

A spectral theorem on the cluster structure of real-world graphs

Sabyasachi Basu ✉

Department of Computer Science and Engineering, University of California, Santa Cruz

Suman Kalyan Bera ✉

Katana Graph

C. Seshadhri ✉

Department of Computer Science and Engineering, University of California, Santa Cruz

Abstract

Partitioning a graph into clusters of vertices is a fundamental problem in computer science and applied mathematics. Arguably, the most important tool for graph partitioning is the Fiedler vector or discrete Cheeger inequality. This result relates the eigenvalues of the normalized adjacency matrix to the low conductance cuts of the graph. However, the Cheeger inequality has little relevance on an important contemporary graph partitioning problem, that of community detection in massive real-world graphs. There are numerous, small, dense clusters in real-world graphs, while Cheeger inequalities focus on partitioning a graph into a few, large clusters. Inspired by the structure of real-world graphs, we define the *spectral transitivity*, a ratio of powers of eigenvalues of the normalized adjacency matrix \mathcal{A} . We discover that constant spectral transitivity implies that a constant fraction of \mathcal{A} is contained in nearly uniform submatrices. Our result is a new spectral theorem that relates the eigenvalues of \mathcal{A} to a cluster structure in \mathcal{A} . The latter structure mimics the observed cluster structure of real-world graphs.

2012 ACM Subject Classification Mathematics of computing → Graph algorithms

Keywords and phrases Graph partitioning, Spectral graph theory, Social networks

Digital Object Identifier 10.4230/LIPIcs...1

1 Introduction

Graph partitioning or clustering is a fundamental problem in theoretical computer science. It has a rich history in the study of algorithms, applied mathematics, and computer science. One of the central tools in graph partitioning is the discrete Cheeger inequality, which goes back to seminal work of Fiedler, and Alon and Milman [9, 1]. This inequality relates the eigenvalues of the graph Laplacian to the graph conductance, showing a connection between the spectrum and graph structure. Consider an undirected graph $G = (V, E)$, where d_v denotes the degree of vertex v . The *normalized adjacency matrix*, denoted \mathcal{A} , is defined as follows: $\mathcal{A}_{uv} = 1/\sqrt{d_u d_v}$ if $(u, v) \in E$, and zero otherwise. The eigenvalues of this matrix are denoted $1 = \lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \lambda_n \geq -1$.

We recall the definition of graph *conductance*. For any subset of vertices $S \subseteq V$, let $\text{Vol}(S) := \sum_{s \in S} d_s$. The conductance of set S is $\Phi(S) := E(S, \bar{S}) / \min(\text{Vol}(S), \text{Vol}(\bar{S}))$. The conductance of the graph G , Φ_G , is defined as $\min_{S \subseteq V} \Phi(S)$. The classic Cheeger inequality relates the spectral gap, $1 - \lambda_2$, to the graph conductance.

► **Theorem 1.1.** (*Cheeger inequality [6, 18]*) $4\sqrt{1 - \lambda_2} \geq \Phi_G \geq (1 - \lambda_2)/4$

This theorem is the foundation of spectral graph theory. The proof also yields an efficient algorithm that finds a low conductance cut.

One of the most important contemporary applications of graph clustering is *community detection* in real-world sparse graphs [25, 24, 23, 10, 11]. Despite the wide applicability of the



© S. Basu, S. K. Bera and C. Seshadhri;
licensed under Creative Commons License CC-BY 4.0
Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

44 Cheeger inequality in general, it is surprisingly irrelevant for network science applications.
 45 Firstly, conductance pertains to breaking the graph into two parts. There are higher order
 46 Cheeger inequalities that deal with k parts, but these are only applicable for $k = O(\log n)$ [20].
 47 Real-world graphs have an extremely large number of small clusters, each of which is internally
 48 dense [21, 29]. Variants of the Cheeger inequalities cannot deal with large k (say, n^δ) and
 49 do not give edge density guarantees about the interior of clusters. We note that there are
 50 local partitioning theorems inspired by the proof of the Cheeger inequality that find small
 51 clusters or give some guarantees on internal structure [37, 20, 26, 27]. But these results do
 52 not connect the graph spectrum to graph structure.

53 Diffusion/PageRank based methods on real-world graphs find a large number of small
 54 sets with conductances around 0.1 or so [21, 13]. For real-world graphs, the connection
 55 between spectral gap and conductance does not seem to be the central theme. In fact,
 56 the commonly observed small world property implies a fairly large spectral gap [19]. Most
 57 real-world networks have a significant fraction of long-range edges or weak ties, that are *not*
 58 part of any community [22, 14, 19]. These edges essentially make the graph be an expander,
 59 in which case the Cheeger inequality has little to say. The spectral gap is sensitive to noise,
 60 so adding (say) a sparse Erdős-Rényi graph (or a set of random edges) on top of an existing
 61 graph could dramatically change the spectral graph and conductances. But that is exactly
 62 what is used for certain models for social networks [29].

63 Our main motivation is:

64 *Can we relate the graph spectrum to the cluster properties of real-world graphs?*

65 1.1 Main result

66 We take inspiration from a central property of real-world graphs, the abundance of tri-
 67 angles [36, 29]. This abundance is widely seen across graphs that come by disparate domains.
 68 Recent work in network science and data mining have used the triangles to effectively cluster
 69 graphs. There is much evidence that the triangle structure aids finding communities in
 70 graphs [28, 33, 3, 34].

71 In network science, the triangle count is often expressed in terms of the *transitivity* or
 72 global clustering coefficient [8, 35]. We define the *spectral transitivity* of the graph G .

73 **► Definition 1.2.** *The spectral transitivity of G , denoted $\tau(G)$, is defined as follows¹. (Recall*
 74 *that the λ_i s are the eigenvalues of the normalized adjacency matrix.)*

$$75 \quad \tau(G) = \frac{\sum_{i \leq n} \lambda_i^3}{\sum_{i \leq n} \lambda_i^2}. \quad (1)$$

76 Standard arguments show that the spectral transitivity is a degree weighted transitivity.
 77 The numerator is a weighted sum over all triangles, while the denominator (squared Frobenius
 78 norm) is a weighted sum over edges (Lemma 3.5).

79 Observe that since $\lambda_i \leq 1$, $\tau \leq 1$. When τ reaches its maximum value of $1 - 1/(n - 1)$,
 80 one can show that G is a clique (Lemma 3.6). We formalize the notion of "clique-like"
 81 submatrices through the concept of uniformity. For a symmetric matrix M and a subset S
 82 of its columns/rows, we use $M|_S$ to denote the square submatrix restricted to S (on both
 83 columns and rows).

¹ If G (or the normalized adjacency matrix \mathcal{A}) are obvious from context, we simply refer to τ instead of $\tau(G)$.

84 ► **Definition 1.3.** Let $\alpha \in (0, 1]$. Let \mathcal{A} be the normalized adjacency matrix of a graph G .
 85 For any subset of vertices S , $\mathcal{A}|_S$ is called α -strongly uniform if at least an α -fraction of
 86 non diagonal entries have values in the range $[\alpha/(|S| - 1), 1/\alpha(|S| - 1)]$.

87 For $s \in S$, let $N(s, S)$ denote the neighborhood of s in S (we define edges by non-zero
 88 entries). An α -uniform matrix is strongly α -uniform if for at least an α -fraction of $s \in S$,
 89 $\mathcal{A}|_{N(s, S)}$ is also α -uniform.

