# Partitioning Well-Clustered Graphs
# with $k$-Means and Heat Kernel

Richard Peng[*]        He Sun[†]        Luca Zanetti[‡]

## Abstract

We study a suitable class of *well-clustered graphs* that admit good $k$-way partitions and present the first almost-linear time algorithm for with almost-optimal approximation guarantees partitioning such graphs. A good *$k$-way partition* is a partition of the vertices of a graph into disjoint clusters (subsets) $\{S_i\}_{i=1}^{k}$, such that each cluster is better connected on the inside than towards the outside. This problem is a key building block in algorithm design, and has wide applications in community detection and network analysis.

Key to our result is a theorem on the multi-cut and eigenvector structure of the graph Laplacians of these well-clustered graphs. Based on this theorem, we give the *first* rigorous guarantees on the approximation ratios of the widely used $k$-means clustering algorithms. We also give an almost-linear time algorithm based on heat kernel embeddings and approximate nearest neighbor data structures.

**Keywords:** graph partitioning, spectral clustering, $k$-means, heat kernel, Cheeger inequalities

---

[*]Massachusetts Institute of Technology, Cambridge, USA. Email: rpeng@mit.edu.

[†]Max Planck Institute for Informatics, Saarbrücken, Germany. Email: hsun@mpi-inf.mpg.de. Part of the work was done during a visit to the Simons Institute for the Theory of Computation, UC Berkeley.

[‡]Max Planck Institute for Informatics, Saarbrücken, Germany. Email: luca.zanetti@mpi-inf.mpg.de.

# Contents

# 1 Introduction

Partitioning a graph into two or more pieces is one of the most fundamental problems in combinatorial optimization, and has wide applications in various disciplines of computer science. One of the most studied graph partitioning problems is the *edge expansion problem*, i.e., finding a cut with few crossing edges normalized by the size of the smaller side of the cut. Formally, let $G = (V, E)$ be an undirected and unweighted graph. For any set $S$, the conductance of set $S$ is defined by

$$\phi_G(S) \triangleq \frac{|E(S, V \setminus S)|}{\text{vol}(S)},$$

where $\text{vol}(S)$ is the total weight of edges incident to vertices in $S$, and the conductance of graph $G$ is

$$\phi(G) \triangleq \min_{S:\text{vol}(S) \leqslant \text{vol}(G)/2} \phi_G(S).$$

The edge expansion problem asks for a set $S \subseteq V$ of $\text{vol}(S) \leqslant \text{vol}(V)/2$ that minimizes $\phi(G)$. This problem is known to be NP-hard [MS90], and assuming the Small Set Expansion Conjecture [RST12], does not admit a polynomial-time algorithm that achieves a constant factor approximation in the worst case.

The *k-way partitioning problem* is a natural generalization of the edge expansion problem. We call subsets of vertices (i.e. *clusters*) $A_1, \ldots, A_k$ a *k-way partition* of $G$ if $A_i \cap A_j = \emptyset$ for different $i$ and $j$, and $\bigcup_{i=1}^{k} A_i = V$. The $k$-way partitioning problem asks for a $k$-way partition of $G$ such that the conductance of any $A_i$ in the partition is at most the *k-way expansion constant*, defined by

$$\rho(k) \triangleq \min_{\text{partition } A_1, \ldots, A_k} \max_{1 \leqslant i \leqslant k} \phi_G(A_i). \tag{1.1}$$

Clusters of low conductance in a real network usually capture the notion of *community*, and algorithms for finding these subsets have applications in various domains such as community detection and network analysis. In computer vision, most image segmentation procedures are based on region-based merge and split [CA79], which in turn rely on partitioning graphs into multiple subsets [SM00]. On a theoretical side, decomposing vertex/edge sets into multiple disjoint subsets is a key technique in the approximation algorithms for Unique Games [Tre08], and also has applications in the design of efficient algorithms [KLOS14, LR99, ST11].

Despite widespread use of various graph partitioning schemes over the past decades, the quantitative relationship between the $k$-way expansion constant and the eigenvalues of the graph Laplacians were unknown until a sequence of very recent results, e.g. [KLL$^+$13, LM14, LOGT12, LRTV12, OGT14]. In particular, Lee et al. [LOGT12] proved the following higher-order Cheeger inequality:

$$\frac{\lambda_k}{2} \leqslant \rho(k) \leqslant O(k^2)\sqrt{\lambda_k}, \tag{1.2}$$

where $0 = \lambda_1 \leqslant \ldots \leqslant \lambda_n \leqslant 2$ are the eigenvalues of the normalized Laplacian matrix of $G$. Informally, the higher-order Cheeger inequality shows that a graph $G$ has a $k$-way partition with low $\rho(k)$ if and only if $\lambda_k$ is small. This implies that a large gap between $\lambda_{k+1}$ and $\rho(k)$ *guarantees* (i) existence of a $k$-way partition $\{S_i\}_{i=1}^{k}$ with bounded $\phi_G(S_i) \leqslant \rho(k)$, and (ii) any $(k+1)$-way partition of $G$ contains a subset with significantly higher conductance compared with $\rho(k)$. That is, a suitable lower bound on the *gap* $\Upsilon$ for some $k$, defined by

$$\Upsilon \triangleq \frac{\lambda_{k+1}}{\rho(k)}, \tag{1.3}$$

implies the existence of a $k$-way partition for which every cluster has low conductance, and that $G$ is a *well-clustered* graph.

**Remark.** *Our gap assumption can be also "informally" interpreted as a gap between $\lambda_{k+1}$ and $\lambda_k$, since a large enough gap between $\lambda_{k+1}$ and $\lambda_k$ implies a lower bound on $\Upsilon$.*

## 1.1 Our Results

In this paper we study spectral properties of graphs satisfying the *gap assumption* $\Upsilon = \Omega(k^3)$. We give structural results that show close connections between the eigenvectors and the characteristic vectors of the clusters. This characterization allows us to show that many variants of $k$-means algorithms, that are based on the spectral embedding and that work "in practice", can be rigorously analyzed "in theory". Moreover, exploiting our gap assumption, we can approximate this spectral embedding using the heat kernel of the graph. Combining this with locality-sensitive hashing, we give an almost-linear time algorithm for the $k$-way partitioning problem.

Our structural results can be summarized as follows. Let $\{S_i\}_{i=1}^{k}$ be a $k$-way partition of $G$ achieving $\rho(k)$ defined in (1.1). We define $\bar{g}_1, \cdots, \bar{g}_k$ to be the normalized characteristic vectors of the clusters $\{S_i\}_{i=1}^{k}$, and $\{f_i\}_{i=1}^{k}$ to be the eigenvectors corresponding to the first $k$ smallest eigenvalues of $\mathcal{L}$. Our first result is about the clusters $S_1, \ldots S_k$ and the structure of $f_1, \ldots, f_k$: under the condition of $\Upsilon = \Omega(k^3)$, the span of $\{\bar{g}_i\}_{i=1}^{k}$ and the span of $\{f_i\}_{i=1}^{k}$ are close to *each other*. It can be stated formally as follows:

**Theorem 1.1** (The Structure Theorem). *Let $\{S_i\}_{i=1}^{k}$ be a $k$-way partition of $G$ achieving $\rho(k)$, and let $\Upsilon = \lambda_{k+1}/\rho(k) > k$. Assume that $\{f_i\}_{i=1}^{k}$ are the first $k$ eigenvectors of matrix $\mathcal{L}$, and $\bar{g}_1, \ldots, \bar{g}_k \in \mathbb{R}^n$ are the characteristic vectors of $\{S_i\}_{i=1}^{k}$, with proper normalization[1]. Then the following statements hold:*

1. *For every $\bar{g}_i$, there is a linear combination of $\{f_i\}_{i=1}^{k}$, called $\hat{f}_i$, such that $\|\bar{g}_i - \hat{f}_i\|^2 \leqslant 1/\Upsilon$.*

2. *For every $f_i$, there is a linear combination of $\{\bar{g}_i\}_{i=1}^{k}$, called $\hat{g}_i$, such that $\|f_i - \hat{g}_i\|^2 \leqslant k/\Upsilon$.*

This theorem generalizes the result shown by Arora et al. ([ABS10], Theorem 2.2), which proves the easier direction (the first statement, Theorem 1.1), and can be thought as a stronger version of the well-known Davis-Kahan theorem [DK70]. We remark that, despite that we use the higher-order Cheeger inequality from (1.2) to motivate the definition of $\Upsilon$, our proof of this structure theorem is self-contained. Specifically, it omits much of the machinery used in the proofs of higher-order and improved Cheeger inequalities [KLL+13, LOGT12].

As a direct application, Theorem 1.1 implies that the set of vectors of $\mathbb{R}^k$ in the span of $\{f_i\}_{i=1}^{k}$ is almost equivalent to the set of vectors of $\mathbb{R}^k$ in the span of $\{\bar{g}_i\}_{i=1}^{k}$. This fact has several interesting consequences. For instance, we look at the well-known spectral embedding $F : V \to \mathbb{R}^k$ defined by

$$F(u) \triangleq \frac{1}{\mathsf{NormalizationFactor}(u)} \cdot (f_1(u), \ldots, f_k(u))^{\mathsf{T}}, \tag{1.4}$$

with a proper normalization factor $\mathsf{NormalizationFactor}(u) \in \mathbb{R}$ for each $u \in V$. We use Theorem 1.1 to prove that (i) *all* points $F(u)$ from the same cluster $u \in S_i$ ($1 \leqslant i \leqslant k$) are close to each other, and (ii) *most pairs* of points $F(u), F(v)$ from two different clusters $S_i, S_j$ are far from each other.

Based on this fact, we analyze the performance of spectral $k$-means algorithms[2], aiming at answering the following longstanding open question: *Why do spectral $k$-means algorithms perform*

---

[1]See the formal definition in Section 3.

[2]For simplicity, we use the word "spectral $k$-means algorithms" to refer to the algorithms which combine a spectral embedding with a $k$-means algorithm in Euclidean space.

*well in practice?* We show that the partition $\{A_i\}_{i=1}^k$ produced by the spectral $k$-means algorithm gives a good approximation of any "optimal" partition $\{S_i\}_{i=1}^k$: every $A_i$ has low conductance, and has large overlap with its correspondence $S_i$. To the best of our knowledge, this is the *first* rigorous guarantee for many practical spectral clustering algorithms. These algorithms have comprehensive applications, and have been the subject of extensive experimental studies (e.g., [AY95, NJW$^+$02, VL07]). Our result also gives an affirmative answer to an open question proposed in [LOGT12]: whether the spectral $k$-means algorithm can be rigorously analyzed in certain general circumstances. Our result is as follows:

**Theorem 1.2** (Approximation Guarantee of Spectral $k$-Means Algorithms)**.** *Let $G$ be a graph satisfying the condition $\Upsilon = \lambda_{k+1}/\rho(k) = \Omega(k^3)$, and $k \in \mathbb{N}$. Let $F : V \to \mathbb{R}^k$ be the embedding defined above. Let $\{A_i\}_{i=1}^k$ be a $k$-way partition by any $k$-means algorithm running in $\mathbb{R}^k$ that achieves an $\mathsf{APT}$-approximation. Then the following statements hold: (i) $\phi_G(A_i) = O\left(\phi_G(S_i) + \mathsf{APT} \cdot k^3 \cdot \Upsilon^{-1}\right)$; (ii) $\mathrm{vol}(A_i \triangle S_i) = O\left(\mathsf{APT} \cdot k^3 \cdot \Upsilon^{-1} \cdot \mathrm{vol}(S_i)\right)$.*

This allows us to apply various $k$-means algorithms (e.g., [KSS04, ORSS12]) in Euclidean space to give spectral clustering algorithms, with different time versus approximation tradeoffs.

Notice that for moderately large values of $k$, e.g. $k \approx n^{0.1}$, the performance of these algorithms becomes super-linear, since most $k$-means algorithms have an $\Omega(nk)$ running time. Moreover, when the number of clusters is $k = \omega(\log n)$, it is not even clear how to obtain the embedding (1.4) in $\widetilde{O}(m)$ time. To obtain a faster algorithm, we introduce another novel technique: we approximate the squared-distance $\|F(u) - F(v)\|^2$ of the embedded points $F(u)$ and $F(v)$ via their *heat-kernel distance*, which allows us to avoid the computation of eigenvectors. Using our gap assumption, we apply approximate nearest-neighbor algorithms, and give an ad hoc variant of the $k$-means algorithm that works in almost-linear time.

**Theorem 1.3** (Almost-Linear Time Algorithm For Partitioning Graphs)**.** *Let $G = (V, E)$ be a graph of $n$ vertices and $m$ edges, and a parameter $k \in \mathbb{N}$. Assume that $\Upsilon = \lambda_{k+1}/\rho(k) = \Omega(k^4 \log^3 n)$, and $\{S_i\}_{i=1}^k$ is a $k$-way partition such that $\phi_G(S_i) \leqslant \rho(k)$. Then there is an algorithm running in $\widetilde{O}(m)$ time[3] that outputs a $k$-way partition $\{A_i\}_{i=1}^k$. Moreover, the following statements hold: (i) $\phi_G(A_i) = O(\phi_G(S_i) + k^3 \log^2 k \cdot \Upsilon^{-1})$; (ii) $\mathrm{vol}(A_i \triangle S_i) = O\left(k^3 \log^2 k \cdot \Upsilon^{-1} \cdot \mathrm{vol}(S_i)\right)$.*

Our algorithm differs from most of the previous spectral clustering algorithms in that it works primarily with distances between the embedded vertices, instead of their coordinates along eigenvectors. This approach closely resembles many practical approaches, partly because it circumvents issues related to the stability of eigenvectors. It also allows us to directly use the heat-kernel embedding, which traditionally is used either as an alternative to pagerank vectors [Chu09] or within the matrix multiplicative weights update frameworks [OSV12]. We believe this distance driven approach has wider applications, especially in graph partitioning settings.

## 1.2 Related Work

There is a large amount of literature on partitioning graphs under various settings. Arora et al. [ABS10] gives an $O(1/\lambda_k)$-approximation algorithm for the sparest cut problem with running time $n^{O(k)}$, by searching for a sparest cut in the $k$-dimensional eigenspace corresponding to the first $k$ eigenvectors. Kwok et al. [KLL$^+$13] shows that spectral partitioning gives a constant factor approximation for the sparest cut problem, when $\lambda_k$ is large for constant values of $k$.

---

[3]The $\widetilde{O}(\cdot)$ term hides a factor of poly $\log n$.

