
Further Clarification on Mantin’s Digraph Repetition Bias in RC4

Pranab Chakraborty · Subhamoy Maitra

Abstract In this note we provide a theoretical argument towards an unsolved question related to Mantin’s Digraph Repetition Bias (2005) that is observed in the key-stream of RC4. The open question, that depends on the observation that arrival of four consecutive same bytes in RC4 key-stream is slightly negatively biased, was posed by Bricout et al [Des. Codes Cryptogr. (2018) 86:743-770] in 2016.

Keywords RC4 · Non-randomness · Sequence · Stream Cipher.

1 Introduction

RC4 is possibly the most popular stream cipher and it attracted huge attention in the domain of cryptanalysis (see for example [4–6] and the references therein). Recently there are evidences of almost practical attacks on this cipher and thus the cipher is not recommended to be deployed in new systems. However, this cipher still handles considerable traffic in different networks and thus of interest to cryptologic community. At the same time, the cipher is a very interesting combinatorial object to study. Even after serious efforts for around four decades, we are still amazed with novel results continuously coming in this domain of research.

The best long term bias observed in RC4 keystream was provided by Mantin long back [3]. It says that the probability of obtaining a substring of the form $ABTAB$ (A, B 8-bit characters and T is a short string of such characters) in RC4 stream is greater than what should be obtained in a true random situation. This bias is famously referred to as the “Digraph Repetition Bias” in RC4. A detailed study in this regard has been presented recently in [1], which is referred as the fine-grained analysis. Through this analysis it has been theoretically argued in [1,

Pranab Chakraborty
Learning and Development, Human Resources, Wipro Limited, Doddakannelli, Sarjapur Road,
Bangalore 560035, India. E-mail: kojagori@gmail.com

Subhamoy Maitra
Applied Statistics Unit, Indian Statistical Institute, 203 B T Road, Kolkata 700108, India,
E-mail: subho@isical.ac.in

Theorem 1] that the bias should be little more in the case when $A = B$. However, all the cases under this situation could not be clarified in [1] and in particular, when $A = B$ and \mathcal{T} is null, then the bias could not be observed at all through experiments. The authors of [1] thus commented,

“However, when $A = B$, we do not see the positive bias behaviour predicted by [1, Theorem 1], but instead a small, negative bias. We do not currently have an explanation for this behaviour.”

In this note, we answer this question with detailed theoretical analysis of RC4 evolution during the pseudo-random key-stream generation process.

Before proceeding further, let us first quickly describe the RC4 algorithm. In RC4, we have an $N = 256$ length array of 8-bit integers 0 to $N - 1$, that works as a permutation. There is also an l length array of bytes K (the secret key), where l may vary from 5 to 32, depending on the key length. There are also two bytes i, j , where i is the deterministic index that increases by 1 in each step and j is updated in a manner so that it behaves pseudo-randomly. The Key Scheduling Algorithm (KSA) of RC4 is as follows:

- $j = 0$; for $i = 0$ to $N - 1$: $S[i] = i$;
- for $i = 0$ to $N - 1$:
 $j = j + S[i] + K[i \bmod l]$; swap($S[i], S[j]$);

Next, the pseudo-random bytes z are generated during the Pseudo Random Generator Algorithm (PRGA) as follows:

- $i = j = 0$;
- for $i = 0$ to $N - 1$:
 $i = i + 1$; $j = j + S[i]$; swap($S[i], S[j]$); $z = S[S[i] + S[j]]$;

All the additions here are modulo N .

The work of [3] presented the first distinguisher for RC4 when any amount of initial keystream bytes are thrown away. This distinguisher is based on the digraph distribution of RC4. The term digraph means a pair of consecutive keystream words. In [3, Section 3], it has been shown that getting strings of the pattern $ABTAB$ (where A, B are bytes and \mathcal{T} is a string of bytes of small length $G \leq 16$), is more probable in RC4 keystream than in random stream. The exact theoretical result [3, Theorem 1] is as follows.

