon_1$ . The probability is at least  $\delta_2$  that there is a  $w'_f$  such that  $\langle u, w'_f \rangle \leq -\epsilon_2$ . Thus with probability at least  $\delta_2 - \delta_1$ , there exist  $w_e, w'_f$  such that  $\langle u, w_e - w'_f \rangle \geq \epsilon_1 + \epsilon_2$ .  $\square$

#### 4.2.2 Proof of Lemma 7

Let  $\sigma = \epsilon\sqrt{d}$ , and assume the stretch  $t = \sigma/l$  of the matching edges is larger than any desired constant.

We set  $D(\sigma, \delta) = 8\sqrt{2 \log(2/\delta)}/\sigma$ , and  $\rho(\delta) = \delta/4$ .

First, we show  $S_1 = X$ . The hypothesis implies that every  $v_i \in X$  is centrally  $(\epsilon, \delta)$ -covered by the set of  $v_j \in \text{Ball}(v_i, 1)$ . We apply Lemma 9 to each of these  $(\epsilon, \delta)$ -covers with  $\gamma = \sigma/4l\sqrt{d}$ . Note that

$$\gamma = \sigma/4l\sqrt{d} > \sqrt{2 \log(2/\delta)}/\sqrt{d} + \sqrt{2 \log(2/\delta)}/\sqrt{d},$$

so we conclude that  $v_i$  is also  $(1 - \delta/2, \epsilon - 2\gamma\ell)$  covered by  $\text{Ball}(v_i, 1)$ . Since  $2\gamma\ell < \sigma/2\sqrt{d} \leq \epsilon/2$ , we have thus shown that every  $v_i \in X$  is also in  $S_1$ .

Assume the induction has worked for  $r$  steps and there is a set  $S_r \subseteq T$  satisfying  $|S_r| \geq \rho^{r-1}|X|$  such that every point  $v_i \in S_r$  is centrally  $(\epsilon_r, 1 - \delta_0/2)$ -covered by the vectors in  $\text{Ball}(v_i, r)$ , where  $\epsilon_r \geq 0.5r\epsilon$ .

For each  $v_i \in S_r$  consider the set of all vectors  $v_j - v_k$  where  $v_j \in \text{Ball}(v_i, r)$  and  $v_k \in \text{Ball}(v_i, 1)$ . Lemma 10 implies that these vectors form an  $(\epsilon_r + \epsilon, \delta/2)$  cover, but unfortunately this is no longer centered at  $v_i$ . Thus we are unable to prove in general that  $v_i \in S_{r+1}$ .

Instead, we argue differently and use an averaging argument to say that if  $S_r$  is large, so is  $S_{r+1}$ . Let  $v_i \in S_r$ . For  $1 - \delta/2$  fraction of directions  $u$ , there is a point  $v_j \in \text{Ball}(v_i, r)$  such that  $\langle v_j - v_i, u \rangle \geq \epsilon_r$ .

Also for  $\delta$  fraction of directions  $u$ , there is a point in  $v_k \in \text{Ball}(v_i, 1)$  such that  $\langle v_k - v_i, u \rangle \leq -\epsilon$  and  $v_k$  is matched to  $v_i$  in the matching  $M_u$ . Thus for a  $\delta/2$  fraction of directions  $u$ , both events happen and thus the pair  $(v_j, v_k)$  satisfies  $\langle v_j - v_k, u \rangle \geq \epsilon_r + \epsilon$ . Since  $v_j \in \text{Ball}(v_k, r + 1)$ , we “assign” this vector  $v_j - v_k$  to point  $v_k$  for direction  $u$ , as a step towards building a cover centered at  $v_k$ . Now we argue that for many  $v_k$ 's, the vectors assigned to it in this way form a  $(\epsilon_r + \epsilon, \rho^r \delta/2)$ -cover.

For each point  $v_i \in S_r$ , for  $\delta/2$  fraction of the directions  $u$  the process above assigns a vector to a point in  $X$  for direction  $u$  according to the matching  $M_u$ . Thus on average for each direction  $u$ , at least  $\delta|S_r|/2$  vectors get assigned it by the process. Equivalently, for a random point in  $X$ , the expected measure of directions for which the point is assigned a vector is at least  $\delta|S_r|/2|X|$ . Furthermore, at most one vector is assigned to any point for a given direction  $u$  (since the assignment is governed by the matching  $M_u$ ). Therefore at least  $\delta|S_r|/4|X|$  fraction of the points in  $X$  must be assigned a vector for  $\delta|S_r|/4|X|$  fraction of the directions.

We will define all such points of  $X$  to be the set  $S_{r+1}$  and note that for  $\rho = \delta/4$  the size is at least  $\rho|S_r|$  as required. However, we have to show another property for  $S_{r+1}$ . Thus far, since  $\delta|S_r|/4|X| = \rho^r$ , we have only shown that each point  $v_k$  in  $S_{r+1}$  is centrally  $(\epsilon_r + \epsilon, \delta\rho^r)$ -covered by  $v_j \in \text{Ball}(v_k, r + 1)$ .

We now invoke measure concentration to show that these centered covers are also centered  $(\epsilon_r + \epsilon/2, 1 - \delta/2)$  covers, so long as  $r = O(\sigma/l)$ . Note that the vectors in the cover have squared length at most  $rl^2$  due to the triangle inequality on the squared lengths. We apply Lemma 9 with  $\ell = \sqrt{rl}$  and

$$\gamma = \epsilon/4\ell = \sigma/4\sqrt{d}\sqrt{rl}.$$

Now, we need

$$\gamma = \frac{\sigma}{4\sqrt{d}\sqrt{rl}} > \frac{2\sqrt{\log \frac{2}{\rho^r}} + \sqrt{\log(\frac{2}{\delta})}}{\sqrt{d}}, \quad (11)$$

to get that  $v_k$  is centrally  $(\epsilon_r + \epsilon - 2\gamma\ell, 1 - \delta/2)$  covered. The condition is satisfied when  $r \leq \sigma/8l\sqrt{2\log(8/\delta)}$ , since  $\rho = \delta/4$ .

By noting that  $2\gamma\ell < \epsilon/2$ , we now observe that each  $v_k \in S_{r+1}$  is  $((r + 1)\epsilon/2, 1 - \delta/2)$ -covered by  $v_j \in \text{Ball}(v_k, r + 1)$  and our induction is complete.  $\square$

Now we state a corollary of the proof of Lemma 7 which will be useful in Section 5. As in Lemma 7, let  $X \subseteq \mathbb{R}^d$  be a pointset that is  $(\epsilon, \delta)$ -matching-covered using edges of squared length at most  $l^2$ .

The Corollary concerns, for some  $s > 0$ , a subset  $T \subseteq X$  of size at least  $|X|/2$  containing every  $v_i$  such that the set  $\{v_j - v_i : |v_i - v_j|^2 \leq s\}$  is an  $(\epsilon_1, 1 - \delta/2)$ -cover. Define  $T'$  to be the set of  $v_i$ 's such that the set  $\{v_j - v_i : |v_i - v_j|^2 \leq s + l^2\}$  is an  $(\epsilon_1 + \epsilon/2, 1 - \delta/2)$ -cover. The corollary assumes  $\epsilon\sqrt{d}$  is some constant, say  $\sigma$ .

**COROLLARY 11 (COVER COMPOSITION)**

*There are constants  $\rho, f$  depending only on  $\sigma, \delta$  such that if  $s + l^2 \leq f$ , then  $|T'| \geq \rho|X|$ .*

**PROOF:** Straightforward from proof of Lemma 7; left to the reader.  $\square$

## 5 Achieving $\Delta = \Omega(1/\sqrt{\log n})$ .

To prove Theorem 1 with  $\Delta = \Omega(1/\sqrt{\log n})$ , we start by invoking SET-FIND with that  $\Delta$  as separation parameter. If SET-FIND succeeds, we are done. Otherwise, as before we end up with matchings  $M_u$  in most directions  $u$ . Now, however, we cannot necessarily show that this leads to a contradiction. The bottleneck in our previous proof lies in the induction of Lemma 7, where the size of the set  $S_r$  decreases geometrically with  $r$  (see the calculation (11)). To bypass that, we describe another algorithm that finds a  $\Delta$ -separated set. This uses the simple observation that if a point is well covered then all points that are close to it are also well covered. Now we formalize this.

LEMMA 12 (COVERING CLOSE POINTS)

Suppose  $v_1, v_2, \dots \in \mathfrak{X}^d$  are vectors such that for some point  $v_0 \in \mathfrak{X}^d$  the vectors  $v_1 - v_0, v_2 - v_0, \dots$ , form an  $(\epsilon, \delta)$ -cover. Then for every point  $v'_0$  such that  $|v_0 - v'_0| = s$ , the vectors  $v_1 - v'_0, v_2 - v'_0, \dots$  form an  $(\epsilon - \frac{ts}{\sqrt{d}}, \delta - e^{-t^2/4})$ -cover.

PROOF: If  $u$  is a random unit vector,  $\Pr_u[\langle u, v_0 - v'_0 \rangle \geq \frac{ts}{\sqrt{d}}] \leq e^{-t^2/4}$ .  $\square$

Armed with this lemma, we construct  $l^2$ -separated sets as follows. Our algorithm is very much like the inductive step in the proof lemma 7.