Lee et al. [LOGT12] studies the higher-order Cheeger inequalities, and shows that every graph can be partitioned into $k$ non-empty subsets such that every subset in the partition has expansion $O(k^3)\sqrt{\lambda_k}$. Oveis Gharan and Trevisan [OGT14] formulate the notion of clusters with respect to the *inner* and *outer* conductance: a cluster $S$ should have low outer conductance, while the conductance of the induced subgraph by $S$ should be high. Under a gap assumption between $\lambda_{k+1}$ and $\lambda_k$, they further present a polynomial-time algorithm that finds a $k$-partition $\{A_i\}_{i=1}^k$ that satisfy the inner- and outer-conductance condition. In order to assure that every $A_i$ has high inner conductance, they assume that $\lambda_{k+1} \geqslant \mathrm{poly}(k)\lambda_k^{1/4}$, which is much stronger than ours. Moreover, their algorithm runs in polynomial-time, in contrast to our almost-linear time algorithm.

Based on the gap between $\lambda_k$ and $\lambda_{k+1}$, Dey et al. [DRS14] proposed a $k$-way partition algorithm, which is based on the $k$-centers problem and on combinatorial arguments. In contrast to our work, their result only holds for bounded-degree graphs, and cannot provide an approximate guarantee for *every* cluster. Moreover, their algorithm runs in almost-linear time only if $k = O(\mathrm{poly}\log n)$.

We also explore the separation between $\lambda_k$ and $\lambda_{k+1}$ from an algorithmic perspective, and show that this assumption interacts well with heat-kernel embeddings. The heat kernel has been used in previous algorithms on local partitioning [Chu09], balanced separators [OSV12]. It also plays a key role in current efficient approximation algorithms for finding low conductance cuts [OSVV08, She09]. However, most of these theoretical guarantees are through the matrix multiplicative weights update framework [AHK12, AK07]. Our algorithm instead directly uses the heat-kernel embedding to find low conductance cuts.

## 1.3   Organization of the Paper

The paper is organized as follows: We first list background knowledge in Section 2. In Section 3, we analyze the structure theorem. Section 4 studies the $k$-means clustering algorithms, and gives a theoretical approximation guarantee of $k$-means clustering on well-clustered graphs. In Section 5, we present an almost-linear time algorithm for partitioning well-clustered graphs.

## 2   Preliminaries

Let $G = (V, E)$ be an undirected and unweighted graph with $n$ vertices and $m$ edges. The set of neighbors of a vertex $u$ is represented by $N(u)$, and its degree is $d(u) = |N(u)|$. Moreover, for any set $S \subseteq V$, let $\mathrm{vol}(S) \triangleq \sum_{u \in S} d_u$. For any set $S, T \subseteq V$, we define $E(S, T)$ to be the set of edges from $S$ to $T$, aka $E(S, T) \triangleq \{\{u, v\} | u \in S \text{ and } v \in T\}$. For simplicity, we write $\partial S = E(S, V \setminus S)$ for any set $S \subseteq V$. For two sets $X$ and $Y$, the symmetric difference of $X$ and $Y$ is defined as $X \triangle Y \triangleq (X \setminus Y) \cup (Y \setminus X)$.

We will work extensively with algebraic objects related to $G$. The adjacency matrix $\mathbf{A}$ of $G$ is given by

$$\mathbf{A}_{u,v} = \begin{cases} 1 & \text{if } \{u, v\} \in E[G], \text{ and} \\ 0 & \text{otherwise.} \end{cases}$$

We will also use $\mathbf{D}$ to denote the $n \times n$ diagonal matrix with $\mathbf{D}_{uu} = d_u$ for $u \in V[G]$. The *Laplacian matrix* of $G$ is defined by $\mathbf{L} \triangleq \mathbf{D} - \mathbf{A}$, and the *normalized Laplacian matrix* of $G$ is defined by

$$\mathcal{L} \triangleq \mathbf{D}^{-1/2}\mathbf{L}\mathbf{D}^{-1/2} = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}.$$

For this matrix, we will denote its $n$ eigenvalues with $0 = \lambda_1 \leqslant \cdots \leqslant \lambda_n \leqslant 2$, and their corresponding eigenvectors with $f_1, \ldots, f_n$. Note that if $G$ is connected, the first eigenvector is $f_1 = \mathbf{D}^{1/2}f$, where $f$ is any non-zero constant vector.

**Figure 1:** Relations among $\{\hat{f}_i\}$, $\{f_i\}$, $\{\bar{g}_i\}$, and $\{\hat{g}_i\}$. Here $\Upsilon$ is the gap defined with respect to $\lambda_{k+1}$ and $\rho(k)$.

For a vector $x \in \mathbb{R}^n$, the 2-norm, or Euclidean norm of $x$ is given by

$$\|x\| = \left( \sum_{i=1}^{n} x_i^2 \right)^{1/2}.$$

The spectral norm of a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined by

$$\|\mathbf{A}\| = \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|=1}} \|\mathbf{A}x\|.$$

If $\mathbf{A}$ is symmetric, $\|\mathbf{A}\| = |\lambda_{\max}(\mathbf{A})|$, where $\lambda_{\max}(\mathbf{A})$ is the largest eigenvalue of $\mathbf{A}$ in absolute value. If $\mathbf{A}$ is not symmetric, then $\|\mathbf{A}\| = \sqrt{|\lambda_{\max}(\mathbf{A}^\mathsf{T}\mathbf{A})|}$.

For any $f : V \to \mathbb{R}$, the *Rayleigh quotient* of $f$ with respect to graph $G$ is then given by

$$\mathcal{R}(f) \triangleq \frac{f^\mathsf{T}\mathcal{L}f}{\|f\|_2^2} = \frac{f^\mathsf{T}\mathbf{L}f}{\|f\|_\mathbf{D}} = \frac{\sum_{\{u,v\}\in E(G)} (f(u) - f(v))^2}{\sum_u d_u f(u)^2},$$

where $\|f\|_\mathbf{D} \triangleq f^\mathsf{T}\mathbf{D}f$.

Throughout the rest of the paper, we will use $S_1, \ldots, S_k$ to express a $k$-way partition of $G$ achieving the minimum conductance, $\rho(k)$. Note that this partition may not be unique.

# 3 Proof of The Structure Theorem

In this section we give a formal description of Theorem 1.1. Recall that the structure theorem states that (i) any normalized characteristic vector $\bar{g}_i$ of cluster $S_i$ can be approximated by a linear combination of the first $k$ eigenvectors, called $\hat{f}_i$, such that $\|\hat{f}_i - \bar{g}_i\| \leqslant 1/\Upsilon$; (ii) any $f_i$ ($1 \leqslant i \leqslant k$) can be approximated by a linear combination of the normalized characteristic vectors $\{\bar{g}_i\}_{i=1}^k$ such that $\|f_i - \hat{g}_i\| \leqslant k/\Upsilon$, see Figure 1 for an illustration.

To formally prove the structure theorem, let $g_i$ be the characteristic vector of cluster $S_i$, defined by

$$g_i(u) = \begin{cases} 1 & \text{if } u \in S_i \\ 0 & \text{if } u \notin S_i \end{cases} \tag{3.1}$$

for any $1 \leqslant i \leqslant k$, and the corresponding normalized vector is defined by

$$\bar{g}_i = \frac{\mathbf{D}^{1/2}g_i}{\|\mathbf{D}^{1/2}g_i\|}. \tag{3.2}$$

5

Notice that the conductance of a set $S_i$ can be expressed as

$$\phi_G(S_i) = \mathcal{R}(\bar{g}_i), \tag{3.3}$$

and hence we can write the gap $\Upsilon$ as

$$\Upsilon = \frac{\lambda_{k+1}}{\rho(k)} = \min_{1 \leqslant i \leqslant k} \frac{\lambda_{k+1}}{\phi_G(S_i)} = \min_{1 \leqslant i \leqslant k} \frac{\lambda_{k+1}}{\mathcal{R}(\bar{g}_i)}. \tag{3.4}$$

We will always assume that

$$\Upsilon \geqslant C \cdot k^3, \tag{3.5}$$

for a large enough constant $C$.

Theorem 3.1 below shows that the normalized characteristic vector of every cluster $S_i$ can be approximated by a linear combination of *the first $k$ eigenvectors*, with respect to the value of $\Upsilon$. We remark that this result is proven implicitly by Arora et al. ([ABS10], Theorem 2.2).

**Theorem 3.1.** *For any $1 \leqslant i \leqslant k$, there is a linear combination of the eigenvectors $f_1, \ldots, f_k$, called vector $\hat{f}_i \in \mathbb{R}^n$, such that*

$$\left\| \bar{g}_i - \hat{f}_i \right\|^2 \leqslant 1/\Upsilon.$$

*Proof.* We write $\bar{g}_i$ as a linear combination of eigenvectors of $\mathcal{L}$, i.e.,

$$\bar{g}_i = \alpha_1^{(i)} f_1 + \cdots + \alpha_n^{(i)} f_n,$$

and let the vector $\hat{f}_i$ be the projection of vector $\bar{g}_i$ on the subspace spanned by $\{f_i\}_{i=1}^k$, i.e.,

$$\hat{f}_i = \alpha_1^{(i)} f_1 + \cdots + \alpha_k^{(i)} f_k.$$

By the definition of Rayleigh quotients, we have that

$$\begin{aligned}
\mathcal{R}(\bar{g}_i) &= \left( \alpha_1^{(i)} f_1 + \cdots + \alpha_n^{(i)} f_n \right)^\mathsf{T} \mathcal{L} \left( \alpha_1^{(i)} f_1 + \cdots + \alpha_n^{(i)} f_n \right) \\
&= \left( \alpha_1^{(i)} \right)^2 \lambda_1 + \cdots + \left( \alpha_n^{(i)} \right)^2 \lambda_n \\
&\geqslant \left( \alpha_2^{(i)} \right)^2 \lambda_2 + \cdots + \left( \alpha_k^{(i)} \right)^2 \lambda_k + \left( 1 - \alpha' - \left( \alpha_1^{(i)} \right)^2 \right) \lambda_{k+1} \\
&\geqslant \alpha' \lambda_2 + \left( 1 - \alpha' - \left( \alpha_1^{(i)} \right)^2 \right) \lambda_{k+1},
\end{aligned}$$

where $\alpha' = \left( \alpha_2^{(i)} \right)^2 + \cdots + \left( \alpha_k^{(i)} \right)^2$. Therefore, we have that

$$1 - \alpha' - \left( \alpha_1^{(i)} \right)^2 \leqslant \mathcal{R}(\bar{g}_i)/\lambda_{k+1} \leqslant 1/\Upsilon,$$

and

$$\| \bar{g}_i - \hat{f}_i \|^2 = \left( \alpha_{k+1}^{(i)} \right)^2 + \cdots + \left( \alpha_n^{(i)} \right)^2 = 1 - \alpha' - \left( \alpha_1^{(i)} \right)^2 \leqslant 1/\Upsilon,$$

which finishes the proof. ∎

Now we will show that the opposite direction holds as well, i.e., any $f_i$ $(1 \leqslant i \leqslant k)$ can be approximated by a linear combination of the normalized characteristic vectors $\{\bar{g}_i\}_{i=1}^k$.

**Theorem 3.2.** *Let $\Upsilon > k$. For any $1 \leqslant i \leqslant k$, there is a vector $\hat{g}_i = \sum_{j=1}^{k} \beta_j^{(i)} \bar{g}_j$, which is a linear combination of $\{\bar{g}_i\}_{i=1}^{k}$, such that*

$$\|f_i - \hat{g}_i\|^2 \leqslant k/\Upsilon.$$

We first discuss the intuition behind proving Theorem 3.2. It is easy to see that, if we could write every $\bar{g}_i$ *exactly* as a linear combination of $\{f_i\}_{i=1}^{k}$, then we could write every $f_i$ ($1 \leqslant i \leqslant k$) as a linear combination of $\{\bar{g}_i\}_{i=1}^{k}$. This is because both of $\{f_i\}_{i=1}^{k}$ and $\{\bar{g}_i\}_{i=1}^{k}$ are sets of linearly independent vectors of the same dimension and span $\{\bar{g}_1, \ldots, \bar{g}_k\} \subseteq$ span $\{f_1, \ldots, f_k\}$.

However, the $\bar{g}_i$'s are only close to a linear combination of the first $k$ eigenvectors. We will denote this combination as $\hat{f}_i$, and use the fact that the errors of approximation are small to show that these $\{\hat{f}_i\}_{i=1}^{k}$ are almost orthogonal between each other. This allows us to show that span $\left\{\hat{f}_1, \ldots, \hat{f}_k\right\} =$ span $\{f_1, \ldots, f_k\}$, which then implies Theorem 3.2.

Based on the fact that

*Proof.* By Theorem 3.1, every $\bar{g}_i$ is approximated by vector $\hat{f}_i$ defined by

$$\hat{f}_i = \alpha_1^{(i)} f_1 + \cdots \alpha_k^{(i)} f_k.$$

Define a $k$ by $k$ matrix $\mathbf{F}$ such that $\mathbf{F}_{i,j} = \alpha_i^{(j)}$, i.e., the $j$th column of matrix $\mathbf{F}$ consists of values $\left\{\alpha^{(j)}\right\}_{j=1}^{k}$ representing $\hat{f}_j$, where

$$\alpha^{(j)} = \left(\alpha_1^{(j)}, \cdots, \alpha_k^{(j)}\right)^{\mathsf{T}}.$$

Notice that (i) each column of $\mathbf{F}$ has almost unit norm, and (ii) different columns are *almost* orthogonal to each other, in the sense that

$$\left|\left\langle \alpha^{(i)}, \alpha^{(j)} \right\rangle\right| \leqslant \max \left\{\mathcal{R}(\bar{g}_i)/\lambda_{k+1}, \mathcal{R}(\bar{g}_j)/\lambda_{k+1}\right\} \leqslant 1/\Upsilon, \qquad \text{for } i \neq j.$$

This implies that $\mathbf{F}$ is almost an orthogonal matrix. Moreover, since $(\mathbf{F}^{\mathsf{T}}\mathbf{F})_{i,i} \geqslant 1 - 1/\Upsilon$ and $|(\mathbf{F}^{\mathsf{T}}\mathbf{F})_{i,j}| \leqslant 1/\Upsilon$ for $i \neq j$, it holds by the Geršgorin Circle Theorem (cf. Theorem A.1) that all the eigenvalues of $\mathbf{F}^{\mathsf{T}}\mathbf{F}$ are at least

$$1 - 1/\Upsilon - (k-1) \cdot 1/\Upsilon = 1 - k/\Upsilon.$$

Therefore, matrix $\mathbf{F}$ has no eigenvalue with value 0 as long as $\Upsilon > k$, i.e., the vectors $\left\{\alpha^{(j)}\right\}_{j=1}^{k}$ are linearly independent. Combining this with the fact that span $\{\hat{f}_1, \ldots, \hat{f}_k\} \subseteq$ span $\{f_1, \ldots, f_k\}$ and $\dim(\text{span}(\{f_1, \ldots, f_k\})) = k$, it holds that span $\{\hat{f}_1, \ldots, \hat{f}_k\} =$ span $\{f_1, \ldots, f_k\}$. Hence, we can write every $f_i$ ($1 \leqslant i \leqslant k$) as a linear combination of $\{\hat{f}_i\}_{i=1}^{k}$, i.e.,

$$f_i = \beta_1^{(i)} \hat{f}_1 + \beta_2^{(i)} \hat{f}_2 + \cdots + \beta_k^{(i)} \hat{f}_k. \tag{3.6}$$

Now define the value of $\hat{g}_i$ as

$$\hat{g}_i = \beta_1^{(i)} \bar{g}_1 + \beta_2^{(i)} \bar{g}_2 + \cdots + \beta_k^{(i)} \bar{g}_k. \tag{3.7}$$

By Theorem 3.1, it is easy to see that

$$\|f_i - \hat{g}_i\|^2 \leqslant k \max_{1 \leqslant j \leqslant k} \|\hat{f}_j - \bar{g}_j\|^2 \leqslant k/\Upsilon. \qquad \blacksquare$$

7

This implies that the first $k$ eigenvectors, normalized by $\mathbf{D}^{-1/2}$, are close (in the $\mathbf{D}$-norm) to a $k$-step function constant on each cluster. Our next lemma shows that, for every pair of clusters, there exists an eigenvector whose coordinates have reasonably different values on two different clusters. This is due to the fact that the first $k$ eigenvectors are able to approximate the characteristic vector of every cluster.

