A one-sided Zimin construction

L. J. Cummings

University of Waterloo ljcummings@math.uwaterloo.ca

M. Mays

West Virginia University mays@math.wvu.edu

Submitted: December 1, 2000; Accepted: July 23, 2001. MR Subject Classifications: 68R15, 20M35

Abstract

A string is Abelian square-free if it contains no Abelian squares; that is, adjacent substrings which are permutations of each other. An Abelian square-free string is maximal if it cannot be extended to the left or right by concatenating alphabet symbols without introducing an Abelian square. We construct Abelian square-free finite strings which are maximal by modifying a construction of Zimin. The new construction produces maximal strings whose length as a function of alphabet size is much shorter than that in the construction described by Zimin.

1 Introduction

Strings are a fundamental data structure. Equivalent names include: sequence, word, vector, codeword, linear array, and list. We take the entries of our strings to be elements of a finite set $A = \{a_0, \ldots, a_m\}$ called the alphabet. The elements of A will be called *entries* or *letters*. Strings may be infinite or finite. Considerable research effort has been directed toward determining those countably infinite strings which do or do not exhibit certain properties, but here we will be concerned with finite strings. Any ordered sequence $\mathbf{x} = b_1 b_2 \cdots b_n$ of elements chosen from A is called a finite *string* of length $|\mathbf{x}| = n$ over A.

In the interest of notational convenience, and without loss of generality, we often choose $A = \{0, \ldots, m\}$ as the alphabet. Every element of the alphabet is also considered to be a string. Two strings $\mathbf{x} = a_1 a_2 \cdots a_p$ and $\mathbf{y} = b_1 b_2 \cdots b_q$ are equal if and only if p = q and $a_i = b_i$ for $i = 1, \ldots, p$. For each $a \in A$ we define the integer-valued function $|\mathbf{x}|_a$ to be the number of times that a appears in the string \mathbf{x} . The (m+1)-tuple $\sharp \mathbf{x} = [|\mathbf{x}|_{a_0}, |\mathbf{x}|_{a_1}, \cdots, |\mathbf{x}|_{a_m}]$ is called the *frequency vector* of \mathbf{x} . We freely concatenate strings and write the concatenation of strings \mathbf{x} and \mathbf{y} as simply \mathbf{xy} . With this operation, A^* , the set of finite strings over A, has an algebraic structure called the free monoid over A but we do not use this fact here. If $1 \leq i \leq j \leq n$ the ordered sequence $a_i a_{i+1} \cdots a_j$ is said to be a *substring* of the string $\mathbf{x} = a_1 a_2 \cdots a_n$.

One of the first questions to ask is whether there are repetitions in a given string; i.e., a substring consisting of a block of letters immediately followed at least once by the same block of letters in the same order. If a string \mathbf{x} contains a substring of the form $\mathbf{y}\mathbf{y}$ then we say \mathbf{x} contains the square $\mathbf{y}\mathbf{y}$. A string without any substrings which are squares is said to be square-free. The string 0102010 is square-free and, moreover, cannot be extended by concatenation over the alphabet $\{0, 1, 2\}$ on either the right or the left without creating a substring which is a square.

Less studied is another kind of repetition that can occur as a substring of a string, called an Abelian square. An *Abelian square* is a string followed by a permutation of itself. Every square is also an Abelian square. Over the alphabet $\{0, 1\}$, 010100 is an Abelian square which contains the squares 0101, 1010, and 00. Thus 010100 contains 4 Abelian squares because the string itself is also an Abelian square.

Erdős [5] first asked what was the minimum alphabet size over which there exist countably infinite strings without Abelian squares. This is a variant of the corresponding problem for squares which was resolved by Thue [8] in 1906. More formally,

Definition 1 An Abelian square over the alphabet A is a non-empty string of the form

$$\mathbf{y}\mathbf{y}^{\sigma} = b_1 \cdots b_k b_{\sigma(1)} \cdots b_{\sigma(k)}$$

where σ is a permutation of A. A string is said to be Abelian square-free if it contains no Abelian squares.

Note that every square is an Abelian square corresponding to the identity permutation. Clearly every Abelian square-free string is square-free.

