

Polar Codes with Balanced Codewords

Utkarsh Gupta*, Han Mao Kiah†, Alexander Vardy‡, and Hanwen Yao‡

*Department of Mathematics, Indian Institute of Technology, Hauz Khas, New Delhi, India

†School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore

‡Department of Electrical & Computer Engineering, University of California San Diego, La Jolla, CA, USA

Emails: utkarshgupta99@gmail.com, hmkiah@ntu.edu.sg, avardy@ucsd.edu, hay125@eng.ucsd.edu

Abstract—The *imbalance* of a binary word refers to the absolute difference between the number of ones and zeros in the word. Motivated by applications in DNA-based data storage and the success of polar codes, we study the problem of reducing imbalance in the codewords of a polar code. To this end, we adapt the technique of Mazumdar, Roth, and Vontobel by considering balancing sets that correspond to low-order Reed-Muller (RM) codes. Such balancing sets are likely to be included as subcodes in polar codes.

Specifically, using the first-order RM code, we show that any message can be encoded into a length- n polar codeword with imbalance at most $o(n)$ in $O(n \log n)$ -time. We then reduce the imbalance even further using two methods. First, we constrain the ambient space \mathbb{X} and analyze the imbalance that the first-order RM code can achieve for words in \mathbb{X} . We demonstrate that for codelengths up to 128, the first-order RM code achieves zero imbalance for appropriate choices of \mathbb{X} that sacrifice only a few message bits. Second, we augment the balancing set by considering higher order RM codes. We give a simple recursive upper bound for the guaranteed imbalance of RM codes. We also prove that the second-order RM code $\text{RM}(2, m)$ balances all even-weight words for $m \leq 5$, while the RM code of order $m - 3$ balances all even-weight words for $m \geq 5$.

I. INTRODUCTION

The *imbalance* of a binary word x refers to the absolute difference between the number of ones and the number of zeros in x . A word is *balanced* if its imbalance is at most one and a code is *balanced* if all its codewords are balanced. Due to their applications in various recording systems, balanced codes have been extensively studied (see [1] for a survey).

Coupled with recent progress in the biotechnology industry, DNA macromolecules are emerging as a next-generation data storage medium with its unprecedented density, durability and replication efficiency [2]. This has rekindled interest in balanced codes. Specifically, a DNA string comprises four bases or letters: A, C, T and G, and a string is *GC-rich* (or *GC-poor*) if a high (or low) proportion of the bases corresponds to either G or C. Since GC-rich or GC-poor DNA strings are prone to both synthesis and sequencing errors [3], [4], we aim to reduce the difference with the number of G and C and the number of A and T on every DNA codeword. This requirement is equivalent to reducing the imbalance of a related binary word.

To further reduce errors, we equip our balanced codes with error-correcting capabilities and previous constructions of balanced error-correcting codes can be found in [5]–[10]. We highlight a few techniques that transform known linear error-correcting codes into balanced ones. In [8], Weber *et al.* take two input codes of distance d : a linear code and a *short* balanced codebook, and constructs a long balanced code of distance d . Later, Chee *et al.* [9] remove the need of a short codebook and construct balanced error-correcting codes from cyclic codes.

Of significance to our work is the technique of Mazumdar *et al.* [7]. At a high level, Mazumdar *et al.* consider a linear code \mathbb{C} of length n and write it as a direct sum of two linear subspaces \mathbb{C}' and \mathbb{B} , in other words, $\mathbb{C} = \mathbb{C}' \oplus \mathbb{B}$. Here, \mathbb{B} is chosen so that it is a *balancing set*. In other words, for any binary word x , there exists a balancing word $b \in \mathbb{B}$ so that $x + b$ is balanced. In this

coding scheme, the codewords are of the form $x + b$, where x belongs to \mathbb{C}' and b is its corresponding balancing word. In the same paper, the authors demonstrated that a random subcode \mathbb{B} of dimension $(3/2) \log n + o(\log n)$ is a balancing set with high probability. However, verifying that \mathbb{B} is indeed a balancing set remains computationally difficult.

In this paper, we apply the technique of Mazumdar *et al.* [7] to a beautiful and important class of codes, the *polar codes*. Invented by Arıkan [11], polar codes achieve capacity for many channels with low encoding and decoding complexities. Recently, they are being adopted in 5G standard [12] and have been adapted for insertion and deletion channels [13], [14]. Given this appeal, we study efficient means of transforming messages into balanced polar codewords while retaining the low complexities of the polar encoding and decoding algorithms.

We remark that techniques to adapt polar codes for processes with memory or constrained systems were studied by Şaşoğlu [15] and Shuval and Tal [16]. However, it is unclear whether the framework in these works can be used efficiently here.

II. PRELIMINARIES

Let x be a binary word of length n . We use $\text{wt}(x)$ and $\mu(x)$ to denote its weight and imbalance, respectively. We have that

$$\mu(x) = |2\text{wt}(x) - n|. \quad (1)$$

We also regard x as a vector in \mathbb{F}_2^n . For a subset $\mathcal{X} \subset \mathbb{F}_2^n$, we use $\text{span}(\mathcal{X})$ to denote the *linear span of vectors in \mathcal{X}* . Given two vectors x and y , we let $\text{wt}(x \cap y)$ denote the number of indices where both x and y are one, and we have the observation:

$$\text{wt}(x + y) = \text{wt}(x) + \text{wt}(y) - 2\text{wt}(x \cap y). \quad (2)$$

We use $\langle x, y \rangle$ to denote the *dot product* of x and y over \mathbb{F}_2 and we have that $\text{wt}(x \cap y) \equiv \langle x, y \rangle \pmod{2}$.

A. Polar Codes

The *polarization kernel matrix* G_1 is given by $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$. For $m \geq 2$, set $n \triangleq 2^m$. Then the *polar transformation matrix* G_m is an $n \times n$ matrix, recursively defined by

$$G_m \triangleq \begin{bmatrix} G_{m-1} & \mathbf{0} \\ G_{m-1} & G_{m-1} \end{bmatrix}.$$

We index the rows of G_m with words in $\mathcal{Z} \triangleq \mathbb{F}_2^m$ and list them in bit-reversal order. To define a polar code, we pick a set $\mathcal{F} \subset \mathcal{Z}$ of *frozen indices* and a *frozen syndrome* $f \in \mathbb{F}_2^{|\mathcal{F}|}$. Then the *polar code* defined by \mathcal{F} and f is the linear code $\mathbb{C}(\mathcal{F}, f) \triangleq \{xG_m : x|_{\mathcal{F}} = f\}$. The index set $\mathcal{I} \triangleq \mathcal{Z} \setminus \mathcal{F}$ corresponds to the *information indices* and $k \triangleq |\mathcal{I}|$ measures the bits of information that the polar code $\mathbb{C}(\mathcal{F}, f)$ can transmit.

