

Achieving DNA Labeling Capacity with Minimum Labels through Extremal de Bruijn Subgraphs

Christoph Hofmeister, Anina Gruica, Dganit Hanania, Rawad Bitar and Eitan Yaakobi

Abstract—DNA labeling is a tool in molecular biology and biotechnology to visualize, detect, and study DNA at the molecular level. In this process, a DNA molecule is labeled by a set of specific patterns, referred to as labels, and is then imaged. The resulting image is modeled as an $(\ell + 1)$ -ary sequence, where ℓ is the number of labels, in which any non-zero symbol indicates the appearance of the corresponding label in the DNA molecule. The labeling capacity refers to the maximum information rate that can be achieved by the labeling process for any given set of labels. The main goal of this paper is to study the minimum number of labels of the same length required to achieve the maximum labeling capacity of 2 for DNA sequences or $\log_2 q$ for an arbitrary alphabet of size q . The solution to this problem requires the study of path unique subgraphs of the de Bruijn graph with the largest number of edges and we provide upper and lower bounds on this value.

I. INTRODUCTION

Labeling of DNA molecules with fluorescent markers is a valuable technique for investigating DNA at the molecular level. This method is extensively employed in the fields of molecular biology and medicine, finding numerous applications in genomics and microbiology [1]–[3]. The labeling can be done in many methods, such as Fluorescence in situ hybridization (FISH) [1], CRISPR [2], [4], and Methyltransferases [5]. The process of optical mapping [6], [7] is one application of the labeling method, and in [8], a mathematical model was developed to optimize labels for specific applications within the optical mapping process.

The study of the information theoretic limits of the labeling process was initiated in [9]. Let $x \in \{A, C, G, T\}^n$ be a DNA sequence of length n and let $\alpha \in \{A, C, G, T\}^m$ be a label of length m . The labeling process marks all the locations where the label α appears in the DNA sequence x . The received sequence from the labeling process of x using the label α is a binary sequence y , also of length n , in which $y_i = 1$ if and only if $(x_i, \dots, x_{i+m-1}) = \alpha$. In general, this process can be extended for any $\ell > 0$ labels and in this case, the resulting sequence will have entries from $\{0, \dots, \ell\}$. For example, for the labels $\alpha_1 = AT$, $\alpha_2 = C$, the labeling of the sequence $x = ATCTGGATATCG$ produces the ternary sequence $y = (1, 0, 2, 0, 0, 0, 1, 0, 1, 0, 2, 0)$. The

labeling capacity, introduced in [9], is the logarithmic ratio between the number of possible outputs by labeling and the length of the sequence. This value expresses the maximum asymptotic information rate that can be achieved through DNA labeling. The labeling capacity was calculated for almost any type of single label and several cases of multiple labels were also found.

Another problem presented in [9] asked for the minimum number of labels of the same length needed to recognize “almost” any DNA sequence according to its labeling output sequence. Here, by almost, we refer to finding the minimum number of labels of the same length reaching the maximum labeling capacity, which is $\log_2(4) = 2$ in case of DNA sequences. In [9], the cases where the length of the labels is 1 or 2 were fully solved. For example, when the length of the labels is 1, the minimum number of labels to reach a labeling capacity of 2 is three. However, for labels of size greater than one, the problem becomes more challenging. Further, a simple extension of the results of [9], [10] shows that the minimum number of required labels of length $m > 1$ necessary to achieve the maximum labeling capacity over an alphabet of size q , which is $\log_2(q)$, equals the minimum number of edges that need to be removed from the $(m - 1)$ -dimensional de Bruijn graph over q symbols to make it path unique. A graph is called path unique, if for any k , between any two vertices, there exists at most one walk of length k . The maximum number of edges in a path unique 1-dimensional de Bruijn graph over q symbols is solved in [11], and from this result, it was observed in [9], that the minimum number of length-2 labels to reach maximum capacity $\log_2(4) = 2$ is 10.

Although [9], [10] provide an exact characterization of the problem of reaching the largest capacity with the minimum number of labels, the problem of finding the largest subgraphs of the de Bruijn graph that are path unique remained open. The main goal of this paper is to study this question.

The rest of this paper is organized as follows. In Section II, we introduce the notation and formally define our problem. We denote by $s(q, m)$ the minimum number of length- m labels over Σ_q that are needed in order to achieve the maximum labeling capacity. To find this value, we study $\gamma(q, d)$, the maximum number of edges in a path unique subgraph of a d -dimensional de Bruijn graph over q symbols. In Section III, we introduce constructions of path unique de Bruijn subgraphs, which provide two lower bounds on $\gamma(q, d)$. The first construction is for any d , and the second is restricted to $d = 2$. In Section IV we prove an upper bound on the value of $\gamma(q, d)$,

CH and RB are with the School of Computation, Information and Technology at the Technical University of Munich, Germany. Emails: {christoph.hofmeister, rawad.bitar}@tum.de

AG is with the Dept. of Mathematics and Computer Science at Eindhoven University of Technology, the Netherlands. Email: a.gruica@tue.nl

DH and EY are with the CS department of Technion—Israel Institute of Technology, Israel. Emails: dganit@campus.technion.ac.il, yaakobi@cs.technion.ac.il

and finally we compare the different bounds in Section V.

II. DEFINITIONS AND PRELIMINARIES

Throughout this paper, q , m , d and n are positive integers, Σ_q denotes the q -ary alphabet $\{0, 1, \dots, q-1\}$, and for a positive integer n we denote by $[n]$ the set $\{1, \dots, n\}$. For $\mathbf{x} = (x_1, \dots, x_n) \in \Sigma_q^n$ and $i, k \in [n]$ with $1 \leq i \leq n-k+1$ we let $\mathbf{x}_{[i;k]} = (x_i, \dots, x_{i+k-1})$ and a *label* $\alpha \in \Sigma_q^m$, $m < n$ is a (relatively short) sequence over Σ_q . For any matrix A of size $n \times m$, we denote by $|A|$ the sum of the entries of A , i.e., $|A| := \sum_{i \in [n], j \in [m]} A_{i,j}$. Next, we mathematically define the labeling process that is studied in this paper.

