

T-shift synchronization codes

R. Ahlswede, B. Balkenhol, C. Deppe, H. Mashurian
and T. Partner

*Fakultät für Mathematik, Universität Bielefeld, Postfach 100131, 33501 Bielefeld,
Germany*

Abstract

In this paper we give a construction of T -shift synchronization codes, i.e. block codes capable of correcting synchronization shifts of length at most T in either direction (left or right). We prove lower and upper bounds on the maximal cardinality of such codes. An infinite number of the constructed codes turn out to be asymptotically optimal.

1 Introduction

Problems of synchronization are of basic importance in coding theory. Such problems arise naturally in situations when, due to some shift in the transmission (caused for example by insertion or deletion of some amount of information), the receiver gets out of synchronism meaning that he does not anymore know the starting points of the codewords. Even in the case of noiseless transmission, asynchronism can cause false decoding, i.e. incorrect word separation. Thus, the problem is to find encoding schemes which enable the two parties to regain synchronism in their communication.

There are two conflicting goals when designing appropriate codes for this purpose. First, one tries to minimize the synchronization delay, i.e. the number s of consecutive symbols which must be read by the receiver for correct decoding of the messages. Second, the code designer tries to maximize the code size. Among different classes of synchronization codes the statistically synchronizable codes are defined by the most relaxed condition on the delay s requiring that $\lim_{S \rightarrow \infty} \Pr(s \leq S) = 1$ holds. As a result, these codes achieve minimal redundancy like the Huffman codes [2].

Synchronization codes with finite delay [3] assume a fixed upper bound S on the random variable s . Already this condition essentially reduces the code size.

It implies the asymptotical upper bound q^n/n for q -ary codes of block length n . Obviously the same bound holds also for more restrictive classes of comma-free codes [9] and overlap-free codes [5]. The latter ones, however, being more redundant than general synchronization codes with finite delay, show better synchronization capability: for the maximal delay we have $S = 2n - 1$ where n is the block length. For the purpose of simple encoding and decoding one also considers synchronizable codes in the family of cosets of cyclic codes [7], [8]. The so-called prefix-synchronized codes [6] also admit easy implementation due to simplicity of the encoding and decoding procedures. However, all these codes have high redundancy.

The model which we discuss in this paper is motivated by the following situation of data storage. Suppose that the data is encoded by means of a code of block length n and the encoded information is stored (for instance, written on a disc) without any disturbance. However, some random words of random length (breaks) may be added between the codewords causing loss of synchronism at the receiver end. We assume that there is an upper threshold T for the length of each inserted word. For the decoder this means that the first n symbols of the received message do not necessarily constitute a codeword but are rather a left-shift of a codeword whereby the shift length t is upper-bounded by T . Shifts in both directions arise when the receiver decodes the message starting at some position in the received sequence (not necessarily the beginning). Thereby he knows with some precision the locations of the codewords separation points. In other words, he knows that the actual starting point of a codeword which is next to the position where the decoding process begins is at most T symbols apart. The aim now is to design maximal (in terms of the code size) codes capable of correcting shifts of the above-mentioned type, T -shift synchronization codes.

This model was first considered in [4], where 1-shift synchronization codes were constructed. Further progress was made in [1] providing a construction of 2-right shift and asymptotically optimal 1-right shift synchronization codes. In this paper we improve the lower bound given in [4] and generalize the results of [1] and [4] to the case of T -shifts. The method used here is a refinement of the method proposed in [1].

The paper is organized as follows. In Section 2 we give some basic notation and necessary definitions. In Section 3 the code construction is presented. Hereby (until Section 5) we restrict ourselves to the case of unidirectional shifts. Theorem 1 proves that the codes thus constructed are indeed right-shift and left-shift synchronizing. In Section 4 lower and upper bounds are given for asymptotical behavior of an optimal unidirectional T -shift synchronization code (Theorem 3). In Section 5 we show that the problem with shifts in both directions reduces to the case of unidirectional shifts. The final Remark 3 indicates some features of the codes constructed. For simplicity of presentation

we discuss only the case of a fixed threshold T and the binary alphabet. The results hold also for arbitrary finite alphabets and any threshold function T depending on the block length n with an order of growth $T(n) = o(n \cdot \log^{-2} n)$ (see Remark 3).