Specifically, given a message $m \in \mathbb{F}_2^k$, we can efficiently encode m to a polar codeword $c \triangleq \text{enc}(m)$ in $O(n \log n)$ time. On the other hand, suppose that the noisy word \tilde{c} is received. Using successive cancellation decoding, we compute $\text{dec}(\tilde{c})$ in $O(n \log n)$ time and recover the message m with high probability.

Now, the choice of \mathcal{F} is channel dependent and typically, the frozen syndrome f is chosen to be the all zero word. When $f = 0$, the polar code $\mathbb{C}(\mathcal{F}, 0)$ is simply the linear span of the rows with indices in \mathcal{I} . In this case, we simply write this linear span as $\text{span}(\mathcal{I})$. For certain choices of frozen indices, the polar code $\mathbb{C}(\mathcal{F}, 0)$ corresponds to the ubiquitous class of *Reed-Muller* (RM) codes [17, Ch. 13]. Specifically, if \mathcal{F}_r is the set of indices whose weight are at most r , then $\mathbb{C}(\mathcal{F}_r, 0)$ defines the class of Reed-Muller code of order $m - r - 1$, or $\text{RM}(m - r - 1, m)$ in short. Equivalently, if \mathcal{J}_{m-r} denote the set of indices whose weight is at least $m - r$, we have that $\text{span}(\mathcal{J}_{m-r}) = \text{RM}(r, m)$. RM codes have been studied extensively and recently, Abbe and Ye demonstrated that the RM codes share similar polarization behaviour as polar codes [18].

In the next subsection, we pick a subset \mathcal{B} from the information set \mathcal{I} for the purposes of reducing the imbalance of codewords. Specifically, we reduce the number of information bits to $k' = n - |\mathcal{F}| - |\mathcal{B}|$ so that a k' -bit message \mathbf{m}' is encoded to a polar codeword \mathbf{c} with small imbalance $\mu(\mathbf{c})$.

B. Linear Balancing Sets

We discuss the technique *ala* Mazumdar *et al.* [7] in the context of polar codes. As before, we fix a set of frozen indices \mathcal{F} and a frozen syndrome f . Let $k = n - |\mathcal{F}|$ and we have the polar code $\mathbb{C}(\mathcal{F}, f)$ with a corresponding pair of encoding $\text{enc} : \mathbb{F}_2^k \rightarrow \mathbb{F}_2^k$ and decoding $\text{dec} : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^k$ algorithms.

We pick a set of *balancing indices* \mathcal{B} so that $\mathcal{F} \cap \mathcal{B} = \emptyset$ and let \mathbb{B} be the linear span of the rows corresponding to \mathcal{B} . We then set $k' \triangleq k - |\mathcal{B}|$ and in this *balancing encoding scheme*, we transmit k' -bit messages (instead of k -bit). Consider $\mathbf{m}' \in \mathbb{F}_2^{k'}$. We first insert $|\mathcal{B}| = k - k'$ zeros to \mathbf{m}' at positions corresponding to \mathcal{B} and compute the corresponding encoding \mathbf{c}' . Next, we find a balancing vector $\mathbf{b} \in \mathbb{B}$ so that the corresponding imbalance $\mu(\mathbf{c}' + \mathbf{b})$ is minimized. Then we transmit the word $\mathbf{c} \triangleq \mathbf{c}' + \mathbf{b}$.

To decode a noisy word $\tilde{\mathbf{c}}$, we simply apply the polar-decoding algorithm and find the k -bit vector $\mathbf{m} = \text{dec}(\tilde{\mathbf{c}})$. The desired message \mathbf{m}' is then the k' -prefix of \mathbf{m} . So, if \mathbf{m} is successfully decoded under the original polar coding scheme, we also successfully recover our message \mathbf{m}' . Therefore, in this encoding / decoding scheme, provided that we can find \mathbf{b} efficiently, we retain the low encoding and decoding complexities and more importantly, the performance guarantees of polar codes.

It remains to understand how small the imbalance $\mu(\mathbf{c}' + \mathbf{b})$ is. To do so, we define the *guaranteed imbalance* as the quantity

$$\mu(\mathbb{X}, \mathbb{B}) \triangleq \max_{\mathbf{c}' \in \mathbb{X}} \left(\min_{\mathbf{b} \in \mathbb{B}} \mu(\mathbf{c}' + \mathbf{b}) \right),$$

where \mathbb{X} is some subset \mathbb{F}_2^n . When $\mathbb{X} = \mathbb{F}_2^n$, we simply write $\mu(\mathbb{B})$, instead of $\mu(\mathbb{F}_2^n, \mathbb{B})$.

In this paper, we are interested in imbalance guarantees, or equivalently, determining upper bounds for $\mu(\mathbb{X}, \mathbb{B})$ for various choices of \mathbb{X} and \mathbb{B} . Before stating our results, we discuss differences with the original scheme of Mazumdar *et al.* [7].

(a) \mathbb{B} need not be a balancing set in strictest sense. As mentioned in the introduction, to reduce errors in DNA data storage, it suffices to have DNA strings with small GC-imbalance. Therefore, instead of requiring $\mu(\mathbb{B}) = 0$, we focus on finding a set \mathbb{B} so that $\mu(\mathbb{B}) = o(n)$. We loosely refer to \mathbb{B} as a balancing set.

(b) \mathbb{X} need not be the whole space \mathbb{F}_2^n . In the balancing encoding scheme, we are only required to reduce the imbalance of the

m	n	$\log \mathbb{X}_1 $	$\mu(\mathbb{X}_1, \mathbb{B}_m)$	$\log \mathbb{X}_2 $	$\mu(\mathbb{X}_2, \mathbb{B}_m)$	$\log \mathbb{X}_3 $	$\mu(\mathbb{X}_3, \mathbb{B}_m)$
2	4	3	0	—	—	—	—
3	8	7	0	—	—	—	—
4	16	15	4	14	0	—	—
5	32	31	≤ 4	30	0	—	—
6	64	63	8	62	≤ 4	≈ 50	0
7	128	127	≤ 8	126	≤ 8	≈ 112	0
8	256	255	16	254	≤ 12	≈ 238	≤ 8
9	512	511	≤ 20	510	≤ 20	≈ 492	≤ 16
10	1024	1023	32	1022	≤ 28	≈ 1002	≤ 24

TABLE I: Guaranteed imbalance $\mu(\mathbb{X}, \mathbb{B}_m)$ for various choices of \mathbb{X} .

words in the polar code $\mathbb{C}(\mathcal{F}, f)$. In other words, it suffices to determine $\mu(\mathbb{C}(\mathcal{F}, f), \mathbb{B})$. However, as the choice of \mathcal{F} and f is channel dependent, it is tedious to determine the guaranteed imbalance for all possibilities. Instead, we determine $\mu(\mathbb{X}, \mathbb{B})$ for a large space \mathbb{X} that is likely to contain $\mathbb{C}(\mathcal{F}, f)$. Then the guaranteed imbalance $\mu(\mathbb{C}(\mathcal{F}, f), \mathbb{B})$ is simply upper bounded by $\mu(\mathbb{X}, \mathbb{B})$. Furthermore, we consider spaces \mathbb{X} that are easily described by the polar transformation matrix. Specifically, for the spaces \mathbb{X} considered in this paper, either we have $\mathbb{X} = \mathbb{C}(\mathcal{F}^*, f^*)$ for some choice of frozen indices or we provide an efficient mapping of messages into $\mathbb{X} \subseteq \mathbb{C}(\mathcal{F}^*, f^*)$ using the polar transformation matrix \mathbf{G}_m . These descriptions allow one to incorporate the balancing technique into the polar encoding. In all cases, the indices in set \mathcal{F}^* have weight at most two and thus, are very likely to belong to an arbitrary frozen set \mathcal{F} .