For each  $r$ , define  $S'_r$  to be the set of points  $v$  which is  $(r\epsilon/4, 1 - 3\delta/4)$ -covered by points that are within length  $\ell_r = \sqrt{2r}/l$  of  $v$ . Consider the smallest  $r$  such that  $S'_{r+1}$  has cardinality less than  $n/2$  (we argue below that such an  $r$  exists). We apply Corollary 11 to  $S'_r$ , to get a set of  $\delta n/4$  points  $T$  where each point  $v \in T$  is  $((r/4 + 1/2)\epsilon, 1 - \delta/2)$  covered by points that are within squared Euclidean length  $\ell_r^2 + l^2$  of  $v$ . It will follow using Lemma 12 that all points within length  $l$  of  $T$  are in  $S_{r+1}$ , and therefore the  $\Omega(n)$  sized sets  $T, S_1 - S_{r+1}$  are at least  $l^2$ -separated.

To argue correctness, we will show that for every  $r \leq r_0$  where  $r_0$  is  $\theta(1/l^2)$ , Corollary 11 implies  $|T| \geq \rho n$  and Lemma 12 implies that the  $l$ -neighborhood of  $T$  is contained in  $S'_{r+1}$ . Now for some  $l$  which is  $\theta(1/\log^{1/4} n)$  we argue that  $|S'_{r_0}| = 0$ . Recall that  $v \in S'_{r_0}$  only if for most directions it participates in a stretched pair of points with stretch  $\Omega(\sqrt{r_0}/l) = \Omega(1/l^2) = \Omega(\sqrt{\log n})$ . In fact, for most directions there are no such stretched pairs.

To finish we argue that  $r_0$  is  $\theta(1/l^2)$ . To use Lemma 12 as above, we need to ensure that the loss in projection  $ts/\sqrt{d}$  (in our context  $s = l$ ) is at most  $\epsilon/4$ , and that the loss in probability  $e^{-t^2/4}$  is at most  $\delta/4$ . Choosing  $t = 2\sqrt{\log 4/\delta}$ , we see that  $l$  must be less than  $\epsilon\sqrt{d}/t$  which is easily satisfied for sufficiently large  $n$ . Moreover, for some constant  $f$ , Corollary 11 can be applied as long as  $\ell_r^2 = 2rl^2 \leq f$ . This implies that the upper limit for  $r$  is  $\theta(1/l^2)$ .

REMARK 7 Simplifying slightly, here is a concrete algorithm.

For each  $r$ , find a set  $\tilde{S}'_r$  where each point in  $\tilde{S}'_r$  is approximately  $(r\epsilon/4, 1 - 3\delta/4)$ -covered by points that are within length  $\sqrt{2r}/l$ . This can be done by sampling  $O(\log n)$  directions. For the first  $\tilde{S}'_r$  with cardinality less than  $|S'_1|/2$ , we take the set  $\tilde{S}'_1 - \tilde{S}'_r$  to be  $S$  and take  $T$  to be all the points that are at least  $l^2$  from  $S$ .

The sets  $S$  and  $T$  are  $l^2$ -separated by construction and the argument above shows that each set is large in the event SET-FIND usually fails.

## 6 $O(\sqrt{\log n})$ Ratio for SPARSEST CUT

Now we describe a rounding technique for the SDP in (8)–(10) that gives a  $O(\sqrt{\log n})$ -approximation to SPARSEST CUT. Note that our results on expander flows in Section 7 given an alternative  $O(\sqrt{\log n})$ -approximation algorithm.

First we see in what sense the SDP in (8)–(10) is a relaxation for SPARSEST CUT. For any cut  $(S, \bar{S})$  consider a vector representation that places all nodes in  $S$  at one point of the sphere of squared radius  $(4|S| \binom{|S|}{|S|})^{-1}$  and all nodes in  $|\bar{S}|$  at the diametrically opposite point. It is easy to verify that this solution is feasible and has value  $|E(S, \bar{S})| / |S| \binom{|S|}{|S|}$ . Since  $\binom{|S|}{|S|} \in [n/2, n]$ , we can treat it as a scaling factor. We conclude that the *optimal* value of the SDP is a lower bound (up to scaling by  $n$ ) for SPARSEST CUT.

The next theorem implies that the integrality gap is  $O(\sqrt{\log n})$ .

THEOREM 13

There is a polynomial-time algorithm that, given a feasible SDP solution with value  $\beta$ , produces a cut  $(S, \bar{S})$  satisfying  $|E(S, \bar{S})| = O(\beta |S| n \sqrt{\log n})$ .

The proof divides into two cases, one of which is similar to that of Theorem 1. The other case is dealt with the following Lemma.

LEMMA 14

There is a polynomial-time algorithm for the following task. Given any feasible SDP solution with  $\beta = \sum_{\{i,j\} \in E} |v_i - v_j|^2$ , and a node  $k$  such that the geometric ball of squared-radius  $1/4n^2$  around  $v_k$  contains at least  $n/2$  vectors, the algorithm finds a cut  $(S, \bar{S})$  with expansion at most  $O(\beta n)$ .

PROOF: Let  $X$  be the subset of nodes that correspond to the vectors in the geometric ball around  $v_k$ . Let  $d(i, j) = |v_i - v_j|^2$  and when  $\{i, j\}$  is an edge  $e$  we write  $d(e)$ .

Since  $\sum_{i < j} d(i, j) = 1$ , the triangle inequality implies that node  $k$  also satisfies:

$$\sum_j d(k, j) \geq 1/2n.$$

Separating terms corresponding to  $j \in X$  and  $j \notin X$  we obtain

$$\sum_{j \notin X} d(j, k) \geq \frac{1}{2n} - \frac{1}{4n^2} \cdot \frac{n}{2} = \frac{3}{8n}. \quad (12)$$

The Lemma's hypothesis also says

$$\sum_{e=\{i,j\}} d(e) \leq \beta. \quad (13)$$

Now, we extend a tree from  $X$  to the rest of the graph by growing along outgoing edges at a uniform rate with respect to  $d(e)$ . (In other words, do breadth-first search on the weighted graph.) Let  $\alpha_{obs}$  to the minimum expansion of any cut induced by this growing boundary. (Notice, the smaller side of the cut lies outside the boundary since  $|X| \geq n/2$ .) In other words, after growing  $\delta$  from the set  $X$  the cut size is at least  $\alpha_{obs} n(\delta)$  where  $n(\delta)$  is the number of nodes we have yet to reach.

This implies that the total edge weight seen in this process is at least the integral of  $\alpha_{obs} n(\delta)$ . From equation (13) we obtain:

$$\beta \geq \int_{\delta > 0}^1 \alpha_{obs} n(\delta) d\delta.$$

Furthermore, we note that the integral of  $n(\delta)$  with respect to  $\delta$  is at least  $\sum_{i \notin X} d(u, i)$ . Thus, from equation (12) we have that

$$\int_{\delta > 0} n(\delta) d\delta \geq \frac{3}{8n}.$$

Combining the above inequalities, we get  $\beta \geq 3\alpha_{obs}/8n$ , or, in other words  $\alpha_{obs} = O(\beta n)$ .

□

Thus we only need to consider the case where the hypothesis of Lemma 14 does not hold for any  $k$ . Namely, for each node  $k$ , less than  $n/2$  vectors lie within a ball of squared radius less than  $1/4n^2$ . That is, the nodes are well spread out. Under this condition, the ideas from Corollary 2 and Section 5 can be used to produce  $c'$ -balanced cuts (where  $c'$  is some constant) of expansion  $O(\beta n)$ , thus showing that the integrality gap is  $O(\sqrt{\log n})$ . Now we sketch how this is done.

First, scale all vectors by  $2n$  so that the squared-length of  $1/4n^2$  becomes 1 and  $\sum_{i < j} |v_i - v_j|^2 = 4n^2$ . Now any sphere of radius 2 contains at most  $1/2$  the points. Furthermore, averaging shows that at least  $9/10$  fraction of points lie inside a sphere of radius 40. Thus  $\Omega(1)$  fraction of nodes lie in a spherical annulus of inner radius 1 and outer radius 40. A version of Theorem 1 applies to such representations with constants appropriately modified for the diameters of the ball. As stated, the Theorem assumed the vector representation of the graph involves unit vectors, but looking over the proofs it is clear that the proofs go through (with the constants not as good) if  $9/10$  of the vectors have

length lowerbounded and upperbounded by some constant. The reason is that by using the Goemans-Williamson analysis as before, we conclude that for some  $c', \sigma$ , the algorithm SET-FIND- $(c', \sigma)$  outputs a  $c'$ -balanced cut with probability  $\Omega(1)$ . Then, the remainder of the proof just uses an upper bound on the vector length in various places.