2 Notation and definitions

For a finite set $\mathcal{X} = \{0, \dots, q-1\}$ called q -ary alphabet, we form $\mathcal{X}^n = \{0, \dots, q-1\}^n$, the words of length n , with letters from \mathcal{X} .

Definition 1 *The word $b^n = (b_1, \dots, b_n)$ is said to be a t -right shift of the word $a^n = (a_1, \dots, a_n)$ for a non-negative integer $t < n$ iff the first $n-t$ symbols of b^n coincide with the last $n-t$ symbols of a^n , that is,*

$$(b_1, \dots, b_{n-t}) = (a_{t+1}, \dots, a_n).$$

In this case, a^n is called a t -left shift of the word b^n .

Note that the 0-shifting (the case $t = 0$) leaves words unchanged. The following notions are central in the paper.

Definition 2 *Let $T < n$ be a positive integer. The numbers t_1 and t_2 below are assumed to be non-negative integers satisfying $t_1, t_2 \leq T$. Consider a block code $C \subseteq \{0, 1\}^n$.*

(a) C is called T -right shift synchronizing iff for all distinct t_1 and t_2 , no t_1 -right shift of any codeword is a t_2 -right shift of any codeword,

(b) Symmetrically, C is called T -left shift synchronizing iff for all distinct t_1 and t_2 , no t_1 -left shift of any codeword is a t_2 -left shift of any codeword,

(c) Finally, C is called T -shift synchronizing iff it is T -right shift synchronizing, T -left shift synchronizing and for all positive t_1 and t_2 , no t_1 -right shift of any codeword is a t_2 -left shift of any codeword.

Remark 1. We infer from Definition 2 that in any code of one of the mentioned types no positive shift of a codeword can be a codeword. This is easily seen by taking $t_1 = 0$.

3 The code construction

Our aim is to construct T -shift synchronization codes of maximal cardinality. It turns out that the problem reduces to the case of shifts in one direction (see Section 5).

We need the following relations between two real numbers a and b . For the binary alphabet we write:

$$a >_T b \quad \text{iff} \quad a - b > T$$

$$a <_{-T} b \quad \text{iff} \quad a - b < -T$$

$$a \sim_T b \quad \text{iff} \quad -T \leq a - b \leq T$$

$$a =_k b \quad \text{iff} \quad a - b = k \text{ with } -T \leq k \leq T$$

We note that in the general case of a q -ary alphabet we have to replace here all appearances of T by $T(q-1)$. Now, choose a natural number m in such a way that $2^{m-1} \leq T < 2^m$ holds. Let n be the block length of the code to be constructed. Since we later let n grow, we can assume that $n \geq 2^m$. So represent n in the form $n = 2^m \cdot r + n_1$ with $0 \leq n_1 < 2^m$. Now, let the last n_1 positions in a codeword take all possible 2^{n_1} values. The construction below shows that the first $2^m \cdot r$ positions already enforce the code properties sought. Therefore, without loss of generality, we can assume $n_1 = 0$, which means that 2^m divides n .

Let (μ_m, \dots, μ_1) be any m -tuple of relations with $\mu_i \in \{>_T, <_{-T}, \sim_T, =_k\}$. The following m (in)equalities for a word (x_1, \dots, x_n) are basic in our code construction:

$$\sum_{i=1}^{n/2^j} x_{2^j \cdot i} \quad \mu_j \quad \sum_{i=1}^{n/2^j} x_{2^j \cdot i - 2^{j-1}} \quad (j^*)$$

for $j = 1, \dots, m$. Expanded they look as follows.

$$(x_2 + x_4 + x_6 + x_8 + \dots) \quad \mu_1 \quad (x_1 + x_3 + x_5 + x_7 + \dots) \quad (1^*)$$

$$(x_4 + x_8 + \dots) \quad \mu_2 \quad (x_2 + x_6 + \dots) \quad (2^*)$$

$$(x_8 + \dots) \quad \mu_3 \quad (x_4 + \dots) \quad (3^*)$$

..... ..

$$(x_{2^m} + x_{2 \cdot 2^m} + \dots) \quad \mu_m \quad (x_{2^{m-1}} + x_{3 \cdot 2^{m-1}} + \dots) \quad (m^*)$$

Now, for each $1 \leq s \leq m$, define

$$C(\mu_s, \dots, \mu_1) := \{(x_1, \dots, x_n) \text{ satisfying } (1^*), (2^*), \dots, (s^*)\}$$

Remark 2. (a) To be more precise, we should have included the word length n in the notation $C(\mu_s, \dots, \mu_1)$, but we leave it out for the sake of simplicity. By convention, we assume that all words in $C(\mu_s, \dots, \mu_1)$ always have the same length n .