(c) \mathbb{B} corresponds to the linear space spanned by rows of \mathbf{G}_m . This requirement is specific to polar codes and allows one to simply read off the original message \mathbf{m}' from the vector obtained from polar decoding. Also, as with the point (b), we study balancing sets that are likely to be *not frozen*.

C. Our Contributions

In this paper, we focus on a family of balancing sets. For $m \geq 2$ and $1 \leq r \leq m - 1$, we let $\mathcal{B}_{r,m}$ be the set of indices with weight between $m - r$ and $m - 1$ (inclusive), and set $\mathbb{B}_{r,m} = \text{span}(\mathcal{B}_{r,m})$. Since $\mathcal{J}_{m-r} = \mathcal{B}_{r,m} \cup 1^m$, we have that $\text{RM}(r, m) = \mathbb{B}_{r,m} \oplus \text{span}(1^n)$. In other words, $\mathbb{B}_{r,m}$ is closely related to the Reed-Muller codes and we make use of certain combinatorial properties of the latter to derive estimates on guaranteed imbalance. Furthermore, when r is fixed, it can be shown that the rows or channels corresponding to $\mathcal{B}_{r,m}$ are “good” asymptotically; thus they are likely to be part of the information indices. We omit the details for lack of space.

In Section III, we study the special case where $r = 1$. Here, we simply write \mathcal{B}_m and \mathbb{B}_m , instead of $\mathcal{B}_{1,m}$ and $\mathbb{B}_{1,m}$, respectively. Our first result states that $\mu(\mathbb{B}_m) \leq \sqrt{n} = o(n)$. Since $|\mathcal{B}_m| = m = \log n$, this means that we can significantly reduce the imbalance of *all codewords* by sacrificing only $\log n$ information bits. Also, we show that we can find the corresponding codeword with the smallest imbalance in $O(n \log n)$ time.

Even though the proportion of imbalance vanishes to zero with codelength, the imbalance is relatively high for small code lengths. Hence, in Section IV, we focus on attention of certain constrained spaces \mathbb{X} and analyse the corresponding guaranteed imbalances. Specifically, we show that for appropriate choices of \mathbb{X} , we can achieve *zero imbalance* for $m \leq 7$ or $n \leq 128$. We summarize our results in Table I.

In Section V, we reduce the guaranteed imbalance by consider $\mathcal{B}_{r,m}$ for $r > 1$. We derive a simple recursive upper bound for $\mu(\mathbb{B}_{r,m})$ and demonstrate the $\mathbb{B}_{m-3,m}$ is able to balance all even-weight codewords for $m \geq 5$.

Example 1. Let $m = 4$. So, $n = 16$ and the polar transformation matrix G_4 is given by

$$\begin{array}{l} 0000 \\ 1000 \\ 0100 \\ \textcolor{magenta}{1100} \\ 0010 \\ \textcolor{magenta}{1010} \\ 0110 \\ \textcolor{magenta}{1110} \\ 0001 \\ 1001 \\ \textcolor{magenta}{0101} \\ 1101 \\ \textcolor{magenta}{0011} \\ 1011 \\ \textcolor{magenta}{0111} \\ 1111 \end{array} \left[\begin{array}{cccccccccccccccc} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{array} \right]$$

The row indices in \mathcal{B}_4 are highlighted in blue. In Section III, we show that the guaranteed imbalance $\mu(\mathcal{B}_4)$ is four, i.e. for any word $\mathbf{x} \in \mathbb{F}_2^n$, we have $\mu(\mathbf{x} + \mathbf{b}) \leq 4$ for some \mathbf{b} in \mathcal{B}_4 . For example, if $\mathbf{x} = 0001000100011110$, we can check that $\mu(\mathbf{x} + \mathbf{b}) = 4$ for all \mathbf{b} in \mathcal{B}_4 .

On the other hand, we consider the space \mathbb{X}_2 whose words are of the form $\mathbf{x}_1\mathbf{x}_2$ with both weights $\text{wt}(\mathbf{x}_1)$ and $\text{wt}(\mathbf{x}_2)$ odd. In other words, $\mathbb{X}_2 = \mathbb{C}(\{0000, 0001\}, (0, 1))$. Now, $\mathbf{x}' \triangleq 0^710^71$ belongs to \mathbb{X}_2 and it is not difficult to see that $\mu(\mathbf{x}' + 1^80^8) = 0$. Furthermore, in Section IV, for all $\mathbf{x} \in \mathbb{X}_2$, we show that it is always possible to find $\mathbf{b} \in \mathcal{B}_m$ so that $\mu(\mathbf{x} + \mathbf{b}) = 0$.

Finally, the row indices in $\mathcal{B}_{2,4}$ (which include \mathcal{B}_4) are highlighted in either blue or magenta. In Section V, we show that $\mathcal{B}_{2,4}$ balances all even-weight words. That is, if \mathbf{x} has even weight, then we can find $\mathbf{b} \in \mathcal{B}_{2,4}$ such that $\mu(\mathbf{x} + \mathbf{b}) = 0$. For instance, when $\mathbf{x} = 0001000100011110$ as above, we can choose $\mathbf{b}' = 1^40^{12} \in \mathcal{B}_{2,4}$ and check that $\mu(\mathbf{x} + \mathbf{b}') = 0$.

III. $\mu(\mathcal{B}_m)$ AND THE WALSH SPECTRUM

For $m \geq 2$, set $n = 2^m$. Recall that \mathcal{B}_m is the set of indices with weight exactly $m-1$ and $\mathcal{B}_m = \text{span}(\mathcal{B}_m)$. In this section, we derive an upper bound for $\mu(\mathcal{B}_m)$.