Over the alphabet $A = \{0, 1, 2\}$, 012201 is an Abelian square while 0102010 is Abelian square-free and cannot be extended on either the right or the left by any element of the alphabet $A = \{0, 1, 2\}$ without introducing Abelian squares.

Dekking showed that Abelian cubes are avoidable over alphabets of size 3 while Abelian fourth powers are avoidable on binary alphabets [4]. Carpi showed that on an alphabet of 4 letters the number of Abelian square-free strings grows exponentially in the length of the string [1].

Definition 2 A finite string \mathbf{x} over an alphabet A is a left (right) maximal Abelian square-free string if, for every $a \in A$, $a\mathbf{x}$ ($\mathbf{x}a$) contains Abelian squares. An Abelian square-free string is maximal if it is both left maximal and right maximal.

On a four letter alphabet, the following is a maximal Abelian square-free string of length 26:

01021302012131203020312010.

Although Abelian square-free implies square-free, a maximal Abelian square-free string need not be a maximal square-free string. A simple example is the string 1020102 over $\{0, 1, 2\}$.

It is an open question as to whether every Abelian square-free string can be extended to a maximal Abelian square-free string. The string 012 is Abelian square-free but is certainly not maximal over $\{0, 1, 2\}$ since it can be embedded in the maximal Abelian square-free string 0201202.

In 1970 Pleasants [7] showed that there existed an infinite Abelian square-free string on an alphabet of 5 elements. This result was sharpened by Keränen [6] who showed the same was true for an alphabet of 4 elements with a computer-aided proof.

Searching strings for Abelian squares is discussed in [3]. It is folklore that any Abelian square-free string over $\{0, 1, 2\}$ has length at most 7. This can be established, say, by diligently constructing the tree of possible Abelian square-free strings starting with 0 and observing that starting with 1 or 2 would yield the same tree. Knowing this allows one to prove there are 117 distinct Abelian square-free finite strings over the alphabet $\{0, 1, 2\}$ [2]. Accepting the result by Keränen [6], the case of just three letters is seen to be important because it is the last case for which all Abelian square-free strings are finite. In what follows we direct attention toward finite strings and the problem of constructing maximal finite Abelian square-free strings of short length.

2 The Recursive Construction of Zimin

We introduce a notation that makes explicit both the alphabet symbols being used and the order of occurrence of the symbols in the string. Consider an alphabet $A = \{a_0, a_1, \ldots, a_m\}$.

Definition 3 Zimin words $\mathbf{z}_k = \mathbf{z}_k(a_0, a_1, \dots, a_k)$ are defined recursively for $k = 0, \dots, m$ by

$$\mathbf{z}_0(a_0) = a_0$$

$$\mathbf{z}_k(a_0, a_1, \cdots, a_k) = \mathbf{z}_{k-1} a_k \mathbf{z}_{k-1}.$$
 (1)

An easy induction proof shows that

$$\sharp \mathbf{z}_k = [2^k, 2^{k-1}, \cdots 2, 1]$$

for each k = 0, ..., m. Summing, one obtains $|\mathbf{z}_k| = 2^{k+1} - 1$ for each k.

Zimin words have many properties. Zimin words were introduced in connection with blocking sets in [9], in the sense that for a pattern \mathbf{p} containing m different letters, then \mathbf{p} is avoidable (on some finite alphabet) if \mathbf{z}_m avoids \mathbf{p} .

Most interesting from our point of view is that not only are Zimin words square-free, but in fact they are maximal Abelian square-free over the alphabet for which they are defined. This is easy to establish by induction, because in the Zimin word \mathbf{z}_k of length $2^{k+1} - 1$, the first $2^k - 1$ entries (and the last $2^k - 1$ entries) form lower order Zimin words, which are Abelian square-free by the induction hypothesis. No Abelian square can span the central entry a_k since $|\mathbf{z}|_{a_k} = 1$ and hence a_k can not appear in two successive subwords. Maximality also follows by an easy induction.

In the next section we consider a variation of Zimin's construction that produces onesided maximal words of shorter length, and build from them two-sided maximal words.