(b) We use the shorthand $C^s(>_T)$ for the set $C(\mu_s, \dots, \mu_1)$ when all the relations μ_1, \dots, μ_s are equal to the same relation $>_T$.

Now we complete our code construction by taking $s = m$ and $\mu_m = \dots = \mu_1 = >_T$. In other words, $C^m(>_T)$ is the constructed code. In Theorem 1 below we will show that $C^m(>_T)$ is unidirectionally T shift synchronizing. It will provide the lower bound in Theorem 3. We need the following key result.

Lemma 1 *Suppose that a t_1 -right shift of the word $a^n = (a_1, \dots, a_n)$ coincides with a t_2 -right shift of the word $b^n = (b_1, \dots, b_n)$ where $t_1, t_2 < n$. If $t_2 - t_1$ is an odd natural number, then*

$$\begin{aligned} -t_2 &\leq \sum_{i=1}^{n/2} b_{2i} - \sum_{i=1}^{n/2} a_{2i-1} \leq t_2 \\ -t_2 &\leq \sum_{i=1}^{n/2} b_{2i-1} - \sum_{i=1}^{n/2} a_{2i} \leq t_2 \end{aligned}$$

Proof. This is immediate, since if the conditions of the lemma are fulfilled, then the entries of the even positions in b^n coincide with the entries of the odd positions in a^n except of at most t_2 positions. The same is true for the odd positions of b^n and even positions of a^n . \square

By symmetry the result of Lemma 1 also holds for left shifts.

Theorem 1 *The code $C^m(>_T)$ is a T -right shift and T -left shift synchronization code.*

Proof. We only prove the result for right shifts. Left shifts are settled completely symmetrically. Suppose that the assertion of the theorem is not true. This means that there are two codewords $a^n, b^n \in C^m(>_T)$ and two distinct non-negative integers $t_1, t_2 \leq T$ such that a t_1 -right shift of a^n coincides with a t_2 -right shift of b^n . Without loss of generality we can assume that $t_1 < t_2$. We come to contradiction by showing that b^n does not satisfy one of the inequalities (1*), (2*), \dots , (m*).

First case: $t_2 - t_1$ is odd. We have that a^n as a codeword satisfies (1*), (2*), \dots , (m*). Then inequality (1*) implies

$$\sum_i a_{2i} - \sum_i a_{2i-1} > T \tag{1}$$

For b^n we obtain

$$\begin{aligned}
\sum_i b_{2i} - \sum_i b_{2i-1} &\leq \sum_i a_{2i-1} + t_2 - \sum_i a_{2i} + t_2 && \text{(Lemma 1)} \\
&< -T + 2t_2 && \text{(Inequality (1))} \\
&\leq T && (t_2 \leq T)
\end{aligned}$$

Thus, inequality (1*) does not hold for b^n and therefore b^n cannot be a code-word.

Second case: $t_2 - t_1$ is even. Separating the maximal power of 2 we represent this difference in the form $t_2 - t_1 = 2^k \cdot t'$ where t' is an odd natural. Since $t_1, t_2 \leq T < 2^m$, we have $1 \leq k \leq m - 1$. Recall that without loss of generality we have assumed the divisibility of n by 2^m . Therefore also 2^k divides n .

The main idea of the proof is already contained in the previous case. We reduce the present case to that one by concentrating on the positions of a^n and b^n which are multiples of 2^k . So consider the following two words:

$$\begin{aligned}
a' &= (a_{2^k} a_{2 \cdot 2^k} a_{3 \cdot 2^k} \cdots a_n) \\
b' &= (b_{2^k} b_{2 \cdot 2^k} b_{3 \cdot 2^k} \cdots b_n)
\end{aligned}$$

We leave it to the reader to easily verify that a t'_1 -right shift of a' coincides with a t'_2 -right shift of b' if we take $t'_1 = \lfloor t_1 / (2^k) \rfloor$ and $t'_2 = t'_1 + t'$. But we know that a' satisfies the inequality $(k+1)^*$, and now $t'_2 - t'_1 = t'$ is odd. As in the previous case we infer that inequality $(k+1)^*$ cannot hold for b' . This means that b^n is not a codeword. This contradiction proves the theorem. \square

4 Lower and upper bounds

To evaluate the size of $C^m(>_T)$ we need some auxiliary results.