90 Observe that the normalized adjacency matrix of a clique is (strongly) 1-uniform. But
 91 submatrices of this matrix are not. Roughly speaking, a constant uniform submatrix
 92 corresponds to a dense subgraph of (say) size k where the *total* degrees of vertices is
 93 $\Theta(k)$. Strong uniformity is closely related to *clustering coefficients*, which is the edge
 94 density of neighborhoods. It is well-known that real-world graphs have high clustering
 95 coefficients [36, 29]. A strongly uniform submatrix essentially exhibits high clustering
 96 coefficients.

97 Our main theorem states that any graph with constant spectral transitivity can be
 98 decomposed into constant uniform blocks. We use $\|M\|_2$ to denote the Frobenius norm of
 99 matrix M .

100 ► **Theorem 1.4 (Spectral Theorem).** There exist absolute constants $\delta > 0$ and $c > 0$ such
 101 that the following holds. Let \mathcal{A} be the normalized adjacency matrix of a graph with spectral
 102 transitivity τ .

103 There exists a collection of disjoint sets of vertices X_1, X_2, \dots, X_k satisfying the following
 104 conditions:

- 105 1. (Cluster structure) For all $i \leq k$, $\mathcal{A}|_{X_i}$ is strongly $\delta\tau^c$ -uniform.
- 106 2. (Coverage) $\sum_{i \leq k} \|\mathcal{A}|_{X_i}\|_2^2 \geq \delta\tau^c \|\mathcal{A}\|_2^2$.

107 We call this output the *spectral triadic decomposition*. Our proof also yields an efficient
 108 algorithm that computes the decomposition, whose running time is dominated by a triangle
 109 enumeration. Details in are given in Theorem 6.1 and §6.

110 1.2 Significance of Theorem 1.4

111 One can think of Theorem 1.4 as a type of Cheeger inequality that is relevant to the structure
 112 of real-world social networks. We explain how it captures many of the salient properties of
 113 clusters in real-world networks. In this discussion, we will assume that τ is a constant.

114 **The spectral transitivity:** We find it remarkable that a bound on a single spectral
 115 quantity, τ , implies such a rich decomposition. The spectral transitivity τ captures a key
 116 property of real-world graphs, the abundance of triangles. While there is a rich body of
 117 empirical work on using triangles to cluster graphs, there is no theory explaining *why* triangles
 118 are so useful. Theorem 1.4 gives a spectral-theoretic explanation.

119 The spectral transitivity is a weighted version of the transitivity, which is typically
 120 around 0.1 for real-world graphs². We also note that the final algorithm that computes the
 121 decomposition focuses on triangle cuts, which is a popular empirical technique for finding
 122 clusters in social networks [3, 34].

123 **The strong uniformity of clusters:** Each cluster X_i of the spectral triadic decomposi-
 124 tion is (constant) strongly uniform. While there is no one definition of a "community" in
 125 real-world graphs, the definition of strong uniformity captures many basic concepts. Most

² Our experiments on these real-world graphs yield similar values for the spectral transitivity.

126 importantly, X_i is internally dense in edges. Let $|X_i| = k$. Then $\Omega(k^2)$ entries in X_i are
 127 $\Omega(1/k)$, which (by averaging) implies that a constant fraction of X_i involves vertices of degree
 128 $\Theta(k)$. Thus, a constant fraction of X_i vertices have a constant fraction of their neighbors
 129 in X_i . Moreover, the submatrix of every neighborhood in X_i is also uniform. This is quite
 130 consistent with the typical notion of a social network community.

131 Crucially, Theorem 1.4 gives a condition on the *internal* structure of the decomposition.
 132 This addresses a key weakness of the Cheeger inequality.

133 **The coverage condition:** It is natural to measure the "mass" of a matrix by the squared
 134 Frobenius norm. The clusters of spectral triadic decomposition of Theorem 1.4 capture a
 135 constant fraction of this squared norm. This is consistent with the fact that a constant
 136 fraction of the edges in a real-world graph are *not* community edges [22, 14, 19, 29]. Any
 137 decomposition into communities would avoid these "long-range" edges, excluding a constant
 138 fraction of the matrix mass.

139 **Robustness to noise:** Taking the above point further, the non-community edges are
 140 often modeled as stochastic (or noisy). The underlying cluster structure of a real-world graph
 141 is robust to such perturbations. Adding (say) an Erdős-Rényi graph with $\Theta(n)$ edges can
 142 only affect the spectral transitivity by a constant factor (by changing the Frobenius norm).
 143 Theorem 1.4 would only be affected by constant factors. Note that the spectral gap, on the
 144 other hand, can dramatically increase by such noise.

145 **Spectral graph theory inspired by real-world graphs:** We consider Theorem 1.4 as
 146 opening up a new direction in spectral graph theory. At a mathematical level, Theorem 1.4
 147 is like a Cheeger inequality, where a spectral condition implies a graph theoretic property.
 148 But all aspects of Theorem 1.4 (the notion of spectral transitivity and the properties of the
 149 decomposition) are inspired by the observed properties of real-world graphs.

150 2 Related Work

151 Spectral graph theory is a deep field of study with much advancement over the past two
 152 decades. We refer the readers to the classic textbook by Chung [7], and the tutorial [31] and
 153 lecture notes [30] by Spielman.

154 The cluster structure of real-world networks has attracted attention from the early days
 155 of network science [12, 23]. Fortunato's (somewhat dated) survey on community detection
 156 has details of the key results [10]. There is no definitive model for social networks, but it
 157 is generally accepted that they have many dense clusters with sparse connections between
 158 them [5, 21, 29]. The study of triangles and neighborhood density goes back to the early days
 159 of social science theory [16, 17, 4, 8]. Early network science papers popularized the notion of
 160 clustering coefficients and transitivity as useful measures [36]. The use of triangles to find
 161 such clusters is a more recent development in network science. A number of contemporary
 162 results explicit use triangle information for algorithmic purposes [28, 33, 3, 34]. Our main
 163 theorem is inspired by these applications.

164 While the Cheeger inequality by itself is not useful for real-world graph clustering, local
 165 versions of spectral clustering are extremely useful [32, 2]. We stress that these results do
 166 not relate the graph spectrum to the partitions. But the algorithm is inspired by the proof
 167 of the Cheeger inequality. Many results on the cluster structure of real-world graphs [21, 13]
 168 use the Personalized PageRank method [2]. Some local partitioning methods yield bounds
 169 on the internal structure of clusters [20, 26, 27].

170 Most relevant to our work is the result of Gupta, Roughgarden, and Seshadhri [15]. They
 171 prove a decomposition theorem for triangle-rich graphs, as measured by graph transitivity.

172 Their main result shows that a triangle-dense graph can be clustered into dense clusters.
 173 The results of [15] do not have any spectral connection, nor do they provide the kind of
 174 uniformity or coverage bounds of Theorem 1.4. Our main insight is in generalizations of their
 175 proof technique, which leads to connections with graph spectrum. We adapt the [15] proof
 176 to deal with normalized adjacency matrix, which adds many complications because of the
 177 non-uniformity of entries.

178 **3 Preliminaries**

179 We use V, E, T to denote the sets of vertices, edges, and triangles of G , respectively. For any
 180 subgraph H of G , we use V_H, E_H, T_H to denote the corresponding sets within H . For any
 181 edge e , let $T_H(e)$ denote the set of triangles in H containing e .

182 For any vertex v , let d_v denote the degree of v (in G).

183 We first define the notion of *weights* for edges and triangles. We will think of edges and
 184 triangles as unordered sets of vertices.