3 The One-Sided Construction

We give a variation of the Zimin construction that depends recursively on previous Zimin words as well as previous values of the construction.

Definition 4 Left Zimin words $\mathbf{l}_k = \mathbf{l}_k(a_0, a_1, \dots, a_k)$ are defined recursively for $k = 0, \dots, m$ by

$$\mathbf{l}_{0}(a_{0}) = a_{0}$$

$$\mathbf{l}_{k}(a_{0}, a_{1}, \cdots, a_{k}) = \mathbf{l}_{k-1}a_{k} \begin{cases} \mathbf{z}_{\lfloor \frac{k-1}{2} \rfloor}(a_{0}, a_{2}, \dots, a_{k-1}) & \text{if } k \text{ is odd} \\ \mathbf{z}_{\lfloor \frac{k-1}{2} \rfloor}(a_{1}, a_{3}, \dots, a_{k-1}) & \text{if } k \text{ is even.} \end{cases}$$
(2)

Right Zimin words can be defined similarly. In fact the construction is symmetric: right Zimin words are the reversals of left Zimin words.

The first few left Zimin words on the alphabet $A = \{0, 1, 2, 3, 4, 5, 6\}$ are

$$\begin{array}{rclrcl} \mathbf{l}_0(0) &=& 0\\ \mathbf{l}_1(0,1) &=& 010\\ \mathbf{l}_2(0,1,2) &=& 01021\\ \mathbf{l}_3(0,1,2,3) &=& 010213020\\ \mathbf{l}_4(0,1,2,3,4) &=& 0102130204131\\ \mathbf{l}_5(0,1,2,3,4,5) &=& 010213020413150204020\\ \mathbf{l}_6(0,1,2,3,4,5,6) &=& 01021302041315020402061315131. \end{array}$$

We note the frequency vectors of the one sided Zimin words in the following lemma, which is easy to establish by induction.

Lemma 1 $\sharp(\mathbf{l}_k) = [2^{\lfloor (k+1)/2 \rfloor}, 2^{\lfloor k/2 \rfloor}, \cdots, 4, 2, 2, 1].$

Observe that the frequency vectors begin with a repeated entry for k even, and a single largest entry when k is odd.

Theorem 1 The string $l_k(a_0, a_1, \dots, a_k)$ is a left maximal Abelian square-free string on the alphabet $\{a_0, a_1, a_2, \dots, a_k\}$, for each $k = 0, \dots, m$.

PROOF. First note that $\mathbf{l}_0(a_0)$ is a single letter, hence Abelian square-free. Using induction, assume \mathbf{l}_{k-1} is Abelian square-free on $\{a_0, a_1, a_2, \dots, a_{k-1}\}$ and write

$$\mathbf{l}_k(a_0, a_1, \cdots, a_k) = \mathbf{l}_{k-1} a_k \mathbf{z}',$$

where \mathbf{z}' is the appropriate lower-order Zimin word as defined in (2). Now $\mathbf{l}_1, \mathbf{l}_2, \ldots, \mathbf{l}_{k-1}$ are Abelian square-free by induction, and \mathbf{z}' is Abelian square-free since it is a Zimin word. No Abelian square substring can contain a_k because a_k occurs only once in the string, hence in at most one factor of a possible Abelian square.

To show that each \mathbf{l}_k is left maximal on the alphabet $\{a_0, a_2, \dots, a_k\}$, we must check for $i = 0, 1, \dots, k$ that each string

$$a_i \mathbf{l}_k(a_0, a_1, \cdots, a_k)$$

contains an Abelian square. Since $\mathbf{l}_k(a_0, a_1, \dots, a_k) = \mathbf{l}_{k-1}a_k\mathbf{z}'$, by induction there is an Abelian square in $a_i\mathbf{l}_{k-1}$ for $1 \le i \le k-1$, hence in \mathbf{l}_k .

