

# Asymmetric Lee Distance Codes for DNA-Based Storage

Ryan Gabrys<sup>\*†</sup>, Han Mao Kiah<sup>‡</sup>, and Olgica Milenkovic<sup>§</sup>

<sup>\*</sup> <sup>§</sup> Coordinated Science Laboratory, University of Illinois, Urbana-Champaign, USA

<sup>†</sup>Spawar Systems Center San Diego, Code 532, USA

<sup>‡</sup> School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore

Emails: <sup>\*</sup>gabrys@illinois.edu, <sup>†</sup>ryan.gabrys@navy.mil, <sup>‡</sup>hmkiah@ntu.edu.sg, <sup>§</sup>milenkov@illinois.edu

**Abstract**—We consider a new family of asymmetric Lee codes that arise in the design and implementation of DNA-based storage systems and systems with parallel string transmission protocols. The codewords are defined over a quaternary alphabet, although the results carry over to other alphabet sizes, and have symbol distances dictated by their underlying binary representation. Our contributions are two-fold. First, we derive upper bounds on the size of the codes under the asymmetric Lee distance measure based on linear programming techniques. Second, we propose code constructions which imply lower bounds.

**Keywords.** Coding for DNA-based storage, Coding theory

## I. INTRODUCTION

Codes for classical channels with single-sequence inputs and single-sequence outputs have been studied extensively, leading to a diverse suite of solutions including algebraic codes [15], codes on graphs, such as LDPC codes [14], and polar codes [16]. Similar advances have been reported for parallel channels [7], with the rather common underlying assumption that the channels introduce uncorrelated errors. The alphabet size of the codes used in both scenarios is usually restricted by the system design, and often, input sequences are de-interleaved or represented as arrays over smaller alphabets in order to enable more efficient transmission. Far less is known about channels that operate on several sequences at the same time and introduce correlated symbol errors, or output ordering errors. The goal of this work is to analyze one such scenario, motivated by emerging applications in DNA-based storage systems.

To motivate our analysis, consider a transmission model in which two binary input sequences are simultaneously passed through *two channels* that introduce substitution errors with some probability  $0 < p < 1/2$  (Figure 1). Simultaneous errors in both strings are less likely than individual string errors. In addition to the substitution errors, the outputs of the channels may be switched – in other words, the label of the channel from which the output symbol originated may be in error. The confusion graph of this type of channel is depicted in Figure 2. In Figure 2, the vertices are indexed by pairs of bits denoting the inputs into the two channels. The labels of the edges denote the channel confusion parameters (weights, distances).

One application of such a model arises in DNA sequencing for archival storage [6], where multiple sequences are read in parallel. The readout errors are rare and it is very uncommon to make simultaneous mistakes in both sequences; nevertheless, the identity of the strands may be confused due to string sorting issues. Another unrelated way to view this model is to assume that the binary encodings represent four symbols of the DNA alphabet  $\{A, T, G, C\}$ , say  $00 \rightarrow G$ ,  $11 \rightarrow C$ ,  $01 \rightarrow A$  and

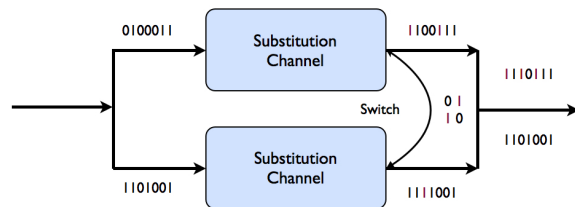


Fig. 1. A pair of channels with individual substitution errors, for which the outputs may also be switched. For the given example, the outputs of the channels at position three are switched.

$10 \rightarrow T$ . In this case, the proposed metric illustrates a DNA readout channel in which the bases  $\{A, T\}$  are very likely to be mutually confused during sequencing, while the basis  $\{G, C\}$  are much less likely to be misinterpreted for each other. Illumina systems and some other sequencing devices have substitution errors that exhibit such “bias” phenomena, and similar effects may be expected for single base sequencing technologies of the third generation [13].

The problem of interest is to design pairs of sequences – henceforth, termed codewords – such that any two codewords are at a sufficiently large “distance” from each other. We subsequently refer to the distance imposed by Figure 2 as the *asymmetric Lee distance* (ALD). The ALD resembles a weighted version of the Lee metric [2], with the exception of *two symbols* being treated differently. These two symbols capture the uncertainty about the actual ordering of the readouts. Given the connection with the Lee metric, a partial analysis of the ALD and related code construction questions may be addressed by invoking results for codes in the Lee metric [1], [2]. Nevertheless, due to the asymmetry of the distance, specialized techniques need to be developed to find bounds on the code size and to construct codes that approach these bounds. To accomplish this task, we formally define the ALD as a judiciously chosen combination of the Lee and the Hamming distance.

The paper is organized as follows. In Section II we introduce the ALD, and then proceed to derive upper and lower bounds on the size of the largest code in this metric using results on codes in the Lee metric. In Section III, we turn our attention on computing tighter upper bounds on the size of codes by applying generalized sphere packing techniques not known from the Lee coding literature. New code constructions and related questions are discussed in Section IV.

## II. PROBLEM FORMULATION

We start by characterizing the channel errors by introducing a new distance measure, which we refer to as the *asymmetric*

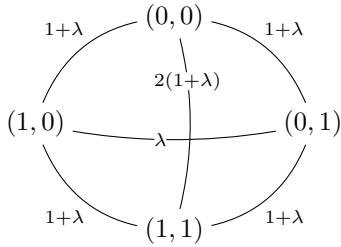


Fig. 2. Weighted confusion graph for codelength  $n = 1$ .

Lee distance (ALD).

