

Studying the Cycle Complexity of DNA Synthesis

Amit Zrihan*, Eitan Yaakobi*, and Zohar Yakhini*†

*Faculty of Computer Science, Technion - Israel Institute of Technology, Israel

†School of Computer Science, RUNI, Israel

{amit.zr@campus.technion.ac.il, yaakobi@cs.technion.ac.il, zohar.yakhini@gmail.com}

Abstract—Storing data in DNA is being explored as an efficient solution for archiving and in-object storage. Synthesis time and cost remain challenging, significantly limiting some applications at this stage. In this paper we investigate efficient synthesis, as it relates to cyclic synchronized synthesis technologies, such as photolithography. We define performance metrics related to the number of cycles needed for the synthesis of any fixed number of bits. We first expand on some results from the literature related to the channel capacity, addressing densities beyond those covered by prior work. This leads us to develop effective encoding achieving rate and capacity that are higher than previously reported. Finally, we analyze cost based on a parametric definition and determine some bounds and asymptotics. We investigate alphabet sizes that can be larger than 4, both for theoretical completeness and since practical approaches to such schemes were recently suggested and tested in the literature.

I. INTRODUCTION

Digital data can be stored in DNA oligonucleotides thanks to the composition of DNA molecules, which consist of four nucleotides: Adenine (A), Cytosine (C), Guanine (G), and Thymine (T). An oligonucleotide (oligo), represents a sequence of these nucleotides in a specific order, and can be denoted as a string over the alphabet $\{A, C, G, T\}$. As it is feasible to chemically synthesize nearly any nucleotide sequence, DNA oligos can serve to store digital data.

The initial phase of DNA data storage involves converting binary data into DNA oligos. This process involves adding redundant bases for indexing and using coding techniques to reduce errors and losses that commonly occur during synthesis and sequencing. The next step consists of synthesis. In this step each oligo in the storage system is typically synthesized thousands to millions of times, resulting in a vast collection of unordered DNA fragments. Storing the DNA which now codes the information may require the use of technology to guarantee its long time integrity [11]. During the reading process, a subset of these fragments is then selected and read using sequencing methods. To reconstruct the stored sequences, the resulting reads are grouped into clusters that are likely originated from the same design oligo. The original designed oligos are then reconstructed using the reads (see e.g [21], [22]). The original data is then retrieved from the reconstructed DNA sequences [4], [6], [10], [12].

Photolithographic DNA synthesis [2], [3], [20] is a technique used in the synthesis step. Photolithography can also be used in other synthesis applications, including for direct synthesis of RNA [17]. In the context of DNA for data storage, the process results in a large number of DNA oligonucleotides, all synthesized in parallel. Under the standard approach, each oligo is grown by at most one nucleotide in each synthesis

cycle. The machine follows a fixed supersequence of possible nucleotides to append to the oligos. As the machine iterates through this supersequence, the next nucleotide is added to a selected subset of the oligos until the machine reaches the end of the supersequence. In particular, each synthesized DNA oligo must be a subsequence of the machine's supersequence. In this paper, we operate under the assumption that there are no constraints on the number of oligos that can be synthesized simultaneously for any practical run of the machine.

Several approaches were suggested and implemented, however, that support larger alphabet [1], [5], [7]–[9], [19], [23]. Throughout the paper, we will work with the abstract alphabet $\Sigma_q \triangleq \{1, 2, \dots, q\}$, where each $\sigma \in \Sigma_q$ is referred to as a symbol. As the synthesis machine's supersequence we will use the alternating sequence A_q , that cyclically repeats all symbols in Σ_q in ascending order. For an alphabet of size q , the alternating sequence is $A_q = (123 \dots q123 \dots q \dots)$. For a given number of synthesis cycles C we will use $A_q[C]$ to denote the length- C prefix of A_q as the supersequence. When synthesizing a predetermined length of oligonucleotides, let L denote that length, and denote also $\rho \triangleq \frac{L}{C}$ ($0 \leq \rho \leq 1$). One of our main objectives in this paper is to minimize the number of synthesis cycles that is required in order to store some N information bits. Hence, our aim is to maximize, for a single oligo, the ratio between the number of information bits and the number of synthesis cycles. This is in the spirit of the term *logical density* used in some of the literature. Another main objective to this work is to investigate the cost of synthesis under various assumptions. These questions will be precisely defined in the next section. The rest of the paper is structured as follows. In Section II, we provide some technical preliminaries and definitions. In Section III, we address encoders that attain previously proven bounds [16]. In Section IV, we analyze the synthesis cost. Finally, a discussion of limitations and future work is given in Section V. Due to space constraints, some of the proofs are omitted.

II. PRELIMINARIES AND PROBLEM FORMULATION

A. Information Theory Limits of Photolithographic Synthesis

The *radius- t deletion sphere* (exactly t deletions) of a string s is denoted by $D(s, t)$, and $D_q(C, t) \triangleq |D(A_q[C], t)|$. As indicated in [13] the recursive formula $D_q(C, t) = \sum_{i=0}^t \binom{C-t}{i} D_{q-1}(t, t-i)$ holds. Let $M_q(C, L)$ be the number of distinct subsequences of length L of $A_q[C]$, and note that $M_q(C, L) = D_q(C, C-L)$.

Lemma 1. For any $q, C, L \in \mathbb{N}$, ($L \leq C$) it holds that

$$M_q(C, L) \leq M_{q+1}(C, L) \leq \binom{C}{L}, \quad (1)$$

where the inequality is strict for $2 \leq q < C$, $2 \leq C$, $0 < L$.