To see that $a_k \mathbf{l}_k$ must contain an Abelian square first suppose k is odd.

$$a_{k}\mathbf{l}_{k} = a_{k}\mathbf{l}_{k-1}a_{k}\mathbf{z}_{\lfloor\frac{k-1}{2}\rfloor}(a_{0}, a_{2}, \dots, a_{k-1})$$

= $a_{k}\mathbf{l}_{k-2}a_{k-1}\mathbf{z}_{\lfloor\frac{k-2}{2}\rfloor}(a_{1}, a_{3}, \dots, a_{k-2})a_{k}\mathbf{z}_{\lfloor\frac{k-1}{2}\rfloor}(a_{0}, a_{2}, \dots, a_{k-1})$

For convenience set

$$\mathbf{z}_1 = \mathbf{z}_{\lfloor \frac{k-2}{2} \rfloor}(a_1, a_3, \dots, a_{k-2})$$
$$\mathbf{z}_2 = \mathbf{z}_{\lfloor \frac{k-1}{2} \rfloor}(a_0, a_2, \dots, a_{k-1})$$

then let

$$\mathbf{u} = a_k \mathbf{I}_{k-2} a_{k-1}$$

and

$$\mathbf{v} = \mathbf{z}_1 a_k \mathbf{z}_2.$$

We establish that \mathbf{uv} is an Abelian square by computing frequency vectors to find that $\sharp \mathbf{u} = \sharp \mathbf{v}$. Since the frequency vector for each \mathbf{l}_k depends on the parity of k, we do the computation in both cases.

For k odd, we have for **u**

$\sharp(a_k)$	=	[0,	0,	0,	0,	···,	0,	0,	0,	0,	1]
$\sharp(\mathbf{l}_{k-2})$	=	[$2^{\frac{k-1}{2}}$,	$2^{\frac{k-3}{2}}$,	$2^{\frac{k-3}{2}}$,	$2^{\frac{k-5}{2}}$,	···,	2,	2,	1,	0,	0]
$\sharp(a_{k-1})$	=	[0,	0,	0,	0,	$\cdots,$	0,	0,	0,	1,	0]
and for \mathbf{v} ,												
$\sharp(\mathbf{z}_1)$	=	[0,	$2^{\frac{k-3}{2}}$,	0,	$2^{\frac{k-5}{2}}$,	···,	2,	0,	1,	0,	0]
$\sharp(a_k)$	=	[0,	0,	0,	0,	···,	0,	0,	0,	0,	1]
$\sharp(\mathbf{z}_2)$	=	[$2^{\frac{k-1}{2}},$	0,	$2^{\frac{k-3}{2}},$	0,	····,	0,	2,	0,	1,	0]

For k even, we have

$$a_{k}\mathbf{l}_{k} = a_{k}\mathbf{l}_{k-1}a_{k}\mathbf{z}_{\lfloor\frac{k-1}{2}\rfloor}(a_{1}, a_{3}, \dots, a_{k-1})$$

= $a_{k}\mathbf{l}_{k-2}a_{k-1}\mathbf{z}_{\lfloor\frac{k-1}{2}\rfloor}(a_{0}, a_{2}, \dots, a_{k-2})a_{k}\mathbf{z}_{\lfloor\frac{k-1}{2}\rfloor}(a_{1}, a_{3}, \dots, a_{k-1}).$

In this case we set

$$\mathbf{z}_1 = \mathbf{z}_{\lfloor \frac{k-1}{2} \rfloor}(a_0, a_2, \dots, a_{k-2})$$
$$\mathbf{z}_2 = \mathbf{z}_{\lfloor \frac{k-1}{2} \rfloor}(a_1, a_3, \dots, a_{k-1})$$

and the decomposition into \mathbf{u} and \mathbf{v} are as before.