We refer to an error that causes a paired-symbol transition from  $(1, 0)$  to  $(0, 1)$  as a Class 1 (switch) error; similarly, we refer to an error that causes a single substitution in one of the input strings as a Class 2 error. An error that causes a paired-symbol transition from  $(0, 0)$  to  $(1, 1)$  is referred to as a Class 3 (simultaneous substitution) error. Note that based on Figure 2, an edge in the confusion graph corresponding to a Class 1 error has weight  $\lambda$ , an edge for a Class 2 error has weight  $1 + \lambda$ , while an edge for a Class 3 error has weight  $2(1 + \lambda)$ .

Let  $m \geq 2$  be an integer. For  $a_1, a_2, \dots, a_m \in \mathbb{Z}_2$ , the indicator function  $\chi$  is such that  $\chi(a_1, a_2, \dots, a_m) = 1$ , if  $a_1 = a_2 = \dots = a_m$  and  $\chi(a_1, a_2, \dots, a_m) = 0$  otherwise. Consider next four sequences  $\mathbf{a} = (a_1, \dots, a_n)$ ,  $\mathbf{b} = (b_1, \dots, b_n)$ ,  $\mathbf{c} = (c_1, \dots, c_n)$ ,  $\mathbf{d} = (d_1, \dots, d_n) \in \mathbb{Z}_2^n$ , paired as  $(\mathbf{a}, \mathbf{b})$  and  $(\mathbf{c}, \mathbf{d})$ , and define two pseudo-metrics:

$$d_S((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) := \sum_{i=1}^n \chi(a_i, b_i) + \chi(c_i, d_i) - 2\chi(a_i, b_i, c_i, d_i),$$

$$d_D((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) := \sum_{i=1}^n 2(\chi(a_i, b_i) + \chi(c_i, d_i)) + \chi(a_i, \bar{b}_i, \bar{c}_i, d_i) - 4\chi(a_i, b_i, c_i, d_i),$$

where  $\bar{x} = 1 - x$ , and  $x \in \mathbb{Z}_2$ . Note that the pseudo-metrics  $d_S((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d}))$  and  $d_D((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d}))$  both depend on how much the sequences within the pairing  $(a_i, b_i), (c_i, d_i)$ ,  $i \in [n]$ , disagree from each other, as well as how much the pairs differ themselves. Furthermore, observe that  $d_S$  is invariant under the change of order in the pairing, i.e.,  $d_S((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) = d_S((\mathbf{b}, \mathbf{a}); (\mathbf{c}, \mathbf{d}))$ , and that it takes the maximum value  $2n$  if and only if the sequences  $\mathbf{a}, \mathbf{b}$  are identical and complements of  $\mathbf{c}, \mathbf{d}$ . Hence, this pseudo-metric does not penalize reorderings of symbols in the same pair of strings. On the other hand, it is easy to see that  $d_D$  is actually a metric which assigns non-zero distances to pairs in which the bits are switched. Consequently, combining  $d_S$  and  $d_D$  in the coding process should allow one to control both the effects of symbol switching events between the channels and the substitution errors in the channels.

For a positive real number  $\lambda$ , define the ALD  $d_\lambda((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d}))$  as a combination of  $d_S$  and  $d_D$ ,

$$\lambda d_D((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) + (1 - \lambda) d_S((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) \quad (1)$$

$$= \sum_{i=1}^n (1 + \lambda) (\chi(a_i, b_i) + \chi(c_i, d_i)) + \lambda \chi(a_i, \bar{b}_i, \bar{c}_i, d_i) - 2(1 + \lambda) \chi(a_i, b_i, c_i, d_i).$$

It can be shown that  $d_\lambda((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d}))$  is a metric. Hence, for any  $\lambda > 0$ ,  $d_\lambda$  is non-negative, symmetric and satisfies the tri-

angle inequality. We henceforth focus our attention on integer  $\lambda$ .

From (1), we observe that the paired symbol distance is asymmetric, in so far that complementary pairs are treated differently than non-complementary pairs. Furthermore, when pairs are complementary, the distance depends on the binary weight of the pairs. The choice of the parameter  $\lambda$  governs the degree of the asymmetry in the distance.

In order to highlight the relationship between  $d_\lambda((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d}))$  and the Lee distance, define the mapping  $\mathcal{Z} : \mathbb{Z}_2^2 \rightarrow \mathbb{Z}_4$  so that  $00 \rightarrow 1, 10 \rightarrow 0, 01 \rightarrow 2$  and  $11 \rightarrow 3$ . Then, by changing the weight between  $(10)$  and  $(01)$  – from  $\lambda$  to  $2(1 + \lambda)$ , we arrive at a scaled version of the Lee distance  $d_L(x, y)$  between symbols  $x, y \in \mathbb{Z}_4$ , which reads as

$$(1 + \lambda) \cdot \min\{4 - |x - y|, |x - y|\} = (1 + \lambda) \cdot d_L(x, y). \quad (2)$$

More precisely, for two sequences  $\mathbf{x}, \mathbf{y}$  over  $\mathbb{Z}_4$ , we have

$$d_\lambda(\mathbf{x}, \mathbf{y}) = (1 + \lambda) \cdot d_L(\mathbf{x}, \mathbf{y}) - (\lambda + 2) \sum_{\{x_i, y_i\} = \{0, 2\}} \chi(x_i, y_i) \leq (1 + \lambda) \cdot d_L(\mathbf{x}, \mathbf{y}). \quad (3)$$

The main difference between the ALD and the Lee metric is that the distance between pairs of symbols, say  $(s, p), (q, r) \in \mathbb{Z}_2^2$  is not completely characterized by the weight of the vector  $(s, p) + (q, r)$ , and that in particular, it depends on the exact values of  $(s, p)$  in a manner analogous to asymmetric error correcting codes. This connection between the ALD, Lee metric, and asymmetric error correcting codes will be used in our subsequent derivations.

