# Polar Coding for Processes with Memory

Eren Şaşoğlu, Ido Tal, *Senior Member, IEEE*

*Abstract*—We study polar coding for stochastic processes with memory. For example, a process may be defined by the joint distribution of the input and output of a channel. The memory may be present in the channel, the input, or both. We show that $\psi$-mixing processes polarize under the standard Arıkan transform, under a mild condition. We further show that the rate of polarization of the *low-entropy* synthetic channels is roughly $O(2^{-\sqrt{N}})$, where $N$ is the blocklength. That is, essentially the same rate as in the memoryless case.

*Index Terms*—Channels with memory, polar codes, mixing, periodic processes, fast polarization, rate of polarization.

## I. Introduction

**P**OLAR codes were invented by Arıkan [1] as a low-complexity method to achieve the capacity of symmetric binary-input memoryless channels. The technique that underlies these codes, called *polarization*, is quite versatile, and has since been applied to numerous classical memoryless problems in information theory.

Many practical sources and channels are not well-described by memoryless models. In wireless communication, for example, memory in the form of intersymbol interference is quite prominent due to multipath propagation, as are slow variations in channel conditions due to mobility. In practice, this type of memory is commonly handled by eliminating it, e.g., by augmenting the transmitter/receiver appropriately to create an overall memoryless channel. Memoryless coding techniques are then used for communication. Channel equalization, interleaving, and OFDM techniques are perhaps the most notable examples of this approach.

In contrast, we are interested here in whether polar coding can be used *directly* on channels and sources with memory. In addition to being of theoretical interest, such results may help simplify the design of communication or compression systems.

Little is known about the theory of polarization for settings with memory. In particular, it was shown in [2] that the successive cancellation decoding complexity of polar codes scales with the number of states of the underlying process, and thus is practical if the amount of memory in the system is modest. It was shown in [3, Chapter 5] that Arıkan's standard transform indeed polarizes a class of mixing processes with finite memory. Whether polarization takes place sufficiently fast to yield a coding theorem has been left open, however, and that is the problem we address here.

We first give a proof of polarization that is both simpler than the one given in [3], and holds for the more general class of $\psi$-mixing processes with finite $\psi_0$ (both concepts are defined in Section II). We further show that the asymptotic rate of polarization of the *low-entropy* synthetic channels is as in the memoryless case. Conversely, we show a simple counterexample of a process that is not $\psi$-mixing and which does not polarize because it is periodic. We remark that in [4], under additional assumptions, fast polarization is shown for the *high-entropy* synthetic channels.

## II. Setting

Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a stationary process, where the $Y_i$ take values in a finite alphabet $\mathcal{Y}$. We assume $X_i \in \{0, 1\}$ to keep the notation simple, but the results here can be generalized to arbitrary finite alphabets using standard techniques. See, for example, [3, Chapter 3]. We think of $X_i$ as a sequence to be estimated, and $Y_i$ as a sequence of observations related to $X_i$. In particular, $X_i$ may be the input sequence to a communication channel, with the corresponding channel output $Y_i$. Alternatively, $X_i$ may be the output of a data source to be compressed, and $Y_i$ may be the side information available to the decompressor.

A key property of the processes we consider is $\psi$-mixing. We follow[1] [5, Page 169] and say that a process $T_i$ is $\psi$-mixing if there exists a nonincreasing sequence $\psi_k \to 1$ as $k \to \infty$ such that

$$\mathbb{P}(A \cap B) \le \psi_k \mathbb{P}(A)\mathbb{P}(B) \qquad (1)$$

for all $A \in \sigma(T_{-\infty}^0)$ and $B \in \sigma(T_{k+1}^\infty)$, where $\sigma(\cdot)$ denotes the sigma-field generated by its argument. Since $\psi_k \to 1$, in a $\psi$-mixing process, any two events $A \in \sigma(T_{-\infty}^0)$ and $B \in \sigma(T_{k+1}^\infty)$ that are sufficiently separated in 'time' are almost independent. Namely, by [6, Definition 3.3, page 67], [6, Proposition 3.11, part a, page 76], and [6, Proposition 5.2, part III.a, page 153]

$$|\mathbb{P}(A \cap B) - \mathbb{P}(A)\mathbb{P}(B)| \le \frac{\psi_k - 1}{2}.$$

In this paper, we require for polarization that a process be $\psi$-mixing with finite $\psi_0$. Since this requirement appears several times, we make the following definition.