Definition 1 (The labeling model [9]). Let $\alpha_1, \dots, \alpha_\ell \in \Sigma_q^m$ be labels of length m . Denote by \mathcal{A} the set $\{\alpha_1, \dots, \alpha_\ell\}$, where $\alpha_1 \leq \dots \leq \alpha_\ell$ are ordered lexicographically.

- (i) The \mathcal{A} -labeling sequence of $\mathbf{x} = (x_1, \dots, x_n) \in \Sigma_q^n$ is the sequence $L_{\mathcal{A}}(\mathbf{x}) = (c_1, \dots, c_n) \in \Sigma_{\ell+1}^n$, in which $c_i = j$ if $\mathbf{x}_{[i;m]} = \alpha_j$ and $i \leq n-m+1$. If such a j does not exist or $i \in \{n-m, \dots, n\}$ then we set $c_i = 0$.
- (ii) A sequence $\mathbf{c} \in \Sigma_{\ell+1}^n$ is called a *valid \mathcal{A} -labeling sequence* if there exists $\mathbf{x} \in \Sigma_q^n$ with $L_{\mathcal{A}}(\mathbf{x}) = \mathbf{c}$.
- (iii) The *labeling capacity* of \mathcal{A} is

$$\text{cap}(\mathcal{A}) = \limsup_{n \rightarrow \infty} \frac{\log_2(|\{L_{\mathcal{A}}(\mathbf{x}) : \mathbf{x} \in \Sigma_q^n\}|)}{q^n}.$$

Example 1. Let $q = 4$, $m = 2$ and $\mathcal{A} = \{\alpha_1, \alpha_2\}$, where $\alpha_1 = (1, 0)$, $\alpha_2 = (2, 2)$. For $\mathbf{x} = (3, 1, 0, 3, 2, 2, 2, 3, 1, 0)$ we have $L_{\mathcal{A}}(\mathbf{x}) = (0, 1, 0, 0, 2, 2, 0, 0, 1, 0)$.

In [9], [10], the labeling capacity of a single label was computed for almost all cases and for several more cases of multiple labels. Another interesting problem presented in [9], [10] asked for the minimum number of labels of the same length that are needed in order to reach the largest labeling capacity of $\log_2 q$. Formally, this problem, which is the main focus of this paper, is stated as follows.

Problem 1. Find the minimum value ℓ of length- m labels over Σ_q such that there exists a set of labels $\mathcal{A} = \{\alpha_1, \dots, \alpha_\ell\}$ achieving the largest capacity. More formally, find the value

$$s(q, m) := \min\{\ell : \exists \mathcal{A} \subseteq \Sigma_q^m, |\mathcal{A}| = \ell, \text{cap}(\mathcal{A}) = \log_2(q)\}.$$

It is easy to see that for all q , $s(q, 1) = q - 1$. In order to solve it for larger values of m , it was observed in [9] that $s(4, 2) = 10$. This result was derived by using a connection between this problem and another problem over graphs that will be explained. The well-studied objects that allow us to do so are the well-known de Bruijn graphs [12], which are defined as follows.

Definition 2 (de Bruijn graphs). The d -dimensional *de Bruijn graph* $\mathcal{B}_{q,d} = (\mathcal{V}_{q,d}, \mathcal{E}_{q,d})$ of q elements is the directed graph with $\mathcal{V}_{q,d} = \Sigma_q^d$ and for $\mathbf{x} = (x_1, \dots, x_d) \in \Sigma_q^d$ and $\mathbf{y} = (y_1, \dots, y_d) \in \Sigma_q^d$ we have $(\mathbf{x}, \mathbf{y}) \in \mathcal{E}_{q,d}$ if and only if $(x_2, \dots, x_d) = (y_1, \dots, y_{d-1})$, i.e., the last $d-1$ entries of \mathbf{x} coincide with the first $d-1$ entries of \mathbf{y} .

Note that for a positive integer k , a *walk* of length k (i.e., a sequence of vertices and connecting edges with a total of k edges) in the de Bruijn graph $\mathcal{B}_{q,d}$ can be seen as a sequence $\mathbf{s} = (s_1, \dots, s_{d+k}) \in \Sigma_q^{d+k}$. In particular, an edge in $\mathcal{B}_{q,d}$ corresponds to a sequence in Σ_q^{d+1} . As shown later in Theorem 1, Problem 1 can be reformulated into a problem concerned with graphs satisfying a specific constraint:

Definition 3 (Path unique graph). Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a directed graph. We call \mathcal{G} *path unique* if for any positive integer k between any two vertices there exists at most one walk of length k .

An *edge-induced subgraph* is a subset of the edges of a graph together with any vertices that are their start- or endpoints. From now on, by subgraph we mean an edge-induced subgraph. The following theorem shows the connection between path unique subgraphs of $\mathcal{B}_{q,d}$ and the value of $s(q, d+1)$. Since it is a generalization of [10, Theorem 12] and it can be proven analogously, we omit the proof here.

Theorem 1. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a subgraph of $\mathcal{B}_{q,d}$ and let $\mathcal{A} = \Sigma_q^{d+1} \setminus \mathcal{E}$. It holds that $\text{cap}(\mathcal{A}) = \log_2(q)$ if and only if \mathcal{G} is path unique.

From Theorem 1 it follows that finding the size of the largest path unique subgraph of $\mathcal{B}_{q,d}$ is equivalent to finding the value of $s(q, d+1)$. This motivates the following reformulation of Problem 1.

Problem 2. For any q and d , let $\Gamma(q, d)$ be the set of subgraphs of $\mathcal{B}_{q,d}$ that are path unique. Find the value

$$\gamma(q, d) := \max\{|\mathcal{E}| : \mathcal{G} = (\mathcal{V}, \mathcal{E}) \in \Gamma(q, d)\}.$$

We call graphs in $\Gamma(q, d)$ that have $\gamma(q, d)$ edges *optimal* subgraphs of $\mathcal{B}_{q,d}$.

An immediate consequence of Theorem 1 is the following.