Let  $d$  be a real positive integer. We say that two pairs of sequences  $(\mathbf{a}, \mathbf{b}), (\mathbf{c}, \mathbf{d}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n$  are  $(d, \lambda)$ -distinguishable if their ALD  $d_\lambda$  is at least  $d$ ; similarly, we say that two pairs of sequences are  $(d, \lambda)$ -indistinguishable if their ALD is less than  $d$ . We use  $A_\lambda(n, d)$  to denote the largest number of  $(d, \lambda)$ -distinguishable sequences of length  $n$ . In the next section, we derive a number of upper bounds on  $A_\lambda(n, d)$ .

### III. UPPER BOUNDS

We start by computing  $A_\lambda(n, d)$  using bounds on the size of codes in the Lee metric. Let  $\mathcal{Z} : \mathbb{Z}_2^2 \rightarrow \mathbb{Z}_4$  denote the mapping introduced in the previous section. Furthermore, let  $A_L(n, d, q)$  denote the size of the largest code of length  $n$  with minimum Lee distance  $d$  over  $\mathbb{Z}_q$ . Since  $\mathcal{Z}$  is invertible, we have the following result which easily follows from (3).

**Proposition 1.** For positive integers  $n$  and  $d$ , and positive integer  $\lambda$ ,  $A_\lambda(n, d) \leq A_L(n, \lfloor \frac{d}{1+\lambda} \rfloor, 4)$ .

A recent Singleton-type bound for codes in the Lee metric and over an even-sized alphabet [1] asserts that

$$A_\lambda(n, d) \leq A_L(n, \lfloor \frac{d}{1+\lambda} \rfloor, 4) \leq 2^{2n - \lfloor \frac{d}{1+\lambda} \rfloor + 1}. \quad (4)$$

Note that the bound is general, as it applies to codes that are not necessarily linear.

### A. Sphere Packing Bounds on $A_\lambda(n, d)$ for large $n$

In what follows, we first derive an expression for the size of a ball of radius  $r$  in the ALD. Using this expression, we proceed to find asymptotic upper bounds on  $A_\lambda(n, 2r + 1)$ .

Given a vector  $\mathbf{x} \in \mathbb{Z}_2^n$ , let  $wt(\mathbf{x})$  denote the Hamming weight of  $\mathbf{x}$ . For a pair of sequences  $(\mathbf{a}, \mathbf{b}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n$ , let

$$\mathcal{B}_{(r,\lambda)}(\mathbf{a}, \mathbf{b}) := \{(\mathbf{c}, \mathbf{d}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n : d_\lambda((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) \leq r\}$$

denote the set of pairs of sequences that are  $(r + 1, \lambda)$ -indistinguishable from  $(\mathbf{a}, \mathbf{b})$ . Note that from (1), the quantity  $|\mathcal{B}_{(r,\lambda)}(\mathbf{a}, \mathbf{b})|$  is a function of the Hamming weight of  $\mathbf{a} - \mathbf{b}$ , which we for simplicity denote by  $w = w(\mathbf{a}, \mathbf{b})$ . Let

$$S_\lambda(n, w, \delta) := |\{(\mathbf{c}, \mathbf{d}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n : d_\lambda((\mathbf{a}, \mathbf{b}); (\mathbf{c}, \mathbf{d})) = \delta\}|,$$

and

$$V_\lambda(n, w, r) := |\mathcal{B}_{(r,\lambda)}(\mathbf{a}, \mathbf{b})|,$$

where we omitted the arguments  $\mathbf{a}, \mathbf{b}$  on the left-hand side of the equations for simplicity of notation. Clearly,  $V_\lambda(n, w, r) = \sum_{j=0}^r S_\lambda(n, w, j)$ .

Using the same techniques as described in [2], we compute the generating function of  $S_\lambda(n, w, \delta)$ ,

$$\sum_{w=0}^{\infty} \sum_{r=0}^{\infty} S_\lambda(n, w, r) x^w z^r = x^w (z^\lambda + (1 + 2z^{1+\lambda}))^w. \quad (5)$$

$$\left( z^{2(1+\lambda)} + (1 + 2z^{1+\lambda}) \right)^{n-w}.$$

From (5), it follows that  $S_\lambda(n, w, r)$  equals the coefficient of  $x^w z^r$ . Since  $V_\lambda(n, w, r) = \sum_{i=0}^r S_\lambda(n, w, i)$ , the next lemma follows after some straightforward algebraic manipulations. For simplicity of notation, we use  $k * \ell * m \leq (r, \lambda)$  to denote  $(2k + \ell)(1 + \lambda) + \lambda m \leq r$ .

**Lemma 2.** For integers  $n, w, r$ ,  $V_\lambda(n, w, r)$  may be written as

$$\sum_{k * \ell * m \leq (r, \lambda)} \binom{w}{m} \binom{n-w}{k} \binom{n-k-m}{\ell} 2^\ell.$$

**Corollary 3.** For positive integers  $n, w$ , and  $r$ , it holds that  $V_\lambda(n, w, r) \geq V_\lambda(n, w - 1, r)$ .

Using Lemma 2, one can derive the following asymptotic upper bound on the size of an ALD, following similar arguments as those used in [11, Theorem 3].

**Lemma 4.** There exists an  $N$  such that for  $n \geq N$ , we have

$$A_\lambda(n, 2r + 1) \leq \frac{4^n}{V_\lambda(n, \lfloor n/2 - \sqrt{5nr \log_2(n)} \rfloor, r)} \cdot (1 + o(1)).$$

From Lemma 4, one recovers the following two special upper bounds.

