

Code Sparsification and its Applications

Sanjeev Khanna*

Aaron (Louie) Putterman†

Madhu Sudan‡

November 3, 2023

Abstract

We introduce a notion of *code* sparsification that generalizes the notion of cut sparsification in graphs. For a (linear) code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k a $(1 \pm \epsilon)$ -*sparsification* of size s is given by a weighted set $S \subseteq [n]$ with $|S| \leq s$ such that for every codeword $c \in \mathcal{C}$ the projection $c|_S$ of c to the set S has (weighted) hamming weight which is a $(1 \pm \epsilon)$ approximation of the hamming weight of c . We show that for every code there exists a $(1 \pm \epsilon)$ -sparsification of size $s = \tilde{O}(k \log(q)/\epsilon^2)$. This immediately implies known results on graph and hypergraph cut sparsification up to polylogarithmic factors (with a simple unified proof) — the former follows from the well-known fact that cuts in a graph form a linear code over \mathbb{F}_2 , while the latter is obtained by a simple encoding of hypergraph cuts. Further, by connections between the eigenvalues of the Laplacians of Cayley graphs over \mathbb{F}_2^k to the weights of codewords, we also give the first proof of the existence of spectral Cayley graph sparsifiers over \mathbb{F}_2^k by Cayley graphs, i.e., where we sparsify the set of generators to nearly-optimal size. Additionally, this work can be viewed as a continuation of a line of works on building sparsifiers for constraint satisfaction problems (CSPs); this result shows that there exist near-linear size sparsifiers for CSPs over \mathbb{F}_p -valued variables whose unsatisfying assignments can be expressed as the zeros of a linear equation modulo a prime p . As an application we give a full characterization of ternary Boolean CSPs (CSPs where the underlying predicate acts on three Boolean variables) that allow for near-linear size sparsification. This makes progress on a question posed by Kogan and Krauthgamer (ITCS 2015) asking which CSPs allow for near-linear size sparsifiers (in the number of variables).

At the heart of our result is a codeword counting bound that we believe is of independent interest. Indeed, extending Karger’s cut-counting bound (SODA 1993), we show a novel decomposition theorem of linear codes: we show that every linear code has a (relatively) small subset of coordinates such that after deleting those coordinates, the code on the remaining coordinates has a smooth upper bound on the number of codewords of small weight. Using the deleted coordinates in addition to a (weighted) random sample of the remaining coordinates now allows us to sparsify the whole code. The proof of this decomposition theorem extends Karger’s proof (and the contraction method) in a clean way, while enabling the extensions listed above without any additional complexity in the proofs.

*School of Engineering and Applied Sciences, University of Pennsylvania, Philadelphia, PA. Email: sanjeev@cis.upenn.edu. Supported in part by NSF awards CCF-1934876 and CCF-2008305.

†School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts, USA. Supported in part by the Simons Investigator Fellowship of Boaz Barak, NSF grant DMS-2134157, DARPA grant W911NF2010021, and DOE grant DE-SC0022199. Supported in part by the Simons Investigator Award of Madhu Sudan and NSF Award CCF 2152413. Email: aputterman@g.harvard.edu.

‡School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts, USA. Supported in part by a Simons Investigator Award and NSF Award CCF 2152413. Email: madhu@cs.harvard.edu.

Contents

1	Introduction	1
1.1	Code Sparsification	1
1.2	Motivation	1
1.3	Main Results	2
1.4	Applications	3
1.5	Proof Techniques	5
1.6	Organization	5
2	A Counting Bound for Codewords	6
2.1	Preliminaries	6
2.2	Overview	6
2.3	A Karger-Style Bound for Codewords	7
3	Preliminaries	9
3.1	Codes	9
3.2	A Probabilistic Bound	9
3.3	Graphs and Graph Sparsification	10
3.4	CSPs and CSP Sparsification	10
4	Near-linear Size Sparsifiers for Polynomially-bounded Codes	11
4.1	Code Decomposition	11
4.2	Code Sparsification Algorithm	12
5	Nearly Linear Size Sparsifiers for Codes of Arbitrary Length	15
5.1	Simple Quadratic Size Sparsifiers	16
5.1.1	Correctness	16
5.1.2	Size Analysis	16
5.2	Removing the $O(\log n)$ factors	17
5.3	Final Algorithm	20
6	Application to Cayley Graph Sparsifiers	23
7	Applications to Sparsifying CSPs	24
7.1	Affine CSPs	24
7.2	Ternary Boolean Predicates	24
8	Application to Hypergraph Cut Sparsifiers	28
9	Conclusions	30
A	A Simpler Construction of Cut Sparsifiers for Graphs	31

1 Introduction

In this work we introduce a notion of sparsification for linear codes, prove that “nearly linear-size” sparsifications exist for every linear code, and give some applications beyond coding theory. We start by recalling the notion of sparsification and give some background.

Sparsification of data with respect to a certain class of queries alludes to a compression of the data that allows all queries in the class to be answered (or estimated) correctly. It has emerged as a fundamental concept in theoretical computer science, with graph sparsification for cut-queries being a central example. In seminal works, Karger [Kar94] and Benczúr and Karger [BK96] showed how graphs may be sparsified by carefully sampling a subset of its edges and assigning weights to them, so that for every cut, the cut-size in the sampled weighted graph gives a good estimate of the cut size in the input graph. Crucially the number of sampled edges was nearly linear-sized in the number of vertices of the graph. This work led to many strengthenings (in particular to allowing the sample to be linear sized [BSS09]), extensions (in particular to more general spectral notions of sparsification [ST11, BSS09], to hypergraphs [KK15, CKN20, KKTY21]), more recently to sparsifying sums of norms [JLLS23], and applications to a host of problems including solvers for max-flow and min-cut [She13, KLOS14, Pen16], as well as to better solvers for structured linear systems [ST11, CKM⁺14, KLP⁺16, LS18, JS20] and applications to clustering [CSWZ16]. In the streaming and sketching settings, small representations of graphs are equally important and work on graph sparsifiers has played a key role.

A natural question that arises is which other classes of structures and queries allow for such sparsification. In this work we explore this question in a new terrain, namely linear codes, where the queries specify a message and the goal is to estimate the Hamming weight of its encoding under the linear code. We describe our setting formally first before motivating the problem.

1.1 Code Sparsification Throughout this paper q will be a prime power and \mathbb{F}_q will denote the finite field on q elements. A linear code \mathcal{C} is a \mathbb{F}_q -linear subspace of \mathbb{F}_q^n , and we will assume it is the image of a linear map $E : \mathbb{F}_q^k \rightarrow \mathbb{F}_q^n$. The Hamming weight of a vector $v \in \mathbb{F}_q^n$, denoted $\text{wt}(v)$ is the number of non-zero coordinates of v . Given a sequence of non-negative (integer) weights $w = (w_1, \dots, w_n)$, the weighted Hamming weight of $v = (v_1, \dots, v_n)$, denoted $\text{wt}_w(v)$, equals $\sum_{i|v_i \neq 0} w_i$. For a vector $v \in \mathbb{F}_q^n$ and set $S \subseteq [n]$ the puncturing of v to the set S , denoted $v|_S$ is the vector $(v_i)_{i \in S}$. The puncturing of a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ to the coordinates $S \subseteq [n]$ is the code $\mathcal{C}|_S \subseteq \mathbb{F}_q^{|S|}$ given by $\mathcal{C}|_S = \{v|_S : v \in \mathcal{C}\}$. We are now ready to define code sparsifiers.

DEFINITION 1.1. (CODE SPARSIFIER) *For integer s , real $\epsilon > 0$ and a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, a $(1 \pm \epsilon)$ -sparsifier of size s for the code \mathcal{C} is a subset $S \subseteq [n]$ with $|S| \leq s$ along with weights $w_S = (w_i)_{i \in S}$ such that for every codeword $v \in \mathcal{C}$, we have*

$$(1 - \epsilon)\text{wt}(v) \leq \text{wt}_w(v|_S) \leq (1 + \epsilon)\text{wt}(v).$$

(In other words the weighted Hamming weight of every codeword in the punctured code roughly equals its weight in the unpunctured code.)

The vanilla representation of a linear code would involve kn elements of \mathbb{F}_q . The sparsification reduces the representation size to sk field elements which may be significantly smaller if $s \ll n$. In several applications we consider later, the code \mathcal{C} itself is obtained by puncturing a known fixed mother code $M \subseteq \mathbb{F}_q^N$. In such cases the vanilla representation of \mathcal{C} (to someone who knows M) would require $n \log N$ bits while the sparsification would require only $s \log N$ bits to describe. Thus in both cases the sparsification definitely compresses the representation of \mathcal{C} . And if we fix any linear encoding scheme $E : \mathbb{F}_q^k \rightarrow \mathbb{F}_q^n$ such that $\mathcal{C} = \{E(m) | m \in \mathbb{F}_q^k\}$, then the sparsification allows us to estimate the hamming weight of the encoding $E(m)$ of every message $m \in \mathbb{F}_q^k$. Thus, the definition of code sparsifiers fits the general notion of sparsification, and so we turn to the motivation for studying this concept.

1.2 Motivation Our initial motivation for studying code sparsification is that it abstracts and generalizes cut sparsification in graphs. Specifically for every graph there is a linear code over \mathbb{F}_2 such that codewords of this code are indicator vectors of the edges crossing cuts in the graph. (This code is obtained by viewing the edge-vertex incidence matrix as the generator of the code.) Thus, a sparsifier for this code corresponds to a cut sparsifier for the associated graph. Existence of graph sparsifiers is typically proved by combinatorial or spectral analysis —

tools that are less amenable to application over codes. Thus the exploration of code sparsification forces us to revisit methods for constructing graph sparsifiers and extract the essential elements in this toolkit.

One broad class of sparsifiers that overlap significantly with code sparsifiers are *CSP sparsifiers*, introduced by Kogan and Krauthgamer [KK15] and studied further by Filtser and Krauthgamer [FK17] and Butti and Zivný [BZ20]. Constraint Satisfaction Problems (CSPs) have as instances n constraints on k variables where each constraint operates on a constant number of variables (called the arity of the constraint). A CSP sparsifier aims to compress an instance of the CSP into a smaller (weighted) one on the same set of variables such that for every assignment to the variables, the sparsified CSP satisfies roughly the same number of constraints as the original one. When the variables take on values in a finite set \mathbb{F}_q and the constraints are linear constraints over \mathbb{F}_q , then the sparsification task is a code sparsification task. (Note that code sparsification allows q as well as the arity of the constraints to be non-constant and so code sparsification is not a subclass of CSP sparsification.) In particular, cut sparsification is also a special case of CSP sparsification. Prior work had shown how to get nearly linear size sparsifiers for CSPs beyond cut sparsifiers. Specifically, in [KK15], it was shown that r -SAT instances on a universe of k Boolean variables admit sparsifiers of size $\tilde{O}(kr/\varepsilon^2)$. This was improved to $\tilde{O}(k/\varepsilon^2)$ by Chen, Khanna and Nagda [CKN20], but the family of CSPs for which this result holds was not broadened. The works [FK17, BZ20] completely classify all binary CSPs (i.e., with arity two), that allow for nearly linear size sparsification, but a classification beyond $r = 2$ remains wide open. Indeed [FK17], pose this as an open question, and [BZ20] highlight the challenge of sparsifying r -XOR CSPs (for $r \geq 3$) as a central problem that remains unaddressed by graph and hypergraph sparsification techniques. Thus, code sparsification seems like the natural next frontier in CSP sparsification and worthy of further attention.

Finally we also give some new applications of code sparsification in this paper itself. In particular, we give a simple reduction from hypergraph cut sparsification to code sparsification that makes the former a special case of the latter (although there may be some polylogarithmic losses in the size of the sparsification). We also show how code sparsification can be used to derive structured sparsification of Cayley graphs over \mathbb{F}_2^k , by other Cayley graphs! We elaborate on these new connections and their implications after describing our main results.

1.3 Main Results Our main theorem in this paper shows that every (possibly weighted) linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k has a nearly-linear sized sparsifier, i.e., one of size $\tilde{O}(\frac{k}{\varepsilon^2} \log q)$.¹

THEOREM 1.1. *For every $\varepsilon > 0$, prime power q , positive integers k and n , every (possibly weighted) linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k has a $(1 \pm \varepsilon)$ -code sparsifier of size $\tilde{O}(\frac{k}{\varepsilon^2} \log q)$.*

Note that the theorem is essentially optimal up to polylogarithmic factors in k (for constant ε and q) in that codes require $\Omega(k)$ -sized sparsifiers. Note also that our result qualitatively reproduces the existential part of the sparsification result in [BK96] while extending it vastly. Of particular interest is the fact that our result does not require the generator matrix of the code to be sparse, something that was evidently true of all previous works on sparsification, and potentially used as a core ingredient in many proofs. In fact, the theorem as stated above is completely independent of the choice of the generator matrix of the code, whereas previous analyses even for the graph-theoretic codes seem to rely on the use of a specific generator matrix.

A central tool in the cut sparsifiers of [Kar94, Kar99, BK96, FHHP11] is “Karger’s cut counting bound” [Kar93, Kar99] which asserts that every graph on k vertices whose minimum cut is c has at most $k^{2\alpha}$ cuts of size at most αc , for every integer α . This bound can be interpreted in coding terms — the minimum cut size is the minimum distance of the corresponding code, and distance is of course a central concept in coding theory. However the lemma is patently false for general codes. Specifically for codes of minimum distance $\Omega(n)$ the bound would suggest that there are at most $n^{\mathcal{O}(1)}$ codewords in the code and every asymptotically good code (those which $k = \Omega(n)$) are counterexamples to this potential extension. In view of the centrality of this bound though, it is natural to ask what weaker bound one can get for general codes. Our next theorem gives a simple weakening that essentially suggests that the only counterexamples come from “good” codes embedded in \mathcal{C} .

THEOREM 1.2. *For every prime power q , parameters d , k and n and every linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, the following holds: there exists a subset $T \subseteq [n]$ with $|T| \leq k \cdot d$ such that for $S = [n] \setminus T$ the code $\mathcal{C}|_S$ satisfies the condition that for every integer $\alpha \geq 1$ the code $\mathcal{C}|_S$ has at most $q^\alpha \cdot \binom{k}{\alpha}$ codewords of weight at most αd .*

¹In this paper we use the notation $\tilde{O}(\cdot)$ to hide poly logarithmic factors in the argument.

Note that in the above theorem, d does not refer to the distance, but rather is a parameter of our choosing.

In other words while \mathcal{C} may have many relatively small weight codewords, they come from a sub-code $\mathcal{C}|_T$ contained on a small set of coordinates while the rest of the code $\mathcal{C}|_S$ has a smooth growth in the number of codewords of small weight. We note that this basic theorem about linear spaces does not seem to have been noticed before and could be of independent interest.

While it is immediate that [Theorem 1.2](#) can be used to get some sparsification for some codes, it is not clear how to use it to go all the way to [Theorem 1.1](#). Indeed in the case of graph sparsification for preserving cuts, known proofs utilize additional notions from graph theory to identify importance of a coordinate (edge). For instance, [BK96] utilizes the notion of edge strengths while [FHHP11] relies on edge connectivity to determine sampling probabilities, and in both cases, the analysis uses graph-theoretic structure to establish correctness of the resulting sparsifiers. We show nevertheless that a simple recursive scheme can be applied to sparsify every code. Indeed even the specialization of this proof to the graph-theoretic case of cut sparsifiers seems new and we describe this simpler proof in [§ A](#).

We remark that one weakness of our results (or a major open question) is that our results are existential and we do not have efficient algorithms to produce the sparsifiers that we show exist. The difficulty roughly emerges from the difficulty of finding and counting low weight codewords in a code which are known hard problems in coding theory.

1.4 Applications

Hypergraph Cut Sparsification. A cut sparsifier for a hypergraph is a simple extension of the notion of a cut sparsifier for graphs. Specifically it is a weighted subgraph of the input hypergraph such that for every 2-coloring of the vertices, the number of bichromatic edges in the original hypergraph is approximately the same as the weight of the bichromatic edges in the subgraph. Previous works by Kogan and Krauthgamer [KK15] (in the constant arity hyperedge case) and ultimately Chen, Khanna and Nagda [[CKN20](#)] (in the unbounded arity case) have given cut-sparsifiers of size $O(k \log(k)/\epsilon^2)$ for every hypergraph on k vertices. We are able to recover their result qualitatively (up to polylogarithmic factors in k and $1/\epsilon$) with a very simple reduction. (Specifically we note that if we choose q to be a large enough prime and associate an r -vertex hyperedge with the vector $(q-r+1, 1, \dots, 1, 0, \dots, 0) \in \mathbb{F}_q^k$ then the only coordinates in encodings of messages in $\{0, 1\}^k$ that are 0 are the monochromatic edges.) See [Remark 8.1](#) for more details. Indeed by applying this reduction to [Theorem 1.2](#) we also obtain a structural decomposition theorem for hypergraphs that does not seem to have been noticed before.

THEOREM 1.3. *For every integer $d \geq 1$, every hypergraph H on k vertices has a set of at most kd hyperedges such that upon their removal, the resulting hypergraph satisfies the condition that for every integer $\alpha \geq 1$ it has at most $(2k)^{2\alpha}$ cuts of size $\leq ad$.*

Indeed, as in the setting of codes, an analog of “Karger’s cut-counting bound” does not hold in the realm of hypergraphs [KK15]. Thus, our analysis of code sparsifiers provides a more universal counting bound which decomposes codes and hypergraphs alike.

Cayley Graph Sparsifiers. A well-studied notion extending that of a cut-sparsifier for graphs is a spectral sparsifier. Formally a spectral sparsifier of a graph is a weighted subgraph whose Laplacian has eigenvalues close to that of Laplacian of the original graph. (The Laplacian of a graph $G = (V, E)$, denoted L_G , is a $|V| \times |V|$ matrix, whose diagonal entries $L_{G,i,i}$ are the degrees of the i th vertex, and whose off diagonal entries $L_{G,i,j}$ are $-w_{i,j}$ where $w_{i,j}$ is the weight of the edge (i, j) in G .) Informally, a spectral sparsifier allows us to estimate the quadratic form $x^T L_G x$ for every real vector x , whereas a cut-sparsifier allows us to estimate this form only for $x \in \{0, 1\}^{|V|}$. Most of the results in this paper only extend the notion of cut-sparsifiers, but not spectral sparsifiers. The only exception is for spectral sparsifiers of “Cayley graphs” on \mathbb{F}_2^k . In this special setting the vertex set of the Cayley graph is \mathbb{F}_2^k , and the edges are specified by a “generating” set $\Gamma \subseteq \mathbb{F}_2^k$. Two vertices $x, y \in \mathbb{F}_2^k$ are adjacent if $x - y \in \Gamma$.

While the general theory of spectral sparsification of course holds for Cayley graphs, this may not lead to a compressed representation of the graph, since the generating set Γ can be much smaller than $|V|$ (and Γ specifies the Cayley graph completely). A natural question in this context would be whether there can be a compression of Cayley graphs that is also a Cayley graph (so leads to a compressed representation of the original graph). While this question remains open for general groups, in the setting of \mathbb{F}_2^k , our main theorem, [Theorem 1.1](#) effectively resolves this positively.

A folklore connection between the eigenvectors of the Cayley graphs and the code generated by Γ (where Γ is viewed an $n \times k$ matrix over \mathbb{F}_2 whose columns generate a code contained in \mathbb{F}_2^n) allows us to show that the weight distribution of a code-sparsifier of Γ closely matches that of Γ , and so leads to a new generating set for the Cayley graph with nearly matching eigenvalue profile. This leads to the following theorem, whose proof may be found in § 6.

THEOREM 1.4. (CAYLEY GRAPH SPECTRAL SPARSIFIER) *For every (possibly weighted) Cayley graph G on \mathbb{F}_2^k with generating set $\Gamma \subseteq \mathbb{F}_2^k$, there exists a weighted sparsifier $\hat{\Gamma} \subseteq \Gamma$, such that for the Cayley graph \hat{G} generated by $\hat{\Gamma}$,*

$$(1 - \varepsilon)L_G \preceq L_{\hat{G}} \preceq (1 + \varepsilon)L_G.$$

Further, $|\hat{\Gamma}| \leq \tilde{O}(k/\varepsilon^2)$.