## 7 Expander Flows: Approximate certificates of expansion

Deciding whether a given graph  $G$  has expansion at least  $\alpha$  is coNP-complete [6] and thus has no short certificate unless the polynomial hierarchy collapses. Jerrum and Sinclair [21] and then Leighton and Rao [22] showed how to use multicommodity flows to give “approximate” certificates; this technique was then clarified by Sinclair [32] and Diaconis and Saloff-Coste [11]. The results of this section represent a continuation of that work but with a better certificate: for any graph with  $\alpha(G) = \alpha$  we can exhibit a certificate to the effect that the expansion is at least  $\Omega(\alpha/\sqrt{\log n})$ . Furthermore, this certificate can be computed in polynomial time. The certificate involves multicommodity flows that are (weighted) expander graphs. (The algorithms of this paper were originally discovered in this setting.)

A word on convention. Weighted graphs in this section will be symmetric, i.e.,  $w_{ij} = w_{ji}$  for all node pairs  $i, j$ . We call  $\sum_j w_{ij}$  the *degree* of node  $i$ . We say that a weighted graph is  $d$ -regular if all degrees are exactly  $d$ . We emphasize that  $d$  can be a fraction.

### 7.1 Multicommodity flows as graph embeddings

A *multicommodity flow* in an unweighted graph  $G = (V, E)$  is an assignment of a *demand*  $f_{ij} \geq 0$  to each node pair  $i, j$  such that we can route  $f_{ij}$  units of flow from  $i$  to  $j$ , and can do this simultaneously for all pairs while satisfying capacity constraints. More formally, for each  $i, j$  and each path  $p \in \mathcal{P}_{ij}$  there exists  $f_p \geq 0$  such that

$$\forall i, j \in V \quad \sum_{p \in \mathcal{P}_{ij}} f_p = f_{ij} \tag{14}$$

$$\forall e \in E \quad \sum_{p \ni e} f_p \leq 1. \tag{15}$$

Note that every multicommodity flow in  $G$  can be viewed as an *embedding* of a weighted graph  $G' = (V, E', f_{ij})$  on the same vertex set such that the weight of edge  $\{i, j\}$  is  $f_{ij}$ . We assume the multicommodity flow is symmetric, i.e.,  $f_{ij} = f_{ji}$ . The following inequality is trivial.

$$\alpha(G) \geq \alpha(G') \tag{16}$$

The following is one way to look at the Leighton-Rao result where  $K_n$  is the complete graph on  $n$  nodes. The embedding mentioned in the theorem is, by (16), a certificate showing that expansion is  $\Omega(\alpha/\log n)$ .

**THEOREM 15 (LEIGHTON-RAO [22])**

*If  $G$  is any  $n$ -node constant-degree graph with  $\alpha(G) = \alpha$ , then it is possible to embed a scaled version of  $K_n$  in it with each  $f_{ij} = \alpha/n \log n$ .*

**REMARK 8** The same theorem is usually stated using all  $f_{ij} = 1$  (i.e., an unweighted copy of  $K_n$ ) and then bumping up the capacity of each edge of  $G$  to  $O(n \log n/\alpha)$ . A similar restatement is possible for our theorem about expander flows (Theorem 17).

We note that the embedding of Theorem 15 can be found in polynomial time using a multicommodity flow computation, and that this embedding is a “certificate” via (16) that  $\alpha(G) = \Omega(\alpha/\log n)$ .

## 7.2 Expanders and Expander Flows

*Expanders* can be defined in more than one way. We use the following definition.

**DEFINITION 7 (EXPANDERS)** *For any  $c > 0$ , a  $d$ -uniform weighted graph  $(w_{ij})$  is a  $\beta$ -expander if for every set of nodes  $S$ ,  $w(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} w_{ij}$  is at least  $\beta d |S|$ .*

If a weighted graph is a  $\beta$ -expander, then the second eigenvalue of its Laplacian lies in the interval  $[\beta^2/2, 2\beta]$ . Usually the word “expander” is reserved for the case when  $\beta > 0$  is a fixed constant independent of the graph size, in which case the second eigenvalue is bigger than some fixed positive constant. Thus an expander family can be recognized in polynomial time by an eigenvalue computation, and this motivates our definition of expander flows.

An *expander flow* is a multicommodity flow that is a  $\beta$ -expander for some constant  $\beta$ . Such flows can be used to certify the expansion of a graph.

**LEMMA 16**

*If a graph  $G$  contains a multicommodity flow  $(f_{ij})$  that is  $d$ -regular and is a  $\beta$ -expander, then  $\alpha(G) \geq \beta d$ .*

**PROOF:** For every set  $S$  of nodes, the amount of flow leaving  $S$  is at least  $\beta d$ , and hence this is a lowerbound on the number of edges leaving  $S$ .  $\square$

The previous Lemma is most useful if  $\beta$  is constant, in which case an eigenvalue computation can be used to verify that the given flow is indeed an expander. Thus a  $d$ -uniform expander flow may be viewed as a certificate that the expansion is  $\Omega(d)$ . The following Theorem says that such flows exist for  $d = \alpha/\sqrt{\log n}$  in every graph  $G$  satisfying  $\alpha(G) = \alpha$ . This is an interesting structural property of every graph, in the same spirit as Theorem 15. It also yields a certificate that  $\alpha(G) = \Omega(\alpha/\sqrt{\log n})$ .

**THEOREM 17**

*There is a constant  $\beta > 0$  such that the following is true for every unweighted graph  $G$ . If  $\alpha(G) = \alpha$ , then  $G$  contains a  $d$ -regular multicommodity flow in it that is a  $\beta$ -expander, where  $d = \alpha/\sqrt{\log n}$ .*

We note that the flow of Theorem 17 is such that all nonzero eigenvalues of its Laplacian are at least  $\lambda_0 = \beta^2/2$ . Thus the flow can be found in polynomial time using an Ellipsoid-like method that, for any, can find the largest  $d$  such that the graph contains a  $d$ -uniform multicommodity flow whose each nonzero eigenvalue is at least  $\lambda_0$ . (See Theorem ??.)

**EXAMPLE 1** We describe a few examples of expander flows. (i)  $G =$  the  $n$ -cycle. The eigenvalue bound is known to be quadratic in the conductance [9]. However, an expander flow is obtained by taking any 3-regular  $n$ -node expander (e.g. from [25]) and setting  $f_{ij} = 1/n$  for each expander edge  $\{i, j\}$ . Clearly, this flow is routable in the  $n$ -cycle since the demand crossing each cut is at most 2. Furthermore in the expander every set  $S$  has  $\Omega(|S|)$  edges leaving it, which corresponds to a flow of  $\Omega(|S|/n)$ . Thus the flow certifies that the expansion of the cycle is  $\Omega(1/n)$ . (ii)  $G =$  the cube connected cycle on  $\{0, 1\}^k$  (namely, the hypercube in which each node is replaced by a  $k$ -cycle.) Embed a random  $n$ -node bounded degree graph for  $n = 2^k$ . If  $\{i, j\}$  is an edge then route  $1/k$  units of flow from  $i$  to  $j$  in the cube-connected cycle using random paths from  $i$  to  $j$ . Simple random graph arguments show that flow capacity constraints are satisfied whp. Thus the flow certifies that the expansion is  $\Omega(1/k)$ . (iii) The graph from (ii), but with the nodes in each  $k$ -cycle permuted arbitrarily. (So this is really a family of graphs, since there are  $(k-1)!$  choices for the permutation at each of the  $2^k$  nodes.) Here we do not know of an explicit expander flow except by the general argument of Theorem 17. We also do not know if the expansion of every graph in the family is always  $\Omega(1/k)$ .

The proof of Theorem 17 uses the following lemma that is a consequence of the von Neumann min-max theorem<sup>3</sup>.

<sup>3</sup>The existence of expander flows can be proved in other ways including SDP duality or plain LP duality, but we think the proof using the min-max theorem is more intuitive.

LEMMA 18

Let  $c < 1/4$ . Let  $G = (V, E)$  be any graph and  $d > 0$  be any number. Let

$$\mathcal{V} = \{(v_1, \dots, v_n) : v_i \text{'s are a unit-}\ell_2^2\text{-representation that is } c\text{-spread}\} \quad (17)$$

$$\mathcal{F} = \{(f_{ij}) : (f_{ij}) \text{ is a } d\text{-regular flow that can be routed in } G\} \quad (18)$$

Then

$$\min_{(v_1, \dots, v_n) \in \mathcal{V}} \max_{(f_{ij}) \in \mathcal{F}} \sum_{ij} f_{ij} |v_i - v_j|^2 = \max_{(f_{ij}) \in \mathcal{F}} \min_{(v_1, \dots, v_n) \in \mathcal{V}} \sum_{ij} f_{ij} |v_i - v_j|^2 \quad (19)$$

PROOF: Consider a zero-sum two-player game in which the players' moves are chosen from  $\mathcal{F}$  and  $\mathcal{V}$  respectively and the payoff from the vector player to the flow player is  $\sum_{ij} f_{ij} |v_i - v_j|^2$ . Clearly, the strategy set  $\mathcal{F}$  is convex. Less obviously, so is  $\mathcal{V}$  once we represent it appropriately. A unit- $\ell_2^2$ -representation  $v_1, v_2, \dots, v_n$  may be represented using the Gram matrix of the  $v_i$ 's, namely,  $M_{ij} = \langle v_i, v_j \rangle$ . Thus  $|v_i - v_j|^2 = 2(1 - M_{ij})$ , which shows that  $\mathcal{V}$  consists of positive semidefinite matrices satisfying some linear constraints; this is a convex set.