Theorem 1 *During RC4 PRGA, for small integer values of $G \geq 0$*

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = B) | (z_r = A, z_{r+1} = B)) = \frac{1}{N^2} \left(1 + \frac{e^{-\frac{8-8G}{N}}}{N}\right).$$

This result is true for most of the cases under some logical assumptions on independence. However, it should be noted that being a deterministic stream cipher on a classical paradigm, the states of RC4 actually dependent on each other, whatever less the influence may be. Thus, there are cases, where the bias is not exactly the same as in Theorem 1. In this direction detailed analysis has been presented in [1] and it has been observed that when $A = B$, the bias is more prominent.

Theorem 2 *During RC4 PRGA, for small integer values of $G \geq 0$*

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = A) | (z_r = A, z_{r+1} = A)) = \frac{1}{N^2} \left(1 + \frac{e^{-\frac{4-6G}{N}}}{N}\right).$$

Interestingly, this bias could not be observed for $G = 0$ in experiments as explained following [1, Figure 2].

In RC4 related research, the biases are generally identified in two ways.

- One can run some experiments to observe the biases and then try to prove them.
- One can theoretically inspect the algorithm to obtain the bias, prove it theoretically and then supplement it with experiments.

As we have commented earlier, the proofs are completed based on certain assumptions. Thus, in specific cases, due to incorrect assumptions, the reported biases may not exist. These are identified later through more disciplined studies. This is exactly what has been pointed out in [1] and left as an open question. In fact, in this case we actually do not concentrate on showing the bias. Rather we try to argue with detailed theoretical analysis that the bias is indeed negligible.

2 Explanation of the small negative bias in the AAAA sequence

In this section we first explain in details the arguments presented by Mantin in Lemma 2 and Theorem 1 in [3] and then describe additional refinements of the result by Bricout et. al. in [1]. While experimenting on the refinements of the proposed results, Bricout et. al. identified a deviation in the observed behavior (from the expected behavior predicted by Theorem 1 in [1]) for a specific form of digraph repetition sequence that has the form AAAA. To the best of our knowledge, this deviation remained unexplained so far. In this section, We present a theoretical explanation to this behavior.

2.1 Revisiting Mantin's result [3]

It appears that due to a possible typographical error, the statement of [3, Theorem 1] is not exactly correct, and it differs slightly from the proof. As given correctly in the proof, during RC4 PRGA, for small integer values of $G \geq 0$

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = B) | (z_r = A, z_{r+1} = B)) = \frac{1}{N^2} \left(1 + \frac{e^{-\frac{8-8G}{N}}}{N}\right).$$

However, the theorem statement says,

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = B) | (z_r = A, z_{r+1} = B)) = \frac{1}{N^2} \left(1 + \frac{e^{-\frac{4-8G}{N}}}{N}\right),$$

where -4 is misprinted instead of -8 in the exponent. This misprint is carried in [1, Result 2] too, where [3, Theorem 1] is referred.

Let us now describe the approach that Mantin had used to prove the stated result and in the process we point out additional clarifications for specific cases that demand refinements of the result. The key observation made by Mantin is the fact that if, with respect to an arbitrary round r of RC4 PRGA, $S_r[i_r + 1] = 1$ and $S_{r-1}[i_r] = x$ is any byte-value other than 1, then at the end of round $r + 1$, the permutation byte pair $(x, 1)$ would move to the location indexed by $(j_r, j_r + 1)$ and if the byte pair remains undisturbed till round $(r + G + 2)$ (where $G \geq 0$

is an integer signifying the gap between the source and destinations pairs), then under an additional condition that $j_{r+G+2} = i_r$, the key-stream byte pair (z_r, z_{r+2}) would repeat as the byte pair (z_{r+G+2}, z_{r+G+3}) .