We stress that this is the first existential result of its kind. For comparison, the work of [BSS09] showed the existence of spectral sparsifiers for any graph on n vertices to size $O(n/\varepsilon^2)$. In the setting of Cayley graphs, this implies the existence of spectral sparsifiers of Cayley graphs over \mathbb{F}_2^k with $O(2^k/\varepsilon^2)$ edges, leading to an average degree which is approximately $O(1/\varepsilon^2)$. However, the key distinction is that the sparsifier returned by [BSS09] is *not* guaranteed to still be a Cayley graph and indeed it can not be. In fact, even just to maintain connectivity, a Cayley graph on \mathbb{F}_2^k requires $\Omega(k)$ generators. Our sparsifiers have degree $\tilde{O}(k)$ (for constant ε), but now our resulting sparsified graph is a Cayley graph.

CSP Sparsification. As described earlier, cut sparsification can be viewed as a special case of CSP sparsification (corresponding to a CSP where constraints apply to two binary variables and require that their XOR be 1). Furthermore when restricted to fields of constant size and generator matrices with exactly r non-zero elements per row (where the columns of the generator matrix generate the code), the code sparsification problem is also a special case of CSP sparsification. Thus this interpretation already leads to a new broad class of CSPs that admit nearly linear sparsifiers.

We say that a predicate $P : \mathbb{F}_q^r \rightarrow \{0, 1\}$ is an *affine* predicate if there exist elements $a_0, a_1, \dots, a_r \in \mathbb{F}_q$ such that $P(b_1, \dots, b_r) = 0$ if and only if $a_0 + \sum_i a_i b_i = 0$ (over \mathbb{F}_q). Equivalently, the predicate $P(b_1, \dots, b_r)$ is evaluating $\sum_i a_i b_i \neq -a_0$.

Now, let \mathcal{P} be a collection of predicates of any arity, and define $\text{CSP}(\mathcal{P})$ to be the family of CSPs where each constraint is a predicate from \mathcal{P} applied to any appropriately sized tuple of variables.

The following theorem asserts that CSPs over affine predicates are sparsifiable.

THEOREM 1.5. *Let $\mathcal{P} = \{P : P \text{ is an affine predicate over } \mathbb{F}_q\}$. Then, any CSP in $\text{CSP}(\mathcal{P})$ admits a nearly linear size $(1 \pm \varepsilon)$ sparsification, namely of size $\tilde{O}_q(k/\varepsilon^2)$, where k denotes the number of variables in the instance (and $O_q(\cdot)$ hides factors of q).*

Note that the above theorem has no dependence on the value r , and in fact, we can sparsify affine predicates even when $r = k$, and simultaneously sparsify affine predicates of different arities as long as they are affine with respect to the same field \mathbb{F}_q .

One immediate application of the above theorem is to any r -XOR constraint. Indeed, the unsatisfying instances of an XOR constraint form a linear subspace over \mathbb{F}_2 , and so we give the first proof of the sparsifiability of XOR constraints to nearly linear size. This addresses one of the open questions of [BZ20], who showed the fundamental inexpressability of XOR constraints in terms of hypergraphs.

To illustrate the power of the theorem above, we also extend a result of [FK17] to predicates on 3 Boolean variables, giving an exact classification of which Boolean predicates on up to three variables are sparsifiable to near-linear size. We say that a predicate $P : \{0, 1\}^r \rightarrow \{0, 1\}$ has an *affine projection to AND* if there exists a function $\pi : [r] \rightarrow \{0, 1, x, \neg x, y, \neg y\}$ such that $\text{AND}(x, y) = P(\pi(1), \dots, \pi(r))$.

THEOREM 1.6. *For a predicate $P : \{0, 1\}^3 \rightarrow \{0, 1\}$, all possible CSPs of P on subsets of k variables are $(1 \pm \varepsilon)$ sparsifiable to size $\tilde{O}(k/\varepsilon^2)$ if and only if P has no affine projection to AND.*

We remark that this classification does not yet extend to arbitrary ternary predicates (specifically over non-Boolean variables), and thus does not extend the result in [BZ20].

1.5 Proof Techniques In our view, one of the strengths of this paper is that the proofs are conceptually simple and short even as they generalize known results and provide some new applications. Indeed once we determine the right form of the weight counting bound, namely [Theorem 1.2](#) (later as [Theorem 2.2](#)), its proof is not very hard.

Roughly, we build a contraction procedure in linear spaces analogous to Karger’s contraction method in [[Kar93](#), [Kar99](#)] in graphs. It is not immediately clear what should be the appropriate analog to Karger’s contraction in our setting. Karger’s method is inherently very reliant on the underlying graph structure, as each contraction “merges” two vertices with an edge between them. In a general linear code, while one can think of each row of the generating matrix as corresponding to an edge, and each column of the generating matrix as corresponding to a vertex, the analogy quickly breaks down as each row can be “involved” with many columns. As such, the perspective we end up taking for “contracting” the generating matrix (and one that extends to the graphical case for edges) is a contraction on a single coordinate j of the generating matrix which is not always 0. We decompose the entire column space of the generating matrix into the entire linear subspace of codewords ($\subseteq \mathbb{F}_q^n$) which is zero on this coordinate (i.e. such that the generating matrix for this subspace is entirely 0 in the j th row), and a single vector which is non-zero in the j th coordinate. The original linear space is the span of these two separate components, and for us, contracting on a coordinate corresponds to keeping only the linear subspace which is zero on this coordinate.

Using this perspective, we show that this contraction procedure when applied to a random, not constantly 0 coordinate of the code (which corresponds to an edge in the underlying graph in Karger’s result) is likely to keep low-weight codewords intact. Thus if we focus on a particular low-weight codeword, repeated application of the contraction procedure should output this codeword with probability at least $q^{-\alpha} \binom{k}{\alpha}^{-1}$ as long as the support of our code never gets too small. We use support here to refer to the total number of remaining not all zero rows of the generating matrix, or equivalently the number of coordinates in the code which are not always 0.

While Karger’s analysis is able to explicitly lower bound the support size (in his case the number of remaining edges) at every iteration in the algorithm by considering the minimum cut size, we can do no such thing as there can exist codes with large dimension and small support contained in our code. Instead, we build a parameterized trade-off saying that whenever our algorithm is unable to make progress in the contraction process, it is because there is a non-trivially high-dimensional code contained in our code that is supported on relatively few coordinates. (See [Theorem 2.1](#) for a precise formulation.) Removing this code and the coordinates it is supported on, and continuing leads to a proof of [Theorem 1.2](#).

The proof of [Theorem 1.1](#) given [Theorem 1.2](#) is also not hard, and we believe this is a simpler proof than previous sparsification results for special cases of graphs and hypergraphs. To prove the existence of near-linear size code sparsifiers, suppose we start with a code \mathcal{C} of dimension k , and length k^2 . We first invoke the above theorem with $d = \sqrt{k}$. Intuitively, this breaks the code into two parts: In one part, we have no guarantee on the distribution of weights of the codewords, but the support size is bounded by $k \cdot d = k^{3/2}$. In the other part, we have a strong bound on the distribution of weights of codewords, but the support can still be as large as k^2 . However, the fact that this code has roughly at most k^α codewords of weight $\alpha\sqrt{k}$ suggests that if we sample $k^{3/2}$ coordinates of this code and give each a weight of \sqrt{k} then we get a pretty good sparsification of this part of the code. Gluing the two parts together gives an $O(k^{3/2})$ size sparsifier of the whole code. To get better sparsifications we can continue to apply this procedure recursively. For instance, if we repeat the process one more level (separately on each of the sparsified codes above), we can now use $d = k^{1/4}$ instead of $d = \sqrt{k}$ and this leads to 4 codes, each of size roughly $k^{5/4}$. Likewise, we can glue these 4 codes together, yielding an $O(k^{5/4})$ size sparsifier. Repeating enough times, and optimizing the parameters carefully leads to the final result. There are additional caveats when moving to codes whose length is super-polynomial in their dimension, and for these we utilize a few additional ideas to get our final sparsification result.

1.6 Organization In § 2, we introduce an analog of Karger’s contraction algorithm using Gaussian elimination to prove our decomposition theorem ([Theorem 1.2](#)) for linear codes. Then, in § 3 we introduce some basic facts and definitions about codes, graphs, and CSPs, which we will utilize in our construction of code sparsifiers and their applications. In § 4 we formalize the argument we gave above regarding sparsifying any code whose length is polynomial in its dimension. Finally, in § 5, we use the algorithm from § 4 as a sub-routine along with several code decomposition tricks to remove any dependence on the length of the code, and present the existence of near-linear size sparsifiers for all linear codes, proving [Theorem 1.1](#). In § 6, we prove [Theorem 1.4](#) by a simple

reduction to codes. In § 7 we prove Theorem 1.5 and Theorem 1.6, extending the classification of near-linearly sparsifiable CSPs again by reducing to the coding case. Finally, in § 8, we interpret hypergraphs in the setting of codes, proving our decomposition result, namely Theorem 1.3.

Finally, as a stand-alone application of our sparsification approach, we provide in § A of the appendix, a self-contained simpler proof of the near-linear size cut sparsifier result of [BK96].

2 A Counting Bound for Codewords

2.1 Preliminaries First, we introduce a few basic definitions. These definitions will be used throughout the paper and are all that is required for this section.

DEFINITION 2.1. *A linear code \mathcal{C} of dimension k and length n is a k -dimensional linear subspace of \mathbb{F}_q^n . We often associate with \mathcal{C} a generator matrix $G \in \mathbb{F}_q^{n \times k}$, which maps vectors $x \in \mathbb{F}_q^k$ to codeword $\in \mathbb{F}_q^n$.*

DEFINITION 2.2. *For a codeword $c \in \mathbb{F}_q^n$, the weight of a codeword is the number of non-zero entries in c . This is will be denoted as $\text{wt}(c)$.*

DEFINITION 2.3. *For a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, we say the support size is*

$$\text{Supp}(\mathcal{C}) = |\{i \in [n] : \exists c \in \mathcal{C} \text{ such that } c_i \neq 0\}|.$$

*In words, it is the number of coordinates of the code which are not always zero. Likewise, we say that for a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ along with a generating matrix $G \in \mathbb{F}_q^{n \times k}$, a coordinate $j \in [n]$ is **non-zero** if there exists a codeword $c \in \mathcal{C}$ such that $c_j \neq 0$.*

DEFINITION 2.4. *For a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, a **subcode** of \mathcal{C} is a linear subspace of \mathcal{C} .*

DEFINITION 2.5. *For a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, we say that the **density** of a code is*

$$\text{Density}(\mathcal{C}) = \frac{\text{Dim}(\mathcal{C})}{\text{Supp}(\mathcal{C})}.$$

We say that the density of the empty code is 0.

2.2 Overview In his seminal work, Karger ([Kar93, Kar99]) showed that for a graph on k vertices with minimum cut-value c , there are at most $k^{2\alpha}$ cuts with size $\leq ac$, for any integer α . As discussed in the introduction, while an analogous statement does not hold for codes in general, we are able to prove a generalization that establishes a bound on the distribution of codewords. Roughly speaking, we show that for a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k and a parameter $d \in \mathbb{Z}^+$ of our choosing, either there exists a dense subcode of \mathcal{C} contained on a small support (i.e. a subcode whose density is $\geq \frac{1}{d}$), or the code satisfies a Karger-style bound on the distribution of codeword weights with parameter d (i.e. there are at most $(qk)^\alpha$ codewords of weight $\leq ad$). An exact statement of this result is given in Theorem 2.1. We will then show that Theorem 2.1 in fact implies Theorem 1.2 described in the introduction; that is, there exists a set of at most kd coordinates, such that upon removing these coordinates (equivalently removing the corresponding rows from the generating matrix), the resulting code satisfies the Karger-style bound with parameter d . An exact statement is provided in Theorem 2.2.

We will prove this bound by describing a “contraction” algorithm akin to that of Karger. Intuitively, the algorithm takes in a generating matrix of a code of dimension k and length n . At random, the algorithm chooses one of the non-zero coordinates of this generating matrix and performs Gaussian elimination to “zero” out all but one of the columns in this coordinate. Then, the algorithm removes the remaining column of the generating matrix that is non-zero in this coordinate, reducing the dimension of the generating matrix by 1. This process is repeated until the dimension of the generating matrix is sufficiently small.

We call this Gaussian elimination step on a random non-zero coordinate of the generating matrix a “contraction”. We will show that as long as the support of the generating matrix is sufficiently large in every iteration, then any low-weight codeword will “survive” (i.e. remain in the span of the generating matrix) after many contractions with high probability. We can then conserve probability mass to argue that in fact, under the condition that the support is never too small, the number of lightweight codewords cannot be too large. This

naturally leads to the statement of Theorem 2.1 as either there exists a subcode of sufficiently high dimension with small support or during our contraction algorithm, the support is always large enough to get a strong bound on the number of codewords.

REMARK 2.1. *When the mother code is the cut-code of a graph, the procedure of choosing a random row and eliminating all but 1 of the non-zero entries is in fact equivalent to Karger's contraction based algorithm ([Kar93, Kar99]).*

First, we will prove some facts about the following algorithm, which takes as input a generating matrix $G \in \mathbb{F}_q^{n \times k}$ for a code of dimension k :

Algorithm 1: Contract(G, α)

```

1 Let  $G_i$  be the  $i$ th column of  $G$ , and let  $k$  be the number of columns of  $G$ , and  $n$  the number of rows.
2 while  $\dim(G) \geq \alpha + 1$  do
3   Choose a random non-zero coordinate  $j \in [n]$  of  $G$ .
4   Let  $G_a$  be the first column of  $G$  where the  $j$ th coordinate is non-zero, and let  $G_{b_1}, \dots, G_{b_p}$  be the
      remaining columns where the  $j$ th coordinate is non-zero.
5   Remove column  $G_a$  from  $G$ , and add  $-G_{j,a}^{-1}G_{b_i,a}G_a$  to each  $G_{b_i}, i \in [p]$ .
6 end

```

CLAIM 2.1. *Given a matrix G of dimension k , after i contractions, the dimension of the column span of G is $k - i$.*

Proof. For the base case, note that after 0 contractions, the dimension is indeed k , as the columns of G are indeed a basis. Now, suppose the claim holds inductively. We will assume that after i contractions, the rank is $k - i$, and show that this holds after the $i + 1$ st contraction. This is clear however, as there are $k - i - 1$ columns after the $i + 1$ st contraction, and if we added the column that we removed back into the matrix, the rank would still be $k - i$ (as every column operation was invertible). If adding one vector can bring the dimension to $k - i$, the dimension before adding the vector must be at least $k - i - 1$. \square

CLAIM 2.2. *For any non-zero coordinate $j \in [n]$ of G that we contract to get G' , the span of G' is exactly all codewords in the column span of G that are zero in coordinate j .*

Proof. First, note that if a code of dimension k' is non-zero in some coordinate j , then it is non-zero in this coordinate in exactly $(q - 1) \cdot q^{k' - 1}$ codewords (and zero in this coordinate in exactly $q^{k' - 1}$ codewords). After we contract on this j th coordinate, we get a new generating matrix G' of dimension $k' - 1$, where the j th row is all 0. Finally, because G' is made by adding columns of G together, the span of G' is contained in the span of G . This means that the span of G' is a $k' - 1$ dimensional subspace of G where the j th row is 0, which is exactly all codewords generated by G that are 0 in their j th coordinate, as the span of G' will have $q^{k' - 1}$ codewords. \square

CLAIM 2.3. *Consider a code \mathcal{C} of dimension k , and a codeword $c \in \mathcal{C}$. If we never contract on any coordinate j where c is non-zero, then c is still in the span of the resulting contracted code.*

Proof. Let G' be the result of performing all the aforementioned contractions. Let the sequence of coordinates we contract on be j_1, \dots, j_k . We know that after each contraction on j_i , the span of the new generator matrix is exactly all remaining codewords of G that are zero in coordinate j_i . Since c is zero on all of j_1, \dots, j_k (because we only ever contract on coordinates where c is 0 by assumption), then c always remains in the span of G after each contraction as the columns we removed are non-zero on these coordinates. So, after all the contractions are performed, c is still in the span of G . \square

2.3 A Karger-Style Bound for Codewords

Here, we will prove the following theorem:

THEOREM 2.1. *For a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k , and any integer $d \geq 1$, at least one of the following is true:*

1. There exists a linear sub-code $\mathcal{C}' \subset \mathcal{C}$ such that $\text{Density}(\mathcal{C}') > \frac{1}{d}$.
2. For all integers α , there are at most $q^\alpha \cdot \binom{k}{\alpha}$ codewords of weight $\leq \alpha d$.

REMARK 2.2. Note that Theorem 2.1 implies Karger's original cut-counting bound as a special case with $q = 2$. For any $c \geq 1$, in any graph with minimum cut size c , the number of edges is necessarily at least $(nc)/2$. So Condition 1 above never arises once we set $d = c/2$, allowing us to recover Karger's original cut-counting bound.

We first prove some sub-claims that will make this easier.

LEMMA 2.1. Suppose we run Algorithm 1, and for some $d \in \mathbb{Z}^+$, after every contraction $\text{Density}(\text{Span}(G)) \leq \frac{1}{d}$. Then, the probability some codeword $c \in \mathcal{C}$ of weight $\leq \alpha d$ (for $\alpha \in \mathbb{Z}^+$) is still in the span of the final contracted G is at least $\binom{k}{\alpha}^{-1}$.

Proof. By Claim 2.3, it follows that if c is in the span of the generating matrix G after i iterations, and we contract on a coordinate where c is non-zero, then c will still be in the span of the contracted generating matrix. So, it follows that the probability c survives is:

$$\Pr[\text{survives first contraction}] \cdots \Pr[\text{survives } (k - \alpha)\text{th contraction}].$$

Using the fact that before the i th contraction, the dimension is $k - i + 1$, we know that the support must be at least $(k - i + 1) \cdot d$ by the assumed density. Because the codeword c survives if we contract on a coordinate where c is 0, the probability of survival in the i th contraction is at least $1 - \frac{\alpha \cdot d}{(k - i + 1) \cdot d}$. Thus,

$$(2.1) \quad \Pr[\text{survives all contractions}] \geq (1 - \alpha/k)(1 - \alpha/(k - 1)) \dots (1 - \alpha/(\alpha + 1))$$

$$(2.2) \quad = \frac{k - \alpha}{k} \cdot \frac{k - 1 - \alpha}{k - 1} \cdots \frac{\alpha + 1 - \alpha}{\alpha + 1} = \binom{k}{\alpha}^{-1}.$$

□

We are now ready to prove our main claim.

Proof. [Proof of Theorem 2.1] Suppose that condition 1 does not hold. Then, every linear sub-code $\mathcal{C}' \subseteq \mathcal{C}$ satisfies $\text{Density}(\mathcal{C}') \leq \frac{1}{d}$. We can then invoke Lemma 2.1 to conclude that any codeword with weight $\leq \alpha \cdot d$ is in the span of the contracted matrix with probability $\geq \binom{k}{\alpha}^{-1}$. Further, note that because the dimension of the contracted matrix is α , there are at most q^α codewords in this span. Because there are q^α codewords in the span of each contracted matrix, the sum of all the probabilities of low weight codewords surviving must be $\leq q^\alpha$. This means there can be at most $q^\alpha \cdot \binom{k}{\alpha}$ codewords of weight $\leq \alpha \cdot d$. □

THEOREM 2.2. For a linear code \mathcal{C} of dimension k and length n over \mathbb{F}_q , for any integer $d \geq 1$, there exists a set of at most $k \cdot d$ coordinates, such that upon their removal, in the resulting code, for any integer $\alpha \geq 1$ there are at most $q^\alpha \cdot \binom{k}{\alpha}$ codewords of weight $\leq \alpha d$.