Thus the payoff function  $\pi$  may be rewritten as

$$\pi((f_{ij}), M) = 2 \sum_{ij} f_{ij} (1 - M_{ij}),$$

which shows that it is linear in each of its two arguments (see Section B in the Appendix). Both strategy sets are bounded convex sets, and thus are finitely approximable at all scales. Then the Lemma follows from Lemma 29.  $\square$

As in many situations involving the minmax theorem, one side of (19) is easier to reason about than the other. The next Lemma shows that the left hand side is  $\Omega(n\alpha(G)/\sqrt{\log n})$ . Its proof appears in Section 7.2.1.

LEMMA 19

Let  $c < 1/4$ . If graph  $G = (V, E)$  satisfies  $\alpha(G) = \alpha$  then for every unit  $\ell_2^2$  representation  $v_1, v_2, \dots, v_n$ , that is  $c$ -spread, there exists a  $d$ -regular multicommodity flow  $(f_{ij})$  for  $d = \alpha/\sqrt{\log n}$  such that

$$\sum_{ij} f_{ij} |v_i - v_j|^2 = \Omega(nd).$$

Since the LHS of (19) is  $\Omega(nd)$  for every choice of the  $v_i$ 's, there is *some* choice of the flow for which the RHS is  $\Omega(nd)$  also. Hence we have proved the following.

COROLLARY 20

Let  $c < 1/4$ . For every graph  $G = (V, E)$  and  $\alpha(G) = \alpha$  there exists a  $d$ -regular flow  $(f_{ij}^*)$  for  $d = \alpha/\sqrt{\log n}$  such that for every unit  $\ell_2^2$ -representation  $v_1, v_2, \dots, v_n$  that that is  $c$ -spread:

$$\sum_{ij} f_{ij}^* |v_i - v_j|^2 = \Omega\left(\frac{\alpha n}{\sqrt{\log n}}\right).$$

Note that the Corollary implies that  $\alpha_c((f^*)) = \Omega(\alpha n/\sqrt{\log n})$ , so we have proven that all large sets expand well in  $f^*$ . To finish the proof of Theorem 17 we need to augment this flow so *all* sets expand well, not just those of size at least  $cn$ . This needs another definition and a lemma. We say that a bipartite weighted graph  $(V_1, V_2, w)$  is a  $\beta$ -*matching* if every weighted degree is at least  $\beta$  and at most  $10\beta$ .

LEMMA 21

Let  $G = (V, E, w)$  be a weighted graph and  $S \subseteq V$  be such that (i) the induced graph  $G|_S$  satisfies  $\alpha(G|_S) \geq \alpha$  and (ii) the induced bipartite weighted graph  $(\bar{S}, S, w)$  is a  $\beta$ -matching. Then  $\alpha(G) \geq \gamma\alpha$ , where  $\gamma = \frac{\beta}{11\beta + \alpha}$ .

PROOF: Let  $U \subseteq V$  be any subset. If  $|U \cap S| \geq \gamma|U|$  then the expansion of  $G|_S$  implies the amount of weight leaving  $U$  is least  $\gamma\alpha|U|$ . So assume  $|U \setminus S| \geq (1 - \gamma)|U|$ . Then  $U \setminus S$  has at least  $(1 - \gamma)\beta|U|$  weight going out of it, of which at most  $\gamma 10\beta|U|$  weight could end up in  $U \cap S$ . Thus at least  $(1 - \gamma)\beta|U| - 10\gamma\beta|U| = \gamma\alpha|U|$  weight leaves  $U$ .  $\square$

Now we prove Theorem 17.

PROOF:(Theorem 17) Let  $c = 1/4 - \epsilon$  for some arbitrarily small  $\epsilon > 0$  and  $(f_{ij}^*)$  be the multicommodity flow given by Corollary 20. Then the minimum  $c$ -balanced cut in the weighted graph  $(f_{ij}^*)$  has capacity  $\Omega(\frac{\alpha n}{\sqrt{\log n}})$ . Thus Lemma 22 implies that there is a subset  $S_1$  of  $(\frac{3}{4} - 2\epsilon)n$  nodes such that the induced (weighted) subgraph of  $(f_{ij}^*)$  on  $S_1$  has expansion  $\Omega(\alpha/\sqrt{\log n})$ .

To extend this subgraph to an expander flow, we take the union of the flow  $f^*$  and the matching given by the next claim for  $d' = \alpha/\sqrt{\log n}$ . By Lemma 21 the union is an expander. In this graph every degree is at most  $2\alpha/\sqrt{\log n}$ , so scaling down by a factor 2 gives a graph with every degree at most  $d = \alpha/\sqrt{\log n}$ . Finally, add a self-loop at each node and give it a weight  $f_{ii}$  equal to  $d - (\text{degree of } i)$ . This uses up no additional capacity and makes the flow  $d$ -regular.

CLAIM: For any  $d' < \alpha/5$  there is a multicommodity flow  $(g_{ij})$  in  $G$  such that the weighted bipartite graph  $(V \setminus S_1, S_1, g)$  is a  $\Theta(d')$ -matching. PROOF: Let  $S_2$  denote  $V \setminus S_1$ . Let  $k = |S_1| / |S_2|$ . Note that  $k < 5$ . Consider a flow problem in which each node of  $S_2$  is the source for  $k$  units of flow and each node of  $T$  is allowed to be the sink for 1 unit of flow and each edge has capacity  $\lceil 1/\alpha \rceil$ . Then by the max-flow min cut theorem, there is a flow of  $|S_1|$  units, namely all available flow gets routed. Now scaling this flow down by  $\alpha$  gives the result.  $\square$   $\square$

LEMMA 22

For each  $c < 1/4$  and  $c_1 < 1$  the following is true for any graph  $G$ . There is an induced subgraph  $G|_U$  of  $G$  with at least  $(1 - c)n$  nodes such that  $\alpha(G|_U) > c_1\alpha_c(G)$ . Furthermore, if  $c \in [1/4, 1/3)$  a similar statement is true for each  $c_1 < \frac{1-2c}{2c}$ .

PROOF: Let  $\alpha$  denote  $\alpha_c(G)$ . Let us iteratively remove sets that do not expand by  $c_1\alpha$ ; clearly the part of the graph that remains at the end (if any!) has expansion at least  $c_1\alpha$ . Let  $S_1, S_2, \dots, S_k \subseteq V$  be the sequence of sets removed at any step and let  $U = V \setminus \cup_{i \leq k} S_i$ . Then the number of edges between  $\cup_i S_i$  and  $U$  is at most  $c_1\alpha(\sum_i |S_i|)$ . Since we know that  $c$ -balanced cuts expand by at least  $\alpha$  and  $\cup_i S_i$  expands by at most  $c_1\alpha < \alpha$ , we conclude that  $\sum_i |S_i| < cn$ .

The proof for  $c \in [1/4, 1/3)$  is similar.  $\square$

**Can one embed expanders integrally instead of fractionally?** We think this should be possible by randomized rounding, though a delicate analysis seems necessary.

### 7.2.1 Proof of Lemma 19

First we state a corollary of our inductive proof of Lemma 7 and the stronger analysis of Section 5.

COROLLARY 23

For every  $c < 1/4$  there are constants  $\gamma, C, \tau > 0$  such that the following is true. If  $G = (V, E)$  is a  $n$ -node graph and  $\alpha(G) = \alpha$  then for every unit- $\ell_2^2$ -representation  $v_1, \dots, v_n$  that is  $c$ -spread, then there is a pair of nodes  $i, j$  such that  $d_G(i, j) \leq C\sqrt{\log n}/\alpha$  and  $|v_i - v_j|^2 \geq \gamma$ . Furthermore this property holds even if we forbid some  $\tau \cdot n$  nodes from playing the role of  $i, j$ .

PROOF: The proof uses the SET-FIND algorithm and appeared in full in the first version of our paper. It used ideas similar to those in our better analysis of Section 5. However, the analysis of our rounding algorithm has now been rewritten with a new geometric viewpoint. So the proof of this corollary needs to be rewritten from scratch. Here is a sketch. Assume for contradiction's sake that endpoints of all paths of length  $O(\sqrt{\log n}/\alpha)$  in the graph are nodepairs  $i, j$  such that  $|v_i - v_j|^2 = o(1)$ . We make a few observations. First, the the matching cover used in Lemma 7 (and other lemmas) can without loss of generality use nodepairs  $(i, j)$  such that the distance of  $i, j$  in the graph is  $O(1/\alpha)$ . This follows from the fact that  $S_u, T_u$  are sets of size  $\Omega(n)$ , and thus have many paths of length  $O(1/\alpha)$  between them.

Next, we note that for any set  $T$  of size  $\rho n$ , breadth-first-search can be used to reach  $\geq n/2$  nodes in  $O(1/\alpha)$  steps. Let  $S$  be the set of nodes reached. Thus if the nodes of  $T$  are well-covered by a set of points, then Lemma 12 can be used to conclude that the nodes of  $S$  are too.