Corollary 1. We have $s(q, d+1) = q^{d+1} - \gamma(q, d)$.

Due to Corollary 1, it is enough to study the value of $\gamma(q, d)$ when solving Problem 1 for labels of length $m = d+1$ (and clearly Problem 2), which will be the focus of the rest of the paper.

The value of $\gamma(q, 1)$ was fully solved in [11], and it concerns the complete graph of q vertices. For example, it holds that $\gamma(4, 1) = 6$ and a corresponding graph using the DNA alphabet is visualized in Fig. 1.

Theorem 2 (see [11, Theorem 1]). We have

$$\gamma(q, 1) = \begin{cases} \frac{(q+1)^2}{4} & \text{if } q \text{ is odd,} \\ \frac{q(q+2)}{4} & \text{if } q \text{ is even.} \end{cases}$$

The approach of [11] relies on the following lemma.

Lemma 1. Let \mathcal{G} be a graph on n vertices, and let A be its adjacency matrix. For any $i, j \in [n]$, and for any positive

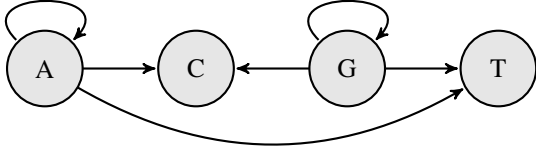


Fig. 1: An optimal path unique graph over $\Sigma_4 = \{A, C, G, T\}$.

integer k the entry in position (i, j) of \mathbf{A}^k equals the number of walks in \mathcal{G} starting at vertex i and ending at vertex j . Further, \mathcal{G} is a path unique graph if and only if for any positive integer k , the entries of \mathbf{A}^k are in $\{0, 1\}$.

We denote the adjacency matrix of a de Bruijn graph $\mathcal{B}_{q,d}$ by $\mathbf{B}_{q,d}$ (with vertices in lexicographic order).

III. CONSTRUCTIONS OF PATH UNIQUE SUBGRAPHS OF $\mathcal{B}_{q,d}$

In this section, we give two constructions of path unique subgraphs of $\mathcal{B}_{q,d}$, which provide lower bounds on $\gamma(q, d)$. The first construction is applicable for any q and d , while the second works for any q and $d = 2$ but gives a tighter lower bound. We start with the first construction.

Construction 1. Let $\mathcal{G}_{q,d}^1 = (\mathcal{V}_{q,d}^1, \mathcal{E}_{q,d}^1)$ be the subgraph of $\mathcal{B}_{q,d} = (\mathcal{V}_{q,d}, \mathcal{E}_{q,d})$ with $\mathcal{V}_{q,d}^1 = \mathcal{V}_{q,d}$ and where the edge $(x_1, \dots, x_{d+1}) \in \mathcal{E}_{q,d}$ is in $\mathcal{E}_{q,d}^1$ if and only if

- I) $x_1 \leq x_2$,
- II) it does **not** hold that $x_1 \leq x_2 \leq \dots \leq x_{d+1} \leq q - 2$.

Lemma 2. The total number of edges in $\mathcal{G}_{q,d}^1$ is

$$\left(\frac{q+1}{2q}\right) \cdot q^{d+1} - \binom{d+q-1}{d+1}.$$

Proof. The proportion of edges $(x_1, \dots, x_{d+1}) \in \Sigma_q^{d+1} = \mathcal{E}_{q,d}$ satisfying $x_1 \leq x_2$ within the set of all edges is $(q+1)/(2q)$. The edges $(x_1, \dots, x_{d+1}) \in \Sigma_q^{d+1}$ for which $x_1 \leq x_2 \leq \dots \leq x_{d+1} \leq q-2$ are clearly a subset of the edges described in I) of Construction 1. There are a total of $\binom{d+q-1}{d+1}$ such edges (see e.g. [13, Chapter II.5]) from which the lemma follows. \square

Example 2. When constructing the graph $\mathcal{G}_{3,3}^1$ we start with the graph shown in Fig. 2 consisting of 54 edges (where $(i, \mathbf{x}) \in \Sigma_q^d$ is the set of vertices whose first entry is equal to $i \in \Sigma_q$). We then remove the edges in the set

$$\{(0, 0, 0, 0), (0, 0, 0, 1), (0, 0, 1, 1), (0, 1, 1, 1), (1, 1, 1, 1)\}$$

resulting in a subgraph of $\mathcal{B}_{3,3}$ with 49 edges, which is $\mathcal{G}_{3,3}^1$.

Theorem 3. The graph $\mathcal{G}_{q,d}^1$ is path unique. In particular,

$$\gamma(q, d) \geq \left(\frac{q+1}{2q}\right) \cdot q^{d+1} - \binom{d+q-1}{d+1}.$$

Proof. Let $\mathbf{s} = (s_1, \dots, s_k) \in \Sigma_q^k$ be a walk on $\mathcal{G}_{q,d}^1$ from $\mathbf{x} = (x_1, \dots, x_d)$ to $\mathbf{y} = (y_1, \dots, y_d)$. If $k \leq 2d$, then there is nothing to show. If $k > 2d$, we show that $s_{d+1} = s_{d+2} =$

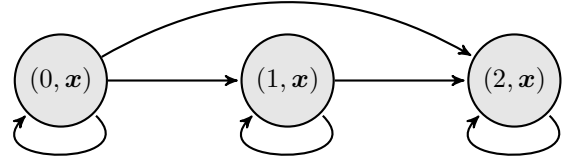


Fig. 2: The edges from part I) of Construction 1 for $q = 3$.

$\dots = s_{k-d} = q - 1$, hence the walk is unique. We split the proof into two parts, and we prove each part separately.

Claim 1. We have $x_1 \leq \dots \leq x_d$.

Proof of Claim 1. If there exists an $\ell \in \{2, \dots, d\}$ with $x_\ell < x_{\ell-1}$, this would mean that there is an edge $(x_{\ell-1}, x_\ell, \mathbf{z})$ in $\mathcal{G}_{q,d}^1$ for some $\mathbf{z} \in \Sigma_q^{d-1}$ with $x_\ell < x_{\ell-1}$, which by I) of Construction 1 is a contradiction. \square

Claim 2. We have $s_{d+1} = q - 1$.