Note that unless specifically mentioned, $S_r[y]$ is the y -th element of the S array, after the swap is done in the r -th round. To be more specific, for a given G , the conditions for the event to occur are as follows:

1. $j_r = i_{r+G+2}$,
2. $S_r[i_r + 1] = 1$,
3. $j_{r+G+2} = i_r$,
4. the locations indexed by i_r, i_{r+1}, i_{r+G+2} and i_{r+G+3} are undisturbed in the intermediate G rounds,
5. the location of the key-stream byte z_r indexed by $S_r[i_r] + S_r[i_{r+G+2}]$ remains unchanged in $G + 2$ rounds ($r + 1, \dots, r + G + 2$) and the location of the key-stream byte z_{r+1} indexed by $S_{r+1}[i_{r+1}] + S_{r+1}[i_{r+G+3}]$ remains unchanged in $G + 2$ rounds ($r + 2, \dots, r + G + 3$).

Incidentally, Mantin stated the conditions of the theorem by referring to the gap with the variable $g = j_r - i_r$ and later introduced the variable $G = (g - 2)$ and called it the real gap. If $g = 0$, $i_r = j_r$ and the permutation byte pair $(x, 1)$ stays in the same locations at the end of round $r + 1$. Similarly, $g = 1$ is forbidden in real RC4 as Finney cycles [2] can't occur. Therefore, we only consider the case $g \geq 2$ or in other words $G \geq 0$.

The probability associated with the conditions (1-3) is clearly $\frac{1}{N^3}$. Using [3, Lemma 1], we find that the probability for condition 4 is around $e^{(-4G)/N}$ for small values of G . Similarly, the probability corresponding to the condition 5 is $(1 - \frac{G+2}{N})^2 \cdot e^{\frac{-2(G+2)}{N}} \approx e^{\frac{-4(G+2)}{N}}$. Hence, the combined probability of the desired event according to conditions 1-5 is $e^{\frac{-8-8G}{N}} \cdot \frac{1}{N^3}$.

On the other hand, for the complimentary scenario with probability $(1 - \frac{1}{N^3})$, in which one or more of the conditions (1-3) do not hold, we consider the event probability as the fair chance of $\frac{1}{N^2}$. So the combined probability for the complimentary scenario is $(1 - \frac{1}{N^3}) \cdot \frac{1}{N^2}$.

By using the above probability values one can obtain the desired result,

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = B) | (z_r = A, z_{r+1} = B)) = \frac{1}{N^2} \left(1 + \frac{e^{\frac{-8-8G}{N}}}{N}\right).$$

2.2 Revisiting Bricout et. al. result [1]

Mantin (in [3]) mentioned that for the sake of simplicity he made certain heuristic assumptions. However, those were not elaborated in [3]. Bricout et. al. [1], while performing a fine grained analysis of the proof, showed that there are certain special cases in which one should not expect any digraph repetition bias. For example, Mantin's result would not be applicable for the following cases:

1. $A = 1$
2. $B = 1$
3. $A = (N - 3)$ and $G = 0$
4. $B = (N - 3)$ and $G = 0$

The reason that these cases do not demonstrate digraph repetition bias as per Mantin's result is due to the fact that in each of these cases the condition 5 as required by Mantin's [3, Theorem 1] as stated in Theorem 1 above, gets violated. In addition to the above cases, Bricout et al. [1] also showed that for a generic pattern of the form $AATAA$, there should be a stronger digraph repetition bias than the bias given by Mantin's result. We here outline the proof of Theorem 2.

The crucial observation for this result is that when $A = B$, the requirement of condition 5 as per the proof given above for [3, Theorem 1] reduces to the condition of non-disturbance of one target permutation byte position instead of two byte positions. Hence the probability increases by a multiplicative factor of $e^{\frac{2G+4}{N}}$. This brings us to the modified result:

$$\Pr((z_{r+G+2} = A) \wedge (z_{r+G+3} = A) | (z_r = A, z_{r+1} = A)) = \frac{1}{N^2} \left(1 + \frac{e^{-\frac{4-6G}{N}}}{N}\right).$$

One should note that the above result is not applicable for certain specific values of A and B [1]. All these deviations, as identified by Bricout et. al. [1], have been experimentally verified in their paper, except one specific class of patterns. For the pattern of the form $AAAA$, the experimental result showed slight negative bias instead of the strong positive bias as expected in Theorem 2. It has been mentioned in [1] that no explanation could be found out for such a deviation. We solve this issue in the next section (Theorem 3) by proving the slight negative bias. We also prove that there is a dependence of the extent of this bias on certain special values of index i_r around which this digraph repetition is observed.