**Lemma 3.3.** *Let $\Upsilon = \Omega(k^3)$. For any $1 \leqslant i \leqslant k$, let*

$$\hat{g}_i = \beta_1^{(i)} \overline{g}_1 + \cdots + \beta_k^{(i)} \overline{g}_k$$

*be such that*

$$\|f_i - \hat{g}_i\| \leqslant \frac{k}{\Upsilon}.$$

*Then, for any $\ell \neq j$, there exists $i \in \{1, \ldots, k\}$ such that*

$$\left| \beta_\ell^{(i)} - \beta_j^{(i)} \right| \geqslant \zeta \triangleq \frac{1}{10\sqrt{k}}. \tag{3.8}$$

*Proof.* Let $\beta^{(i)} = \left( \beta_1^{(i)}, \ldots, \beta_k^{(i)} \right)^\mathsf{T}$, for $1 \leqslant i \leqslant k$. Since $\overline{g}_i \perp \overline{g}_j$ for any $i \neq j$, we have that $\langle \hat{g}_i, \hat{g}_j \rangle = \langle \beta^{(i)}, \beta^{(j)} \rangle$, and therefore

$$\begin{aligned}
\left| \left\langle \beta^{(i)}, \beta^{(j)} \right\rangle \right| &= |\langle \hat{g}_i, \hat{g}_j \rangle| \leqslant |\langle f_i - (f_i - \hat{g}_i), f_j - (f_j - \hat{g}_j) \rangle| \\
&= |\langle f_i, f_j \rangle - \langle f_i - \hat{g}_i, f_j \rangle - \langle f_j - \hat{g}_j, f_i \rangle + \langle f_i - \hat{g}_i, f_j - \hat{g}_j \rangle| \\
&\leqslant \|f_i - \hat{g}_i\| + \|f_j - \hat{g}_j\| + \|f_i - \hat{g}_i\|\|f_j - \hat{g}_j\| \\
&\leqslant 2\sqrt{k/\Upsilon} + k/\Upsilon,
\end{aligned}$$

where the last inequality follows from Theorem 3.2. This implies that $\beta^{(i)}$'s are almost orthogonal to each other.

Now we construct a $k$ by $k$ matrix $\mathbf{B}$, where the $j$th column of $\mathbf{B}$ is $\beta^{(j)}$. Using the same technique as in Theorem 3.2, we know that, for any eigenvalue $\lambda$ of matrix $\mathbf{B}$ with the corresponding normalized eigenvector $x$, it holds that

$$|\lambda|^2 x^\mathsf{T} x = (\mathbf{B}x)^\mathsf{T} \mathbf{B}x = x^\mathsf{T} \mathbf{B}^\mathsf{T} \mathbf{B}x \in \left( 1 - k(2\sqrt{k/\Upsilon} + k/\Upsilon), 1 + k(2\sqrt{k/\Upsilon} + k/\Upsilon) \right), \tag{3.9}$$

i.e., matrix $\mathbf{B}$ is almost orthogonal and its eigenvalues have modulus close to 1.

We can now show that $\beta_\ell^{(i)}$ and $\beta_j^{(i)}$ are far from each other by contradiction. Suppose there exist $\ell \neq j$ such that

$$\zeta' \triangleq \max_{1 \leqslant i \leqslant k} \left| \beta_\ell^{(i)} - \beta_j^{(i)} \right| < \frac{1}{10\sqrt{k}}.$$

This implies that the $j$th row and $\ell$th row of matrix $\mathbf{B}$ are somewhat close to each other. Let us now define matrix $\mathbf{E} \in \mathbb{R}^{k \times k}$, where

$$\mathbf{E}_{\ell, i} \triangleq \beta_j^{(i)} - \beta_\ell^{(i)},$$

and $\mathbf{E}_{t, i} = 0$ for any $t \neq \ell$ and $1 \leqslant i \leqslant k$. Moreover, let $\mathbf{Q} = \mathbf{B} + \mathbf{E}$. Notice that $\mathbf{Q}$ has two identical rows, and rank at most $k - 1$. Therefore $\mathbf{Q}$ has an eigenvalue with value 0, and the spectral norm $\|\mathbf{E}\|$ of $\mathbf{E}$, the largest singular value of $\mathbf{E}$, is at most $\sqrt{k}\zeta'$. By definition of matrix $\mathbf{Q}$ we have that

$$\mathbf{Q}^\mathsf{T} \mathbf{Q} = \mathbf{B}^\mathsf{T} \mathbf{B} + \mathbf{B}^\mathsf{T} \mathbf{E} + \mathbf{E}^\mathsf{T} \mathbf{B} + \mathbf{E}^\mathsf{T} \mathbf{E}.$$

Since $\mathbf{B}^\intercal\mathbf{B}$ is symmetric and 0 is an eigenvalue of $\mathbf{Q}^\intercal\mathbf{Q}$, by Theorem A.2 we know that, if $\hat{\lambda}$ is an eigenvalue of $\mathbf{Q}^\intercal\mathbf{Q}$, then there is an eigenvalue $\lambda$ of $\mathbf{B}^\intercal\mathbf{B}$ such that

$$
\begin{aligned}
|\hat{\lambda} - \lambda| &\leqslant \|\mathbf{B}^\intercal\mathbf{E} + \mathbf{E}^\intercal\mathbf{B} + \mathbf{E}^\intercal\mathbf{E}\| \\
&\leqslant \|\mathbf{B}^\intercal\mathbf{E}\| + \|\mathbf{E}^\intercal\mathbf{B}\| + \|\mathbf{E}^\intercal\mathbf{E}\| \\
&\leqslant 4\sqrt{k}\zeta' + k\zeta'^2,
\end{aligned}
$$

which implies that

$$
\hat{\lambda} \geqslant \lambda - 4\sqrt{k}\zeta' - k\zeta'^2 \geqslant 1 - k(2\sqrt{k/\Upsilon} + k/\Upsilon) - 4\sqrt{k}\zeta' - k\zeta'^2,
$$

due to (3.9). By setting $\hat{\lambda} = 0$, we have that

$$
1 - k(2\sqrt{k/\Upsilon} + k/\Upsilon) - 4\sqrt{k}\zeta' - k\zeta'^2 \leqslant 0.
$$

By the condition of $\Upsilon$ in (3.5), the inequality above implies that $\zeta' \geqslant \frac{1}{10\sqrt{k}}$, which leads to a contradiction. ∎

**Remark 3.4.** *It was shown in [KLL$^+$13] that the first $k$ eigenvectors can be approximated by a $(2k + 1)$-step function. The quality of the approximation is the same as the one given by our structure theorem. However, a $(2k + 1)$-step function is not enough to show that the entire cluster is concentrated around a certain point.*

# 4  Analysis of Spectral $k$-Means Algorithms

In this section we analyze an algorithm based on the classical spectral clustering paradigm, and give an approximation guarantee of this method on well-clustered graphs. We will show that any the approximation guarantee of any $k$-means algorithm $\mathsf{AlgoMean}(\mathcal{X}, k)$ can be translated to one for the $k$-way partitioning problem. Furthermore, it suffices to call $\mathsf{AlgoMean}$ in a black-box manner with a point set $\mathcal{X} \subseteq \Re^d$.

This section is structured as follows. We first give a quick overview of spectral and $k$-means clustering in Section 4.1. In Section 4.2, we use the structure theorem to analyze the spectral embedding. Section 4.3 gives a general result about the $k$-means algorithm when applied to this embedding, and a formal proof of Theorem 1.2.

## 4.1  $k$-Means Clustering

Given a set of points $\mathcal{X} \subseteq \mathbb{R}^d$, a *k-means algorithm* $\mathsf{AlgoMean}(\mathcal{X}, k)$ seeks to find a set $\mathcal{K}$ of $k$ centers $c_1, \cdots, c_k$ to minimize the sum of the squared-distance between $x \in \mathcal{X}$ and the center to which it is assigned. Formally, for any partition $\mathcal{X}_1, \cdots, \mathcal{X}_k$ of the set $\mathcal{X} \subseteq \mathbb{R}^d$, we define the cost function by

$$
\mathsf{COST}(\mathcal{X}_1, \ldots, \mathcal{X}_k) \triangleq \min_{c_1, \ldots, c_k \in \mathbb{R}^d} \sum_{i=1}^{k} \sum_{x \in \mathcal{X}_i} \|x - c_i\|^2,
$$

i.e., the $\mathsf{COST}$ function minimizes the total squared-distance between the points $x$'s and their individually closest center point $c_i$, where $c_1, \ldots, c_k$ are chosen arbitrarily in $\mathbb{R}^d$. We further define the optimal clustering cost by

$$
\Delta_k^2(\mathcal{X}) \triangleq \min_{\text{partition } \mathcal{X}_1, \ldots, \mathcal{X}_k} \mathsf{COST}(\mathcal{X}_1, \ldots, \mathcal{X}_k). \tag{4.1}
$$

9

A typical spectral $k$-means algorithm on graphs can be described as follows: (i) Compute the first $k$ eigenvectors $f_1, \cdots, f_k$ of the normalized Laplacian matrix[4] of graph $G$. (ii) Map every vertex $u \in V[G]$ to a point $F(u) \in \mathbb{R}^k$ according to

$$F(u) = \frac{1}{\mathsf{NormalizationFactor}(u)} \cdot (f_1(u), \ldots, f_k(u))^{\mathsf{T}}, \tag{4.2}$$

with a proper normalization factor $\mathsf{NormalizationFactor}(u) \in \mathbb{R}$ for each $u \in V$. (iii) Let $\mathcal{X} \triangleq \{F(u) : u \in V\}$ be the set of embedded points from vertices in $G$. Run $\mathsf{AlgoMean}(\mathcal{X}, k)$, and group vertices of $G$ into $k$ clusters, according to the output of $\mathsf{AlgoMean}(\mathcal{X}, k)$. This approach that combines a $k$-means algorithm with a spectral embedding has been widely used in practice for a long time, although there is a lack of rigorous analyses of its performance prior to our results.

## 4.2 Analysis of the Spectral Embedding

The first step of the $k$-means clustering technique described above is to map vertices of a graph into points in Euclidean space, through the spectral embedding (1.4). This subsection analyzes the properties of this embedding. Let us define the normalization factor to be

$$\mathsf{NormalizationFactor}(u) \triangleq \sqrt{d_u}.$$

We will show that the embedding (4.2) with the normalization factor above has very nice properties: embedded points from different clusters of $G$ are far from each other, while embedded points from the same cluster $S_i$ are concentrated around their center $c_i \in \mathbb{R}^k$. These properties imply that a simple $k$-means algorithm is able to produce a good clustering[5].

We first define $k$ points $p^{(i)} \in \mathbb{R}^k$ ($1 \leqslant i \leqslant k$), where

$$p^{(i)} \triangleq \frac{1}{\sqrt{\mathrm{vol}\,(S_i)}} \left( \beta_i^{(1)}, \ldots, \beta_i^{(k)} \right)^{\mathsf{T}}, \tag{4.3}$$

i.e., $p^{(i)}$ can be expressed as

$$p^{(i)} = \left( \mathbf{D}^{-1/2} \hat{g}_1(u), \ldots, \mathbf{D}^{-1/2} \hat{g}_k(u) \right),$$

where $u$ is any vertex in $S_i$. We will show in Lemma 4.1 that all embedded points $\mathcal{X}_i \triangleq \{F(u) : u \in S_i\}$ ($1 \leqslant i \leqslant k$) are concentrated around $p^{(i)}$. Moreover, we bound the total squared-distance between vertices in $\mathcal{X}_i$ and $p^{(i)}$, which is proportional to $1/\Upsilon$: the bigger the value of $\Upsilon$, the higher concentration the points within the same cluster have. Notice that we *do not* claim that $p^{(i)}$ is the actual center of $\mathcal{X}_i$. However, these approximated points $p^{(i)}$'s suffice for our analysis.

**Lemma 4.1.** *It holds that*

$$\sum_{i=1}^{k} \sum_{u \in S_i} d_u \left\| F(u) - p^{(i)} \right\|^2 \leqslant k^2/\Upsilon. \tag{4.4}$$

---

[4]Other graph matrices (e.g. the adjacency matrix, and the Laplacian matrix) are also widely used in practice. Notice that, with proper normalization, the choice of these matrices does not substantially influence the performance of $k$-means algorithms.

[5]Notice that this embedding is similar with the one used in [LOGT12], with the only difference that $F(u)$ is not normalized and so it is not necessarily a unit vector. This difference, though, is crucial for our analysis.