Proof of Claim 2. We show that the only outgoing edge from \mathbf{x} is to $(x_2, \dots, q-1)$. Towards a contradiction, assume that there is an edge $(x_1, \dots, x_d, i) \in \mathcal{E}_{q,d}^1$ with $i \leq q-2$. First, note that $x_d \leq i$, as otherwise, similar to the proof of Claim 1, there is an edge (x_d, i, \mathbf{z}) in $\mathcal{G}_{q,d}^1$ for some $\mathbf{z} \in \Sigma_q^{d-1}$ with $x_d > i$, contradicting I) of Construction 1. Combined with Claim 1, it follows that $x_1 \leq \dots \leq x_d \leq i \leq q-2$ contradicting II) of Construction 1. \square

Similar to the proof of Claim 2, if $k-2d \geq 2$, by considering the walk going from $(x_2, \dots, x_d, s_{d+1})$ to \mathbf{y} , we have that $s_{d+2} = q-1$. We can show inductively that $s_{d+1} = s_{d+2} = \dots = s_{k-d} = q-1$. \square

Theorem 3 gives the following asymptotic lower bounds.

Corollary 2. The following statements hold.

- (i) For every fixed q , we have

$$\lim_{d \rightarrow +\infty} \frac{\gamma(q, d)}{q^{d+1}} \geq \frac{q+1}{2q}.$$

- (ii) For every fixed d , we have

$$\lim_{q \rightarrow +\infty} \frac{\gamma(q, d)}{q^{d+1}} \geq \frac{1}{2} - \frac{1}{(d+1)!}.$$

Next, we present our second construction.

Construction 2. Let $\mathcal{G}_q^2 = (\mathcal{V}_q^2, \mathcal{E}_q^2)$ be the subgraph of $\mathcal{B}_{q,2} = (\mathcal{V}_{q,2}, \mathcal{E}_{q,2})$ with $\mathcal{V}_q^2 = \mathcal{V}_{q,2}$ where the edge $(x_1, x_2, x_3) \in \mathcal{E}_{q,2}$ is in \mathcal{E}_q^2 if and only if it meets one of the following (mutually exclusive) conditions:

- I) $x_1 = x_2 = x_3$,
- II) $0 < x_1 = x_2 < x_3$,
- III) $x_1 > x_2$ and $x_2 \leq x_3$,
- IV) $0 = x_1 < x_2 \leq x_3$, or
- V) $q-1 = x_1 > x_2 > x_3 > 0$.

We partition the vertices in $\mathcal{V}_{q,2}^2$ into q blocks indexed between 0 and $q-1$. Within each block, we divide the vertices

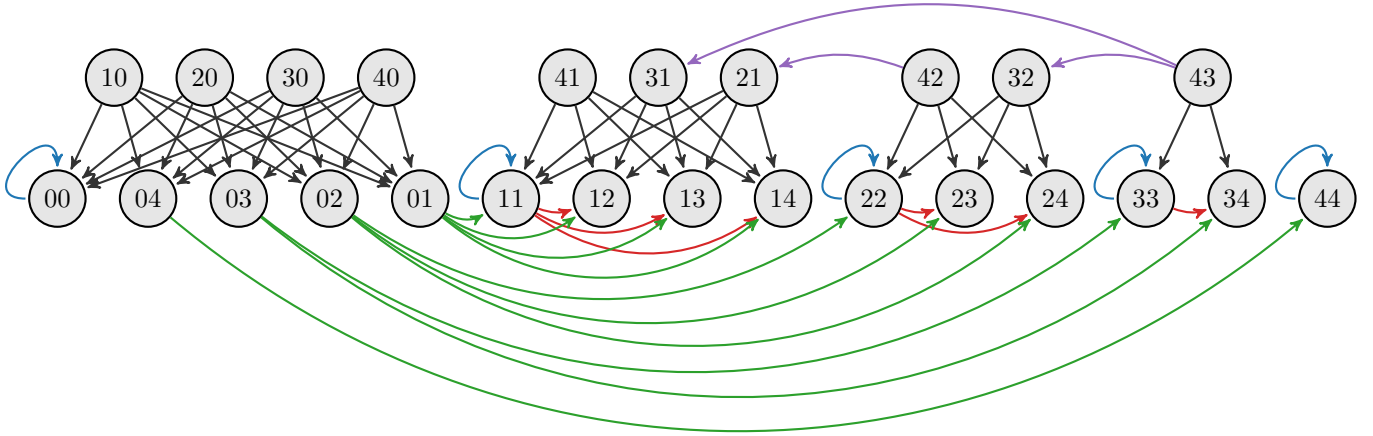


Fig. 3: \mathcal{G}_5^2 from Construction 2. For compactness, each vertex $(x_1, x_2) \in \Sigma_q^2$ is labeled by x_1x_2 , e.g., $(0, 1)$ by 01. Edges satisfying Condition I) are blue, edges satisfying Condition II) are red, edges satisfying Condition III) are black, edges satisfying Condition IV) are green, and edges satisfying Condition V) are purple.

into two parts. Each vertex (x_1, x_2) is in block $\min(x_1, x_2)$. Further, we say that (x_1, x_2) is in the *top part* if $x_1 > x_2$ and in the *bottom part* otherwise. For clarity, we assign a color to each edge $(x_1, x_2, x_3) \in \mathcal{E}_q^2$ based on which of the above conditions it meets, cf. Fig. 3 for a depiction.

More precisely:

- (i) Condition I) is met by self-loops which we color blue.
- (ii) Condition II) is met by edges from the constant vertex (x_1, x_1) in the bottom part of a block x_1 , with $x_1 > 0$, to another vertex in the bottom part of block x_1 . We color these edges in red.
- (iii) Condition III) is met by edges from the top part of a block to the bottom part of the same block. We color these edges in black.
- (iv) Condition IV) is met by edges from the bottom part in block 0 to the bottom part of another cluster. These edges are green.
- (v) Finally, Condition V) is met by edges from the top part of a block x_2 , with $1 < x_2 < q - 1$, to the top part of another block x_3 , with $0 < x_3 < x_2$. We assign to these edges the color purple.