**Corollary 5.** For  $n \geq N$  and  $L(n) = \lfloor n/2 - \sqrt{10n \log_2(n)} \rfloor$ ,

$$A_1(n, 5) \leq \frac{4^n}{1 + 2n + \sum_{j=1}^2 \binom{L(n)}{j}}.$$

**Corollary 6.** For  $n \geq N$  and  $L(n) = \lfloor n/2 - \sqrt{15n \log_2(n)} \rfloor$ ,

$$A_1(n, 7) \leq \frac{4^n}{1 + \sum_{j=1}^3 \binom{L(n)}{j} + 2L(n)(n-1) + 2n}.$$

### B. Non-Asymptotic Upper Bounds

We start by introducing our notation, following [4]. Fix  $r, \lambda$ ,  $\mathbf{x}$  and recall the definition of  $\mathcal{B}_{(r,\lambda)}(\mathbf{x})$ . Define a directed graph  $\mathcal{G}_{(r,\lambda)}$  on the vertex set  $\mathbb{Z}_2^n \times \mathbb{Z}_2^n$  such that there exists an arc from  $\mathbf{x}$  to  $\mathbf{y}$  if  $\mathbf{y} \in \mathcal{B}_{(r,\lambda)}(\mathbf{x})$ . Let  $\mathbf{A}_{(r,\lambda)}$  be an adjacency matrix of dimension  $4^n \times 4^n$ , indexed by the elements in  $\mathbb{Z}_2^n \times \mathbb{Z}_2^n$ . Then  $\mathbf{A}_{(r,\lambda)}(i, j) = 1$  if  $\mathbf{x}_i \in \mathcal{B}_{(r,\lambda)}(\mathbf{x}_j)$ , and zero otherwise. Let  $\mathbb{R}_+$  denote the set of non-negative reals, and define

$$\tau^*(\mathbf{A}_{(r,\lambda)}) = \min \left\{ \sum_{i=1}^{4^n} w_i : \mathbf{w} \in \mathbb{R}_+^{4^n}, \mathbf{A}_{(r,\lambda)}^T \cdot \mathbf{w} \geq \mathbf{1} \right\}. \quad (6)$$

The results of [10] show that one may use  $\tau^*(\mathbf{A}_{(r,\lambda)})$  as an upper bound for  $A_\lambda(n, 2r + 1)$ . However, (6) is a linear program involving  $4^n$  equations and we may significantly reduce this number by observing certain symmetries of  $\mathcal{G}_{(r,\lambda)}$ .

An automorphism of  $\mathcal{G}_{(r,\lambda)}$  is a permutation of its vertices that preserves adjacency. Let  $\mathbb{S}_n$  be the symmetric group on  $n$  symbols. The set of all automorphisms of  $\mathcal{G}_{(r,\lambda)}$  is defined as  $\text{Aut}(\mathcal{G}_{(r,\lambda)}) = \{\pi \in \mathbb{S}_n | \pi \text{ is an automorphism of } \mathcal{G}\}$ . Given a subgroup  $H$  of  $\text{Aut}(\mathcal{G}_{(r,\lambda)})$ , let  $H$  partition the vertex set into  $n_H$  equivalence classes  $\{X_1, \dots, X_{n_H}\}$ . Let  $\mathbf{A}_{H,(r,\lambda)}$  be an  $n_H \times n_H$  adjacency matrix corresponding to the subgroup  $H$ , such that for  $1 \leq i, j \leq n_H$ ,

$$\mathbf{A}_{H,(r,\lambda)}(i, j) = \frac{|\{(\mathbf{x}, \mathbf{y}) : \mathbf{x} \in X_j, \mathbf{y} \in X_i, \mathbf{y} \in \mathcal{B}_{(r,\lambda)}(\mathbf{x})\}|}{|X_i|}.$$

The authors of [4] demonstrated that solving (6) is equivalent to solving a linear program involving only  $\mathbf{A}_{H,(r,\lambda)}$ . Combining this result with those of [10], we have the following.

**Theorem 7.** (c.f. [4], [10]) Let  $H$  be a subgroup of  $\text{Aut}(\mathcal{G}_{(r,\lambda)})$  and define  $\mathbf{A}_{H,(r,\lambda)}$  as above. Then,

$$\tau^*(\mathbf{A}_{(r,\lambda)}) = \min \left\{ \sum_{i=1}^{n_H} |X_i| \cdot w_i : \mathbf{w} \in \mathbb{R}_+^{n_H}, \mathbf{A}_{H,(r,\lambda)}^T \cdot \mathbf{w} \geq \mathbf{1} \right\}.$$

and  $A_\lambda(n, 2r + 1) \leq \tau^*(\mathbf{A}_{(r,\lambda)})$ .

Next, we define a set of automorphisms on  $\mathcal{G}_{(r,\lambda)}$  as follows. We first introduce a mapping denoted  $\pi_{\sigma, \mathbf{x}}$ . For every permutation  $\sigma$  in the symmetric group  $\mathbb{S}_n$  and any  $\mathbf{x} \in \mathbb{Z}_2^n$ , let  $\pi_{\sigma, \mathbf{x}} : \mathbb{Z}_2^n \times \mathbb{Z}_2^n \rightarrow \mathbb{Z}_2^n \times \mathbb{Z}_2^n$  be a mapping such that for all  $(\mathbf{a}, \mathbf{b}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n$ ,  $(\pi_{\sigma, \mathbf{x}}(\mathbf{a}, \mathbf{b}))_i = (\mathbf{a}_{\sigma(i)}, \mathbf{b}_{\sigma(i)})$  if  $x_i = 0$  and  $(\pi_{\sigma, \mathbf{x}}(\mathbf{a}, \mathbf{b}))_i = (\mathbf{b}_{\sigma(i)}, \mathbf{a}_{\sigma(i)})$  otherwise. For instance, if  $(\mathbf{a}, \mathbf{b}) = ((0, 0, 1), (1, 0, 1))$ ,  $\mathbf{x} = (0, 1, 0)$ , and  $\sigma = (3, 1, 2)$ , then  $\pi_{\sigma, \mathbf{x}}(\mathbf{a}, \mathbf{b}) = ((1, 1, 0), (1, 0, 0))$ . It can be shown that the set  $H = \{\pi_{\sigma, \mathbf{x}} : \sigma \in \mathbb{S}_n, \mathbf{x} \in \mathbb{Z}_2^n\}$  is a subgroup of  $\text{Aut}(\mathcal{G}_{(r,\lambda)})$ . As a result of the previous discussion, a bound on  $A_\lambda(n, 2r + 1)$  may be obtained by considering the quantity