*Proof.* Since $\|x\|^2 = \|\mathbf{D}^{-1/2}x\|_{\mathbf{D}}$ holds for any $x \in \mathbb{R}^n$, by Theorem 3.2 we have for any $1 \leqslant j \leqslant k$ that

$$\sum_{i=1}^{k} \sum_{u \in S_i} d_u \left( F(u)_j - p_j^{(i)} \right)^2 = \left\| \mathbf{D}^{-1/2} f_j - \mathbf{D}^{-1/2} \hat{g}_j \right\|_{\mathbf{D}} \leqslant k/\Upsilon.$$

Summing over all $j$ for $1 \leqslant j \leqslant k$ implies that

$$\sum_{i=1}^{k} \sum_{u \in S_i} d_u \left\| F(u) - p^{(i)} \right\|^2 = \sum_{i=1}^{k} \sum_{j=1}^{k} \sum_{u \in S_i} d_u \left( F(u)_j - p_j^{(i)} \right)^2 \leqslant k^2/\Upsilon. \qquad \blacksquare$$

**Lemma 4.2.** *It holds for every $1 \leqslant i \leqslant k$ that*

$$\frac{9}{10\,\mathrm{vol}(S_i)} \leqslant \left\| p^{(i)} \right\|^2 \leqslant \frac{11}{10\,\mathrm{vol}(S_i)}.$$

*Proof.* By (4.3), we have that

$$\left\| p^{(i)} \right\|^2 = \frac{1}{\mathrm{vol}(S_i)} \left\| \left( \beta_i^{(1)}, \dots, \beta_i^{(k)} \right)^{\mathsf{T}} \right\|^2.$$

Notice that $p^{(i)}$ is just the $i$th row of matrix $\mathbf{B}$ defined in Lemma 3.3, normalized by $\sqrt{\mathrm{vol}(S_i)}$. Taking the transpose of $\mathbf{B}$ and $x = \mathbf{1}$, we apply the same argument as in (3.9) and obtain that

$$\left\| \left( \beta_i^{(1)}, \dots, \beta_i^{(k)} \right)^{\mathsf{T}} \right\|^2 \in [9/10, 11/10], \tag{4.5}$$

which implies the statement. $\blacksquare$

Lemma 4.2 shows that $\left\| p^{(i)} \right\|^2$ is proportional to $1/\mathrm{vol}(S_i)$. We will further show in Lemma 4.3 that these points $p^{(i)}(1 \leqslant i \leqslant k)$ exhibit another excellent property: the distance between $p^{(i)}$ and $p^{(j)}$ is inversely proportional to the volume of the *smaller* cluster between $S_i$ and $S_j$. Therefore, embedded points in $\mathcal{X}_i$ from $S_i$ of smaller $\mathrm{vol}(S_i)$ are far from embedded points in $\mathcal{X}_j$ of bigger $\mathrm{vol}(S_j)$. Notice that, if this was not the case, a small misclassification of points in a bigger cluster $S_j$ could introduce a large error in the cluster of smaller volume.

**Lemma 4.3.** *For every $i \neq j$, it holds that*

$$\left\| p^{(i)} - p^{(j)} \right\|^2 \geqslant \frac{\zeta^2}{10 \min \left\{ \mathrm{vol}(S_i), \mathrm{vol}(S_j) \right\}},$$

*where $\zeta$ is defined in (3.8).*

*Proof.* By Lemma 3.3, there exists $1 \leqslant \ell \leqslant k$ such that

$$\left| \beta_i^{(\ell)} - \beta_j^{(\ell)} \right| \geqslant \zeta.$$

By the definition of $p^{(i)}$ and $p^{(j)}$ it follows that

$$\left\| \frac{p^{(i)}}{\|p^{(i)}\|} - \frac{p^{(j)}}{\|p^{(j)}\|} \right\|^2 \geqslant \left( \frac{\beta_i^{(\ell)}}{\sqrt{\sum_{t=1}^{k} \left( \beta_i^{(t)} \right)^2}} - \frac{\beta_j^{(\ell)}}{\sqrt{\sum_{t=1}^{k} \left( \beta_j^{(t)} \right)^2}} \right)^2.$$

11

By Lemma 4.2, we know that

$$\sum_{\ell=1}^{k} \left(\beta_j^{(\ell)}\right)^2 = \left\|\left(\beta_j^{(1)}, \ldots, \beta_j^{(k)}\right)^{\mathsf{T}}\right\|^2 \in [9/10, 11/10].$$

Therefore, we have that

$$\left\|\frac{p^{(i)}}{\|p^{(i)}\|} - \frac{p^{(j)}}{\|p^{(j)}\|}\right\|^2 \geqslant \frac{1}{2} \cdot \left(\beta_i^{(\ell)} - \beta_j^{(\ell)}\right)^2 \geqslant \frac{1}{2} \cdot \zeta^2,$$

and

$$\left\langle \frac{p^{(i)}}{\|p^{(i)}\|}, \frac{p^{(j)}}{\|p^{(j)}\|} \right\rangle \leqslant 1 - \zeta^2/4.$$

Without loss of generality, we assume that $\left\|p^{(i)}\right\|^2 \geqslant \left\|p^{(j)}\right\|^2$. By Lemma 4.2, it holds that

$$\left\|p^{(i)}\right\|^2 \geqslant \frac{9}{10 \cdot \operatorname{vol}(S_i)},$$

and

$$\left\|p^{(i)}\right\|^2 \geqslant \left\|p^{(j)}\right\|^2 \geqslant \frac{9}{10 \cdot \operatorname{vol}(S_j)}.$$

Hence, it holds that

$$\left\|p^{(i)}\right\|^2 \geqslant \frac{9}{10 \min\left\{\operatorname{vol}(S_i), \operatorname{vol}(S_j)\right\}}.$$

We can now finish the proof by considering two cases based on $\left\|p^{(i)}\right\|$.

*Case 1:* Suppose that $\left\|p^{(i)}\right\|^2 \geqslant 4\left\|p^{(j)}\right\|^2$. We have that

$$\left\|p^{(i)} - p^{(j)}\right\| \geqslant \left\|p^{(i)}\right\| - \left\|p^{(j)}\right\| \geqslant \frac{1}{2}\left\|p^{(i)}\right\|,$$

which implies that

$$\left\|p^{(i)} - p^{(j)}\right\|^2 \geqslant \frac{1}{4}\left\|p^{(i)}\right\|^2 \geqslant \frac{1}{5 \min\left\{\operatorname{vol}(S_i), \operatorname{vol}(S_j)\right\}}.$$

*Case 2:* Suppose $\left\|p^{(j)}\right\| = \alpha\left\|p^{(i)}\right\|$ for $\alpha \in (\frac{1}{4}, 1]$. In this case, we have that

$$\begin{aligned}
\left\|p^{(i)} - p^{(j)}\right\|^2 &= \left\|p^{(i)}\right\|^2 + \left\|p^{(j)}\right\|^2 - 2\left\langle \frac{p^{(i)}}{\|p^{(i)}\|}, \frac{p^{(j)}}{\|p^{(j)}\|} \right\rangle \left\|p^{(i)}\right\| \left\|p^{(j)}\right\| \\
&\geqslant \left\|p^{(i)}\right\|^2 + \left\|p^{(j)}\right\|^2 - 2(1 - \zeta^2/4) \cdot \left\|p^{(i)}\right\| \left\|p^{(j)}\right\| \\
&= (1 + \alpha^2)\left\|p^{(i)}\right\|^2 - 2(1 - \zeta^2/4)\alpha \cdot \left\|p^{(i)}\right\|^2 \\
&= (1 + \alpha^2 - 2\alpha + \alpha\zeta^2/2)\left\|p^{(i)}\right\|^2 \\
&\geqslant \frac{\alpha\zeta^2}{2} \cdot \left\|p^{(i)}\right\|^2 \geqslant \zeta^2 \cdot \frac{1}{10 \min\left\{\operatorname{vol}(S_i), \operatorname{vol}(S_j)\right\}},
\end{aligned}$$

and the lemma follows. $\blacksquare$

## 4.3 Approximation Analysis of Spectral $k$-Means Algorithms

We now give an explanation of why spectral $k$-means algorithms perform well for solving the $k$-way partitioning problem. Throughout the whole subsection, we assume that $A_1, \ldots, A_k$ is any $k$-way partition of $G$ that is returned by a $k$-means algorithm with an approximation ratio of $\mathsf{APT}$.

We first map every vertex $u$ to $d_u$ identical points in $\mathbb{R}^k$. This "trick" allows us to bound the volume of the overlap between the clusters retrieved by a $k$-means algorithm and the optimal ones. For this reason, the cost function of partition $A_1, \ldots, A_k$ of $V[G]$ is defined by

$$\mathsf{COST}(A_1, \ldots, A_k) \triangleq \min_{c_1, \ldots, c_k \in \mathbb{R}^k} \sum_{i=1}^{k} \sum_{u \in A_i} d_u \|F(u) - c_i\|^2,$$

and the optimal clustering cost is defined by

$$\Delta_k^2 \triangleq \min_{\text{partition } A_1, \ldots, A_k} \mathsf{COST}(A_1, \ldots, A_k),$$

i.e., we define the optimal clustering cost in the same way as in (4.1), except that we look at the embedded points from vertices of $G$ in the definition. From now on, we always refer to $\mathsf{COST}$ and $\Delta_k^2$ as the $\mathsf{COST}$ and optimal $\mathsf{COST}$ values of points $\{F(u)\}_{u \in V}$, where for technical reasons every point is counted $d_u$ times.

**Lemma 4.4.** *The optimal solution of a $k$-means clustering satisfies $\Delta_k^2 \leqslant k^2/\Upsilon$.*

*Proof.* Since $\Delta_k^2$ is obtained by minimizing over all partitions $A_1, \ldots, A_k$ and $c_1, \ldots, c_k$, we have that

$$\Delta_k^2 \leqslant \sum_{i=1}^{k} \sum_{u \in S_i} d_u \left\| F(u) - p^{(i)} \right\|^2. \tag{4.6}$$

Hence the statement follows by applying Lemma 4.1. ∎

By Lemma 4.4 and the assumption that $A_1, \cdots, A_k$ is an $\mathsf{APT}$-approximation of an optimal clustering, we have that $\mathsf{COST}(A_1, \ldots, A_k) \leqslant \mathsf{APT} \cdot k^2/\Upsilon$. In the following, we show that this upper bound of $\mathsf{APT} \cdot k^2/\Upsilon$ suffices to show that this approximate clustering $A_1, \ldots, A_k$ is close to the "actual" clustering $S_1, \ldots, S_k$, in the sense that, (i) every $A_i$ has low conductance, and (ii) under a proper permutation $\sigma : \{1, \ldots, k\} \to \{1, \ldots, k\}$, the symmetric difference between $A_i$ and $S_{\sigma(i)}$ is low.

**Lemma 4.5.** *Let $A_1, \ldots, A_k$ be a partition of $V$. Suppose that, for every permutation of the indices $\sigma : \{1, \ldots, k\} \to \{1, \ldots, k\}$, there exists $i$ such that $\mathrm{vol}\left(A_i \triangle S_{\sigma(i)}\right) \geqslant 2\varepsilon \, \mathrm{vol}\left(S_{\sigma(i)}\right)$ for $\varepsilon \geqslant \frac{1000k^2}{\zeta^2 \Upsilon}$. Then, it holds that*

$$\mathsf{COST}(A_1, \ldots, A_k) \geqslant \min\left\{\frac{\varepsilon \zeta^2}{100}, \frac{\zeta^2}{100k}\right\}.$$

We will give a complete proof of Lemma 4.5 in the next subsection. Now we are ready to prove Theorem 1.2.

**Lemma 4.6.** *Let $A_1, \ldots, A_k$ be a $k$-way partition that achieves an approximation ratio of $\mathsf{APT}$. Then, there exists a permutation $\sigma$ of the indices such that*

$$\mathrm{vol}\left(A_i \triangle S_{\sigma(i)}\right) \leqslant \frac{2000k^2 \cdot \mathsf{APT}}{\zeta^2 \Upsilon} \mathrm{vol}(S_{\sigma(i)})$$

*for any $1 \leqslant i \leqslant k$.*

*Proof.* The proof is by contradiction. Assume that there is $i \in \{1, \ldots, k\}$ such that

$$\mathrm{vol}(A_i \triangle S_{\sigma(i)}) > \frac{2000 k^2 \cdot \mathsf{APT}}{\zeta^2 \Upsilon} \, \mathrm{vol}(S_{\sigma(i)}).$$

This implies by Lemma 4.5 that

$$\mathsf{COST}(A_1, \ldots, A_k) > 10 \cdot \mathsf{APT} \cdot k^2 / \Upsilon,$$

which contradicts to the fact that $A_1, \ldots, A_k$ is an $\mathsf{APT}$-approximation to a $k$-way partition, whose optimal cost is at most $\mathsf{APT} \cdot k^2 / \Upsilon$. ∎

**Lemma 4.7.** *Let $A_1, \ldots, A_k$ be a $k$-way partition that achieves an approximation ratio of $\mathsf{APT}$, and $\sigma : \{1, \cdots, k\} \to \{1, \cdots, k\}$ be the permutation defined in Lemma 4.6. Let*

$$\varepsilon = \frac{2000 k^2 \cdot \mathsf{APT}}{\zeta^2 \Upsilon} = O\left(\frac{k^3 \cdot \mathsf{APT}}{\Upsilon}\right).$$

*Then, it holds for every $1 \leqslant i \leqslant k$ that*

$$\phi_G(A_i) = O(\phi_G(S_{\sigma(i)}) + \mathsf{APT} \cdot k^3 / \Upsilon).$$

*Proof.* For any $1 \leqslant i \leqslant k$, the number of leaving edges of $A_i$ is upper bounded by

$$
\begin{aligned}
|\partial(A_i)| &\leqslant |\partial(A_i \setminus S_{\sigma(i)})| + |\partial(A_i \cap S_{\sigma(i)})| \\
&\leqslant |\partial(A_i \triangle S_{\sigma(i)})| + |\partial(A_i \cap S_{\sigma(i)})| \\
&\leqslant \varepsilon \, \mathrm{vol}(S_{\sigma(i)}) + \phi_G(S_{\sigma(i)}) \, \mathrm{vol}(S_{\sigma(i)}) \\
&= (\varepsilon + \phi_G(S_{\sigma(i)})) \, \mathrm{vol}(S_{\sigma(i)}),
\end{aligned}
$$

where the third inequality follows from Lemma 4.6 and the fact that we use the same $\sigma$ as in Lemma 4.6. On the other hand, we have that

$$\mathrm{vol}(A_i) \geqslant \mathrm{vol}(A_i \cap S_{\sigma(i)}) \geqslant (1 - \varepsilon) \, \mathrm{vol}(S_{\sigma(i)}).$$

Hence,

$$\phi_G(A_i) \leqslant \frac{(\varepsilon + \phi_G(S_{\sigma(i)})) \, \mathrm{vol}(S_{\sigma(i)})}{(1 - \varepsilon) \, \mathrm{vol}(S_{\sigma(i)})} = \frac{\varepsilon + \phi_G(S_{\sigma(i)})}{1 - \varepsilon} = O(\phi_G(S_{\sigma(i)}) + \mathsf{APT} \cdot k^3 / \Upsilon). \qquad \blacksquare$$

Theorem 1.2 follows by combining Lemma 4.6 and Lemma 4.7.