Lemma 3. The number of edges in \mathcal{G}_q^2 is

$$\frac{1}{3}q^3 + \frac{3}{2}q^2 - \frac{23}{6}q + 4.$$

Proof. We count the number of triples $(x_1, x_2, x_3) \in \Sigma_q^3$ that fulfill each condition. Condition I) is met by q triples, Condition II) by $\binom{q-1}{2}$ triples, Condition III) by $\sum_{i=1}^{q-1} i^2 + i = 2\binom{q}{2} + 2\binom{q}{3}$ triples, Condition IV) by $\binom{q}{2}$ triples, and Condition V) by $\binom{q-2}{2}$ triples. Simplifying gives the desired expression. \square

Example 3. The graph \mathcal{G}_5^2 from Construction 2 is shown in Fig. 3. The color of each edge corresponds to the condition in Construction 2 it fulfills. The vertices are partitioned according

to the blocks of the construction, with the vertices in the top part arranged above the vertices in the bottom part.

Theorem 4. The graph \mathcal{G}_q^2 is path unique. In particular, we have

$$\gamma(q, 2) \geq \frac{1}{3}q^3 + \frac{3}{2}q^2 - \frac{23}{6}q + 4.$$

Before we prove the theorem, we compare the bounds obtained from Theorem 3 and Theorem 4.

Observation 1. The lower bound derived in Theorem 4 is tighter than the one derived in Theorem 3 for $d = 2$ and for all $q \geq 3$. For $d = 2$ and $q = 2$ the two bounds are the same.

To prove Theorem 4, we require the following two claims, which are proven in Appendix B and A, respectively.

Claim 3. For any walk of length three $(x_1, x_2, x_3, x_4, x_5)$ in \mathcal{G}_q^2 it holds that $x_3 = x_4$.

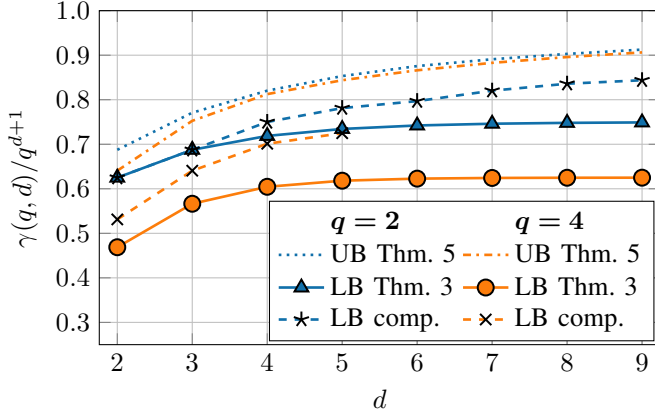
Claim 4. For any $x_1, x_2, x_4, x_5 \in \Sigma_q$, there exists an $x_3 \in \Sigma_q$ such that $(x_1, x_2, x_3, x_4, x_5)$ is a walk of length three in \mathcal{G}_q^2 if and only if (x_1, x_2, x_4, x_5) is a walk in \mathcal{G}_q^2 . Furthermore, if x_3 exists, it is unique.

Proof of Theorem 4. Let \mathbf{A} be the adjacency matrix of \mathcal{G}_q^2 . From Claim 4 it follows that $\mathbf{A}^2 = \mathbf{A}^3$, which implies $\mathbf{A}^2 = \mathbf{A}^i$ for all integers $i \geq 2$. Since $\mathbf{B}_{q,d}$ and $\mathbf{B}_{q,d}^2$ have no entry greater than 1, neither have \mathbf{A} and \mathbf{A}^2 . \square

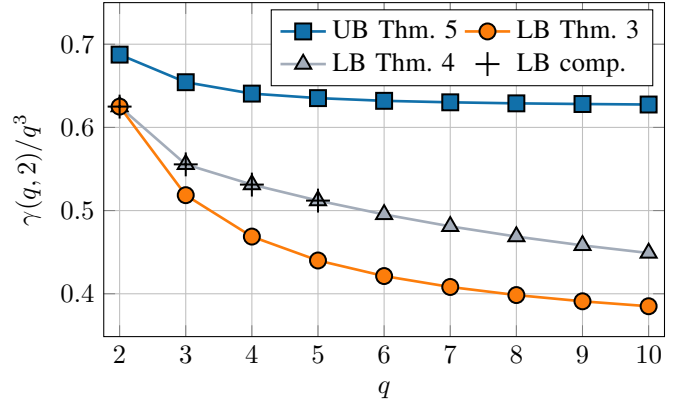
IV. AN UPPER BOUND ON $\gamma(q, d)$

Inspired by the upper bound on $\gamma(q, 1)$ in [11], in this section we give a general upper bound on $\gamma(q, d)$. We will need the following notation.

Notation 1. For a positive integer k we denote by $\eta(q, d, k)$ the maximum number of distinct walks in $\mathcal{B}_{q,d}$ of length k that have (at least) one edge in common.



(a) Alphabet size $q = 2$ and $q = 4$.



(b) de Bruijn graph dimension $d = 2$.

Fig. 4: Upper and lower bounds on the maximum number of edges in a path unique subgraph of $\mathcal{B}_{q,d}$ relative to the number of edges in the full de Bruijn graph $\mathcal{B}_{q,d}$. Results found by computer search are denoted as “LB comp.”.