$$\tau^*(\mathbf{A}_{(r,\lambda)}) = \min \left\{ 2^n \cdot \sum_{\ell=0}^n \binom{n}{\ell} \cdot w_\ell : \mathbf{w} \in \mathbb{R}_+^{n+1}, \mathbf{A}_{H,(r,\lambda)}^T \cdot \mathbf{w} \geq \mathbf{1} \right\}. \quad (7)$$

Note that using the automorphism group, the dimension of the weight vector  $\mathbf{w}$  has been reduced from  $O(4^n)$  (as in (6)) to

$n+1$  in (7). The next theorem provides a feasible weight vector for the optimization problem in (7).

**Theorem 8.** Let  $V_\lambda(n, w, r) = V_\lambda(n, 0, r)$  for  $w < 0$ . Suppose that  $n, r, \lambda > 0$  and that  $\mu = \lfloor \frac{r}{1+\lambda} \rfloor$ . Let  $\mathbf{w} = (w_0, w_1, \dots, w_n) \in \mathbb{R}_+^{n+1}$  be defined as  $w_i = \frac{1}{V_\lambda(n, i - \mu, r)}$ . Then,

$$A_\lambda(n, 2r + 1) \leq 2^n \cdot \sum_{i=0}^n \binom{n}{i} \cdot \frac{1}{V_\lambda(n, i - \mu, r)}.$$

*Proof:* The result follows directly from Theorem 7 and (7), provided that we can show that  $\mathbf{A}_{H, (r, \lambda)}^T \cdot \mathbf{w} \geq 1$ . For the given weight assignment, and for an arbitrary choice of  $(\mathbf{a}, \mathbf{b}) \in \mathbb{Z}_2^n \times \mathbb{Z}_2^n$ , we have

$$\begin{aligned} \sum_{(\mathbf{c}, \mathbf{d}) \in \mathcal{B}_{(r, \lambda)}(\mathbf{a}, \mathbf{b})} w_{w(\mathbf{c}, \mathbf{d})} &= \sum_{(\mathbf{c}, \mathbf{d}) \in \mathcal{B}_{(r, \lambda)}(\mathbf{a}, \mathbf{b})} \frac{1}{V_\lambda(n, w(\mathbf{c}, \mathbf{d}) - \mu, r)} \\ &\geq \sum_{(\mathbf{c}, \mathbf{d}) \in \mathcal{B}_{(r, \lambda)}(\mathbf{a}, \mathbf{b})} \frac{1}{V_\lambda(n, w(\mathbf{a}, \mathbf{b}), r)} \\ &= \frac{V_\lambda(n, w(\mathbf{a}, \mathbf{b}), r)}{V_\lambda(n, w(\mathbf{a}, \mathbf{b}), r)} = 1. \end{aligned}$$

Since  $(\mathbf{c}, \mathbf{d}) \in \mathcal{B}_{(r, \lambda)}(\mathbf{a}, \mathbf{b})$ ,  $w(\mathbf{c}, \mathbf{d}) \leq w(\mathbf{a}, \mathbf{b}) + \mu$  and now the inequality follows from Corollary 3.  $\blacksquare$

Throughout the remainder of this section, we consider more sophisticated weight assignments that in many cases improve the upper bound of Theorem 8. Notice first that if  $r = \lambda$ , then the choice of weights  $\mathbf{w}$  from the previous theorem produces the best possible upper bound on  $A_\lambda(n, 3)$  achievable via Theorem 7. This can be seen by noting when  $r = \lambda$ , one may write

$$\tau^*(\mathbf{A}_{(\lambda, \lambda)}) = \min \left\{ 2^n \cdot \sum_{\ell=0}^n \binom{n}{\ell} \cdot w_\ell : w_\ell \cdot (\ell + 1) \geq 1, 0 \leq \ell \leq n, w_\ell \geq 0 \right\}.$$

**Proposition 9** Suppose that  $r = \lambda$ . Then

$$\tau^*(\mathbf{A}_{(r, \lambda)}) = 2^n \cdot \sum_{\ell=0}^n \frac{\binom{n}{\ell}}{\ell + 1} = \frac{2^n(2^{n+1} - 1)}{n + 1}.$$

Hence,  $A_\lambda(n, 2\lambda + 1) \leq 2^n(2^{n+1} - 1)/(n + 1)$ .

