

Reconstructing Words from Right-Bounded-Block Words

Pamela Fleischmann¹, Marie Lejeune^{2*}, Florin Manea^{3**}, Dirk Nowotka¹, and
Michel Rigo²

¹ Kiel University, Germany {fpa,dn}@informatik.uni-kiel.de

² University of Liège, Belgium {m.lejeune,m.rigo}@uliege.be

³ University of Göttingen, Germany florin.manea@informatik.uni-goettingen.de

Abstract. A reconstruction problem of words from scattered factors asks for the minimal information, like multisets of scattered factors of a given length or the number of occurrences of scattered factors from a given set, necessary to uniquely determine a word. We show that a word $w \in \{\mathbf{a}, \mathbf{b}\}^*$ can be reconstructed from the number of occurrences of at most $\min(|w|_{\mathbf{a}}, |w|_{\mathbf{b}}) + 1$ scattered factors of the form $\mathbf{a}^i \mathbf{b}$, where $|w|_{\mathbf{a}}$ is the number of occurrences of the letter a in w . Moreover, we generalize the result to alphabets of the form $\{1, \dots, q\}$ by showing that at most $\sum_{i=1}^{q-1} |w|_i (q - i + 1)$ scattered factors suffices to reconstruct w . Both results improve on the upper bounds known so far. Complexity time bounds on reconstruction algorithms are also considered here.

1 Introduction

The general scheme for a so-called *reconstruction problem* is the following one: given a sufficient amount of information about substructures of a hidden discrete structure, can one uniquely determine this structure? In particular, what are the fragments about the structure needed to recover it all. For instance, a square matrix of size at least 5 can be reconstructed from its principal minors given in any order [20].

In graph theory, given some subgraphs of a graph (these subgraphs may share some common vertices and edges), can one uniquely rebuild the original graph? Given a finite undirected graph $G = (V, E)$ with n vertices, consider the multiset made of the n induced subgraphs of G obtained by deleting exactly one vertex from G . In particular, one knows how many isomorphic subgraphs of a given class appear. Two graphs leading to the same multiset (generally called a *deck*) are said to be *hypomorphic*. A conjecture due to Kelly and Ulam states that two hypomorphic graphs with at least three vertices are isomorphic [14, 21]. A similar conjecture in terms of edge-deleted subgraphs has been proposed by Harary [11]. These conjectures are known to hold true for several families of graphs.

* Supported by a FNRS fellowship.

** Supported by the DFG grant MA 5725/2-1.

A *finite word*, i.e., a finite sequence of letters of some given alphabet, can be seen as an edge- or vertex-labeled linear tree. So variants of the graph reconstruction problem can be considered and are of independent interest. Participants of the Oberwolfach meeting on Combinatorics on Words in 2010 [2] gave a list of 18 important open problems in the field. Amongst them, the twelfth problem is stated as *reconstruction from subwords of given length*. In the following statement and all along the paper, a *subword* of a word is understood as a subsequence of not necessarily contiguous letters from this word, i.e., subwords can be obtained by deleting letters from the given word. To highlight this latter property, they are often called *scattered subwords* or *scattered factors*, which is the notion we are going to use.

Definition 1. *Let k, n be natural numbers. Words of length n over a given alphabet are said to be k -reconstructible whenever the multiset of scattered factors of length k (or k -deck) uniquely determines any word of length n .*

Notice that the definition requires multisets to store the information how often a scattered factor occurs in the words. For instance, the scattered factor **ba** occurs three times in **baba** which provides more information for the reconstruction than the mere fact that **ba** is a scattered factor.

The challenge is to determine the function $f(n) = k$ where k is the least integer for which words of length n are k -reconstructible. This problem has been studied by several authors and one of the first trace goes back to 1973 [13]. Results in that direction have been obtained by M.-P. Schützenberger (with the so-called *Schützenberger’s Guessing game*) and L. Simon [25]. They show that words of length n sharing the same multiset of scattered factors of length up to $\lfloor n/2 \rfloor + 1$ are the same. Consequently, words of length n are $(\lfloor n/2 \rfloor + 1)$ -reconstructible. In [15], this upper bound has been improved: Krasikov and Roditty have shown that words of length n are k -reconstructible for $k \geq \lfloor 16\sqrt{n}/7 \rfloor + 5$. On the other hand Dudik and Schulmann [6] provide a lower bound: if words of length n are k -reconstructible, then $k \geq 3^{(\sqrt{2/3} - o(1)) \log_3^{1/2} n}$. Bounds were also considered in [19]. Algorithmic complexity of the reconstruction problem is discussed, for instance, in [5]. Note that the different types of reconstruction problems have application in phylogenetic networks, see, e.g., [12], or in the context of molecular genetics [7] and coding theory [16].

Another motivation, close to combinatorics on words, stems from the study of k -binomial equivalence of finite words and k -binomial complexity of infinite words (see [23] for more details). Given two words of the same length, they are k -binomially equivalent if they have the same multiset of scattered factors of length k , also known as *k -spectrum* ([1], [18], [24]). Given two words x and y of the same length, one can address the following problem: decide whether or not x and y are k -binomially equivalent? A polynomial time decision algorithm based on automata and a probabilistic algorithm have been addressed in [10]. A variation of our work would be to find, given k and n , a minimal set of scattered factors for which the knowledge of the number of occurrences in x and y permits to decide k -binomial equivalence.

Over an alphabet of size q , there are q^k pairwise distinct length- k factors. If we relax the requirement of only considering scattered factors of the same length, another interesting question is to look for a minimal (in terms of cardinality) multiset of scattered factors to reconstruct entirely a word. Let the *binomial coefficient* $\binom{u}{x}$ be the number of occurrences of x as a scattered factor of u . The general problem addressed in this paper is therefore the following one.

Problem 2. Let Σ be a given alphabet and n a natural number. We want to reconstruct an unknown word $w \in \Sigma^n$. To that aim, we are allowed to ask questions of the type "What is the value of $\binom{w}{u_i}$?". Based on the answers to questions related to $\binom{w}{u_1}, \dots, \binom{w}{u_i}$, we can decide which will be the next question (i.e. decide which word will be u_{i+1}). We want to have the shortest sequence (u_1, \dots, u_k) uniquely determining w by knowing the values of $\binom{w}{u_1}, \dots, \binom{w}{u_k}$.

In this new context, we naturally look for a value of k less than the upper bound for k -reconstructibility.

In this paper, we firstly recall the use of Lyndon words in the context of reconstructibility. A word w over a totally ordered alphabet is called *Lyndon word* if it is the lexicographically smallest amongst all its rotations, i.e., $w = xy$ is smaller than yx for all non trivial factorisations $w = xy$. Every binomial coefficient $\binom{w}{x}$ for arbitrary words w and x over the same alphabet can be deduced from the values of the coefficients $\binom{u}{u}$ for Lyndon words u that are lexicographically less than or equal to x . This result is presented in Section 2 along with the basic definitions. We consider an alphabet equipped with a total order on the letters. Words of the form $a^n b$ with letters $a < b$ and a natural number n are a special form of Lyndon words, the so-called *right-bounded-block* words.

We consider the reconstruction problem from the information given by the occurrences of right-bounded-block words as scattered factors of a word of length n . In Section 3 we show how to reconstruct a word uniquely from $m + 1$ binomial coefficients of right-bounded-block words where m is the minimum number of occurrences of a and b in the word. We also prove that this is less than the upper bound given in [15]. In Section 4 we reduce the problem for arbitrary finite alphabets $\{1, \dots, q\}$ to the binary case. Here we show that at most $\sum_{i=1}^{q-1} |w|_i (q - i + 1) \leq q|w|$ binomial coefficients suffice to uniquely reconstruct w with $|w|_i$ being the number of occurrences of letter i in w . Again, we compare this bound to the best known one for the classical reconstruction problem (from words of a given length). In the last section of the paper we also propose several results of algorithmic nature regarding the efficient reconstruction of words from given scattered factors.

Due to space restrictions some proofs (marked with $*$) can be found in [9].

2 Preliminaries

Let \mathbb{N} be the set of natural numbers, $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$, and let $\mathbb{N}_{\geq k}$ be the set of all natural numbers greater than or equal to k . Let $[n]$ denote the set $\{1, \dots, n\}$ and $[n]_0 = [n] \cup \{0\}$ for an $n \in \mathbb{N}$.

An *alphabet* $\Sigma = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots\}$ is a finite set of letters and a *word* is a finite sequence of letters. We let Σ^* denote the set of all finite words over Σ . The *empty word* is denoted by ε and Σ^+ is the free semigroup $\Sigma^* \setminus \{\varepsilon\}$. The length of a word w is denoted by $|w|$. Let $\Sigma^{\leq k} := \{w \in \Sigma^* \mid |w| \leq k\}$ and Σ^k be the set of all words of length exactly $k \in \mathbb{N}$. The number of occurrences of a letter $\mathbf{a} \in \Sigma$ in a word $w \in \Sigma^*$ is denoted by $|w|_{\mathbf{a}}$. The i^{th} letter of a word w is given by $w[i]$ for $i \in [|w|]$. The powers of $w \in \Sigma^*$ are defined recursively by $w^0 = \varepsilon$, $w^n = ww^{n-1}$ for $n \in \mathbb{N}$. A word $u \in \Sigma^*$ is a *factor* of $w \in \Sigma^*$, if $w = xuy$ holds for some words $x, y \in \Sigma^*$. Moreover, u is a *prefix* of w if $x = \varepsilon$ holds and a *suffix* if $y = \varepsilon$ holds. The factor of w from the i^{th} to the j^{th} letter will be denoted by $w[i..j]$ for $1 \leq i \leq j \leq |w|$. Two words $u, v \in \Sigma^*$ are called *conjugates* or *rotations* of each other if there exist $x, y \in \Sigma^*$ with $u = xy$ and $v = yx$. Additional basic information about combinatorics on words can be found in [17].

Definition 3. Let $<$ be a total ordering on Σ . A word $w \in \Sigma^*$ is called *right-bounded-block word* if there exist $\mathbf{x}, \mathbf{y} \in \Sigma$ with $\mathbf{x} < \mathbf{y}$ and $\ell \in \mathbb{N}_0$ with $w = \mathbf{x}^\ell \mathbf{y}$.

Definition 4. A word $u = \mathbf{a}_1 \cdots \mathbf{a}_n \in \Sigma^n$, for $n \in \mathbb{N}$, is a *scattered factor* of a word $w \in \Sigma^+$ if there exist $v_0, \dots, v_n \in \Sigma^*$ with $w = v_0 \mathbf{a}_1 v_1 \cdots v_{n-1} \mathbf{a}_n v_n$. For words $w, u \in \Sigma^*$, define $\binom{w}{u}$ as the number of occurrences of u as a scattered factor of w .