**Definition 1** (Promptly $\psi$-mixing). Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a stationary process, where $X_i \in \{0, 1\}$ and the $Y_i$ take values in a finite alphabet $\mathcal{Y}$. Such a process is called *promptly $\psi$-mixing* if it is $\psi$-mixing and $\psi_0 < \infty$.

Many source and channel models of practical importance satisfy our requirements of being promptly $\psi$-mixing. Specifically, this holds for a class of models with memory that have an

[1]To the best of our understanding, the first displayed equation on page 169 of [5] should be "$\sum_v \mu(uvw) \le \cdots$".

underlying ergodic Markov structure, as shown in [4, Lemma 5]. There, these processes are termed Finite-state, Aperiodic, Irreducible (hidden) Markov processes, or FAIM for short. The parameter $\psi_0$ plays an important role in this paper, and can be computed easily if the underlying process is FAIM [4, Equation 19].

We are interested in the effects of Arıkan's standard polar transform on stationary processes with memory. For this purpose, we let $U_1^N = X_1^N \mathsf{B}_N \mathsf{G}_N$, where the matrix multiplications are over the binary field, $N = 2^n$ for positive integers $n$, $\mathsf{G}_N$ is the $n$th Kronecker power of $\left( \begin{smallmatrix} 1 & 0 \\ 1 & 1 \end{smallmatrix} \right)$, and $\mathsf{B}_N$ is the $N \times N$ bit-reversal matrix. The conditional entropy rate of $X_i$ is defined as

$$\mathcal{H}_{X|Y} = \lim_{N \to \infty} \frac{1}{N} H(X_1^N | Y_1^N)$$
$$= \lim_{N \to \infty} \frac{1}{N} H(X_1^N, Y_1^N) - \lim_{N \to \infty} \frac{1}{N} H(Y_1^N).$$

The limits on the right-hand-side exist due to stationarity [7, Theorem 4.2.1]. Also useful for the analysis is the parameter

$$Z(A|B) = 2 \sum_{b \in \mathcal{B}} \sqrt{p_{A,B}(0,b) p_{A,B}(1,b)}$$

for random variables $A \in \{0, 1\}$ and $B \in \mathcal{B}$. Sometimes called the Bhattacharyya parameter, $Z(A|B)$ upper-bounds the error probability of optimally guessing $A$ by observing $B$. See, for example, [3, Proposition 2.2].

## III. MAIN RESULTS

The following two theorems relate to the polarization of promptly $\psi$-mixing process.

**Theorem 1** (Polarization). *Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a promptly $\psi$-mixing process, then for all $\epsilon > 0$*

$$\lim_{N \to \infty} \frac{1}{N} \left| \left\{ i : H(U_i | U_1^{i-1}, Y_1^N) > 1 - \epsilon \right\} \right| = \mathcal{H}_{X|Y},$$
$$\lim_{N \to \infty} \frac{1}{N} \left| \left\{ i : H(U_i | U_1^{i-1}, Y_1^N) < \epsilon \right\} \right| = 1 - \mathcal{H}_{X|Y}.$$

**Theorem 2** (Fast polarization of the low-entropy set). *Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a promptly $\psi$-mixing process, then for all $\beta < 1/2$*

$$\lim_{N \to \infty} \frac{1}{N} \left| \left\{ i : Z(U_i | U_1^{i-1}, Y_1^N) < 2^{-N^\beta} \right\} \right| = 1 - \mathcal{H}_{X|Y}.$$

We conjecture that an analog of Theorem 2 holds for the high-entropy set.

**Conjecture 3** (Fast polarization of the high-entropy set). *Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a promptly $\psi$-mixing process, then for all $\beta < 1/2$*

$$\lim_{N \to \infty} \frac{1}{N} \left| \left\{ i : Z(U_i | U_1^{i-1}, Y_1^N) > 1 - 2^{-N^\beta} \right\} \right| = \mathcal{H}_{X|Y}.$$

Resolving the above conjecture would be an important step for polar codes. We refer the reader to [4, Theorem 13], which shows that the conjecture indeed holds if the process is FAIM. To recap, assuming that the process $(X_i, Y_i)$ is governed by an underlying state sequence having a certain structure allows one to prove Conjecture 3. However, we will *not* assume an underlying state sequence when proving Theorems 1 and 2.
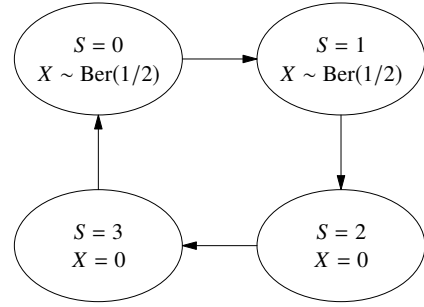


Fig. 1. A periodic data source that does not polarize. The source output is Bernoulli $1/2$ for two consecutive states and zero for next two consecutive states. There is no side information, i.e., $Y_i$ is constant.