**Proposition 10** For integers  $r, \lambda$ , where  $\lambda | r$ ,

$$A_\lambda(n, 2r + 1) \leq 2^n \cdot \sum_{\ell=0}^n \frac{\binom{n}{\ell}}{\sum_{j=0}^{r/\lambda} \binom{\ell}{j}}.$$

We now consider the case where  $r \geq 2$  and  $\lambda = 1$ . To ease the notation, we introduce the function

$$V(n, w, r, \ell) = \sum_{4k+m \leq r-2\ell} \binom{n-w}{k} \binom{w}{m} 2^\ell.$$

We also assume that  $V(n, w, r, \ell) = 0$  if  $r - 2\ell < 0$ . Given this setup, we may write  $\mathbf{A}_{H, (r, 1)}^T = (a_{i,j})_{i=1, j=1}^{n+1}$ , where  $a_{i,j} = V(n, j - 1, r, |i - j|)$ . We produce a weight assignment for (7)

by considering another matrix  $\hat{\mathbf{A}}(n, r)$  related to  $\mathbf{A}_{H, (r, 1)}^T$ . The weights are given as  $\mathbf{w} = \hat{\mathbf{A}}(n, r)^{-1} \cdot \mathbf{1}$ . Theorem 11 states that, indeed,  $\mathbf{w} = \hat{\mathbf{A}}(n, r)^{-1} \cdot \mathbf{1}$  is a feasible weight assignment for (7).

We introduce the matrix  $\hat{\mathbf{A}}(n, r)$  as follows: Let  $\hat{\mathbf{A}}(n, r) = (\hat{a}_{i,j})_{i=1, j=1}^{n+1}$  be defined so that for  $j \neq i$ ,  $\hat{a}_{i,j} = \min(V(n, j - 1, r, |i - j|), \frac{1}{r}(V(n, j - 1, r, 0) - 1))$  and  $\hat{a}_{j,j} = V(n, j - 1, r, 0)$ .

**Theorem 11.** For integers  $n, r$ ,

$$A_1(n, 2r + 1) \leq 2^n \cdot \sum_{\ell=0}^n \binom{n}{\ell} \cdot \hat{w}_\ell,$$

where  $\hat{\mathbf{w}} = (\hat{w}_0, \dots, \hat{w}_{n+1}) = \hat{\mathbf{w}} = \hat{\mathbf{A}}(n, r)^{-1} \cdot \mathbf{1}$ .

The results of our bounds on  $A_1(n, 5)$ , when  $n \leq 20$ , are listed in Table I. As can be seen from the table, for every code length considered, Theorem 11 provides a tighter bound than Theorem 8.

TABLE I  
COMPARISON OF UPPER BOUNDS FOR  $A_1(n, 5)$

| Length | Bound from (7) | Theorem 8   | Proposition 10 | Theorem 11  |
|--------|----------------|-------------|----------------|-------------|
| 5      | 65             | 716         | 254            | 197         |
| 6      | 209            | 2348        | 793            | 589         |
| 7      | 681            | 7545        | 2508           | 1771        |
| 8      | 2285           | 23959       | 8048           | 5396        |
| 9      | 7723           | 75688       | 26190          | 16719       |
| 10     | 27137          | 239112      | 86393          | 52906       |
| 11     | 95480          | 758457      | 288649         | 170584      |
| 12     | 340889         | 2422954     | 975954         | 562157      |
| 13     | 1233644        | 7812585     | 3336118        | 1885717     |
| 14     | 4471386        | 25462344    | 11518362       | 6425947     |
| 15     | 16320256       | 83943512    | 40130869       | 22271529    |
| 16     | 59909131       | 279998120   | 140971957      | 78091743    |
| 17     | 220589555      | 944741909   | 498899141      | 276648820   |
| 18     | 815168373      | 3222862985  | 1777507455     | 991500693   |
| 19     | 3022921187     | 11108080869 | 6371682078     | 3578006784  |
| 20     | 11241799535    | 38650901357 | 22966595378    | 12983261249 |

#### IV. CODE CONSTRUCTIONS

In what follows, we present constructions of linear codes under the ALD for the case where  $\lambda = 1$ . We first address the case where the minimum ALD is equal to three.

For a positive integer  $v$ , let  $H'_3 \in \mathbb{F}_2^{v \times (2^v - 2)}$  be a matrix which has as its columns all non-zero vectors from  $\mathbb{F}_2^v$ , except for the all-ones vector. Write  $H'_3 = (\mathbf{h}'_1, \mathbf{h}'_2, \dots, \mathbf{h}'_{2^v - 2})$ , where for  $i \in \{1, \dots, 2^v - 2\}$ ,  $\mathbf{h}'_i$  represents the  $i$ -th column of  $H'_3$ . Furthermore, let  $\mathbf{1}_v \in \mathbb{F}_2^{v \times 1}$  be the all-ones vector.

In this setting, let  $\mathcal{C}(2^v - 2, 3) \subseteq \mathbb{Z}_2^{2^v - 2} \times \mathbb{Z}_2^{2^v - 2}$  be equal to

$$\mathcal{C}(2^v - 2, 3) := \left\{ (\mathbf{a}, \mathbf{b}) \in \mathbb{F}_2^{2^v - 2} \times \mathbb{F}_2^{2^v - 2} : \sum_{i=1}^{2^v - 2} a_i \cdot \mathbf{h}'_i + \sum_{i=1}^{2^v - 2} b_i \cdot \mathbf{1}_v = \mathbf{0} \right\}.$$

Note that a code with minimum ALD distance three (for  $\lambda = 1$ ) can either:

- 1) Correct a single Class 1 error, or
- 2) Detect a single Class 2 error.

In the following lemma, we show that the code  $\mathcal{C}(2^v - 2, 3)$  can perform either 1) or 2).