Remark 5. Notice that $|w|_{\mathbf{x}} = \binom{w}{\mathbf{x}}$ for all $\mathbf{x} \in \Sigma$.

The following definition addresses Problem 2.

Definition 6. A word $w \in \Sigma^n$ is called *uniquely reconstructible/determined* by the set $S \subset \Sigma^*$ if for all words $v \in \Sigma^n \setminus \{w\}$ there exists a word $u \in S$ with $\binom{w}{u} \neq \binom{v}{u}$.

Consider $S = \{\mathbf{ab}, \mathbf{ba}\}$. Then $w = \mathbf{abba}$ is not uniquely reconstructible by S since $\left[\binom{w}{\mathbf{ab}}, \binom{w}{\mathbf{ba}}\right] = [2, 2]$ is also the 2-vector of binomial coefficients of \mathbf{baab} . On the other hand $S = \{\mathbf{a}, \mathbf{ab}, \mathbf{ab}^2\}$ reconstructs w uniquely. The following remark gives immediate results for binary alphabets.

Remark 7. Let $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and $w \in \Sigma^n$. If $|w|_{\mathbf{a}} \in \{0, n\}$ then w contains either only \mathbf{b} or \mathbf{a} and by the given length n of w , w is uniquely determined by $S = \{\mathbf{a}\}$. This fact is in particular an equivalence: $w \in \Sigma^n$ can be uniquely determined by $\{\mathbf{a}\}$ iff $|w|_{\mathbf{a}} \in \{0, n\}$. If $|w|_{\mathbf{a}} \in \{1, n-1\}$, w is not uniquely determined by $\{\mathbf{a}\}$ as witnessed by \mathbf{ab} and \mathbf{ba} for $n = 2$. It is immediately clear that the additional information $\binom{w}{\mathbf{ab}}$ leads to unique determinism of w .

Lyndon words play an important role regarding the reconstruction problem. As shown in [22] only scattered factors which are Lyndon words are necessary to determine a word uniquely, i.e., S can always be assumed to be a set of Lyndon words.

Definition 8. Let $<$ be a total ordering on Σ . A word $w \in \Sigma^*$ is a *Lyndon word* iff for all $u, v \in \Sigma^+$ with $w = uv$, we have $w <_{\text{lex}} vu$ where $<_{\text{lex}}$ is the lexicographical ordering on words induced by $<$.

Proposition 9 ([22]). *Let w and u be two words. The binomial coefficient $\binom{w}{u}$ can be computed using only binomial coefficients of the type $\binom{w}{v}$ where v is a Lyndon word of length up to $|u|$ such that $v \leq_{lex} u$.*

To obtain a formula to compute the binomial coefficient $\binom{w}{u}$ for $w, u \in \Sigma^*$ by binomial coefficients $\binom{w}{v_i}$ for Lyndon words v_1, \dots, v_k with $v_i \in \Sigma^{\leq |u|}$, $i \in [k]$, and $k \in \mathbb{N}$ the definitions of shuffle and infiltration are necessary [17].

Definition 10. *Let $n_1, n_2 \in \mathbb{N}$, $u_1 \in \Sigma^{n_1}$, and $u_2 \in \Sigma^{n_2}$. Set $n = n_1 + n_2$. The shuffle of u_1 and u_2 is the polynomial $u_1 \sqcup u_2 = \sum_{I_1, I_2} w(I_1, I_2)$ where the sum has to be taken over all pairs (I_1, I_2) of sets that are partitions of $[n]$ such that $|I_1| = n_1$ and $|I_2| = n_2$. If $I_1 = \{i_{1,1} < \dots < i_{1,n_1}\}$ and $I_2 = \{i_{2,1} < \dots < i_{2,n_2}\}$, then the word $w(I_1, I_2)$ is defined such that $w[i_{1,1}]w[i_{1,2}] \cdots w[i_{1,n_1}] = u_1$ and $w[i_{2,1}]w[i_{2,2}] \cdots w[i_{2,n_2}] = u_2$ hold.*

The infiltration is a variant of the shuffle in which equal letters are merged.

Definition 11. *Let $n_1, n_2 \in \mathbb{N}$, $u_1 \in \Sigma^{n_1}$, and $u_2 \in \Sigma^{n_2}$. Set $n = n_1 + n_2$. The infiltration of u_1 and u_2 is the polynomial $u_1 \downarrow u_2 = \sum_{I_1, I_2} w(I_1, I_2)$, where the sum has to be taken over all pairs (I_1, I_2) of sets of cardinality n_1 and n_2 respectively, for which the union is equal to the set $[n']$ for some $n' \leq n$. Words $w(I_1, I_2)$ are defined as in the previous definition. Note that some $w(I_1, I_2)$ are not well defined if $i_{1,j} = i_{2,k}$ but $u_1[j] \neq u_2[k]$. In that case they do not appear in the previous sum.*

Considering for instance $u_1 = \text{aba}$ and $u_2 = \text{ab}$ gives the polynomials

$$\begin{aligned} u_1 \sqcup u_2 &= 2\text{ababa} + 4\text{aabba} + 2\text{aabab} + 2\text{abaab}, \\ u_1 \downarrow u_2 &= \text{aba} \sqcup \text{ab} + \text{aba} + 2\text{abba} + 2\text{aaba} + 2\text{abab}. \end{aligned}$$

Based on Definitions 10 and 11, we are able to give a formula to compute a binomial coefficient from the ones making use of Lyndon words. This formula is given implicitly in [22, Theorem 6.4]: Let $u \in \Sigma^*$ be a non-Lyndon word. By [22, Corollary 6.2] there exist non-empty words $x, y \in \Sigma^*$ and with $u = xy$ and such that every word appearing in the polynomial $x \sqcup y$ is lexicographically less than or equal to u . Then, for all word $w \in \Sigma^*$, we have

$$\binom{w}{u} = \frac{1}{(x \sqcup y, u)} \left[\binom{w}{x} \binom{w}{y} - \sum_{v \in \Sigma^*, v \neq u} (x \downarrow y, v) \binom{w}{v} \right],$$

where (P, v) is a notation giving the coefficient of the word v in the polynomial P . One may apply recursively this formula until only Lyndon factors are considered. Some examples can be found in the appendix.

3 Reconstruction from Binary Right-Bounded-Block Words

In this section we present a method to reconstruct a binary word uniquely from binomial coefficients of right-bounded-block words. Let $n \in \mathbb{N}$ be a natural number and $w \in \{\mathbf{a}, \mathbf{b}\}^n$ a word. Since the word length n is assumed to be known, $|w|_{\mathbf{a}}$ is known if $|w|_{\mathbf{b}}$ is given and vice versa. Set for abbreviation $k_u = \binom{w}{u}$ for $u \in \Sigma^*$. Moreover we assume w.l.o.g. $k_{\mathbf{a}} \leq k_{\mathbf{b}}$ and that $k_{\mathbf{a}}$ is known (otherwise substitute each \mathbf{a} by \mathbf{b} and each \mathbf{b} by \mathbf{a} , apply the following reconstruction method and revert the substitution). This implies that w is of the form

$$\mathbf{b}^{s_1} \mathbf{a} \mathbf{b}^{s_2} \dots \mathbf{b}^{s_{k_{\mathbf{a}}}} \mathbf{a} \mathbf{b}^{s_{k_{\mathbf{a}}+1}} \quad (1)$$

for $s_i \in \mathbb{N}_0$ and $i \in [|w|_{\mathbf{a}} + 1]$ with $\sum_{i \in [k_{\mathbf{a}}+1]} s_i = n - k_{\mathbf{a}} = k_{\mathbf{b}}$ and thus we get for $\ell \in [k_{\mathbf{a}}]_0$

$$k_{\mathbf{a}^\ell \mathbf{b}} = \binom{w}{\mathbf{a}^\ell \mathbf{b}} = \sum_{i=\ell+1}^{k_{\mathbf{a}}+1} \binom{i-1}{\ell} s_i. \quad (2)$$

Remark 12. Notice that for fixed $\ell \in [k_{\mathbf{a}}]_0$ and $c_i = \binom{i-1}{\ell}$ for $i \in [k_{\mathbf{a}} + 1] \setminus \{\ell\}$, we have $c_i < c_{i+1}$ and especially $c_{\ell+1} = 1$ and $c_{\ell+2} = \ell + 1$.

Equation (2) shows that reconstructing a word uniquely from binomial coefficients of right-bounded-block words equates to solve a system of Diophantine equations. The knowledge of $k_{\mathbf{b}}, \dots, k_{\mathbf{a}^\ell \mathbf{b}}$ provides $\ell + 1$ equations. If the equation of $k_{\mathbf{a}^\ell \mathbf{b}}$ has a unique solution for $\{s_{\ell+1}, \dots, s_{k_{\mathbf{a}}+1}\}$ (in this case we say, by language abuse, that $k_{\mathbf{a}^\ell \mathbf{b}}$ is *unique*), then the system in row echelon form has a unique solution and thus the binary word is uniquely reconstructible. Notice that $k_{\mathbf{a}^{k_{\mathbf{a}}}} \mathbf{b}$ is always unique since $k_{\mathbf{a}^{k_{\mathbf{a}}}} \mathbf{b} = s_{k_{\mathbf{a}}+1}$.

Consider $n = 10$ and $k_{\mathbf{a}} = 4$. This leads to $w = \mathbf{b}^{s_1} \mathbf{a} \mathbf{b}^{s_2} \mathbf{a} \mathbf{b}^{s_3} \mathbf{a} \mathbf{b}^{s_4} \mathbf{a} \mathbf{b}^{s_5}$ with $\sum_{i \in [5]} s_i = 6$. Given $k_{\mathbf{a} \mathbf{b}} = 4$ we get $4 = s_2 + 2s_3 + 3s_4 + 4s_5$. The s_i are not uniquely determined. If $k_{\mathbf{a}^2 \mathbf{b}} = 2$ is also given, we obtain the equation $2 = s_3 + 3s_4 + 6s_5$ and thus $s_3 = 2$ and $s_4 = s_5 = 0$ is the only solution. Substituting these results in the previous equation leads to $s_2 = 0$ and since we only have six \mathbf{b} , we get $s_1 = 4$. Hence $w = \mathbf{b}^4 \mathbf{a}^2 \mathbf{b}^2 \mathbf{a}^2$ is uniquely reconstructed by $S = \{\mathbf{a}, \mathbf{a} \mathbf{b}, \mathbf{a}^2 \mathbf{b}\}$.

The following definition captures all solutions for the equation defined by $k_{\mathbf{a}^\ell \mathbf{b}}$ for $\ell \in [k_{\mathbf{a}}]_0$.