 $\begin{aligned} & \sharp(a_k) &= \begin{bmatrix} 0, 0, 0, 0, 0, \cdots, 0, 0, 0, 0, 1 \end{bmatrix} \\ & \sharp(\mathbf{l}_{k-2}) &= \begin{bmatrix} 2^{\frac{k-2}{2}}, 2^{\frac{k-2}{2}}, 2^{\frac{k-4}{2}}, 2^{\frac{k-4}{2}}, 2^{\frac{k-4}{2}}, \cdots, 2, 2, 1, 0, 0 \end{bmatrix} \\ & \sharp(a_{k-1}) &= \begin{bmatrix} 0, 0, 0, 0, 0, \cdots, 0, 0, 0, 1, 0 \end{bmatrix} \\ & \text{and for } \mathbf{v}, \\ & \sharp(\mathbf{z}_1) &= \begin{bmatrix} 2^{\frac{k-2}{2}}, 0, 2^{\frac{k-4}{2}}, 0, \cdots, 2, 0, 1, 0, 0 \end{bmatrix} \\ & \sharp(a_k) &= \begin{bmatrix} 0, 0, 0, 0, 0, \cdots, 0, 0, 0, 0, 1 \end{bmatrix} \\ & \sharp(\mathbf{z}_2) &= \begin{bmatrix} 0, 2^{\frac{k-2}{2}}, 0, 2^{\frac{k-4}{2}}, \cdots, 0, 2, 0, 1, 0 \end{bmatrix} \end{aligned}$

In both cases, **uv** is an Abelian square.

We obtain maximal words from the one-sided maximal words of the construction as follows.

Theorem 2 For m > 0 the string

$$\mathbf{l}'_{m} = \mathbf{l}'_{m}(a_{0}, a_{1}, \cdots, a_{m}) = \mathbf{l}_{m-1}(a_{0}, a_{1}, \cdots, a_{m-1})a_{m}(\mathbf{l}_{m-1}(a_{0}, a_{1}, \cdots, a_{m-1}))^{r}$$

is maximal Abelian square-free over the alphabet $\{a_0, a_1, \dots, a_m\}$, where \mathbf{x}^r denotes the reversal of a string \mathbf{x} .

PROOF Since $\mathbf{l}_{m-1}(a_0, a_1, \dots, a_{m-1})$ is left maximal over $\{a_0, a_1, \dots, a_{m-1}\}$ by Theorem 1, none of the symbols $\{a_0, a_1, \dots, a_{m-1}\}$ can be prepended to $\mathbf{l}_{m-1}(a_0, a_1, \dots, a_{m-1})$ or appended to $(\mathbf{l}_{m-1}(a_0, a_1, \dots, a_{m-1}))^r$ without creating an Abelian square. Note that a_m can not be prepended because

$$\sharp(a_m \mathbf{l}_{m-1}(a_0, a_1, \cdots, a_{m-1})) = \sharp(a_m (\mathbf{l}_{m-1}(a_0, a_1, \cdots, a_{m-1}))^r).$$

There can be no Abelian square in either $\mathbf{l}_{m-1}(a_0, a_1, \dots, a_{m-1})$ or in $(\mathbf{l}_{m-1}(a_0, a_1, \dots, a_{m-1}))^r$, and no Abelian square can include the single occurrence of a_m in $\mathbf{l}_{m-1}a_m\mathbf{l}_{m-1}^r$.

We obtain (m + 1)! other maximal Abelian squarefree strings by permuting the underlying alphabet of m + 1 letters.

4 Calculations and Asymptotics

Reference [2] provides a complete catalog of Abelian square-free words on an alphabet of size 3. From this we can isolate the maximal Abelian square-free words. If we assume that the alphabet symbols $\{0, 1, 2\}$ have their first occurrences in a word in that order,

the possibilities can be summarized in a tree diagram in which one edge is included for each possible extension of a word by a letter.

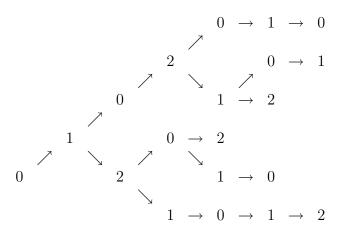


Figure 1: Right maximal Abelian square-free words

We observe $\mathbf{z}_2(0, 1, 2) = \mathbf{l}'_2(0, 1, 2)$ occurring in the topmost path in the tree.

With permutations of the alphabet included there are 6×1 right maximal words of length 5, 6×2 right maximal words of length 6, and 6×3 right maximal words of length 7 in the catalog, corresponding to the 6 leaves of the tree in Figure 1.