As a concrete example of the distinction between promptly $\psi$-mixing and FAIM processes, consider the family of processes given in [8, Example 3]. Each such process $(X_i')$, $i \in \mathbb{Z}$, is $\psi$-mixing, with $\psi_0 < \infty$. Also, the support of each $X_i'$ is $[0, 1)$. Next, fix such a process, and let $B$ be some Borel set on $[0, 1)$. For example, $B = [0, 1/2)$. Define the process $(X_i)$, $i \in \mathbb{Z}$, such that $X_i = 1$ if $X_i' \in B$, and $X_i = 0$ otherwise. Since the process $(X_i)$ is a marginalization of $(X_i')$, we deduce from (1) that $(X_i)$ is also $\psi$-mixing, with finite $\psi_0$. That is, we deduce that $(X_i)$ is promptly $\psi$-mixing, and hence Theorems 1 and 2 are applicable. However, since the underlying process $(X_i')$ is not finite state, it is not FAIM, and thus it is not clear if Conjecture 3 holds for $(X_i)$.

The following theorem shows an example of a process that has memory and that *does not* polarize because it is periodic.

**Theorem 4** (Periodic processes may not polarize). *The stationary periodic Markov process described in Figure 1 does not polarize. Indeed, for all $\frac{5N}{8} < i \le \frac{6N}{8}$,*

$$\left| H(U_i | U_1^{i-1}) - \frac{1}{2} \right| \le \epsilon_N, \quad \lim_{N \to \infty} \epsilon_N = 0. \quad (2)$$

## IV. NOTATION

We will prove the above theorems in the following sections. Throughout, we will use the shorthand

$$H^{\mathbf{b}} = H(U_i | U_1^{i-1}, Y_1^N),$$
$$Z^{\mathbf{b}} = Z(U_i | U_1^{i-1}, Y_1^N),$$

where $\mathbf{b} \in \{0, 1\}^n$ is the $n$-bit binary expansion of $i - 1 \in \{0, \dots, N - 1\}$. We will omit the ranges of indices when they are clear from context. The following are immediate from the definition of $\mathsf{B}_N \mathsf{G}_N$:

$$H^{\mathbf{b}0} = H(U_{2i-1} | U_1^{2i-2}, Y_1^{2N})$$
$$H^{\mathbf{b}1} = H(U_{2i} | U_1^{2i-1}, Y_1^{2N})$$

for all $\mathbf{b} \in \{0, 1\}^n$. These identities also hold when the $H$'s are replaced by $Z$'s. Further, if we let $B_1, B_2, \dots$ be a sequence of i.i.d. $\mathrm{Ber}(1/2)$ random variables, then it is easy to see that the random variables $H_n = H^{B_1 \dots B_n}$ and $Z_n = Z^{B_1 \dots B_n}$ are uniformly distributed over the sets of $H^{\mathbf{b}}$'s and $Z^{\mathbf{b}}$'s, respectively. Theorems 1 and 2 are then equivalent to

**Theorem 5.** *Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a promptly $\psi$-mixing process, then for all $\epsilon > 0$*

$$\lim_{n \to \infty} \mathbb{P}\left(H_n > 1 - \epsilon\right) = \mathcal{H}_{X|Y} \,,$$
$$\lim_{n \to \infty} \mathbb{P}\left(H_n < \epsilon\right) = 1 - \mathcal{H}_{X|Y} \,.$$

**Theorem 6.** *Let $(X_i, Y_i)$, $i \in \mathbb{Z}$, be a promptly $\psi$-mixing process, then for all $\beta < 1/2$*

$$\lim_{n \to \infty} \mathbb{P}\left(Z_n < 2^{-N^\beta}\right) = 1 - \mathcal{H}_{X|Y} \,.$$