**Lemma 12.** For any positive integer  $n$ ,  $d_1(\mathcal{C}(2^v - 2, 3)) \geq 3$ .

*Proof:* Let  $n = 2^v - 2$ . Suppose that  $(\mathbf{a}, \mathbf{b}) \in \mathcal{C}(n, 3)$  was transmitted and that the vector  $(\mathbf{c}, \mathbf{d}) \in \mathbb{F}_2^n \times \mathbb{F}_2^n$  was received, where  $(\mathbf{c}, \mathbf{d}) = ((c_1, d_1), (c_2, d_2), \dots, (c_n, d_n))$  is the result of at most a single Class 2 error occurring in  $(\mathbf{a}, \mathbf{b})$ . For the remainder of this proof, let  $\mathbf{s} = \sum_{i=1}^n c_i \cdot \mathbf{h}'_i + \sum_{i=1}^n d_i \cdot \mathbf{1}_v$ . If a single Class 2 error has occurred in position  $j \in \{1, \dots, n\}$  of  $\mathbf{a}$ , then  $\mathbf{s} = \mathbf{h}'_j \neq \mathbf{0}$ . Otherwise, if the Class 2 error occurred in position  $j$  of  $\mathbf{b}$ , then  $\mathbf{s} = \mathbf{1}_v \neq \mathbf{0}$  holds as well. Clearly, if no Class 2 errors occurred, we have  $\mathbf{s} = \mathbf{0}$ . Thus,  $\mathcal{C}(n, 3)$  can detect whether a single Class 2 error has occurred by checking if  $\mathbf{s}$  is non-zero.

Suppose that  $(\mathbf{a}, \mathbf{b}) \in \mathcal{C}(n, 3)$  was transmitted and that the vector  $(\mathbf{c}, \mathbf{d}) \in \mathbb{F}_2^n \times \mathbb{F}_2^n$  was received, where  $(\mathbf{c}, \mathbf{d}) = ((c_1, d_1), (c_2, d_2), \dots, (c_n, d_n))$  is the result of at most one Class 1 error in  $(\mathbf{a}, \mathbf{b})$ . We describe next how to recover  $(\mathbf{a}, \mathbf{b})$  from  $(\mathbf{c}, \mathbf{d})$ . Let  $\mathbf{s}_2 = \mathbf{s} + \mathbf{1}_v$ . Note that if  $(\mathbf{c}, \mathbf{d})$  is the result of a single Class 1 error in  $(\mathbf{a}, \mathbf{b})$ , then

$$\mathbf{s}_2 = \left( \sum_{i=1}^n c_i \cdot \mathbf{h}'_i + \sum_{i=1}^n d_i \cdot \mathbf{1}_v \right) + \mathbf{1}_v = \mathbf{h}'_j + \mathbf{1}_v + \mathbf{1}_v = \mathbf{h}'_j,$$

where  $j \in \{1, \dots, n\}$  is the position of the error in  $(\mathbf{a}, \mathbf{b})$ . Otherwise, if no error occurred,  $\mathbf{s}_2 = \mathbf{s} + \mathbf{1}_v = \mathbf{1}_v$ . From the above discussion (recall  $\mathbf{h}'_j \neq \mathbf{1}_v$  by design) it is clear that a  $\mathcal{C}(n, 3)$  decoder can recover  $(\mathbf{a}, \mathbf{b})$  from  $(\mathbf{c}, \mathbf{d})$  from  $\mathbf{s}_2$  as follows. If  $\mathbf{s}_2 = \mathbf{1}_v$ , then the decoder concludes that no errors have occurred. Otherwise, if  $\mathbf{s}_2 = \mathbf{h}'_j$  for some  $j \in \{1, \dots, n\}$ , then the decoder corrects a Class 1 error at position  $j$ . ■

As a consequence of Lemma 12 and Proposition 9, we have

$$\frac{4^n}{n+2} \leq A_1(n, 3) \leq \frac{2^n(2^{n+1} - 1)}{n+1}.$$

We now turn our attention to the problem of constructing codes with minimum ALD equal to  $d$ . We first describe the code construction, and then proceed to provide a proof of its correctness. Let  $H'_d \in \mathbb{F}_2^{s \times 2n}$  be a parity check matrix for a code  $\mathcal{C}$  with Hamming distance  $d$ . Write  $H'_d = (\mathbf{h}'_1, \dots, \mathbf{h}'_{2n})$ , where, as before,  $\mathbf{h}'_i$  denotes the  $i$ -th column of  $H'_d$ . Let  $\mathcal{C}(n, d) \subseteq \mathbb{Z}_2^n \times \mathbb{Z}_2^n$  be defined as

$$\mathcal{C}(n, d) := \left\{ (\mathbf{a}, \mathbf{b}) = ((a_1, b_1), \dots, (a_n, b_n)) \in \mathbb{F}_2^n \times \mathbb{F}_2^n : \sum_{i=1}^n a_i \cdot \mathbf{h}'_i + \sum_{i=1}^n b_i \cdot (\mathbf{h}'_i + \mathbf{h}'_{n+i}) = \mathbf{0} \right\}.$$

**Lemma 13.** For a positive integer  $n$ ,  $d_1(\mathcal{C}(n, d)) \geq d$ .

*Proof:* Since the code  $\mathcal{C}(n, d)$  is linear, we only need to show that for any  $(\mathbf{c}, \mathbf{d}) \in \mathbb{F}_2^n \times \mathbb{F}_2^n$ , we have  $\mathbf{s} = \sum_{i=1}^n c_i \cdot \mathbf{h}'_i + \sum_{i=1}^n d_i \cdot (\mathbf{h}'_i + \mathbf{h}'_{n+i}) \neq \mathbf{0}$  when  $d_1((\mathbf{0}, \mathbf{0}); (\mathbf{c}, \mathbf{d})) < d$ . Suppose, in particular, that  $(\mathbf{c}, \mathbf{d})$  is the result of  $m$  Class 1 errors,  $\ell$  Class 2 errors, and  $k$  Class 3 errors, where  $m + 2\ell + 4k < d$ . Then,  $\mathbf{s}$  is the sum of at most  $m + 2\ell + k$  columns of  $H'_d$ . Since  $m + 2\ell + 4k < d$ , then  $m + 2\ell + k < d$ , and since  $H'_d$  is a parity-check matrix of a code with Hamming distance  $d$ , we conclude that  $\mathbf{s} \neq \mathbf{0}$ . ■

We now comment on the size of the code  $\mathcal{C}(n, d)$  for the case when  $r = 2$  or  $d = 5$ . We compare the ALD code  $\mathcal{C}(n, 5)$  with

codes constructed directly for the Hamming distance, assumed to be large enough. In particular, we consider (a) binary codes of length  $2n$ ; and (b) the binary image of quaternary codes of length  $n$ .