185 ► **Definition 3.1.** For any edge $e = (u, v)$, define the weight $\text{wt}(e)$ to be $\frac{1}{d_u d_v}$. For any
 186 triangle $t = (u, v, w)$, define the weight $\text{wt}(t)$ to be $\frac{1}{d_u d_v d_w}$.

187 For any set S consisting solely of edges or triangles, define $\text{wt}(S) = \sum_{s \in S} \text{wt}(s)$.

188 We state some basic facts that relate the sum of weights to sum of eigenvalue powers.
 189 Let $S \subset V$ be any subset of vertices, and let $\mathcal{A}|_S$ denote the submatrix of \mathcal{A} restricted to S .
 190 We use $\lambda_i(S)$ to denote the i th largest eigenvalue of the symmetric submatrix $\mathcal{A}|_S$. Abusing
 191 notation, we use E_S and T_S to denote the edges and triangles contained in the graph induced
 192 on S .

193 ► **Claim 3.2.** $\sum_{i \leq |S|} \lambda^2(S)_i = 2 \sum_{e \in E(S)} \text{wt}(e)$

194 **Proof.** By the properties of the Frobenius norm of matrices, $\sum_{i \leq |S|} \lambda_i^2 = \sum_{s, t \in S} \mathcal{A}_{st}^2$. Note
 195 that $\mathcal{A}_{st} = A_{st}/\sqrt{d_s d_t}$. Hence, $\sum_{s, t} \mathcal{A}_{s, t}^2 = 2 \sum_{e=(u, v) \in E(S)} 1/d_u d_v$. (We get a 2-factor
 196 because each edge (u, v) appears twice in the adjacency matrix.) ◀

197 ► **Claim 3.3.** $\sum_{i \leq |S|} \lambda^3(S)_i = 6 \sum_{t \in T(S)} \text{wt}(t)$.

198 **Proof.** Note that $\sum_{i \leq |S|} \lambda^3(S)_i$ is the trace of $(\mathcal{A}|_S)^3$. The diagonal entry $(\mathcal{A}|_S)_{ii}^3$ is precisely
 199 $\sum_{s \in S} \sum_{s' \in S} \mathcal{A}_{is} \mathcal{A}_{ss'} \mathcal{A}_{s'i}$. Note that $\mathcal{A}_{is} \mathcal{A}_{ss'} \mathcal{A}_{s'i}$ is non-zero iff (i, s, s') form a triangle. In
 200 that case, $\mathcal{A}_{is} \mathcal{A}_{ss'} \mathcal{A}_{s'i} = 1/\sqrt{d_i d_s} \cdot 1/\sqrt{d_s d_{s'}} \cdot 1/\sqrt{d_{s'} d_i} = \text{wt}((i, s, s'))$. We conclude that
 201 $(\mathcal{A}|_S)_{ii}^3$ is $2 \sum_{t \in T(S), t \ni i} \text{wt}(t)$. (There is a 2 factor because every triangle is counted twice.)

202 Thus, $\sum_{i \leq n} \lambda^3(S)_i = \sum_i 2 \sum_{t \in T, t \ni i} \text{wt}(t) = 2 \sum_{t \in T} \sum_{i \in t} \text{wt}(t) = 6 \sum_{t \in T} \text{wt}(t)$. (The
 203 final 3 factor appears because a triangle contains exactly 3 vertices.) ◀

204 ► **Claim 3.4.** $\sum_{t \in T(S)} \text{wt}(t) \leq \|\mathcal{A}|_S\|_2^2/6$.

205 **Proof.** By Claim 3.3 $\sum_{t \in T(S)} \text{wt}(t) = \sum_{i \leq |S|} \lambda^3(S)_i/6$. The maximum eigenvalue of \mathcal{A}
 206 is 1, and since $\mathcal{A}|_S$ is a submatrix, $\lambda(S)_1 \leq 1$ (Cauchy's interlacing theorem). Thus,
 207 $\sum_{i \leq |S|} \lambda^3(S)_i \leq \sum_{i \in |S|} \lambda^2(S)_i = \|\mathcal{A}|_S\|_2^2$. ◀

208 As a direct consequence of the previous claims applied on \mathcal{A} , we get the following
 209 characterization of the spectral triadic content in terms of the weights.

210 ► **Lemma 3.5.** $\tau = \frac{3 \sum_{t \in T} \text{wt}(t)}{\sum_{e \in E} \text{wt}(e)}$.

1:6 A spectral theorem on the cluster structure of real-world graphs

211 While the following bound is not necessary for our main result, it is instructive to see the
212 largest possible value of the spectral transitivity.

213 ► **Lemma 3.6.** *Consider normalized adjacency matrices \mathcal{A} with n vertices. The maximum
214 value of $\tau(\mathcal{A})$ is $1 - 1/(n-1)$. This value is attained for the unique strongly 1-uniform
215 matrix, the normalized adjacency matrix of the n -clique.*

Proof. First, consider the normalized adjacency matrix \mathcal{A} of the n -clique. All off-diagonal
entries are precisely $1/(n-1)$ and \mathcal{A} can be expressed as $(n-1)^{-1}(\mathbf{1}\mathbf{1}^T - I)$. The matrix \mathcal{A}
is 1-regular. The largest eigenvalue is 1 and all the remaining eigenvalues are $-1/(n-1)$.
Hence, $\sum_i \lambda_i^3 = 1 - (n-1)/(n-1)^3 = 1 - 1/(n-1)^2$. The sum of squares of eigenvalue is
 $\sum_i \lambda_i^2 = 1 + (n-1)/(n-1)^2 = 1 + 1/(n-1)$. Dividing,

$$\frac{\sum_{i \leq n} \lambda_i^3}{\sum_{i \leq n} \lambda_i^2} = 1 - 1/(n-1).$$

216 Since the matrix has zero diagonal, the trace $\sum_i \lambda_i$ is zero. We will now prove the following
217 claim.

218 ► **Claim 3.7.** Consider any sequence of numbers $1 = \lambda_1 \geq \lambda_2 \dots \geq \lambda_n$ such that $\forall i, |\lambda_i| \leq 1$
219 and $\sum_i \lambda_i = 0$. If $\sum_i \lambda_i^3 \geq (1 - 1/(n-1)) \sum_i \lambda_i^2$, then $\forall i > 1, \lambda_i = -1/(n-1)$.

220 **Proof.** Let us begin with some basic manipulations.

$$221 \quad \sum_i \lambda_i^3 \geq [1 - 1/(n-1)] \sum_i \lambda_i^2 \quad (2)$$

$$222 \quad \implies 1 + \sum_{i>1} \lambda_i^3 \geq [1 - 1/(n-1)] \cdot (1 + \sum_{i>1} \lambda_i^2)$$

$$223 \quad \implies \sum_{i>1} \lambda_i^3 \geq [1 - 1/(n-1)] \sum_{i>1} \lambda_i^2 - 1/(n-1). \quad (3)$$

224 For $i > 1$, define $\delta_i := \lambda_i + 1/(n-1)$. Note that $\sum_{i>1} \lambda_i = -1$, so $\sum_{i>1} \delta_i = 0$. Moreover,
225 $\forall i > 1, \delta_i \leq 1 + 1/(n-1)$. We plug in $\lambda_i = \delta_i - 1/(n-1)$ in (3).