As is usual in proofs of polarization, we will analyze how the entropies and Bhattacharyya parameters evolve in a single recursion of the polarization transform. That is, when two smaller polarization blocks are combined to form a larger block. Due to the dependence between the combined blocks, we will need to keep track of more random variables than is required in the analysis of the memoryless case. The following shorthand will then be useful:

$$U_1^N = X_1^N \mathsf{B}_N \mathsf{G}_N \,, \tag{3a}$$
$$V_1^N = X_{N+1}^{2N} \mathsf{B}_N \mathsf{G}_N \,, \tag{3b}$$
$$Q_i = (U_1^{i-1}, Y_1^N) \,, \tag{3c}$$
$$R_i = (V_1^{i-1}, Y_{N+1}^{2N}) \,. \tag{3d}$$

## V. PROOF OF THEOREM 1

Throughout this section, we assume that $(X_i, Y_i)$, $i \in \mathbb{Z}$, is a promptly $\psi$-mixing process. We will prove Theorem 1 by showing that $H_n$ converges almost surely (a.s.) and in $L^1$ to a $\{0, 1\}$-valued random variable $H_\infty$. As in [1], we first show that $H_\infty \in [0, 1]$.

**Lemma 7.** *The sequence $H_n$ converges a.s. and in $L^1$ to a random variable $H_\infty \in [0, 1]$.*

*Proof:* Recall that for $i = 1 + (B_1 \dots B_n)_2$ we have that

$$H_n = H(U_i | U_1^{i-1}, Y_1^N) = H(U_i | Q_i) \,.$$

Also, for $i$ as above,

$$H_{n+1} = \begin{cases} H(U_i + V_i | Q_i, R_i) \,, & \text{if } B_{n+1} = 0 \,, \\ H(V_i | Q_i, R_i, U_i + V_i) \,, & \text{if } B_{n+1} = 1 \,. \end{cases} \tag{4}$$

Next, note that

$$\begin{aligned} H(U_i &+ V_i | Q_i, R_i) + H(V_i | Q_i, R_i, U_i + V_i) \\ &= H(U_i, V_i | Q_i, R_i) \\ &\leq H(U_i | Q_i) + H(V_i | R_i) \\ &= 2H(U_i | Q_i) \,, \end{aligned}$$

where the inequality follows since conditioning reduces entropy, and the last step follows from stationarity. Thus, since $B_{n+1}$ is uniform, $\mathbb{E}[H_{n+1} | H_1, \dots, H_n] \leq H_n$. The entropy is bounded, $H_n \in [0, 1]$, and thus it follows that $H_1, H_2, \dots$ is a bounded supermartingale. We conclude by [9, Theorem 9.4.5] that it converges almost surely and in $L^1$ to a $[0, 1]$-valued random variable $H_\infty$. $\blacksquare$

Our approach to proving that $H_\infty \in \{0, 1\}$ shares similarities with the proof in [10, Section 2.2] for the memoryless case. In essence, the proof there hinges on [10, Lemma 2.2], which

shows that if $H(U_i | V_i)$ is bounded away from both 0 and 1, then $H(U_i + V_i | Q_i, R_i) - H(U_i | Q_i)$ is bounded away from 0. Informally, if $H_n$ has not polarized, then it has not converged. Thus, our main focus now is on $H(U_i + V_i | Q_i, R_i)$.

Recalling the definitions of $Q_i$ and $R_i$ in (3), we see that $Y_N \in Q_i$ and $Y_{N+1} \in R_i$. Since $Y_N$ and $Y_{N+1}$ are generally dependent, we deduce that $Q_i$ and $R_i$ are generally dependent as well. However, suppose that $U_i$ and $V_i$ were independent given $Q_i$ and $R_i$. This is not generally true, but if it were, we would be closer to the memoryless setting and our task of analyzing $H(U_i + V_i | Q_i, R_i)$ would be simpler. Informally, inequality (5) in the next lemma shows that this is "almost true".

**Lemma 8.** *For any $\epsilon > 0$, the fraction of indices $i$ for which*

$$I(U_i; V_i | Q_i, R_i) < \epsilon \,, \tag{5}$$
$$I(U_i; R_i | Q_i) < \epsilon \,, \tag{6}$$
$$I(V_i; Q_i | R_i) < \epsilon \,, \tag{7}$$