Proof. As long as Condition 1 of Theorem 2.1 continues to hold, we can write the matrix in the form $\begin{bmatrix} A & B \\ 0 & C \end{bmatrix}$, where the coordinates corresponding to A, B are the subcode of density $> \frac{1}{d}$. We can then remove all these coordinates corresponding to this subspace, and continue repeating this process until Condition 1 no longer holds. Once condition 2 holds, the dimension of the new code will be some value $\leq k$, so the number of codewords of weight $\leq \alpha d$ will be at most $q^\alpha \cdot \binom{k}{\alpha}$. To see why the number of coordinates we remove is at most kd , whenever the subspace we remove has dimension k' , we remove at most $k'd$ coordinates, so after this removal, the resulting code is of dimension $k - k'$. It follows that we can remove at most kd coordinates before the dimension of the code is 0.

Note that after removing coordinates, the resulting code will have dimension $\leq k$, so in particular, the generating matrix will have a non-trivial nullspace. This means that there will be several messages that map to the same codeword. However, if the encoding of two messages is the same, i.e. yielding the same codeword, we do not count these as separate instances. Instead, this bound treats this as a single codeword. □

3 Preliminaries

3.1 Codes

DEFINITION 3.1. For a code $\mathcal{C} \subseteq \mathbb{F}_q^n$, its **distance** is

$$\min_{c \in \mathcal{C}: c \neq 0} \text{wt}(c).$$

DEFINITION 3.2. For a code $\mathcal{C} \subseteq \mathbb{F}_q^n$, the **coordinates** of the code are $[n]$. When we refer to the number of coordinates, this is interchangeable with the length of the code, which is exactly n .

In this work, we will be concerned with code sparsifiers as defined below.

DEFINITION 3.3. For a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ with associated generating matrix G , a $(1 \pm \varepsilon)$ -sparsifier for \mathcal{C} is a subset $S \subseteq [n]$, along with a set of weights $w_S : S \rightarrow \mathbb{R}^+$ such that for any $x \in \mathbb{F}_q^k$

$$(1 - \varepsilon)\text{wt}(Gx) \leq \text{wt}_S(G|_S x) \leq (1 + \varepsilon)\text{wt}(Gx).$$

Here, wt_S is meant to imply that if the codeword is non-zero in its coordinate corresponding to an element $i \in S$, then it contributes $w_S(i)$ to the weight. We will often denote $G|_S$ with the corresponding weights as \tilde{G} .

We next present a few simple results for code sparsification that we will use frequently.

CLAIM 3.1. For a vertical decomposition of a generating matrix,

$$G = \begin{bmatrix} G_1 \\ G_2 \\ \vdots \\ G_k \end{bmatrix},$$

if we have a $(1 \pm \varepsilon)$ sparsifier to codeword weights in each G_i , then their union is a $(1 \pm \varepsilon)$ sparsifier for G .

Proof. Consider any codeword $c \in \text{Span}(G)$. Let c_i denote the restriction to each G_i in the vertical decomposition. It follows that if in the sparsifier $\text{wt}(\hat{c}_i) \in (1 \pm \varepsilon)\text{wt}(c_i)$, then $\text{wt}(\hat{c}) = \sum_i \text{wt}(\hat{c}_i) \in (1 \pm \varepsilon) \sum_i \text{wt}(c_i) = (1 \pm \varepsilon)\text{wt}(c)$. \square

CLAIM 3.2. Suppose \mathcal{C}' is $(1 \pm \delta)$ sparsifier of \mathcal{C} , and \mathcal{C}'' is a $(1 \pm \varepsilon)$ sparsifier of \mathcal{C}' , then \mathcal{C}'' is a $(1 - \varepsilon)(1 - \delta), (1 + \varepsilon)(1 + \delta)$ approximation to \mathcal{C} (i.e. preserves the weight of any codeword to a factor $(1 - \varepsilon)(1 - \delta)$ below and $(1 + \varepsilon)(1 + \delta)$ above).

Proof. Consider any codeword Cx . We know that $(1 - \varepsilon)\text{wt}(\mathcal{C}x) \leq \text{wt}(\mathcal{C}'x) \leq (1 + \varepsilon)\text{wt}(\mathcal{C}x)$. Additionally, $(1 - \delta)\text{wt}(\mathcal{C}'x) \leq \text{wt}(\mathcal{C}''x) \leq (1 + \delta)\text{wt}(\mathcal{C}'x)$. Composing these two facts, we get our claim. \square

Finally, we will use the following claim many times implicitly in our arguments, as we will freely change the generating matrix of the code we are looking at.

CLAIM 3.3. Suppose generating matrices G and G' both generate the same dimension k code \mathcal{C} of length n . Then, if some weighted subset of the rows of G yields a $(1 \pm \varepsilon)$ sparsifier $G|_S = \tilde{G}$, the same weighted subset of the rows of G' yields a $(1 \pm \varepsilon)$ sparsifier $G'|_S = \tilde{G}'$.

Proof. Consider any codeword $c \in \mathcal{C}$. By construction, there is an x, x' such that $Gx = c, G'x' = c$. Now, $G|_S x = c|_S$ and $G'|_S x' = c|_S$. Hence, if \tilde{G} is a $(1 \pm \varepsilon)$ sparsifier of codewords in \mathcal{C} , then so too is \tilde{G}' . \square

3.2 A Probabilistic Bound We will frequently utilize the following probabilistic bound.

CLAIM 3.4. ([FHHP11]) Let X_1, \dots, X_ℓ be random variables such that X_i takes on value $1/p_i$ with probability p_i , and is 0 otherwise. Also, suppose that $\min_i p_i \geq p$. Then, with probability at least $1 - 2e^{-0.38\varepsilon^2\ell p}$,

$$\sum_i X_i \in (1 \pm \varepsilon)\ell.$$

3.3 Graphs and Graph Sparsification We recall here a few basic concepts about graphs and graph sparsification.

DEFINITION 3.4. A cut in a graph $G = (V, E)$ is a subset $S \subseteq V$. We will specify the size of a cut $|\delta_G(S)|$ to be the number of edges that cross from S to $V - S$ (i.e. the number of edges that go from a set S to the rest of the vertices). When a graph $G = (V, E)$ also has an associated weight function $w : E \rightarrow \mathbb{R}^+$, we take $|\delta_G(S)|$ to be the sum over all the edges that cross from S to $V - S$ of the weights of these crossing edges.

DEFINITION 3.5. A $(1 \pm \varepsilon)$ cut-sparsifier for a weighted graph $G = (V, E)$ is a new, weighted graph $\hat{G} = (V, \hat{E})$, with associated weight function w , such that

1. $\hat{E} \subseteq E$.

2. For every cut $S \subseteq V$,

$$(1 - \varepsilon)|\delta_G(S)| \leq |\delta_{\hat{G}}(S)| \leq (1 + \varepsilon)|\delta_G(S)|.$$

The famous result of Karger relates the size of the minimum cut to the number of cuts of other sizes.

DEFINITION 3.6. The minimum cut of a graph $G = (V, E)$ is

$$\min_{S \subseteq V: S \neq V, S \neq \emptyset} |\delta_G(S)|.$$

THEOREM 3.1. (KARGER'S CUT-COUNTING BOUND) [Kar93, Kar99] Suppose a graph $G = (V, E)$ has n vertices, and minimum cut value c . Then, for any integer α , the number of cuts of size at most αc is at most $n^{2\alpha}$.

DEFINITION 3.7. (CUT CODE) For intuition, we will several times refer to the “cut code” of a corresponding graph $G = (V, E)$. Intuitively, this is the code over \mathbb{F}_2 with a generating matrix on $|V|$ columns, where for every edge $e = (u, v) \in E$, we add a row in the generating matrix which has a 1 in the columns corresponding to u and v . If we denote this generating matrix by G' , it can be verified that for any cut $(S, V - S)$,

$$|\delta_G(S)| = \text{wt}(G' \mathbf{1}_S).$$

That is, the weight of the codeword corresponding to the encoding of the indicator vector of S is exactly the number of edges crossing the cut $(S, V - S)$.

3.4 CSPs and CSP Sparsification We formally define here CSPs and CSP sparsification. We start by introducing the notion of a predicate.

DEFINITION 3.8. A predicate P of arity r is a function going from $\{0, 1\}^r \rightarrow \{0, 1\}$.

We will look at CSPs which are defined using a single predicate.

DEFINITION 3.9. A CSP over k variables with predicate P , is a collection of constraints of the form $P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)})$ where r is the arity of P , and i ranges from 1 to m .

DEFINITION 3.10. The value obtained by the CSP on an assignment x is correspondingly

$$\sum_{i=1}^m P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)})$$

In some cases, the CSP has a corresponding weight $w_i \in \mathbb{R}^+$ for each constraint. In this case, the value of the CSP on assignment x is

$$\sum_{i=1}^m w_i \cdot P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)}).$$

DEFINITION 3.11. A $(1 \pm \varepsilon)$ -sparsifier for a CSP C on m constraints, is a new CSP \hat{C} specified by a subset $T \subseteq [m]$, along with weights $(w_i)_{i \in T}$, such that for any assignment $x \in \{0, 1\}^k$,

$$(1 - \varepsilon) \sum_{i=1}^m P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)}) \leq \sum_{i \in T} w_i P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)}) \leq (1 + \varepsilon) \sum_{i=1}^m P(x_1^{(i)}, x_2^{(i)}, \dots, x_r^{(i)}).$$

In words, we choose a weighted subset of the constraints of C , such that for any assignment x , the value of the CSP is preserved to a $(1 \pm \varepsilon)$ factor.

Note that in some works, sparsifying a CSP is meant to only preserve satisfiability while reducing the number of constraints. In our setting, the goal is to approximately preserve the value of the satisfied constraints.

DEFINITION 3.12. For a universe of variables $x \in \{0, 1\}^k$, an affine projection is a restriction of the variables of the form $x_i = 0, x_i = 1, x_i = x_j$ or $x_i = \neg x_j$.

REMARK 3.1. We say that an affine projection of a predicate P of arity r yields an AND of arity 2 if there exists a function $\pi : [r] \rightarrow \{0, 1, x, \neg x, y, \neg y\}$ such that $\text{AND}(x, y) = P(\pi(1), \dots, \pi(r))$.

For instance, the predicate $P : \{0, 1\}^3 \rightarrow \{0, 1\}$, with the only satisfying assignments 000, 001 is equal to an AND of arity 2 under affine projections. If we consider the restriction R which sets the third variable equal to 0, we get a new predicate $P|_R$ whose only satisfying assignment is 00. Thus, this predicate $P|_R(y_1, y_2) = \neg y_1 \wedge \neg y_2$.

4 Near-linear Size Sparsifiers for Polynomially-bounded Codes

In this section, we will prove the existence of near-linear size sparsifiers for codes of length $k^{O(1)}$, where k is the dimension of the code. That is, when \mathcal{C} is a linear code of dimension k and length $k^{O(1)}$. We will also assume that \mathcal{C} is unweighted (or that every coordinate has the same weight). The main theorem to be proved in this section is Theorem 4.1 showing the result for polynomial length (which is a specific case of the more general Theorem 4.2 which applies to codes of any length, but loses extra factors, also proved here). In particular, this is a special case of Theorem 1.1 proved in the introduction.

Intuitively, the proof will use Theorem 2.2 to repeatedly decompose the code \mathcal{C} . In each application, we invoke Theorem 2.2 with a specific choice of parameter d to decompose the generator matrix for the code \mathcal{C} into the form

$$\begin{bmatrix} A & B \\ 0 & C \end{bmatrix},$$

where A is exactly the low-dimensional subcode with bounded support. We will argue that we can keep all of the coordinates of A because the support of A is sufficiently bounded. In turn, this means that B is effectively preserved as well, so as long as we can get a $(1 \pm \varepsilon)$ approximation to C of sufficiently small size, we will be okay. To deal with C , we argue that the distribution of the weights of codewords in C is sufficiently smooth, such that we may subsample the coordinates at rate roughly $1/d$.

This now yields two separate codes, each with size that is strictly smaller than the starting size. We operate inductively, and recursively break down these two smaller codes. Ultimately, in each recursive step, we take a code of size $k \cdot k^\gamma$, and return two codes of size roughly $k \cdot k^{\gamma/2}$. For an initial code of length $k^{O(1)}$ (i.e. whose length is polynomial in the dimension), after $\log \log k$ levels of recursion, we have $\log k$ codes, each of length roughly k . In turn, we can glue these codes back together, and return a sparsifier for our original code.

4.1 Code Decomposition

First, we consider Algorithm 2:

Algorithm 2: CodeDecomposition(\mathcal{C}, d)

- 1 Let k be the dimension of \mathcal{C} .
 - 2 Let S be the set of coordinates to be removed as specified by Theorem 2.2.
 - 3 Let \mathcal{C}' be the code \mathcal{C} after removing the set of coordinates S .
 - 4 **return** S, \mathcal{C}'
-

CLAIM 4.1. 1. After the termination of the Algorithm 2, $|S| \leq k \cdot d$.

2. The final resulting \mathcal{C}' of Algorithm 2 satisfies condition 2 of Theorem 2.1.

Proof. For item 1: This follows exactly from Theorem 2.2.

For item 2: Because we exited the while loop, Condition 1 no longer holds for \mathcal{C}' . Thus, by Theorem 2.1, Condition 2 must hold. \square

CLAIM 4.2. To get a $(1 \pm \varepsilon)$ code sparsifier for \mathcal{C} , it suffices to get S, \mathcal{C}' from Algorithm 2, and then sample all of the indices in S with probability 1, and get a $(1 \pm \varepsilon)$ -sparsifier for \mathcal{C}' .

Proof. This follows because we are creating a “vertical” decomposition of the code. The coordinates of the code corresponding to S contain a dimension k' subspace, and the remaining coordinates of the code define \mathcal{C}' . Thus, we conclude by using Claim 3.1. \square

4.2 Code Sparsification Algorithm In this section we will present the code sparsification algorithm, and argue its correctness and its sparsity.

Algorithm 3: CodeSparsify($\mathcal{C} \subseteq \mathbb{F}_q^n, k, \varepsilon, \eta$)

```

1 Let  $n$  be the length of  $\mathcal{C}$ .
2 if  $n \leq 100 \cdot k \cdot \eta \log(k) \log(q) / \varepsilon^2$  then
3   | return  $\mathcal{C}$ 
4 end
5 Let  $d = \frac{n\varepsilon^2}{\eta \cdot k \log(k) \log(q)}$ .
6 Let  $S, \mathcal{C}' = \text{CodeDecomposition}(\mathcal{C}, \sqrt{d} \cdot \eta \cdot \log(k) \log(q) / \varepsilon^2)$ . Let  $\mathcal{C}_1 = \mathcal{C}|_S$ . Let  $\mathcal{C}_2$  be the result of
   sampling every coordinate of  $\mathcal{C}'$  at rate  $1/\sqrt{d}$ .
7 return CodeSparsify( $\mathcal{C}_1, k, \varepsilon, \eta$ )  $\cup \sqrt{d} \cdot \text{CodeSparsify}(\mathcal{C}_2, k, \varepsilon, \eta)$ 

```

We will use the following fact, which is a simple extension of a result from Karger [Kar94]. In Karger’s work, it was noted that for a graph with minimum cut value c , one can roughly sample the edges at rate $\log(n)/(c\varepsilon^2)$ and scale the weights of the sampled edges up by $c\varepsilon^2/\log(n)$ while still maintaining a $(1 \pm \varepsilon)$ approximation to the cuts in the graph. In the following claim, we adapt this fact to codes which satisfy a smooth bound on the number of codewords of a given weight.

CLAIM 4.3. Suppose \mathcal{C} is a code of dimension k over \mathbb{F}_q , and let $b \geq 1$ be an integer such that for any integer $\alpha \geq 1$, the number of codewords of weight $\leq \alpha b$ is at most $(qk)^\alpha$. Suppose further that the minimum distance of the code \mathcal{C} is b . Then, sampling the coordinates of \mathcal{C} at rate $\frac{\log(k) \log(q) \eta}{b\varepsilon^2}$ with weights $\frac{b\varepsilon^2}{\log(k) \log(q) \eta}$ yields a $(1 \pm \varepsilon)$ sparsifier with probability $1 - 2^{-(0.19\eta - 110) \log k} \cdot k^{-101}$.

Proof. Consider any codeword c of weight $[\alpha b/2, \alpha b]$ in \mathcal{C} . We know that there are at most $(qk)^\alpha$ codewords that have weight in this range. The probability that our sampling procedure fails to preserve the weight of c up to a $(1 \pm \varepsilon)$ fraction can be bounded by Claim 3.4. Indeed,

$$\Pr[\text{fail to preserve weight of } c] \leq 2e^{-0.38 \cdot \varepsilon^2 \cdot \frac{\alpha b}{2} \cdot \frac{\eta \log(k) \log(q)}{\varepsilon^2 b}} = 2e^{-0.19\alpha\eta \log(k) \log(q)}.$$

Now, let us take a union bound over the at most $(qk)^\alpha$ codewords of weight between $[\alpha b/2, \alpha b]$. Indeed,

$$\begin{aligned} \Pr[\text{fail to preserve any } c \text{ of weight } [\alpha b/2, \alpha b]] &\leq 2^{\alpha \log(qk)} \cdot 2e^{-0.19\alpha\eta \log(k) \log(q)} \\ &\leq 2^{\alpha \cdot (-0.19\eta + 1) \log(k) \log(q)} \\ &\leq 2^{\alpha \cdot (-0.19\eta + 1) \log(k)} \\ &\leq 2^{-(0.19\eta - 110)\alpha \log k} \cdot 2^{-109\alpha \log k} \\ &\leq 2^{-(0.19\eta - 110) \log k} \cdot k^{-109\alpha}, \end{aligned}$$

where we have chosen η to be sufficiently large. Now, by integrating over $\alpha \geq 1$, we can bound the failure probability for any integer choice of α by $2^{-(0.19\eta-110)\log k} \cdot k^{-101}$. \square

LEMMA 4.1. *In Algorithm 3, starting with a code \mathcal{C} of size $dk \log(k) \log(q)/\varepsilon^2$, after i levels of recursion, with probability $1 - 2^i \cdot 2^{-\eta k}$, the code being sparsified at level i , $\mathcal{C}^{(i)}$ has at most*

$$(1 + 1/2 \log \log(k))^i \cdot d^{1/2^i} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2$$

surviving coordinates.

Proof. Let us prove the claim inductively. For the base case, note that in the 0th level of recursion the number of surviving coordinates in $\mathcal{C}^{(0)} = \mathcal{C}$ is $d \cdot k \log(k) \log(q)/\varepsilon^2$, so the claim is satisfied trivially.

Now, suppose the claim holds inductively. Let $\mathcal{C}^{(i)}$ denote a code that we encounter in the i th level of recursion, and suppose that it has at most

$$(1 + 1/2 \log \log(k))^i \cdot d^{1/2^i} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2$$

coordinates. Denote this number of coordinates by ℓ . Now, if this number is smaller than $100k\eta \log(k) \log(q)/\varepsilon^2$, we will simply return this code, and there will be no more levels of recursion, so our claim holds vacuously. Instead, suppose that this number is larger than $100k\eta \log(k) \log(q)/\varepsilon^2$. Let $d' = \frac{\ell\varepsilon^2}{\eta k \log(k) \log(q)} \leq (1 + 1/2 \log \log(k))^i \cdot d^{1/2^i}$.

Then, we decompose $\mathcal{C}^{(i)}$ into two codes, \mathcal{C}_1 and \mathcal{C}_2 . \mathcal{C}_1 contains coordinates of $\mathcal{C}^{(i)}$ that contain some k' dimensional code on a support of size $\leq k' \cdot \sqrt{d'} \cdot \eta \cdot \log(k) \log(q)/\varepsilon^2$. Because $k' \leq k$, it follows that the number of coordinates in \mathcal{C}_1 , the first code we recurse on, is at most $(1 + 1/2 \log \log(k))^i \cdot d'^{1/2^{i+1}} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2$ as we desire.