Lemma 2 *The sets $C(>_T, \mu_s, \dots, \mu_1)$, $C(<_{-T}, \mu_s, \dots, \mu_1)$, $C(\sim_T, \mu_s, \dots, \mu_1)$ are pairwise disjoint and their union is equal to $C(\mu_s, \dots, \mu_1)$.*

Proof. Obvious, since $>_T$, $<_{-T}$, and \sim_T are mutually excluding and complement each other. \square

Lemma 3

$$|C(>_T, \mu_s, \dots, \mu_1)| = |C(<_{-T}, \mu_s, \dots, \mu_1)|$$

Proof. By symmetry. For each sequence $(x_1, \dots, x_n) \in C(>_T, \mu_s, \dots, \mu_1)$, we exchange the positions $x_{2^{s+1}, i}$ and $x_{2^{s+1}, i-2^s}$ for all $1 \leq i \leq n/2^{s+1}$ and keep all other positions unchanged. The obtained sequence is in $C(<_{-T}, \mu_s, \dots, \mu_1)$ and this correspondence is a bijection. \square

Lemma 4 *Let s be a non-negative integer such that 2^{s+1} divides the block length n , so $n = 2^{s+1} \cdot r$. For any integer k satisfying $-T \leq k \leq T$ it holds*

$$|C(=_{k, \mu_s, \dots, \mu_1})| \leq \binom{2r}{r+k} 2^{n-2r}. \quad (2)$$

Proof. Each word from $C(=_{k, \mu_s, \dots, \mu_1})$ satisfies inequalities $(1^*), \dots, (s^*), (s+1)^*$ where relation $=_k$ is substituted for μ_{s+1} . Ignoring the first s of these restrictions would increase the size of $C(=_{k, \mu_s, \dots, \mu_1})$. So now we upper-bound the number of words of length n satisfying inequality $(s+1)^*$. This inequality refers to the $2r$ positions $2^s, 2 \cdot 2^s, 3 \cdot 2^s, \dots, n$. Since the remaining $n - 2r$ positions have no any restriction, they contribute with the factor 2^{n-2r} in (2). Therefore the problem reduces to the following one.

Denote

$$A_r = \{a^{2r} = (a_1, a_2, \dots, a_{2r}) : \sum_{i=1}^r a_{2i} - \sum_{i=1}^r a_{2i-1} = k\}.$$

Show that

$$|A_r| \leq \binom{2r}{r+k}. \quad (3)$$

We proceed as follows. Each word a^{2r} from A_r we transform into a new word b^{2r} by negating the letters in odd positions of a^{2r} . It is easily seen that b^{2r} has $r+k$ ones. Since every a^{2r} produces a different b^{2r} , inequality (3) follows. The lemma is proved. \square

Lemma 5

$$\frac{|C(\sim_T, \mu_s, \dots, \mu_1)|}{2^n} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Proof. We evaluate the size of $C(\sim_T, \mu_s, \dots, \mu_1)$ from above.

$$\begin{aligned}
|C(\sim_T, \mu_s, \dots, \mu_1)| &= \left| \bigcup_{k=-T}^T C(=k, \mu_s, \dots, \mu_1) \right| \\
&\leq \sum_{k=-T}^T |C(=k, \mu_s, \dots, \mu_1)| \\
&\leq \sum_{k=-T}^T \binom{2r}{r+k} 2^{n-2r} \quad (\text{Lemma 4}) \\
&\leq (2T+1) \binom{2r}{r} 2^{n-2r}
\end{aligned}$$

The well-known formula

$$\frac{\binom{2r}{r}}{2^{2r}} = \frac{1}{\sqrt{\pi r}} (1 + \alpha(r))$$

with vanishing $\alpha(r)$ as $r \rightarrow \infty$ completes the proof. \square

As a consequence we obtain the following result about the asymptotical behavior of the code $|C^m(>_T)|$.