Proof. The first inequality can be shown by induction on q and the recursive formula of $D_q(C, t)$. For the second inequality,

This research was funded by the European Union (DiDAX, 101115134). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

$\binom{C}{L}$ represents the number of ways to choose L indices out of string of length C . However, since these options may not necessarily generate distinct strings, the inequality arises. \square

Lemma 2. It holds that $\binom{v}{u}^n \leq M_v(n \cdot v, n \cdot u)$, ($u, v, n \in \mathbb{N}$).

Proof. Consider a scenario where for each window of length v we select u symbols for the subsequence, yielding a subsequence of length $n \cdot u$. Every distinct choice made in every window results in a distinct subsequence. Consequently, the total number of ways to choose u symbols out of v , n times is $\binom{v}{u}^n$. However, since this is represent only a subset of all the possible subsequences, the inequality holds. \square

We define, analogous to [16]

$$\begin{aligned} \text{cap}(q, \rho) &\triangleq \limsup_{C \rightarrow \infty} \frac{\log_2 M_q(C, \lfloor \rho \cdot C \rfloor)}{C}, \\ \text{cap}(q) &\triangleq \limsup_{C \rightarrow \infty} \frac{\log_2 \left(\sum_{L=0}^C M_q(C, L) \right)}{C}. \end{aligned}$$

These values represent the maximum number of information bits that can be stored per synthesis cycle for a fixed oligonucleotide length $L = \lfloor \rho \cdot C \rfloor$ and for flexible L , respectively. All when synthesizing a single oligo under A_q as the supersequence of length C .

Theorem 1. It holds that

$$\text{cap}(q, \rho) \leq \text{cap}(q+1, \rho) \leq \lim_{q \rightarrow \infty} \text{cap}(q, \rho) = H(\rho) \quad (2)$$

$$\text{cap}(q) \leq \text{cap}(q+1) \leq \lim_{q \rightarrow \infty} \text{cap}(q) = 1 \quad (3)$$

where H is the binary entropy function.

Proof. As for (2), both inequalities follows from Lemma 1. As for the the limit, we start from an upper bound

$$\begin{aligned} \lim_{q \rightarrow \infty} \text{cap}(q, \rho) &= \lim_{q \rightarrow \infty} \limsup_{C \rightarrow \infty} \frac{\log_2 M_q(C, \lfloor \rho \cdot C \rfloor)}{C} \\ &\stackrel{\text{Lemma 1}}{\leq} \lim_{q \rightarrow \infty} \limsup_{C \rightarrow \infty} \frac{\log_2 \binom{C}{\lfloor \rho \cdot C \rfloor}}{C} = H(\rho). \end{aligned}$$

We showed that the sequence $\text{cap}(q, \rho)$ is monotonously increasing and bounded from above therefore it converges. As for the lower bound, let $\rho = \frac{s}{t} \in \mathbb{Q}$. It holds that

$$\begin{aligned} \lim_{q \rightarrow \infty} \text{cap}(q, \rho) &= \lim_{n \rightarrow \infty} \text{cap}(n \cdot t, \rho) \\ &= \lim_{n \rightarrow \infty} \limsup_{C \rightarrow \infty} \frac{\log_2 M_{n \cdot t}(C, \lfloor \rho \cdot C \rfloor)}{C} \\ &\geq \lim_{n \rightarrow \infty} \lim_{m \rightarrow \infty} \frac{\log_2 M_{n \cdot t}(m \cdot nt, m \cdot ns)}{m \cdot nt} \\ &\stackrel{\text{Lemma 2}}{\geq} \lim_{n \rightarrow \infty} \lim_{m \rightarrow \infty} \frac{\log_2 \binom{nt}{ns}^m}{m \cdot nt} \\ &= \lim_{n \rightarrow \infty} \frac{\log_2 \binom{nt}{ns}}{nt} = \lim_{n \rightarrow \infty} \frac{\log_2 \binom{nt}{\rho \cdot nt}}{nt} = H(\rho). \end{aligned}$$

Note that the above proof formally holds for $\rho \in \mathbb{Q}$. As for (3) it holds that

$$\begin{aligned} 1 &= H\left(\frac{1}{2}\right) = \lim_{q \rightarrow \infty} \text{cap}\left(q, \frac{1}{2}\right) \leq \lim_{q \rightarrow \infty} \text{cap}(q) \\ &\stackrel{\text{Lemma 1}}{\leq} \limsup_{C \rightarrow \infty} \frac{\log_2 \left(\sum_{L=0}^C \binom{C}{L} \right)}{C} = \limsup_{C \rightarrow \infty} \frac{\log_2 (2^C)}{C} = 1. \end{aligned}$$

From [16], it is known that \square

$$\text{cap}(q, \rho) = \begin{cases} \rho \log_2 q, & 0 \leq \rho \leq \frac{2}{q+1} \\ \rho \log_2 \left(\sum_{i=1}^q x_q(\rho)^{i-\frac{1}{\rho}} \right), & \frac{2}{q+1} < \rho < 1, \end{cases} \quad (4)$$

where $x_q(\rho)$ is the unique solution to $\sum_{i=1}^q (1-\rho i)x^i = 0$ in the range $0 < x < 1$. Furthermore, as per [15], for any q we have $\text{cap}(q) = -\log_2 x_q$, where x_q , ($0 < x_q < 1$) is the unique real solution of the polynomial $\sum_{i=1}^q x^i = 1$. Fig. 1 shows $\text{cap}(q, \rho)$ for some values of q in relation to $H(\rho)$. Note that as q grows, capacity is maximized for ρ around $\frac{1}{2}$.