As mentioned earlier, $\text{span}(\mathcal{B}_m \cup \{1^n\})$ is equal to the first order RM codes. Hence, since the covering radius of the latter is at most $(n - \sqrt{n})/2$ [19], it can be argued that $\mu(\mathcal{B}_m)$ is at most \sqrt{n} . There are many derivations of the upper bound for the covering radius of $\text{RM}(1, m)$ and we briefly outline one proof that relies on the *Hadamard transform* and *Walsh spectrum*. The proof illustrates certain combinatorial properties that we exploit to reduce the imbalance in later sections.

Definition 2. For $m \geq 2$, set $n = 2^m$. The *Hadamard matrix of order m* is an integer-valued symmetric $n \times n$ matrix defined by the following recursion. Set $\mathbf{H}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ and for $m \geq 3$, $\mathbf{H}_m \triangleq \begin{bmatrix} \mathbf{H}_{m-1} & \mathbf{H}_{m-1} \\ \mathbf{H}_{m-1} & -\mathbf{H}_{m-1} \end{bmatrix}$. For $\mathbf{x} = x_1x_2 \dots x_n \in \mathbb{F}_2^n$, we define the *Walsh spectrum of \mathbf{x}* , denoted by $\mathbf{W}(\mathbf{x})$, to be the integer vector $\mathbf{H}_m \mathbf{u}$, where $\mathbf{u} = ((-1)^{x_i})_{i=1}^n$.

From the definition of Hadamard matrices, we observe that

$$\mathbf{H}_m \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \end{bmatrix} = \begin{bmatrix} \mathbf{H}_m \mathbf{u}_1^T + \mathbf{H}_m \mathbf{u}_2^T \\ \mathbf{H}_m \mathbf{u}_1^T - \mathbf{H}_m \mathbf{u}_2^T \end{bmatrix} \quad (3)$$

Using (3) and standard divide-and-conquer techniques, we can compute the Walsh spectrum efficiently.

Theorem 3 ([17, Ch. 14]). *Let $\mathbf{x} \in \mathbb{F}_2^n$. Then the Walsh spectrum of \mathbf{x} can be computed in $O(n \log n)$ time.*

The next theorem provides the connection between Walsh spectrum and imbalances.

Theorem 4 ([17, Ch. 14]). *Let $\mathbf{x} \in \mathbb{F}_2^n$ and write the Walsh spectrum of \mathbf{x} as $\mathbf{W}(\mathbf{x}) = (W_1, W_2, \dots, W_n)$. Then the n absolute values in the Walsh spectrum $|W_i|$, $i = 1, 2, \dots, n$, correspond to the n imbalances $\mu(\mathbf{x} + \mathbf{b})$, $\mathbf{b} \in \mathcal{B}_m$.*

Therefore, finding the smallest imbalance is equivalent to determining the smallest absolute value in the Walsh spectrum of \mathbf{x} . To provide an upper bound for the smallest value, we consider the square of the norm of $\mathbf{W}(\mathbf{x})$, in other words, $\|\mathbf{W}(\mathbf{x})\|^2 \triangleq \sum_{i=1}^n W_i^2$. Since $\mathbf{H}_m^2 = n\mathbf{I}_n$, we have that

$$\|\mathbf{W}(\mathbf{x})\|^2 = \mathbf{W}(\mathbf{x})^T \mathbf{W}(\mathbf{x}) = \mathbf{u}^T \mathbf{H}_m^T \mathbf{H}_m \mathbf{u} = n \mathbf{u}^T \mathbf{u} = n^2. \quad (4)$$

Therefore, the average squared value in the Walsh spectrum is n and so, the smallest absolute value is at most \sqrt{n} . When m is even, it can be shown that there exist words \mathbf{x} with $\mathbf{W}(\mathbf{x}) = (\sqrt{n})_{i=1}^n$. Thus the upper bound is essentially sharp.

We summarize our discussion in the following theorem.

Theorem 5. $\mu(\mathcal{B}_m) \leq \lfloor \sqrt{n} \rfloor = \lfloor 2^{m/2} \rfloor$. Furthermore, we have equality when m is even.

Remark 6. Here, we point out that the reduction in imbalance is nontrivial. Assuming a symmetric channel, if a polar code is properly designed, then each bit in a random codeword is either 0 or 1 with probability 1/2. Thus, Hoeffding's inequality guarantees that the probability that the imbalance of a random codeword exceeding tn is at most $2e^{-2nt^2}$. Now, taking $t = 1/\sqrt{n}$, we have the probability that a random codeword has imbalance at least \sqrt{n} is about about $2e^{-2} \approx 0.27$. In contrast, the method in this section ensures all transmitted codewords have imbalance at most \sqrt{n} and this comes with a small rate penalty of $\log n$ bits.

IV. REDUCING THE IMBALANCE FOR CERTAIN SPACES

In this section, we demonstrate that if we constrain the space \mathbb{X} , we are able to reduce the guaranteed imbalance. Our first result makes use of the following lemma which is implied by (1).

Lemma 7. *Let $n \equiv 0 \pmod{4}$. Then*

$$\mu(\mathbf{x}) \equiv \begin{cases} 0 & \text{when } \text{wt}(\mathbf{x}) \text{ is even,} \\ 2 & \text{when } \text{wt}(\mathbf{x}) \text{ is odd.} \end{cases}$$

Now, for $m \geq 2$, n is always divisible by four. Hence, if \mathbf{x} has even weight, it is necessary that its imbalance is divisible by four and the following result is straightforward.

Theorem 8. *Set $\mathbb{X}_1 \triangleq \mathbb{C}(\{\mathbf{0}\}, (0))$. Then $\mu(\mathbb{X}_1, \mathcal{B}_m) \leq 4 \lfloor \sqrt{n}/4 \rfloor$ for $m \geq 2$.*

Next, we consider a smaller space $\mathbb{X}_2 \subset \mathbb{X}_1$ and show that the guaranteed imbalance can be reduced further.

Theorem 9. *Set $\mathbb{X}_2 \triangleq \mathbb{C}(\{\mathbf{0}, 0^{m-1}1\}, (0, 1))$. If $\delta(n) = \lfloor (\sqrt{n} - 2)/4 \rfloor$, then $\mu(\mathbb{X}_2, \mathcal{B}_m) \leq 4\delta(n)$ for $m \geq 3$.*

Proof. Let $\mathbf{x} = x_1x_2 \dots x_n \in \mathbb{X}_2$. Set $\mathbf{u}_1 = ((-1)^{x_i})_{i=1}^{n/2}$, $\mathbf{u}_2 = ((-1)^{x_i})_{i=n/2+1}^n$ and $\mathbf{W}^{(j)} = \mathbf{H}_{m-1} \mathbf{u}_j$ for $j \in \{1, 2\}$. It follows from the definition of \mathbb{X}_2 that both $x_1x_2 \dots x_{n/2}$ and $x_{n/2+1}x_{n/2+2} \dots x_n$ have odd weight. Therefore, Lemma 7 states that both Walsh spectra $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ have entries congruent to two modulo four.