## 4.4 Proof of Lemma 4.5

The proof of Lemma 4.5 is based on the following high-level idea: suppose by contradiction that there is a cluster $S_j$ which is very different from every cluster $A_\ell$, where $A_1, \ldots, A_k$ is an $\mathsf{APT}$-approximate $k$-way partition. Then there is a cluster $A_i$ with significant overlaps with two different clusters $S_j$ and $S_{j'}$. However, Lemma 4.3 gives that any two clusters are far from each other. This implies that the $\mathsf{COST}$ value of $A_1, \ldots, A_k$ is high, giving a contradiction.

**Lemma 4.8.** *Suppose for every permutation $\pi : \{1, \ldots, k\} \to \{1, \ldots, k\}$ there exists index $i$ such that*

$$\mathrm{vol}(A_i \triangle S_{\pi(i)}) \geqslant 2\varepsilon \, \mathrm{vol}(S_{\pi(i)}).$$

*Then one of the following statements holds:*

- *For any index $i$ there are indices $i_1 \neq i_2$ and $\varepsilon_i \geqslant 0$ such that*

$$\operatorname{vol}(A_i \cap S_{i_1}) \geqslant \operatorname{vol}(A_i \cap S_{i_2}) \geqslant \varepsilon_i \min\left\{\operatorname{vol}(S_{i_1}), \operatorname{vol}(S_{i_2})\right\},$$

  *and $\sum_{i=1}^{k} \varepsilon_i \geqslant \varepsilon$.*

- *There are indices $i', j, \ell$ such that*

$$\operatorname{vol}(A_{i'} \cap S_j) \geqslant \operatorname{vol}(A_{i'} \cap S_\ell) \geqslant \operatorname{vol}(S_\ell)/k.$$

*Proof.* Let $\sigma : \{1, \ldots, k\} \to \{1, \ldots, k\}$ be the function defined by

$$\sigma(i) = \operatorname*{argmax}_{1 \leqslant j \leqslant k} \frac{\operatorname{vol}(A_i \cap S_j)}{\operatorname{vol}(S_j)}.$$

We first assume that $\sigma$ is one-to-one, i.e. $\sigma$ is a permutation. By the hypothesis of the lemma, there exists an index $i$ such that $\operatorname{vol}(A_i \triangle S_{\sigma(i)}) \geqslant 2\varepsilon \operatorname{vol}(S_{\sigma(i)})$. Without loss of generality, we assume that $i = 1$. Notice that

$$\operatorname{vol}\left(A_1 \triangle S_{\sigma(1)}\right) = \sum_{j \neq 1} \operatorname{vol}\left(A_j \cap S_{\sigma(1)}\right) + \sum_{j \neq \sigma(1)} \operatorname{vol}\left(A_1 \cap S_j\right). \tag{4.7}$$

Hence, one of the summations on the right hand side of (4.7) is at least $\varepsilon \operatorname{vol}(S_{\sigma(1)})$. Now the proof is based on the case distinction.

*Case 1:* Assume that $\sum_{j \neq 1} \operatorname{vol}\left(A_j \cap S_{\sigma(1)}\right) \geqslant \varepsilon \operatorname{vol}(S_{\sigma(1)})$. We define $\tau_j$ for $1 \leqslant j \leqslant k, j \neq 1$, to be

$$\tau_j = \frac{\operatorname{vol}\left(A_j \cap S_{\sigma(1)}\right)}{\operatorname{vol}\left(S_{\sigma(1)}\right)}.$$

We have that

$$\sum_{j \neq 1} \tau_j \geqslant \varepsilon,$$

and by the definition of $\sigma$ we have that

$$\operatorname{vol}\left(A_j \cap S_{\sigma(j)}\right) \geqslant \tau_j \cdot \operatorname{vol}\left(S_{\sigma(j)}\right)$$

for any $1 \leqslant j \leqslant k$.

*Case 2:* Assume that

$$\sum_{j \neq \sigma(1)} \operatorname{vol}\left(A_1 \cap S_j\right) \geqslant \varepsilon \operatorname{vol}(S_{\sigma(1)}). \tag{4.8}$$

Let us define $\tau'_j$ for $1 \leqslant j \leqslant k, j \neq \sigma(1)$, to be

$$\tau'_j = \frac{\operatorname{vol}(A_1 \cap S_j)}{\operatorname{vol}\left(S_{\sigma(1)}\right)}.$$

By (4.8) we have that

$$\sum_{j \neq \sigma(1)} \tau'_j \geqslant \varepsilon.$$

This case holds by assuming $\operatorname{vol}\left(A_1 \cap S_{\sigma(1)}\right) \geqslant \varepsilon \operatorname{vol}\left(S_{\sigma(1)}\right)$, since otherwise we have

$$\sum_{j \neq 1} \operatorname{vol}\left(A_j \cap S_{\sigma(1)}\right) \geqslant \varepsilon' \operatorname{vol}(S_{\sigma(1)})$$

15

**Figure 2:** We use the fact that $\|p^{(i_1)} - c_i\| \geqslant \|p^{(i_2)} - c_i\|$, and lower bound the value of COST function by only looking at the contribution of points $u \in B_i$ for all $1 \leqslant i \leqslant k$.

for $\varepsilon' = 1 - \varepsilon$, and this case was proven in *Case 1*.

Let us now consider the case that $\sigma$ as defined earlier is *not* one-to-one. Hence, there is $j$ $(1 \leqslant j \leqslant k)$ such that $j \notin \{\sigma(1), \ldots, \sigma(k)\}$. Since $\{A_1, \ldots, A_k\}$ is a partition, there exists $i'$ such that $\mathrm{vol}(A_{i'} \cap S_j) \geqslant \mathrm{vol}(S_j)/k$. However, by the definition of $\sigma$, we have that $\mathrm{vol}\left(A_{i'} \cap S_{\sigma(i)}\right) \geqslant \mathrm{vol}\left(S_{\sigma(i')}\right)/k$ for $\sigma(i) \neq j$, which completes the proof. ∎

*Proof of Lemma 4.5.* By Lemma 4.8 for every $i$ there exist $i_1 \neq i_2$ such that

$$\begin{aligned}
\mathrm{vol}(A_i \cap S_{i_1}) &\geqslant \varepsilon_i \min\left\{\mathrm{vol}(S_{i_1}), \mathrm{vol}(S_{i_2})\right\}, \\
\mathrm{vol}(A_i \cap S_{i_2}) &\geqslant \varepsilon_i \min\left\{\mathrm{vol}(S_{i_1}), \mathrm{vol}(S_{i_2})\right\},
\end{aligned} \tag{4.9}$$

for some $\varepsilon \geqslant 0$, and

$$\sum_{i=1}^{k} \varepsilon_i \geqslant \min\{\varepsilon, 1/k\}.$$

Let $c_i$ be the center of $A_i$. Let us assume without loss of generality that $\|c_i - p^{(i_1)}\| \geqslant \|c_i - p^{(i_2)}\|$, which implies $\|p^{(i_1)} - c_i\| \geqslant \|p^{(i_1)} - p^{(i_2)}\|/2$. However, points in $B_i = A_i \cap S_{i_1}$ are far away from $c_i$, see Figure 2. We lower bound the value of $\mathsf{COST}(A_1, \ldots, A_k)$ by only looking at the contribution of points in the $B_i$s . Notice that by Lemma 4.1 the sum of the squared-distances between points in $B_i$ and $p^{(i_1)}$ is at most $k^2/\Upsilon$, while the distance between $p^{(i_1)}$ and $p^{(i_2)}$ is large (Lemma 4.3). Therefore, we have that

$$\begin{aligned}
\mathsf{COST}(A_1, \ldots, A_k) &= \sum_{i=1}^{k} \sum_{u \in A_i} d_u \|F(u) - c_i\|^2 \\
&\geqslant \sum_{i=1}^{k} \sum_{u \in B_i} d_u \|F(u) - c_i\|^2
\end{aligned}$$

By applying the inequality $a^2 + b^2 \geqslant (a - b)^2/2$, we have that

16

$$\mathsf{COST}(A_1, \ldots, A_k) \geqslant \sum_{i=1}^{k} \sum_{u \in B_i} d_u \left( \frac{\left\| p^{(i_1)} - c_i \right\|^2}{2} - \left\| F(u) - p^{(i_1)} \right\|^2 \right)$$

$$\geqslant \sum_{i=1}^{k} \sum_{u \in B_i} d_u \frac{\left\| p^{(i_1)} - c_i \right\|^2}{2} - \sum_{i=1}^{k} \sum_{u \in B_i} d_u \left\| F(u) - p^{(i_1)} \right\|^2$$

$$\geqslant \sum_{i=1}^{k} \sum_{u \in B_i} d_u \frac{\left\| p^{(i_1)} - c_i \right\|^2}{2} - \frac{k^2}{\Upsilon} \tag{4.10}$$

$$\geqslant \sum_{i=1}^{k} \sum_{u \in B_i} d_u \frac{\left\| p^{(i_1)} - p^{(i_2)} \right\|^2}{8} - \frac{k^2}{\Upsilon}$$

$$\geqslant \sum_{i=1}^{k} \frac{\zeta^2 \operatorname{vol}(B_i)}{80 \min \left\{ \operatorname{vol}(S_{i_1}), \operatorname{vol}(S_{i_2}) \right\}} - \frac{k^2}{\Upsilon} \tag{4.11}$$

$$\geqslant \sum_{i=1}^{k} \frac{\zeta^2 \varepsilon_i \min \left\{ \operatorname{vol}(S_{i_1}), \operatorname{vol}(S_{i_2}) \right\}}{80 \min \left\{ \operatorname{vol}(S_{i_1}), \operatorname{vol}(S_{i_2}) \right\}} - \frac{k^2}{\Upsilon}$$

$$\geqslant \sum_{i=1}^{k} \frac{\zeta^2 \varepsilon_i}{80} - \frac{k^2}{\Upsilon}$$

$$\geqslant \min \left\{ \frac{\zeta^2 \varepsilon}{80}, \frac{\zeta^2}{80k} \right\} - \frac{k^2}{\Upsilon} \geqslant \min \left\{ \frac{\zeta^2 \varepsilon}{100}, \frac{\zeta^2}{100k} \right\}$$

where (4.10) follows from Lemma 4.1, (4.11) follows from Lemma 4.3 and the last inequality follows from the assumption that $\varepsilon \geqslant \frac{1000k^2}{\zeta \Upsilon}$. ∎

# 5 Partitioning Well-Clustered Graphs in Almost-Linear Time

In this section we present the first almost-linear time algorithm for partitioning well-clustered graphs. Our algorithm is motivated by the heat kernel embedding, which allows us to approximate distances between $F(u)$ in nearly-linear time. We introduce these objects in Section 5.1, present an overview of our algorithm in Section 5.2, and give its analysis in Section 5.3.

## 5.1 Heat Kernel Embedding

The heat kernel is the fundamental solution of the heat equation

$$\frac{\partial u}{\partial t} = -\mathcal{L}u.$$

Through the heat kernel, the Laplacian is associated with the rate of dissipation of heat. In the discrete case, we can define the heat kernel of a graph. Formally, for any graph $G$ with the normalized Laplacian matrix $\mathcal{L}$, the heat kernel of $G$ is defined by

$$\mathbf{H}_t \triangleq \mathrm{e}^{-t\mathcal{L}}, \tag{5.1}$$

for a *temperature* $t \geqslant 0$. By the definition of the matrix exponential, we can rewrite (5.1) as

$$\mathbf{H}_t = \sum_{i=1}^{n} \mathrm{e}^{-t\lambda_i} f_i f_i^{\mathsf{T}}, \tag{5.2}$$

where $\lambda_1, \ldots, \lambda_n$ are the eigenvalues of matrix $\mathcal{L}$, with the corresponding eigenvectors $f_1, \ldots, f_n$. It is known that the heat kernel on a graph defined in (5.1) or similar forms relates to a geometric embedding, and continuous random walks [LPW09]. We refer the reader to [Chu97] for further details on the heat kernels.

In this work we view the heat kernel as a geometric embedding from vertices of $G$ to points in $\mathbb{R}^n$. Formally, we define the heat kernel embedding $\psi_t : V \to \mathbb{R}^n$ for any fixed $t \geqslant 0$ by

$$\psi_t(u) \triangleq \frac{1}{\sqrt{d_u}} \cdot \left( \mathrm{e}^{-(t/2)\cdot\lambda_1} f_1(u), \cdots, \mathrm{e}^{-(t/2)\cdot\lambda_n} f_n(u) \right). \tag{5.3}$$

This means the squared-distance between the embedded points $\psi(u)$ and $\psi(v)$ can be written as

$$\eta_t(u, v) \triangleq \|\psi_t(u) - \psi_t(v)\|^2. \tag{5.4}$$

Notice that, in contrast to the spectral embedding (1.4) that maps vertices of $G$ to points in $\mathbb{R}^k$, the heat kernel embedding maps vertices to points in $\mathbb{R}^n$. We will resolve this issue with We first show that under the condition of $k = \Omega(\log n)$ and the gap assumption of $\Upsilon$, there is a wide range of $t$ for which $\eta_t(u, v)$ gives a good approximation of $\|F(u) - F(v)\|^2$.