Theorem 2 *For $2^{m-1} \leq T < 2^m$ we have*

$$\lim_{n \rightarrow \infty} \frac{|C^m(>_T)|}{2^n} = \frac{1}{2^m}$$

Proof. Lemmas 2, 3 and 5 imply that the size of the code $C(>_T, \mu_s, \dots, \mu_1)$ has the same asymptotical behavior as the half of $|C(\mu_s, \dots, \mu_1)|$. By iterating this m times for $C^m(>_T)$, we get the result. \square

Theorem 3 *For any optimal T -right shift synchronization code of block length n we have*

$$\frac{1}{2^m} \leq \lim_{n \rightarrow \infty} \frac{|C_{opt}|}{2^n} \leq \frac{1}{T+1} \quad (4)$$

where $2^{m-1} \leq T < 2^m$. The same is true for optimal T -left shift synchronization codes.

Proof. Again, by symmetry, we only need to settle the case of right shifts. The lower bound follows directly from Theorem 2. To show the upper bound, let C be any T -right shift synchronization code. For each $1 \leq t \leq T$ and $x^t \in \{0, 1\}^t$ we consider the function f_{x^t} acting on the set C as follows. The f_{x^t} -image of a codeword $c^n \in C$ is obtained by removing the first t symbols of c^n and appending x^t on the right. By C_{x^t} we denote the f_{x^t} -image of C . Obviously, for distinct words x^t the corresponding sets C_{x^t} are disjoint. Now

put

$$C_t := \bigcup_{x^t \in \{0,1\}^t} C_{x^t}.$$

In other words, C_t is the set of all t -right shifts of codewords from C . To evaluate the size of C_t we notice the following. If two codewords from C (they are words of length n) have the same f_{x^t} -image, then their last $n - t$ positions must coincide. But there are at most 2^t possibilities for the first t positions. This implies that no more than 2^t words from C can have the same image under the mapping f_{x^t} . Consequently,

$$|C_{x^t}| \geq 2^{-t} \cdot |C|,$$

and therefore

$$|C_t| = \sum_{x^t \in \{0,1\}^t} |C_{x^t}| \geq |C|.$$

Now, according to the definition of a T -right shift synchronization code, the code C and the sets C_t for $1 \leq t \leq T$ must be pairwise disjoint. We obtain

$$\left| \bigcup_{t=1}^T C_t \cup C \right| \leq 2^n \quad (\text{the whole space})$$

This implies

$$(T + 1) \cdot |C| \leq 2^n$$

and hence the assertion. \square

When $T = 1$ and $T = 2$, we obtain the results of [1] as a consequence of the lower bound in (4).

5 Shifts in both directions

It turns out that the synchronization problem for shifts in both directions can be easily reduced to the case of unidirectional shifts. The next lemma shows the connection.

Lemma 6 *If a t_1 -right shift of a word a^n is equal to a t_2 -left shift of a word b^n , where $t_1 + t_2 < n$, then b^n is a $(t_1 + t_2)$ -right shift of a^n .*

Proof. Let z^n be a t_1 -right shift of a^n and at the same time a t_2 -left shift of b^n . Then b^n is a t_2 -right shift of z^n . Therefore we can obtain b^n from a^n by moving to the right: first by t_1 steps (obtaining z^n) and then by t_2 steps. \square

This implies

Lemma 7 *Each code which is simultaneously $2T$ -right shift and $2T$ -left shift synchronizing, is a T -shift synchronization code.*

Now, as a consequence of this lemma and Theorem 3, we obtain corresponding lower and upper bounds for T -shift synchronization codes.

Theorem 4 *For any optimal T -shift synchronization code we have*

$$\frac{1}{2^m} \leq \lim_{n \rightarrow \infty} \frac{|C_{opt}|}{2^n} \leq \frac{1}{2T+1}$$

where $2^{m-1} \leq 2T < 2^m$.