Now, let $j \in \{1, 2\}$. We claim that strictly more than half the entries in $\mathbf{W}^{(j)}$ have absolute values at most $4\delta(n) + 2$. We

prove by contradiction. Suppose otherwise that half of the entries in $W^{(j)}$ have absolute values strictly greater than $4\delta(n) + 2$. Then

$$\|W^{(j)}\|^2 \geq \frac{n}{4}(4\delta(n) + 6)^2 > \frac{n}{4} \left(4 \left(\frac{\sqrt{n} - 2}{4} - 1 \right) + 6 \right)^2 = \frac{n^2}{4}.$$

However, (4) states that $\|W^{(j)}\|^2 = (n/2)^2 = n^2/4$ and hence, we obtain a contradiction.

As such, we have a common index i where the absolute values of both $W_i^{(1)}$ and $W_i^{(2)}$ are at most $4\delta(n) + 2$. Since $W(x) = (W^{(1)} + W^{(2)}, W^{(1)} - W^{(2)})$, either $W_i^{(1)} + W_i^{(2)}$ or $W_i^{(1)} - W_i^{(2)}$ has absolute value at most $4\delta(n) + 2 - 2 = 4\delta(n)$. \square

Theorem 9 guarantees that $\mu(\mathbb{X}_2, \mathbb{B}_m) = 0$ for $m \in \{4, 5\}$. For larger values of m , we consider another choice of constrained space to reduce the guaranteed imbalance. For $m \geq 5$, we consider the extended Hamming code $\mathbb{H}_m \triangleq \text{RM}(m-2, m)$. The constrained space \mathbb{X}_3 is defined as follows:

$$\mathbb{X}_3 \triangleq \{x_1 x_2 : x_j \in \mathbb{H}_{m-1}, \text{wt}(x_j) \equiv 2 \pmod{4} \text{ for } j \in \{1, 2\}\}. \quad (5)$$

To analyse $\mu(\mathbb{X}_3, \mathbb{B}_m)$, we consider the dual relation of \mathbb{H}_m with \mathbb{B}_m and consider the imbalances modulo eight.

Lemma 10. *Let $m \geq 5$. If $x \in \mathbb{H}_m$ and $b \in \mathbb{B}_m$, then*

$$\mu(x + b) \equiv \begin{cases} 4 \pmod{8}, & \text{when } \text{wt}(x) \equiv 2 \pmod{4}, \\ 0 \pmod{8}, & \text{when } \text{wt}(x) \equiv 0 \pmod{4}. \end{cases}$$

Proof. From (2), we have that $\text{wt}(x + b) = \text{wt}(x) + \text{wt}(b) - 2\text{wt}(x \cap b)$. Since $b \in \mathbb{B}_m \subset \text{RM}(1, m)$ and $\text{RM}(1, m)$ is the dual code for \mathbb{H}_m , we have that $\langle x, b \rangle = 0$ and so, $2\text{wt}(x \cap b) \equiv 0 \pmod{4}$. For $m \geq 5$, since $\text{wt}(b) \equiv 0 \pmod{4}$, we have $\text{wt}(x + b) \equiv \text{wt}(x) \pmod{4}$. The lemma follows from (1). \square

Therefore, for $x_1 x_2 \in \mathbb{X}_3$, the values in the Walsh spectra of both x_1 and x_2 are congruent to four modulo eight. Proceeding in a similar manner as the proof of Theorem 9, we have the following result on the guaranteed imbalance.

Theorem 11. *Define \mathbb{X}_3 as in (5). If $\epsilon(n) = \lfloor(\sqrt{n} - 4)/8\rfloor$, then $\mu(\mathbb{X}_3, \mathbb{B}_m) \leq 8\epsilon(n)$ for $m \geq 6$.*

It remains to estimate the size of $|\mathbb{X}_3|$.

Proposition 12. $|\mathbb{H}_m| \geq 2^{n-m-2} - 2^{n/2-1}$. Therefore,

$$|\mathbb{X}_3| \geq \left(2^{n/2-m-1} - 2^{n/4-1}\right)^2 \quad (6)$$

The proof of Proposition 12 is technical and relies on certain known properties of RM codes. First, we consider the span of certain rows in the polar matrix G_m . For $m \geq 2$, define

$$\mathcal{D}_m \triangleq \begin{cases} \mathcal{J}_{(m+1)/2}, & \text{if } m \text{ is odd,} \\ \mathcal{J}_{m/2+1} \cup \{z'1 : \text{wt}(z') = m/2 - 1\}, & \text{if } m \text{ is even.} \end{cases} \quad (7)$$

It is straightforward to verify that $|\mathcal{D}_m| = 2^{m-1} = n/2$. Let $\mathbb{D}_m \triangleq \text{span}(\mathcal{D}_m)$. The next lemma states that \mathbb{D}_m is self-dual.

Lemma 13. \mathbb{D}_m is a self-dual code.

Proof. When m is odd, \mathbb{D}_m is $\text{RM}((m-1)/2, m)$ and the latter is known to be self-dual. When m is even, \mathbb{D}_m is the direct sum $\text{RM}(m/2-1, m) \oplus \mathbb{D}'$, where $\mathbb{D}' \subset \text{RM}(m/2, m)$. It is also known that the dual of $\text{RM}(m/2-1, m)$ is $\text{RM}(m/2, m)$ [17, Ch. 13].

Now, to show that \mathbb{D}_m is self-dual, we need to demonstrate that $\langle x, x' \rangle = 0$ for any vectors x, x' in \mathbb{D}_m . We write $x = u + v$ and $x' = u' + v'$ with $u, u' \in \text{RM}(m/2-1, m)$ and $v, v' \in \mathbb{D}'$. \square

Then $\langle x, x' \rangle = \langle u, u' \rangle + \langle u, v' \rangle + \langle v, u' \rangle + \langle v, v' \rangle = \langle v, v' \rangle$. Note that $\langle u, u' \rangle = \langle u, v' \rangle = \langle v, u' \rangle = 0$ because of the duality of $\text{RM}(m/2-1, m)$ and $\text{RM}(m/2, m)$. Next, since $v = yy$ and $v' = y'y'$, we have that $\langle v, v' \rangle = 2\langle y, y' \rangle = 0$. \square

Proof of Proposition 12. Consider the set of binary words:

$$\mathcal{M} \triangleq \left\{ \mathbf{m}_1 \mathbf{m}_2 : \mathbf{m}_1 \in \mathbb{F}_2^{n/2-m-1} \setminus \{\mathbf{0}\}, \mathbf{m}_2 \in \mathbb{F}_2^{n/2-1} \right\}.$$

For our proof, we present an injective map / encoding $\phi : \mathcal{M} \rightarrow \mathbb{H}_m$ such that $\text{wt}(\phi(\mathbf{m}_1 \mathbf{m}_2)) \equiv 2 \pmod{4}$. This then implies the lower bound as $|\mathcal{M}| = 2^{n/2-1}(2^{n/2-m-1} - 1) = 2^{n-m-2} - 2^{n/2-1}$. To compute $\phi(\mathbf{m}_1 \mathbf{m}_2)$, we proceed as follows. Recall \mathcal{F}_1 is the set of indices with weight ≤ 1 and we set $\mathcal{I} = \mathcal{Z} \setminus (\mathcal{F}_1 \cup \mathcal{D}_m)$. Hence, $|\mathcal{F}_1| = m + 1$, $|\mathcal{I}| = n/2 - m - 1$, $|\mathcal{D}_m| = n/2$.