Now we use an induction similar to the one in the proof of Lemma 7 and that in Section 5 to derive a contradiction to the assumption that all paths of length  $O(\sqrt{\log n}/\alpha)$  in the graph correspond to  $i, j$  where  $|v_i - v_j|^2 = o(1)$ . The "robustness" property follows from the observation that all the arguments are essentially unchanged if we set aside  $\tau n$  nodes from participating in the matching cover.  $\square$

Now we prove Lemma 19. First, we note that it suffices to find flows of max-degree  $d$ . Such a flow can be made  $d$ -regular by augmenting it with a self-loop at each node and with a weight  $f_{ii}$  equal to  $d - (\text{degree of } i)$ . This uses up no additional capacity.

To prove the existence of the desired flow is tantamount to proving that the optimum value of the following LP is  $\Omega(nd)$ . In this LP, for  $i \neq j$ ,  $\mathcal{P}_{ij}$  denotes the set of paths  $p$  connecting  $i$  and  $j$  in graph  $G$ . Note that  $v_i$ 's appear as "constants" in this LP; the variables are the  $f_p$ 's.

$$\max \frac{1}{4} \sum_{i,j} \sum_{p \in \mathcal{P}_{ij}} f_p |v_i - v_j|^2 \quad (20)$$

$$\forall e \in E \quad \sum_{p \ni e} f_p \leq 1 \quad (21)$$

$$\forall i \quad \sum_j \sum_{p \in \mathcal{P}_{ij}} f_p \leq d \quad (22)$$

$$f_p \geq 0 \quad (23)$$

We write the dual LP. It involves finding a weighted graph  $\{w_e\}$  whose underlying edge set is the same as  $G$ , and an assignment of non-negative weights  $\{s_i\}$  to the vertices.

$$\min \sum_{e \in E} w_e + d \sum_i s_i \quad (24)$$

$$\forall i, j \quad \forall p \in \mathcal{P}_{ij}, \quad \sum_{e \in p} w_e + s_i + s_j \geq \frac{1}{4} |v_i - v_j|^2. \quad (25)$$

$$\forall e \in E \quad w_e \geq 0 \quad (26)$$

Now we show that for  $d = \alpha/\sqrt{\log n}$  the primal optimum is  $\Omega(\alpha n/\sqrt{\log n}) = \Omega(nd)$ . By LP duality, it suffices to exhibit that every feasible dual has objective value  $\Omega(nd)$ . Fix any feasible  $w_e$ 's and  $s_i$ 's for the dual. If  $d(\sum_i s_i) = \Omega(nd)$  then there is nothing to prove so from now on assume  $d(\sum_i s_i) = o(nd)$ . The goal is to show that then  $W = \sum_e w_e$  is  $\Omega(nd)$ , and this will be done in Claim 1 below.

First, we define a new unweighted graph  $G'$  that is related to  $G$  and the given weights  $w_e$ 's. Let  $\epsilon > 0$  be some small enough constant to be specified later. Graph  $G' = (V', E')$  is obtained by replacing each edge  $e = \{k, l\}$  of  $G$  by a path  $k, k_1, k_2, \dots, k_m, l$  of length  $m + 1 = \lceil \frac{\epsilon w_e n}{2W} \rceil$ . (Here  $k_1, k_2, \dots$ , are indices

of new vertices that are added to the graph.) Let  $V'$  be the set of nodes of  $G'$ . Then the set of new nodes  $V' \setminus V$  satisfies (since  $|E| \leq 2n$ )

$$|V' \setminus V| \leq \sum_e \frac{\epsilon w_e n}{2W} + \sum_{e: w_e > 2W/n\epsilon} 1 \quad (27)$$

$$\leq 2\epsilon n \quad (28)$$

Thus  $G'$  has at most  $n(1 + 2\epsilon)$  nodes.

Furthermore,  $G'$  inherits a natural unit- $\ell_2^2$  representation from  $G$ . Using the above edge  $\{k, l\}$  as an example again, we can use the vector  $v_k$  for  $k, k_1, k_2, \dots, k_{\lfloor m/2 \rfloor}$  and  $v_l$  for the other nodes of the path.

In  $G'$  every set  $U$  of at least  $4\epsilon n$  nodes contains at least  $|U| - 2\epsilon n$  original nodes and hence has at least  $\alpha(|U| - 2\epsilon n) \geq \alpha|U|/2$  edges leaving it. Thus  $\alpha_{4\epsilon}(G') \geq \alpha/2$  and so Lemma 22 implies that  $G'$  contains a subgraph  $G''$  with at least  $(1 - 4\epsilon)n$  original nodes such that  $\alpha(G'') \geq \alpha/3$ . Furthermore, the inherited  $\ell_2^2$  representation is  $c'$ -spread where  $c'$  is some other constant, namely,  $g$  such that  $g(1 - g) = c(1 - c) - (4\epsilon)^2$ . Let  $\tau > 0$  denote maximum allowed fraction of the “forbidden set” in Corollary 23 for this constant  $g$ , and let  $\gamma > 0$  be the constant such that the Theorem yields  $i, j$  such that  $\frac{|v_i - v_j|^2}{4} \geq \gamma$ . Assume  $\epsilon$  is small enough that  $\sum_i s_i < (\tau - 2\epsilon)\gamma n/3$ . (Recall that we assumed  $\sum_i s_i = o(n)$ .)

The next claim completes the proof of Lemma 19.

CLAIM 1:  $W = \Omega(nd)$

PROOF: Let  $B$  be the set of nodes  $i$  satisfying  $s_i \geq \gamma/3$ . Since  $\sum_i s_i < (\tau - 2\epsilon)\gamma n/3$ , we have  $|B| \leq (\tau - 2\epsilon)n$ .

Now treat  $B \cup (V' \setminus V)$  as the “forbidden set” (which we can do since it has size at most  $(\tau - 2\epsilon)n + 2\epsilon n \leq \tau n$ ) and apply Corollary 23 to conclude that there is a node pair  $i, j \in G''$  that are not forbidden and  $d_{G''}(i, j) = O(\sqrt{\log n}/\alpha) = O(1/d)$  and  $\frac{|v_i - v_j|^2}{4} \geq \gamma$ . Since  $i, j$  are not forbidden, they are in  $V$  and  $s_i, s_j < \gamma/3$ . Thus we have shown the existence of  $i, j$  in  $G$  and a path  $p$  connecting them such that  $\sum_{e \in p} \lceil \frac{\epsilon w_e n}{2W} \rceil = O(1/d)$  but at the same time satisfies (thanks to constraint (25) on the dual solutions):

$$\sum_{e \in p} w_e + 2 \cdot \frac{\gamma}{3} \geq \gamma.$$

We conclude that  $W = \Omega(dn)$ , and the Claim is proved.  $\square$

### 7.3 Alternative Approximation Algorithm for SPARSEST CUT

Now we show that Theorem 17 leads to an alternative  $O(\sqrt{\log n})$ -approximation for SPARSEST CUT that does not use SDP.

THEOREM 24

*There is a  $\beta_0 > 0$  and a polynomial-time algorithm that, given a graph  $G = (V, E)$  and a degree bound  $d$  either finds a  $d$ -regular  $\beta_0$ -expander flow in  $G$  or else finds a cut of expansion  $O(d\sqrt{\log n})$ .*

PROOF: We write an LP expressing the existence of a  $\beta$ -expander flow. For all paths  $p$  in the graph, we have a non-negative variable  $f_p$ . Let  $\mathcal{D}$  denote the polytope of demand vectors  $\bar{d} = (d_{ij})$  that correspond to multicommodity flows routable in  $G$  without exceeding edge capacities. (We omit the detailed description of  $\mathcal{D}$ , which is standard.) For any such demand vector  $\bar{d}$  and cut  $(S, \bar{S})$  we denote by  $d(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} d_{ij}$  the amount of demand crossing the cut. We denote the degree of node  $i$ ,  $\sum_j d_{ij}$  by  $d_i$ .

The LP is the following:

$$\bar{d} \in \mathcal{D} \quad (29)$$

$$d_i = d \quad \forall i \quad (30)$$

$$d(S, \bar{S}) \geq \beta |S| \quad \forall S \subseteq V \quad (31)$$

First, we relax the LP so that (34) is only required to hold for  $|S| \geq n/10$  and  $d_i = d$  is relaxed to  $d_i \leq d$ .

To determine feasibility of this relaxed LP by the Ellipsoid method, we need a polynomial-time separation oracle for the constraints. The separation oracle for the first two constraints is trivial, so assume those constraints are satisfied by the current demand vector. An exact separation oracle for the third constraint probably does not exist since the expansion problem is co-NP complete, as noted. So we use an approximate oracle relying on the connection (due to Cheeger) between eigenvalues and expansion. We determine  $\lambda_2(\mathcal{L}(\bar{d}))$ , the second eigenvalue of the Laplacian of the weighted graph  $(d_{i,j})$  (turned into a  $d$ -regular graph by adding self-loops). If this is less than  $\beta^2/2$ , we can use the corresponding eigenvector to find a cut  $(S_1, \bar{S}_1)$  of expansion less than  $\beta$  in  $G$ . If  $|S_1| \geq n/10$  we have found a violated constraint. Otherwise iterate on the demand vector restricted to  $V \setminus S_1$  to find other sets  $S_2, S_3, \dots$ , of expansion less than  $\beta$  until the second eigenvalue of the remaining demand graph rises above  $\beta^2/2$  and no more such sets can be found. Now if  $|S_1 \cup S_2 \cup S_3 \dots| > n/10$  then we have found a violated constraint, and failing that, we have convinced ourselves that  $G$  contains a subgraph of at least  $9n/10$  vertices where  $\bar{d}$  is a  $\beta^2/2$ -expander.