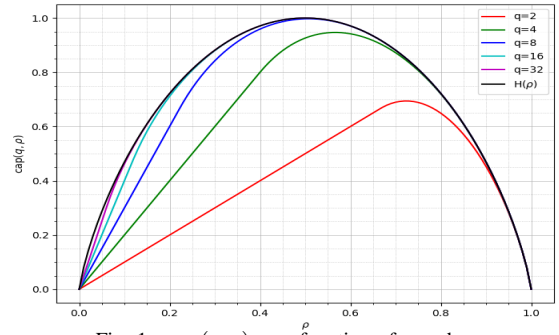


Fig. 1. $\text{cap}(q, \rho)$ as a function of q and ρ .

B. Capacity Achieving Encoder

Having explored the information-theoretical limits of photographic synthesis, our focus shifts to formulating an efficient encoder-decoder scheme capable of achieving the capacity. For the special case of $\rho = \frac{2}{q+1}$, [15] presents an efficient encoder-decoder using a single redundancy symbol. The time complexity of the encoder is linear and it encodes any string of length $L-1$ over an alphabet of size q to a length- L string over the same alphabet, such that its synthesis time is at most $C = \lfloor \frac{q+1}{2} \cdot L \rfloor$. We demonstrate that this encoder's rate achieves the capacity. The rate R of this encoder satisfies

$$\begin{aligned} R &= \lim_{L \rightarrow \infty} \frac{\log_2 q^{L-1}}{\lfloor \frac{q+1}{2} \cdot L \rfloor} = \lim_{L \rightarrow \infty} \frac{(L-1) \cdot \log_2 q}{\lfloor \frac{q+1}{2} \cdot L \rfloor} \\ &= \frac{2}{q+1} \cdot \log_2 q = \rho \cdot \log_2 q = \text{cap}\left(q, \rho = \frac{2}{q+1}\right). \end{aligned}$$

For any other value of ρ we want to achieve the best possible rate using a variation of this encoder. This problem is stated formally as follows.

Problem 1. Given $0 < \rho < 1$ and q as the alphabet size to be used, our objective is to design a supersequence with an efficient encoder-decoder that achieves $\text{cap}(q, \rho)$.

Problem 1 is addressed in Section III.

C. Optimizing The Synthesis Cost

In this section we define a *cost function* for photolithographic synthesis of symbols from Σ_q . Let $\alpha, \beta \in \mathbb{R}^+$ be the cost of cycle synthesis and symbol synthesis, respectively. To derive a general cost function for storing N bits using C cycles, s.t. $\rho = \frac{L}{C}$ we are using the expression

$$\text{cost}(N, C, q, \rho) = \alpha C + \beta \cdot (\#\text{synthesized symbols}).$$

Further, assume that we can achieve the capacity $\text{cap}(q, \rho)$ for any q, ρ , and so we can thus store $C \cdot \text{cap}(q, \rho) = \frac{L}{\rho} \cdot \text{cap}(q, \rho)$ bits per oligo of length $L = \rho C$. We determine the number of oligos to be synthesized to be $\frac{N}{\frac{L}{\rho} \cdot \text{cap}(q, \rho)}$ which implies that $\frac{\rho \cdot N}{\text{cap}(q, \rho)}$ symbols are synthesized. Note that we should use the floor and ceil functions, but for the sake of simplicity, we have omitted them. We therefore see that, under our assumptions, the optimal achievable cost is given by the formula

$$\text{cost}^*(N, C, q, \rho) \triangleq \alpha C + \beta N \cdot \frac{\rho}{\text{cap}(q, \rho)}. \quad (5)$$

Problem 2. Our objective is to find the values of

- 1) $\arg \min_{\rho \in (0,1)} \text{cost}^*(N, C, q, \rho)$, for fixed N, C, q .
- 2) $\arg \min_{\rho \in (0,1), q \in \mathbb{N}} \text{cost}^*(N, C, q, \rho)$, for fixed N, C .

Remark. Changes in ρ will also impact the oligo length L , consequently altering the number of oligos to be synthesized. This will also affect the barcode lengths required for indexing the oligos, influencing the overall cost. While we acknowledge this effect, we do not address its specific implications.

III. CAPACITY ACHIEVING ENCODER

While the capacity of photolithographic synthesis was calculated in prior work, an efficient encoder-decoder scheme that achieves the capacity, was only described, to our knowledge, for $\rho = \frac{2}{q+1}$. We now address Problem 1, by introducing two families of efficient encoders applicable for a wider range of ρ values, and achieving higher information rate.

A. Lookup Table

We start with a naive encoding approach by using a lookup table. Let $d \in \mathbb{N}, \rho \in (0,1)$ s.t. $\rho \cdot d \cdot q \in \mathbb{N}$ and denote $L = \rho \cdot d \cdot q, C = d \cdot q$. For $B = \lfloor \log_2(M_q(C, L)) \rfloor$, we pre-calculate a lookup table of 2^B entries, while each entry corresponds to a distinct subsequence of length L out of $A_q[d \cdot q]$. The encoding process iteratively encodes every B bits into a string of length L over $A_q[d \cdot q]$ using the lookup table. We attain that the rate R of this encoder satisfies $R = \frac{\lfloor \log_2(M_q(dq, \rho \cdot dq)) \rfloor}{dq}$. Note that $\lim_{d \rightarrow \infty} \frac{\lfloor \log_2(M_q(dq, \rho \cdot dq)) \rfloor}{dq} = \text{cap}(q, \rho)$. Two major drawbacks of this approach is that the lookup table requires a substantial use of memory, and its calculation is time consuming. Fig. 2 depicts the rate of this encoder using a lookup table of up to 2^{32} entries alongside with the capacity.