Proof. According to Theorem 1, the code $C^m(>_{2T})$ obtained by our construction is $2T$ -right shift and $2T$ -left shift synchronizing. Therefore it is a T -shift synchronization code. Due to Theorem 2 this gives the lower bound. To show the upper bound, let C be an arbitrary T -shift synchronization code. Consider for all values $t = 1, 2, \dots, T$ the sets

$$C_t^{right} := \{\text{All } t\text{-right shifts of codewords from } C\}$$

and

$$C_t^{left} := \{\text{All } t\text{-left shifts of codewords from } C\}$$

Like in the proof of Theorem 3, these sets and C are pairwise disjoint and each of them has size greater or equal to $|C|$. Hence the upper bound. \square

Theorem 4 improves the lower bound on the size of an optimal 1-shift synchronization code (the case $T = 1$) stated in [4].

The next remark summarizes some features of the codes constructed and the results obtained in the paper.

Remark 3.

(a) Although formulated for the binary case, all theorems hold also for arbitrary q -ary alphabets. Only the denominator 2^n has to be replaced then by q^n , the size of the whole space. Everything else remains unchanged.

(b) By a slight modification of the code construction, the same arguments show that the results are true also for a more general situation, namely when the maximal shift length $T(n)$ depending on the block length n has order of growth

$$T(n) = o(n \cdot \log^{-2} n). \quad (5)$$

For this, we just have to consider the new relations $>_{T \cdot 2^{-i}}$, $<_{-T \cdot 2^{-i}}$, $\sim_{T \cdot 2^{-i}}$ for $i = 0, 1, \dots, m-1$. The relations considered earlier correspond to the case $i = 0$. Now, we obtain the modified code by putting in $C(\mu_m, \dots, \mu_1)$ the values $\mu_{i+1} = >_{T \cdot 2^{-i}}$ for $i = 0, 1, \dots, m-1$. Note that equation (5) is fulfilled for any function $T(n) = O(n^a)$ with $a < 1$. For $T(n) = n-1$, the class

of T -right shift synchronizable codes coincides with the class of overlap-free codes.

(c) The smaller the positive difference between 2^m and T , the smaller is the gap between the lower and upper bounds. Especially, in the model of unidirectional shifts, the bounds coincide for $T = 2^m - 1$, $m = 1, 2, \dots$ (see Theorem 3), thus providing asymptotically optimal synchronization codes.

(d) Like for comma-free codes and overlap-free codes, the maximal synchronization delay is $S = 2n - 1$, which means that at most $2n - 1$ consecutive symbols have to be read by the decoder in order to regain synchronism.

(e) The constructed codes allow easy encoding and decoding, since no look-up table is needed in order to decide whether a sequence of symbols is a codeword or not. One has to verify the equations (1^*) , (2^*) , \dots , (m^*) .

References

- [1] R. Ahlswede, B. Balkenhol and T. Partner: Shift synchronization codes, *Preprint* 00 – 118, Universität Bielefeld, 2000.
- [2] R.M. Capocelli, A.A. De Santis, L. Gargano, and U. Vaccaro: On the construction of statistically synchronizable codes, *IEEE Trans. on Inf. Theory*, vol.38, no.2, 407–414, 1992.
- [3] S.W. Golomb and B. Gordon: Codes with bounded synchronization delay, *Information and Control*, 8, 355–372, 1965.
- [4] G. Khachatryan and S. Sargissian: Almost optimal single-shift self-synchronized codes, *Proceedings of the first INTAS International Seminar on Coding Theory and Combinatorics*, Tsahkadzor, Armenia, 122–124, October 1996.
- [5] V.I. Levenshtein: Maximum number of words in codes without overlaps, *Problems of Inf. Transm.*, 7(3), 215–222, 1971.
- [6] H. Morita, A.J. van Wijngaarden and A.J. Han Vinck: On the construction of maximal prefix-synchronized codes, *IEEE Trans. on Inf. Theory*, vol.42, no.6, 2158–2166, 1996.
- [7] W.W. Peterson and E.J. Weldon, Jr.: *Error-Correcting Codes*, MIT Press, Cambridge, MA and London, England, 1972.
- [8] J.J. Stiffler: *Theory of Synchronous Communications*, Prentice Hall, Englewood Cliffs, NJ, 1971.
- [9] B. Tang, S.W. Golomb and R.L. Graham: A new result on comma-free codes of even word-length, *Can. J. Math.*, vol.39, no.3, 513–526, 1987.