$$226 \quad \sum_{i>1} \left[\delta_i - 1/(n-1) \right]^3 \geq [1 - 1/(n-1)] \sum_{i>1} \left[\delta_i - 1/(n-1) \right]^2 - 1/(n-1)$$

$$227 \quad \implies \sum_{i>1} \left[\delta_i^3 - 3\delta_i^2/(n-1) + 3\delta_i/(n-1)^2 - 1/(n-1)^3 \right]$$

$$228 \quad \geq [1 - 1/(n-1)] \sum_{i>1} \left[\delta_i^2 - 2\delta_i/(n-1) + 1/(n-1)^2 \right] - 1/(n-1).$$

229 Recall that $\sum_{i>1} \delta_i = 0$. Hence, we can simplify the above inequality.

$$230 \quad \sum_{i>1} \delta_i^3 - (3/(n-1)) \sum_{i>1} \delta_i^2 - 1/(n-1)^2$$

$$231 \quad \geq [1 - 1/(n-1)] \sum_{i>1} \delta_i^2 + 1/(n-1) - 1/(n-1)^2 - 1/(n-1)$$

$$232 \quad \implies \sum_{i>1} \delta_i^3 \geq [1 + 2/(n-1)] \sum_{i>1} \delta_i^2. \quad (\text{Canceling terms and rearranging})$$

233 Since $\delta_i \leq (1 + 1/(n-1))$, we get that $\sum_{i>1} \delta_i^3 \leq [1 + 1/(n-1)] \sum_{i>1} \delta_i^2$. Combining with
234 the above inequality, we deduce that $[1 + 2/(n-1)] \sum_{i>1} \delta_i^2 \leq [1 + 1/(n-1)] \sum_{i>1} \delta_i^2$. This
235 can only happen if $\sum_{i>1} \delta_i^2$ is zero, implying all δ_i values are zero. Hence, for all $i > 1$,
236 $\lambda_i = -1/(n-1)$. ◀

237 With this claim, we conclude that any matrix \mathcal{A} maximizing the ratio of cubes and squares
 238 of eigenvalues has a fixed spectrum. It remains to prove that a unique normalized adjacency
 239 matrix has this spectrum. We use the rotational invariance of the Frobenius norm: sum of
 240 squares of entries of \mathcal{A} is the same as the sum of squares of eigenvalues. Thus,

$$241 \quad \sum_{(u,v) \in E} \frac{2}{d_u d_v} = 1 + \frac{1}{n-1} = \frac{n}{n-1}. \quad (4)$$

242 Observe that $\frac{2}{d_u d_v} \geq 1/(d_u(n-1)) + 1/(d_v(n-1))$, since all degrees are at most $n-1$.
 243 Summing this inequality over all edges,

$$244 \quad \sum_{(u,v) \in E} \frac{2}{d_u d_v} \geq \sum_{v \in V} \sum_{u \in N(v)} \frac{1}{d_v(n-1)} = \sum_{v \in V} \frac{d_v}{d_v(n-1)} = \frac{n}{n-1}. \quad (5)$$

245 Hence, for (4) to hold, for all edges (u, v) , we must have the equality $\frac{2}{d_u d_v} = 1/(d_u(n-1)) +$
 246 $1/(d_v(n-1))$. That implies that for all edge (u, v) , $d_u = d_v = n-1$. So all vertices have
 247 degree $(n-1)$, and the graph is an n -clique. ◀

248 We will need the following “reverse Markov inequality” for some intermediate proofs.

249 ▶ **Lemma 3.8.** *Consider a random variable Z taking values in $[0, b]$. If $\mathbf{E}[Z] \geq \sigma b$, then*
 250 $\Pr[Z \geq \sigma b/2] \geq \sigma/2$.

251 **Proof.** In the following calculations, we will upper bound the conditional expectation by the
 252 maximum value (under that condition).

$$253 \quad \sigma b \leq \mathbf{E}[Z] = \Pr[Z \geq \sigma b/2] \cdot \mathbf{E}[Z|Z \geq \sigma b/2] + \Pr[Z$$

$$254 \quad \leq \sigma b/2] \cdot \mathbf{E}[Z|Z \leq \sigma b/2] \quad (7)$$

$$255 \quad \leq \Pr[Z \geq \sigma b/2] \cdot b + \sigma b/2 \quad (8)$$

256 We rearrange to complete the proof. ◀

257 4 Cleaned graphs and extraction

258 For convenience, we set $\varepsilon = \tau/6$.

259 ▶ **Definition 4.1.** *A connected subgraph H is called clean if $\forall e \in E(H)$, $\text{wt}(T_H(e)) \geq \varepsilon \text{wt}(e)$.*

Algorithm 1 Extract(H)

-
- 1: Pick $v \in V(H)$ that minimizes d_v .
 - 2: Construct the set $L := \{u | (u, v) \in E(H), d_u \leq 2\varepsilon^{-1}d_v\}$ (L is the set of low degree neighbors of v in H .)
 - 3: For every vertex $w \in V(H)$, define ρ_w to be the total weight of triangles of the form (w, u, u') where $u, u' \in L$.
 - 4: Sort the vertices in decreasing order of ρ_w , and construct the “sweep cut” C to be the smallest set satisfying $\sum_{w \in C} \rho_w \geq (1/2) \sum_{w \in V(H)} \rho_w$.
 - 5: Output $X := \{v\} \cup L \cup C$.
-

260 The main theorem of this section follows.

► **Theorem 4.2.** *Suppose the subgraph H is connected and clean. Let X denote the output of the procedure $\text{Extract}(H)$. Then*

$$\sum_{t \in T(H), t \subseteq X} \text{wt}(t) \geq (\varepsilon^8/2000) \sum_{t \in T(H), t \cap X \neq \emptyset} \text{wt}(t)$$

261 (The triangle weight contained inside X is a constant fraction of the triangle weight incident
262 to X .)

263 Moreover, $\mathcal{A}|_X$ is strongly $\delta\varepsilon^{12}$ -uniform.

264 We will need numerous intermediate claims to prove this theorem. We use v , L , and C as
265 defined in $\text{Extract}(H)$. We use N to denote the neighborhood of v in H . Note that $L \subseteq N$.

266 For any vertex $u \in N$, we define the set of partners $P(u)$ to be $\{w : (u, v, w) \in T_H\}$.

267 The following lemma is an important tool in our analysis.

268 ► **Lemma 4.3.** *For any $u \in N$, $\sum_{w \in P(u) \cap L} d_w^{-1} \geq \varepsilon/2$.*

269 **Proof.** Let $e = (u, v)$. Since H is clean, $\text{wt}(T_H(e)) \geq \varepsilon \text{wt}(e)$. Expanding out the definition
270 of weights,

$$271 \sum_{w: (u, v, w) \in T_H} \frac{1}{d_u d_v d_w} \geq \frac{\varepsilon}{d_u d_v} \implies \sum_{w \in P(u)} d_w^{-1} \geq \varepsilon. \quad (9)$$

272 Note that L (as constructed in $\text{Extract}(H)$) is the subset of N consisting of vertices with
273 degree at most $2\varepsilon^{-1}d_v$. For $w \in N \setminus L$, we have the lower bound $d_w \geq 2\varepsilon^{-1}d_v$. Hence,

$$274 \sum_{w \in N \setminus L} d_w^{-1} \leq |N \setminus L|(\varepsilon/2)d_v^{-1} \leq d_v \times (\varepsilon/2)d_v^{-1} = \varepsilon/2. \quad (10)$$

275 In the calculation below, we split the sum of (9) into the contribution from L and from
276 outside L . We apply (10) to bound the latter contribution.

$$277 \varepsilon \leq \sum_{w \in P(u)} d_w^{-1} \leq \sum_{w \in P(u) \cap L} d_w^{-1} + \sum_{w \in N \setminus L} d_w^{-1} \leq \sum_{w \in P(u) \cap L} d_w^{-1} + \varepsilon/2. \quad (11)$$

278 ◀

279 ► **Claim 4.4.** $|L| \geq \varepsilon d_v/2$

280 **Proof.** Since H is connected, there must exist some edge $e = (u, v) \in E(H)$. By Lemma 4.3,
281 $\sum_{w \in P(u) \cap L} d_w^{-1} \geq \varepsilon/2$. Hence, $\sum_{w \in L} d_w^{-1} \geq \varepsilon/2$. Since v is the vertex in $V(H)$ minimizing
282 d_v , for any $w \in V(H)$, $d_w \geq d_v$. Thus,

$$283 \varepsilon/2 \leq \sum_{w \in L} d_w^{-1} \leq \sum_{w \in L} d_v^{-1} = |L|d_v^{-1}. \quad (12)$$

284 ◀

285 ► **Claim 4.5.** $\sum_{e \in E(H), e \subseteq L} \text{wt}(e) \geq \varepsilon^2/8$.