2.3 Our result

We first prove two results (Lemma 1 and Lemma 2) that will be referred to prove our final result, i.e., Theorem 3.

Lemma 1 *During RC4 PRGA,*

1. $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = 0, i_r \neq 1) = \frac{2}{N^2}$ and
2. $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = 0, i_r = 1) = \frac{3}{N^2}$.

Proof Let us assume $S_r[i_r] = p$ and $S_r[j_r] = q$, where p and q are two arbitrary byte values. In case j_r coincides with i_r , it is evident that p equals q . As $S_r[i_r + 1] = 0$, p can't be 0. Clearly, $S_r[p + q] = z_r$.

We now investigate what happens in round $r+1$. Since $S_r[i_r + 1] = 0$, $j_{r+1} = j_r$. This implies $S_{r+1}[j_{r+1}] = q$ immediately before the swap step. After the swap operation, $S_{r+1}[i_{r+1}] = q$ and $S_{r+1}[j_{r+1}] = 0$. Therefore, $S_{r+1}[q + 0] = S_{r+1}[q] = z_{r+1}$. If $z_r = z_{r+1}$, then $S_r[p + q] = S_{r+1}[q]$. We now identify the situations that lead to these conditions.

As $p \neq 0$, $(p + q) \neq q$. That means $(p + q)$ and q must point to two different array byte positions of S . Thus, the only way $S_r[p + q]$ can be equal to $S_{r+1}[q]$ is when $z_r (= S_r[p + q])$ gets swapped in round $r + 1$ and moves to a new position pointed to by q . This may happen in the following situations:

1. $(p + q)$ is same as i_{r+1} and q is same as j_{r+1} ,
2. $(p + q)$ is same as j_{r+1} and q is same as i_{r+1} ,

3. If $i_r = 1$, we assume $j_r = 1$ as well as $S_r[i_r] = S_r[j_r] = 1$.

By considering the first two situations, we obtain $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = 0, i_r \neq 1) = \frac{2}{N^2}$ and by considering all the three situations we get $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = 0, i_r = 1) = \frac{3}{N^2}$. \square

Lemma 2 *During RC4 PRGA,*

1. $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = (N - 1), i_r \notin \{(N - 2), (N - 1)\}) = \frac{2}{N} - \frac{1}{N^2}$ and
2. $\Pr(z_r = z_{r+1} | S_r[i_r + 1] = (N - 1), i_r \in \{(N - 2), (N - 1)\}) = \frac{2}{N} - \frac{1}{N^3}$.

Proof We first investigate a configuration that is applicable for all values of i_r and then focus on two special cases one corresponding to $i_r = (N - 1)$ and the other corresponding to $i_r = (N - 2)$.