The work of Keränen [6] suggests that there exist infinitely many maximal Abelian square-free words over an alphabet with 4 letters. A search reveals that the shortest maximal word length is 11. The maximal words of length 11 are determined by the three classes represented by

01021032030 01021302101 01021312010.

Each class contains words for the 4! permutations of the alphabet, for a total of 72 words.

The last of these words is $\mathbf{l}'_3(0, 1, 2, 3)$.

All 312 of the words of length 12 come from the 13 words

All 792 of the words of length 13 come from the 33 words

0102010302030	0102010321020	0102032012030	0102032102030
0102101302101	0102101312010	0102101320121	0102103012010
0102103021020	0102103102101	0102103120121	0102103201023
0102103201202	0102103201203	0102123010212	0102123020121
0102123101202	0102123121020	0102123202101	0102123212010
0102130102101	0102130121012	0102130201020	0102130201202
0102131012010	0102302012030	0102302102030	0120131021013
0121012310212	0121321012312	0123021023012	0123120132131
0123210231232.			

Continuing the search, we find 37 classes of words of length 14, 47 classes of length 15 (one of which is $\mathbf{z}_3(0, 1, 2, 3)$), 49 classes of length 16, 81 of length 17, and 203 of length 18.

In the calculations above, $l'_3(0, 1, 2, 3)$ occurs in a maximal word of minimal length, whereas $\mathbf{z}_3(0, 1, 2, 3)$ is a bit longer.

The difference in length is exaggerated as the alphabet size grows, since

$$\begin{aligned} |\mathbf{z}_m(a_0, a_1, \cdots, a_m)| &= 2^{m+1} - 1, \text{ and} \\ |\mathbf{l}'_m(a_0, a_1, \cdots, a_m)| &= \begin{cases} 4 \cdot 2^{m+1/2} - 5 & , m \text{ odd} \\ 6 \cdot 2^{m/2} - 5 & , m \text{ even}. \end{cases} \end{aligned}$$

Both lengths grow geometrically with m, but the modified words are better for a given m since

$$\lim_{m \to \infty} \frac{|\mathbf{l}'(a_1, a_2, \cdots, a_m)|}{|\mathbf{l}'(a_1, a_2, \cdots, a_{m-1})|} = \sqrt{2}$$

rather than 2, which is the limiting ratio for $\mathbf{z}_m/\mathbf{z}_{m-1}$.

The string constructed by the one-sided Zimin technique generates a two-sided string of the minimum possible length 11 for alphabet size 4. For alphabet size 3, the one-sided Zimin technique produces 0102010 of length 7, but a shorter maximal word 010212 exists. It would be interesting to know if, as the alphabet size grows, the length of the maximal word produced by the one-sided Zimin technique remains close to minimal.

References

- A. Carpi, On the number of Abelian square-free words on four letters Discrete Applied Mathematics, 81(1998) pp. 155-167.
- [2] L. J. Cummings, Strongly Square-Free Strings on Three Letters The Australasian Journal of Combinatorics, 14(1996), 259–266.
- [3] L. J. Cummings and W. F. Smyth, Weak repetitions in strings, J. Combinatorial Mathematics and Combinatorial Computing, 24(1997), 33-48.
- [4] F. M. Dekking, Strongly nonrepetitive sequences and progression-free sets, J. Combinatorial Theory, A 27(1979), 181-185.

- [5] P. Erdős, Some unsolved problems, Hungarian Academy of Sciences Mat. Kutató Intézet Közl, 6(1961) 221–254.
- [6] V. Keränen, Abelian squares are avoidable on 4 letters, Lecture Notes in Computer Science, No.623, 1992, 41–52.
- [7] P. A. B. Pleasants, Non-repetitive sequences, Proc. Cambridge Phil. Soc. 68(1970), 267–274.
- [8] A. Thue, Über unendliche Zeichenreihen, Norske Vid. Selsk. Skr. I, Mat. Nat. Kl., Christiana, 7(1906), 1–22
- [9] A. I. Zimin, Blocking sets of terms, Math. USSR Sbornik, 47(1984), No. 2 353– 364.