286 **Proof.** By Lemma 4.3, $\forall w \in L$, $\sum_{w' \in P(w) \cap L} d_{w'}^{-1} \geq \varepsilon/2$. We multiply both sides by d_w^{-1} and
287 sum over all $w \in L$.

$$288 \sum_{w \in L} \sum_{w' \in P(w) \cap L} (d_w d_{w'})^{-1} \geq (\varepsilon/2) \sum_{w' \in L} d_{w'}^{-1}. \quad (13)$$

289 By Lemma 4.3, $\sum_{w' \in L} d_{w'}^{-1} \geq \varepsilon/2$. Note that $w' \in P(w)$ only if $(w, w') \in E(H)$. Hence,

290 $\sum_{w \in L} \sum_{w' \in L, (w, w') \in E(H)} \text{wt}((w, w')) \geq \varepsilon^2/4$. Note that the summation counts all edges
291 twice, so we divide by 2 to complete the proof. ◀

292 We now come to the central calculations of the main proof. Recall, from the description
 293 of **Extract**, that ρ_w is the total triangle weight of the triangles (w, u, u') , where $u, u' \in L$.
 294 We will prove that $\sum_w \rho_w$ is large; moreover, there are a few entries that dominate the sum.
 295 The latter bound is crucial to arguing that the sweep set C is not too large.

296 \triangleright **Claim 4.6.** $\sum_{w \in V(H)} \rho_w \geq \varepsilon^3/8$.

297 **Proof.** Note that $\sum_{w \in V(H)} \rho_w$ is equal to $\sum_{e \in E(H), e \subset L} \text{wt}(T_H(e))$. Both these expressions
 298 give the total weight of all triangles in H that involve two vertices in L . Since H is
 299 clean, for all edges $e \in E(H)$, $\text{wt}(T_H(e)) \geq \varepsilon \text{wt}(e)$. Hence, $\sum_{e \in E(H), e \subset L} \text{wt}(T_H(e)) \geq$
 300 $\varepsilon \sum_{e \in E(H), e \subset L} \text{wt}(e)$. Applying Claim 4.5, we can lower bound the latter by $\varepsilon^3/8$. \blacktriangleleft

301 We now show that a few ρ_w values dominate the sum, using a somewhat roundabout
 302 argument. We upper bound the sum of square roots.

303 \triangleright **Claim 4.7.** $\sum_{w \in V(H)} \sqrt{\rho_w} \leq 2\varepsilon^{-1} \sqrt{d_v}$

304 **Proof.** Let c_w be the number of vertices in L that are neighbors (in H) of w . Note that for
 305 any triangle (u, u', w) where $u, u' \in L$, both u and u' are common neighbors of w and v . The
 306 number of triangles (u, u', w) where $u, u' \in L$ is at most c_w^2 . The weight of any triangle in H
 307 is at most d_v^{-3} , since d_v is the lowest degree (in G) of all vertices in H . As a result, we can
 308 upper bound $\rho_w \leq d_v^{-3} c_w^2$.

309 Taking square roots and summing over all vertices,

$$310 \sum_{w \in V(H)} \sqrt{\rho_w} \leq d_v^{-3/2} \sum_{w \in V(H)} c_w \quad (14)$$

311 Note that $\sum_{w \in V(H)} c_w$ is exactly the sum over $u \in L$ of the degrees of u in the subgraph
 312 H . (Every edge incident to $u \in L$ gives a unit contribution to the sum $\sum_{w \in V(H)} c_w$.) By
 313 definition, every vertex in L has degree in H at most $2\varepsilon^{-1} d_v$. The size of L is at most d_v .

314 Hence, $\sum_{w \in V(H)} c_w \leq 2\varepsilon^{-1} d_v^2$. Plugging into (14), we deduce that $\sum_{w \in V(H)} \sqrt{\rho_w} \leq$
 315 $2\varepsilon^{-1} \sqrt{d_v}$. \blacktriangleleft

316 We now prove that the sweep cut C is small, which is critical to proving Theorem 4.2.

317 \triangleright **Claim 4.8.** $|C| \leq 144\varepsilon^{-5} d_v$.

318 **Proof.** For convenience, let us reindex vertices so that $\rho_1 \geq \rho_2 \geq \rho_3 \dots$. Let $r \leq n$ be
 319 an arbitrary index. Because we index in non-increasing order, note that $\sum_{j \leq n} \rho_j \geq r \rho_r$.
 320 Furthermore, $\forall j > r, \rho_j \leq \rho_r$.

$$321 \sum_{j > r} \rho_j \leq \sqrt{\rho_r} \sum_{j > r} \sqrt{\rho_j} \leq \sqrt{\frac{\sum_{j \leq n} \rho_j}{r}} \sum_{j \leq n} \sqrt{\rho_j} = \left[\frac{\sum_{j \leq n} \sqrt{\rho_j}}{\sqrt{r} \cdot \sqrt{\sum_{j \leq n} \rho_j}} \right] \sum_{j \leq n} \rho_j \quad (15)$$

322 Observe that Claim 4.7 gives an upper bound on the numerator, while Claim 4.6 gives a
 323 lower bound on (a term in) the denominator. Plugging those bounds in (15),

$$324 \sum_{j > r} \rho_j \leq \frac{2\varepsilon^{-1} \sqrt{d_v}}{\sqrt{r} \cdot \varepsilon^{3/2} / \sqrt{8}} \sum_{j \leq n} \rho_j \leq \frac{1}{\sqrt{r}} \cdot \frac{6\sqrt{d_v}}{\varepsilon^{5/2}} \cdot \sum_{j \leq n} \rho_j. \quad (16)$$

325 Suppose $r > 144\varepsilon^{-5} d_v$. Then $\sum_{j > r} \rho_j < (1/2) \sum_{j \leq n} \rho_j$. The sweep cut C is constructed
 326 with the smallest value of r such that $\sum_{j > r} \rho_j < (1/2) \sum_{j \leq n} \rho_j$. Hence, $|C| \leq 144\varepsilon^{-5} d_v$. \blacktriangleleft

1:10 A spectral theorem on the cluster structure of real-world graphs

327 An additional technical claim we need bounds the triangle weight incident to a single
328 vertex.

329 \triangleright **Claim 4.9.** For all vertices $u \in V(H)$, $\text{wt}(T_H(u)) \leq (2d_v)^{-1}$.