*approaches 1 as $N \to \infty$.*

*Proof:* We only prove the first and the third inequalities, since the second inequality follows by symmetry. We have

$$\begin{aligned} &\log(\psi_0) \\ &\geq \mathbb{E}\left[\log \frac{p_{X_1^{2N}, Y_1^{2N}}(X_1^{2N}, Y_1^{2N})}{p_{X_1^N, Y_1^N}(X_1^N, Y_1^N) \cdot p_{X_{N+1}^{2N}, Y_{N+1}^{2N}}(X_{N+1}^{2N}, Y_{N+1}^{2N})}\right] \\ &= I(X_1^N, Y_1^N; X_{N+1}^{2N}, Y_{N+1}^{2N}) \\ &= I(U_1^N, Y_1^N; V_1^N, Y_{N+1}^{2N}) \\ &= I(Y_1^N; V_1^N, Y_{N+1}^{2N}) + I(U_1^N; V_1^N, Y_{N+1}^{2N} | Y_1^N) \\ &\geq I(U_1^N; V_1^N, Y_{N+1}^{2N} | Y_1^N) \\ &= \sum_{i=1}^N I(U_i; V_1^N, Y_{N+1}^{2N} | Y_1^N, U_1^{i-1}) \\ &= \sum_{i=1}^N I(U_i; R_i, V_i, V_{i+1}^N | Q_i) \,. \end{aligned}$$

The first inequality above follows from the definition of $\psi_0$. Since all terms inside the last sum are non-negative, it follows that at most $\sqrt{\log(\psi_0)N}$ (a vanishing fraction) of them are at least $\sqrt{\log(\psi_0)/N}$ (a vanishing quantity). Thus, to conclude the proof, it suffices to show that the $i$th term is greater than both $I(U_i; R_i | Q_i)$ and $I(U_i; V_i | Q_i, R_i)$. Indeed,

$$\begin{aligned} I(U_i; &R_i, V_i, V_{i+1}^N | Q_i) \\ &= I(U_i; R_i | Q_i) + I(U_i; V_i | Q_i, R_i) + I(U_i; V_{i+1}^N | Q_i, V_i, R_i), \end{aligned}$$

and all the terms are non-negative. $\blacksquare$

In fact (5) is the only inequality we will need from Lemma 8. We have stated (6) and (7) to serve as motivation for the following. Namely, for index $1 \leq i \leq N$, we now introduce the random variables $\tilde{U}_i$ and $\tilde{V}_i$. The joint distribution of $(X_1^{2N}, Y_1^{2N}, U_1^N, V_1^N, Q_1^N, R_1^N, \tilde{U}_1^N, \tilde{V}_1^N)$ is defined as follows. First $X_1^{2N}$ and $Y_1^{2N}$ are picked according to the process distribution. This uniquely determines the values of $U_1^N, V_1^N, Q_1^N$, and $R_1^N$, according to (3). Finally, for each $i = 1, 2, \dots, n$ we pick $\tilde{U}_i$ and $\tilde{V}_i$ independently according to the marginal

distributions $p_{U_i|Q_i}(\cdot|q_i)$ and $p_{V_i|R_i}(\cdot|r_i)$, where $q_i$ and $r_i$ are the realizations of $Q_i$ and $R_i$. The key property to note is that the joint distribution of $(\tilde{U}_i, \tilde{V}_i)$ with $(Q_i, R_i)$ is of the form

$$p_{\tilde{U}_i, \tilde{V}_i, Q_i, R_i}(\tilde{u}_i, \tilde{v}_i, q_i, r_i)$$
$$= p_{U_i|Q_i}(\tilde{u}_i|q_i) \cdot p_{V_i|R_i}(\tilde{v}_i|r_i) \cdot p_{Q_i, R_i}(q_i, r_i) . \quad (8)$$

Thus, by definition, $\tilde{U}_i$ and $\tilde{V}_i$ are independent given $Q_i$ and $R_i$. In fact, more is true: if we replace $U_i$ and $V_i$ by $\tilde{U}_i$ and $\tilde{V}_i$, respectively, in (5)–(7), then all the mutual informations become zero. See (37)–(39) in the appendix for a proof of this fact.

As explained, it will be easier to analyze $H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i)$ in place of $H(U_i + V_i | Q_i, R_i)$. The following corollary to Lemma 8 serves as justification for this shift, since it shows that the two quantities are "close". It is proved in the appendix and will be used later on.

**Corollary 9.** *For any $\epsilon > 0$, the fraction of indices $i$ for which*

$$|H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(U_i + V_i | Q_i, R_i)| < \epsilon \quad (9)$$

*approaches* 1 *as* $N \to \infty$.