Theorem 5. For any positive integer k , we have

$$\gamma(q, d) \leq q^{d+1} - \frac{q^{d+k} - \gamma(q^d, 1)}{\eta(q, d, k)}.$$

Proof. Let $\mathcal{G} = (\mathcal{V}_{q,d}, \mathcal{E}) \in \Gamma(q, d)$ be an arbitrary (but fixed) path unique subgraph of $\mathcal{B}_{q,d} = (\mathcal{V}_{q,d}, \mathcal{E}_{q,d})$ with adjacency matrix $\mathbf{A} \in \{0, 1\}^{q^d \times q^d}$. By Lemma 1, for any positive integer k it holds that $\mathbf{A}^k \in \{0, 1\}^{q^d \times q^d}$. In particular, for any positive integer i we have $\mathbf{A}^{ki} \in \{0, 1\}^{q^d \times q^d}$, i.e., \mathbf{A}^k is also an adjacency matrix of a path unique graph. Theorem 1 implies $|\mathbf{A}^k| \leq \gamma(q^d, 1)$. In the rest of the proof we upper bound the number of edges in \mathcal{G} using the fact that $|\mathbf{A}^k| \leq \gamma(q^d, 1)$.

Let $\mathcal{S} = \mathcal{E}_{q,d} \setminus \mathcal{E}$ be the set of edges that are in $\mathcal{B}_{q,d}$ but not in \mathcal{G} . The set \mathcal{W} of walks of length k in \mathcal{G} is the set of walks of length k in $\mathcal{B}_{q,d}$ that do not traverse any edge in \mathcal{S} . We have $|\mathcal{W}| = |\mathbf{A}^k| \leq \gamma(q^d, 1)$ and $|\mathcal{S}| = |\mathcal{E}_{q,d}| - |\mathcal{E}| = q^{d+1} - |\mathbf{A}|$. Since any edge in \mathcal{S} is traversed by at most $\eta(q, d, k)$ distinct walks of length k , it holds that $|\mathcal{W}| \geq q^{d+k} - \eta(q, d, k)|\mathcal{S}|$, which gives $|\mathcal{S}| \geq \frac{(q^{d+k} - |\mathcal{W}|)}{\eta(q, d, k)}$. Combining the above, we have

$$\gamma(q, d) \leq |\mathbf{A}| = q^{d+1} - |\mathcal{S}| \leq q^{d+1} - \frac{q^{d+k} - \gamma(q^d, 1)}{\eta(q, d, k)}. \quad \square$$

In order to explicitly evaluate the bound of Theorem 5, we need Theorem 2 and the following result, which we prove in Appendix C.

Lemma 4. We have

$$\eta(q, d, k) = \sum_{i=1}^{\lfloor \frac{d+k}{d+1} \rfloor} (-1)^{i+1} q^{d+k-i(d+1)} \binom{k - (i-1)d}{i}.$$

Corollary 3. By setting $k = d$ and $k = d + 1$, respectively, in Theorem 5 we get the following

$$\gamma(q, d) \leq q^{d+1} \left(1 - \frac{3}{4d} + \frac{1}{2dq^d} \right),$$

$$\gamma(q, d) \leq q^{d+1} \left(1 - \frac{1}{d+1} + \frac{1}{4q(d+1)} + \frac{1}{2(d+1)q^{d+1}} + \frac{1}{4(d+1)q^{2d+1}} \right).$$

In particular, for every fixed d it holds that

$$\lim_{q \rightarrow \infty} \frac{\gamma(q, d)}{q^{d+1}} \leq \min \left\{ 1 - \frac{1}{d+1}, 1 - \frac{3}{4d} \right\}.$$

Note that the asymptotic upper bound for fixed q and $d \rightarrow \infty$ for the ratio $\gamma(q, d)/q^{d+1}$ is (trivially) 1.

V. COMPARISON AND CONCLUSION

Generalizing a problem initially presented in [9], we analyzed the value of $\gamma(q, d)$, which is the maximum number of edges in a path unique subgraph of the de Bruijn graph $\mathcal{B}_{q,d}$. We presented two constructions. The first provides a lower bound on $\gamma(q, d)$ for any q and d , while the second is restricted to $d = 2$. Additionally, we proved an upper bound on $\gamma(q, d)$.

We conclude this paper by comparing the lower bounds derived in Theorems 3 and 4 and the upper bound derived in Theorem 5. Additionally, we include results from a computer search for small values of q and d , cf. Appendix D.

Fig. 4a displays the cases $q = 2$ and $q = 4$ for various values of d while Fig. 4b compares the results for $d = 2$ and various values of q . It can be seen that Theorem 4 gives a tighter bound than Theorem 3 and that Theorem 4 gives the same values as the computer search where the latter is available.

The results of this paper provide a significant step towards finding the value of $\gamma(q, d)$. Nevertheless, the plots in Fig. 4a and Fig. 4b illustrate that (for all parameters) there is an interesting gap between our lower and upper bounds. More progress on closing this gap, hence determining the exact values of $\gamma(q, d)$, remains for future work.