For \mathcal{C}_2 , we define random variables $X_1 \dots X_\ell$ for each coordinate in the support of \mathcal{C}_2 . X_i will take value 1 if we sample coordinate i , and it will take 0 otherwise. Let $X = \sum_{i=1}^\ell X_i$, and let $\mu = \mathbb{E}[X]$. Note that

$$\frac{\mu^2}{\ell} = \left(\frac{\ell}{\sqrt{d'}} \right)^2 / \ell = \frac{\ell}{d'} \geq \eta \cdot k \cdot \log(k) \log(q)/\varepsilon^2.$$

Now, using Chernoff,

$$\Pr[X \geq (1 + 1/2 \log \log(k))\mu] \leq e^{\frac{-2}{4 \log^2 \log(k)} \cdot \eta \cdot k \cdot \log(k) \log(q)/\varepsilon^2} \leq 2^{-\eta k},$$

as we desire. Since $\mu = \ell/\sqrt{d'} \leq (1 + 1/2 \log \log(k))^i \cdot d^{1/2^{i+1}} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2$, we conclude our result.

Now, to get our probability bound, we also operate inductively. Suppose that up to recursive level $i-1$, all sub-codes have been successfully sparsified to their desired size. At the i th level of recursion, there are at most 2^{i-1} codes which are being probabilistically sparsified. Each of these does not exceed its expected size by more than the prescribed amount with probability at most $2^{-\eta k}$. Hence, the probability all codes will be successfully sparsified up to and including the i th level of recursion is at least $1 - 2^{i-1}2^{-\eta k} - 2^{i-1}2^{-\eta k} = 1 - 2^i2^{-\eta k}$. \square

LEMMA 4.2. *For any iteration of Algorithm 3 called on a code \mathcal{C} , $\mathcal{C}_1 \cup \sqrt{d} \cdot \mathcal{C}_2$ is a $(1 \pm \varepsilon)$ approximation to \mathcal{C} with probability at least $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$.*

Proof. First, we note that $\mathcal{C}', \mathcal{C}|_S$ (as returned from Algorithm 2) form a *vertical* decomposition of \mathcal{C} . Hence, it suffices to show that \mathcal{C}_1 is a $(1 \pm \varepsilon)$ -sparsifier to $\mathcal{C}|_S$, and \mathcal{C}_2 is a $(1 \pm \varepsilon)$ -sparsifier to \mathcal{C}' .

Seeing that \mathcal{C}_1 is a $(1 \pm \varepsilon)$ -sparsifier to $\mathcal{C}|_S$ is trivial, as $\mathcal{C}_1 = \mathcal{C}|_S$, since we preserve every coordinate with probability 1.

To see that \mathcal{C}_2 is a $(1 \pm \varepsilon)$ -sparsifier to \mathcal{C}' , first note that every codeword in \mathcal{C}' is of weight at least $\sqrt{d} \cdot \eta \cdot \log(k) \log(q)/\varepsilon^2$. This is because if there were a codeword of weight smaller than this, there would exist a subcode of \mathcal{C}' with dimension 1, and support bounded by $\sqrt{d} \cdot \eta \cdot \log(k) \log(q)/\varepsilon^2$. But, because we used Algorithm 2, we know that there can be no such sub-code remaining in \mathcal{C}' . Thus, every codeword in \mathcal{C}' is of weight at least $\sqrt{d} \cdot \eta \cdot \log(k) \log(q)/\varepsilon^2$.

Now, we can invoke Claim 4.3 with $b = \sqrt{d}\eta \log(k) \log(q)/\varepsilon^2$. Note that the hypothesis of Claim 4.3 is satisfied by virtue of our code decomposition. Indeed, we removed coordinates of the code such that in the resulting \mathcal{C}_2 ,

for any $\alpha \geq 1$, there are at most $(qk)^\alpha$ codewords of weight $\leq \alpha\sqrt{d}\eta \log(k) \log(q)/\varepsilon^2$. Using the concentration bound of Claim 4.3 yields that with probability at least $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$, the resulting sparsifier for \mathcal{C}_2 is a $(1 \pm \varepsilon)$ sparsifier, as we desire.

□

COROLLARY 4.1. *If Algorithm 3 achieves maximum recursion depth ℓ when called on a matrix \mathcal{C} , and $\eta > 600$, then the result of the algorithm is a $(1 \pm \varepsilon)^\ell$ sparsifier to \mathcal{C} with probability $\geq 1 - (2^\ell - 1) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$*

Proof. We prove the claim inductively. Clearly, if the maximum recursion depth reached by the algorithm is 0, then we have simply returned the code itself. This is by definition a $(1 \pm \varepsilon)^0$ sparsifier to itself.

Now, suppose the claim holds for maximum recursion depth $i - 1$. We will show it holds for maximum recursion depth i . Let the code we are sparsifying be \mathcal{C} . We break this into \mathcal{C}_1 , \mathcal{C}_2 , and sparsify these. By our inductive claim, with probability $1 - (2^{i-1} - 1) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$ each of the sparsifiers for $\mathcal{C}_1, \mathcal{C}_2$ are $(1 \pm \varepsilon)^{i-1}$ sparsifiers. Now, by Lemma 4.2 and our value of η , $\mathcal{C}_1, \mathcal{C}_2$ themselves together form a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} with probability $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$. So, by using Claim 3.2, we can conclude that with probability $1 - (2^i - 1) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$, the result of sparsifying $\mathcal{C}_1, \mathcal{C}_2$ forms a $(1 \pm \varepsilon)^i$ approximation to \mathcal{C} , as we desire. □

We can then state the main theorem from this section:

THEOREM 4.1. *For a code \mathcal{C} on alphabet \mathbb{F}_q of dimension k , and length $k^{O(1)}$, Algorithm 3 creates a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} of size $O(k\eta \log^2(k) \log(q)(\log \log(k))^2/\varepsilon^2)$ with probability $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-100}$.*

Proof. For a code of dimension k , and length $k^{O(1)}$, this means that our value of d as specified in the first call to Algorithm 3 is at most $k^{O(1)}$ as well. As a result, after only $\log \log k$ iterations, $d = k^{O(1)/2^{\log \log k}} = k^{O(1)/\log k} = O(1)$. So, by Corollary 4.1, because the maximum recursion depth is only $\log \log k$, it follows that with probability at least $1 - (2^{\log \log k} - 1) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-101} \geq 1 - (\log k) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$, the returned result from Algorithm 3 is a $(1 \pm \varepsilon)^{\log \log k}$ sparsifier for \mathcal{C} .

Now, by Lemma 4.1, with probability $\geq 1 - 2^{\log \log k} \cdot 2^{-\eta k} \geq 1 - \log(k)2^{-\eta k} \geq 1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-101}$, every code at recursive depth $\log \log k$ has at most

$$(1 + 1/2 \log \log(k))^{\log \log k} \cdot d^{1/\log k} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2 = O(k\eta \log(k) \log(q)/\varepsilon^2)$$

coordinates. Because the ultimate result from calling our sparsification procedure is the *union* of all of the leaves of the recursive tree, the returned result has size at most

$$2^{\log \log k} \cdot O(k\eta \log(k) \log(q)/\varepsilon^2) = O(k\eta \log^2(k) \log(q)/\varepsilon^2),$$

with probability at least $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-100}$.

Finally, note that we can replace ε with a value $\varepsilon' = \varepsilon/2 \log \log k$. Thus, the resulting sparsifier will be a $(1 \pm \varepsilon')^{\log \log k} \leq (1 \pm \varepsilon)$ sparsifier, with the same high probability.

Taking the union bound of our errors, we can conclude that with probability $1 - 2^{-(0.19\eta-110)\log k} \cdot k^{-100}$, Algorithm 3 returns a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} that has at most $O(k\eta \log^2(k) \log(q)(\log \log(k))^2/\varepsilon^2)$ coordinates. □

THEOREM 4.2. *For a code \mathcal{C} of dimension k , and length n over \mathbb{F}_q , Algorithm 3 creates a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} with probability $1 - \log(n) \cdot 2^{-(0.19\eta-110)\log k} \cdot k^{-100}$ with at most*

$$O(k\eta \log(k) \log(q) \log^2(n)(\log \log(n))^2/\varepsilon^2)$$

coordinates.

Proof. For a code of dimension k , and length n , this means that our value of d as specified in the first call to Algorithm 3 is at most n as well. As a result, after only $\log \log n$ iterations, $d = n^{1/2^{\log \log n}} = n^{1/\log n} = O(1)$. So, by Corollary 4.1, because the maximum recursion depth is only $\log \log n$, it follows that with probability at least

$1 - (2^{\log \log n} - 1) \cdot 2^{-(0.19\eta - 110)\log k} \cdot k^{-101}$, the returned result from Algorithm 3 is a $(1 \pm \varepsilon)^{\log \log n}$ sparsifier for \mathcal{C} .

Now, by Lemma 4.1, with probability $\geq 1 - 2^{\log \log n} \cdot 2^{-\eta k} \geq 1 - \log(n) \cdot 2^{-(0.19\eta - 110)\log k} \cdot 2^{-k}$, every code at recursive depth $\log \log n$ has at most

$$(1 + 1/2 \log \log(k))^{\log \log n} \cdot n^{1/\log n} \cdot \eta \cdot k \log(k) \log(q)/\varepsilon^2 = O(k\eta \log(k) \log(q) \cdot e^{\frac{\log \log n}{\log \log k}}/\varepsilon^2)$$

coordinates. Because the ultimate result from calling our sparsification procedure is the *union* of all of the leaves of the recursive tree, the returned result has size at most

$$\log(n) \cdot e^{\frac{\log \log n}{\log \log k}} \cdot O(k\eta \log(k) \log(q)/\varepsilon^2) = O(k\eta \log(k) \log(q) \log^2(n)/\varepsilon^2),$$

with probability at least $1 - \log(n) \cdot 2^{-(0.19\eta - 110)\log k} \cdot k^{-101}$.

Finally, note that we can replace ε with a value $\varepsilon' = \varepsilon/2 \log \log n$. Thus, the resulting sparsifier will be a $(1 \pm \varepsilon')^{\log \log n} \leq (1 \pm \varepsilon)$ sparsifier, with the same high probability.

Taking the union bound of our errors, we can conclude that with probability $1 - \log(n) \cdot 2^{-(0.19\eta - 110)\log k} \cdot k^{-100}$, Algorithm 3 returns a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} that has at most

$$O(k\eta \log(k) \log(q) \log^2(n)(\log \log(n))^2/\varepsilon^2)$$

coordinates. \square

However, as we will address in the next section, this result is not perfect:

1. For large enough n , there is no guarantee that this probability is ≥ 0 unless η depends on n .
2. For large enough n , $\log^2(n)$ may even be larger than k .

5 Nearly Linear Size Sparsifiers for Codes of Arbitrary Length

In this section, our goal is to prove Theorem 1.1 (the exact version proved will be Theorem 5.1). We will do this by using Theorem 4.2 as a sub-routine in another algorithm.

In the previous section, we saw an algorithm which produces a near-linear size sparsifier for codes of length polynomial in the dimension. However, if we start with a code of arbitrary length n , simply applying the algorithm from the previous section led to spurious $\log(n)$ factors, which unfortunately can dwarf k . In this section, we will show how we can be a little more careful with our sparsifier to avoid these extra $\log(n)$ factors. To do this, in § 5.1, we will showcase a simple one-shot algorithm that returns a size $O(k^2 \log(q)/\varepsilon^2)$ weighted $(1 \pm \varepsilon)$ sparsifier of any code of dimension k and length n . Ideally, we could simply use this sparsification and compose on top of it Algorithm 3. However, as written, Algorithm 3 only works for *unweighted* codes (or codes where every coordinate has the same weight).

Further, the ratio of the weights that are returned by this sparsifier is unbounded in k (and in many cases will be as large as $\Omega(n)$). Naive notions of turning the weighted code into an unweighted code unfortunately do not work, as if we try to replace coordinates of weight w with w unweighted coordinates, we will no longer be guaranteed that length of the code is polynomial, and we will not have gained anything.

Instead, we will show that once we have a weighted sparsifier of size polynomial in the dimension, we can group coordinates together by their weights. That is, we set a parameter $\alpha = \text{poly}(k/\varepsilon)$ sufficiently large, and set the i th group to contain coordinates with weights between $[\alpha^{i-1}, \alpha^i]$. Our key observation is that if a codeword is non-zero in any coordinate in the i th group, then for an appropriately chosen α , the total weighted contribution from any coordinates in groups $i-2, i-3, \dots$ is much less than an $\varepsilon/100$ fraction of the weight coming from group i . Thus if we let i be the largest integer such that a codeword is non-zero in the i th weight group, we can effectively ignore all the coordinates corresponding to weight groups $i-2, \dots$ when sparsifying this codeword. Now, starting with the largest i , we decompose the code into codewords which are non-zero in group i and those which are zero in group i . For those which are non-zero, we can effectively ignore all the coordinates from groups $i-2, i-3, \dots$. This means that all the coordinates we are concerned with have weights in the range $[\alpha^{i-2}, \alpha^i]$. To turn this into an unweighted code, we simply pull out a factor of α^{i-2} , and now for a coordinate of weight w , we can repeat it roughly w times. Because the weights are polynomial in the dimension, the resulting unweighted code is also polynomial in dimension, and we can invoke the results from the previous section.

5.1 Simple Quadratic Size Sparsifiers In this section, we will introduce a one-shot method for sparsifying a code of dimension k and length n on alphabet \mathbb{F}_q that maintains $O(k^2 \log(q)/\varepsilon^2)$ indices of the original code. We state the algorithm here, and then analyze the space complexity and correctness of this algorithm.

Algorithm 4: QuadraticSparsify($\mathcal{C} \subseteq \mathbb{F}_q^n, k, \varepsilon$)

```

1 Let  $n$  be the length of  $\mathcal{C}$ . for  $i = 1, \dots, n$  do
2   | Let  $w_i$  be  $\min_{c \in \mathcal{C}: c_i \neq 0} \text{wt}(c)$ .
3 end
4 Let  $\mathcal{C}'$  be the result of sampling every coordinate of  $\mathcal{C}$  with probability  $\min(1, a \cdot k \log(q)/(\varepsilon^2 w_i))$ , and
   weight  $1/\min(1, a \cdot k/(\varepsilon^2 w_i))$ .
5 return  $\mathcal{C}'$ 
```

5.1.1 Correctness

First, we prove the correctness of this algorithm.

LEMMA 5.1. *For a code \mathcal{C} of dimension k , and a fixed codeword $c \in \mathcal{C}$, Algorithm 4 returns a sparsifier \mathcal{C}' for \mathcal{C} , such that the new weight of c is a $(1 \pm \varepsilon)$ approximation to the old weight with probability at least $1 - 2^{-2k}$.*

Proof. Consider any codeword $c \in \mathcal{C}$ of weight ℓ . Then, by Claim 3.4,

$$\Pr[\mathcal{C}' \text{ does not make a } (1 \pm \varepsilon) \text{ approximation to } c] \leq 2e^{-0.38\varepsilon^2 \ell \frac{a \cdot k \log(q)}{\varepsilon^2 \ell}} = 2e^{-0.38ak \log(q)}.$$

Here, we have used that because w_i is the minimum weight codeword for which a specific coordinate is 1, in a codeword of weight ℓ , every index in the support has $w_i \leq \ell$. Now, by choosing $a = 10$, and taking a union bound over all q^k codewords, we get that with probability $\geq 1 - 2^{-2k}$, Algorithm 4 returns a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} . \square

5.1.2 Size Analysis

Next, we bound the space taken by this algorithm. To do this, we will first need to take advantage of some structural results about codes.

FACT 5.1. *For any linear code, there exists a basis such that codewords of weight ℓ can be written as the sum of codewords of weight $\leq \ell$ from the basis.*

Proof. Fix a basis b_1, \dots, b_k , which maximizes the number of codewords which can be written as the sum of codewords of weight less than itself. Suppose that this basis cannot write some codeword c as a sum of codewords of weight $\leq \text{wt}(c)$. We know that $c = b_1 + \dots + b_k$ for some basis elements. Without loss of generality, assume that the weights of b_1, \dots, b_k are all ordered. Further, assume that all basis elements starting at index $1 \leq j \leq k$ are of weight $\geq \text{wt}(c)$. It follows then that we can swap c with b_k . For any old codeword which required basis element b_k in its decomposition, we can substitute $b_k = b_1 + \dots + b_{k-1} + c$. All of these basis elements are of weight $\leq \text{wt}(b_k)$, so if previously the codeword could be written as the sum of basis elements of weight less than itself, this will still hold true. Further, the new code can express c in its basis with codewords of weight $\leq \text{wt}(c)$ (just by taking itself). Hence, the new basis is better than the original, which is a contradiction. So, the original basis must be able to write every codeword as the sum of codewords of smaller weight. \square

REMARK 5.1. *Fix such a basis b_1, \dots, b_k as specified by the previous fact. Now, consider any coordinate of this code. If we look at the weights of all codewords that are non-zero in this coordinate, the minimum weight codeword must be with one of the basis vectors. This follows easily: any codeword non-zero in its i th coordinate must be the sum of at least 1 basis vector which is non-zero in its i th coordinate. This means that the weight of the codeword must be greater than the weight of the corresponding basis vectors in its sum.*

By the previous remark, when we try to bound the possible values attained by $\max_x \frac{\langle r_i, x \rangle}{\text{wt}(Dx)}$, it will suffice to analyze the possible value attained just by looking at the basis codewords for this special basis. This simplifies things, as we do not have to look at what can happen with possible linear combinations of the codewords. Instead of looking at a basis, we may simply consider a matrix of dimension $n \times k$.

CLAIM 5.1. Fix an $n \times k$ matrix A . Let c_i denote the i th coordinate of c , and let A_j denote the j th column of A . Then,

$$\sum_{i=1}^n \max_{j \in [k]: (A_j)_i = 1} \frac{1}{\text{wt}(c_j)} \leq k.$$

Proof. This follows because each column c_j can only be the “minimizing” column $\text{wt}(c_j)$ times. Each such time, it contributes $\frac{1}{\text{wt}(c_j)}$. There are k columns, so the total contribution is thus bounded by k . This is in fact tight, as the identity matrix will achieve k . \square

CLAIM 5.2. For any linear code \mathcal{C} of dimension k and length n , with a generator matrix G consisting of rows r_i ,

$$\sum_{i=1}^n \max_{c \in \mathcal{C}: c_i = 1} \frac{1}{\text{wt}(c)} \leq k.$$

Proof. This follows by taking the specified codeword basis from Remark 5.1 and invoking Claim 5.1. \square

Finally, we can prove a bound on the size of the sketch.

LEMMA 5.2. With probability $1 - 2^{-k}$, Algorithm 4 does not sample more than $O(k^2 \log(q)/\varepsilon^2)$ coordinates.

The expected number of coordinates that are sampled by Algorithm 4 is

$$\sum_{i=1}^n \min(1, a \cdot k \log(q)/(\varepsilon^2 w_i)) \leq \sum_{i=1}^n a \cdot k \log(q)/(\varepsilon^2 w_i) \leq \frac{ak \log(q)}{\varepsilon^2} \cdot \sum_{i=1}^n 1/w_i \leq \frac{ak^2 \log(q)}{\varepsilon^2}.$$

Note that here we have used Claim 5.2. Finally, by using a Chernoff bound, we can argue that the size of the sketch is not more than double its expected size with probability $1 - 2^{-k}$. Hence, the size of the sketch is at most $O(k^2 \log(q)/\varepsilon^2)$ with probability $1 - 2^{-k}$.

5.2 Removing the $O(\log n)$ factors Similar to [CKN20], we want to remove the extra factors of $\log n$. To this end, we suggest the following procedure upon being given a code of length n and dimension k in Algorithm 5.