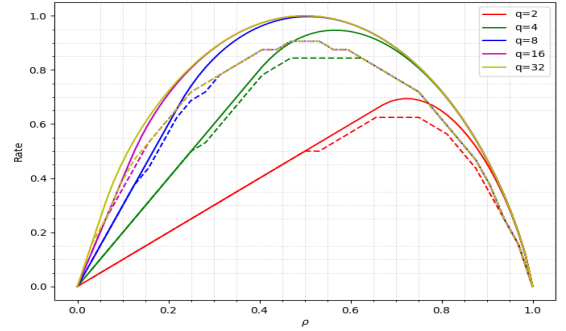


Fig. 2. Information rates achievable by the construction of Section III-A, compared to the related capacities $\text{cap}(q, \rho)$, for selected values of q . The dashed line represents the maximum achievable rate of the encoder, while the solid line represents the capacity. The value of d for every pair of (q, ρ) was chosen to be the largest integer s.t. $B \leq 32$.

B. Multi-size Alphabet Encoder

In this section we present a variation of the encoder presented in [15] (described above) where we use multiple alphabet sizes. Initially, we analyze a scenario where a fraction of the synthesized oligos is encoded over an alphabet of size q_1 , and the remainder is encoded over an alphabet of size q_2 . Using C_1 cycles over the alternating sequence A_{q_1} and C_2 cycles over A_{q_2} produces the machine's supersequence $A_{q_1}[C_1] \circ A_{q_2}[C_2]$, which concatenates the first C_1 symbols of A_{q_1} with the first C_2 symbols of A_{q_2} . First we consider a case where half of the synthesized oligos is over q_1 and the other half is over q_2 . In this case, to synthesize oligos of length L it holds that $C_1 = \frac{q_1+1}{2} \cdot \frac{L}{2}, C_2 = \frac{q_2+1}{2} \cdot \frac{L}{2}, C = C_1 + C_2, \rho = \frac{L}{C} = \frac{4}{q_1+q_2+2}$. The rate R of this encoder satisfies

$$\begin{aligned} R &= \lim_{L \rightarrow \infty} \frac{\log_2 \left(q_1^{\frac{L}{2}-1} \cdot q_2^{\frac{L}{2}-1} \right)}{\frac{q_1+1}{2} \cdot \frac{L}{2} + \frac{q_2+1}{2} \cdot \frac{L}{2}} = \lim_{L \rightarrow \infty} \frac{(\frac{L}{2}-1) \log_2 q_1 + (\frac{L}{2}-1) \log_2 q_2}{\frac{L}{4} \cdot (q_1 + q_2 + 2)} \\ &= \frac{4 \cdot (\frac{1}{2} \log_2 q_1 + \frac{1}{2} \log_2 q_2)}{q_1 + q_2 + 2} = \rho \cdot \left(\frac{1}{2} \log_2 q_1 + \frac{1}{2} \log_2 q_2 \right). \end{aligned}$$

A more general approach is when using all q alphabet sizes $1, 2, \dots, q$, where each alphabet size i is used on a fraction α_i of the synthesized string ($\sum_{i=1}^q \alpha_i = 1$). Similarly to the previous case we get that $C = \sum_{i=1}^q C_i = \sum_{i=1}^q \frac{i+1}{2} \cdot \alpha_i L, \rho = \frac{L}{C} = \frac{2}{1 + \sum_{i=1}^q i \alpha_i}$. The rate R of this encoder satisfies

$$\begin{aligned} R &= \lim_{L \rightarrow \infty} \frac{\log_2 \left(\prod_{i=1}^q i^{\alpha_i L - 1} \right)}{\sum_{i=1}^q \frac{i+1}{2} \cdot \alpha_i L} = \lim_{L \rightarrow \infty} \frac{\sum_{i=1}^q (\alpha_i L - 1) \log_2 i}{\sum_{i=1}^q \frac{i+1}{2} \cdot \alpha_i L} \\ &= \frac{2 \sum_{i=1}^q \alpha_i \log_2 i}{1 + \sum_{i=1}^q i \alpha_i} = \rho \cdot \sum_{i=1}^q \alpha_i \log_2 i. \end{aligned}$$

Consider a fixed $\rho \in (0,1)$, our goal is to maximize the rate R that can be achieved for this ρ . We also assume a fixed q but we can control $\vec{\alpha} \triangleq (\alpha_1, \alpha_2, \dots, \alpha_q)$. Also recall that relevant values of $\vec{\alpha}$ satisfy $\sum_{i=1}^q i \alpha_i = \frac{2}{\rho} - 1$.

This yields the following linear programming (LP) system.

$$\max_{\vec{\alpha}} \{ \mathbf{c}^T \vec{\alpha} \} \text{ s.t. } \mathbf{A} \vec{\alpha} = \mathbf{b}, \vec{\alpha} \geq 0$$

$$\begin{aligned} \mathbf{A} &= \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 2 & \dots & q \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 1 \\ \frac{2}{\rho} - 1 \end{pmatrix}, \\ \mathbf{c}^T &= (\log_2 1, \log_2 2, \dots, \log_2 q). \end{aligned}$$

Denoting by $\vec{\alpha}^*$ an optimal solution leads us to the following.