**Lemma 5.1.** *Let $t \in (c \log n/\lambda_{k+1}, 1/\lambda_k)$, for a constant $c > 1$. Then, it holds for every $u, v \in V$ that*

$$\frac{1}{\mathrm{e}} \cdot \|F(u) - F(v)\|^2 \leqslant \eta_t(u, v) \leqslant \|F(u) - F(v)\|^2 + \frac{1}{n^{c-1}}.$$

*Proof.* By the definition of the heat kernel distance in (5.4), we have that

$$\begin{aligned}
\eta_t(u, v) &= \sum_{i=1}^{n} \mathrm{e}^{-t\lambda_i} \left( \frac{f_i(u)}{\sqrt{d_u}} - \frac{f_i(v)}{\sqrt{d_v}} \right)^2 \\
&= \sum_{i=1}^{k} \mathrm{e}^{-t\lambda_i} \left( \frac{f_i(u)}{\sqrt{d_u}} - \frac{f_i(v)}{\sqrt{d_v}} \right)^2 + \sum_{i=k+1}^{n} \mathrm{e}^{-t\lambda_i} \left( \frac{f_i(u)}{\sqrt{d_u}} - \frac{f_i(v)}{\sqrt{d_v}} \right)^2.
\end{aligned} \tag{5.5}$$

Notice that it holds for $1 \leqslant i \leqslant k$ that

$$1 \geqslant \mathrm{e}^{-t\lambda_i} \geqslant \mathrm{e}^{-\lambda_i/\lambda_k} \geqslant \frac{1}{\mathrm{e}}, \tag{5.6}$$

while it holds for $k + 1 \leqslant i \leqslant n$ that

$$\mathrm{e}^{-t\lambda_i} \leqslant \mathrm{e}^{-c \log n \lambda_i/\lambda_{k+1}} \leqslant \mathrm{e}^{-c \log n \lambda_{k+1}/\lambda_{k+1}} = \frac{1}{n^c}. \tag{5.7}$$

By (5.6), the first summation in (5.5) is $[1/\mathrm{e}, 1] \cdot \|F(u) - F(v)\|^2$, and by (5.7) the second summation in (5.5) is at most $n^{-c+1}$. Hence, the statement holds. ∎

The proof above shows why heat kernel embedding can be used to approximate the spectral embedding used in the spectral $k$-means algorithms: Under the condition of $k = \Omega(\log n)$ and $\Upsilon = \Omega(k^3)$, there is $t \in (c \log n/\lambda_{k+1}, 1/\lambda_k)$, such that, when viewing $\|\psi_t(u) - \psi_t(v)\|^2$, the contribution

to $\|\psi_t(u) - \psi_t(v)\|^2$ from the first $k$ coordinates of $\psi_t(u)$ and $\psi_t(v)$ gives a $(1/e)$-approximation of $\|F(u) - F(v)\|^2$, while the contribution to $\|\psi_t(u) - \psi_t(v)\|^2$ from the remaining $n - k$ coordinates of $\psi_t(u)$ and $\psi_t(v)$ is $O(n^{-c})$, for a constant $c$. We remark that a similar intuition which views the heat kernel embedding as a weighted combination of multiple eigenvectors was discussed in [OSV12]. The main reason to use the heat kernel embedding instead of the spectral embedding given by the first $k$ eigenvectors is that there is an almost-linear time algorithm approximating $e^{-\mathbf{A}}x$ for any SDD matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and any vector $x \in \mathbb{R}^n$.

**Theorem 5.2** ([OSV12])**.** *Given an $n \times n$ SDD matrix $\mathbf{A}$ with $m_\mathbf{A}$ nonzero entries, a vector $v$ and a parameter $\delta > 0$, there is an algorithm that can compute a vector $x$ such that $\|e^{-\mathbf{A}}y - x\| \leqslant \delta\|y\|$ in time $\widetilde{O}((m_\mathbf{A} + n)\log(2 + \|\mathbf{A}\|))$[6]. Moreover, this algorithm corresponds to a linear operator realized by a matrix $\mathbf{Z}$ such that for any vector $x$, its output is $\mathbf{Z}x$.*

The following lemma shows that we can obtain an embedding in almost-linear time, and this embedding can be used to approximate the heat kernel distance between vertices.

**Lemma 5.3.** *Let $G$ be a graph with $n$ vertices and $m$ edges. Let $k = \Omega(\log n)$, and $\Upsilon = \Omega(k^3)$. Then, for any parameters $t, \varepsilon > 0$, we can compute an embedding of the vertices, $x_t(u) \in \mathbb{R}^{O(\varepsilon^{-2} \cdot \log n)}$, in $\widetilde{O}(\varepsilon^{-2} \cdot (m + n)\log(2 + t))$ time[7] such that with high probability it holds for all vertices $u$ and $v$ that*

$$(1 - \varepsilon)\eta_t(u, v) \leqslant \|x_t(u) - x_t(v)\|^2 \leqslant \eta_t(u, v) + n^{-c},$$

*for some $c > 1$. In other words, the $\ell_2^2$-distance given by $x_t$ is a good approximation to the heat kernel distance.*

*Proof.* Since $\mathbf{H}_t = \mathbf{H}_{t/2}\mathbf{H}_{t/2}$, we have that

$$\eta_t(u, v) = \left\|\mathbf{H}_{t/2}\left(\xi_u - \xi_v\right)\right\|^2.$$

Replacing $\mathbf{H}_{t/2}$ with an operator $\mathbf{Z}$ of error $\delta$, we get

$$\left|\|\mathbf{Z}\left(\xi_u - \xi_v\right)\| - \eta_t^{1/2}(u, v)\right| \leqslant \delta\|\xi_u - \xi_v\| \leqslant \delta,$$

where the last inequality follows from $d_u, d_v \geqslant 1$. This is equivalent to

$$\eta_t^{1/2}(u, v) - \delta \leqslant \|\mathbf{Z}\left(\xi_u - \xi_v\right)\| \leqslant \eta_t^{1/2}(u, v) + \delta. \tag{5.8}$$

We invoke the Johnson-Lindenstrauss transform in a way analogous to the computation of effective resistances from [SS11] and [KLP12]. For an $O(\varepsilon^{-2} \cdot \log n) \times n$ Gaussian matrix $\mathbf{Q}$, with high probability it holds for all $u, v$ that

$$(1 - \varepsilon)\|\mathbf{Z}\left(\xi_u - \xi_v\right)\| \leqslant \|\mathbf{QZ}\left(\xi_u - \xi_v\right)\| \leqslant (1 + \varepsilon)\|\mathbf{Z}\left(\xi_u - \xi_v\right)\|. \tag{5.9}$$

Combining (5.8) and (5.9) gives us that

$$(1 - \varepsilon)\left(\eta_t^{1/2}(u, v) - \delta\right) \leqslant \|\mathbf{QZ}\left(\xi_u - \xi_v\right)\| \leqslant (1 + \varepsilon)\left(\eta_t^{1/2}(u, v) + \delta\right).$$

---

[6]The $\widetilde{O}$ notation here hides $\mathrm{poly}(\log n)$ and $\mathrm{poly}(\log(1/\delta))$ factors.

[7]The $\widetilde{O}$ notation here hides factors of $\log(n/\varepsilon)$.

Square both sides, and invoking the inequality

$$(1 - \varepsilon)\alpha^2 - (1 + \varepsilon^{-1})b^2 \leqslant (a + b)^2 \leqslant (1 + \varepsilon)\alpha^2 + (1 + \varepsilon^{-1})b^2,$$

then gives

$$(1 - 5\varepsilon)\,\eta_t(u, v) - 2\delta^2\varepsilon^{-1} \leqslant \|\mathbf{QZ}\,(\xi_u - \xi_w)\|^2 \leqslant (1 + 5\varepsilon)\,\eta_t(u, v) + 2\delta^2\varepsilon^{-1}.$$

Scaling $\mathbf{QZ}$ by a factor of $(1 + 5\varepsilon)^{-1}$, and appending an extra entry in each vector to create an additive distortion of $2\delta\varepsilon^{-1}$ then gives the desired bounds when $\delta$ is set to $\varepsilon n^{-c}$. The running time then follows from $\|\mathcal{L}\| \leqslant 2$ and the performance of the approximate exponential algorithm from [OSV12] described in Theorem 5.2. ∎

Combing Lemma 5.1 with Lemma 5.3, we obtain the following result:

**Lemma 5.4.** *Let $G$ be a graph with $n$ vertices and $m$ edges. Let $k = \Omega(\log n)$, and $\Upsilon = \Omega(k^3)$. Then, there is an embedding of vertices $x_t(u) \in \mathbb{R}^{O(\varepsilon^{-2} \cdot \log n)}$, which is computable in $\widetilde{O}(\varepsilon^{-2} \cdot (m+n))$ time, such that with high probability it holds for all $u$, $v$ that*

$$(1 - \varepsilon) \cdot \frac{1}{\mathrm{e}} \cdot \|F(u) - F(v)\|^2 \leqslant \|x_t(u) - x_t(v)\|^2 \leqslant \|F(u) - F(v)\|^2 + \frac{2}{n^{c-1}}.$$

## 5.2 Algorithm Overview

Conceptually, our algorithm follows the general framework of $k$-means algorithms, which consists of two key steps: a seeding step and a grouping step. The seeding step chooses $k$ candidate centers such that, with good probability, each one is close to the actual center of a different cluster. The grouping step assigns each of the remaining vertices to the candidate center closest to it.

We emphasize that choosing good candidate centers is crucial for most $k$-means algorithms, and has been studied extensively in literature (e.g. [AV07, ORSS12]). Recent results show that good initial centers can be obtained by iteratively picking vertices from a *non-uniform* distribution, leading to algorithms running in $\Omega(nk)$ time. The additional structure of our embedding allows for a simpler sampling scheme motivated by these routines. Since $\|F(u)\|^2$ is approximately equal to $1/\operatorname{vol}(S_i)$ for most vertices $u \in S_i$, we can simply sample vertices with probabilities proportional to $d_u \cdot \|F(u)\|^2$ to ensure that we sample from the clusters uniformly. This allows us to show that $|C| = \Theta(k \log k)$ samples ensure that, for every cluster, we pick a vertex close to its center. The well-separation property of the embedded points also allows us to remove the vertices in $C$ which are close to each other. This removal process ensures that at the end of the seeding step there is exactly one vertex left from every cluster, forming a set $C^\star$.

After obtaining $C^\star$, we can proceed with the grouping step. Thanks once again to the well-separation property of our points, we just need to assign every vertex to its nearest sampled center in $C^\star$, which is much simpler than most $k$-means algorithms, e.g. [ORSS12]. Naively this takes $\Omega(nk)$ time. We speed this up further by showing that for most vertices, the correct center is an $\varepsilon$-approximate nearest neighbor even for moderate values of $\varepsilon$. This allows us to obtain an almost-linear time routine based on approximate nearest neighbor data structures [IM98].

When $k = O(\operatorname{poly} \log n)$, this framework directly gives an almost-linear time algorithm when combined with algorithms for computing the first $k$ eigenvectors. However, this becomes more expensive as $k$ becomes large. Note however that our algorithm only needs distance information between the points $\{F(u)\}_{u \in V[G]}$: in fact, we can check that any constant factor approximation of these distances suffices. This means we can use embeddings in lower dimensional spaces that

approximate the distances given by $F$. Furthermore, we can compute such an embedding directly using the heat kernel embedding given by the matrix $e^{-t\mathcal{L}}$, which can be approximated in nearly-linear time [OSV12]. In the case of larger $k$, the gap assumption allows us to show that there is $t \in (c \log n / \lambda_{k+1}, 1/\lambda_k)$ for which heat kernel distances approximate distances in $F$ well. Moreover, if we consider all $t$ of the form $t = 2^i, i = O(\log n)$, we will have considered a $t$ in this range due to the gap assumption. The minimum cost partition returned at these values of $t$ will then give a good clustering. Our overall algorithm framework for $k = \omega(\log n)$ is described in Figure 3.

---

$\text{CLUSTER}(G, k)$

1. For $1 \leqslant i \leqslant k$ do $A_i' := \emptyset$

2. $\text{COST}(A_1', \ldots, A_k') := \infty$

3. For $t = 2, 4, 8, \ldots, \text{poly}(n)$ do

   (a) $N \leftarrow \Theta(k \log k)$

   (b) $(c_1, \ldots, c_k) \leftarrow \text{SEEDANDTRIM}(G, N, k, t)$

   (c) Compute a partition $A_1, \ldots, A_k$ of $V$: for every $v \in V$ assign $v$ to its nearest center $c_i$ using the algorithm of the $\varepsilon$-NNS problem with $\varepsilon = \log k$.

   (d) If $\text{COST}(A_1, \ldots, A_k) \leqslant \text{COST}(A_1', \ldots, A_k')$ $\text{SET} A_i' := A_i$ FOR $1 \leqslant i \leqslant k$

4. $\text{RETURN}(A_1', \cdots, A_k')$

---

**Figure 3:** Clustering Algorithm

**Remark 5.5.** *Notice that both the algorithm and analysis in the case of $k = \Omega(\log n)$ are more involved, as additional approximation to the spectral embedding is needed. Hence, in the rest of this section, we only focus on the case of $k = \Omega(\log n)$. However, we still use the embedding $F(u)$, and due to Lemma 5.4, we can get a constant factor approximation guarantee when we use $x_t(u)$ instead of $F(u)$.*

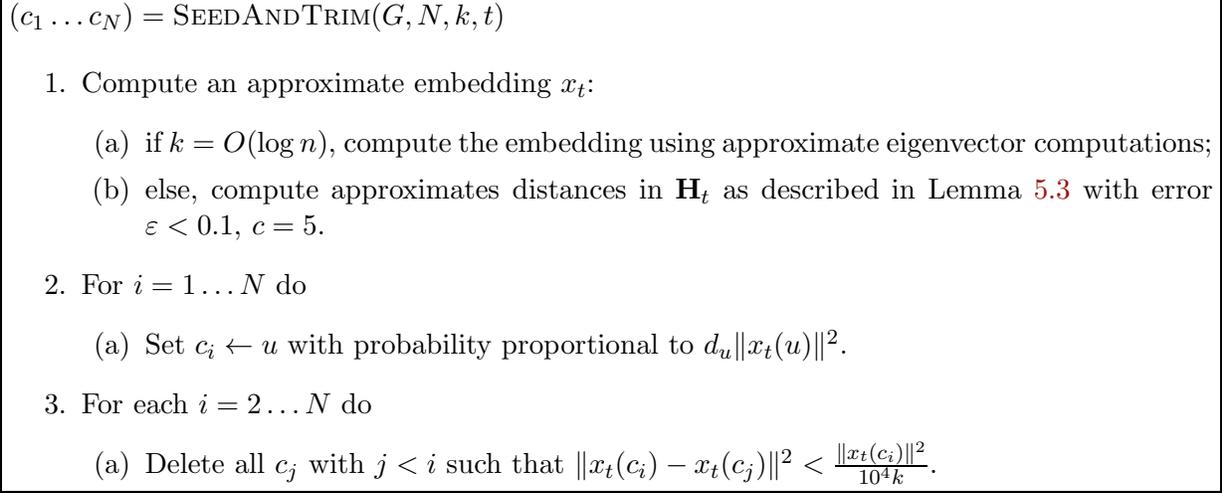## 5.3 Analysis of the Algorithm

In this subsection, we analyze the seeding step, and the group step, as well as the approximate guarantees of the clusters returned by our algorithm. Throughout this section we will assume that $\Upsilon = \Omega(k^4 \log^3 n)$.