Algorithm 5: WeightClassDecomposition($\mathcal{C}, \varepsilon, k$)

- 1 Let $\mathcal{C}' = \text{QuadraticSparsify}(\mathcal{C}, k, \varepsilon/4)$.
 - 2 Let $\alpha = \frac{k^3 \log(q)}{\varepsilon^3}$.
 - 3 Let E_i be all coordinates of \mathcal{C}' that have weight between $[\alpha^{i-1}, \alpha^i]$.
 - 4 Let $\mathcal{D}_{\text{odd}} = E_1 \cup E_3 \cup E_5 \cup \dots$, and let $\mathcal{D}_{\text{even}} = E_2 \cup E_4 \cup E_6 \cup \dots$
 - 5 **return** $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}}$.
-

Next, we prove some facts about this algorithm.

LEMMA 5.3. Consider a code \mathcal{C} of dimension k and length n . Let

$$\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}} = \text{WeightClassDecomposition}(\mathcal{C}, \varepsilon, k).$$

To get a $(1 \pm \varepsilon)$ -sparsifier for \mathcal{C} , it suffices to get a $(1 \pm \varepsilon/4)$ sparsifier to each of $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}}$.

Proof. First, note that in Algorithm 5, we let $\mathcal{C}' = \text{QuadraticSparsify}(\mathcal{C}, t, \varepsilon/4)$. From before, we know that this will return a $(1 \pm \varepsilon/4)$ sparsifier \mathcal{C}' to \mathcal{C} , of size $O(k^2 \log(q)/\varepsilon^2)$ with probability $\geq 1 - 2^{-k-1}$. Now, the creation of $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}}$ forms a *vertical* decomposition of the code \mathcal{C}' . Thus, by Claim 3.1, if we have a $(1 \pm \varepsilon/4)$ sparsifier for each of $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}}$, we have a $(1 \pm \varepsilon/4)$ sparsifier to \mathcal{C}' , so by Claim 3.2 we have a $(1 \pm \varepsilon)$ approximation to \mathcal{C} (with probability $1 - 2^{k-1}$). \square

Because of the previous claim, it is now our goal to create sparsifiers for $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}}$. Without loss of generality, we will focus our attention only on $\mathcal{D}_{\text{even}}$, as the procedure for \mathcal{D}_{odd} is exactly the same (and the proofs will be the same as well). At a high level, we will take advantage of the fact that

$$\mathcal{D}_{\text{even}} = E_2 \cup E_4 \cup \dots,$$

where each E_i contains edges of weights $[\alpha^{i-1}, \alpha^i]$, for $\alpha = \frac{k^3 \log(q)}{\varepsilon^3}$. Because the returned result from Quadratic sparsify has at most $O(k^2 \log(q)/\varepsilon^2)$ edges with high probability, whenever a codeword $c \in \mathcal{C}'$ has a 1 in a coordinate corresponding to E_i , we can effectively ignore all coordinates of lighter weights E_{i-2}, E_{i-4}, \dots . This is because any coordinate in $E_{\leq i-2}$ has weight at most a $\frac{\varepsilon^3}{k^3 \log(q)}$ fraction of any single coordinate in E_i . Because there are at most $O(k^2 \log(q)/\varepsilon^2)$ coordinates in \mathcal{C}' , it follows that the total possible weight of all coordinates in $E_{\leq i-2}$ is still at most a $O(\varepsilon/k)$ fraction of the weight of a single coordinate in E_i . Thus, we will argue that when we are creating a sparsifier for codewords that have a 1 in a coordinate corresponding to some E_i , we will be able to effectively ignore all coordinates corresponding to $E_{\leq i-2}$. To argue this, we will first have to show how to decompose the code into blocks that are non-zero in coordinates in E_i . So, consider the following algorithm:

Algorithm 6: SingleSpanDecomposition($\mathcal{D}_{\text{even}}, \alpha, i$)

- 1 Let E_i be all coordinates of $\mathcal{D}_{\text{even}}$ with weights between α^{i-1} and α^i .
 - 2 Let G be a generating matrix for $\mathcal{D}_{\text{even}}$.
 - 3 Let k' be the rank of $G|_{E_i}$.
 - 4 Let $b_1, \dots, b_{k'}$ be k' linearly independent columns in $G|_{E_i}$.
 - 5 Permute the columns of $\mathcal{D}_{\text{even}}$ so $b_1, \dots, b_{k'}$ become the first k' columns of $G|_{E_i}$.
 - 6 Perform column operations on G to cancel out all remaining columns in $G|_{E_i}$.
 - 7 **return** $G|_{E_i}, G|_{\bar{E}_i}, k'$
-

CLAIM 5.3. *Line 6 in Algorithm 6 is always possible.*

Proof. Because the rank of $G|_{E_i}$ is k' , and the first k' columns are said to be linearly independent, it follows that there exists a sequence of column operations we can do to zero-out all the remaining $k - k'$ columns of $G|_{E_i}$. \square

CLAIM 5.4. *Line 6 in Algorithm 6 does not change the span of the overall generating matrix G .*

Proof. Because the first k' columns remain untouched, and then are added to the remaining $k - k'$ columns, it follows that all these operations can be undone. Since the starting generator matrix G was rank k , it follows that the new generating matrix is also rank k , so the span has not changed. \square

We now describe some structural properties of this decomposition:

CLAIM 5.5. *After a single iteration of Algorithm 6, suppose the result of running the algorithm looks like*

$$G = \begin{bmatrix} A & 0 \\ B & C \end{bmatrix},$$

where A corresponds to the first k' columns of $G|_{E_i}$, and B, C are the remaining coordinates of the decomposition. Then, if a codeword $c \in \mathcal{D}_{\text{even}}$ is 0 on all of the coordinates corresponding to E_i in the above matrix, then c lives entirely in the span of the final $k - k'$ columns.

Proof. Suppose to the contrary that c is 0 on all of the coordinates corresponding to E_i , but requires a non-zero linear combination including the first k' coordinates. Then c would be non-zero in the coordinates corresponding to E_i because A is a set of linearly independent columns for E_i . Thus, any non-zero linear combination including the columns of A will be non-zero in E_i , so c can not include any of the first k' columns. \square

Now, we can continue to decompose the code by repeating Algorithm 6 multiple times.

Algorithm 7: SpanDecomposition($\mathcal{D}_{\text{even}}, \alpha$)

```

1 Let  $\mathcal{D}'_{\text{even}} = \mathcal{D}_{\text{even}}$ .
2 Let  $S = \{\}$ .
3 while  $\mathcal{D}'_{\text{even}}$  is not empty do
4   Let  $k$  be the dimension of  $\mathcal{D}'_{\text{even}}$ . Let  $i$  be the largest integer such that  $E_i$  is non-empty in  $\mathcal{D}'_{\text{even}}$ .
5   Let  $G|_{E_i}, G|_{\bar{E}_i}, k' = \text{SingleSpanDecomposition}(\mathcal{D}'_{\text{even}}, \alpha, i)$ .
6   Let  $H_i$  be the first  $k'$  columns of  $G|_{E_i}$ .
7   Let  $\mathcal{D}'_{\text{even}}$  be the span of the final  $k - k'$  columns of  $G|_{\bar{E}_i}$ .
8   Add  $i$  to  $S$ .
9 end
10 return  $S, H_i$  for every  $i \in S$ 

```

CLAIM 5.6. *Let S, H_i be as returned by Algorithm 7. Then, $\sum_{i \in S} \text{rank}(H_i) = \text{rank}(\mathcal{D}_{\text{even}})$.*

Proof. This follows because in line 6 of Algorithm 7 we set H_i to be the first k' columns of $G|_{E_i}$, and recurse on $\mathcal{D}'_{\text{even}}$ being the span of the remaining $k - k'$ columns. Hence, the total rank is conserved in every inner loop. \square

LEMMA 5.4. *Suppose we have a code of the form $\mathcal{D}_{\text{even}}$ created by Algorithm 5. Then, if we run Algorithm 7 on $\mathcal{D}_{\text{even}}$, to get $S, H_i \forall i \in S$, it suffices to get a $(1 \pm \varepsilon/2)$ sparsifier for each of the H_i in order to get a $(1 \pm \varepsilon)$ sparsifier for $\mathcal{D}_{\text{even}}$.*

Proof. Consider any codeword c in the span of $\mathcal{D}_{\text{even}}$. Let j be the largest integer such that c is non-zero in the coordinates of $\mathcal{D}_{\text{even}}$ corresponding to E_j .

First, we will show that it suffices to approximate the weight of c to $(1 \pm \varepsilon/2)$ on the coordinates in E_j to approximate its weight to $(1 \pm \varepsilon/2)$ overall. Indeed, this follows because any single coordinate in E_j has more weight than all the combined coordinates of E_{j-2}, E_{j-4}, \dots . This is because there are only $O(k^2 \log(q)/\varepsilon^2)$ coordinates in the code total, and by our choice of α in Algorithm 5, any coordinate in E_j is at least $k^3 \log(q)/\varepsilon^3$ of the fraction of the weight of a coordinate in E_{j-2}, E_{j-4}, \dots . So, any single coordinate in E_j contributes $\Omega(k/\varepsilon)$ more weight than all coordinates in E_{j-2}, \dots combined. If we approximate the weight of c in E_j to a $(1 \pm \varepsilon/2)$ fraction then, we are never overestimating the weight of c in the code by more than $(1 + \varepsilon/2)$, and we never underestimate by more than $(1 - \varepsilon/2)(1 - O(\varepsilon/k)) \geq (1 - \varepsilon)$. Hence, this does indeed yield a $(1 \pm \varepsilon)$ approximation to the weight of a codeword c .

Next, we must argue that by creating sparsifiers for each of the H_i , we are indeed approximating the weight of any codeword c to a $(1 \pm \varepsilon/2)$ fraction on the coordinates of $\mathcal{D}_{\text{even}}$ corresponding to E_j . To see why this is true, let us look at a single iteration of Algorithm 6. If WLOG we assume the coordinates of $\mathcal{D}_{\text{even}}$ are sorted by weight, the result of running the algorithm looks like

$$\mathcal{D}_{\text{even}} = \begin{bmatrix} A & 0 \\ B & C \end{bmatrix}.$$

In this case, A is a dimension k' code corresponding to the coordinates of E_j in $\mathcal{D}_{\text{even}}$. We then disregard the matrix B , and iteratively decompose C in the same manner. The key fact is from Claim 5.5. Any codeword which is zero on E_j will in fact live entirely in the span of the final $k - k'$ columns. Thus, it suffices to estimate the weight of c on these final $k - k'$ columns, but because the coordinates of E_j are 0 in these columns, it in fact suffices to simply build a sparsifier for the matrix C in the above decomposition. Thus, inductively, it suffices to continue decomposing C in the above manner. To conclude, if for some c , j is the largest integer for which c is non-zero on E_j , then in each iteration when we decompose the generator into

$$G^{(i)} = \begin{bmatrix} A & 0 \\ B & C \end{bmatrix},$$

c will continue living in the span of the bottom right matrix C , until we finally call Algorithm 6 with parameter j . Then, we are indeed approximating c on the coordinates of $\mathcal{D}_{\text{even}}$ corresponding to E_j , so our argument is complete. \square

Dealing with Bounded Weights Let us consider any H_i that is returned by Algorithm 7, when called with $\alpha = k^3 \log(q)/\varepsilon^3$. By construction, H_i will contain weights only in the range $[\alpha^{i-1}, \alpha^i]$ and will have at most $O(k^2 \log(q)/\varepsilon^2)$ coordinates. In this subsection, we will show how we can turn H_i into an unweighted code with at most $O(k^5 \log^2(q)/\varepsilon^6)$ coordinates. First, note however, that we can simply pull out a factor of α^{i-1} , and treat the remaining graph as having weights in the range of $[1, \alpha]$. Because multiplicative approximation does not change under multiplication by a constant, this is valid. Formally, consider the following algorithm:

Algorithm 8: MakeUnweighted($\mathcal{C}, \alpha, i, \varepsilon$)

- 1 Divide all edge weights in \mathcal{C} by α^{i-1} .
 - 2 Make a new unweighted code \mathcal{C}' by duplicating every coordinate of \mathcal{C} $\lfloor 10w(r)/\varepsilon \rfloor$ times.
 - 3 **return** $\mathcal{C}, \alpha^{i-1} \cdot \varepsilon/10$
-

LEMMA 5.5. *Consider a code \mathcal{C} with weights bounded in the range $[1, \alpha]$. To get a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} it suffices to return a $(1 \pm \varepsilon/10)$ sparsifier for \mathcal{C}' weighted by $\varepsilon/10$, where \mathcal{C}' is the result of calling Algorithm 8 on $\mathcal{C}, \alpha, 1, \varepsilon$.*

Proof. It suffices to show that \mathcal{C}' is $(1 \pm \varepsilon/10)$ sparsifier for \mathcal{C} , as our current claim will then follow by Claim 3.2. Now, to show that \mathcal{C}' is $(1 \pm \varepsilon/10)$ sparsifier for \mathcal{C} , we will use Claim 3.1. Indeed, for every coordinate r in \mathcal{C} , consider the corresponding $\lfloor 10w(r)/\varepsilon \rfloor$ coordinates in \mathcal{C}' . We will show that the contribution from these coordinates in \mathcal{C}' , when weighted by $\varepsilon/10$, is a $(1 \pm \varepsilon/10)$ approximation to the contribution from r .

So, consider an arbitrary coordinate r , and let its weight be w . Then,

$$\frac{10w}{\varepsilon} - 1 \leq \lfloor 10w/\varepsilon \rfloor \leq \frac{10w}{\varepsilon}.$$

When we normalize by $\frac{\varepsilon}{10}$, we get that the combined weight of the new coordinates w' satisfies

$$w - \varepsilon/10 \leq w' \leq w.$$

Because $w \geq 1$, it follows that this yields a $(1 \pm \varepsilon/10)$ sparsifier, and we can conclude our statement. \square

CLAIM 5.7. *Suppose a code \mathcal{C} of length n has weight ratio bounded by α , and minimum weight α^{i-1} . Then, calling Algorithm 8 with error parameter ε yields a new unweighted code of length $O(n\alpha/\varepsilon)$.*

Proof. Each coordinate is repeated at most $O(\alpha/\varepsilon)$ times. \square

5.3 Final Algorithm Finally, we state our final algorithm in Algorithm 9, which will create a $(1 \pm \varepsilon)$ sparsifier for any code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k preserving only $\tilde{O}(k \log(q)/\varepsilon^2)$ coordinates.

First, we analyze the space complexity. WLOG we will prove statements only with respect to $\mathcal{D}_{\text{even}}$, as the proofs will be identical for \mathcal{D}_{odd} .

CLAIM 5.8. *Suppose we are calling Algorithm 9 on a code \mathcal{C} of dimension k . Let $k_{\text{even},i} = \text{rank}(\widehat{H}_{\text{even},i})$ from each call to the for loop in line 5.*

For each call $\widetilde{H}_{\text{even},i} = \text{CodeSparsify}(\widehat{H}_{\text{even},i}, \text{rank}(\widehat{H}_{\text{even},i}), \varepsilon/10, 100(\log(k/\varepsilon) \log \log(q))^2)$ in Algorithm 9, the resulting sparsifier has

$$O(k_{\text{even},i} \log(k_{\text{even},i}) \log^2(k/\varepsilon) \cdot \log^2(k/\varepsilon) \log(q) (\log \log(k/\varepsilon) \log \log(q))^2 / \varepsilon^2)$$

coordinates with probability at least $1 - \log(k \log(q)/\varepsilon) \cdot 2^{-\Omega(\log^2(k/\varepsilon)(\log \log(q))^2)}$.

Proof. We use several facts. First, we use Theorem 4.2. Note that we have replaced the n in the statement of Theorem 4.2 with $k^5 \log^2(q)/\varepsilon^6$ by using Claim 5.7. Indeed, because $\alpha = k^3 \log(q)/\varepsilon^3$, and we started with a weighted code of length $O(k^2 \log(q)/\varepsilon^2)$, it follows that after using Algorithm 8, the support size is bounded by $O(k^5 \log^2(q)/\varepsilon^6)$. We've also added the fact that η is no longer a constant, and instead carries $O((\log(k/\varepsilon) \log \log(q))^2)$, and carried this through to the probability bound. \square

Algorithm 9: FinalCodeSparsify(\mathcal{C}, ε)

```

1 Let  $k$  be the dimension of  $\mathcal{C}$ .
2 Let  $\alpha = k^3 \log(q)/(\varepsilon/2)^3$ , and  $\mathcal{D}_{\text{odd}}, \mathcal{D}_{\text{even}} = \text{WeightClassDecomposition}(\mathcal{C}, \varepsilon, k)$ .
3 Let  $S_{\text{even}}, \{H_{\text{even},i}\} = \text{SpanDecomposition}(\mathcal{D}_{\text{even}}, \alpha)$ .
4 Let  $S_{\text{odd}}, \{H_{\text{odd},i}\} = \text{SpanDecomposition}(\mathcal{D}_{\text{odd}}, \alpha)$ .
5 for  $i \in S_{\text{even}}$  do
6   Let  $\hat{H}_{\text{even},i}, w_{\text{even},i} = \text{MakeUnweighted}(H_{\text{even},i}, \alpha, i, \varepsilon/8)$ .
7   Let  $\tilde{H}_{\text{even},i} = \text{CodeSparsify}(\hat{H}_{\text{even},i}, \text{rank}(\hat{H}_{\text{even},i}), \varepsilon/80, 100(\log(k/\varepsilon) \log \log(q))^2)$ .
8 end
9 for  $i \in S_{\text{odd}}$  do
10  Let  $\hat{H}_{\text{odd},i}, w_{\text{odd},i} = \text{MakeUnweighted}(H_{\text{odd},i}, \alpha, i, \varepsilon/8)$ .
11  Let  $\tilde{H}_{\text{odd},i} = \text{CodeSparsify}(\hat{H}_{\text{odd},i}, \text{rank}(\hat{H}_{\text{odd},i}), \varepsilon/80, 100(\log(k/\varepsilon) \log \log(q))^2)$ .
12 end
13 return  $\bigcup_{i \in S_{\text{even}}} (w_{\text{even},i} \cdot \tilde{H}_{\text{even},i}) \cup \bigcup_{i \in S_{\text{odd}}} (w_{\text{odd},i} \cdot \tilde{H}_{\text{odd},i})$ 

```

LEMMA 5.6. *In total, the combined number of coordinates over $i \in S_{\text{even}}$ of all of the $\tilde{H}_{\text{even},i}$ is at most $\tilde{O}(k \log(q)/\varepsilon^2)$ with probability at least $1 - \log(k \log(q)/\varepsilon) \cdot 2^{-\Omega(\log^2(k/\varepsilon) (\log \log(q))^2)}$.*

Proof. First, we use Claim 5.6 to see that

$$\sum_{i \in S_{\text{even}}} k_{\text{even},i} \leq k,$$

where $k_{\text{even},i} = \text{rank}(\hat{H}_{\text{even},i})$. Thus, in total, the combined length (total number of coordinates preserved) of all the $\tilde{H}_{\text{even},i}$ is

$$\begin{aligned}
& \sum_{i \in S_{\text{even}}} \text{number of coordinates in } \hat{H}_{\text{even},i} \\
& \leq \sum_{i \in S_{\text{even}}} O(k_{\text{even},i} \log(k_{\text{even},i}) \log^2(k/\varepsilon) \cdot \log^2(k/\varepsilon) \log(q) (\log \log(k/\varepsilon) \log \log(q))^2 / \varepsilon^2) \\
& \leq \sum_{i \in S_{\text{even}}} (k_{\text{even},i}) \cdot \tilde{O}(\log^4(k) \log(q) / \varepsilon^2) \\
& = k \cdot \tilde{O}(\log^4(k) \log(q) / \varepsilon^2) \\
& = \tilde{O}(k \log(q) / \varepsilon^2).
\end{aligned}$$

To see the probability bound, we simply take the union bound over all at most k distinct $\tilde{H}_{\text{even},i}$, and invoke Claim 5.8. \square

Now, we will prove that we also get a $(1 \pm \varepsilon)$ sparsifier for $\mathcal{D}_{\text{even}}$ when we run Algorithm 9.