Thus by the end of the Ellipsoid method either we conclude the relaxed LP is infeasible or we find a demand graph that is almost an expander flow.

CASE 1: We find a demand graph  $\bar{d}$  such that for some subgraph  $G|_S$  with at least  $9n/10$  vertices it is a  $\beta^2/2$ -expander. As shown in the proof of Theorem 17, specifically, the Claim in that proof, given any feasible solution to even this relaxed LP we can try to extend it into an expander flow using a single  $s$ - $t$  maximum flow computation. However, instead of the weight  $1/\alpha$  used there (recall that  $\alpha$  is not known) use a weight of  $1/d$  on each edge. If the max-flow computation does yield a bipartite matching, then we have obtained a  $d$ -regular flow that is an  $\Omega(\beta^2)$ -expander. If the max-flow computation does not yield a bipartite matching, then the  $s$ - $t$ -minimum cut must have expansion less than  $d$  and we have found a cut.

CASE 2: during the execution of the Ellipsoid algorithm we discovered  $\text{poly}(n)$  cuts  $S_1, S_2, S_3, \dots$ , each containing at least  $n/10$  nodes and such that the LP with expansion constraints for just these sets is infeasible. Namely, the following LP:

$$\bar{d} \in \mathcal{D} \quad (32)$$

$$d_i \leq d \quad \forall i \quad (33)$$

$$d(S, \bar{S}) \geq \beta |S| \quad \forall S \in \{S_1, S_2, \dots\} \quad (34)$$

Using the polytopal characterization of  $\mathcal{D}$  and Farkas' Lemma this is infeasible iff there is an assignment of weight  $w_e \geq 0$  to each edge  $e$ ,  $s_i \geq 0$  on each node  $i$ , and  $z_S \geq 0$  to each  $S \in \{S_1, S_2, \dots\}$  such that  $\sum_S z_S \leq 10$

$$\sum_{e \in E} w_e + d \sum_i s_i - \sum_S |S| z_S < 0 \quad (35)$$

$$\forall i, j \in V, p \in \mathcal{P}_{ij} s_i + s_j + \sum_{e \in p} w_e \geq \text{cross}(i, j) \quad (36)$$

where  $\text{cross}(i, j)$  is shorthand for  $\sum_{S: i \in S, j \in \bar{S}} z_S$ .

Determine such weights by linear programming. Now we turn the  $z_S$ 's into a unit  $\ell_2^2$ -representation by considering the distribution of cuts given by  $\frac{z_S}{Z}$  where  $Z = \sum_S z_S$ . Recall that every cut has an obvious  $\ell_2^2$  representation, hence so does every distribution on cuts. Now  $|v_i - v_j|^2$  is another way to write  $\text{cross}(i, j)$ .

Then go from the weighted graph given by  $w_e$ 's to an unweighted graph as in the proof of Lemma 19 and apply Corollary 25 below to conclude that if such  $w_e$ 's,  $s_i$ 's,  $z_S$ 's exist then the algorithm must have found a cut of expansion at most  $O(d\sqrt{\log n})$ .

□

Above, we used the following “algorithmic” version of Corollary 23.

**COROLLARY 25**

*For every  $c < 1/4$  there are constants  $\gamma, C, \tau$  and a  $O(mn)$ -time randomized algorithm such that the following is true. Given a graph  $G$  with a unit- $\ell_2^2$ -representation  $v_1, v_2, \dots, v_n$  that is  $c$ -spread, the algorithm finds a pair of vertices  $i, j$  and a cut of expansion  $\alpha_{\text{obs}}$  such that their graph distance is  $\leq C\sqrt{\log n}/\alpha_{\text{obs}}$  and  $|v_i - v_j|^2 \geq \gamma$ . Furthermore, the algorithm can find such a pair even if we forbid some  $\tau n$  vertices from playing the role of  $i, j$ .*

## 8 The four conjectures

We feel that our techniques can be strengthened to upperbound the integrality gap by a smaller function than  $\sqrt{\log n}$ , and possibly even by  $O(1)$  (as conjectured in [15]). As noted earlier, our main theorem (Theorem 1) about the existence of large  $\Delta$ -separated subsets cannot be improved. However, the alternative approach using expander flows seems to not suffer from such a limitation. In particular, we do not know if Corollary 23 is tight and our four conjectures below stem from our failure both to strengthen it and to rule out such a strengthening. (Aside: Our discussion of these conjectures is somewhat sketchy in this version because the paper has been recently overhauled to focus on the purely geometric viewpoint represented by Theorem 1. The conjectures fitted in better with the original presentation of our ideas, which involve a mix of geometry and graph theory as represented by Corollary 23.)

Now we list the first three conjectures. All three concern any constant degree graph  $G = (V, E)$  and any unit- $\ell_2^2$ -representation  $v_1, v_2, \dots, v_n$ , that is  $c$ -spread.

**Conjecture 1:** There is an edge  $\{i, j\} \in E$  such that  $|v_i - v_j|^2 = \Omega(\alpha)$ .

**Conjecture 2:** There are pairs of vertices  $i, j$  such that  $d_G(i, j) = O(1/\alpha)$  and  $|v_i - v_j|^2 = \Omega(1)$ .

**Conjecture 3:** Version of Conjectures 2 whereby  $\tau n$  fraction of nodes are forbidden from being chosen as  $i, j$  and nevertheless these  $i, j$  exist.

Now we summarise the implications of these conjectures. All conjectures are interesting in their own right.

**LEMMA 26**

1. *If the integrality gap of the SDP is  $O(1)$  then Conjecture 1 holds.*
2. *Proving Conjecture 2 suffices to prove that the integrality gap is  $O(1)$ .*
3. *Conjecture 3 implies the existence of optimal (upto a constant factor) expander flows in graphs. Namely, Theorem 17 is true with  $d = \alpha$ .*

**PROOF:**

1. If the integrality gap of the SDP is  $O(1)$  then the value of the objective is  $\Omega(\alpha n)$  and hence at least one edge has value  $\Omega(\alpha)$ .
2. Uses a modification of our original proof that the integrality gap is  $O(\sqrt{\log n})$ . That proof was modified during the switch to the geometric viewpoint. It will be put in again here in the journal version.
3. Conjecture 3 is a stronger version of Corollary 23. One can then mimic the proof of Theorem 17 using this stronger corollary.

□

Note that we have proved weaker versions of Conjectures 2 and 3 by replacing  $\Omega(1)$  by  $\Omega(1/\sqrt{\log n})$ . Replacing  $\sqrt{\log n}$  by any weaker function of  $n$  also gives analogous bounds on integrality gaps and expander flows.

### 8.1 Bringing the conjectures “down” to $\ell_1$

A unit  $\ell_2^2$ -representation  $v_1, v_2, \dots, v_n$  of a graph is said to be  $\ell_1$  if there is a set of vectors  $u_1, u_2, \dots, u_n$  such that  $|v_i - v_j|^2 = |u_i - u_j|_1$ . We say that it is  $\ell_1$  upto distortion  $c$  if  $|u_i - u_j|_1 \leq |v_i - v_j|^2 \leq c \cdot |u_i - u_j|_1$ . It had been conjectured that every  $\ell_2^2$ -representation is  $\ell_1$  upto a distortion factor  $O(1)$ . If true, this would imply an integrality gap of  $O(1)$  for the SDP in (4) to (7), using the following characterization of  $\ell_1$  representations.

LEMMA 27 (WELL-KNOWN)

A unit  $\ell_2$ -representation  $v_1, v_2, \dots, v_n$  is  $\ell_1$  iff there is a  $\alpha_S \geq 0$  associated with each cut  $(S, \bar{S})$  such that

$$|v_i - v_j|^2 = \sum_S \alpha_S d_S(i, j), \quad (37)$$

where  $d_S(i, j) = 0$  if  $i, j$  are on the same side of the cut and 4 otherwise. (In other words,  $\ell_1$  representations correspond exactly to the cut cone.)

PROOF: One direction is trivial from convexity since every cut metric is  $\ell_1$ .

For a proof of the other direction see Shmoys’ survey [31] or Matousek’s book [26]. □

Thus minimization over  $\ell_1$  metrics is exactly equivalent to minimizing over cuts (and thus NP-hard).