Claim 1. *The optimal solution $\bar{\alpha}^*$ has no two positive entries, with indices that are more than one index apart (i.e. no index i so that $\alpha_i > 0$ and so that for some $k \geq 2$ also $\alpha_{i+k} > 0$).*

Proof. Assume by contradiction that $\alpha_i, \alpha_{i+k} > 0$ for some $i \in \{1, 2, \dots, q-2\}$ and $k \geq 2$. Denote $\alpha_i = x$, $\alpha_{i+k} = y$, $\alpha_{i+1} = u$, $\alpha_{i+k-1} = v$. Assume $y \leq x$ (the opposite case is similar). Define $\bar{\alpha}_{new} = (\alpha'_1, \alpha'_2, \dots, \alpha'_q)$ by $\alpha'_i = x - y$, $\alpha'_{i+1} = u + y$, $\alpha'_{i+k-1} = v + y$, $\alpha'_{i+k} = 0$, and $\forall j \neq i, i+1, i+k-1, i+k : \alpha'_j = \alpha_j$. Considering only the changed indices $i, i+1, i+k-1, i+k$, it can be shown that the constraints still hold for $\bar{\alpha}_{new}$. As for the target function, let $r = x \log_2(i) + u \log_2(i+1) + v \log_2(i+k-1) + y \log_2(i+k)$ and $r_{new} = (x-y) \log_2(i) + (u+y) \log_2(i+1) + (v+y) \log_2(i+k-1)$. We observe that the difference $r_{new} - r$ is

$$\begin{aligned} & y(\log_2(i+1) + \log_2(i+k-1) - \log_2(i) - \log_2(i+k)) \\ &= y \log_2 \frac{(i+1)(i+k-1)}{i(i+k)} = y \log_2 \frac{i^2 + ik + k - 1}{i^2 + ik} \\ &= y \log_2 \left(1 + \frac{k-1}{i^2 + ik}\right), \end{aligned}$$

which is positive for all $i \geq 1, k \geq 2$. Therefore $r_{new} > r$ contradicting the optimality of $\bar{\alpha}^*$. \square

Corollary 1. *For the optimal solution $\bar{\alpha}^*$ there is an index $i \in \{1, 2, \dots, q-1\}$ such that for all $j \neq i, i+1$, $\alpha_j^* = 0$.*

Theorem 2. *Let $s \in \{1, 2, \dots, q-1\}$ and $\rho \in [\frac{2}{(s+1)+1}, \frac{2}{s+1}]$. The optimal rate R^* for this configuration, obtained by the approach described above, is*

$$R^* = (\rho(s+2) - 2) \log_2(s) + (2 - \rho(s+1)) \log_2(s+1).$$

Proof. Given that $\rho \in [\frac{2}{(s+1)+1}, \frac{2}{s+1}]$ it holds that $\frac{2}{\rho} - 1 \in [s, s+1]$. Let i be the index of $\bar{\alpha}^*$ from corollary 1. From the first row of the matrix, it holds that $i\alpha_i^* + (i+1)\alpha_{i+1}^* \in [i, i+1]$. From the second row of the matrix it holds that $i\alpha_i^* + (i+1)\alpha_{i+1}^* \in [s, s+1]$. Therefore $i = s$.

Since $\alpha_{s+1}^* = 1 - \alpha_s^*$ we get $s\alpha_s^* + (s+1)(1 - \alpha_s^*) = \frac{2}{\rho} - 1$, yielding $\alpha_s^* = s + 2 - \frac{2}{\rho}$. Substituting into the rate function gives

$$\begin{aligned} R^* &= \rho \cdot (\alpha_s^* \log_2(s) + (1 - \alpha_s^*) \log_2(s+1)) \\ &= \rho \cdot \left((s + 2 - \frac{2}{\rho}) \log_2(s) + (\frac{2}{\rho} - s - 1) \log_2(s+1) \right) \\ &= (\rho(s+2) - 2) \log_2(s) + (2 - \rho(s+1)) \log_2(s+1). \end{aligned}$$

C. Effective Encoder for $\rho = \frac{1}{2}$

In this section we address the specific case $\rho = \frac{1}{2}$. In particular, we introduce an effective encoder-decoder for ρ around $\frac{1}{2}$ applicable to all alphabet sizes. Let q be the alphabet size to be used, and let $f_q \triangleq \max\{n \in \mathbb{N} : n + \lceil \log_2 n \rceil + 1 \leq q\}$. In our encoder we will use alphabet of size $q^* = f_q + \lceil \log_2(f_q) \rceil + 1$. Note that $q^* \leq q$ so this is possible. Our encoder achieves $\rho = \frac{\lfloor \frac{q^*}{2} \rfloor}{q^*}$ and rate $R = \frac{f_q}{q^*}$.

For $u \in \{0, 1\}^*$ let $W_H(u)$ be the Hamming weight of u . We define $I(u) \triangleq \{i : u_i = 1\}$, e.g. $I(0110) = \{2, 3\}$. Let s be a string of the same length as u , we define $\pi(s, u)$, the projection of s on u , to be obtained by keeping the indices $I(u)$ of s and deleting the rest. $\pi(s, u) \triangleq t_1 t_2 \dots t_{H(u)}$ such that $t_i = s_{I(u)_i}$, where $I(u)$ is sorted by ascending order. For example, $\pi(abcd, 0110) = bc$. Let $K(u)$ be the output

of applying a variation of Knuth algorithm [14] using the redundancy symbols encoded with Gray code on u [18]. We first define the basic encoder $E_1 : \{0, 1\}^{f_q} \rightarrow \Sigma_{q^*}^{q^*}$ to be

$$E_1(u) \triangleq \pi(12 \dots q^*, K(u))$$

For example, let $q = 6$. We then have $f_q = 3$. Consider $u = 100$. We get that $K(u) = 010110$ and $E_1(u) = 245$.

Note that since K is injective so is the encoder E_1 .