Definition 13. Set $M(k_{\mathbf{a}^\ell \mathbf{b}}) = \{(r_{\ell+1}, \dots, r_{k_{\mathbf{a}}+1}) \mid k_{\mathbf{a}^\ell \mathbf{b}} = \sum_{i=\ell+1}^{k_{\mathbf{a}}+1} \binom{i-1}{\ell} r_i\}$ for fixed $\ell \in [k_{\mathbf{a}}]_0$. We call $k_{\mathbf{a}^\ell \mathbf{b}}$ *unique* if $|M(k_{\mathbf{a}^\ell \mathbf{b}})| = 1$.

By Remark 12 the coefficients of each equation of the form (2) are strictly increasing. The next lemma provides the range each $k_{\mathbf{a}^\ell \mathbf{b}}$ may take under the constraint $\sum_{i=1}^{k_{\mathbf{a}}+1} s_i = n - k_{\mathbf{a}}$.

Lemma 14. Let $n \in \mathbb{N}$, $k \in [n]_0$, $j \in [k + 1]$ and $c_1, \dots, c_{k+1}, s_1, \dots, s_{k+1} \in \mathbb{N}_0$ with $c_i < c_{i+1}$, for $i \in [k]$, and $\sum_{i=1}^{k+1} s_i = n - k$. The sum $\sum_{i=j}^{k+1} c_i s_i$ is maximal iff $s_{k+1} = n - k$ (and consequently $s_i = 0$ for all $i \in [k]$).

Proof. The case $k = 0$ is trivial. Consider the case $n = k$, i.e., $\sum_{i=1}^{k+1} s_i = 0$. This implies immediately $s_i = 0$ for all $i \in [k+1]$ and the equivalence holds. Assume for the rest of the proof $k < n$. If $s_{k+1} = n - k$, then $s_i = 0$ for all $i \leq k$ and $\sum_{i=j}^{k+1} c_i s_i = c_{k+1}(n - k)$. Let us assume that the maximal value for $\sum_{i=j}^{k+1} c_i s_i$ can be obtained in another way and that there exist $s'_1, \dots, s'_{k+1} \in \mathbb{N}_0$, $\ell \in [n - k]$ such that $\sum_{i=1}^{k+1} s'_i = n - k$ and $s'_{k+1} = n - k - \ell$. Thus

$$c_{k+1}(n - k) \leq \sum_{i=j}^{k+1} c_i s'_i = \left(\sum_{i=j}^k c_i s'_i \right) + c_{k+1}(n - k - \ell).$$

This implies $\sum_{i=j}^k c_i s'_i \geq c_{k+1}\ell$. Since the coefficients are strictly increasing we get $\sum_{i=j}^k c_i s'_i \leq c_k \sum_{i=j}^k s'_i < c_{k+1}\ell$, hence the contradiction. \square

Corollary 15. *Let $k_a \in [n]_0$, $\ell \in [k_a]_0$, and $s_1, \dots, s_{k_a+1} \in \mathbb{N}_0$ with $\sum_{i=1}^{k_a+1} s_i = n - k_a$. Then $\binom{w}{a^{\ell b}} \in \left[\binom{k_a}{\ell}(n - k_a) \right]_0$.*

Proof. It follows directly from Equation (2) and Lemma 14. \square

The following lemma shows some cases in which $k_{a^{\ell b}}$ is unique.

Lemma 16. *Let $k_a \in [n]$, $\ell \in [k_a]_0$ and $s_1, \dots, s_{k_a+1} \in \mathbb{N}_0$ with $\sum_{i=1}^{k_a+1} s_i = n - k_a$. If $k_{a^{\ell b}} \in [\ell]_0 \cup \left\{ \binom{k_a}{\ell}(n - k_a) \right\}$ or $k_{a^{\ell b}} = \binom{k_a-1}{\ell}r + \binom{k_a}{\ell}(n - k_a - r)$ for $r \in [k_b]_0$ then $k_{a^{\ell b}}$ is unique.*

Proof. Consider firstly $k_{a^{\ell b}} \in [\ell]_0$. By Remark 12 we have $c_{\ell+1} = 1$ and $c_{\ell+2} = \ell + 1$. By $c_i < c_{i+1}$ we obtain immediately $s_i = 0$ for $i \in [k_a + 1] \setminus [\ell + 1]$. By setting $s_{\ell+1} = k_{a^{\ell b}}$ the claim is proven. If $k_{a^{\ell b}} = \binom{k_a}{\ell}(n - k_a)$, $s_{k_a+1} = (n - k_a)$ and $s_i = 0$ for $i \in [k_a]_0$ is the only possibility. Let secondly be $r \in [k_b]_0$ and $k_{a^{\ell b}} = \binom{k_a-1}{\ell}r + \binom{k_a}{\ell}(n - k_a - r)$ and suppose that $k_{a^{\ell b}}$ is not unique. This implies $s_{k_a+1} < n - k_a - r$. Assume that $s_{k_a+1} = n - k_a - r'$ for $r' \in [k_b]_{>r}$. Thus there exists $x \in \mathbb{N}$ with $\binom{k_a}{\ell}(n - k_a - r') + x = \frac{(k_a-1)!(k_a(n-k_a)-\ell r)}{\ell!(k_a-\ell)!}$, i.e., $x = \frac{(k_a-1)!(k_a r' - \ell r)}{\ell!(k_a-\ell)!}$. By $k_b = n - k_a$ we have $x \leq \binom{k_a-1}{\ell}r' = \frac{(k_a-1)!(k_a r' - \ell r)}{\ell!(k_a-\ell)!}$ (we only have r' occurrences of b left to distribute). By $r' > r$ we have $\frac{(k_a-1)!(k_a r' - \ell r)}{\ell!(k_a-\ell)!} = x < \frac{(k_a-1)!(k_a r' - \ell r)}{\ell!(k_a-\ell)!}$ - a contradiction. \square

Since we are not able to fully characterise the uniquely determined values for each $k_{a^{\ell b}}$ for arbitrary n and ℓ , the following proposition gives the characterisation for $\ell \in \{0, 1\}$. Notice that we use k_a immediately since it is determinable by n and $k_{a^0 b} = k_b$.

Proposition 17 (*). *The word $w \in \Sigma^n$ is uniquely determined by k_a and k_{ab} iff one of the following occurs*

- $k_a = 0$ or $k_a = n$ (and obviously $k_{ab} = 0$),

- $k_a = 1$ or $k_a = n - 1$ and k_{ab} is arbitrary,
- $k_a \in [n - 2]_{\geq 2}$ and $k_{ab} \in \{0, 1, k_a(n - k_a) - 1, k_a(n - k_a)\}$.

In all cases not covered by Proposition 17 the word cannot be uniquely determined by $\binom{w}{a}$ and $\binom{w}{ab}$. The following theorem combines the reconstruction of a word with the binomial coefficients of right-bounded-block words.

Theorem 18. *Let $j \in [k_a]_0$. If $k_{a^j b}$ is unique, then the word $w \in \Sigma^n$ is uniquely determined by $\{b, ab, a^2b, \dots, a^j b\}$.*

Proof. If $k_{a^j b}$ is unique, the coefficients $s_{j+1}, \dots, s_{k_a+1}$ are uniquely determined. Substituting backwards the known values in the first $j - 1$ equations (2) (for $\ell = 1, \dots, j - 1$) we can now obtain successively the values for s_j, \dots, s_1 . \square

Corollary 19. *Let ℓ be minimal such that $k_{a^\ell b}$ is unique. Then w is uniquely determined by $\{a, ab, a^2b, \dots, a^\ell b\}$ and not uniquely determined by any $\{a, ab, a^2b, \dots, a^i b\}$ for $i < \ell$.*

Proof. It follows directly from Theorem 18. \square

By [15] an upper bound on the number of binomial coefficients to uniquely reconstruct the word $w \in \Sigma^n$ is given by the amount of the binomial coefficients of the $(\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5)$ -spectrum. Notice that implicitly the full spectrum is assumed to be known. As proven in Section 2, Lyndon words up to this length suffice. Since there are $\frac{1}{n} \sum_{d|n} \mu(d) \cdot 2^{\frac{n}{d}}$ Lyndon words of length n , the combination of both results presented in [15, 22] states that, for $n > 6$,

$$\sum_{i=1}^{\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5} \frac{1}{i} \sum_{d|i} \mu(d) \cdot 2^{\frac{n}{d}} \quad (3)$$

binomial coefficients are sufficient for a unique reconstruction with the Möbius function μ . Up to now, it was the best known upper bound.

Theorem 18 shows that $\min\{k_a, k_b\} + 1$ binomial coefficients are enough for reconstructing a binary word uniquely. By Proposition 17 we need exactly one binomial coefficient if $n \in [3]$ and at most two if $n = 4$. For $n \in \{5, 6\}$ we need at most $n - 2$ different binomial coefficients. The following theorem shows that by Theorem 18 we need strictly less binomial coefficients for $n > 6$.

Theorem 20 (*). *Let $w \in \Sigma^n$. We have that $\min\{k_a, k_b\} + 1$ binomial coefficients suffice to uniquely reconstruct w . If $k_a \leq k_b$, then the set of sufficient binomial coefficients is $S = \{b, ab, a^2b, \dots, a^h b\}$ where $h = \lfloor \frac{n}{2} \rfloor$. If $k_a > k_b$, then the set is $S = \{a, ba, b^2a, \dots, b^h a\}$. This bound is strictly smaller than (3).*

Remark 21. By Lemma 16 we know that $k_{a^\ell b}$ is unique if it is in $[\ell]_0$ or exactly $\binom{k_a}{\ell}(n - k_a)$. The probability for the latter is $\frac{1}{2^n}$ for $w \in \{a, b\}^n$. If $k_{a^\ell b} = m \in [\ell]_0$ we get by (2) immediately $s_{\ell+1} = m$ and $s_i = 0$ for $\ell + 2 \leq i \leq k_a + 1$. Hence, the values for s_j for $j \in [\ell]$ are not determined. By $\sum_{i \in [\ell]} s_i = n - k_a - m$ there are $d = \sum_{i \in [\ell]_0} \binom{\ell}{i-1} \binom{n - k_a - m - 1}{i-1}$ possibilities to fulfill the constraints, i.e., we have a probability of $\frac{d}{2^n}$ to have such a word.

4 Reconstruction for Arbitrary Alphabets

In this section we address the problem of reconstructing words over arbitrary alphabets from their scattered factors. We begin with a series of results of algorithmic nature. Let $\Sigma = \{\mathbf{a}_1, \dots, \mathbf{a}_q\}$ be an alphabet equipped with the ordering $\mathbf{a}_i < \mathbf{a}_j$ for $1 \leq i < j \leq q \in \mathbb{N}$.