Proving the equivalence (upto  $O(1)$  distortion) of  $\ell_2^2$  and  $\ell_1$  seems difficult, however. The only known result is due to Goemans [16] who shows that  $\ell_2^2$  metrics derived from vectors in  $\mathbb{R}^d$  embed into  $\ell_1$  with distortion  $O(\sqrt{d})$ , which for  $d = o(\log^2 n)$  improves Bourgain [7]’s more general bound of  $O(\log n)$ , but helps little for the general case  $d = n$ . Furthermore, many researchers (including some of the authors) do not believe that all  $\ell_2^2$  metrics embed in  $\ell_1$  with constant distortion. Thus it is conceivable that the conjectured equivalence between  $\ell_2^2$  and  $\ell_1$  is false, and yet Conjecture 2 is true (and thus the SDP has a constant integrality gap).

However, many researchers find  $\ell_1$  metrics much easier to grasp than  $\ell_2^2$  metrics, so here we present a conjecture about  $\ell_1$  metrics that would also suffice to prove an integrality gap of  $O(1)$  for the SDP.

[Conjecture 4:] If the unit- $\ell_2^2$ -representation is actually  $\ell_1$ , then there is a vertex pair  $i, j$  such that  $d_G(i, j) = O(1/\alpha)$  and  $|v_i - v_j|^2 = \Omega(1)$ .

In light of the preceding discussion, the following equivalence between our conjectures for  $\ell_2^2$  and  $\ell_1$  may be surprising.

LEMMA 28

Conjecture 2 and Conjecture 4 are equivalent.

PROOF: Conjecture 2  $\Rightarrow$  Conjecture 4: Trivial since Conjecture 4 is a subcase of Conjecture 2.

Conjecture 4  $\Rightarrow$  Conjecture 2: Given any unit- $\ell_2^2$  representation  $v_1, v_2, \dots, v_n$  of the graph, consider the uniform distribution on all hyperplane cuts (a la [GW]). Represent this distribution using Lemma 27 by an  $\ell_2^2$  representation that is  $\ell_1$ , namely, a set of unit vectors  $u_1, u_2, \dots, u_n$  where

$$\left| u_i - u_j \right|^2 = \Pr[i, j \text{ are separated in hyperplane cut produced using } v_1, \dots, v_n].$$

Let  $i, j$  be the pair of nodes in  $G$  whose existence follows from Conjecture 4, namely, with  $\left| u_i - u_j \right|^2 = \Omega(1)$ . Then these also suffice for Conjecture 2 since if

$$\Pr[i, j \text{ are separated in hyperplane cut on } v_1, \dots, v_n] = \Omega(1),$$

then  $\left| v_i - v_j \right|^2 = \Omega(1)$  also.  $\square$

REMARK 9 If we weaken Conjectures 2 and 4 to require  $\left| v_i - v_j \right|^2 = \Omega(1/s(n))$  then the equivalence doesn't quite hold. The weakened Conjecture 2 still implies the weakened Conjecture 4, but the weakened Conjecture 4 only implies an even weaker Conjecture 2 with  $s(n)$  replaced by  $s(n)^2$ .

REMARK 10 We can make a robust form of Conjecture 4 whereby the pair exists even if we forbid  $\tau n$  nodes from playing the role of  $i, j$ . This is a stronger version of Corollary 23. Then this conjecture is *equivalent* to the conjecture that the graph as  $d$ -regular expander flows for  $d = \alpha$ . (We argued the more interesting “only if” direction but “if” is easy to prove as well.)

The optimistic view of Lemma 28 is that if we feel  $\ell_1$  metrics are easier to work with than  $\ell_2^2$ , then our hopes of resolving Conjecture 2 are bolstered. (Indeed, some would see Lemma 28 as supporting the conjecture that  $\ell_2^2$  and  $\ell_1$  are equivalent upto distortion  $O(1)$ .)

The pessimistic view of Lemma 28 is that that trying to strengthen our SDP with other constraints (in addition to the triangle inequality) may not help. Though this could conceivably give a class of metrics that is a proper subset of  $\ell_2^2$ , proving our conjectures for that other class of metrics will not be any easier, since that class of metrics (being a relaxation of cut metrics) would always contain  $\ell_1$ . This is why the title of the current subsection announces that the conjectures have been brought “down” to  $\ell_1$ .

## 9 Conclusions

We feel it should be possible to show that the SDP has integrality gap  $O(1)$  or  $O(s(n))$  for some slowly growing function like  $s(n) = \log \log n$ . Our conjectures provide a roadmap to this task. However, some newer rounding algorithm may be required since on hypercubes and related graphs, our rounding algorithm produces cuts whose value is  $O(\sqrt{\log n})$  times the SDP value.

In this connection —especially to derive approximations for generalizations of SPARSEST CUT— it will also be useful to resolve the conjecture about low-distortion embeddings of  $\ell_2^2$  into  $\ell_1$ . As mentioned in the introduction, our geometric results may be a starting point.

Our approximation algorithms are fairly inefficient (though polynomial time) because they use SDPs or related convex optimization, and solving the SDP of (4)–(7) takes about  $n^{4.5}$  time for an  $n$ -node graph using interior point methods. Do more efficient (combinatorial?) approximation algorithms exist possibly using expander flows? One loose analogy would be combinatorial versions of the Leighton-Rao multicommodity flow algorithm (see [28, 13], two papers in a long line of research). Such algorithms may be useful in practice. Many practitioners continue to prefer eigenvalue methods over Leighton-Rao because the geometric meaning of eigenvalues (e.g., the connection to stretched rubberbands and such) has relevance in their application —computer vision, for example. Since the SDP relaxation may

be viewed as a higher-dimensional analogue of eigenvalue computation, it may well turn out to share these properties of eigenvalues, and hence also their practical appeal.

Extending our ideas to other problems should be possible though it doesn't seem to be immediate. The problems in [31] would be a good list to try, especially minimum multicut, for which an  $O(\log k)$ -approximation was designed in [14].

Finally, we would love to see an application of expander flows to something other than estimating the conductance/expansion.

## Acknowledgements

This project evolved over several years. A partial list of people who gave us useful feedback about our ideas and/or thought about our conjectures (apologies to people we forgot): Farid Alizadeh, Dorit Aharonov, Noga Alon, Moses Charikar, Michel Goemans, Eran Halperin, Mohammed Hajiaghayi, Subhash Khot, Robi Krauthgamer, Tom Leighton, Laci Lovász, Ran Raz, Madhu Sudan, Vijay Vazirani, Santosh Vempala, David Williamson.

We are very grateful to James Lee for comments on our first manuscript. He drew our attention to the geometric core of our argument, which helped us greatly improve the presentation.

## References

- [1] F. Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM J. Optimization*, 5:13–51, 1995.
- [2] N. Alon. Eigenvalues and expanders. *Combinatorica* 6:83–96, 1986.
- [3] N. Alon and V. Milman.  $\lambda_1$ , isoperimetric inequalities for graphs and superconcentrators. *J. Combin. Theory B* 38:73–88, 1985.
- [4] Y. Aumann and Y. Rabani. An  $O(\log k)$  approximate min-cut max-flow theorem and approximation algorithms. *SIAM J. Comp*
- [5] K. Ball. An elementary introduction to modern convex geometry, in *Flavors of Geometry*, S. Levy (ed.), Cambridge University Press, 1997.
- [6] M. Blum, R. Karp, O. Vornberger, C. Papadimitriou, M. Yannakakis. The complexity of testing whether a graph is a superconcentrator. *Inf. Proc. Letters* 13:164-167, 1981.
- [7] J. Bourgain. On Lipschitz embeddings of finite metric spaces in Hilbert space. *Israel J. Mathematics* 52:46–52, 1985.
- [8] J. Cheeger. A lower bound for the smallest eigenvalue of the Laplacian, in *Problem in Analysis*, 195-199, Princeton Univ. Press, (1970),
- [9] F. Chung. Spectral graph theory. *CBMS Regional Conference Series in Mathematics*, 92, American Mathematical Society, 1997.
- [10] L. Danzer and B. Grünbaum. On two problems of P. Erdős and V. L. Klee concerning convex bodies (in German). *Math. Zeitschrift* 79:95–99, 1962.
- [11] P. Diaconis and L. Saloff-Coste. Comparison theorems for reversible Markov chains. *Annals of Applied Probability*, 3:696–730, 1993.