- Let x_1 be such that $x_1|_{\mathcal{I}} = \mathbf{m}_1$ and $x_1|_{\mathcal{F}_1 \cup \mathcal{D}_m} = \mathbf{0}$. Set $c_1 \triangleq x_1 G_m$.
- Since $x_1|_{\mathcal{I}} \neq \mathbf{0}$, we have that $c_1 \notin \mathbb{D}_m$. Now, as \mathbb{D}_m is self-dual, there exists a row v with index in \mathcal{D}_m such that $\langle c_1, v \rangle = 1$. Set j^* be the smallest such index.
- Let x_2 be such that $x_2|_{\mathcal{D}_m \setminus \{j^*\}} = \mathbf{m}_2$ and $x_2|_{\mathcal{F}_1 \cup \mathcal{I} \cup \{j^*\}} = \mathbf{0}$. Set $c_2 \triangleq x_2 G_m$.
- If $\text{wt}(c_1 + c_2) \equiv 2 \pmod{4}$, we return $c_1 + c_2$. Note that $c_1 + c_2$ belongs to \mathbb{H}_m as $x_1|_{\mathcal{F}_1} = x_2|_{\mathcal{F}_1} = \mathbf{0}$.
- If $\text{wt}(c_1 + c_2) \equiv 0 \pmod{4}$, we return $c_1 + c_2 + v$. Since $v \in \mathbb{D}_m \subseteq \mathbb{H}_m$, the word $c_1 + c_2 + v$ belongs to \mathbb{H}_m too. To check its weight, recall that $\text{wt}(c_1 + c_2 + v) = \text{wt}(c_1 + c_2) + \text{wt}(v) - 2\text{wt}((c_1 + c_2) \cap v)$. Now, $\text{wt}(c_1 + c_2) \equiv \text{wt}(v) \equiv 0 \pmod{4}$. Also, $\text{wt}((c_1 + c_2) \cap v) \equiv \langle c_1 + c_2, v \rangle = \langle c_1, v \rangle + \langle c_2, v \rangle \equiv 1 \pmod{2}$. Therefore, $\text{wt}(c_1 + c_2 + v) \equiv 2 \pmod{4}$.

To show injectivity, we consider $\phi(\mathbf{m}_1 \mathbf{m}_2) = c$ and demonstrate how to recover \mathbf{m}_1 and \mathbf{m}_2 from c . First, compute $\mathbf{x} = c G_m$. Since $G_m G_m = I_n$, we have that $\mathbf{x} = c_1 + c_2$ or $\mathbf{x} = c_1 + c_2 + v$. In either case, $\mathbf{x}|_{\mathcal{I}} = \mathbf{m}_1$ and we can then figure out j^* . Subsequently, we have that $\mathbf{x}|_{\mathcal{D}_m \setminus \{j^*\}} = \mathbf{m}_2$. \square

The lower bound (6) can be improved as long as we have a code \mathbb{D}^* satisfying certain properties. Specifically, we have the following proposition and its proof follows from the proof of Proposition 12 by replacing \mathbb{D}_m with \mathbb{D}^* . However, while $|\mathbb{D}^*| \leq 2^{n/2}$, it is unclear how to find \mathbb{D}^* efficiently.

Proposition 14. *Suppose that \mathbb{D}^* is a maximal linear subcode of \mathbb{H}_m whose words have weight divisible by four. That is, if $c \in \mathbb{H}_m \setminus \mathbb{D}^*$, then there is a word in $\text{span}(\mathbb{D}^* \cup \{c\})$ whose weight is not divisible by four. Then the number of words in \mathbb{H}_m is at least $2^{n-m-2} - |\mathbb{D}^*|/2$.*

V. REDUCING THE IMBALANCE USING $\mathbb{B}_{r,m}$ FOR $r > 1$

In this section, instead of constraining the space \mathbb{X} , we reduce the guaranteed imbalance by augmenting the balancing set. Specifically, we provide upper bounds for the guaranteed imbalance for $\mathbb{B}_{r,m}$ when $r > 1$. Now, since $\text{RM}(r, m) = \mathbb{B}_{r,m} \oplus \text{span}(1^n)$ and $\mu(\mathbf{x}) = \mu(\mathbf{x} + 1^n)$, we have that $\mu(\mathbb{B}_{r,m}) = \mu(\text{RM}(r, m))$. Hence, it suffices to compute the guaranteed imbalance for RM codes and this is convenient as RM codes have many well-known properties.

First, we have the following recursive upper bound.

Proposition 15. $\mu(\mathbb{B}_{r+1, m+1}) \leq \mu(\mathbb{B}_{r,m})$.

Proof. We show that $\mu(\text{RM}(r+1, m+1)) \leq \mu(\text{RM}(r, m))$. Let \mathbf{x} be a binary word of length n whose imbalance is to be reduced.

In other words, if we write $\mu(\text{RM}(r, m)) = \mu_0$, our task is to find $\mathbf{b} \in \text{RM}(r+1, m+1)$ such that $\mu(\mathbf{x} + \mathbf{b}) \leq \mu_0$.

If $\text{wt}(\mathbf{x}) = n/2$, then $\mu(\mathbf{x}) = 0 \leq \mu_0$ and we choose $\mathbf{b} = \mathbf{0}$. Otherwise, we have that $\text{wt}(\mathbf{x} + \mathbf{b}) \notin \{0, n\}$ for all $\mathbf{b} \in \mathbb{B}_{1, m+1} \subseteq \text{RM}(r+1, m+1)$. We compute the Walsh spectrum of \mathbf{x} . Again, if $\mu(\mathbf{x} + \mathbf{b}) = 0$ for some $\mathbf{b} \in \mathbb{B}_{1, m+1}$, we simply choose the corresponding \mathbf{b} .