Definition 22. *Let $w_1, \dots, w_k \in \Sigma^*$ for $k \in \mathbb{N}$, and $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$ a sequence of $|\Sigma|$ natural numbers. A K -valid marking of w_1, \dots, w_k is a mapping $\psi : [k] \times \mathbb{N} \rightarrow \mathbb{N}$ such that for all $j \in [k]$, $i, \ell \in [|w_j|]$, and $\mathbf{a} \in \Sigma$ there holds*

- if $w_j[i] = \mathbf{a}$ then $\psi(j, i) \leq k_{\mathbf{a}}$,
- if $i < \ell \leq |w_j|$ and $w_j[i] = w_j[\ell] = \mathbf{a}$ then $\psi(j, i) < \psi(j, \ell)$.

A K -valid marking of w_1, \dots, w_k is represented as the string $w_1^\psi, w_2^\psi, \dots, w_k^\psi$, where $w_j^\psi[i] = (w_j[i])_{\psi(j,i)}$ for fresh letters $(w_j[i])_{\psi(j,i)}$.

For instance, let $k = 2$, $\Sigma = \{\mathbf{a}, \mathbf{b}\}$, and $w_1 = \mathbf{aab}$, $w_2 = \mathbf{abb}$. Let $k_{\mathbf{a}} = 3$, $k_{\mathbf{b}} = 2$ define the sequence K . A K -valid marking of w_1, w_2 would be $w_1^\psi = (\mathbf{a})_1(\mathbf{a})_3(\mathbf{b})_1$, $w_2^\psi = (\mathbf{a})_2(\mathbf{b})_1(\mathbf{b})_2$ defining ψ implicitly by the indices. We used parentheses in the marking of the letters in order to avoid confusions.

We recall that a topological sorting of a directed graph $G = (V, E)$, with $V = \{v_1, \dots, v_n\}$, is a linear ordering $v_{\sigma(1)} < v_{\sigma(2)} < \dots < v_{\sigma(n)}$ of the nodes, defined by the permutation $\sigma : [n] \rightarrow [n]$, such that there exists no edge in E from $v_{\sigma(i)}$ to $v_{\sigma(j)}$ for any $i > j$ (i.e., if v_a comes after v_b in the linear ordering, for some $a = \sigma(i)$ and $b = \sigma(j)$, then we have $i > j$ and there should be no edge between v_a and v_b). It is a folklore result that any directed graph G has a topological sorting if and only if G is acyclic.

Definition 23. *Let $w_1, \dots, w_k \in \Sigma^*$ for $k \in \mathbb{N}$, $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$ a sequence of $|\Sigma|$ natural numbers, and ψ a K -valid marking of w_1, \dots, w_k . Let G_ψ be the graph that has $\sum_{\mathbf{a} \in \Sigma} k_{\mathbf{a}}$ nodes, labelled with the letters $(\mathbf{a})_1, \dots, (\mathbf{a})_{k_{\mathbf{a}}}$, for all $\mathbf{a} \in \Sigma$, and the directed edges $((w_j[i])_{\psi(j,i)}, (w_j[i+1])_{\psi(j,i+1)})$, for all $j \in [k]$, $i \in [|w_j|]$, and $((\mathbf{a})_i, (\mathbf{a})_{i+1})$, for all occurring i and $\mathbf{a} \in \Sigma$. We say that there exists a valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k if there exists a topological sorting of the nodes of G_ψ , i.e., G_ψ is a directed acyclic graph.*

The graph associated with the K -valid marking of w_1, w_2 from above would have the five nodes $(\mathbf{a})_1, (\mathbf{a})_2, (\mathbf{a})_3, (\mathbf{b})_1, (\mathbf{b})_2$ and the six directed edges $((\mathbf{a})_1, (\mathbf{a})_3)$, $((\mathbf{a})_3, (\mathbf{b})_1)$, $((\mathbf{a})_2, (\mathbf{b})_1)$, $((\mathbf{b})_1, (\mathbf{b})_2)$, $((\mathbf{a})_1, (\mathbf{a})_2)$, $((\mathbf{a})_2, (\mathbf{a})_3)$ (where the direction of the edge is from the left node to the right node of the pair defining it). This graph has the topological sorting $(\mathbf{a})_1(\mathbf{a})_2(\mathbf{a})_3(\mathbf{b})_1(\mathbf{b})_2$.

Theorem 24 (*). *For $w_1, \dots, w_k \in \Sigma^*$ and a sequence $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$ of $|\Sigma|$ natural numbers, there exists a word w such that w_i is a scattered factor of w with $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $i \in [k]$ and all $\mathbf{a} \in \Sigma$, if and only if there exist a K -valid marking ψ of the words w_1, \dots, w_k and a valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k .*

Next we show that in Theorem 24 uniqueness propagates in the \leftarrow -direction.

Corollary 25. *Let $w_1, \dots, w_k \in \Sigma^*$ and $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$ a sequence of $|\Sigma|$ natural numbers. If the following hold*

- *there exists a unique K -valid marking ψ of the words w_1, \dots, w_k ,*
- *in the unique K -valid marking ψ we have that for each $\mathbf{a} \in \Sigma$ and $\ell \in [k_{\mathbf{a}}]$ there exists $i \in [k]$ and $j \in [|w_i|]$ with $\psi(i, j) = \ell$, and*
- *there exists a unique valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k*

then there exists a unique word w such that w_i is a scattered factor of w , for all $i \in [k]$ and $|w|_{\mathbf{a}} = k_{\mathbf{a}}$ for all $\mathbf{a} \in \Sigma$.

Proof. Let w be the word obtained by writing in order the letters of the unique valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k and removing their markings. It is clear that w' has w_i as a scattered factor, for all $i \in [k]$, and that $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $\mathbf{a} \in \Sigma$. The word w is uniquely defined (as there is no other K -valid marking nor valid topological sorting of the ψ -marked letters), and $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $\mathbf{a} \in \Sigma$. \square

In order to state the second result, we need the projection $\pi_S(w)$ of a word $w \in \Sigma^*$ on $S \subseteq \Sigma$: $\pi_S(w)$ is obtained from w by removing all letters from $\Sigma \setminus S$.

Theorem 26. *Set $W = \{w_{\mathbf{a}, \mathbf{b}} \mid \mathbf{a} < \mathbf{b} \in \Sigma\}$ such that*

- *$w_{\mathbf{a}, \mathbf{b}} \in \{\mathbf{a}, \mathbf{b}\}^*$ for all $\mathbf{a}, \mathbf{b} \in \Sigma$,*
- *for all $w, w' \in W$ and all $\mathbf{a} \in \Sigma$, if $|w|_{\mathbf{a}} \cdot |w'|_{\mathbf{a}} > 0$, then $|w|_{\mathbf{a}} = |w'|_{\mathbf{a}}$.*

Then there exists at most one $w \in \Sigma^$ such that $w_{\mathbf{a}, \mathbf{b}}$ is $\pi_{\{\mathbf{a}, \mathbf{b}\}}(w)$ for all $\mathbf{a}, \mathbf{b} \in \Sigma$.*

Proof. Notice firstly $|W| = \frac{q(q-1)}{2}$. Let $k_{\mathbf{a}} = |w_{\mathbf{a}, \mathbf{b}}|_{\mathbf{a}}$, for $\mathbf{a} < \mathbf{b} \in \Sigma$. These numbers are clearly well defined, by the second item in our hypothesis. Let $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$. It is immediate that there exists a unique K -valid marking ψ of the words $(w_{\mathbf{a}, \mathbf{b}})_{\mathbf{a} < \mathbf{b} \in \Sigma}$. As each two marked letters $(\mathbf{a})_i$ and $(\mathbf{b})_j$ (i.e., each two nodes $(\mathbf{a})_i$ and $(\mathbf{b})_j$ of G_{ψ}) appear in the marked word $w_{\mathbf{a}, \mathbf{b}}^{\psi}$, we know the order in which these two nodes should occur in a topological sorting of G_{ψ} . This means that, if G_{ψ} is acyclic, then it has a unique topological sorting. Our statement follows now from Corollary 25. \square

Remark 27. Given the set $W = \{w_{\mathbf{a}, \mathbf{b}} \mid \mathbf{a} < \mathbf{b} \in \Sigma\}$ as in the statement of Theorem 26, with $k_{\mathbf{a}} = |w_{\mathbf{a}, \mathbf{b}}|_{\mathbf{a}}$, for $\mathbf{a} < \mathbf{b} \in \Sigma$, and $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$, we can produce the unique K -valid marking ψ of the words $(w_{\mathbf{a}, \mathbf{b}})_{\mathbf{a} < \mathbf{b} \in \Sigma}$ in linear time $O(\sum_{\mathbf{a} < \mathbf{b} \in \Sigma} |w_{\mathbf{a}, \mathbf{b}}|) = O((q-1) \sum_{\mathbf{a} \in \Sigma} k_{\mathbf{a}})$: just replace the i^{th} letter \mathbf{a} of $w_{\mathbf{a}, \mathbf{b}}$ by $(\mathbf{a})_i$, for all \mathbf{a} and i . The graph G_{ψ} has $O((q-1) \sum k_{\mathbf{a}})$ edges and $O(\sum k_{\mathbf{a}})$ vertices and can be constructed in linear time $O((q-1) \sum k_{\mathbf{a}})$. Sorting G_{ψ} topologically takes $O((q-1) \sum k_{\mathbf{a}})$ time (see, e.g., the handbook [4]). As such, we conclude that reconstructing a word $w \in \Sigma^*$ from its projections over all two-letter-subsets of Σ can be done in linear time w.r.t. the total length of the respective projections.

Theorem 26 is in a sense optimal: in order to reconstruct a word over Σ uniquely, we need all its projections on two-letter-subsets of Σ . That is, it is always the case that for a strict subset U of $\{\{\mathbf{a}, \mathbf{b}\} \mid \mathbf{a} < \mathbf{b} \in \Sigma\}$, with $|U| = \frac{q(q-1)}{2} - 1$, there exist two words $w' \neq w$ such that $\{\pi_p(w') \mid p \in U\} = \{\pi_p(w) \mid p \in U\}$. We can, in fact, show the following results:

Theorem 28. *Let S_1, \dots, S_k be subsets of Σ . The following hold:*

1. *If each pair $\{\mathbf{a}, \mathbf{b}\} \subseteq \Sigma$ is included in at least one of the sets S_i , then we can reconstruct any word uniquely from its projections $\pi_{S_1}(\cdot), \dots, \pi_{S_k}(\cdot)$.*
2. *If there exists a pair $\{\mathbf{a}, \mathbf{b}\}$ that is not contained in any of the sets S_1, \dots, S_k , then there exist two words w and w' such that $w \neq w'$ and $\pi_{S_1}(w) = \pi_{S_1}(w'), \dots, \pi_{S_k}(w) = \pi_{S_k}(w')$.*

Proof. The first part is, once again, a consequence of Corollary 25. The second part can be shown by assuming that $\Sigma = \{\mathbf{a}_1, \dots, \mathbf{a}_q\}$ and the pair $\{\mathbf{a}_1, \mathbf{a}_2\}$ is not contained in any of the sets S_1, \dots, S_k . Then, for $w = \mathbf{a}_1 \mathbf{a}_3 \mathbf{a}_4 \dots \mathbf{a}_q$ and $w' = \mathbf{a}_2 \mathbf{a}_3 \mathbf{a}_4 \dots \mathbf{a}_q$, we have that $\pi_{S_1}(w) = \pi_{S_1}(w'), \dots, \pi_{S_k}(w) = \pi_{S_k}(w')$. \square

In this context, we can ask how efficiently can we decide if a word is uniquely reconstructible from the projections $\pi_{S_1}(\cdot), \dots, \pi_{S_k}(\cdot)$ for $S_1, \dots, S_k \subset \Sigma$.