- **Scenario 1:** Let us consider a configuration in which $j_r = i_r + 1$. In this case if $S_r[i_r] = p$ where p is any arbitrary byte value (other than $(N - 1)$), we get $z_r = S_r[p + (N - 1)]$. In the next round (i.e., $(r + 1)$ -th round), i_{r+1} and j_{r+1} interchange their positions as compared to those of the r -th round. Hence, after the swap operation, $S_{r+1}[i_{r+1}] = p$ and $S_{r+1}[j_{r+1}] = (N - 1)$. Therefore, $z_{r+1} = S_{r+1}[(N - 1) + p] = z_r$. Here, we ignore the two cases where the position of z_r happens to coincide with either i_r or i_{r+1} as their effect in probability calculation is negligible. Thus, we consider that the probability associated with this configuration is $\frac{1}{N}$.
- **Scenario 2:** We now consider a special configuration that is applicable only for $i_r = (N - 2)$. In this case, we take $j_r = i_r + 2$ and if $S_r[i_r] = p$ where p is any arbitrary byte value (other than $(N - 1)$) then we assume $S_r[j_r] = (N - 2) - p$. This implies $z_r = S_r[p + (N - 2) - p] = S_r[N - 2] = p$. In the $(r + 1)$ -th round j_{r+1} becomes same as i_{r+1} which leads to the fact that $z_{r+1} = S_{r+1}[(N - 1) + (N - 1)] = S_{r+1}[N - 2] = p = z_r$. Since we have assumed three conditions in this configuration, where the first condition has probability of $\frac{1}{N}$, the second condition has the probability $\frac{1}{2}$ and the third condition has probability $\frac{2}{N}$, considering independence, the probability associated would be $\frac{1}{N^2}$.
- **Scenario 3:** We now consider a special configuration that is applicable only for $i_r = (N - 1)$. In this case, we take $j_r = i_r$ and if $S_r[i_r - 1] = p$ where p is any arbitrary byte value (other than $(N - 1)$) such that $p + (N - 1)$ is an even number (say $2k$ for a positive integer k), we assume $S_r[i_r] = S_r[j_r] = k$ or $S_r[i_r] = S_r[j_r] = \frac{N}{2} + k$. This implies $z_r = S_r[2k]$. In the $(r + 1)$ -th round j_{r+1} becomes same as $i_r - 1$ which leads to the fact that $z_{r+1} = S_{r+1}[p + (N - 1)] = S_{r+1}[2k]$. Assuming that the byte value indexed by $2k$ has not changed place in round $(r + 1)$, $z_{r+1} = z_r$ as desired. Since we have assumed two conditions in this configuration where the first condition has a probability of $\frac{1}{2N}$ and the second condition has a probability of $\frac{2}{N}$, the combined probability associated would be $\frac{1}{N^2}$.

Scenario 1 leads to the result

$$\Pr(z_r = z_{r+1} | S_r[i_r + 1] = (N - 1), i_r \notin \{(N - 2), (N - 1)\}) = \frac{1}{N} + (1 - \frac{1}{N}) \cdot \frac{1}{N} = \frac{2}{N} - \frac{1}{N^2}.$$

Scenarios 1, 2 and 3, when combined, lead to the result

$$\Pr(z_r = z_{r+1} | S_r[i_r + 1] = (N - 1), i_r \in \{(N - 2), (N - 1)\}) = \frac{1}{N} + \frac{1}{N^2} + (1 - \frac{1}{N} - \frac{1}{N^2}) \cdot \frac{1}{N} = \frac{2}{N} - \frac{1}{N^3}. \quad \square$$

Now let us present our main theorem. In this regard, we refer to a comment from [1]:

“Aside from the special case of $A = B$ and $G = 0$, we did not observe any additional significant deviations from the behaviour predicted by Result 2 [1] and our refinements of that result. However, a larger-scale computation might well reveal further fine structure. For example, as suggested by a reviewer, it is possible that there is a dependence of biases on i . Since i is known to the attacker, if such biases were present and of significant size, then this would result in exploitable behaviour.”

In fact that is what we study in this paper. The result our Theorem 3 actually points out the values of i for which significant biases do not exist and thus not exploitable.

Theorem 3 *During RC4 PRGA, assuming that the RC4 state is in a random permutation in the r -th round,*

1. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \notin \{0, 1, (N-2), (N-3)\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{4}{N^4},$
2. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \in \{(N-2), (N-3)\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{5}{N^4},$
3. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \in \{0, 1\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{6}{N^4}.$

Proof We prove this lemma by analyzing the following four cases.

[Case 1:] This corresponds to the configuration that was originally used by Mantin in [1] to prove the $ABTAB$ bias. The conditions are as follows -

- (i) $S_r[i_r + 1] = 1,$
- (ii) $j_r = i_r + 2$ and
- (iii) $j_{r+2} = i_r.$

These conditions lead to the desired outcome of $z_r = z_{r+2}$ and $z_{r+1} = z_{r+3}$. Since, $z_r = z_{r+1}$ (given condition), we get the pattern $AAAA$. The probability of obtaining the configuration is $\frac{1}{N^3}$ and using Mantin's result [1, Lemma 2] the probability of occurrence of the desired event (given that configuration) is $e^{-\frac{8}{N}}$. Hence, the combined probability associated with this configurations is $e^{-\frac{8}{N}} \cdot \frac{1}{N^3}$.