LEMMA 5.7. *After combining the $\tilde{H}_{\text{even},i}$ from Lines 5-8 in Algorithm 9, the result is a $(1 \pm \varepsilon/4)$ -sparsifier for $\mathcal{D}_{\text{even}}$ with probability at least $1 - \log(k \log(q)/\varepsilon) \cdot 2^{-\Omega(\log^2(k/\varepsilon) (\log \log(q))^2)}$.*

Proof. We use Lemma 5.4, which states that to sparsify $\mathcal{D}_{\text{even}}$ to a factor $(1 \pm \varepsilon/4)$, it suffices to sparsify each of the $H_{\text{even},i}$ to a factor $(1 \pm \varepsilon/8)$, and then combine the results.

Then, we use Lemma 5.5, which states that to sparsify any $H_{\text{even},i}$ to a factor $(1 \pm \varepsilon/8)$, it suffices to sparsify $\hat{H}_{\text{even},i}$ to a factor $(1 \pm \varepsilon/80)$, where again, $\hat{H}_{\text{even},i}$ is the result of calling Algorithm 8. Then, we must multiply $\hat{H}_{\text{even},i}$ by a factor $\alpha^{i-1} \cdot \varepsilon/10$.

Finally, the resulting code $\hat{H}_{\text{even},i}$ is now an unweighted code, whose length is bounded by $O(k^5 \log^2(q)/\varepsilon^6)$, with rank $k_{\text{even},i}$. The accuracy of the sparsifier then follows from Theorem 4.2 called with parameter $\varepsilon/80$.

The failure probability follows from noting that we take the union bound over at most $k H_{\text{even},i}$. By Theorem 4.2, our choice of η , and the bound on the length of the support being $O(k^5 \log^2(q)/\varepsilon^5)$, the probability bound follows. \square

The reason for our choice of η is a little subtle. For Theorem 4.2, the failure probability is characterized in terms of the dimension of the code that is being sparsified. However, when we call Algorithm 3 as a sub-routine in Algorithm 9, we have no guarantee that the rank is $\omega(1)$. Indeed, it is certainly possible that the decomposition in H_i creates k different matrices all of rank 1. Then, choosing η to only be a constant, as stated in Theorem 4.2, the failure probability could be constant, and taking the union bound over k choices, we might not get anything meaningful. To amend this, instead of treating η as a constant in Algorithm 3, we set $\eta = 100(\log(k/\varepsilon) \log \log(q))^2$, where now k is the *overall* rank of the code \mathcal{C} , *not* the rank of the current code that is being sparsified H_i . With this modification, we can then attain our desired probability bounds.

THEOREM 5.1. *For any code \mathcal{C} of dimension k and length n , Algorithm 9 returns a $(1 \pm \varepsilon)$ sparsifier to \mathcal{C} with $\tilde{O}(k \log(q)/\varepsilon^2)$ coordinates with probability $\geq 1 - 2^{-\Omega((\log(k/\varepsilon) \log \log(q))^2)} - 2^{-k}$.*

Proof. First, we use Lemma 5.3. This Lemma states that in order to get a $(1 \pm \varepsilon)$ sparsifier to a code \mathcal{C} , it suffices to get a $(1 \pm \varepsilon/4)$ sparsifier to each of $\mathcal{D}_{\text{even}}$, \mathcal{D}_{odd} , and then combine the results.

Then, we invoke Lemma 5.7 to conclude that with probability $\geq 1 - 2^{-\Omega((\log(k/\varepsilon) \log \log(q))^2)}$, Algorithm 9 will produce $(1 \pm \varepsilon/4)$ sparsifiers for $\mathcal{D}_{\text{even}}$, \mathcal{D}_{odd} .

Further, to argue the sparsity of the algorithm, we use Lemma 5.6. This states that with probability $\geq 1 - 2^{-\Omega((\log(k/\varepsilon) \log \log(q))^2)}$, Algorithm 9 will produce code sparsifiers of size $\tilde{O}(k \log(q)/\varepsilon^2)$ for $\mathcal{D}_{\text{even}}$, \mathcal{D}_{odd} .

Finally, we can bound the error probability of the subcall to WeightClassDecomposition in Algorithm 9 by 2^{-k} .

Thus, in total, the failure probability is at most $2^{-k} + 2^{-\Omega((\log(k/\varepsilon) \log \log(q))^2)}$, the total size of the returned code sparsifier is at most $\tilde{O}(k \log(q)/\varepsilon^2)$, and the returned code is indeed a $(1 \pm \varepsilon)$ sparsifier for \mathcal{C} , as we desire.

Note that the returned sparsifier may have some duplicate coordinates because of Algorithm 8. Even when counting duplicates of the same coordinate separately, the size of the sparsifier will be at most $\tilde{O}(k \log(q)/\varepsilon^2)$. We can remove duplicates of coordinates by adding their weights to a single copy of the coordinate. \square

Extension to the weighted case. Finally, note that as stated, this section proves the existence of near-linear size code sparsifiers for *unweighted* codes of any length. However, this extends simply to *weighted* codes of any length, as we can simply repeat coordinates in accordance with their weights, and then sparsify the resulting unweighted code.

We make this more rigorous below:

CLAIM 5.9. *Suppose we are given a weighted code \mathcal{C} of dimension k , where the weight of coordinate i is $w_i \in \mathbb{R}^+$, and the smallest weight of any coordinate is w . Then, if we create a new unweighted code \mathcal{C}' by repeating coordinate i $\lfloor \frac{100w_i}{\varepsilon w} \rfloor$ times, then for any $c \in \mathcal{C}$, the corresponding $c' \in \mathcal{C}'$ satisfies*

$$\frac{w\varepsilon}{100} \cdot \text{wt}(c') \in (1 \pm \varepsilon/10)\text{wt}(c).$$

Further, a $(1 \pm \varepsilon/10)$ -sparsifier to \mathcal{C}' when weighted by $\frac{w\varepsilon}{100}$ yields a $(1 \pm \varepsilon)$ -sparsifier for \mathcal{C} . Hence, there exist sparsifiers of size $\tilde{O}(k \log(q)/\varepsilon^2)$ for any weighted code \mathcal{C} .

Proof. We will consider an arbitrary coordinate i from \mathcal{C} and compare its weight contribution to c in \mathcal{C} versus $c' \in \mathcal{C}'$. In \mathcal{C} , the contribution is w_i . In \mathcal{C}' (after weighting the entire code by $\frac{w\varepsilon}{100}$), the contribution from the duplicates of coordinate i is $\frac{w\varepsilon}{100} \cdot \lfloor \frac{100w_i}{\varepsilon w} \rfloor$. Now, it follows that

$$\frac{w\varepsilon}{100} \cdot \left(\frac{100w_i}{\varepsilon w} \right) \leq \frac{w\varepsilon}{100} \cdot \lfloor \frac{100w_i}{\varepsilon w} \rfloor \leq \frac{w\varepsilon}{100} \cdot \left(\frac{100w_i}{\varepsilon w} + 1 \right),$$

and hence

$$w_i \leq \frac{w\varepsilon}{100} \cdot \lfloor \frac{100w_i}{\varepsilon w} \rfloor \leq w_i \cdot (1 + \frac{w\varepsilon}{w_i 100}) \leq w_i(1 + \varepsilon/100).$$

Hence, (after weighting the entire code by $\frac{w\varepsilon}{100}$), the contribution from the duplicates of coordinate i is within $(1 \pm \varepsilon/100)$ of the desired weight, and hence overall, the corresponding weight of c' is preserved to a $(1 \pm \varepsilon/100)$.

Now, suppose we get a $(1 \pm \varepsilon/10)$ -sparsifier for \mathcal{C}' . Then, this sparsifier preserves the weight of any codeword $c' \in \mathcal{C}'$ to a $(1 \pm \varepsilon/10)$ factor. Hence, for a codeword $c \in \mathcal{C}$, the sparsifier for \mathcal{C}' (when weighted by $w\varepsilon/100$) preserves the weight of c to a factor $(1 \pm \varepsilon)$ by Claim 3.2 (composing approximations).

The size of our resulting sparsifier is $\tilde{O}(k \log(q)/\varepsilon^2)$. This follows because duplicating coordinates does not change the dimension or our field size, only the length. Because after duplicating the coordinates, the code is unweighted, we can invoke Theorem 4.2 to create a $(1 \pm \varepsilon/10)$ -sparsifier of size $\tilde{O}(k \log(q)/\varepsilon^2)$. Note that this sparsifier may still have duplicate coordinates, which we can remedy by combining them together (adding their weights). \square

6 Application to Cayley Graph Sparsifiers

In this section, we will explore the connection between the weights of codewords for an error correcting code \mathcal{C} of dimension k and the eigenvalues of the Laplacian of a Cayley graph on \mathbb{F}_2^k . At a high level, it is well known that there exist Cayley graph expanders on \mathbb{F}_2^k with constant expansion when the degree of the graph is $\Omega(k)$. Any such expander H can be viewed as a sparsifier to the complete Cayley graph G on \mathbb{F}_2^k (where we take the set of generators S to be exactly \mathbb{F}_2^k), under the constraint that the resulting sparsifier still has a Cayley graph structure, and $L_G \approx_\varepsilon L_H$. At a high level, our result says that for *any* Cayley graph G over \mathbb{F}_2^k , there exists a sparsifier H with at most $\tilde{O}(k/\varepsilon^2)$ edges, such that H is still a Cayley graph, and $L_H \approx_\varepsilon L_G$.

Note that by prior work [BSS09], we know that there exist sparsifiers H for any Cayley graph G such that $L_H \approx_\varepsilon L_G$, however this is the first work which proves the existence of such graphs under the restriction that H is also a Cayley graph and of nearly-linear size.

Preliminaries First, we introduce some definitions related to Cayley graphs. In this section, we will fix a binary linear code \mathcal{C} , as well as a generating matrix for \mathcal{C} , denoted by $G_{\mathcal{C}}$. Let r_i denote the i th row of $G_{\mathcal{C}}$.

DEFINITION 6.1. A Cayley graph G is a graph with algebraic structure; its vertex set is defined to be a group, and the edges correspond to a set of generators S , along with weight $(w_i)_{i \in S}$. For every element in $s \in S$, and for every vertex v , there is an edge from v to $v + s$ of weight w_s .

DEFINITION 6.2. Let $\chi_x(r) = (-1)^{\langle x, r \rangle}$, where the inner product is taken modulo 2, and $x, r \in \mathbb{F}_2^k$.

FACT 6.1. For a Cayley graph G defined over \mathbb{F}_2^k with generating set S , there is exactly one eigenvalue of G for every vertex x in its vertex set (which is \mathbb{F}_2^k). The corresponding eigenvalue (of the adjacency matrix) is

$$\lambda_x(G) = \sum_{r \in S} w_r \chi_x(r).$$

The corresponding eigenvector is χ_x . Note that this means that any Cayley graph defined on the same vertex set has the same eigenvectors.

Going forward, we will let G be a Cayley graph over \mathbb{F}_2^k , where its set of generators S is exactly $\{r_1, \dots, r_n\}$, where these are the rows of the generating matrix $G_{\mathcal{C}}$.

FACT 6.2. For a message $x \in \mathbb{F}_2^k$, we have that

$$\text{Bias}(G_{\mathcal{C}}x) = \mathbb{E}_{r \in S} \chi_x(r).$$

The above statement is very intuitive. $\chi_x(r_i)$ is 1 if the i th bit in the codeword corresponding to x is 0, and is -1 if the i th bit is 1. As a result, this expectation is exactly measuring how many more 0's there are than 1's.

FACT 6.3. We can generalize the previous fact to the eigenvalues of the Laplacian. In this way we get that

$$\lambda_x(L_G) = n - \sum_{r \in S} \chi_x(r) = n(1 - \text{Bias}(G_{\mathcal{C}}x)) = 2 \cdot \text{wt}(G_{\mathcal{C}}x).$$

CLAIM 6.1. By preserving the weight of every codeword of \mathcal{C} to a $(1 \pm \varepsilon)$ factor, we preserve the eigenvalues of L_G to a $(1 \pm \varepsilon)$ factor.

Proof. This follows exactly from Fact 6.3. If we have a code sparsifier \hat{C} for \mathcal{C} , then for any $x \in \mathbb{F}_2^k$, $\text{wt}(G_{\hat{\mathcal{C}}}) \in (1 \pm \varepsilon)\text{wt}(G_{\mathcal{C}}x)$. Because the codeword weights of the generating set and the eigenvalues of the Cayley graph are exactly equal, this $(1 \pm \varepsilon)$ approximation to the codeword weights implies that a Cayley graph with the same weighted generating set as used by the code sparsifier would be a $(1 \pm \varepsilon)$ spectral sparsifier by Fact 6.3. \square

7 Applications to Sparsifying CSPs

In this section, we show how to use our result on the sparsifiability of codes in the setting of CSPs. Specifically we first show that all affine predicates are sparsifiable, thus proving Theorem 1.5. Then we use this theorem to classify all Boolean ternary CSPs, CSPs on variables that take values in $\{0, 1\}$ where the predicate applies on three variables. This leads to a proof of Theorem 1.6.

7.1 Affine CSPs Recall that a predicate $P : \mathbb{F}_q^r \rightarrow \{0, 1\}$ is an *affine* predicate if there exist elements $a_0, a_1, \dots, a_r \in \mathbb{F}_q$ such that $P(b_1, \dots, b_r) = 0$ if and only if $a_0 + \sum_i a_i b_i = 0$ (over \mathbb{F}_q). We say further that P is *linear* if $a_0 = 0$.

The proof of Theorem 1.5 is completely straightforward if P is linear, given the definition of a linear code. The extension to the affine case uses a simple reduction from the affine case to the linear case (with one extra variable).

Proof. [Proof of Theorem 1.5] Given an instance Φ of $\text{CSP}(\mathcal{P})$ with variables x_1, \dots, x_k and constraints C_1, \dots, C_n where $C_j = P^{(j)}(x_{(j),1}, \dots, x_{(j),r_j})$, we will create a code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k generated by the matrix $G \in \mathbb{F}_q^{n \times k}$, where each row of the generating matrix corresponds to a single constraint. Let $P^{(j)} \in \mathcal{P}$ denote the predicate showing up in the j th constraint of our CSP instance. We start with the case that $P^{(j)}$ is linear with elements $a_{(j),1}, \dots, a_{(j),r_j} \in \mathbb{F}_q$ being such that $P^{(j)}(b_1, \dots, b_{r_j}) = 0$ if and only if $\sum_i a_{(j),i} b_i = 0$. Then, in the corresponding j th row of the generating matrix, for each $i \in [r_j]$, we place $a_{(j),i}$ in the column corresponding to variable $x_{(j),i}$, and leave all other entries in the row to be 0. It is straightforward to verify that for an assignment $x \in \mathbb{F}_q^k$, $(Gx)_j = 0$ if and only if C_j is unsatisfied. Thus $\text{wt}(Gx)$ counts the number of satisfied constraints for assignment x and thus a code sparsifier for \mathcal{C} is a sparsifier for the instance Φ of $\text{CSP}(\mathcal{P})$.

Now considering the case of a general affine $P^{(j)}$ given by $P^{(j)}(b_1, \dots, b_r) = 0$ if and only if $a_{(j),0} + \sum_i a_{(j),i} b_i = 0$. Now let $\widehat{P^{(j)}}(b_0, \dots, b_r) = \sum_{i=0}^r a_{(j),i} b_i$. Note that $\widehat{P^{(j)}}$ is linear. Given an instance Φ of $\text{CSP}(\mathcal{P})$ on variables x_1, \dots, x_k with constraints C_1, \dots, C_n where $C_j = P^{(j)}(x_{(j),1}, \dots, x_{(j),r_j})$, let $\hat{\Phi}$ be the instance of $\text{CSP}(\mathcal{P})$ on variables x_0, \dots, x_k with constraints $\hat{C}_1, \dots, \hat{C}_n$ given by $\hat{C}_j = \widehat{P^{(j)}}(x_0, x_{(j),1}, \dots, x_{(j),r_j})$. We note that an assignment $x \in \mathbb{F}_q^k$ for Φ corresponds to the assignment $(1, x) \in \mathbb{F}_q^{k+1}$ to $\hat{\Phi}$ and so a sparsifier for $\hat{\Phi}$ (available from the previous paragraph) also sparsifies Φ . This concludes the proof. \square

A major open question from the work of [BZ20] was the sparsifiability of XOR predicates. Even on 3 variables, it was not known if the predicate $P(x_1, x_2, x_3) = x_1 \oplus x_2 \oplus x_3$ was sparsifiable to near-linear size. As a consequence of Theorem 1.5, we get the following result that resolves this question.

COROLLARY 7.1. *On a universe of k variables, any CSP with r -XOR predicates for $1 \leq r \leq k$ is $(1 \pm \varepsilon)$ sparsifiable to size $\tilde{O}(k/\varepsilon^2)$.*

7.2 Ternary Boolean Predicates We now turn to the classification of ternary Boolean predicates. Recall that a predicate $P : \{0, 1\}^r \rightarrow \{0, 1\}$ has an *affine projection to AND* if there exists a function $\pi : [r] \rightarrow \{0, 1, x, \neg x, y, \neg y\}$ such that $\text{AND}(x, y) = P(\pi(1), \dots, \pi(r))$. We wish to prove that $P : \{0, 1\}^3 \rightarrow \{0, 1\}$ is sparsifiable nearly linear size if and only if it has no affine projection to AND (Theorem 1.6).

The hardness result follows from a well-known result showing that the dicut problem is not sparsifiable to subquadratic size [FK17], which in our language is equivalent to saying that the binary AND predicate is not sparsifiable. (We include a precise statement and proof below for completeness — see Lemma 7.1.) Extending this to all predicates that have an affine projection to AND is simple (and holds for general r). This is stated and proved as Lemma 7.2 below. The bulk of the section then does a case analysis and shows that all ternary Boolean predicates that do not have an affine projection to AND can be sparsified by appealing to Theorem 1.5.

Non-sparsifiability of predicates with affine projection to AND.

LEMMA 7.1. ([FK17]) *For every $\epsilon \in [0, 1]$, every $(1 \pm \epsilon)$ -sparsifier of size s for $CSP(AND)$ on k variables requires $s = \Omega(k^2)$.*

The proof is actually more general and shows that any “sketch” of an instance Φ of $CSP(AND)$ requires $\Omega(k^2)$ bits.