Theorem 29 (*). *Given the sets $S_1, \dots, S_k \subset \Sigma$, we decide whether we can reconstruct any word uniquely from its projections $\pi_{S_1}(\cdot), \dots, \pi_{S_k}(\cdot)$ in $O(q^2 k)$ time. Moreover, under the Strong Exponential Time Hypothesis (see the survey [3] and the references therein), there is no $O(q^{2-d} k^c)$ algorithm for solving the above decision problem, for any $d, c > 0$.*

Coming now back to combinatorial results, we use the method developed in Section 3 to reconstruct a word over an arbitrary alphabet. We show that we need at most $\sum_{i \in [q]} |w|_i (q+1-i)$ different binomial coefficients to reconstruct w uniquely for the alphabet $\Sigma = \{1, \dots, q\}$. In fact, following the results from the first part of this section, we apply this method on all combinations of two letters. Consider for an example that for $w \in \{\mathbf{a}, \mathbf{b}, \mathbf{n}\}^6$ the following binomial coefficients $\binom{w}{\mathbf{a}^0 \mathbf{b}} = 1$, $\binom{w}{\mathbf{a}^0 \mathbf{n}} = 2$, $\binom{w}{\mathbf{a}^1 \mathbf{b}} = 0$, $\binom{w}{\mathbf{a}^1 \mathbf{n}} = 3$, $\binom{w}{\mathbf{b}^1 \mathbf{n}} = 2$, and $\binom{w}{\mathbf{a}^2 \mathbf{n}} = 1$ are given. By $|w| = 6$, $|w|_{\mathbf{b}} = 1$, and $|w|_{\mathbf{n}} = 2$, we get $|w|_{\mathbf{a}} = 3$. Applying the method from Section 3 for $\{\mathbf{a}, \mathbf{b}\}$, $\{\mathbf{a}, \mathbf{n}\}$, and $\{\mathbf{b}, \mathbf{n}\}$ we obtain the scattered factors \mathbf{ba}^3 , \mathbf{anana} , and \mathbf{bn}^2 . Combining all these three scattered factors gives us uniquely \mathbf{banana} . Notice that in this example we only needed six binomial coefficients instead of ten, which is the worst case.

Remark 30. As seen in the example we have not only the word length but also $\binom{w}{\mathbf{x}}$ for all $\mathbf{x} \in \Sigma$ but one. Both information give us the remaining single letter binomial coefficient and hence we will assume that we know all of them.

For convenience in the following theorem consider $\Sigma = \{1, \dots, q\}$ for $q > 2$ and set $\alpha := \lfloor \frac{16}{7} \sqrt{n} \rfloor + 5$. In the general case the results by [22] and [15] yield

that

$$\sum_{i \in [\alpha]} \frac{1}{i} \frac{(q+1)^{\frac{i}{2}} - 1}{q} \quad (4)$$

is smaller than the best known upper bound on the number of binomial coefficients sufficient to reconstruct a word uniquely.

The following theorem generalises Theorem 20 on an arbitrary alphabet.

Theorem 31 (*). *For reconstructing a word $w \in \Sigma^*$ of length at least $q-1$ uniquely, $\sum_{i \in [q]} |w|_i (q+1-i)$ binomial coefficients suffice, which is strictly smaller than (4).*

Remark 32. Since the estimation in Theorem 31 depends on the distribution of the letters in contrast to the method of reconstruction, it is wise to choose an order $<$ on Σ such that $x < y$ if $|w|_x \leq |w|_y$. In the example we have chosen the *natural* order $\mathbf{a} < \mathbf{b} < \mathbf{n}$ which leads in the worst case to fourteen binomial coefficients that has to be taken into consideration. If we chose the order $\mathbf{b} < \mathbf{n} < \mathbf{a}$ the formula from Theorem 31 provides that ten binomial coefficients suffice. This observation leads also to the fact that less binomial coefficients suffice for a unique determinism if the letters are not distributed equally but some letters occur very often and some only a few times.

Remark 33. Let's note that the number of binomial coefficients we need is at most qn . Indeed, we will prove that $\sum_{i \in [q]} |w|_i (q+1-i) \leq qn$. We have $qn = qn + n - n = q \sum_{i \in [q]} |w|_i + \sum_{i \in [q]} |w|_i - \sum_{i \in [q]} |w|_i \geq q \sum_{i \in [q]} |w|_i + \sum_{i \in [q]} |w|_i - \sum_{i \in [q]} (|w|_i i) = \sum_{i \in [q]} |w|_i (q+1-i)$.

5 Conclusion

In this paper we have proven that a relaxation of the so far investigated reconstruction problem from scattered factors from k -spectra to arbitrary sets yields that less scattered factors than the best known upper bound are sufficient to reconstruct a word uniquely. Not only in the binary but also in the general case the distribution of the letters plays an important role: in the binary case the amount of necessary binomial coefficients is smaller the larger $|w|_{\mathbf{a}} - |w|_{\mathbf{b}}$ is. The same observation results from the general case - if all letters are equally distributed in w then we need more binomial coefficients than in the case where some letters rarely occur and others occur much more often. Nevertheless the restriction to right-bounded-block words (that are intrinsically Lyndon words) shows that a word can be reconstructed by fewer binomial coefficients if scattered factors from different spectra are taken. Further investigations may lead into two directions: firstly a better characterisation of the uniqueness of the $k_{\mathbf{a} \neq \mathbf{b}}$ would be helpful to understand better in which cases less than the worst case amount of binomial coefficients suffices and secondly other sets than the right-bounded-block words could be investigated for the reconstruction problem.

References

1. J. Berstel, J. Karhumäki, Combinatorics on Words – A Tutorial, *Bull. Eur. Assoc. Theor. Comput. Sci. EATCS* **79** (2003), 178–228.
2. V. Berthé, J. Karhumäki, D. Nowotka, J. Shallit, Mini-Workshop: Combinatorics on Words, Oberwolfach Rep. **7** (2010), 2195–2244. doi: 10.4171/OWR/2010/37.
3. K. Bringmann: Fine-Grained Complexity Theory (Tutorial). 36th International Symposium on Theoretical Aspects of Computer Science (STACS 2019), R. Niedermeier, C. Paul Eds., *Leibniz International Proceedings in Informatics* **4** (2019), 1–7.
4. T. H. Cormen, C. E. Leiserson, R. L. Rivest, C. Stein, *Introduction to Algorithms*, 3rd Edition. MIT Press (2009).
5. A. W. M. Dress, P. L. Erdős, Reconstructing words from subwords in linear time, *Ann. Comb.* **8** (2004), 457–462.
6. M. Dudik, L. J. Schulman, Reconstruction from subsequences, *J. Combin. Theory, Ser. A* **103** (2003), 337–348.
7. P. L. Erdős, P. Ligeti, P. Sziklai, D. C. Torney, Subwords in reverse-complement order, *Ann. Comb.* **10** (2006), 415–430.
8. M. Ferov, Irreducible Polynomial modulo p , Bachelor Thesis at Charles University Prague (2008).
9. P. Fleischmann, M. Lejeune, F. Manea, D. Nowotka, M. Rigo, Reconstructing words from right-bounded-block words (2020), arXiv:2001.11218, 21 pages.
10. D. D. Freydenberger, P. Gawrychowski, J. Karhumäki, F. Manea, W. Rytter, Testing k -binomial equivalence. Multidisciplinary Creativity: homage to G. Păun on his 65th birthday, 239–248, Ed. Spandugino, Bucharest, Romania (2015). Available as: arXiv:1509.00622.
11. F. Harary, On the reconstruction of a graph from a collection of subgraphs. In *Theory of Graphs and its Applications* (Proc. Sympos. Smolenice, 1963). Publ. House Czechoslovak Acad. Sci., Prague (1964), 47–52.
12. L. van Iersel, V. Moulton, Leaf-reconstructibility of phylogenetic networks, *SIAM J. Discrete Math.* **32** (2018), 2047–2066.
13. L. I. Kalashnik, The reconstruction of a word from fragments, in “Numerical Mathematics and Computer Technology”, 56–57, Akad. Nauk Ukrain. SSR Inst. Mat., Preprint IV, (1973).
14. P. J. Kelly, A congruence theorem for trees, *Pacific J. Math.* **7** (1957), 961–968.
15. I. Krasikov, Y. Roditty, On a Reconstruction Problem for Sequences, *J. Combin. Theory, Ser. A* **77** (1997), 344–348.
16. V.I. Levenshtein, On perfect codes in deletion and insertion metric, *Discrete Math. Appl.* **2** (1992), 241–258.
17. M. Lothaire, *Combinatorics on Words*, Cambridge University Press (1997).
18. J. Manüch, Characterization of a Word by its Subwords, in *Developments in Language Theory* (1999), 210–219, World Scientific.
19. B. Manvel, A. Meyerowitz, A. Schwenk, K. Smith, P. Stockmeyer, Reconstruction of sequences, *Discrete Math.* **94** (1991), 209–219.
20. B. Manvel, P. K. Stockmeyer, On reconstruction of matrices, *Math. Mag.* **44** (1971), 218–221.
21. P. V. O’Neil, Ulam’s conjecture and graph reconstructions, *Amer. Math. Monthly* **77** (1970), 35–43.
22. C. Reutenauer, *Free Lie Algebras*, London Mathematical Society Monographs New Series, P. M. Cohn, H. G. Dales Eds., (1993).

23. M. Rigo, P. Salimov, Another generalization of abelian equivalence: Binomial complexity of infinite words, *Theoret. Comp. Sci.* **601** (2015), 47–57.
24. G. Rozenberg, A. Salomaa, *Handbook of Formal Languages* (3 volumes), Springer (1997).
25. I. Simon, Piecewise testable events, In: Automata Theory and Formal Languages, H. Brakhage ed., *Lect. Notes Comp. Sci.* **33** (1975), 214–222.