REFERENCES

- [1] A. Moter and U. B. Göbel, “Fluorescence in situ hybridization (FISH) for direct visualization of microorganisms,” *Journal of microbiological methods*, vol. 41, no. 2, pp. 85–112, 2000.
- [2] B. Chen, W. Zou, H. Xu, Y. Liang, and B. Huang, “Efficient labeling and imaging of protein-coding genes in living cells using CRISPR-tag,” *Nature communications*, vol. 9, no. 1, p. 5065, 2018.
- [3] D. Gruszka, J. Jeffet, S. Margalit, Y. Michaeli, and Y. Eberstein, “Single-molecule optical genome mapping in nanochannels: Multidisciplinarity at the nanoscale,” *Essays in Biochemistry*, vol. 65, no. 1, pp. 51–66, 2021.
- [4] H. Ma, A. Naseri, P. Reyes-Gutierrez, S. A. Wolfe, S. Zhang, and T. Pederson, “Multicolor crispr labeling of chromosomal loci in human cells,” *Proceedings of the National Academy of Sciences*, vol. 112, no. 10, pp. 3002–3007, 2015.
- [5] J. Deen, C. Vranken, V. Leen, R. K. Neely, K. P. Janssen, and J. Hofkens, “Methyltransferase-directed labeling of biomolecules and its applications,” *Angewandte Chemie International Edition*, vol. 56, no. 19, pp. 5182–5200, 2017.
- [6] M. Levy-Sakin and Y. Eberstein, “Beyond sequencing: Optical mapping of DNA in the age of nanotechnology and nanoscopy,” *Current opinion in biotechnology*, vol. 24, no. 4, pp. 690–698, 2013.
- [7] V. Müller and F. Westerlund, “Optical DNA mapping in nanofluidic devices: Principles and applications,” *Lab on a Chip*, vol. 17, no. 4, pp. 579–590, 2017.
- [8] Y. Nogin, D. Bar-Lev, D. Hanania, T. Detinis Zur, Y. Eberstein, E. Yaakobi, N. Weinberger, and Y. Shechtman, “Design of optimal labeling patterns for optical genome mapping via information theory,” *Bioinformatics*, vol. 39, no. 10, 2023.
- [9] D. Hanania, D. Bar-Lev, Y. Nogin, Y. Shechtman, and E. Yaakobi, “On the capacity of DNA labeling,” in *2023 IEEE International Symposium on Information Theory (ISIT)*, 2023, pp. 567–572.
- [10] —, “On the capacity of DNA labeling,” *arXiv preprint arXiv:2305.07992*, 2024.
- [11] Z. Huang and X. Zhan, “Extremal digraphs whose walks with the same initial and terminal vertices have distinct lengths,” *Discrete Mathematics*, vol. 312, no. 15, pp. 2203–2213, Aug. 2012.
- [12] N. G. de Bruijn, “A combinatorial problem,” *Proceedings of the Section of Sciences of the Koninklijke Nederlandse Akademie van Wetenschappen te Amsterdam*, vol. 49, no. 7, pp. 758–764, 1946.
- [13] W. Feller, *An introduction to probability theory and its applications, Volume 2*. John Wiley & Sons, 1991, vol. 81.
- [14] L. J. Guibas and A. M. Odlyzko, “String overlaps, pattern matching, and nontransitive games,” *Journal of Combinatorial Theory, Series A*, vol. 30, no. 2, pp. 183–208, Mar. 1981.
- [15] D. J. MacKay, *Information theory, inference and learning algorithms*. Cambridge university press, 2003.

APPENDIX

For convenience we restate results before the corresponding proofs.

A. Proof of Claim 3

Claim 3. For any walk of length three $(x_1, x_2, x_3, x_4, x_5)$ in \mathcal{G}_q^2 it holds that $x_3 = x_4$.

Proof. Let $(x_1, x_2, x_3, x_4, x_5)$ be a walk of length three in \mathcal{G}_q^2 . We perform a case distinction based on the color of the last edge (x_3, x_4, x_5) . If (x_3, x_4, x_5) is red or blue then $x_3 = x_4$ holds. It remains to show that (x_3, x_4, x_5) cannot be black, green or purple.

- Assume (x_3, x_4, x_5) is black. Since we have $x_3 > x_4$ the edge (x_2, x_3, x_4) can only be purple, i.e., $q - 1 = x_2 > x_3 > x_4 > 0$. However, then there exists no edge (x_1, x_2, x_3) .
- Assume (x_3, x_4, x_5) is green. We have $0 = x_3 < x_4 = x_5$. Then, (x_2, x_3, x_4) can only possibly be black implying $x_2 \neq 0$ and $x_3 = 0$ in which case there exists no edge (x_1, x_2, x_3) .
- Assume (x_3, x_4, x_5) is purple. If the edge (x_3, x_4, x_5) is purple then $q - 1 = x_3 > x_4 > x_5 > 0$ and there exists no possible edge (x_2, x_3, x_4) .

□

B. Proof of Claim 4

Claim 4. For any $x_1, x_2, x_4, x_5 \in \Sigma_q$, there exists an $x_3 \in \Sigma_q$ such that $(x_1, x_2, x_3, x_4, x_5)$ is a walk of length three in \mathcal{G}_q^2 if and only if (x_1, x_2, x_4, x_5) is a walk in \mathcal{G}_q^2 . Furthermore, if x_3 exists, it is unique.

Proof. Claim 3 shows that for any walk of length three $(x_1, x_2, x_3, x_4, x_5)$ in \mathcal{G}_q^2 it holds that $x_3 = x_4$. From the claim it directly follows that for any walk of length two with the same start and end vertex (x_1, x_2, x_4, x_5) there can be at most one x_3 such that $(x_1, x_2, x_3, x_4, x_5)$ is a walk of length three in \mathcal{G}_q^2 . It remains to show that if (x_1, x_2, x_3) , (x_2, x_3, x_4) , and (x_3, x_4, x_5) are edges in \mathcal{G}_q^2 , then so are (x_1, x_2, x_4) and (x_2, x_4, x_5) .

Since $x_3 = x_4$, if (x_1, x_2, x_3) is an edge in \mathcal{G}_q^2 , then so is (x_1, x_2, x_4) . We proceed to show that $(x_2, x_4, x_5) = (x_2, x_3, x_5)$ is an edge in \mathcal{G}_q^2 by case distinction on the color of the edges (x_3, x_3, x_5) and (x_2, x_3, x_3) . An edge of the form (x_3, x_3, x_5) can be blue or red.

- If (x_3, x_3, x_5) is blue, then $x_3 = x_4 = x_5$ and $(x_2, x_4, x_5) = (x_2, x_3, x_4)$.
- If (x_3, x_3, x_5) is red, then $x_3 < x_5$ and (x_2, x_3, x_3) can be blue, black, or green.
 - If (x_2, x_3, x_3) is blue, then $x_2 = x_3$ and $(x_2, x_4, x_5) = (x_3, x_4, x_5)$.
 - If (x_2, x_3, x_3) is black, then $x_3 > x_4, x_4 < x_5$ and (x_2, x_4, x_5) is also black.
 - If (x_2, x_3, x_3) is green, then $x_2 < x_4 < x_5$ and (x_2, x_4, x_5) is also green.