Claim 2. *For any $u \in \{0, 1\}^{f_q}$, $E_1(u)$ is a subsequence of length $\lfloor \frac{q^*}{2} \rfloor$ of $A_{q^*}[q^*]$.*

Next, for all $u = u_1 u_2 \dots u_n$ such that $u_i \in \{0, 1\}^{f_q}$, we define the encoder $E_2 : \{0, 1\}^{n \cdot f_q} \rightarrow \Sigma_{q^*}^{n q^*}$ to be

$$E_2(u) \triangleq E_1(u_1) E_1(u_2) \dots E_1(u_n).$$

Observation 1. *For any $u \in \{0, 1\}^{n \cdot f_q}$, $E_2(u)$ is a subsequence of length $n \lfloor \frac{q^*}{2} \rfloor$ of $A_q[nq^*]$, and therefore can be synthesized with nq^* cycles.*

Theorem 3. *For all $q \in \mathbb{N}$, the encoder E_2 achieves $\rho = \frac{\lfloor \frac{q^*}{2} \rfloor}{q^*}$ and rate $R = \frac{f_q}{q^*}$.*

Theorem 3 follows from Observation 1 and the fact that $R = \lim_{n \rightarrow \infty} \frac{\log_2(2^{n \cdot f_q})}{n \cdot q^*} = \frac{f_q}{q^*}$. Also note that

$$\lim_{q \rightarrow \infty} \frac{f_q}{q^*} = \lim_{n \rightarrow \infty} \frac{n}{n + \lceil \log_2 n \rceil + 1} = 1 = \lim_{q \rightarrow \infty} \text{cap}(q, \frac{1}{2}).$$

For $q = 4, 8, 16, 32$, we get that the rate of this encoder is 0.5, 0.5714, 0.6875, 0.8125 respectively, with the capacities $\text{cap}(q, 0.5) = 0.92, 0.996, 0.99998, 1 - 10^{-9}$, respectively.

Note that by allowing a flexible $\rho \leq \frac{1}{2}$, in each window of length q we can encode any binary word of length $q-1$ using 1 redundancy symbols. For instance, if the weight of the binary word exceeds $\frac{q}{2}$, flipping the bits can be done, with the redundancy indicating whether the flipping occurred. The rate R of this encoder would be $R = \frac{q-1}{q} = 1 - \frac{1}{q}$.

D. Summary and Comparison Between the Approaches

In Section III-A we presented a naive encoding approach for a wide range of values of ρ using a lookup table. The rate of this encoding approaches the capacity as d gets large. But the calculation of the lookup table is time consuming and requires a substantial memory usage.

Next, in Section III-B we presented a linear time complexity encoding scheme, which is a generalization of the encoding approach presented in [15]. While this encoding achieves the capacity only for $\rho = \frac{2}{q+1}$, it is applicable for all $\rho \geq \frac{2}{q+1}$.

Finally, in Section III-C we presented a linear time complexity encoding scheme for $\rho = \frac{1}{2}$ that approaches the capacity as q gets large. Although this scheme, strictly speaking, is applicable only for $\rho = \frac{1}{2}$ (or $\rho = \frac{1}{2} - \frac{1}{2q}$), note that values close to $\frac{1}{2}$ achieve the highest capacities. Fig. 3 depicts the capacity and the maximum achievable rate using the methods from Section III-B and Section III-C. In summary, this section introduced two families of efficient encoders that achieve the capacity for certain values of ρ , while providing good performance in other configurations.

IV. OPTIMIZING THE SYNTHESIS COST

In this section we address Problem 2, seeking to identify the optimal parameters that minimize the synthesis cost as defined above by $\text{cost}^*(N, C, q, \rho) = \alpha C + \beta N \cdot \frac{\rho}{\text{cap}(q, \rho)}$.

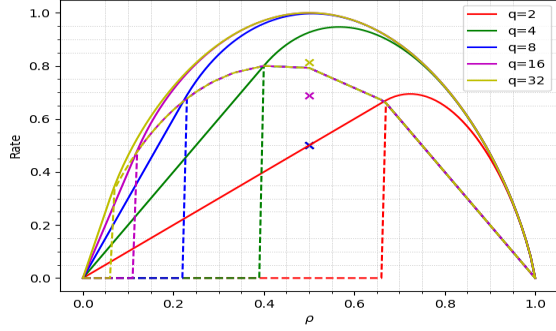


Fig. 3. Information rates achievable by the construction of Section III-B, compared to the related capacities $\text{cap}(q, \rho)$, for selected values of q and $\rho \in [0, 1]$. Dashed lines represent the maximum achievable rate of the encoder of Section III-B, the x markers represent the rate achievable of the encoder of Section III-C and solid lines represent capacities.

Observation 2. $\arg \min_{\rho \in [0, 1]} \text{cost}^*(N, C, q, \rho) = \arg \min_{\rho \in [0, 1]} \frac{\rho}{\text{cap}(q, \rho)}$.

This follows directly from the definition of the cost function, assuming all variables except ρ are fixed.

Claim 3. $\arg \min_{\rho \in [0, 1]} \text{cost}^*(N, C, q, \rho) = \arg \min_{\rho \in [\frac{2}{q+1}, 1]} \text{cost}^*(N, C, q, \rho)$.

Proof. For all $\rho \in (0, \frac{2}{q+1}]$ it holds that $\text{cap}(q, \rho) = \rho \log_2 q$ (see eq. (4)). Hence,

$$\text{cost}^*(N, C, q, \rho) = \alpha \cdot C + \beta \cdot N \frac{\rho}{\text{cap}(q, \rho)} = \alpha \cdot C + \beta \cdot \frac{N}{\log_2 q}.$$

Since the cost function is independent of ρ within this interval, we obtain the result. \square

Claim 4. $\min_{\rho \in (0, 1)} \text{cost}^*(N, C, q, \rho) > \min_{\rho \in (0, 1)} \text{cost}^*(N, C, q+1, \rho)$.