Appendix for Section 2

Example for Binomial Coefficients computed by Lyndon Words.

Let us recall that, for computing $\binom{w}{u}$, x and y are words such that $u = xy$ and such that every word appearing in the polynomial $x \sqcup y$ is lexicographically less than or equal to u .

Example 34. Considering $\Sigma = \{a, b\}$ the binomial coefficient $\binom{w}{ba}$ can be computed using the Lyndon words a , and b by

$$\binom{w}{ba} = \frac{1}{(b \sqcup a, ba)} \left[\binom{w}{b} \binom{w}{a} - (b \downarrow a, ab) \binom{w}{ab} \right] = \binom{w}{b} \binom{w}{a} - \binom{w}{ab}.$$

Regarding word length three, the Lyndon words are aab and abb . Let us give formulas to compute $\binom{w}{aba}$, $\binom{w}{baa}$, $\binom{w}{bab}$ and $\binom{w}{bba}$. Having $x = ab$ and $y = a$, we obtain

$$\binom{w}{aba} = \binom{w}{ab} \left[\binom{w}{a} - 1 \right] - 2 \binom{w}{aab}.$$

For $u = baa$, we can either choose $x = b$ and $y = aa$ or $x = ba$ and $y = a$. In the first case, we get

$$\binom{w}{baa} = \binom{w}{b} \binom{w}{aa} - \binom{w}{aba} - \binom{w}{aab}$$

and by reinjecting formulas for $\binom{w}{aa}$ and $\binom{w}{aba}$, obtained recursively,

$$\binom{w}{baa} = \left[\binom{w}{a} - 1 \right] \left[\frac{1}{2} \binom{w}{a} \binom{w}{b} - \binom{w}{ab} \right] + \binom{w}{aab}.$$

Finally, the last two formulas are quite similar to what we already had:

$$\binom{w}{bab} = \binom{w}{ab} \left[\binom{w}{b} - 1 \right] - 2 \binom{w}{abb}$$

and

$$\binom{w}{bba} = \left[\binom{w}{b} - 1 \right] \left[\frac{1}{2} \binom{w}{a} \binom{w}{b} - \binom{w}{ab} \right] + \binom{w}{abb}.$$

Appendix for Section 3

Proposition 17. *The word $w \in \Sigma^n$ is uniquely determined by k_a and k_{ab} iff one of the following occurs*

- $k_a = 0$ or $k_a = n$ (and obviously $k_{ab} = 0$),
- $k_a = 1$ or $k_a = n - 1$ and k_{ab} is arbitrary,
- $k_a \in [n - 2]_{\geq 2}$ and $k_{ab} \in \{0, 1, k_a(n - k_a) - 1, k_a(n - k_a)\}$.

Proof. Let us first prove that w is uniquely determined in these cases. It is obvious if $k_a = 0$ or $k_a = n$ since the word is composed of the same letter repeated n times. If $k_a = 1$, then $w = \mathbf{b}^{s_1} \mathbf{a} \mathbf{b}^{n-1-s_1}$ and $\binom{w}{\mathbf{a}\mathbf{b}} = n-1-s_1 = k_{\mathbf{a}\mathbf{b}}$. Therefore w is uniquely determined. If $k_a = n-1$, then $w = \mathbf{b}^{s_1} \mathbf{a} \mathbf{b}^{s_2} \dots \mathbf{a} \mathbf{b}^{s_n}$ with exactly one of the s_i being non zero and, in fact, equal to one. We have $\binom{w}{\mathbf{a}\mathbf{b}} = \sum_{i=2}^n (i-1)s_i$ and, if $k_{\mathbf{a}\mathbf{b}}$ is given (between 0 and $n-1$), then $s_{k_{\mathbf{a}\mathbf{b}}+1} = 1$ is the only non zero exponent. Consider now $k_a \in [n-2]_{\geq 2}$, i.e. $w = \mathbf{b}^{s_1} \mathbf{a} \mathbf{b}^{s_2} \dots \mathbf{b}^{s_{k_a}} \mathbf{a} \mathbf{b}^{s_{k_a+1}}$. Thus $k_{\mathbf{a}\mathbf{b}} = 0$ implies $s_1 = n - k_a$ and $s_2 = 0, \dots, s_{k_a+1} = 0$ while $k_{\mathbf{a}\mathbf{b}} = 1$ implies $s_2 = 1, s_1 = n - k_a - 1$ and $s_3 = 0, \dots, s_{k_a+1} = 0$. By Lemma 14, we know that (2) is maximal if and only if $s_{k_a+1} = n - k_a$ and all the other s_i are equal to zero. In that case, the value of the sum equals $k_a(n - k_a)$. Therefore, if $\binom{w}{\mathbf{a}\mathbf{b}} = k_a(n - k_a)$, the word w is uniquely determined. Finally, if $k_{\mathbf{a}\mathbf{b}} = k_a(n - k_a) - 1$, we must have $s_{k_a+1} \leq n - k_a - 1$. If we choose $s_{k_a+1} = n - k_a - 1$, it remains that $\sum_{i=1}^{k_a} s_i = 1$ and $\sum_{i=2}^{k_a} (i-1)s_i = k_a - 1$. We must have $s_{k_a} = 1$ and the other ones equal to zero. In fact, choosing $s_{k_a+1} = n - k_a - 1$ is the only possibility: if otherwise $s_{k_a+1} = n - k_a - \ell$ with $\ell > 1$, we obtain that $\sum_{i=2}^{k_a} (i-1)s_i \geq \ell k_a - 1$ with $\sum_{i=1}^{k_a} s_i = \ell$. It is easy to check with Lemma 14 that these conditions are incompatible.

We now need to prove that w cannot be uniquely determined if $k_a \in [n-2]_{\geq 2}$ and $k_{\mathbf{a}\mathbf{b}} \in [k_a(n - k_a) - 2]_{\geq 2}$. To this aim we will give two different sets of values for the s_i . The first decomposition is the greedy one. Let us put $s_{k_a+1} = \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor$, $s_{(k_{\mathbf{a}\mathbf{b}} \bmod k_a)+1} = 1$ and the other s_i equal to 0. Let us finally modify the value of s_1 (which is, at this stage, equal to 0 or 1) by adding the value needed. By $\sum_{i=1}^{k_a+1} s_i = n - k_a$ we get $s_1 \leftarrow s_1 + (n - k_a) - \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor - 1$. This implies $\sum_{i=1}^{k_a+1} s_i = 1 + (n - k_a) - \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor - 1 + \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor = n - k_a$ and $s_i \geq 0$ for all i . Moreover we have $\sum_{i=2}^{k_a+1} (i-1)s_i = (k_{\mathbf{a}\mathbf{b}} \bmod k_a) + k_a \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor = k_{\mathbf{a}\mathbf{b}}$.

Now we provide a second decomposition for the s_i . First, let us assume that $2 \leq k_{\mathbf{a}\mathbf{b}} < k_a$. In that case, the greedy algorithm sets $s_{k_{\mathbf{a}\mathbf{b}}+1} = 1$, $s_1 = n - k_a - 1$ and the other s_i to 0. Let us now set $s_1 = n - k_a - 2$ and all the other s_i to 0. Then, update $s_{k_{\mathbf{a}\mathbf{b}}} \leftarrow s_{k_{\mathbf{a}\mathbf{b}}} + 1$ and $s_2 \leftarrow s_2 + 1$ (in the case where $k_{\mathbf{a}\mathbf{b}} = 2$, s_2 will be equal to 2 after these manipulations). We have that the sum in (2) is equal to $1 + (k_{\mathbf{a}\mathbf{b}} - 1)$ as needed. Finally, if $k_{\mathbf{a}\mathbf{b}} \geq k_a$, then s_{k_a+1} was non zero in the greedy decomposition, and the idea is to reduce it of a value 1. Let us set $s_{k_a+1} = \lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor - 1$ and the other s_i to 0. Then, let us update some values: $s_{(k_{\mathbf{a}\mathbf{b}} \bmod k_a)+2} \leftarrow s_{(k_{\mathbf{a}\mathbf{b}} \bmod k_a)+2} + 1$ and $s_{k_a} \leftarrow s_{k_a} + 1$ if $(k_{\mathbf{a}\mathbf{b}} \bmod k_a) \neq k_a - 1$, and $s_{k_a} = 2, s_2 = 1$ otherwise. Finally, set s_1 to the right value, i.e., $n - k_a - \sum_{i=2}^{k_a+1} s_i$. It can be easily checked that, in both cases, $s_1 \geq 0$ (notice that $(k_{\mathbf{a}\mathbf{b}} \bmod k_a) = k_a - 1$ implies that $\lfloor \frac{k_{\mathbf{a}\mathbf{b}}}{k_a} \rfloor \leq n - k_a - 2$) and that all s_i sum up to $n - k_a$. Similarly, we can check that $\sum_{i=2}^{k_a+1} (i-1)s_i$ is equal to $k_{\mathbf{a}\mathbf{b}}$ in both cases.

To sum up, we gave two different decompositions for the s_i in cases where $k_a \in [n-2]_{\geq 2}$ and $k_{\mathbf{a}\mathbf{b}} \in [k_a(n - k_a) - 2]_{\geq 2}$. That implies that w cannot be uniquely determined in those cases. \square

Theorem 20. Let $w \in \Sigma^n$. We have that $\min\{k_a, k_b\} + 1$ binomial coefficients suffice to uniquely reconstruct w . If $k_a \leq k_b$, then the set of sufficient binomial coefficients is $S = \{\mathbf{b}, \mathbf{ab}, \mathbf{a}^2\mathbf{b}, \dots, \mathbf{a}^h\mathbf{b}\}$ where $h = \lfloor \frac{n}{2} \rfloor$. If $k_a > k_b$, then the set is $S = \{\mathbf{a}, \mathbf{ba}, \mathbf{b}^2\mathbf{a}, \dots, \mathbf{b}^h\mathbf{a}\}$. This bound is strictly smaller than (3).

Proof. Assume w.l.o.g. $k_a \leq k_b$. Then $k_a \leq \frac{n}{2}$ and Theorem 18 shows that words in the set $\{\mathbf{b}, \mathbf{ab}, \dots, \mathbf{a}^{\lfloor \frac{n}{2} \rfloor} \mathbf{b}\}$ can reconstruct w uniquely. If $k_a > k_b$, the set S is obtained by replacing the letter \mathbf{a} by \mathbf{b} and vice-versa.