Note that by (8),

$$H(\tilde{U}_i | Q_i, R_i) = H(\tilde{U}_i | Q_i) = H(U_i | Q_i) . \quad (10)$$

Thus, in light of this and Corollary 9, we will consider the difference $H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(\tilde{U}_i | Q_i)$ as a proxy for our ultimate quantity of interest, $H(U_i + V_i | Q_i, R_i) - H(U_i | Q_i)$. Note that in order to save space, we will usually prefer writing $H(\tilde{U}_i | Q_i)$ in place of the longer but more informative $H(\tilde{U}_i | Q_i, R_i)$. The same remark applies to $H(\tilde{V}_i | Q_i)$ versus $H(\tilde{V}_i | Q_i, R_i)$, which are also equal due to (8).

Recall that we aim to mimic the memoryless proof in [10, Section 2.2] as much as possible. Hence our informal strategy will soon be the following: show that if $H(\tilde{U}_i | Q_i)$ is bounded away from both 0 and 1, then $H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(\tilde{U}_i | Q_i)$ is bounded away from 0.

We now motivate the following lemma. Namely, we will now introduce an apparent difficulty, which the following lemma will resolve. Recall that we prefer analyzing $\tilde{U}_i$ and $\tilde{V}_i$ over $U_i$ and $V_i$, since the former are independent given $(Q_i, R_i)$. In contrast, as we have already mentioned, $Q_i$ and $R_i$ are generally dependent. This presents an apparent problem with the strategy outlined in the previous paragraph: suppose $H(\tilde{U}_i | Q_i)$ is bounded away from both 0 and 1. Suppose further that for every value $q_i$ that $Q_i$ can take, we have that $H(\tilde{U}_i | Q_i = q_i)$ is either 0 or 1. That is, imagine what is effectively an erasure channel, mapping $\tilde{U}_i$ to $Q_i$. By stationarity, the same property must hold for $H(\tilde{V}_i | R_i = r_i)$. Now, since $Q_i$ and $R_i$ are *not* independent, it is conceivable that they collude, i.e., that it is always the case that the values $q_i$ and $r_i$ that the random variables $Q_i$ and $R_i$ respectively take are such that either $H(\tilde{U}_i | Q_i = q_i) = H(\tilde{V}_i | R_i = r_i) = 0$ or $H(\tilde{U}_i | Q_i = q_i) = H(\tilde{V}_i | R_i = r_i) = 1$. In other words, in two consecutive uses of the above channel, we always have either two non-erasures or two erasures. In such a case, it is easy to see that $H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(\tilde{U}_i | Q_i)$ is identically 0. That is, if the above assumptions are valid, our plan is

doomed to fail: we have an apparent counter-example in which $H(\tilde{U}_i | Q_i)$ is bounded away from both 0 and 1, yet the difference $H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(\tilde{U}_i | Q_i)$ is not bounded away from 0. Informally, an important corollary of the following lemma is that such synchronized erasures cannot happen. That is, as intuition for the following lemma, think of $A = 1$ ($B = 1$) as indicating that $Q_i$ ($R_i$) corresponds to an erasure of $\tilde{U}_i$ ($\tilde{V}_i$).

**Lemma 10.** *For all $\xi > 0$, there exists $N_0$ and $\delta(\xi) > 0$ such that for all $N > N_0$ and all $\{0, 1\}$-valued random variables $A = f(X_1^N, Y_1^N)$ and $B = f(X_{N+1}^{2N}, Y_{N+1}^{2N})$,*

$$p_A(1) \in (\xi, 1 - \xi) \quad \text{implies} \quad p_{A,B}(1, 0) > \delta(\xi) .$$

*Proof:* Let us start by explaining informally why the claim is true. Define $C = f(X_{2N+1}^{3N}, Y_{2N+1}^{3N})$, and suppose to the contrary that $B$ equals $A$ with very high probability. Hence, by stationarity, $C$ equals $B$ with very high probability. We conclude that $A$ equals $C$ with probability very close to 1, a contradiction to the mixing property.

Let us now give a formal proof. First, clearly, we may assume that $\xi \leq 1/2$, or else the claim is vacuous. We have

$$2p_{A,B}(1, 0) = p_{A,B}(1, 0) + p_{B,C}(1, 0)$$
$$\geq p_{A,B,C}(1, 0, 0) + p_{A,B,C}(1, 1, 0)$$
$$= p_{A,C}(1, 0)$$
$$= p_A(1) - p_{A,C}(1, 1)$$
$$\geq p_A(1)(1 - \psi_N p_C(1))$$
$$= p_A(1)(1 - \psi_N p_A(1))$$

where the first and last equalities are due to stationarity. Recall that $\psi_N$ converges to 1 from above. We now commit to an $N_0$ such that $\psi_N < \frac{1}{\sqrt{1-\xi}}$ for all $N > N_0$. Recalling that $p_A(1) \in (\xi, 1 - \xi)$, we can bound the last term in the above displayed equation as