Observation 3. For all $q \in \mathbb{N}, \rho \in (0, 1) : \frac{\rho}{\text{cap}(q, \rho)} \geq \frac{\rho}{H(\rho)}$.

Lemma 3. The function $f(\rho) = \frac{\rho}{H(\rho)}$ is monotonously increasing for all $\rho \in (0, 1)$, and $\lim_{\rho \rightarrow 0} f(\rho) = 0$.

Theorem 4. Let $\rho^* \in (0, 1)$ be the solution to the equation $\frac{\rho}{H(\rho)} = \frac{1}{\log_2 q}$. Then,

$$\arg \min_{\rho \in [0, 1]} \frac{\rho}{\text{cap}(q, \rho)} = \arg \min_{\rho \in [\frac{2}{q+1}, \rho^*]} \frac{\rho}{\text{cap}(q, \rho)}.$$

Proof. According to Claim 3 it is enough to consider only $\rho \in [\frac{2}{q+1}, 1)$. As for the right end of the interval, for all $\rho \geq \rho^*$

$$\frac{\rho}{\text{cap}(q, \rho)} > \frac{\rho}{H(\rho)} \geq \frac{\rho^*}{H(\rho^*)} = \frac{1}{\log_2 q} = \frac{2}{\text{cap}(q, \frac{2}{q+1})},$$

where the first and second inequalities is by Observation 3 and Lemma 3. The last step follows from Equation (4). \square

Fig. 4 depicts ρ^* and $\frac{2}{q+1}$ as a function of q . Note that ρ^* is an implicit quantity. Using the monotonicity stated in Lemma 3 it can be easily approximated to any desirable degree. Turning our attention to Problem 2.2, let $\text{cost}_{OPT}^*(N, C) \triangleq \inf \{ \text{cost}^*(N, C, q, \rho) : \rho \in (0, 1), q \in \mathbb{N} \}$.

Claim 5. It holds that $\text{cost}_{OPT}^*(N, C) = \alpha C$.

Proof. From Claim 4 we get

$$\begin{aligned} \text{cost}_{OPT}^*(N, C) &= \inf \{ \lim_{q \rightarrow \infty} \text{cost}^*(N, C, q, \rho) : \rho \in (0, 1) \} \\ &= \inf \{ \alpha C + \beta N \frac{\rho}{H(\rho)} : \rho \in (0, 1) \}. \end{aligned}$$

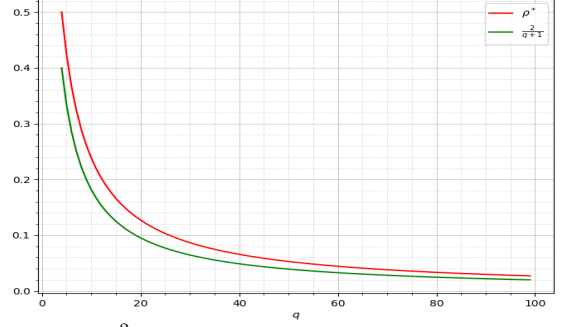


Fig. 4. Values of $\frac{2}{q+1}$ and ρ^* , as a function of q . The space between the lines is the interval of interest, as proven in Theorem 4.

From Lemma 3 we get that

$$\inf \left\{ \alpha C + \beta N \frac{\rho}{H(\rho)} : \rho \in (0, 1) \right\} = \lim_{\rho \rightarrow 0} \alpha C + \beta N \frac{\rho}{H(\rho)} = \alpha C.$$

When the alphabet size is bounded by Q we get

$$\begin{aligned} \text{cost}_{OPT}^*(N, C) &= \inf \{ \text{cost}^*(N, C, q, \rho) : \rho \in (0, 1), q \leq Q \} \\ &= \inf \{ \text{cost}^*(N, C, Q, \rho) : \rho \in (0, 1) \} \\ &\in \left\{ \alpha C + \beta N \frac{\rho}{\text{cap}(Q, \rho)} : \rho \in [\frac{2}{Q+1}, \rho^*] \right\}, \end{aligned}$$

assuming a characterization based on the implicit quantity ρ^* . \square

V. DISCUSSION

We study computational aspects of photolithographic synthesis. We start by expanding on literature results and presenting efficient coding and decoding at values of ρ that were not previously addressed. The number $\rho \in [0, 1]$ represents the fraction of synthesis cycles when using a total of C cycles to synthesize oligos of length $L < C$. Previous work calculated the exact capacity at a single value of ρ , for every alphabet size. This value, however, does not yield the most efficient scheme. We expand on these results, presenting effective coding and show, in particular, that for large enough values of q (that can be achieved, e.g. by composite alphabets [1], [19]), $\rho = \frac{1}{2}$, or close to it, is the most interesting value, achieving rates that approach capacity. We defined cost using (5). We note two key considerations. First, this cost is calculated assuming that capacity is achieved, and therefore is denoted cost^* . Further investigation is needed in order to estimate the cost at lower rates. Second, while β is assumed constant in this paper, larger q will result in higher cost per every added nucleotide. Therefore, it is natural to formulate our questions in terms of β as a function of q . This would require a deeper understanding of β in relation to q , which is currently lacking. We hope to address this question in future work. We note again that $\rho^* = \rho^*(q)$, an important parameter for the development presented in this paper and that ρ^* is implicit. Nonetheless, it can be well approximated thanks to the monotonicity of $\frac{\rho}{H(\rho)}$ in $(0, 1)$ (see Section IV). Finally - we note that efficient coding for values of q that are not large can be obtained by a lookup table, using the method described in Section III.