Proof. Consider edge $(u, w) \in E(H)$. We will prove that $\text{wt}(T_H((u, w))) \leq d_u^{-1}d_v^{-1}$. Recall that d_v is the smallest degree among vertices in H . Furthermore, $|T_H((u, w))| \leq d_w$, since the third vertex in a triangle containing (u, w) is a neighbor of w .

$$\text{wt}(T_H((u, w))) = \sum_{z:(z,u,w) \in T(H)} \frac{1}{d_u d_w d_z} \leq \frac{1}{d_u d_v} \sum_{z:(z,u,w) \in T(H)} \frac{1}{d_w} \leq \frac{1}{d_u d_v} \times \frac{d_w}{d_w} = \frac{1}{d_u d_v}$$

330 We now bound $\text{wt}(T_H(u))$ by summing over all neighbors of u in H .

$$\begin{aligned} 331 \quad \text{wt}(T_H(u)) &= (1/2) \sum_{w:(u,w) \in E(H)} \text{wt}(T_H((u, w))) \\ 332 \quad &\leq (1/2) \sum_{w:(u,w) \in E(H)} \frac{1}{d_u d_v} = \frac{1}{2d_v} \sum_{w:(u,w) \in E(H)} \frac{1}{d_u} \\ 333 \quad &\leq \frac{1}{2d_v} \times \frac{d_u}{d_u} = \frac{1}{2d_v}. \end{aligned}$$

334 ◀

335 4.1 The proof of Theorem 4.2

336 **Proof.** (of Theorem 4.2) By construction of X as $\{v\} \cup L \cup C$, all the triangles of the form
337 (w, u, u') , where $w \in C$ and $u, u' \in L$, are contained in X . The total weight of such triangles
338 is at least $\sum_{v \leq n} \rho_v/2$, by the construction of C . By Claim 4.6, $\sum_{v \leq n} \rho_v/2 \geq \varepsilon^3/16$.

339 Let us now bound that total triangle weight incident to X in \bar{H} . Observe that $|X| =$
340 $1 + |L| + |C|$ which is at most $1 + d_v + \varepsilon^{-5}144d_v$, by Claim 4.8. We can further bound
341 $|X| \leq \varepsilon^{-5}146d_v$. By Claim 4.9, the total triangle weight incident to a vertex is at most
342 $(2d_v)^{-1}$. Hence, the total triangle weight incident to all of X is at most $73\varepsilon^{-5}$.

343 Thus, the triangle weight contained in X is at least $\frac{\varepsilon^3/16}{73\varepsilon^{-5}}$ times the triangle weight
344 incident to X . The ratio is at least $\varepsilon^8/2000$, completing the proof of the first statement.

345 **Proof of uniformity of $\mathcal{A}|_X$:** We first prove a lower bound on the uniformity of $\mathcal{A}|_X$.
346 For convenience, let B denote the set $\{e \in E(H), e \subseteq L\}$. By Claim 4.5, $\sum_{e \in B} \text{wt}(e) \geq \varepsilon^2/8$.
347 There are at most $\binom{d_v}{2} \leq d_v^2/2$ edges in B . For every edge e , $\text{wt}(e) \leq 1/d_v^2$. Let k denote the
348 number of edges in B whose weight is at least $\varepsilon^2/16$.

$$\begin{aligned} 349 \quad \frac{\varepsilon^2}{8} &\leq \sum_{\substack{e \in B \\ \text{wt}(e) \leq \varepsilon^2 d_v^{-2}/16}} \text{wt}(e) + \sum_{\substack{e \in B \\ \text{wt}(e) \geq \varepsilon^2 d_v^{-2}/16}} \text{wt}(e) \\ 350 \quad &\leq |B| \times \varepsilon^2 d_v^{-2}/16 + k d_v^{-2} \\ 351 \quad &\leq d_v^2 \times \varepsilon^2 d_v^{-2}/16 + k d_v^{-2} \\ 352 \quad &= \varepsilon^2/16 + k d_v^{-2}. \end{aligned}$$

353 Rearranging, $k \geq \varepsilon^2 d_v^2/16$.

354 Hence, there are at least $\varepsilon^2 d_v^2/16$ edges contained in X with weight at least $\varepsilon^2 d_v^{-2}/16$.
355 Consider the random variable Z that is the weight of a uniform random edge contained in
356 X . Since $|X| \leq \varepsilon^{-5}144d_v$, the number of edges in X is at most $\varepsilon^{-10}(144)^2 d_v^2$. So,

$$357 \quad \mathbf{E}[Z] \geq \frac{\varepsilon^2 d_v^2/16}{\varepsilon^{-10}(144)^2 d_v^2} \times \varepsilon^2 d_v^{-2}/16 \geq 2\delta \varepsilon^{14} d_v^{-2}. \quad (17)$$

358 The maximum value of Z is the largest possible weight of an edge in $E(H)$, which is at most
 359 d_v^{-2} . Applying the reverse Markov bound of Lemma 3.8, $\Pr[Z \geq \delta \varepsilon^{14} d_v^{-2}] \geq \delta \varepsilon^{14}$. Thus, an
 360 ε^{14} fraction of edges in $|X|$ have weight at least $\delta \varepsilon^{14} d_v^{-2} \geq \delta \varepsilon^c / |X|^2$. Moreover, every edge
 361 has weight at most $d_v^{-2} \leq 1 / (\delta \varepsilon^c |X|^2)$. So we prove the uniformity of $\mathcal{A}|_X$.

362 The largest possible weight for any edge in $E(H)$ is d_v^{-2} . The size of $|X|$ is at least d_v
 363 and at most $\varepsilon^{-5} 144 d_v$. Hence, $\mathcal{A}|_X$ is at least $\delta \varepsilon^{12}$ -uniform.

364 **Proof of strong uniformity:** For strong uniformity, we need to repeat the above
 365 argument within neighborhoods in X . We prove in the beginning of this proof that the total
 366 triangle weight inside X is at least $\varepsilon^3 / 16$. We also proved that $|X| \leq 146 \varepsilon^{-5} d_v$. Consider
 367 the random variable Z that is the triangle weight contained in X incident to a uniform
 368 random vertex in X . Note that $\mathbf{E}[Z] \geq (\varepsilon^3 / 16) / (146 \varepsilon^{-5} d_v) \geq 2 \delta' \varepsilon^8 d_v^{-1}$. By Claim 4.9, Z
 369 is at most $(2d_v)^{-1}$. Applying Lemma 3.8, $\Pr[Z \geq \delta' \varepsilon^8 d_v^{-1}] \geq \delta \varepsilon^8$. This means that at least
 370 $\delta' \varepsilon^8 |X|$ vertices in X are incident to at least $\delta' \varepsilon^8 d_v^{-1}$ triangle weight inside X .

371 Consider any such vertex u . Let $N(u)$ be the neighborhood of u in X . Every edge e in
 372 $N(u)$ forms a triangle with u with weight $\text{wt}(e) / d_u$. Hence, noting that $d_u \geq d_v$,

$$373 \quad \sum_{e \subseteq N(u)} \text{wt}(e) d_u^{-1} \geq \delta' \varepsilon^8 d_v^{-1} \implies \sum_{e \subseteq N(u)} \text{wt}(e) \geq \delta' \varepsilon^8. \quad (18)$$

374 There are at most $|X|^2 \leq \varepsilon^{-10} (146)^2 d_v^2$ edges in $N(u)$. Let Z denote the weight of a uniform
 375 random edge in $N(u)$. Note that $\mathbf{E}[Z] \geq \delta' \varepsilon^8 / (\varepsilon^{-10} (146)^2 d_v^2) \geq 2 \delta \varepsilon^{18} d_v^{-2}$. The maximum
 376 weight of an edge is at most d_v^{-2} . By Lemma 3.8, at least $\delta \varepsilon^{18}$ fraction of edges in $N(u)$ have
 377 a weight of at least $\delta \varepsilon^{18} d_v^{-2}$. Since $|N(u)| \leq |X| \leq \varepsilon^{-5} 146 d_v$, this implies that $N(u)$ is also
 378 $\delta \varepsilon^c$ -uniform. Hence, we prove strong uniformity as well.