For (a), since a code with minimum ALD equal to five must be able to correct four errors (resulting from two Class 1 errors), we require a binary code  $\mathcal{C}_2$  that has minimum Hamming distance at least nine. By definition,  $|\mathcal{C}(n, 5)|$  is strictly larger than  $|\mathcal{C}_2|$ , and the direct method offers worse code rates than our construction. As an example, for  $2n = 2^v - 2$ , where  $v \geq 2$  is a positive integer,  $|\mathcal{C}(n, 5)| \geq \frac{4^n}{(2n+2)^2}$ , if shortened binary BCH codes are used as defining codes. Proposition 10 for the case  $\lambda = 1$  and  $r = 2$ , and the lower bound on  $|\mathcal{C}(n, 5)|$ , imply that for  $2n = 2^v - 2$ , and integer  $v > 2$ ,

$$\frac{4^n}{(2n+2)^2} \leq A_1(n, 5) \leq \frac{3 \cdot 2^n}{(n+1)(n+2)} (2^{n+2} - n - 3).$$

Next, consider the case (b). For a quaternary code  $\mathcal{C}_4$  to have minimum ALD distance five, the code should have Hamming distance at least five. Writing the sphere packing bound for quaternary codes with Hamming distance five, we have that  $|\mathcal{C}_4| \leq \frac{4^n}{\sum_{j=0}^4 \binom{n}{j} 3^j}$ . This value is strictly smaller than  $\frac{4^n}{(2n+2)^2}$ , which is a lower bound for the value of  $|\mathcal{C}(n, 5)|$ , i.e.,

$$|\mathcal{C}_4| \leq \frac{4^n}{\sum_{j=0}^2 \binom{n}{j} 3^j} < \frac{4^n}{(2n+2)^2} \leq |\mathcal{C}(n, 5)|,$$

where  $v \geq 5$ . Hence, our construction outperforms the direct approach for the case of a quaternary alphabet as well.

**Acknowledgment:** This work was funded by the Strategic Research Initiative (SRI) program at the University of Illinois, a CIA Postdoctoral Fellowship, by the NSF STC on Science of Information, and by the NISE program at SSC Pacific. This work was completed when H. M. Kiah was at University of Illinois, Urbana-Champaign.

#### REFERENCES

- [1] T. Alderson and S. Huntemann, "On maximum Lee distance codes," *J. of Discrete Mathematics*, 2013.
- [2] J. Astola, "The theory of Lee-codes," *Lappeenranta University of Technology, Department of Physics and Mathematics, Research Report*, Jan. 1982.
- [3] G. M. Church, Y. Gao, and S. Kosuri, "Next-generation digital information storage in DNA," *Science*, vol. 337, no. 6102, pp. 1628-1628, Sep. 2012.
- [4] A. Fazeli, A. Vardy, E. Yaakobi, "Generalized Sphere Packing Bound," available at <http://arxiv.org/abs/1401.6496>, 2014.
- [5] R. Feynman, "There's plenty of room at the bottom," Caltech, Pasadena. 29 Dec. 1959. Lecture.
- [6] N. Goldman et al., "Towards practical, high-capacity, low-maintenance information storage in synthesized DNA," *Nature*, Jan. 2013.
- [7] E. Hof et al., "Capacity-achieving polar codes for arbitrarily permuted parallel channels," *IEEE Trans. on Info. Theory*, vol. 59, no. 3, pp. 1505-1516, March 2013.
- [8] K.A.S. Immink, *Codes for mass data storage systems*, Shannon Foundation Publisher, 2004.
- [9] M. Kaykobad, "Positive solutions of positive linear systems," *Lin. Alg. and its App.*, vol. 64, pp. 133-140, Jan. 1985.
- [10] A.A. Kulkarni, and N. Kiyavash, "Nonasymptotic upper bounds for deletion correcting codes," *IEEE Trans. on Info. Theory*, vol. 59, no. 8, pp. 5115-5130, April 2013.
- [11] A. Mazumdar, A. Barg, and N. Kashyap, "Coding for high-density recording on a 1-D granular magnetic medium," *IEEE Trans. on Info. Theory*, vol. 57, no. 11, pp. 7403-7417, June 2011.
- [12] O. Milenkovic and N. Kashyap, "On the design of codes for DNA computing," *Coding and Cryptography*, Springer Berlin Heidelberg, pp. 100-119, 2006.
- [13] K. Nakamura, et al., "Sequence-specific error profile of Illumina sequencers," *Nucleic acids research*, voll. 39, no. 13, 2011.
- [14] T. Richardson and R. Urbanke, "The capacity of low-density parity-check codes under message-passing decoding," *IEEE Trans. on Info. Theory*, vol. 47, no. 2, pp. 599-618, August 2002.
- [15] Roth, R., *Introduction to coding theory*, Cambridge University Press, 2006.
- [16] Tal, I., Vardy, A., "How to construct polar codes," *IEEE Trans. on Info. Theory*, vol. 59, no. 10, pp. 6562 - 6582, September 2013.