Hence, in the worst case, $\text{wt}(\mathbf{x} + \mathbf{b}) \notin \{0, n/2, n\}$ for all $\mathbf{b} \in \mathbb{B}_{1, m+1}$. Since there are $(n+1)-3 = n-2$ possible weights and the number of words in $\mathbb{B}_{1, m+1}$ is $2^{m+1} = n$, we must have $\text{wt}(\mathbf{x} + \mathbf{a}) = \text{wt}(\mathbf{x} + \mathbf{b})$ for some $\mathbf{a}, \mathbf{b} \in \mathbb{B}_{1, m+1}$. Using the automorphism group of the RM codes [17, Ch. 13], we can assume that

$$\mathbf{a} = 0^{n/4}1^{n/4}1^{n/4}0^{n/4}, \text{ and } \mathbf{b} = 1^{n/4}0^{n/4}1^{n/4}0^{n/4}.$$

Now, we divide \mathbf{x} into four blocks of equal length (i.e. $n/4$) and let w_i denote the weight of the i th block for $i \in \{1, 2, 3, 4\}$. Since $\text{wt}(\mathbf{x} + \mathbf{a}) = w_1 + w_4 + (n - w_2 - w_3)$, and $\text{wt}(\mathbf{x} + \mathbf{b}) = w_2 + w_4 + (n - w_1 - w_3)$, we have that $w_1 = w_2$.

We then consider $\mathbf{c} = 1^{n/4}0^{3n/4} \in \mathbb{B}_{2, m+1} \subseteq \text{RM}(r+1, m+1)$ and check that the number of ones in the first $n/2$ coordinates of $\mathbf{x} + \mathbf{c}$ is exactly $(n/4 - w_1) + w_2 = n/4$. In other words, the prefix of length $n/2$ of $\mathbf{x} + \mathbf{c}$ has zero imbalance.

Let \mathbf{x}' be the suffix of \mathbf{x} of length $n/2$. Since $\mu(\text{RM}(r, m)) = \mu_0$, we can find \mathbf{d} in $\text{RM}(r, m)$ such that $\mu(\mathbf{x}' + \mathbf{d}) \leq \mu_0$. Since $\text{RM}(r+1, m+1) = \{(\mathbf{u}, \mathbf{u} + \mathbf{v}) : \mathbf{u} \in \text{RM}(r+1, m), \mathbf{v} \in \text{RM}(r, m)\}$ [17, Ch. 13], the word $\mathbf{0d}$ belongs to $\text{RM}(r+1, m+1)$. Therefore, if we choose the balancing word to be $\mathbf{c} + \mathbf{0d} \in \text{RM}(r+1, m+1)$, then the resulting imbalance $\mu(\mathbf{x} + \mathbf{c} + \mathbf{0d}) = \mu(\mathbf{x} + \mathbf{c}) + (\mathbf{x}' + \mathbf{d}) \leq \mu_0$, as desired. \square

We apply Proposition 15 recursively to obtain the next result.

Corollary 16. $\mu(\mathbb{B}_{r, m}) \leq \lfloor 2^{(m-r+1)/2} \rfloor$ for $1 \leq r \leq m$.

An immediate consequence of Corollary 16 is that $\mu(\mathbb{B}_{2, 4}) \leq 2$. Recall that \mathbb{X}_1 is the set of all words with even weight and it is now immediate that $\mu(\mathbb{X}_1, \mathbb{B}_{2, 4}) = 0$. In other words, $\mathbb{B}_{2, 4}$ or $\text{RM}(2, 4)$ balances all words of even weight. The next proposition states that $\text{RM}(2, 5)$ also balances all words of even weight.

Proposition 17. $\mu(\mathbb{X}_1, \mathbb{B}_{2, 4}) = \mu(\mathbb{X}_1, \mathbb{B}_{2, 5}) = 0$.

We defer the proof to the next subsection where we allude to certain geometric properties of RM codes. Nevertheless, a direct consequence of Propositions 15 and 17 is the following corollary.

Corollary 18. $\mu(\mathbb{X}_1, \mathbb{B}_{m-3, m}) = 0$ for $m \geq 5$.

Now, for all $m \geq 2$, define the quantity

$$\rho(m) \triangleq \min\{r : \mu(\mathbb{X}_1, \mathbb{B}_{r, m}) = 0\}. \quad (8)$$

That is, $\rho(m)$ is the smallest value r such that $\mathbb{B}_{r, m}$ balances all even-weight words. The results of the paper state that

$$\rho(2) = \rho(3) = 1, \rho(4) = \rho(5) = 2, 1 < \rho(m) \leq m-3 \text{ for } m \geq 6.$$

For $m \geq 4$, we ran computer simulations and it seems that $\mathbb{B}_{2, m}$ is always able to balance an even-weight word. Hence, we make the following conjecture.

Conjecture. $\rho(m) = 2$ for all $m \geq 4$.

A. Proof of Proposition 17

In this subsection, we show that $\text{RM}(2, 5)$ balances all words of even weight. To do so, we recall certain classical characterizations

of RM codewords. Fix $m \geq 2$ and index the codewords in $\text{RM}(r, m)$ with \mathbb{F}_2^m . For $X \subseteq \mathbb{F}_2^m$, we use $\mathbf{c}(X)$ to denote the binary word whose support is exactly X . In other words, if $\mathbf{c}(X) = (c_{\alpha})_{\alpha \in \mathbb{F}_2^m}$, we have that $c_{\alpha} = 1$ if and only if $\alpha \in X$.

Lemma 19 ([17, Ch. 13]). *If $X \subseteq \mathbb{F}_2^m$ is an affine space with dimension at least $m-r$, then $\mathbf{c}(X)$ belongs to $\text{RM}(r, m)$.*

Using this characterization, the next lemma tells us when a word in $\text{RM}(1, m)$ is able to balance a certain even-weight word.

Lemma 20. *Let $X \subseteq \mathbb{F}_2^m$ with $|X|$ even. If U is a linear space of dimension $m-1$ such that $|U \cap X| = |X|/2$, then $\mathbf{c}(U) \in \text{RM}(1, m)$ balances $\mathbf{c}(X)$.*

The next lemma follows from the covering radius of $\text{RM}(2, 5)$.

Lemma 21 ([19]). *For $\mathbf{x} \in \mathbb{X}_1$, there exists $\mathbf{c} \in \text{RM}(2, 5)$ such that $\text{wt}(\mathbf{x} + \mathbf{c}) \in \{0, 2, 4, 6\}$.*

We are now ready to show that $\text{RM}(2, 5)$ balances all words of even weight. Our broad strategy is as follows.

- Let X be the support of the word to be balanced. From Lemma 21, it suffices to assume that $|X| \in \{0, 2, 4, 6\}$.
- To apply Lemma 20, we find a four-dimensional linear space U such that $|U \cap X| = |X|/2$. In most cases, this is always possible and in what follows, we simply state the space U and omit the detailed verification of the intersection size. When it is not possible to find U , we modify X so that Lemma 20 is applicable.
- We assume that X contains $\mathbf{0}$. Otherwise, we can shift it using the automorphism group of $\text{RM}(2, 5)$. Also, the vectors $\alpha, \beta, \gamma, \delta$, and ϵ are assumed to be linearly independent.

We consider the following cases according to the size of X .