- [12] U. Feige and R. Krauthgamer. A polylogarithmic approximation of the minimum bisection. In *IEEE FOCS 2001* pp 105-115.
- [13] N. Garg and J. Könemann. Faster and Simpler Algorithms for Multicommodity Flow and other Fractional Packing Problems. In *IEEE FOCS 1997*.
- [14] N. Garg and V. V. Vazirani and M. Yannakakis. Approximate max-flow min-(multi)cut theorems and their applications. *SIAM J. Computing*, **25**(2):235-251, 1996. *Prelim. version in Proc. ACM STOC'93*.
- [15] M.X. Goemans. Semidefinite programming in combinatorial optimization. *Math. Programming*, **79**:143-161, 1997.
- [16] M. X. Goemans. *unpublished note*.
- [17] M.X. Goemans and D.P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *JACM*, **42**(6):1115-1145, 1995.
- [18] M. Grötschel, L. Lovász, and A. Schrijver. *Geometric algorithms and combinatorial optimization*. Springer-Verlag, 1993.
- [19] H. Karloff and U. Zwick. A 7/8-approximation algorithm for MAX 3SAT? *Proc. of 38th IEEE FOCS* (1997), 406-415.
- [20] D. Karger, R. Motwani, and M. Sudan. Approximate graph coloring by semidefinite programming. *JACM*, **45**(2):246-265, 1998.
- [21] M. Jerrum and A. Sinclair. Approximating the permanent. *SIAM J. Comput.*, **18**(6):1149-1178, 1989.
- [22] T. Leighton and S. Rao. Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms. *JACM* **46** 1999. *Prelim. version in ACM STOC 1988*.
- [23] N. Linial, E. London, and Y. Rabinovich. The geometry of graphs and some of its algorithmic applications. *Combinatorica* **15**:215-246, 1995.
- [24] L. Lovász. On the Shannon capacity of a graph. *IEEE Trans. on Info. Theory* **IT-25**:1-7, 1979.
- [25] A. Lubotzky, R. Philips, and P. Sarnak. Ramanujan graphs. *Combinatorica*, **8**:261-277, 1988.
- [26] J. Matousek. *Lectures on Discrete Geometry*. Springer Verlag, 2002.
- [27] Y. Nesterov and A. Nemirovskii. *Interior point polynomial methods in convex programming*. SIAM, Philadelphia, PA 1994.
- [28] S. Plotkin and D. B. Shmoys and E. Tardos. Fast approximation algorithms for fractional packing and covering problems. *Math. Operations Res.* **20**:257-301, 1995. *Prelim. version IEEE Foundations of Computer Science, 1991, 495-504*.
- [29] G. Schechtman. Concentration, results and applications. *Handbook of the Geometry of Banach Spaces*, volume 2, W.B. Johnson and J. Lindenstrauss (eds.), North Holland, 2003. Draft version available from Schechtman's website.
- [30] F. Shahrokhi and D.W. Matula. The maximum concurrent flow problem. *Journal of the ACM*, **37**:318-334, 1990.
- [31] D. S. Shmoys. Cut problems and their application to divide and conquer. *Approximation Algorithms for NP-hard problems*, D.S. Hochbaum (ed.), PWS Publishing, 1995.
- [32] A. Sinclair. Improved bounds for mixing rates of Markov chains and multicommodity flow. *Combinatorics, Prob., Comput.* **1**:351-370, 1992.

[33] V. Vazirani. Approximation algorithms. Springer Verlag, 2002.

[34] K. Zatloukal. *Personal communication*, November 2003.

## A Reduction to bounded degree graphs

For completeness we give approximation-preserving reductions from the SPARSEST CUT and  $c$ -BALANCED CUT problems to the bounded degree versions of these problems.

Given any instance  $G = (V, E)$  of the  $c$ -BALANCED CUT problem we replace each node  $i$  by a gadget  $H_i$  that is a *strong expander* on  $m = O(n^3)$  nodes. A *strong expander* is a graph in which for each subset of nodes  $S$  of fewer than half the nodes, the number of edges leaving it is at least  $|S| + 1$ ; these were constructed in [25] (though weaker expander constructions would also suffice for us). If  $d_i$  was the degree of  $i$  in  $G$ , then we designate  $d_i$  nodes of  $H_i$  as *special*, and if  $\{i, j\}$  was an edge in  $G$  then we connect two special nodes of  $H_i, H_j$ . These special nodes will not be used for any other edges. Let  $G'$  denote this new graph on  $mn = O(n^4)$  nodes.

Now if  $(S', \overline{S'})$  is the optimum cut in  $G'$ , then it must be that every gadget  $H_i$  lies entirely in  $S'$  or  $\overline{S'}$ . (For, if not, then moving it to one side or the other must strictly decrease the number of edges in the cut.) Thus the cut corresponds to a  $c$ -balanced cut in  $G$ .

The same reduction works for sparsest cut.

## B Version of Min-Max Theorem

Section 7 uses the following consequence of the von Neumann min-max theorem that we prove for sake of completeness. Let  $\pi : A \times B \rightarrow \mathfrak{R}$  be a payoff function where  $A, B$  are finite sets. According to the min-max theorem:

$$\min_{D_1} \max_{b \in B} E_{a \in D_1}[\pi(a, b)] = \max_{D_2} \min_{a \in A} E_{b \in D_2}[\pi(a, b)], \quad (38)$$

where  $D_1, D_2$  are distributions on  $A, B$  respectively.

We wish to identify sufficient conditions under which the optimum strategies in the game are deterministic. Furthermore, we wish to allow  $A, B$  to be infinite sets. We will express the optimum as a limit of optima of finite games.

Set  $A$  is *convex* if for each  $a_1, a_2 \in A$  and  $t \in [0, 1]$  there is a way to define  $ta_1 + (1 - t)a_2 \in A$ .

We say that the payoff function is *linear* in each coordinate if for all positive  $t_1, t_2$  satisfying  $t_1 + t_2 = 1$ :  $\pi(t_1 a_1 + t_2 a_2, b) = t_1 \pi(a_1, b) + t_2 \pi(a_2, b)$  and  $\pi(a, t_1 b_1 + t_2 b_2) = t_1 \pi(a, b_1) + t_2 \pi(a, b_2)$ .

Note that if  $A, B$  are convex and the payoff function is linear then for any distribution  $D_1$  that gives probability  $p_a$  to  $a \in A$  the element  $a^* = \sum_a p_a \cdot a$  is also a member of  $A$  and for all  $b \in B$ :  $E_{a \in D_1}[\pi(a, b)] = \pi(a^*, b)$ .

Finally,  $(A, B, \pi)$  are *finitely approximable at all scales* if for every  $\epsilon > 0$ , there exist finite  $A_\epsilon \subseteq A, B_\epsilon \subseteq B$  such that:

$$\forall a \in A \exists a_\epsilon \in A \text{ such that } \pi(a, b) \in [\pi(a_\epsilon, b) - \epsilon, \pi(a_\epsilon, b) + \epsilon],$$

and a similar property is true for  $B_\epsilon$ .

LEMMA 29

*If  $\pi$  is linear in each coordinate, and  $A, B$  are convex sets that are finitely approximable at all scales, then*

$$\min_{a \in A} \max_{b \in B} \pi(a, b) = \max_{b \in B} \min_{a \in A} \pi(a, b)$$

PROOF: Consider the finite game on  $A_\epsilon, B_\epsilon$  and take the limit as  $\epsilon \rightarrow 0$ .  $\square$

## C Eigenvalues and expansion

DEFINITION 8 Let  $G$  be a weighted complete graph on  $n$  nodes with  $w_{ij}$  denoting the (nonnegative) weight on edge  $\{i, j\}$ . The expansion is defined as

$$\nu(w) = \min_{S \subseteq [n]} \frac{w(S, \bar{S})}{\min\{|S|, |\bar{S}|\}}. \quad (39)$$

and the conductance is defined as

$$c(w) = \min_{S \subseteq [n]} \frac{w(S, \bar{S})}{\min\{w(S), w(\bar{S})\}}. \quad (40)$$

Let  $D$  be the diagonal matrix in which  $D_{ii} = \sum_j w_{ij}$ . The analog of the adjacency matrix for weighted graphs is  $W = (w_{ij})$ , and the *Combinatorial Laplacian*  $C$  is  $D - W$ . The *Laplacian*  $\mathcal{L}$  is  $D^{-1/2}CD^{-1/2}$ . This is positive semidefinite and its smallest eigenvalue is 0. Denoting the second smallest eigenvalues by  $\lambda_{\mathcal{L}}$  we have the standard facts

$$\lambda_{\mathcal{L}} = \min_{\vec{x}: \sum_i D_{ii} x_i = 0} \frac{\sum_{ij} w_{ij} (x_i - x_j)^2}{\sum_i D_{ii} x_i^2} \quad (41)$$

In the special case when  $D_{ii} = d$ .

The connection to expansion is as follows:

$$2c(w) \geq \lambda_{\mathcal{L}} \geq \frac{c(w)^2}{2} \quad (\text{Alon-Cheeger and Alon-Milman inequalities}) \quad (42)$$

**An  $O(1/\alpha)$  bound on the integrality gap.** We sketch a proof that if  $\alpha(G) = \alpha$  then the sparsest cut relaxation of (8)–(10) has optimum value  $\Omega(\alpha^2/n)$ . (For an alternative proof see Goemans [15].) We can prove this using an analogue of Conjecture 4 but using a pair of nodes whose distance is  $O(1/\alpha^2)$  rather than  $O(1/\alpha)$ . This uses the fact that the random walk on such graphs mixes in  $O(1/\lambda)$  time, which is  $O(1/\alpha^2)$ . Recall that an  $\ell_1$  representation corresponds to a probability distribution on cuts; see Lemma 27. It actually suffices (by arguments similar to our existence proof of expander flows) to consider distributions on balanced cuts. Then the upperbound on mixing times shows that there is a path  $p \in \mathcal{P}_{ij}$  of length  $O(1/\alpha^2)$  that crosses a constant fraction of cuts, i.e.,  $\|v_i - v_j\|_1 = \Omega(1)$ .