Set $N_2(i) := \frac{1}{i} \sum_{d|i} \mu(d) 2^{\frac{i}{d}}$ for all $i \in [\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5]$, i.e., $\sum_{i=1}^{\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5} N_2(i)$, which is Equation (3), binomial coefficients suffice. By [8, Lemma 2.4] we have

$$N_2(i) \geq \frac{1}{i} \left(2^i - \frac{2^{\frac{i}{2}} - 1}{2 - 1} \right) = \frac{1}{i} \left(2^i - 2^{\frac{i}{2}} + 1 \right) = \frac{1}{i} \left(2^{\frac{i}{2}} (2^{\frac{i}{2}} - 1) + 1 \right) \geq \frac{2^{\frac{i}{2}}}{i}.$$

This results in

$$\sum_{i=1}^{\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5} N_2(i) \geq \sum_{i=1}^{\lfloor \frac{16}{7} \sqrt{n} \rfloor + 5} \frac{2^{\frac{i}{2}}}{i} \geq \frac{1}{\frac{16}{7} \sqrt{n} + 5} \frac{\sqrt{2}^{\frac{16}{7} \sqrt{n} + 5} - 1}{\sqrt{2} - 1}.$$

We want to show that this quantity is at least equal to $\frac{n+1}{2}$. Let us define

$$f(x) = \frac{1}{\frac{16}{7} \sqrt{x} + 5} \frac{\sqrt{2}^{\frac{16}{7} \sqrt{x} + 5} - 1}{\sqrt{2} - 1} - \frac{x + 1}{2}$$

for all $x > 0$, which is the continuous extension on \mathbb{R}^+ of the quantity we are interested in. It is easy to verify by hand that $f(1), f(2), f(3)$ and $f(4)$ are positive. Let us formally show that $f(x) > 0$ for all $x \geq 5$. Since this function is differentiable, we get with $y = \frac{16}{7} \sqrt{x} + 5$

$$f'(x) = \frac{1}{y} \sqrt{2}^y \frac{\ln(\sqrt{2})}{\sqrt{2} - 1} \frac{8}{7\sqrt{x}} - \frac{1}{y^2} \frac{8}{7\sqrt{x}} \frac{\sqrt{2}^y - 1}{\sqrt{2} - 1} - \frac{1}{2}.$$

We thus have $\sqrt{x} = \frac{7y-35}{16}$ and $y \geq 7$ for all $x \geq 1$. By injecting y in the previous expression, and reducing to the common denominator, we have to show that

$$\begin{aligned} & 2y\sqrt{2}^y 128 \ln(\sqrt{2}) - 256(\sqrt{2}^y - 1) - 7(7y - 35)y^2(\sqrt{2} - 1) \\ &= \sqrt{2}^y (128 \ln(2)y - 256) + 256 - 49y^3(\sqrt{2} - 1) + 245y^2(\sqrt{2} - 1) \\ &\geq \sqrt{2}^y 365 + 256 - 49y^3(\sqrt{2} - 1) + 245y^2(\sqrt{2} - 1) \end{aligned}$$

is strictly positive. Let us call the last quantity $g(y)$. We will show that it is positive for all $y \geq 10.05$, which means that $f(x)$ is positive for all x such that $\frac{16}{7} \sqrt{x} + 5 \geq 10.05$, i.e., for all $x \geq 5$. We have

$$\begin{aligned} g'(y) &= 365\sqrt{2}^y \ln(\sqrt{2}) - 147(\sqrt{2} - 1)y^2 + 490(\sqrt{2} - 1)y, \\ g''(y) &= 365\sqrt{2}^y (\ln(\sqrt{2}))^2 - 294(\sqrt{2} - 1)y + 490(\sqrt{2} - 1), \\ g'''(y) &= 365\sqrt{2}^y (\ln(\sqrt{2}))^3 - 294(\sqrt{2} - 1), \end{aligned}$$

and $g'''(7) > 50$, $g''(8.5) > 2$, $g'(10.05) > 8$ and finally $g(10.05) > 1787$. Since $g'''(y)$ is increasing and positive in 7, $g''(y)$ is increasing for $y \geq 7$. Therefore $g'(y)$ is increasing for $y \geq 8.5$ and finally $g(y)$ is increasing for $y \geq 10.05$ and positive. \square

Appendix for Section 4

Theorem 24. *For $w_1, \dots, w_k \in \Sigma^*$ and a sequence $K = (k_{\mathbf{a}})_{\mathbf{a} \in \Sigma}$ of $|\Sigma|$ natural numbers, there exists a word w such that w_i is a scattered factor of w with $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $i \in [k]$ and all $\mathbf{a} \in \Sigma$, if and only if there exist a K -valid marking ψ of the words w_1, \dots, w_k and a valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k .*

Proof. If w is such that w_i is a scattered factor of w , for all $i \in [k]$, and $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $\mathbf{a} \in \Sigma$, then we can mark the i^{th} occurrence of \mathbf{a} as $(\mathbf{a})_i$, for all $\mathbf{a} \in \Sigma$ and $i \in [k_{\mathbf{a}}]$. This induces a K -valid marking ψ of the words w_i , and, moreover, the linear ordering of the nodes of G_ψ induced by the order in which the marked letters (i.e., nodes of G_ψ) occur in w is a topological sorting of G_ψ .

Let us now assume that there exists a K -valid marking ψ of the words w_1, \dots, w_k , and there exists a valid topological sorting of the ψ -marked letters of the words w_1, \dots, w_k . Let w' be the word obtained by writing the nodes of G_ψ in the order given by its topological sorting and removing their markings. It is clear that w' has w_i as a scattered factor, for all $i \in [k]$, and that $|w'|_{\mathbf{a}} \leq k_{\mathbf{a}}$, for all $\mathbf{a} \in \Sigma$. Let now $w = w' \prod_{\mathbf{a} \in \Sigma} \mathbf{a}^{k_{\mathbf{a}} - |w'|_{\mathbf{a}}}$, where $\prod_{\mathbf{a} \in \Sigma} \mathbf{a}^{k_{\mathbf{a}} - |w'|_{\mathbf{a}}}$ is the concatenation of the factors $\mathbf{a}^{k_{\mathbf{a}} - |w'|_{\mathbf{a}}}$, for $\mathbf{a} \in \Sigma$ in some fixed order. Now w has w_i as a scattered factor, for all $i \in [k]$, and $|w|_{\mathbf{a}} = k_{\mathbf{a}}$, for all $\mathbf{a} \in \Sigma$. \square

Theorem 29. *Given the sets $S_1, \dots, S_k \subset \Sigma$, we decide whether we can reconstruct any word uniquely from its projections $\pi_{S_1}(\cdot), \dots, \pi_{S_k}(\cdot)$ in $O(q^2 k)$ time. Moreover, under the Strong Exponential Time Hypothesis (see the survey [3] and the references therein), there is no $O(q^{2-d} k^c)$ algorithm for solving the above decision problem, for any $d, c > 0$.*

Proof. We begin with a series of preliminaries. Let us recall the *Orthogonal Vectors* problem: Given sets A, B consisting of n vectors in $\{0, 1\}^k$, decide whether there are vectors $a \in A$ and $b \in B$ which are orthogonal (i.e., for any $i \in [k]$ we have $a[i]b[i] = 0$). This problem can be solved naïvely in $O(n^2 k)$ time, but under the *Strong Exponential Time Hypothesis* there is no $O(n^{2-d} k^c)$ algorithm for solving it, for any $d, c > 0$ (once more, see the survey [3] and the references therein).

We show that our problem is equivalent to the *Orthogonal Vectors* problem.

Let us first assume that we are given the sets $S_1, \dots, S_k \subset \Sigma$, and we want to decide whether we can reconstruct any word uniquely from its projections $\pi_{S_1}(\cdot), \dots, \pi_{S_k}(\cdot)$. This is equivalent, according to Theorem 28, to checking whether each pair $\{\mathbf{a}, \mathbf{b}\} \subseteq \Sigma$ is included in at least one of the sets S_i . For each letter \mathbf{a} of Σ we define the k -dimensional vectors $x_{\mathbf{a}}$ where $x_{\mathbf{a}}[i] = 1$ if $\mathbf{a} \in S_i$ and

$x_{\mathbf{a}}[i] = 0$ if $\mathbf{a} \notin S_i$. This can be clearly done in $O(qk)$ time. Now, there exists a pair $\{\mathbf{a}, \mathbf{b}\}$ that is not contained in any of the sets S_1, \dots, S_k if and only if there exists a pair of vectors $\{x_{\mathbf{a}}, x_{\mathbf{b}}\}$ such that $x_{\mathbf{a}}[i]x_{\mathbf{b}}[i] = 0$ for all $i \in [k]$. We can check whether there exists a pair of vectors $\{x_{\mathbf{a}}, x_{\mathbf{b}}\}$ such that $x_{\mathbf{a}}[i]x_{\mathbf{b}}[i] = 0$ for all $i \in [k]$ by solving the *Orthogonal Vectors* problem by using for both input sets of vectors the set $\{x_{\mathbf{a}} \mid \mathbf{a} \in \Sigma\}$. As such, we can check whether there exists a pair $\{\mathbf{a}, \mathbf{b}\}$ that is not contained in any of the sets S_1, \dots, S_k in $O(q^2k)$ time.

Let us now assume that we are given two sets A, B consisting of n vectors in $\{0, 1\}^k$, and we want to decide whether there are vectors $a \in A$ and $b \in B$ which are orthogonal (i.e., for any $i \in [k]$ we have $a[i]b[i] = 0$). We can compute the set of $(k+2)$ -dimensional vectors A' containing the vectors of A extended with two new positions (position $k+1$ and position $k+2$) set to 10 and the vectors of B extended with two new positions (position $k+1$ and position $k+2$) set to 01. To decide whether there are vectors $a \in A$ and $b \in B$ which are orthogonal is equivalent to decide whether there are vectors $a, b \in A'$ which are orthogonal (if two such vectors exist, they must be different on their last two positions, so one must come from A and one from B). Assume that $A' = \{x_1, x_2, \dots, x_{2n}\}$. Now we define an alphabet $\Sigma = \{\mathbf{a}_1, \dots, \mathbf{a}_{2n}\}$ of size $2n$ and the sets S_1, \dots, S_k , where $\mathbf{a}_j \in S_i$ if and only if $x_j[i] = 1$. Computing A' and then the alphabet Σ and the sets S_i , for $i \in [k]$, takes $O(nk)$ time. Now, to decide whether there are vectors $x_i, x_j \in A'$ which are orthogonal is equivalent to decide whether there exists a pair of letters $\{\mathbf{a}_i, \mathbf{a}_j\}$ of Σ that is not contained in any of the sets S_1, \dots, S_k . The conclusion of the theorem now follows. \square

Theorem 31. *For reconstructing a word $w \in \Sigma^*$ of length at least $q-1$ uniquely, $\sum_{i \in [q]} |w|_i(q+1-i)$ binomial coefficients suffice, which is strictly smaller than (4).*

Proof. The claim that $\sum_{i \in [q]} |w|_i(q+1-i)$ binomial coefficients suffice to reconstruct w uniquely follows by Theorem 26: for each pair of letters we apply the method of the binary case. We are thus going to reconstruct words $w_{\mathbf{a}, \mathbf{b}}$ for all pairs of letters $\mathbf{a} < \mathbf{b}$. To determine such a word uniquely, $\min(k_{\mathbf{a}}, k_{\mathbf{b}}) + 1 \leq k_{\mathbf{a}} + 1$ binomial coefficients suffice. In total, we thus need $\sum_{i \in [q]} (|w|_i + 1)(q-i)$ binomial coefficients. Among them are the ones related to letters of $\Sigma \setminus \{q\}$. The letter i is counted $(q-i)$ times in the above formula. So in fact $\sum_{i \in [q]} |w|_i(q-i) + (q-1)$ binomial coefficients suffice. This quantity is less than or equal to $\sum_{i \in [q]} |w|_i(q+1-i)$ for every w of length at least $q-1$.