□

C. Proof of Lemma 4

Lemma 4. We have

$$\eta(q, d, k) = \sum_{i=1}^{\lfloor \frac{d+k}{d+1} \rfloor} (-1)^{i+1} q^{d+k-i(d+1)} \binom{k - (i-1)d}{i}.$$

Proof. We use the correspondence between walks in $\mathcal{B}_{q,d}$ and sequences over Σ_q . To determine $\eta(q, d, k)$, we count the maximum number of sequences $\mathbf{b} \in \Sigma_q^{d+k}$ that contain a $\mathbf{a} \in \Sigma_q^{d+1}$ as a contiguous subsequence, i.e.,

$$\eta(q, d, k) = \max_{\mathbf{a} \in \Sigma_q^{d+1}} |\{ \mathbf{b} \in \Sigma_q^{d+k} : \mathbf{b}_{[\ell, d+1]} = \mathbf{a} \text{ for some } \ell \in [k] \}|.$$

In [14, Section 7] it is shown that the number of sequences not containing a string monotonically increases with the autocorrelation of the string. Hence, the maximum is attained for a sequence \mathbf{a} with the lowest possible autocorrelation, i.e., a string that does not overlap itself, e.g. $\mathbf{a} = (1, 0, 0, \dots, 0)$. For $d+k < 2d$ we have $\eta(q, d, k) = kq^{d+k-i(d+1)}$ since there are k positions where \mathbf{a} can be placed in \mathbf{b} and $q^{d+k-(d+1)}$ possibilities for choosing the remaining symbols of \mathbf{b} . In the general case, to avoid overcounting the sequences \mathbf{b} that contain \mathbf{a} more than once, we use the Inclusion-Exclusion Principle to obtain

$$\eta(q, d, k) = \sum_{i=1}^{\lfloor \frac{d+k}{d+1} \rfloor} (-1)^{i+1} q^{d+k-i(d+1)} P_i,$$

where P_i denotes the number of possibilities to place i copies of \mathbf{a} into \mathbf{b} . It is multiplied by the number of ways to choose the remaining symbols of \mathbf{b} , which is $q^{d+k-i(d+1)}$. Note that P_i is the same as the number of integer vectors (x_1, \dots, x_i) with $1 \leq x_1 \leq \dots \leq x_i \leq d+k-i(d+1)$ (each x_j can be seen as the position for which the length- $(d+k)$ vector \mathbf{b} satisfies $\mathbf{b}_{[(j-1)(d+1)+x_j, d+1]} = \mathbf{a}$, i.e., the position of \mathbf{b} in which \mathbf{a} starts). Similar to the proof of Lemma 2 (also here see e.g. [13, Chapter II.5] for a reference), we have $P_i = \binom{d+k-i(d+1)+i}{i}$, which gives the formula in the lemma. □

D. Computer Search and Further Comparison

We conduct a computer search for path unique subgraphs of de Bruijn graphs $\mathcal{B}_{q,d}$ for small values of q and d .

Loosely based on [11, Lemma 5], we search for permutation matrices $\mathbf{\Pi}$ such that the upper triangle of the permuted adjacency matrix $\text{triu}(\mathbf{\Pi}^{-1} \mathbf{B}_{q,d} \mathbf{\Pi})$ is path unique and contains a high number of ones. The search algorithm is in the spirit of simulated annealing, cf. [15, Chapter 30]. Starting with a random permutation, at each iteration a random swap is performed on the permutation matrix. The new permutation matrix is accepted if it results in a path unique matrix with a higher number of ones. If it does not, it is still going to be accepted with a small probability.

The results of the computer search are reported in Table I along with explicit evaluations of our theoretical results. For $q = 2$ with $d \in \{2, 3, 4\}$ and $q = 3$ with $d = 2$ direct exhaustive search for path unique allocation matrices is feasible. The

results for these values coincide with the values found through the aforementioned search algorithm and reported in the table.

q	d	LB comp.	LB Thm. 3	LB Thm. 4	UB Thm. 5
2	2	5	5	5	5
2	3	11	11	-	12
2	4	24	23	-	26
2	5	50	47	-	54
2	6	102	95	-	112
2	7	210	191	-	228
2	8	428	383	-	462
2	9	864	767	-	934
3	2	15	14	15	17
3	3	53	49	-	61
3	4	174	156	-	197
3	5	553	479	-	617
3	6	1712	1450	-	1900
4	2	34	30	34	41
4	3	164	145	-	192
4	4	718	619	-	832
4	5	2969	2532	-	3456
5	2	64	55	64	79
5	3	372	340	-	469

TABLE I: Evaluation of the lower bound obtained experimentally through the computer search (denoted by “LB comp.”), the lower bound in Theorem 3 (denoted by “LB Thm. 3”), the lower bound in Theorem 4 (denoted by “LB Thm. 4”), and the upper bound in Theorem 5 (denoted by “UB Thm. 5”) for various values of q and d .

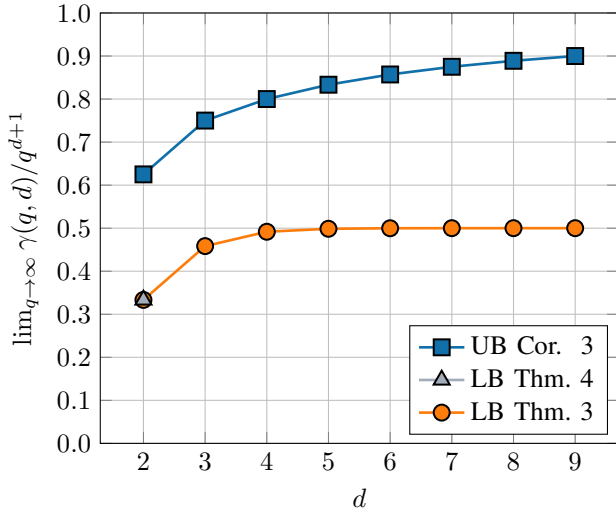


Fig. 5: Upper and lower bounds on the relative number of edges in an optimal path unique subgraph of the de Bruijn graph when the alphabet size tends to infinity.

Fig. 5 displays the asymptotic results for large alphabet sizes q . For $d = 2$, both lower bounds tend to $\frac{1}{3}$, while the upper bound from Corollary 3 tends to $\frac{5}{8}$. For large d , the lower bound from Theorem 3 tends to $\frac{1}{2}$ while the upper bound tends to 1.