**Analysis of the Seeding Step.** In the seeding step, we sample $N \triangleq \Theta(k \log k)$ vertices, each with probability proportional to $d_u \|F(u)\|^2$. After that we delete the sampled vertices that are close to each other until there are exactly $k$ vertices left. A formal description of this routine is in Figure 4.

Now we analyze the seeding step. For any $1 \leqslant i \leqslant k$, we define $\mathcal{E}_i$ to be

$$\mathcal{E}_i \triangleq \sum_{u \in S_i} d_u \left\| F(u) - p^{(i)} \right\|^2,$$

$(c_1 \ldots c_N) = \textsc{SeedAndTrim}(G, N, k, t)$

1. Compute an approximate embedding $x_t$:

    (a) if $k = O(\log n)$, compute the embedding using approximate eigenvector computations;

    (b) else, compute approximates distances in $\mathbf{H}_t$ as described in Lemma 5.3 with error $\varepsilon < 0.1$, $c = 5$.

2. For $i = 1 \ldots N$ do

    (a) Set $c_i \leftarrow u$ with probability proportional to $d_u \|x_t(u)\|^2$.

3. For each $i = 2 \ldots N$ do

    (a) Delete all $c_j$ with $j < i$ such that $\|x_t(c_i) - x_t(c_j)\|^2 < \frac{\|x_t(c_i)\|^2}{10^4 k}$.

**Figure 4:** Seeding Algorithm. For simplicity in Step 2 and 3, we only write the case of $k = \Omega(\log n)$. When $k = O(\log n)$, we can simply use $\|F(u)\|$ to replace $x_t(u)$.

and define the radius of $S_i$ to be

$$R_i^\alpha \triangleq \frac{\alpha \cdot \mathcal{E}_i}{\mathrm{vol}(S_i)}$$

for some parameter $\alpha$, i.e., $R_i^\alpha$ is the approximate mean square error in cluster $S_i$. We define $\mathsf{CORE}_i^\alpha \subseteq S_i$ to be the set of vertices whose $\ell_2^2$-distance to $p^{(i)}$ is at most $R_i^\alpha$, i.e.,

$$\mathsf{CORE}_i^\alpha \triangleq \left\{ u \in S_i : \left\| F(u) - p^{(i)} \right\|^2 \leqslant R_i^\alpha \right\}.$$

By the averaging argument it holds that

$$\mathrm{vol}(S_i \setminus \mathsf{CORE}_i^\alpha) \leqslant \frac{\sum_{u \in S_i} d_u \left\| F(u) - p^{(i)} \right\|^2}{R_i^\alpha} = \frac{\mathrm{vol}(S_i)}{\alpha},$$

and therefore $\mathrm{vol}(\mathsf{CORE}_i^\alpha) \geqslant \left(1 - \frac{1}{\alpha}\right) \mathrm{vol}(S_i)$. From now on, we assume that $\alpha = \Theta(N \log N)$.

**Lemma 5.6.** *For each cluster $S_i$, it holds that*

$$\sum_{u \in \mathsf{CORE}_i^\alpha} d_u \cdot \|F(u)\|^2 \geqslant \frac{9}{10} \left(1 - \frac{1}{100N}\right),$$

*and also the sum over the vertices not in the cores satisfies*

$$\sum_{i=1}^{k} \sum_{u \notin \mathsf{CORE}_i^\alpha} d_u \cdot \|F(u)\|^2 \leqslant \frac{k}{100N}.$$

22

*Proof.* By the definition of $\mathsf{CORE}_i^\alpha$, we have that

$$\sum_{i=1}^{k} \sum_{u \in \mathsf{CORE}_i^\alpha} d_u \cdot \|F(u)\|^2 \geqslant \frac{1}{\alpha} \sum_{i=1}^{k} \int_0^\alpha \sum_{u \in \mathsf{CORE}_i^\rho} d_u \cdot \|F(u)\|^2 d\rho$$

$$\geqslant \frac{1}{\alpha} \sum_{i=1}^{k} \int_0^\alpha \left( \left\| p^{(i)} \right\| - \sqrt{R_i^\rho} \right)^2 \mathrm{vol}(\mathsf{CORE}_i^\rho) d\rho \tag{5.10}$$

$$\geqslant \frac{1}{\alpha} \sum_{i=1}^{k} \int_0^\alpha \left( \left\| p^{(i)} \right\|^2 - 2\sqrt{R_i^\rho} \cdot \left\| p^{(i)} \right\| \right) \left( 1 - \frac{1}{\rho} \right) \mathrm{vol}(S_i) d\rho \tag{5.11}$$

$$\geqslant \frac{1}{\alpha} \int_0^\alpha \left( k - 2 \sum_{i=1}^{k} \sqrt{\frac{11 \cdot \mathcal{E}_i \rho}{10}} \right) \left( 1 - \frac{1}{\rho} \right) d\rho \tag{5.12}$$

where (5.10) follows from the fact that for all $u \in \mathsf{CORE}_i^\rho$, $\|F(u)\| \geqslant \|p^{(i)}\| - \sqrt{R_i^\rho}$, (5.11) from $\mathrm{vol}(\mathsf{CORE}_i^\rho) \geqslant \left( 1 - \frac{1}{\rho} \right) \mathrm{vol}(S_i)$, and (5.12) from the definition of $R_i^\rho$ and the fact that $\sum_{1=1}^{k} \|p^{(i)}\|^2 \cdot \mathrm{vol}(S_i) = k$. By the Cauchy-Schwarz inequality, we have that

$$\sum_{i=1}^{k} \sqrt{\frac{11 \mathcal{E}_i \rho}{10}} \leqslant \sqrt{k \cdot \sum_{i=1}^{k} \frac{11}{10} \mathcal{E}_i \rho} \leqslant \sqrt{\frac{11 \cdot k^3 \rho}{10 \cdot \Upsilon}},$$

and combing this with (5.12) gives us that

$$\sum_{i=1}^{k} \sum_{u \in \mathsf{CORE}_i^\alpha} d_u \cdot \|F(u)\|^2 \geqslant \frac{1}{\alpha} \int_0^\alpha \left( k - 2\sqrt{\frac{11 \cdot k^3 \rho}{10 \Upsilon}} \right) \left( 1 - \frac{1}{\rho} \right) d\rho \tag{5.13}$$

$$\geqslant \frac{1}{\alpha} \int_0^\alpha \left( k - 2\sqrt{\frac{11 \cdot k^3 \rho}{10 \Upsilon}} - \frac{k}{\rho} \right) d\rho$$

$$\geqslant k \left( 1 - \sqrt{\frac{k \cdot \alpha}{\Upsilon}} - \frac{\ln \alpha}{\alpha} \right)$$

$$\geqslant k \left( 1 - \frac{1}{100N} \right),$$

where the last inequality holds by assuming $\alpha = \Theta(N \log N)$ and $\Upsilon \geqslant 100c^2 k N^3 \log N$ for a sufficiently large constant $c$. Combing this with the fact

$$\sum_{u \in V[G]} d_u \|F(u)\|^2 = \sum_{u \in V[G]} \sum_{i=1}^{k} f_i^2(u) = k$$

yields the second statement of the lemma.

Using a similar argument we can show that

$$\sum_{u \in \mathsf{CORE}_i^\alpha} d_u \cdot \|F(u)\|^2 \geqslant \frac{9}{10} \left( 1 - \frac{1}{cN} \right),$$

which finishes the proof. ∎

The next lemma shows that, after sampling $\Theta(k \log k)$ vertices, with constant probability all the sampled vertices are from the cores of $k$ clusters, and every core contains at least one sampled vertex.

**Lemma 5.7.** *Assume that $N = \Omega(k \log k)$ vertices are sampled, in which every vertex is sampled with probability proportional to $d_u \cdot \|F(u)\|^2$. Then, with constant probability the set $Z = \{c_1 \ldots c_N\}$ of sampled vertices has the following properties:*

*1. Set $Z$ only contains vertices from the cores, i.e. $Z \subseteq \bigcup_{i=1}^{k} \mathsf{CORE}_i^{\alpha}$, and*

*2. Set $Z$ contains at least one vertex from each cluster, i.e.*

$$Z \cap S_i \neq \emptyset \qquad \forall 1 \leqslant i \leqslant k.$$

*Proof.* By Lemma 5.4, it holds for every vertex $u$ that

$$\frac{1}{2\mathrm{e}} \cdot \|F(u)\|^2 \leqslant \|x_t(u)\|^2 \leqslant \|F(u)\|^2 + \frac{1}{n^5}.$$

Since

$$\sum_{u \in V[G]} d_u \|F(u)\|^2 = \sum_{u \in V[G]} \sum_{i=1}^{k} f_i^2(u) = k,$$

the total probability mass that we use to sample vertices, i.e. $\sum_{u \in V[G]} d_u \|x_t(u)\|^2$, is between $\frac{1}{2\mathrm{e}} \cdot k$ and $k + 1$. We first bound the probability that we sample at least one vertex from every core. For every $1 \leqslant j \leqslant k$, we have that the probability of each sample coming from $\mathsf{CORE}_j^{\alpha}$ is at least

$$\frac{\sum_{u \in \mathsf{CORE}_i^{\alpha}} d_u \cdot \|x_t(u)\|^2}{k+1} \geqslant \frac{\sum_{u \in \mathsf{CORE}_i^{\alpha}} d_u \cdot \|F(u)\|^2}{2\mathrm{e} \cdot (k+1)} \geqslant \frac{\frac{9}{10}\left(1 - \frac{1}{100N}\right)}{2\mathrm{e} \cdot (k+1)} \geqslant \frac{1}{10k}.$$

Therefore, the probability that we never encounter a vertex from sampling $N$ vertices is at most

$$\left(1 - \frac{1}{10k}\right)^N \leqslant \frac{1}{10k}.$$

Also, the probability that a sampled vertex is outside the cores of the clusters is at most

$$\frac{\sum_{u \notin \mathsf{CORE}_i^{\alpha}, \forall i} d_u \cdot \|x_t(u)\|^2}{k/2} \leqslant \frac{\sum_{u \notin \mathsf{CORE}_i^{\alpha}, \forall i} d_u \cdot \left(\|F(u)\|^2 + n^{-5}\right)}{k/2}$$

$$\leqslant \frac{\frac{k}{100N} + n^{-4}}{k/2} \leqslant \frac{1}{n^2} + \frac{2}{100N}.$$

Taking a union bound over all these events gives that the total probability of undesired events is bounded by

$$k \cdot \frac{1}{10k} + N \cdot \left(\frac{1}{n^2} + \frac{2}{100N}\right) \leqslant \frac{1}{2}. \qquad \blacksquare$$

We now show that points from the same core are much closer between each other than points from different cores. In other words, we show that the procedure SEEDANDTRIM succeeds with constant probability.

**Lemma 5.8.** *For any two vertices $u, v \in \mathsf{CORE}_i^\alpha$, it holds that*

$$\|x_t(u) - x_t(v)\|^2 \leqslant \frac{12\alpha k^2}{\Upsilon \operatorname{vol}(S_i)} < \frac{1}{2 \cdot 10^4 k}.$$

*Proof.* By the definition of $\mathsf{CORE}_i^\alpha$, it holds for any $u \in \mathsf{CORE}_i^\alpha$ that

$$\left\| F(u) - p^{(i)} \right\| \leqslant \sqrt{R_i^\alpha}$$

By the triangle inequality, it holds for any $u \in \mathsf{CORE}_i^\alpha$ and $v \in \mathsf{CORE}_i^\alpha$ that

$$\|F(u) - F(v)\| \leqslant 2\sqrt{R_i^\alpha},$$

or

$$\|F(u) - F(v)\|^2 \leqslant 4R_i^\alpha = \frac{4\alpha \mathcal{E}_i}{\operatorname{vol}(S_i)} \leqslant \frac{4\alpha k^2}{\Upsilon \operatorname{vol}(S_i)},$$

where the last inequality follows from the fact that $\sum_{i=1}^k \mathcal{E}_i \leqslant k^2/\Upsilon$. On the other hand, we also have

$$\|F(u)\|^2 \geqslant \left( \left\| p^{(i)} \right\| - \sqrt{R_i} \right)^2 \geqslant \frac{9}{10} \cdot \left( 1 - \frac{1}{cN} \right) \cdot \left\| p^{(i)} \right\|^2 \geqslant \frac{4}{5 \operatorname{vol}(S_i)}.$$

and

$$\|F(u)\|^2 \leqslant \left( \left\| p^{(i)} \right\| + \sqrt{R_i} \right)^2 \leqslant \frac{11}{10} \left( 1 + \frac{1}{cN} \right) \cdot \left\| p^{(i)} \right\|^2 \leqslant \frac{6}{5 \operatorname{vol}(S_i)}.$$

Therefore we can incorporate the conditions on $x_t(u)$ to give

$$\|x_t(u) - x_t(v)\|^2 \leqslant \|F(u) - F(v)\|^2 + \frac{1}{n^3}$$

$$\leqslant \frac{4\alpha k^2}{\Upsilon \operatorname{vol}(S_i)} + \frac{1}{n^3}$$

$$\leqslant \frac{10\alpha k^2}{\Upsilon} \|F(u)\|^2$$

$$\leqslant \frac{12\alpha k^2}{\Upsilon \operatorname{vol}(S_i)}.$$

By the conditions on $\alpha$ and $\Upsilon$, and the fact that $\|x_t(u)\|^2 \leqslant 2\|F(u)\|^2$, it also holds

$$\|x_t(u) - x_t(v)\|^2 \leqslant \frac{10\alpha k^2}{\Upsilon} \|F(u)\|^2 < \frac{\|x_t(u)\|^2}{2 \cdot 10^4 k}.$$

∎

**Lemma 5.9.** *For any $u \in \mathsf{CORE}_i^\alpha$ and $v \in \mathsf{CORE}_j^\alpha$ where $i \neq j$, we have*

$$\|x_t(u) - x_t(v)\|^2 \geqslant \frac{1}{1000k \operatorname{vol}(S_i)} > \frac{\|x_t(u)\|^2}{10^4 k}.$$

*Proof.* By the triangle inequality, it holds for any pair of $u \in \mathsf{CORE}_i^\alpha$ and $v \in \mathsf{CORE}_j^\alpha$ that

$$\|F(u) - F(v)\| \geqslant \left\|p^{(i)} - p^{(j)}\right\| - \left\|F(u) - p^{(i)}\right\| - \left\|F(v) - p^{(j)}\right\|.$$

By Lemma 4.3, we have for any $i \neq j$,

$$\left\|p^{(i)} - p^{(j)}\right\|^2 \geqslant \frac{1}{10k \min\left\{\mathrm{vol}(S_i), \mathrm{vol}(S_j)\right\}}.$$

Combing this with the fact that

$$\left\|F(u) - p^{(i)}\right\| \leqslant \sqrt{R_i^\alpha} \leqslant \sqrt{\frac{\alpha \cdot k^2}{\Upsilon \, \mathrm{vol}(S_i)}},$$

we obtain that

$$\|F(u) - F(v)\| \geqslant \left\|p^{(i)} - p^{(j)}\right\| - \left\|F(u) - p^{(i)}\right\| - \left\|F(v) - p^{(j)}\right\|$$

$$\geqslant \sqrt{\frac{1}{10k \min\left\{\mathrm{vol}(S_i), \mathrm{vol}(S_j)\right\}}} - \sqrt{\frac{\alpha k^2}{\Upsilon \, \mathrm{vol}(S_i)}} - \sqrt{\frac{\alpha k^2}{\Upsilon \, \mathrm{vol}(S_j)}}$$

$$\geqslant \sqrt{\frac{1}{100k \min\left\{\mathrm{vol}(S_i), \mathrm{vol}(S_j)\right\}}}$$

Hence, we have that

$$\|x_t(u) - x_t(v)\|^2 \geqslant \frac{1}{2\mathrm{e}} \|F(u) - F(v)\|^2 \geqslant \frac{1}{1000k \, \mathrm{vol}(S_i)} > \frac{\|x_t(u)\|^2}{10^4 k}. \qquad \blacksquare$$

Combing Lemma 5.8 and Lemma 5.9 directly gives us the following result:

**Lemma 5.10.** *The procedure* SEEDANDTRIM *returns in* $\widetilde{O}(m + k^2)$ *time a set of centers* $c_1 \dots c_k$ *such that each* $\mathsf{CORE}_i^\alpha$ *contains exactly one* $c_i$.