- If $|X| = 0$, we pick any four-dimensional space U .
- If $|X| = 2$, then $X = \{\mathbf{0}, \alpha\}$. Pick the three-dimensional space U such that $\alpha \notin U$. For example, we can set $U = \text{span}(\{\beta, \gamma, \delta, \epsilon\})$.
- When $|X| = 4$, we have the following two cases. If the rank of X is two, then $X = \{\mathbf{0}, \alpha, \beta, \alpha + \beta\}$ and we pick $U = \text{span}(\{\alpha, \gamma, \delta, \epsilon\})$. If the rank of X is three, then $X = \{\mathbf{0}, \alpha, \beta, \gamma\}$ and we pick $U = \text{span}(\{\alpha + \beta, \gamma, \delta, \epsilon\})$.
- When $|X| = 6$, we have the following four cases.
 - (i) If the rank of X is three, we pick a three-dimensional space V that contains X . So, $\mathbf{c}(V)$ belongs to $\text{RM}(2, 5)$ and $\mathbf{c}(V) + \mathbf{c}(U)$ has weight two. We then proceed as before to balance the word $\mathbf{c}(V) + \mathbf{c}(U)$.
 - (ii) Suppose that the rank of X is four and $X = \{\mathbf{0}, \alpha, \beta, \gamma, \delta, \alpha + \beta + \gamma + \delta\}$. Here, we set $V = \text{span}(\{\alpha, \beta, \epsilon\})$ and so, $\mathbf{c}(V)$ belongs to $\text{RM}(2, 5)$. Let X' be the support of $\mathbf{c}(X) + \mathbf{c}(X)$ and here, $|X'| = 8$. Now, if we pick the four-dimensional space $U = \text{span}(\{\alpha, \beta, \gamma, \delta\})$, then $|X' \cap U| = 4$. Therefore, Lemma 20 states that $\mathbf{c}(X') + \mathbf{c}(U) = \mathbf{c}(X) + \mathbf{c}(V) + \mathbf{c}(U)$ is balanced.
 - (iii) Suppose that the rank of X is four and X contains $\{\mathbf{0}, \alpha, \beta, \gamma, \delta\}$ with the sixth vector not equal to the sum of the five nonzero vectors. Then we may assume that the sixth vector is $\beta + f(\gamma, \delta)$ where $f(\gamma, \delta)$ is a linear combination of γ and δ , not involving α . In this case, we pick $U = \text{span}(\{\alpha + \beta, \gamma, \delta, \epsilon\})$.
 - (iv) If the rank of X is five, then $X = \{\mathbf{0}, \alpha, \beta, \gamma, \delta, \epsilon\}$ and we pick $U = \text{span}(\{\alpha + \beta, \alpha + \gamma, \delta, \epsilon\})$.

ACKNOWLEDGEMENT

We thank the reviewers for the insightful comments. In particular, Remark 6 was highlighted by an anonymous reviewer.

REFERENCES

- [1] K. A. S. Immink, *Codes for Mass Data Storage Systems*. Shannon Foundation Publisher, 2004.
- [2] S. Yazdi, H. M. Kiah, E. R. Garcia, J. Ma, H. Zhao, and O. Milenkovic, “DNA-based storage: Trends and methods,” *IEEE Trans. Molecular, Biological, Multi-Scale Commun.*, vol. 1, no. 3, pp. 230–248, 2015.
- [3] M. G. Ross, C. Russ, M. Costello, A. Hollinger, N. J. Lennon, R. Hegarty, C. Nusbaum, D. B. Jaffe, “Characterizing and Measuring Bias in Sequence Data”, *Genome Biol.* 14, R51, 2013.
- [4] K. A. S. Immink and K. Cai, “Properties and constructions of constrained codes for DNA-based data storage,” *arXiv preprint arXiv:1812.06798*, 2018.
- [5] H. Van Tilborg and M. Blaum, “On error-correcting balanced codes,” *IEEE Trans. Inform. Theory*, vol. 35, no. 5, pp. 1091–1095, 1989.
- [6] S. Al-Bassam and B. Bose, “Design of efficient error-correcting balanced codes,” *IEEE Trans. Computers*, vol. 42, no. 10, pp. 1261–1266, 1993.
- [7] A. Mazumdar, R. M. Roth, and P. O. Vontobel, “On linear balancing sets,” *Adv. Math. Commun.* vol. 4, no. 3, pp. 345–361, 2010.
- [8] J. H. Weber, K. A. S. Immink, and H. C. Ferreira, “Error-correcting balanced knuth codes,” *IEEE Trans. Inform. Theory*, vol. 58, no. 1, pp. 82–89, 2012.
- [9] Y. M. Chee, H. M. Kiah, and H. J. Wei, “Efficient and Explicit Balanced Primer Codes,” *arXiv preprint arXiv:1901.01023*, 2019.
- [10] K. Cai, Y. M. Chee, R. Gabrys, H. M. Kiah, and T. T. Nguyen, “Optimal Codes Correcting a Single Indel / Edit for DNA-Based Data Storage”, *arXiv preprint arXiv:1910.06501*, 2019.
- [11] E. Arikan. “Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels”, *IEEE Trans. Inform. Theory*, vol. 55, no. 7, pp. 3051–3073, 2009.
- [12] D. Kane, *Samsung Licenses 5G Polar Coding Technology Developed by UC San Diego Engineers*, UC San Diego Press Release, October 2018. [Online]. Available as: https://ucsdnews.ucsd.edu/pressrelease/samsung_licenses_5g_polar_coding_technology_developed_by_uc_san_diego_engineers
- [13] K. Tian, A. Fazeli, and A. Vardy, “Polar coding for deletion channels: Theory and implementation”, in *Proc. IEEE Int. Symp. Inform. Theory*, 2018, pp. 1869–1873.
- [14] I. Tal, H. D. Pfister, A. Fazeli, and A. Vardy, “Polar Codes for the Deletion Channel: Weak and Strong Polarization”, *arXiv preprint arXiv:1904.13385*, 2019.
- [15] E. Şaşoğlu, “Polarization in the presence of memory,” in *Proc. IEEE Symp. Inform. Theory*, 2011, pp. 189–193.
- [16] B. Shuval and I. Tal. “Fast Polarization for Processes with Memory.” *IEEE Trans. Inf. Theory*, vol. 65, no. 4, pp. 2004–2020, 2019.
- [17] F. J. MacWilliams, and N. J. A. Sloane. *The Theory of Error-Correcting Codes*. North-Holland, 1977.
- [18] E. Abbe and M. Ye, “Reed-Muller codes polarize,” [Online]. *arXiv preprint arXiv:1901.11533*, 2019.
- [19] G. D. Cohen, and S. N. Litsyn, “On the covering radius of Reed-Muller codes,” *Discrete Mathematics*, 106, pp. 147–155, 1992.