We show the second claim by induction on q where the binary case in Theorem 20 serves as induction basis. This implies

$$\begin{aligned} \sum_{i \in [q]} |w|_i(q+1-i) &= \left(\sum_{i \in [q-1]} |w|_i(q+1-i) \right) + |w|_q(q+1-q) \\ &= \left(\sum_{i \in [q-1]} |w|_i(q-i) \right) + \sum_{i \in [q-1]} |w|_i + |w|_q \end{aligned}$$

and therefore

$$\begin{aligned} \sum_{i \in [q]} |w|_i (q+1-i) &\leq \left(\sum_{i \in [\alpha]} \frac{1}{i} \frac{q^{\frac{i}{2}} - 1}{q-1} \right) + n \\ &= \left(\sum_{i \in [\alpha]} \frac{1}{i} \frac{q(q^{\frac{i}{2}} - 1)}{q(q-1)} \right) + n. \end{aligned}$$

On the other hand, we have to compare this quantity with (4) which can be rewritten as

$$\sum_{i \in [\alpha]} \frac{1}{i} \frac{(q+1)^{\frac{i}{2}} - 1}{q} = \sum_{i \in [\alpha]} \frac{1}{i} \frac{(q-1)((q+1)^{\frac{i}{2}} - 1)}{q(q-1)}.$$

Thus the claim is proven, if the subtraction of the latter one and the previous one is greater than zero, i.e., we show that

$$\left(\sum_{i \in [\alpha]} \frac{1}{i} \frac{(q-1)((q+1)^{\frac{i}{2}} - 1) - q(q^{\frac{i}{2}} - 1)}{q(q-1)} \right) - n > 0, \text{ i.e.} \quad (5)$$

$$\left(\sum_{i \in [\alpha]} \frac{1}{i} \frac{(q-1)(q+1)^{\frac{i}{2}} - qq^{\frac{i}{2}} + 1}{q(q-1)} \right) - n > 0. \quad (6)$$

With $f(i) = \frac{(q-1)(q+1)^{\frac{i}{2}} - qq^{\frac{i}{2}} + 1}{iq(q-1)}$ for all $i \in [\alpha]$, the proof of (6) contains the following steps

1. For all $i \geq 2$ we have $f(i) \geq 0$,
2. $f(5) + f(1) \geq 0$,
3. $f(\alpha) - n > 0$.

Step 1.: For $i = 2$ we have

$$f(2) = \frac{1}{2} \frac{(q-1)(q+1) - q^2 + 1}{q(q-1)} = \frac{q^2 - 1 - q^2 + 1}{2q(q-1)} = 0.$$

For $i = 3$ we have

$$f(3) = \frac{1}{3} \frac{(q-1)(q+1)\sqrt{q+1} - q^2\sqrt{q} + 1}{q(q-1)}.$$

Consider the function $g : \mathbb{R} \rightarrow \mathbb{R}; q \mapsto q^4 - 2q^3 - 2q^2 + q + 1$. This function has two minima (between -0.75 and -0.5 as well as between 1.75 and 2) and one maximum (between 0.125 and 0.25). Since g has only two inflexion points and g is strictly greater than zero at the first minima, g has only two roots. The first root is between 0.7 and 0.8 and the second root is between 2.5 and 2.75 . Thus

for all $q \geq 2.75$ we have $g(q) > 0$. This implies $q^5 + q^4 - 2q^3 - 2q^2 + q + 1 > q^5$. Hence equivalently we get $(q+1)(q^4 - 2q^2 + 1) > q^5$, i.e., $(q+1)(q^2 - 1)^2 > qq^4$. This implies $\sqrt{q+1}(q^2 - 1) > \sqrt{q}q^2$ which proves that the numerator of $f(3)$ is positive and hence $f(3) > 0$. Before we prove the claim for $i \geq 4$, we will prove that $(q-1)(q+1)^j \geq q^{j+1}$ for $j \geq 2$. Firstly we get

$$(q-1)(q+1)^j = \left(\sum_{k \in [j]} \left(\binom{j}{k-1} - \binom{j}{k} \right) q^k \right) + q^{j+1} - 1.$$

Due to the central symmetry of each row of the Pascal triangle and since the distribution of the binomial coefficient is unimodal, for $k \leq \lfloor j/2 \rfloor$, we have

$$\binom{j}{j-k} - \binom{j}{j-k+1} = - \left(\binom{j}{k-1} - \binom{j}{k} \right) > 0$$

and thus

$$(q-1)(q+1)^j = \left(\sum_{k \in [\lfloor j/2 \rfloor]} \left(\binom{j}{k} - \binom{j}{k-1} \right) (q^{j-k+1} - q^k) \right) + q^{j+1} - 1.$$

Since $k \leq \lfloor j/2 \rfloor$, we have $j-k+1 > k$ and each term of the above sum is thus positive. This shows that $(q-1)(q+1)^j \geq q^{j+1}$. This leads to the following estimations for $f(i)$. For $i = 2j$ and $j \geq 2$ we get

$$f(i) = \frac{(q-1)(q+1)^j - qq^j + 1}{iq(q-1)} \geq \frac{q^{j+1} - q^{j+1} + 1}{iq(q-1)} > 0.$$

Finally for $i = 2j + 1$ and $j \geq 2$ we get

$$f(i) = \frac{(q-1)(q+1)^j \sqrt{q+1} - qq^j \sqrt{q} + 1}{iq(q-1)} \geq \frac{q^{j+1}(\sqrt{q+1} - \sqrt{q}) + 1}{iq(q-1)} > 0.$$

Step 2.: Notice that $\alpha \geq 7$ holds and thus $f(5)$ is always a summand. For $f(5) + f(1)$ we have to prove

$$\frac{(q-1)(q+1)^2 \sqrt{q+1} - qq^2 \sqrt{q} + 1}{5q(q-1)} + \frac{(q-1)\sqrt{q+1} - q\sqrt{q} + 1}{q(q-1)} \geq 0$$

Thus we get for the numerator

$$\begin{aligned} & (q-1)(q+1)^2 \sqrt{q+1} - qq^2 \sqrt{q} + 1 + 5(q-1)\sqrt{q+1} - 5q\sqrt{q} + 5 \\ &= (q-1)\sqrt{q+1}((q+1)^2 + 5) - q\sqrt{q}(q^2 + 5) + 6 \\ &= (q-1)\sqrt{q+1}(q^2 + 2q + 6) - q\sqrt{q}(q^2 + 5) + 6 \\ &= q^3 \sqrt{q+1} + q^2 \sqrt{q+1} + 4q\sqrt{q+1} - 6\sqrt{q+1} - q^3 \sqrt{q} - 5q\sqrt{q} + 6. \end{aligned}$$

We have $q^3\sqrt{q+1} > q^3\sqrt{q}$ and, since $q \geq 3$,

$$q^2\sqrt{q+1} \geq (6+q)\sqrt{q+1}.$$

Therefore $q^2\sqrt{q+1} + 4q\sqrt{q+1} \geq 6\sqrt{q+1} + 5q\sqrt{q}$ and the numerator is positive.

Step 3.: Notice that for fixed i , $f(i)$ is monotonically increasing for increasing q . This implies

$$f(\alpha) \geq \frac{2 \cdot 4^{\frac{\alpha}{2}} - 3 \cdot 3^{\frac{\alpha}{2}} + 1}{6\alpha} = \frac{2^{\alpha+1} - 3^{\frac{\alpha}{2}+1} + 1}{6\alpha}.$$

We are going to prove that

$$2^{\alpha+1} - 3^{\frac{\alpha}{2}+1} + 1 > 6\alpha n. \quad (7)$$

Recall that α is a function of n , given by $\alpha = \lfloor \frac{16}{7}\sqrt{n} \rfloor + 5$.

First, we have

$$2^{\alpha+1} - 3^{\frac{\alpha}{2}+1} > 2^{\alpha-1} - 2^{\frac{\alpha}{2}}.$$

Indeed, this inequality is equivalent to

$$2^{\frac{\alpha}{2}} (3 \cdot 2^{\frac{\alpha}{2}-1} + 1) > 3^{\frac{\alpha}{2}+1} \Leftrightarrow 2^{\frac{\alpha}{2}-1} + \frac{1}{3} > \left(\frac{3}{2}\right)^{\frac{\alpha}{2}}.$$

We can check that this last inequality is true by taking the logarithm of both sides, since $\alpha > 5$.

Therefore, it is sufficient for (7) to show that

$$2^{\alpha-1} - 2^{\frac{\alpha}{2}} = 2^{\frac{\alpha}{2}} (2^{\frac{\alpha}{2}-1} - 1) > 6\alpha n.$$

Note that $2^{\frac{\alpha}{2}} > n$ (indeed, $\lfloor \frac{16}{7}\sqrt{n} \rfloor + 5 > 2\sqrt{n} + 5$, thus $2^{\frac{\alpha}{2}} > 2^{\frac{5}{2}} \cdot 2^{\sqrt{n}}$). Once again, taking the logarithms, one can check that $2^{\frac{5}{2}} \cdot 2^{\sqrt{n}} > n$ holds.

To verify (7), it remains to show that $2^{\frac{\alpha}{2}-1} - 1 \geq 6\alpha$ or that $2^{\frac{\alpha}{2}-1} > 6\alpha$. Taking the logarithms, it is equivalent to

$$\begin{aligned} \frac{\alpha}{2} - 1 &> \log(6) + \log(\alpha) \\ \Leftrightarrow \alpha - 2\log(\alpha) &> 2\log(6) + 2, \end{aligned}$$

which is true for $\alpha \geq 15$, that is for $n \geq 16$. Equation (7) can be verified by a computer for $q-1 \leq n < 16$.

By 1., 2., and 3. Equation (6) is proven and this proves the claim. \square