$$p_A(1)(1 - \psi_N p_A(1)) > p_A(1) \left(1 - \frac{1}{\sqrt{1-\xi}} p_A(1)\right)$$
$$> p_A(1) \left(1 - \frac{1}{\sqrt{1-\xi}}(1 - \xi)\right)$$
$$= p_A(1) \left(1 - \sqrt{1-\xi}\right)$$
$$> \xi \left(1 - \sqrt{1-\xi}\right) ,$$

under the assumption that $N > N_0$. That is, for all $N > N_0$, we deduce that $2p_{A,B}(1, 0) > \xi \left(1 - \sqrt{1-\xi}\right)$. Thus, we take $\delta(\xi) = \xi \left(1 - \sqrt{1-\xi}\right)/2$. $\blacksquare$

The next lemma will be instrumental in the following setting. Let $q_i$ and $r_i$ be given. Assume that $H(\tilde{U}_i | Q_i = q_i)$ and $H(\tilde{V}_i | R_i = r_i)$ are not both close to 0, nor are they both close to 1. To emphasize: we only rule out the case where both entropies are close to each other and extremal. Then, we will deduce from the following lemma that the entropy $H(\tilde{U}_i + \tilde{V}_i | Q_i = q_i, R_i = r_i)$ is non-negligibly greater than the mean of $H(\tilde{U}_i | Q_i = q_i)$ and $H(\tilde{V}_i | R_i = r_i)$. The proof is given in the appendix.

**Lemma 11.** *Let A and B be independent binary random variables. For every $\xi > 0$, there exists $\Delta(\xi) > 0$ such that*

$$\max\{H(A), H(B)\} > \xi$$

*and*

$$\min\{H(A), H(B)\} < 1 - \xi$$

*imply*

$$H(A + B) > \frac{H(A) + H(B)}{2} + \Delta(\xi) \ .$$

We are now ready to state and prove the cardinal lemma of this section. Informally, we now show that if $H_n = H(\tilde{U}_i|Q_i)$ has not polarized, then it has not converged.

**Lemma 12.** *For all $\xi > 0$ there exist $\theta(\xi) > 0$ and $N_0$ such that for all $N > N_0$ and all $1 \le i \le N$,*

$$H(\tilde{U}_i|Q_i) \in (3\xi, 1 - 3\xi) \ implies$$
$$H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i) - H(\tilde{U}_i|Q_i) > 2\theta(\xi) \ . \quad (11)$$

*Proof:* For a given $\xi > 0$, let $\theta(\xi) = \delta(\xi)\Delta(\xi)/2$, where $\delta(\xi)$ and $\Delta(\xi)$ are as in Lemmas 10 and 11. Also, let $N_0$ be as in Lemma 10. The motivation for these choices will soon become apparent. Set $N > N_0$ and let $i$ be given. We must show that (11) holds.

Let us first introduce some notation. Let $X$ and $Y$ be generic random variables in this paragraph. Note that $H(X|Y = y)$ is a function of $y$, which we denote in this paragraph as $g(y)$. We shall denote $g(Y)$ as $H(X|\underline{Y})$. We emphasize: the underline in $H(X|\underline{Y})$ signifies that we are dealing with a random variable, which is a function of the underlined quantity.[2] A simple and concise result of this definition is that

$$H(X|Y) = \mathbb{E}\left[H(X|\underline{Y})\right] \ .$$

Assume that

$$H(\tilde{U}_i|Q_i) \in (3\xi, 1 - 3\xi) \ , \quad (12)$$

otherwise the claim is vacuous. Together with our assumption that $\xi$ is positive, the above trivially implies that

$$0 < \xi < \frac{1}{6} \ . \quad (13)$$

Recall that $H(\tilde{U}_i|Q_i) = \mathbb{E}[H(\tilde{U}_i|\underline{Q_i})]$. In order to keep the notation light, we further denote

$$\alpha = \mathbb{P}\left(H(\tilde{U}_i|\underline{Q_i}) \le \xi\right) \ , \quad (14)$$

$$\beta = \mathbb{P}\left(H(\tilde{U}_i|\underline{Q_i}) \in (\xi, 1 - \xi)\right) \ , \quad (15)$$

$$\gamma = \mathbb{P}\left(H(\tilde{U}_i|\underline{Q_i}) \ge 1 - \xi\right) \ . \quad (16)$$

We will prove (11) for two cases, $\beta < \xi$ and $\beta \ge \xi$.