[Case 2:] In the first sub-case we assume the configuration

- (i) $S_{r+1}[i_{r+1} + 1] = 0.$

The probability associated with this configuration is $\frac{1}{N}$. Based on Lemma 1, we know that $\Pr(z_{r+2} = z_{r+1} | S_{r+1}[i_{r+1} + 1] = 0, i_{r+1} \neq 1) = \frac{2}{N^2}$ and $\Pr(z_{r+2} = z_{r+1} | S_{r+1}[i_{r+1} + 1] = 0, i_{r+1} = 1) = \frac{3}{N^2}$. Subsequently, we consider $\Pr(z_{r+3} = z_{r+2}) = \frac{1}{N}$ under the fair chance assumption. Hence, if $i_{r+1} \neq 1$, the probability of getting the desired outcome of $AAAA$ is $\frac{2}{N^4}$ and if $i_{r+1} = 1$, the probability of getting the desired outcome of $AAAA$ is $\frac{3}{N^4}$.

In the next sub-case we consider the configuration

$$(i) S_{r+2}[i_{r+2} + 1] = 0$$

The probability associated with this configuration is $\frac{1}{N}$. Based on Lemma 1, we know that $\Pr(z_{r+3} = z_{r+2} | S_{r+2}[i_{r+2} + 1] = 0, i_{r+2} \neq 2) = \frac{2}{N^2}$ and $\Pr(z_{r+3} = z_{r+2} | S_{r+2}[i_{r+2} + 1] = 0, i_{r+2} = 2) = \frac{3}{N^2}$. We also consider a fair chance of $\frac{1}{N}$ for $z_{r+1} = z_{r+2}$. Hence, if $i_{r+2} \neq 2$, the probability of obtaining the desired outcome of AAAA is $\frac{2}{N^4}$ and if $i_{r+2} = 2$, the probability of getting the desired outcome of AAAA is $\frac{3}{N^4}$.

[Case 3:] In the first sub-case we consider the configuration

$$(i) S_{r+1}[i_{r+1} + 1] = (N - 1)$$

$$(ii) j_{r+1} = i_{r+1} + 1$$

The probability associated with this configuration is $\frac{1}{N^2}$. This configuration ensures that $z_{r+2} = z_{r+1}$. So we need to investigate what happens in round $r + 3$.

Based on this configuration, we can say that in round r , j_r must have equal to $i_r + 3$. Let $S_r[i_r] = p$ and $S_r[j_r] = q$ where p and q are two arbitrary byte-values. In that case $z_r = S_r[p + q]$. The given configuration also implies that in round $(r + 2)$, $j_{r+2} = i_{r+1}$ where $S_{r+2}[j_{r+2}] = (N - 1)$. Therefore, in round $(r + 3)$, it would not be possible to have $S_{r+3}[i_{r+3}] + S_{r+3}[j_{r+3}] = (p + q)$, instead it would become $(q + s)$ for some arbitrary byte value $s \neq p$ in position $S_{r+3}[j_{r+3}]$ before the swap operation or $S_{r+3}[i_{r+3}]$ after the swap operation. So $z_{r+3} = S_{r+3}[q + s]$. The only way for z_{r+3} to be equal to z_r is if the permutation array byte indexed by $(p + q)$ moves to the new position indexed by $(q + s)$ in round $r + 3$. Using the argument similar to that used in Lemma 1 we get the probability of this event (given the configuration of the sub-case) as $\frac{2}{N^2}$. Hence, the probability of getting the desired outcome of AAAA in this sub-case is $\frac{2}{N^4}$.

By using the argument presented in Lemma 2, one can argue that for $i_{r+1} \in \{(N - 2), (N - 1)\}$, there can be an additional configuration (with the associated probability of $\frac{1}{N^3}$) in which the probability of getting the desired outcome of AAAA will be $\frac{2}{N^5}$. We consider any term of the order of $\frac{1}{N^5}$ as negligible in this theorem.