379

◀

380 5 Obtaining the decomposition

■ Algorithm 2 Decompose(G)

-
- 1: Initialize \mathbf{X} to be an empty family of sets, and initialize subgraph $H = G$.
 - 2: **while** H is non-empty **do**
 - 3: **while** H is not clean **do**
 - 4: Remove an edge $e \in E(H)$ from H such that $\text{wt}(T_H(e)) < (\varepsilon) \text{wt}(e)$.
 - 5: **end while**
 - 6: Add output $\text{Extract}(H)$ to \mathbf{X} .
 - 7: Remove these vertices from H .
 - 8: **end while**
 - 9: Output \mathbf{X} .
-

381 We first describe the algorithm that obtains the decomposition promised in Theorem 1.4.

382 We partition all the triangles of G into three sets depending on how they are affected by
 383 $\text{Decompose}(G)$. (i) The set of triangles removed by the cleaning step of Step 4, (ii) the set
 384 of triangles contained in some $X_i \in \mathbf{X}$, or (iii) the remaining triangles. Abusing notation,
 385 we refer to these sets as T_C , T_X , and T_R respectively. Note that the triangles of T_R are the
 386 triangles "cut" when X_i is removed.

387 ▷ **Claim 5.1.** $\text{wt}(T_C) \leq (\tau/6) \sum_{e \in E} \text{wt}(e)$.

1:12 A spectral theorem on the cluster structure of real-world graphs

388 **Proof.** Consider an edge e removed at Step 4 of **Decompose**. Recall that ε is set to $\tau/6$. At
 389 that removal, the total weight of triangles removed (cleaned) is at most $(\tau/6)\text{wt}(e)$. An edge
 390 can be removed at most once, so the total weight of triangles removed by cleaning is at most
 391 $(\tau/6) \sum_{e \in E} \text{wt}(e)$. ◀

392 **Proof.** (of Theorem 1.4) Let us denote by H_1, H_2, \dots, H_k the subgraphs of which **Extract**
 393 is called. Let the output of **Extract**(H_i) be denoted X_i . By the uniformity guarantee of
 394 Theorem 4.2, each $\mathcal{A}|_{X_i}$ is $\delta\tau^c$ -uniform.

395 It remains to prove the coverage guarantee. We now sum the bound of Theorem 4.2
 396 over all X_i . (For convenience, we expand out ε as $\tau/6$ and let δ' denote a sufficiently small
 397 constant.)

$$398 \sum_{i \leq k} \sum_{t \in T(H), t \subseteq X} \text{wt}(t) \geq (\delta'\tau^8) \sum_{i \leq k} \sum_{t \in T(H), t \cap X \neq \emptyset} \text{wt}(t). \quad (19)$$

399 The LHS is precisely $\text{wt}(T_X)$. Note that a triangle appears at most once in the double
 400 summation in the RHS. That is because if $t \cap X_i \neq \emptyset$, then t is removed when X_i is removed.
 401 Since H_i is always clean, the triangles of T_C cannot participate in this double summation.
 402 Hence, the RHS summation is $\text{wt}(T_X) + \text{wt}(T_R)$ and we deduce that

$$403 \text{wt}(T_X) \geq \delta'\tau^8(\text{wt}(T_X) + \text{wt}(T_R)) \quad (20)$$

404 Note that $\text{wt}(T_c) + \text{wt}(T_x) + \text{wt}(T_r) = \sum_{t \in T} \text{wt}(t)$. There is where the definition of τ makes
 405 its appearance. By Lemma 3.5, we can write the above equality as $\text{wt}(T_c) + \text{wt}(T_x) + \text{wt}(T_r) =$
 406 $(\tau/3) \sum_{e \in E} \text{wt}(e)$. Applying Claim 5.1, (20), and the relation of edge weights to the Frobenius
 407 norm (Claim 3.2),

$$408 (\delta'\tau^8)^{-1} \text{wt}(T_X) \geq (\tau/6) \sum_{e \in E} \text{wt}(e) \implies \text{wt}(T_X) \geq \delta\tau^c \|\mathcal{A}\|_2^2 \quad (\text{by Claim 3.2}) \quad (21)$$

409 By Claim 3.4, $\sum_{i \leq k} \|\mathcal{A}|_{X_i}\|_2^2 \geq \text{wt}(T_X)$, completing the proof of the coverage bound. ◀

410 6 Algorithmics and implementation

411 We discuss theoretical and practical implementations of the procedures computing the
 412 decomposition of Theorem 1.4. The main operation required is a triangle enumeration of
 413 G ; there is a rich history of algorithms for this problem. The best known bound for sparse
 414 graph is the classic algorithm of Chiba-Nishizeki that enumerates all triangles in $O(m\alpha)$
 415 time, where α is the graph degeneracy.

416 We first provide a formal theorem providing a running time bound. We do not explicitly
 417 describe the implementation through pseudocode, and instead explain the main details in
 418 the proof.

419 ► **Theorem 6.1.** *There is an implementation of **Decompose**(G) whose running time is
 420 $O(R + (m + n + T) \log n)$, where R is the running time of listing all triangles. The space
 421 required is $O(T)$ (where T is the triangle count).*

422 **Proof.** We assume an adjacency list representation where each list is stored in a dictionary
 423 data structure with logarithmic time operations (like a self-balancing binary tree).

424 We prepare the following data structure that maintains information about the current
 425 subgraph H . We initially set $H = G$. We will maintain all lists as hash tables so that
 426 elementary operations on them (insert, delete, find) can be done in $O(1)$ time.

- 427 1. A list of all triangles in $T(H)$ indexed by edges. Given an edge e , we can access a list of
428 triangles in $T(H)$ containing e .
- 429 2. A list of $\text{wt}(T_H(e))$ values for all edges $e \in E(H)$.
- 430 3. A list U of all (unclean) edges such that $\text{wt}(T_H(e)) < \varepsilon \text{wt}(e)$.
- 431 4. A min priority queue Q storing all vertices in $V(H)$ keyed by degree d_v . We will assume
432 pointers from v to the corresponding node in Q .

433 These data structures can be initialized by enumerating all triangles, indexing them, and
434 preparing all the lists. This can be done in $O(R)$ time.

435 We describe the process to remove an edge from H . When edge e is removed, we go over
436 all the triangles in $T(H)$ containing e . For each such triangle t and edge $e' \in t$, we remove t
437 from the triangle list of e' . We then update $\text{wt}(T_H(e'))$ by reducing it by $\text{wt}(t)$. If $\text{wt}(T_H(e'))$
438 is less than $\text{wt}(e)$, we add it to U . Finally, if the removal of e removes a vertex v from $V(H)$,
439 we remove v from the priority queue Q . Thus, we can maintain the data structures. The
440 running time is $O(|T_H(e)|)$ plus an additional $\log n$ for potentially updating Q . The total
441 running time for all edge deletes is $O(T + n \log n)$.