Proof. Let S and w_S be a $(1 \pm \epsilon)$ sparsifier of size s of a CSP instance Φ on variables x_1, \dots, x_k . Given a sparsifier, i.e., a subset of the constraints S and a pair $(i, j) \in \binom{k}{2}$, consider the weight of the constraints satisfied by the assignment x^{ij} given by $x_k^{ij} = 1$ if $k \in \{i, j\}$ and 0 otherwise. This weight is positive in Φ if and only if the constraint $x_i \wedge x_j$ appears in Φ . Thus this weight is positive in the weighted sparsified instance using constraints from S if and only if the constraint $x_i \wedge x_j$ appears in Φ (since $\epsilon < 1$); and furthermore the weight is positive in the unweighted sparsified instance on S if and only if the constraint $x_i \wedge x_j$ appears in Φ . (The presence of the weights only affect the weight of the constraints that are satisfied, but not whether the number is positive or not.) Since this is true to every $(i, j) \in \binom{k}{2}$, it follows that S allows us to reconstruct $\Omega(k^2)$ independent bits of information about Φ and thus by the pigeonhole principle $|S| \geq \binom{k}{2} = \Omega(k^2)$ (for some instance Φ). \square

LEMMA 7.2. *For every r , if a predicate $P : \{0, 1\}^r \rightarrow \{0, 1\}$ has an affine projection to AND, then for every $\epsilon \in [0, 1]$, every $(1 \pm \epsilon)$ -sparsifier of size s for $CSP(P)$ requires $s = \Omega_r(k^2)$.*

Proof. Let $\pi : [r] \rightarrow \{0, 1, x, \neg x, y, \neg y\}$ be such that $AND(x, y) = P(\pi(1), \dots, \pi(r))$. Given an instance Φ of $CSP(AND)$ on k variables x_1, \dots, x_k we create an instance Ψ of $CSP(P)$ on $2rk + 2$ variables denoted $G_0, G_1, Y_{i,t,b}$ for $i \in [k], t \in [r], b \in \{0, 1\}$ as follows: For every constraint $AND(x_i, x_j)$ in Φ , we introduce the constraint $P(v_1, \dots, v_r)$ where for $t \in [r]$,

$$v_t = \begin{cases} G_0 & \text{if } \pi(t) = 0 \\ G_1 & \text{if } \pi(t) = 1 \\ Y_{i,t,0} & \text{if } \pi(t) = x \\ Y_{i,t,1} & \text{if } \pi(t) = \neg x \\ Y_{j,t,0} & \text{if } \pi(t) = y \\ Y_{j,t,1} & \text{if } \pi(t) = \neg y \end{cases}.$$

Given an assignment to x_1, \dots, x_k , it can be verified that the resulting predicate simulates $AND(x_i, x_j)$ if $G_0 = 0$, $G_1 = 1$, and $Y_{\ell,t,0} = x_\ell$ and $Y_{\ell,t,1} = \neg x_\ell$ for all $\ell \in [k]$ and $t \in [r]$. Thus a sparsification (S, w_S) of Ψ yields a sparsification of Φ . From the lower bound in Lemma 7.1, we get that $s = \Omega(k^2)$. Relative to k' the number of variables of Ψ we get that $s = \Omega((k'/r)^2) = \Omega_r(k'^2)$ as desired. \square

Sparsifying 3-CSPs with no affine projections to AND. Finally we show that if a predicate $P : \{0, 1\}^3 \rightarrow \{0, 1\}$ has no affine projection to AND, then there exists linear-size sparsifiers for P . We will make use of the following corollary of Theorem 1.5:

COROLLARY 7.2. *Suppose there exists a linear equation $E(x_1, x_2, x_3) = ax_1 + bx_2 + cx_3 + d \pmod{p}$ (for p a prime) such that the unsatisfying assignments to a predicate $P(x_1, x_2, x_3) : \{0, 1\}^3 \rightarrow \{0, 1\}$ are exactly the assignments to x_1, x_2, x_3 such that E evaluates to 0, then a valued CSP containing this predicate can be sparsified to size $\tilde{O}(k \log(p)/\epsilon^2)$.*

We note that the corollary is immediate from Theorem 1.5, which even allows the variables to take values in all of \mathbb{Z}_p while we only need them to take values in $\{0, 1\}$. Restricting the set of assignments preserves sparsification and so the sparsifier from Theorem 1.5 certainly suffices to get Corollary 7.2.

With this in hand, we will now use a case by case analysis to show how to write any predicate $P : \{0, 1\}^3 \rightarrow \{0, 1\}$ with no affine projection to AND as a linear equation modulo some prime. We state four claims that cover the different cases, and prove them in turn. Given the four claims, and Lemma 7.2, the proof of Theorem 1.6 is immediate.

CLAIM 7.1. *If $P : \{0,1\}^3 \rightarrow \{0,1\}$ has zero, six, seven or eight satisfying assignments then P is sparsifiable to nearly linear size.*

CLAIM 7.2. *If $P : \{0,1\}^3 \rightarrow \{0,1\}$ has five satisfying assignments and P has no affine projections to AND then P is sparsifiable to nearly linear size.*

CLAIM 7.3. *If $P : \{0,1\}^3 \rightarrow \{0,1\}$ has four satisfying assignments and P has no affine projections to AND then P is sparsifiable to nearly linear size.*

CLAIM 7.4. *If $P : \{0,1\}^3 \rightarrow \{0,1\}$ has one, two or three satisfying assignments then P has an affine projection to AND.*

Proof. [Proof of [Theorem 1.5](#)] If P has an affine projection to AND, then by [Lemma 7.2](#), P has no subquadratic sized sparsifiers. So assume P has no affine projections to AND. Then by [Claim 7.4](#) (in contrapositive form) P has at least four satisfying assignments. And [Claim 7.1-Claim 7.3](#) show that in all remaining cases P is sparsifiable to nearly linear size. \square

Thus all that remains is to prove [Claim 7.1-Claim 7.4](#). Before turning to the proofs of these claims we mention some basic symmetries that allows us to simplify all the cases.

LEMMA 7.3. *For every predicate $P : \{0,1\}^r \rightarrow \{0,1\}$, permutation $\pi : [r] \rightarrow [r]$ and index $b \in \{0,1\}^r$, let $P'(z_1, \dots, z_r) = P(z_{\pi(1)} \oplus b_1, \dots, z_{\pi(r)} \oplus b_r)$. Then P has a nearly linear sparsifier if and only if P' does.*

Proof. The proof when P' is just a permutation of the variables of P (i.e., when $b = 0^r$) is straightforward - we map any constraint of an instance of $\text{CSP}(P)$ to the constraint obtained by permuting the sequence of variables according to π . It thus suffices to prove the lemma for the case where $b \neq 0^r$ and π is the identity. (By composing the two steps we get the full lemma). For this case we use an idea similar to the idea in the proof of [Lemma 7.2](#).

We first note that we can assume w.l.o.g that variables of the instances to be sparsified come in r blocks and each constraint application applies constraints in which the t th variable comes from the t th block. To see this suppose we have an instance Φ with variables x_1, \dots, x_k and suppose some constraint is $P(x_{i_1}, \dots, x_{i_r})$. We create a new instance Φ' on variables $x_{i,t}$ for $i \in [k]$ and $t \in [r]$ and replace the constraint above by the constraint $P(x_{i_1,1}, \dots, x_{i_r,r})$. We claim a sparsification of Φ' yields a sparsification of Φ . (In Φ we are only interested in assignments in which the r copies of a variable all take on the same value. The sparsification of Φ' yields an estimate of the number of satisfied constraints for all assignments including these).

Once the instances apply constraints to variables from distinct blocks, we can now negate any subset of variables of P . Fix $b \in \{0,1\}^r$ and let $P'(z) = P(z \oplus b)$. Given a canonical instance Φ' of $\text{CSP}(P')$ as above, we can simply let $\hat{\Phi}$ be the $\text{CSP}(P')$ instance where every constraint application applies P' instead of P to the same sequence of variables. To compute the sparsification of Φ' we simply use a sparsification of $\hat{\Phi}$ and use the fact that the weight of constraints satisfied by an assignment $\{a_{i,t}\}_{i,t}$ in $\hat{\Phi}$ is the same as the weight of constraints satisfied by the assignment $\{a_{i,t} \oplus b_t\}_{i,t}$ in Φ' . This allows us to use a sparsification of $\text{CSP}(P')$ to get a sparsification of $\text{CSP}(P)$ (and the other direction follows similarly). \square

We now proceed by cases. By [Corollary 7.2](#), in order to show near-linear size sparsifiability, it suffices to show that the unsatisfying assignments to these predicates can be written as solutions to a linear equation $\mod p$ for some prime p .

Proof. [Proof of [Claim 7.1](#)]

1. P has 0 or 8 Satisfying assignments: in both cases, the predicate is a constant function. So, we can in fact sparsify to a single constraint (just a single constraint with weight equal to the number of total constraints).

2. P has 7 Satisfying assignments: if there are 7 satisfying assignments to P , then this is simply an OR on 3 variables. By the reduction provided in [\[KK15\]](#), along with known hypergraph sparsification results [\[CKN20\]](#), this can be sparsified to size $\tilde{O}(k/\varepsilon^2)$. This can also be shown using [Corollary 7.2](#): W.l.o.g. P is unsatisfied by the all zeros assignment. So $P(x,y,z) = 0$ if and only if $x + y + z = 0 \mod 5$ which fits within the framework of [Corollary 7.2](#).

3. 6 Satisfying assignments to P : if there are 6 satisfying assignments to P , then by replacing x_i with $\neg x_i$, we can assume WLOG that one of the unsatisfying assignments is 000, as argued earlier. There are then 3 cases:

- (a) The other unsatisfying assignment is at distance 3 from 000, i.e. 111 is the other unsatisfying assignment. Then, these are exactly the unsatisfying assignments to $x_1 + x_2 + x_3 \bmod 3$.
- (b) The other unsatisfying assignment is at distance 2 from 000. Then, by our argument before, we can permute the bits of this other unsatisfying assignment, i.e. up to re-ordering 011 is the other unsatisfying assignment. Then, these are exactly the unsatisfying assignments to $x_1 + 2x_2 + 3x_3 \bmod 5$.
- (c) The other unsatisfying assignment is at distance 1 from 000. Then, up to re-ordering / negation, the other unsatisfying assignment is 001. Then, our expression is exactly $x_1 \vee x_2$, which is known to be sparsifiable by [FK17], or similarly, by viewing it as the equation $x_1 + x_2 \bmod 3$.

This concludes the proof of [Claim 7.1](#). □

Proof. [Proof of [Claim 7.2](#)] We consider P that has five satisfying assignments. Again, we break into several cases. By replacing variables with their negations, and possibly reordering the variables, we can always assume one unsatisfying assignment is 000. We consider some cases on the distance to the nearest other satisfying assignment.

1. Suppose one other unsatisfying assignment is at distance 1 from 000. WLOG, let this other unsatisfying assignment be 001. Note that this means the final unsatisfying assignment must start with 11, as otherwise it will contain an AND of arity 2 (to see this, if we assume WLOG for the sake of contradiction that the first bit is a 0, consider the affine projection where $x_1 = 0$, there are 3 unsatisfying assignments remaining under this projection, and hence an AND). Further, by negating the third variable, we can get either 110, or 111 without changing the other two unsatisfying assignments. Thus, we assume the final unsatisfying assignment is 111. Indeed, the only case we have to deal with here is when the predicate P has unsatisfying assignments which are 000, 001, 111. However, note that under the affine restriction where $x_1 = x_2$, this is exactly an AND, as the unsatisfying assignments will be 00, 01, 11. Hence, sparsifying this expression can require $\Omega(k^2)$ constraints.
2. Suppose one other unsatisfying assignment is at distance 2 from 000. WLOG let this other unsatisfying assignment be 011. Note that if we included an unsatisfying assignment that was at distance 1 from 011, then by variable negation and re-ordering, we would be back in the previous case. Hence, the only other case which has not yet been considered is when the other unsatisfying assignment is also at distance 2 from 011 and 000. WLOG let this other unsatisfying assignment be 110. Then, note that we can express this as the zeros to $x_1 + x_2 + 2x_3 \bmod 3$.

□

Proof. [Proof of [Claim 7.3](#)] Let P have 4 satisfying assignments. Note then on the 3-dimensional hypercube, every face must have either 0, 2, or 4 satisfying assignments (as otherwise P has an affine projection to AND). If any of the faces has all 4 satisfying assignments, then P is exactly just x_i or $\neg x_i$ (which is easy to sparsify). So, suppose every face has 2 satisfying assignments. Without loss of generality, suppose one of the satisfying assignments is 111. Then, there are two cases:

1. There is a satisfying assignment at distance 1 from 111, WLOG let this be 110. Then the other two satisfying assignments must be 000 and 001, as otherwise, there exists a face with more than 2 satisfying assignments. If this is the case, then constraint can be expressed as $x_1 + x_2 + 1 \bmod 2$.
2. There is no satisfying assignment at distance 1 from 111. This means there are only satisfying assignments at distance 2 from 111, WLOG let this be 100. Note that if 000 is also a satisfying assignment in this case, then by negating all the variables, we are back in the previous case. Hence, the only other case is when 000 is not a satisfying assignment, which means that all the satisfying assignments are 100, 001, 010, 111, which is exactly $x_1 + x_2 + x_3 \bmod 2$.

□

Proof. [Proof of [Claim 7.4](#)] We now consider the case where P has one, two or three satisfying assignments and prove that in each case it has an affine projection to AND.

1. 3 Satisfying assignments: Suppose the satisfying assignments are $a_1a_2a_3$, $b_1b_2b_3$ and $c_1c_2c_3$. Choose a coordinate such that not all the strings are equal on this coordinate. This means that by restricting this coordinate, there are either 1 or 2 satisfying assignments on the face. So, fix the coordinate to make 1 satisfying assignment, which yields an AND.
2. 2 Satisfying assignments: Suppose the satisfying assignments are $a_1a_2a_3$ and $b_1b_2b_3$. Choose a coordinate such that not all the strings are equal on this coordinate. By restricting this coordinate, we create a predicate with one satisfying assignment, and hence an AND.
3. 1 Satisfying assignment: Let the satisfying assignment be $a_1a_2a_3$. Restrict $x_1 = a_1$. This yields a predicate with 1 satisfying assignment, and hence an AND.

□

This concludes the proof of [Claim 7.1](#)-[Claim 7.4](#) and thus the proof of [Theorem 1.6](#).

8 Application to Hypergraph Cut Sparsifiers

In this section, we will show how our result implies the existence of near-linear size hypergraph cut sparsifiers. First, we introduce the definition of a hypergraph.

DEFINITION 8.1. A hypergraph $G = (V, E)$ is a set of n vertices V , along with a set of hyperedges E . Each hyperedge e is a subset of V , of any size.

Next, we introduce the definition of a *cut* in a hypergraph.

DEFINITION 8.2. For a hypergraph $G = (V, E)$, a cut in the hypergraph is a non-empty subset $S \subset V$. The size of the cut S is the number of hyperedges in E that are not completely contained in S or $V - S$. Intuitively, this is the number of edges that cross between $S, V - S$. We use $\delta_S(G)$ to denote the crossing edges corresponding to S in G .

With this, we can then state a consequence of our main result in the setting of hypergraphs.

COROLLARY 8.1. For a hypergraph $G = (V, E)$ on k vertices, there exists a weighted sub-hypergraph \hat{G} of G with $\tilde{O}(k/\varepsilon^2)$ hyperedges, such that for any subset $S \subseteq V$,

$$(1 - \varepsilon)\text{wt}(\delta_G(S)) \leq \text{wt}(\delta_{\hat{G}}(S)) \leq (1 + \varepsilon)\text{wt}(\delta_G(S)).$$

At a high level, our proof takes advantage of the fact that we showed the existence of code sparsifiers over *any* arbitrary field \mathbb{F}_q . In particular, for a hypergraph on k vertices, we will choose a prime q between k and $2k$. Then, we will create a generating matrix for a code over \mathbb{F}_q where each row of the generating matrix corresponds to a hyperedge in the hypergraph. To start, we create a generating matrix with k columns. Now, for any hyperedge by e , we denote its size by $|e|$. If we analyze the row of the generating matrix corresponding to edge e , we then place a 1 in the columns corresponding to vertices $e_1, \dots, e_{|e|-1}$. For $e_{|e|}$ (i.e. the final vertex contained in the hyperedge), we place the value $q - |e| + 1$. In doing so, the row-sum of any row of the generating matrix will be exactly 0. Indeed, for any $\{0, 1\}$ weighted linear combination of the columns of this generating matrix, a row is identically 0 if either *all* of the vertices corresponding to the hyperedge are included in the linear combination, or none of them are. This is exactly the definition of a hypergraph cut.

First, we use a basic number theoretic result:

FACT 8.1. (BERTRAND'S POSTULATE) For any positive integer n , there exists a prime between n and $2n$.

Thus, for any hypergraph on k vertices, we can find a prime number between $k, 2k$. Next, we define more specifically the generating matrix corresponding to a hypergraph:

DEFINITION 8.3. For a hypergraph $H = (V, E)$ on k vertices, let q be a prime between $k, 2k$ as guaranteed by [Fact 8.1](#). Let G be a generating matrix of a code defined over $\mathbb{F}_q^{|E|}$, and let G have k columns. Now, for any hyperedge $e = v_1, \dots, v_{|e|}$, let G_{e,v_i} be 1 if $i \leq |e| - 1$, and $q - |e| + 1$ if $i = |e|$. All other entries in the row are zero. Call this generating matrix G the generating matrix associated with H .

REMARK 8.1. Let G be the generating matrix associated with a hypergraph H . Let $S \subseteq [k]$, and let x be the indicator vector for S . Then,

$$\text{wt}(\delta_H(S)) = \text{wt}(Gx).$$

Proof. Consider any such $S \subseteq [k]$. $\text{wt}(\delta_H(S))$ is exactly the number of hyperedges that are not completely contained in S or $V - S$. Now, $\text{wt}(Gx)$ is the number of non-zero entries in Gx . By our construction of G , the only way for a $\{0, 1\}$ row-sum of G to be zero is when a $\{0, 1\}$ vector assigns either all 0's, or all 1's to the corresponding vertices of the hyperedge. This is exactly the same as the hyperedge being completely contained in S or $V - S$. \square

CLAIM 8.1. Let G be the generating matrix associated with a hypergraph H . If, there exists a sparsifier \hat{G} such that for every message $x \in \mathbb{F}_q^k$, $\text{wt}(\hat{G}x) \in (1 \pm \varepsilon)\text{wt}(Gx)$, then if we select the corresponding edges of H with the same weights as in \hat{G} , we will recover a $(1 \pm \varepsilon)$ hypergraph cut sparsifier for H .

Proof. By the previous Remark, for any set $S \subseteq [k]$, and x being the indicator vector for S ,

$$\text{wt}(\delta_H(S)) = \text{wt}(Gx).$$

Further, given \hat{G} , we can create the corresponding hypergraph \hat{H} . It is still true that

$$\text{wt}(\delta_{\hat{H}}(S)) = \text{wt}(Gx).$$

Thus, we conclude that for any $S, x = \mathbf{1}[S]$,

$$(1 - \varepsilon)\text{wt}(\delta_H(S)) = (1 - \varepsilon)\text{wt}(Gx) \leq \text{wt}(\hat{G}x) = \text{wt}(\delta_{\hat{H}}(S)) \leq (1 + \varepsilon)\text{wt}(Gx) = (1 + \varepsilon)\text{wt}(\delta_H(S)).$$

It follows that the hypergraph associated with \hat{G} is indeed a hypergraph cut-sparsifier for H . \square

COROLLARY 8.2. For any hypergraph H on k vertices, there exists a hypergraph cut sparsifier of H with $\tilde{O}(k/\varepsilon^2)$ weighted hyperedges.

Proof. Let G be the generating matrix associated with H . Let \hat{G} be the $(1 \pm \varepsilon)$ code-sparsifier for the code generated by G . Note that because $q \leq 2k$, the number of rows in \hat{G} is

$$\tilde{O}(k \log q/\varepsilon^2) = \tilde{O}(k/\varepsilon^2).$$

By the previous claim, we can then let \hat{H} be the hypergraph associated with \hat{G} . \square

Finally, by using Corollary [2.2](#), we can state a novel fact about the decomposition of hypergraphs.

COROLLARY 8.3. For any hypergraph H on n vertices, for any integer $d \geq 1$, there exists a set of at most nd hyperedges, such that upon their removal, the resulting hypergraph has at most $(2n)^{2\alpha}$ cuts of size $\leq \alpha d$.

Proof. Let G be the generating matrix associated with H . By the previous corollary, it follows that there exists a prime $q \leq 2n$, and a set of at most nd rows we can remove from G such that the number of codewords of weight $\leq \alpha d$ is at most $q^\alpha \cdot \binom{k}{\alpha}$. However, each such codeword corresponds with a possible cut in the graph of size at most αd . Hence, the number of possible cuts in the graph of size $\leq \alpha d$ is at most $(2n)^{2\alpha}$.