ACKNOWLEDGMENT

We thank the Yaakobi and Yakhini research groups as well as many DiDAX members for useful discussions.

REFERENCES

- [1] L. Anavy, I. Vaknin, O. Atar, R. Amit, and Z. Yakhini. Data storage in DNA with fewer synthesis cycles using composite DNA letters. *Nature Biotechnology*, 37:1229 – 1236, 2019.
- [2] P. L. Antkowiak, J. Lietard, M. Z. Darestani, M. M. Somoza, W. J. Stark, R. Heckel, and R. N. Grass. Low cost DNA data storage using photolithographic synthesis and advanced information reconstruction and error correction. *Nature Communications*, 11, 2020.
- [3] J. Behr, T. Michel, M. Giridhar, S. Santhosh, A. Das, H. Sabzalipoor, et al. An open-source advanced maskless synthesizer for light-directed chemical synthesis of large nucleic acid libraries and microarrays. *ChemRxiv*, 2024. This content is a preprint and has not been peer-reviewed.
- [4] M. Blawat, K. Gaedke, I. Hütter, X.-M. Chen, B. Turczyk, S. Inverso, B. W. Pruitt, and G. M. Church. Forward error correction for dna data storage. *Procedia Computer Science*, 80:1011–1022, 2016. International Conference on Computational Science 2016, ICCS 2016, 6-8 June 2016, San Diego, California, USA.
- [5] Y. Choi, T. Ryu, A. C. Lee, H. Choi, H.-B. Lee, J. Park, S.-H. Song, S. Kim, H. Kim, W. Park, and S. Kwon. High information capacity dna-based data storage with augmented encoding characters using degenerate bases. *Scientific Reports*, 9, 2019.
- [6] G. Church, Y. Gao, and S. Kosuri. Next-generation digital information storage in DNA. *Science (New York, N.Y.)*, 337:1628, 08 2012.
- [7] A. L. et al. Molecular data storage systems and methods, 2021. US Patent Application.
- [8] A. L. et al. Molecular data storage systems and methods, 2023. US Patent Application.
- [9] R. N. et al. Systems for nucleic acid-based data storage, 2024. US Patent Application.
- [10] N. Goldman, P. Bertone, S. Chen, C. Dessimoz, E. M. LeProust, B. Sipos, and E. Birney. Towards practical, high-capacity, low-maintenance information storage in synthesized DNA. *Nature*, 494(7435):77–80, 2013.
- [11] R. N. Grass, R. Heckel, M. Puddu, D. Paunescu, and W. J. Stark. Robust chemical preservation of digital information on dna in silica with error-correcting codes. *Angewandte Chemie*, 54 8:2552–5, 2015.
- [12] R. N. Grass, R. Heckel, M. Puddu, D. Paunescu, and W. J. Stark. Robust chemical preservation of digital information on DNA in silica with error-correcting codes. *Angewandte Chemie International Edition*, 54(8):2552–2555, feb 2015.
- [13] D. S. Hirschberg and M. Régnier. Tight bounds on the number of string subsequences daniel s. *Journal of Discrete Algorithms*, 2000.
- [14] D. E. Knuth. Efficient balanced codes. *IEEE Trans. Inf. Theory*, 32:51–53, 1986.
- [15] A. Lenz, Y. Liu, C. Rashtchian, P. H. Siegel, A. Wachter-Zeh, and E. Yaakobi. Coding for efficient dna synthesis. In *2020 IEEE International Symposium on Information Theory (ISIT)*, pages 2885–2890, 2020.
- [16] A. Lenz, S. Melczer, C. Rashtchian, and P. H. Siegel. Multivariate analytic combinatorics for cost constrained channels and subsequence enumeration. *ArXiv*, abs/2111.06105, 2021.
- [17] J. Lietard, D. Ameer, M. J. Damha, and M. M. Somoza. High-density rna microarrays synthesized in situ by photolithography. *Angewandte Chemie (International Ed. in English)*, 57:15257 – 15261, 2018.
- [18] E. N. Mambou and T. G. Swart. Encoding and decoding of balanced q-ary sequences using a gray code prefix. *2016 IEEE International Symposium on Information Theory (ISIT)*, pages 380–384, 2016.
- [19] I. Preuss, M. Rosenberg, Z. Yakhini, and L. Anavy. Efficient dna-based data storage using shortmer combinatorial encoding. *Scientific Reports*, 14(1):7731, Apr 2024.
- [20] E. Schaudy, J. Lietard, and M. M. Somoza. Enzymatic synthesis of high-density rna microarrays. *Current Protocols*, 3, 2023.
- [21] R. Shafir, O. Sabary, L. Anavy, E. Yaakobi, and Z. Yakhini. Sequence reconstruction under stutter noise in enzymatic dna synthesis. *2021 IEEE Information Theory Workshop (ITW)*, pages 1–6, 2021.
- [22] R. Shafir, O. Sabary, L. Anavy, E. Yaakobi, and Z. Yakhini. Sequence design and reconstruction under the repeat channel in enzymatic dna synthesis. *IEEE Transactions on Communications*, 72:675–691, 2024.
- [23] Y. Yan, N. Pinnamaneni, S. Chalapati, C. Crosbie, and R. Appuswamy. Scaling logical density of dna storage with enzymatically-ligated composite motifs. *Scientific Reports*, 13(1):15978, Sep 2023.