In the next sub-case we assume the configuration

$$(i) S_{r+2}[i_{r+2} + 1] = (N - 1),$$

$$(ii) j_{r+2} = i_{r+2} + 1.$$

The probability associated with this configuration is $\frac{1}{N^2}$. The configuration readily implies that $z_{r+2} = z_{r+3}$. By considering a fair chance of $\frac{1}{N}$ for the event of $z_{r+1} = z_{r+2}$, we get the probability of getting the desired outcome of AAAA in this sub-case as $\frac{1}{N^3}$.

By using the argument used in Lemma 2, for $i_{r+1} \in \{(N - 3), (N - 2)\}$, there can be an additional configuration (with the associated probability of $\frac{1}{N^3}$) in which the probability of getting the desired outcome of AAAA will be $\frac{1}{N^4}$.

[Case 4:] In this case we consider the rest of the configurations (i.e., complimentary to the combination of Cases 1, 2 and 3). The probability associated with this case would be $(1 - \frac{2}{N} - \frac{2}{N^2} - \frac{1}{N^3})$ for $i_{r+1} \notin \{(N - 3), (N - 2)\}$ and it would be $(1 - \frac{2}{N} - \frac{2}{N^2} - \frac{2}{N^3})$ for $i_{r+1} \in \{(N - 3), (N - 2)\}$. For each of these situations we consider that the desired configuration of $z_{r+1} = z_{r+2}$ and $z_{r+2} = z_{r+3}$ has the fair chance of $\frac{1}{N^2}$.

By combining all the cases we get the result for the desired outcome as

1. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \notin \{0, 1, (N-2), (N-3)\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{4}{N^4},$
2. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \in \{(N-2), (N-3)\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{5}{N^4},$
3. $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1}, i_{r+1} \in \{0, 1\})$
 $\approx \frac{1}{N^2} - \frac{(1-e^{-\frac{8}{N}})}{N^3} + \frac{6}{N^4}.$

□

For $N = 256$, all these three probabilities are less than $\frac{1}{N^2}$, where $\frac{1}{N^2}$ corresponds to the uniform random case.

3 Conclusion

In this note, we solve an open question that is related to Mantin's bias [3] in RC4 key-stream. This bias is till date the most significant long term one to distinguish RC4 key-stream from uniform random distribution. However, this is mostly a generic result with a few logical assumptions. Unfortunately, in a very few cases, the assumptions are not correct and such issues have been studied in great detail in [1]. The theoretical analysis could be formalized in [1], except one experimental observation, that could not be supported by theoretical argument. This is related to $\Pr((z_r, z_{r+1}) = (z_{r+2}, z_{r+3}) | z_r = z_{r+1})$, that is for the sub-string of the form AAAA. While the analysis of [1] could only point out to a positive bias, the experiments show that it is actually slightly negative in such a case. In this note, we prove this result with proper theoretical justification.

References

1. R. Bricout, S. Murphy, K. G. Paterson, T. van der Merwe. Analysing and exploiting the Mantin biases in RC4. *Designs Codes and Cryptography*, 86:743-770, 2018.
2. H. Finney. An RC4 cycle that can't happen. Post in sci.crypt, September 1994.
3. I. Mantin. Predicting and Distinguishing Attacks on RC4 Keystream Generator. *EUROCRYPT 2005*, pages 491506, vol. 3494, *Lecture Notes in Computer Science*, Springer.
4. K. G. Paterson, B. Poettering and J. C. N. Schuldt. Big Bias Hunting in Amazonia: Large-scale Computation and Exploitation of RC4 Biases. *ASIACRYPT 2014. LNCS, Part 1*, pp. 398-419, Vol. 8873, 2014.
5. S. SenGupta, S. Maitra, G. Paul, S. Sarkar. (Non-)Random Sequences from (Non-)Random Permutations – Analysis of RC4 stream cipher. *Journal of Cryptology*, 27(1):67-108, 2014
6. P. Sepehrdad, S. Vaudenay, and M. Vuagnoux. Statistical Attack on RC4 - Distinguishing WPA. *EUROCRYPT 2011. LNCS pp. 343-363, Vol. 6632*, 2011.