Note that again, if two separate vertex cuts $S_1, V - S_1$ and $S_2, V - S_2$ lead to exactly the same hyperedges being cut, we do not consider these to be separate cuts. Indeed, when we bound the number of cuts, we are bounding the number of distinct sets of hyperedges being cut of a given size. \square

9 Conclusions

In this work, we showed that for any linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$ of dimension k , there exists a weighted set S of $\tilde{O}(k \log(q)/\varepsilon^2)$ coordinates, such that for any codeword $c \in \mathcal{C}$, the quantity $\text{wt}(c|_S)$ is a $(1 \pm \varepsilon)$ approximation to $\text{wt}(c)$. This result provides a unified approach to recover known results about existence of near-linear size graph and hypergraph cut-sparsifiers, as well as some new results that include near-linear size Cayley-graph sparsifiers of Cayley graphs over \mathbb{F}_2^k , and near-linear size sparsifiers for a broader class of CSPs than were previously known.

The existential nature of our sparsification result raises the following natural question. Is there a poly-time algorithm to find the decomposition of Theorem 2.1? Alternately, is there a poly-time algorithm to compute code sparsifiers of near-linear size? The NP-hardness of the Minimum Distance Problem (MDP) for linear codes makes it particularly hard to implement Algorithm 4 (which explicitly calculates the minimum distance) and likewise finding the decomposition of Theorem 2.1 for arbitrary d can in *some* cases solve the MDP as well.

Another interesting direction for future work is to extend our classification theorem for sparsifiability of CSPs to predicates of arity greater than 3.

Acknowledgments

We thank Salil Vadhan for pointing out the connection between codes and the eigenvalues of the Laplacians of Cayley graphs over \mathbb{F}_2^k .

References

- [BK96] András A. Benczúr and David R. Karger. Approximating s - t minimum cuts in $\tilde{O}(n^2)$ time. In Gary L. Miller, editor, *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing, Philadelphia, Pennsylvania, USA, May 22-24, 1996*, pages 47–55. ACM, 1996.
- [BSS09] Joshua D. Batson, Daniel A. Spielman, and Nikhil Srivastava. Twice-ramanujan sparsifiers. In Michael Mitzenmacher, editor, *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, pages 255–262. ACM, 2009.
- [BZ20] Silvia Butti and Stanislav Zivný. Sparsification of binary csp. *SIAM J. Discret. Math.*, 34(1):825–842, 2020.
- [CKM⁺14] Michael B. Cohen, Rasmus Kyng, Gary L. Miller, Jakub W. Pachocki, Richard Peng, Anup B. Rao, and Shen Chen Xu. Solving SDD linear systems in nearly $m \log^{1/2} n$ time. In David B. Shmoys, editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 343–352. ACM, 2014.
- [CKN20] Yu Chen, Sanjeev Khanna, and Ansh Nagda. Near-linear size hypergraph cut sparsifiers. In Sandy Irani, editor, *61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020*, pages 61–72. IEEE, 2020.
- [CSWZ16] Jiecao Chen, He Sun, David P. Woodruff, and Qin Zhang. Communication-optimal distributed clustering. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 3720–3728, 2016.
- [FHHP11] Wai Shing Fung, Ramesh Hariharan, Nicholas J.A. Harvey, and Debmalya Panigrahi. A general framework for graph sparsification. In *Proceedings of the Forty-Third Annual ACM Symposium on Theory of Computing, STOC '11*, page 71–80, New York, NY, USA, 2011. Association for Computing Machinery.
- [FK17] Arnold Filtser and Robert Krauthgamer. Sparsification of two-variable valued constraint satisfaction problems. *SIAM J. Discret. Math.*, 31(2):1263–1276, 2017.
- [JLLS23] Arun Jambulapati, James R. Lee, Yang P. Liu, and Aaron Sidford. Sparsifying sums of norms. *CoRR*, abs/2305.09049, 2023.
- [JS20] Arun Jambulapati and Aaron Sidford. Ultrasparse ultrasparsifiers and faster laplacian system solvers. *CoRR*, abs/2011.08806, 2020.
- [Kar93] David R. Karger. Global min-cuts in rnc, and other ramifications of a simple min-cut algorithm. In Vijaya Ramachandran, editor, *Proceedings of the Fourth Annual ACM/SIGACT-SIAM Symposium on Discrete Algorithms, 25-27 January 1993, Austin, Texas, USA*, pages 21–30. ACM/SIAM, 1993.
- [Kar94] David R. Karger. Using randomized sparsification to approximate minimum cuts. In Daniel Dominic Sleator, editor, *Proceedings of the Fifth Annual ACM-SIAM Symposium on Discrete Algorithms. 23-25 January 1994, Arlington, Virginia, USA*, pages 424–432. ACM/SIAM, 1994.
- [Kar99] David R. Karger. Random sampling in cut, flow, and network design problems. *Math. Oper. Res.*, 24(2):383–413, 1999.

- [KK15] Dmitry Kogan and Robert Krauthgamer. Sketching cuts in graphs and hypergraphs. In Tim Roughgarden, editor, *Proceedings of the 2015 Conference on Innovations in Theoretical Computer Science, ITCS 2015, Rehovot, Israel, January 11-13, 2015*, pages 367–376. ACM, 2015.
- [KKT^Y21] Michael Kapralov, Robert Krauthgamer, Jakab Tardos, and Yuichi Yoshida. Towards tight bounds for spectral sparsification of hypergraphs. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC ’21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021*, pages 598–611. ACM, 2021.
- [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 Chandra Chekuri, editor, *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2014, Portland, Oregon, USA, January 5-7, 2014*, pages 217–226. SIAM, 2014.
- [KLP⁺16] Rasmus Kyng, Yin Tat Lee, Richard Peng, Sushant Sachdeva, and Daniel A. Spielman. Sparsified cholesky and multigrid solvers for connection laplacians. In Daniel Wichs and Yishay Mansour, editors, *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 842–850. ACM, 2016.
- [LS18] Yin Tat Lee and He Sun. Constructing linear-sized spectral sparsification in almost-linear time. *SIAM J. Comput.*, 47(6):2315–2336, 2018.
- [Pen16] Richard Peng. Approximate undirected maximum flows in $O(m \text{polylog}(n))$ time. In Robert Krauthgamer, editor, *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 1862–1867. SIAM, 2016.
- [She13] Jonah Sherman. Nearly maximum flows in nearly linear time. In *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA*, pages 263–269. IEEE Computer Society, 2013.
- [ST11] Daniel A. Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM J. Comput.*, 40(4):981–1025, 2011.

A A Simpler Construction of Cut Sparsifiers for Graphs

Let $G = (V, E)$ be a graph on n vertices, and $\varepsilon \in (0, 1)$ be some given accuracy parameter. Our goal in this section is to output a $(1 \pm \varepsilon)$ cut-sparsifier for G . The first proof showing the existence of cut-sparsifiers of near-linear size was given by [BK96], using the notion of a strength decomposition of a graph. Subsequently, a different proof based on edge connectivities was given in [FHHP11]. While the starting point for these proofs is Karger’s cut counting bound [Kar93, Kar99], they both require additional ideas and are somewhat involved. We present here a simpler proof that directly utilizes Karger’s cut counting bound to recursively sparsify G while preserving its cuts.

At a high level, the algorithm is very simple: starting with a graph G that has $d \cdot n \log(n)/\varepsilon^2$ edges, we first remove all edges involved in cuts in this graph that have size $\leq \sqrt{d} \log(n)/\varepsilon^2$, and place these edges in a graph G_1 . By a simple argument, we can show that G_1 will have at most $n\sqrt{d} \log(n)/\varepsilon^2$ edges. Now, the resulting graph from removing all the small cuts from G is a graph we call G' . G' has the special property that all non-zero cuts in G' have size at least $\sqrt{d} \log(n)/\varepsilon^2$. So, we can decompose G' into its connected components, and using Karger’s cut-counting bound for each of these connected components, we can argue that sampling each edge with probability $1/\sqrt{d}$ and weight \sqrt{d} will preserve cut-sizes with high probability. As a result, with high probability G' has $\sqrt{dn} \log(n)/\varepsilon^2$ edges, and cut-sizes are preserved. Thus, starting with a graph on $d \cdot n \log(n)/\varepsilon^2$, we get two graphs on $\sqrt{d} \cdot n \log(n)/\varepsilon^2$ edges. By repeating this procedure recursively, we can then get near-linear size cut-sparsifiers.

Note that, as stated, the algorithm only works for unweighted graphs, and indeed our proof here is intended only for unweighted graphs in the name of simplicity. However, this proof can be extended to weighted graphs by using the more advanced decomposition techniques from § 5.

In this section, the algorithm we propose will create cut sparsifiers with

$$O\left(\frac{n(\log(n) \log \log(n))^2}{\varepsilon^2}\right)$$

edges. The only non-trivial fact we will use is Theorem 3.1. Our algorithm is presented in Algorithm 10.

Note that in the algorithm, we are using $\gamma(n) = C \cdot \log n$, where C is a constant (depending on ε). We will need some claims in order to prove this result.

CLAIM A.1. *For any n -vertex graph $G(V, E)$ and a positive integer $c \in [1..(n - 1)]$, to remove all non-empty cuts*

Algorithm 10: GraphSparsify(G, ε, n, i)

```

1 if  $i = \log \log n$  or  $G$  has  $\leq n \log(n)/\varepsilon^2$  edges then
2   | return  $G$ 
3 end
4 Initialize empty graphs  $G_1, G_2$  on  $n$  vertices.
5 while  $\exists$  a non-empty cut  $(S, V - S)$  in  $G$  of size  $\leq \gamma(n) \cdot n^{1/2^i}$  do
6   | Remove the edges in  $\delta_G(S)$  from  $G$ , and place them in  $G_1$ .
7 end
8 for the remaining edges  $e \in G$  do
9   | Sample each  $e$  with probability  $\frac{1}{n^{1/2^i}}$ , and if sampled, add it to  $G_2$ .
10 end
11 return GraphSparsify( $G_1, \varepsilon, n, i + 1$ )  $\cup n^{1/2^i} \cdot$  GraphSparsify( $G_2, \varepsilon, n, i + 1$ ).

```

of size $\leq c$ in the graph, we require removing only $(n - 1)c$ edges.

Proof. Fix a $c \in [1..(n - 1)]$. Let $T(n)$ denote the maximum number of edges that are involved in cuts of size at most c in an n -vertex graph.

We prove this claim by induction on number of vertices. For our base case, consider any graph on $p \leq (c + 1)$ vertices. There are at most $\binom{p}{2} \leq \frac{(c+1)(p-1)}{2} \leq (p-1) \cdot c$ edges in this graph, and all of them are in cuts of size at most c , so $T(p) \leq (p-1)c$. Now, consider an arbitrary n -vertex graph. We will repeatedly remove from G edges in cuts of size at most c . Upon finding a cut $(S, V - S)$ of size at most c , we remove all edges involved in the cut, and recursively continue on $G[S]$ and $G[V - S]$. So assuming $|S| = s$, we get that

$$T(n) \leq c + T(s) + T(n - s),$$

where $1 \leq s \leq n - 1$. Invoking the inductive hypothesis, we get $T(n) \leq c + (s - 1)c + (n - s - 1)c = (n - 1)c$, completing the proof. \square

CLAIM A.2. At the i th level of recursion, with high probability, each non-empty G_i has at most $(1 + \frac{1}{\log \log n})^i \cdot n^{1+1/2^i} \cdot \gamma(n)$ edges, where $\gamma(n) = C \log n$.

Proof. We prove the claim by induction. Consider the first level of recursion of the algorithm. Let $n^{1/2} \cdot \gamma(n)$ be the minimum cut value on which we set our threshold. Then, the number of edges that we keep corresponding to the minimum cuts is at most $n^{3/2} \cdot \gamma(n)$. So, the first of the graphs in the subcall is of size at most $\gamma(n) \cdot n^{3/2}$.

Now, we sample the remaining edges with probability $\frac{1}{n^{1/2}}$. Because the support size is $\geq n \log(n)/\varepsilon^2$, we can use a Chernoff bound to argue that with high probability, at most $(1 + 1/\log \log n)$ the expected number of edges will be included by the sampling procedure. This means that at most

$$n^2 \cdot \frac{(1 + 1/\log \log n)}{n^{1/2}} = (1 + 1/\log \log n)n^{3/2} \leq (1 + 1/\log \log n) \cdot \gamma(n) \cdot n^{3/2}$$

edges will be included with high probability. So, the number of edges in the second graph on which we recurse is at most $(1 + 1/\log \log n)\gamma(n) \cdot n^{3/2}$.

Now, we suppose the claim holds by induction. Let us analyze the j th level of recursion. So, let the graph at this level be G , and have at most $(1 + 1/\log \log n)^j \cdot n^{1+1/2^j} \cdot \gamma(n)$ edges. At this level of recursion, the first subgraph we make has all cuts of size at most $\gamma(n)n^{1/2^{j+1}}$. By our previous claim, there are at most $n \cdot \gamma(n)n^{1/2^{j+1}}$ many edges involved in those cuts, so our first graph we recurse on will satisfy the desired bound. To construct the second graph we recurse on, we sample all remaining edges with probability $\frac{1}{n^{1/2^{j+1}}}$. By Chernoff again, we know the number of edges sampled is with very high probability at most $(1 + 1/\log \log n)$ the expected number, because the graph has at least $n \log(n)/\varepsilon^2$ edges. Hence, with high probability, the size of the second graph we recurse on has at most

$$\frac{(1 + 1/\log \log n)}{n^{1/2^{j+1}}} \cdot (1 + 1/\log \log n)^j \cdot n^{1+1/2^j} \cdot \gamma(n) = \gamma(n) \cdot (1 + 1/\log \log n)^{j+1} n^{1+1/2^{j+1}}$$

edges. \square

CLAIM A.3. Suppose a graph G is decomposed into two graphs $G = G_1 \cup G_2$, and we are given a $(1 \pm \varepsilon)$ cut-sparsifier H_1 for G_1 , H_2 for G_2 . Then, $H_1 \cup H_2$ provides a $(1 \pm \varepsilon)$ approximation to every cut in G .

Proof. Consider any cut $S, V - S$ in G . The claim follows because

$$\begin{aligned} (1 - \varepsilon)\text{wt}(\delta_G(S)) &= (1 - \varepsilon)\text{wt}(\delta_{G_1}(S)) + (1 - \varepsilon)\text{wt}(\delta_{G_2}(S)) \leq \text{wt}(\delta_{H_1}(S)) + \text{wt}(\delta_{H_2}(S)) \\ &\leq (1 + \varepsilon)\text{wt}(\delta_{G_1}(S)) + (1 + \varepsilon)\text{wt}(\delta_{G_2}(S)) = (1 + \varepsilon)\text{wt}(\delta_G(S)). \end{aligned}$$

\square

CLAIM A.4. If a cut-sparsifier G'' is a $(1 \pm \varepsilon)$ cut-approximation to G' , and G' is a $(1 \pm \delta)$ cut-approximation to G , then G'' is a $(1 - \delta)(1 - \varepsilon), (1 + \delta)(1 + \varepsilon)$ approximation to G .

Proof. Consider any cut $(S, V - S)$. We know that $(1 - \varepsilon)\text{wt}(\delta_{G''}(S)) \leq \text{wt}(\delta_{G'}(S)) \leq (1 + \varepsilon)\text{wt}(\delta_{G''}(S))$. Additionally, $(1 - \delta)\text{wt}(\delta_{G'}(S)) \leq \text{wt}(\delta_G(S)) \leq (1 + \delta)\text{wt}(\delta_{G''}(S))$. Composing these two facts, we get our claim. \square

LEMMA A.1. For any $j \in [1, \dots, \log \log n]$, the output of our algorithm when called with $i = \log \log n - j$ on a graph G is a $(1 \pm \varepsilon)^j$ cut-sparsifier for G with probability at least $1 - 4^j/n^8$.

Proof. In the base case, we consider when $j = 0$. Clearly then, we return G , which is indeed a 1-approximation.

Now, suppose the claim holds by induction. Now, let $i = \log \log(n) - j$. The algorithm, after receiving G , breaks G in G_1, G_2 , where G_1 contains all edges that are in cuts of size $\leq \gamma(n) \cdot n^{1/2^i}$, and G_2 contains all remaining edges. Note that from our previous claims, it suffices to argue that we get $(1 \pm \varepsilon)$ approximations to G_1, G_2 , as the returned sparsifier for G_1, G_2 will be $(1 \pm \varepsilon)^{j-1}$ approximations by induction.

Our algorithm completely preserves G_1 , so this is not an issue. Instead, we focus on G_2 . The algorithm samples every edge from G_2 with probability $\frac{1}{n^{1/2^i+1}}$. Because the minimum cut size in each component of G_2 is $\geq \gamma(n)n^{1/2^i}$, we know that the number of cuts in the graph of size $\leq \alpha \cdot \gamma(n) \cdot n^{1/2^i}$ is at most $n^{2\alpha}$. This is because if we denote the sizes of the components as x_1, \dots, x_r , $\sum_{i=1}^r x_i^{2\alpha} \leq n^{2\alpha}$. Now, if we preserve every cut in each component of size $\leq \alpha c$, it follows that we preserve every cut in G_2 of size $\leq \alpha c$ (since the empty cuts are preserved trivially). For a cut of size $[\alpha/2\gamma(n) \cdot n^{1/2^i}, \alpha\gamma(n) \cdot n^{1/2^i}]$, if we sample with probability $\frac{1}{n^{1/2^i+1}}$, then the probability that we do not get a $(1 \pm \varepsilon)$ approximation to the cut is at most

$$2 \cdot 2^{-0.38\varepsilon^2 \frac{\alpha}{2}\gamma(n) \cdot n^{1/2^i} \cdot n^{-1/2^{i+1}}} = 2 \cdot 2^{-0.19\varepsilon^2 \alpha\gamma(n) \cdot n^{1/2^{i+1}}}.$$

Taking the union bound over the at most $n^{2\alpha}$ cuts of this size, the probability that we fail for cuts of size between $[\alpha/2, \alpha] \cdot C \cdot n^{1/2^i}$ is at most

$$2 \cdot 2^{-0.19\alpha\gamma(n)\varepsilon^2 \cdot n^{1/2^{i+1}}} \cdot 2^{2\alpha \log n} \leq 2^{\alpha(-0.19\varepsilon^2 C \cdot n^{1/2^{i+1}} \log n + 3 \log n)}.$$

Setting $C = \frac{100}{\varepsilon^2}$, we get that with probability at most $1/n^{10}$ the sparsifier for cuts of size $[\alpha/2, \alpha] \cdot \gamma(n) \cdot n^{1/2^i}$ will fail. Now, we can take a union bound over $\alpha = 1, 2, \dots, n^2$ to conclude that all cuts in G_2 will be preserved to a $(1 \pm \varepsilon)$ fraction with probability $1 - 1/n^8$.

Now, by induction, each of the recursive calls will return cut sparsifiers with probability $1 - 4^{j-1}/n^8$. So, the total failure probability at level j is bounded by $4^{j-1}/n^8 + 4^{j-1}/n^8 + 1/n^8 \leq 4^j/n^8$, as we desire.

Continuing to $j = \log \log n$, we get our desired result. \square

LEMMA A.2. After recursing to depth $\log \log n$, the final returned graph has at most

$$O(n \log^2(n) (\log \log(n))^2 / \varepsilon^2)$$

edges with high probability.

Proof. We set $\varepsilon' = \frac{\varepsilon}{2\log\log n}$. Then, the error in approximation for our starting graph is $(1 - \frac{\varepsilon}{2\log\log n})^{\log\log n}$, $(1 + \frac{\varepsilon}{2\log\log n})^{\log\log n}$, which yields a $(1 \pm \varepsilon)$ cut-sparsifier with high probability.

At the $\log\log n$ th level of recursion, each graph is of size at most

$$(1 + 1/\log\log n)^{\log\log n} \cdot n \cdot n^{1/2^{\log\log n}} \cdot \gamma(n) = O(n \cdot \log(n)(\log\log(n))^2/\varepsilon^2),$$

where we use that $\gamma(n) = O(\frac{\log n}{\varepsilon^2})$. Now, we take the union bound over all $2^{\log\log n}$ graphs at level $\log\log n$ in the recursion to conclude that there are at most $O(n \cdot \log^2(n)(\log\log(n))^2/\varepsilon^2)$ edges in the sparsified graph with high probability. \square