**Analysis of the Grouping Step.** After the seeding step, we obtain $k$ vertices $c_1, \cdots, c_k$. The analysis about the seeding step assures that these $k$ vertices belong to $k$ different clusters. The next step is to assign the remaining $n - k$ vertices to different clusters. Based on the well-separation property, we can simply ask every vertex to choose its nearest point in the embedded space. Hence we reduce this step to the following $\varepsilon$-approximate nearest neighbor problem ($\varepsilon$-NNS).

**Problem 1** ($\varepsilon$-approximate nearest neighbor Problem). *Given a set of point* $P \in \mathbb{R}^d$ *and a point* $q \in \mathbb{R}^d$, *find a point* $p \in P$ *such that, for all* $p' \in P$, $\|p - q\| \leqslant (1 + \varepsilon)\|p' - q\|$.

The grouping step of our algorithm uses the algorithm in [IM98] for the $\varepsilon$-NNS problem.

**Theorem 5.11** (Proposition 1 of [IM98]). *Given a set of points* $P \subset \mathbb{R}^d$, *and* $\varepsilon > 0$, *there is an algorithm for* $\varepsilon$-NNS *which uses* $\widetilde{O}\left(|P|^{1 + \frac{1}{1+\varepsilon}} + d|P|\right)$ *preprocessing and requires* $\widetilde{O}\left(d|P|^{\frac{1}{1+\varepsilon}}\right)$ *query time.*

By applying Theorem 5.11 and setting $\varepsilon = \Theta(\log k)$, the grouping step takes $\widetilde{O}(nd)$ time in total.

**Approximation Analysis of the Algorithm.** Now we analyze the approximation ratios of the returned $k$-way partition.

**Lemma 5.12.** *Let $(A_1, \ldots, A_k) = \text{CLUSTER}(G, k)$ be the partition of $V[G]$ computed by the algorithm of Figure 3. Then, under a proper permutation of the indices, for any $1 \leqslant i \leqslant k$ it holds*

$$\text{vol}(A_i \triangle S_i) = O\left(\frac{k^3 \log^2 k}{\Upsilon} \text{vol}(S_i)\right)$$

*and*

$$\phi_G(A_i) = O\left(\phi_G(S_i) + \frac{k^3 \log^2 k}{\Upsilon}\right).$$

*Proof.* First we bound the symmetric difference between $A_i$ and its correspondence $S_i$ $(1 \leqslant i \leqslant k)$.

$$
\begin{aligned}
\text{vol}(A_i \triangle S_i) &\leqslant \sum_{i \neq j} \text{vol}\left(\left\{v \in S_i : \|c_i - x_t(v)\| \geqslant \frac{\|c_j - x_t(v)\|}{\log k}\right\}\right) \\
&\quad + \sum_{i \neq j} \text{vol}\left(\left\{v \in S_j : \|c_j - x_t(v)\| \geqslant \frac{\|c_i - x_t(v)\|}{\log k}\right\}\right) \quad (5.14) \\
&\leqslant \text{vol}\left(\left\{v \in S_i : \|c_i - x_t(v)\|^2 \geqslant \frac{1}{1000k \log^2 k \, \text{vol}(S_i)}\right\}\right) \\
&\quad + \sum_{i \neq j} \text{vol}\left(\left\{v \in S_j : \|c_j - x_t(v)\|^2 \geqslant \frac{1}{1000k \log^2 k \, \text{vol}(S_i)}\right\}\right) \quad (5.15) \\
&\leqslant \frac{2000k^3 \log^2 k}{\Upsilon} \text{vol}(S_i). \quad (5.16)
\end{aligned}
$$

where (5.14) follows from Theorem 5.11 by setting $\varepsilon = \log k - 1$, (5.15) follows from Lemma 5.9, and (5.16) follows by Lemma 4.1 and the fact that

$$\sum_{u \in S_j} d_v \|x_t(u) - c_j\|^2 \leqslant 2 \sum_{u \in S_j} d_v \left\|F(u) - p^{(j)}\right\|^2$$

for any $j$.

The bound in the outer conductance of the $A_i$'s follows from the same argument of Lemma 4.7. ∎

# References

[ABS10] Sanjeev Arora, Boaz Barak, and David Steurer. Subexponential algorithms for unique games and related problems. In *51st Annual IEEE Symposium on Foundations of Computer Science (FOCS'10)*, pages 563–572, 2010. 2, 3, 6

[AHK12] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012. 4

[AK07]  Sanjeev Arora and Satyen Kale. A combinatorial, primal-dual approach to semidefinite programs. In *39th Annual ACM Symposium on Theory of Computing (STOC'07)*, pages 227–236, 2007. 4

[AV07]  David Arthur and Sergei Vassilvitskii. *k*-means++: The advantages of careful seeding. In *18th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'07)*, pages 1027–1035, 2007. 20

[AY95]  Charles J. Alpert and So-Zen Yao. Spectral partitioning: The more eigenvectors, the better. In *Discrete Applied Mathematics*, pages 195–200, 1995. 3

[CA79]  Guy B. Coleman and Harry C. Andrews. Image segmentation by clustering. *Proceedings of the IEEE*, 67(5):773–785, 1979. 1

[Chu97]  Fan R. K. Chung. Spectral graph theory. *Regional Conference Series in Mathematics, American Mathematical Society*, 92:1–212, 1997. 18

[Chu09]  Fan R. K. Chung. A local graph partitioning algorithm using heat kernel pagerank. *Internet Mathematics*, 6(3):315–330, 2009. 3, 4

[DK70]  Chandler Davis and William M. Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970. 2

[DRS14]  Tamal K. Dey, Alfred Rossi, and Anastasios Sidiropoulos. Spectral concentration, robust *k*-center, and simple clustering. *CoRR*, abs/1404.1008, 2014. 4

[HJ12]  Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2012. 30

[IM98]  Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *30th Annual ACM Symposium on Theory of Computing (STOC'98)*, pages 604–613, 1998. 20, 26

[KLL$^+$13]  Tsz Chiu Kwok, Lap Chi Lau, Yin Tat Lee, Shayan Oveis Gharan, and Luca Trevisan. Improved cheeger's inequality: analysis of spectral partitioning algorithms through higher order spectral gap. In *45th Annual ACM Symposium on Theory of Computing (STOC'13)*, pages 11–20. ACM, 2013. 1, 2, 3, 9

[KLOS14]  Jonathan A. Kelner, Yin Tat Lee, Lorenzo Orecchia, and Aaron Sidford. An almost-linear-time algorithm for approximate max flow in undirected graphs, and its multicommodity generalizations. In *25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'14)*, pages 217–226, 2014. 1

[KLP12]  Ioannis Koutis, Alex Levin, and Richard Peng. Improved Spectral Sparsification and Numerical Algorithms for SDD Matrices. In *29th International Symposium on Theoretical Aspects of Computer Science (STACS'12)*, volume 14, pages 266–277, Dagstuhl, Germany, 2012. 19

[KSS04]  Amit Kumar, Yogish Sabharwal, and Sandeep Sen. A simple linear time (1+ $\varepsilon$)-approximation algorithm for geometric *k*-means clustering in any dimensions. In *45th Annual IEEE Symposium on Foundations of Computer Science (FOCS'04)*, pages 454–462, 2004. 3

[LM14]  Anand Louis and Konstantin Makarychev. Approximation algorithm for sparsest *k*-partitioning. In *25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'14)*, pages 1244–1255, 2014. 1

[LOGT12]  James R. Lee, Shayan Oveis Gharan, and Luca Trevisan. Multi-way spectral partitioning

and higher-order cheeger inequalities. In *44th Annual ACM Symposium on Theory of Computing (STOC'12)*, pages 1117–1130, 2012. 1, 2, 3, 4, 10

[LPW09] David A. Levin, Yuval Peres, and Elizabeth L. Wilmer. *Markov chains and mixing times*. Providence, R.I. American Mathematical Society, 2009. 18

[LR99] Frank T. Leighton and Satish Rao. Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms. *J. ACM*, 46(6):787–832, 1999. 1

[LRTV12] Anand Louis, Prasad Raghavendra, Prasad Tetali, and Santosh Vempala. Many sparse cuts via higher eigenvalues. In *44th Annual ACM Symposium on Theory of Computing (STOC'12)*, pages 1131–1140, 2012. 1

[MS90] David W. Matula and Farhad Shahrokhi. Sparsest cuts and bottlenecks in graphs. *Discrete Applied Mathematics*, 27(1-2):113–123, 1990. 1

[NJW+02] Andrew Y. Ng, Michael I. Jordan, Yair Weiss, et al. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 2:849–856, 2002. 3

[OGT14] Shayan Oveis Gharan and Luca Trevisan. Partitioning into expanders. In *25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'14)*, pages 1256–1266, 2014. 1, 4

[ORSS12] Rafail Ostrovsky, Yuval Rabani, Leonard J. Schulman, and Chaitanya Swamy. The effectiveness of lloyd-type methods for the $k$-means problem. *J. ACM*, 59(6):28, 2012. 3, 20

[OSV12] Lorenzo Orecchia, Sushant Sachdeva, and Nisheeth K Vishnoi. Approximating the exponential, the lanczos method and an $\widetilde{O}(m)$-time spectral algorithm for balanced separator. In *44th Annual ACM Symposium on Theory of Computing (STOC'12)*, pages 1141–1160, 2012. 3, 4, 19, 20, 21

[OSVV08] Lorenzo Orecchia, Leonard J. Schulman, Umesh V. Vazirani, and Nisheeth K. Vishnoi. On partitioning graphs via single commodity flows. In *40th Annual ACM Symposium on Theory of Computing (STOC'08)*, pages 461–470, 2008. 4

[RST12] Prasad Raghavendra, David Steurer, and Madhur Tulsiani. Reductions between expansion problems. In *27th Conference on Computational Complexity (CCC'12)*, pages 64–73, 2012. 1

[She09] Jonah Sherman. Breaking the multicommodity flow barrier for $O(\sqrt{\log n})$-approximations to sparsest cut. In *50th Annual IEEE Symposium on Foundations of Computer Science (FOCS'09)*, pages 363–372, 2009. 4

[SM00] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(8):888–905, 2000. 1

[SS11] Daniel A. Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. *SIAM J. Comput.*, 40(6):1913–1926, 2011. 19

[ST11] Daniel A. Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on Computing*, 40(4):981–1025, 2011. 1

[Tre08] Luca Trevisan. Approximation algorithms for unique games. *Theory of Computing*, 4(1):111–128, 2008. 1

[VL07] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007. 3

# A    Auxiliary Results

**Theorem A.1** (Geršgorin Circle Theorem). *Let* $\mathbf{A}$ *be an* $n \times n$ *matrix , and let* $R_i(\mathbf{A}) = \sum_{j \neq i} |\mathbf{A}_{i,j}|$, *for* $1 \leqslant i \leqslant n$. *Then, all eigenvalues of* $\mathbf{A}$ *are in the union of Geršgorin Discs defined by*

$$\bigcup_{i=1}^{n} \{z \in \mathbb{C} : |z - \mathbf{A}_{i,i}| \leqslant R_i(\mathbf{A})\} .$$

**Theorem A.2** (Corollary 6.3.4 [HJ12]). *Let* $\mathbf{A}$ *be an* $n \times n$ *normal matrix with eigenvalues* $\lambda_1, \ldots, \lambda_n$ *and* $\mathbf{E}$ *be an* $n \times n$ *matrix. If* $\hat{\lambda}$ *is an eigenvalue of* $\mathbf{A} + \mathbf{E}$, *then there is some eigenvalue* $\lambda_i$ *of* $\mathbf{A}$ *for which* $|\hat{\lambda} - \lambda_i| \leqslant \|\mathbf{E}\|$.

# B    Generalization For Weighted Graphs

Our result can be easily generalized to more general graphs, i.e., undirected weighted graphs for which the edge weights are polynomially bounded. Formally, for any weighted graph $G = (V, E, w : E \to \mathbb{R})$, we define the weighted adjacency matrix $\mathbf{A}$ of $G$ by

$$\mathbf{A}_{u,v} = \begin{cases} w(u,v) & \text{if } \{u,v\} \in E, \\ 0 & \text{otherwise.} \end{cases}$$

where $w(u,v) = w(v,u)$ is the weight on the edge $\{u,v\}$. For every vertex $u \in V$ we define the *weighted degree* of $u$ as $d_u = \sum_{\{u,v\} \in E} w(u,v)$, and the degree matrix $\mathbf{D}$ is defined by $\mathbf{D}_{u,u} = d_u$. We can define the Laplacian matrix $\mathcal{L}$ and the heat kernel $\mathbf{H}_t$ in the same way as in the case of unweighted graphs. Then, it is easy to verify that all the results in Section 5 hold.