442 With this setup in place, we discuss how to implement **Decompose**. The cleaning operation
443 in **Decompose** can be implemented by repeatedly deleting edges from the list U , until it is
444 empty.

445 We now discuss how to implement **Extract**. We will maintain a max priority queue R
446 maintaining the values $\{\rho_w\}$. Using Q as defined earlier, we can find the vertex v of minimum
447 degree. By traversing its adjacency list in H , we can find the set L . We determine all edges
448 in L by traversing the adjacency lists of all vertices in L . For each such edge e , we enumerate
449 all triangles in H containing e . For each such triangle t and $w \in t$, we will update the value
450 of ρ_w in R .

451 We now have the total $\sum_w \rho_w$ as well. We find the sweep cut by repeatedly deleting from
452 the max priority queue R , until the sum of ρ_w values is at least half the total. Thus, we can
453 compute the set X to be extracted. The running time is $O((|X| + |E(X)| + |T(X)|) \log n)$,
454 where $E(X), T(X)$ are the set of edges and triangles incident to X .

455 Overall, the total time for all the extractions and resulting edge removals is $O((n + m +$
456 $T) \log n)$. The initial triangle enumeration takes R time. We add to complete the proof. ◀

457 ——— References ———

- 458 1 N Alon and V.D Milman. λ_1 , isoperimetric inequalities for graphs, and superconcentrators.
459 *Journal of Combinatorial Theory, Series B*, 38(1):73–88, 1985.
- 460 2 Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using pagerank vectors.
461 In *Foundations of Computer Science (FOCS)*, pages 475–486, 2006.
- 462 3 A. Benson, D. F. Gleich, and J. Leskovec. Higher-order organization of complex networks.
463 *Science*, 353(6295):163–166, 2016.
- 464 4 Ronald S. Burt. Structural holes and good ideas. *American Journal of Sociology*, 110(2):349–
465 399, 2004.
- 466 5 Deepayan Chakrabarti and Christos Faloutsos. Graph mining: Laws, generators, and al-
467 gorithms. *ACM Comput. Surv.*, 38(1):2–es, jun 2006.
- 468 6 Jeff Cheeger. *A Lower Bound for the Smallest Eigenvalue of the Laplacian*, pages 195–200.
469 Princeton University Press, 1971.
- 470 7 F. R. K. Chung. *Spectral Graph Theory*. American Mathematical Society, 1997.
- 471 8 Katherine Faust. Comparing social networks: Size, density, and local structure. *Metodoloski*
472 *zvezki*, 3:185–216, 07 2006.

- 473 9 M. Fiedler. Algebraic connectivity of graphs. *Czechoslovak Mathematical Journal*, 23(98):298–
474 305, 1973.
- 475 10 Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3):75–174, 2010.
- 476 11 Santo Fortunato and Mark E. J. Newman. 20 years of network community detection. *Nature*
477 *Physics*, 18(8):848–850, jul 2022.
- 478 12 M. Girvan and M. Newman. Community structure in social and biological networks. *Proceedings*
479 *of the National Academy of Sciences*, 99(12):7821–7826, 2002.
- 480 13 D. Gleich and C. Seshadhri. Vertex neighborhoods, low conductance cuts, and good seeds for
481 local community methods. In *SIGKDD Conference on Knowledge Discovery and Data Mining*
482 *(KDD)*, pages 597–605, 2012.
- 483 14 M. Granovetter. The strength of weak ties: A network theory revisited. *Sociological Theory*,
484 1:201–233, 1983.
- 485 15 Rishi Gupta, Tim Roughgarden, and C. Seshadhri. Decompositions of triangle-dense graphs.
486 *Innovations in Theoretical Computer Science*, pages 471–482, 2014.
- 487 16 P. Holland and S. Leinhardt. A method for detecting structure in sociometric data. *American*
488 *Journal of Sociology*, 76:492–513, 1970.
- 489 17 Paul W. Holland and Samuel Leinhardt. Local structure in social networks. *Sociological*
490 *Methodology*, 7:1–45, 1976.
- 491 18 Mark Jerrum and Alistair Sinclair. Approximating the permanent. *SIAM Journal on Comput-*
492 *ing*, 18(6):1149–1178, 1989.
- 493 19 J. Kleinberg. Navigation in a small world. *Nature*, 406(6798), 2000.
- 494 20 James R. Lee, Shayan Oveis Gharan, and Luca Trevisan. Multiway spectral partitioning and
495 higher-order cheeger inequalities. *J. ACM*, 61(6):1–30, 2014.
- 496 21 J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney. Community structure in large
497 networks: Natural cluster sizes and the absence of large well-defined clusters. *Internet*
498 *Mathematics*, 6(1):29–123, 2009.
- 499 22 S. Milgram. The small world problem. *Psychology Today*, 1(1):60–67, 1967.
- 500 23 M. E. J. Newman. Properties of highly clustered networks. *Phys. Rev. E*, 68:026121, Aug
501 2003.
- 502 24 M. E. J. Newman. Finding community structure in networks using the eigenvectors of matrices.
503 *Phys. Rev. E*, 74:036104, Sep 2006.
- 504 25 M. E. J. Newman. Modularity and community structure in networks. *Proceedings of the*
505 *National Academy of Sciences*, 103:8577–8582, 2006.
- 506 26 Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an
507 algorithm. In *Neural Information Processing Systems: Natural and Synthetic*, NIPS’01, page
508 849–856. MIT Press, 2001.
- 509 27 Richard Peng, He Sun, and Luca Zanetti. Partitioning well-clustered graphs: Spectral clustering
510 works! In *Conference on Learning Theory (COLT)*, volume 40 of *Proceedings of Machine*
511 *Learning Research*, pages 1423–1455, 2015.
- 512 28 A. Erdem Sariyuce, C. Seshadhri, A. Pinar, and U. Catalyurek. Finding the hierarchy of dense
513 subgraphs using nucleus decompositions. In *World Wide Web (WWW)*, pages 927–937, 2015.
- 514 29 C. Seshadhri, Tamara G. Kolda, and Ali Pinar. Community structure and scale-free collections
515 of Erdos-Renyi graphs. *Physical Review E*, 85:056109, 2012.
- 516 30 Daniel A. Spielman. *Spectral and Algebraic Graph Theory*. [http://cs-www.cs.yale.edu/
517 homes/spielman/sagt/](http://cs-www.cs.yale.edu/homes/spielman/sagt/).
- 518 31 Daniel A. Spielman. Spectral graph theory and its applications. In *IEEE Symposium on*
519 *Foundations of Computer Science (FOCS)*, pages 29–38, 2007.
- 520 32 Daniel A. Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs
521 and its application to nearly linear time graph partitioning. *SIAM Journal on Computing*,
522 42(1):1–26, 2013.
- 523 33 Charalampos E. Tsourakakis. The k-clique densest subgraph problem. In *The Web Conference*
524 *(WWW)*, pages 1122–1132, 2015.

- 525 **34** Charalampos E. Tsourakakis, Jakub Pachocki, and Michael Mitzenmacher. Scalable motif-
526 aware graph clustering. In *The Web Conference (WWW)*, pages 1451–1460, 2017.
- 527 **35** Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications*.
528 Structural Analysis in the Social Sciences. Cambridge University Press, 1994.
- 529 **36** D. Watts and S. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393:440–442,
530 1998.
- 531 **37** Zeyuan Allen Zhu, Silvio Lattanzi, and Vahab S. Mirrokni. A local algorithm for finding
532 well-connected clusters. In *International Conference on Machine Learning (ICML)*, volume 28,
533 pages 396–404, 2013.