*Case 1*: Consider first the case in which

$$\beta < \xi \ . \quad (17)$$

---

[2]One might benefit from verbalizing $H(X|\underline{Y})$ as "the conditional entropy of $X$, as a function of $Y$". Note that this definition is similar to the definition of $\mathbb{E}[X|Y]$, which is usually taken to be a random variable that is a function of $Y$.

In words: the probability that $Q_i$ equals a value $q_i$ for which $H(\tilde{U}_i|Q_i = q_i) \in (\xi, 1 - \xi)$ is denoted $\beta$, and is less than $\xi$. Informally, for $\xi > 0$ small, this means that a typical realization of $Q_i$ implies either an "almost certainty" regarding the value of $\tilde{U}_i$ or an "almost erasure".

Informally, we next show that for $\xi$ "small", and under the assumptions (12) and (17), the probability of an "almost erasure", $\gamma$, is not trivial. That is, for a lower bound on $\gamma$, we employ (12)–(17) and deduce that

$$3\xi < H(\tilde{U}_i|Q_i) \le \alpha \cdot \xi + \beta \cdot (1 - \xi) + \gamma \cdot 1$$
$$< \alpha \cdot \xi + \xi \cdot (1 - \xi) + \gamma \cdot 1$$
$$\le (1 - \gamma) \cdot \xi + \xi \cdot (1 - \xi) + \gamma \cdot 1 \ ,$$

where the last inequality follows from $\alpha \le 1 - \gamma$ (since $\alpha$, $\beta$ and $\gamma$ are probabilities summing to 1). Rearranging the above gives

$$\gamma > \frac{\xi + \xi^2}{1 - \xi} \ . \quad (18)$$

For an upper bound on $\gamma$, we again use (12)–(17) to show that

$$1 - 3\xi > H(\tilde{U}_i|Q_i) \ge \alpha \cdot 0 + \beta \cdot \xi + \gamma \cdot (1 - \xi)$$
$$\ge \gamma \cdot (1 - \xi) \ .$$

Rearranging gives

$$\gamma < \frac{1 - 3\xi}{1 - \xi} \ . \quad (19)$$

By (13), (18), (19), and some simple algebra, we deduce that

$$\gamma \in (\xi, 1 - \xi) \ . \quad (20)$$

Recall that by (3), $Q_i$ is a deterministic function of $X_1^N$ and $Y_1^N$. Thus, there clearly exists a $\{0, 1\}$-valued function $f$ such that $f(X_1^N, Y_1^N)$ equals 1 iff $H(\tilde{U}_i|\underline{Q_i}) \ge 1 - \xi$. That is, for $\xi$ "small", $f(X_1^N, Y_1^N)$ equals 1 iff $Q_i$ corresponds to an "almost erasure" of $\tilde{U}_i$. By the symmetry of definitions in (3) and (8), the above $f$ also satisfies that $f(X_{N+1}^{2N}, Y_{N+1}^{2N}) = 1$ iff $H(\tilde{V}_i|\underline{R_i}) \ge 1 - \xi$. Recalling (16), (20), and our definition of $f$, we get from Lemma 10 that

$$\mathbb{P}\left(H(\tilde{U}_i|\underline{Q_i}) \ge 1 - \xi \ , \ H(\tilde{V}_i|\underline{R_i}) < 1 - \xi\right) > \delta(\xi) \ . \quad (21)$$

Let us now define the "good" (with respect to Lemma 11) set $G$ of pairs $(q_i, r_i)$ as

$$G = \Big\{(q_i, r_i) :$$
$$\max\{H(\tilde{U}_i|Q_i = q_i), H(\tilde{V}_i|R_i = r_i)\} > \xi$$
$$\text{and}$$
$$\min\{H(\tilde{U}_i|Q_i = q_i), H(\tilde{V}_i|R_i = r_i)\} < 1 - \xi\Big\} \ .$$

By (13) and (21),

$$\mathbb{P}((Q_i, R_i) \in G) > \delta(\xi) \ . \quad (22)$$

We are now ready to show (11). We claim that

$$H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i) - H(\tilde{U}_i|Q_i)$$

$$= H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i) - \frac{H(\tilde{U}_i|Q_i) + H(\tilde{V}_i|Q_i)}{2}$$

$$= \sum_{(q_i, r_i)} p_{Q_i, R_i}(q_i, r_i)\left(H(\tilde{U}_i + \tilde{V}_i|Q_i = q_i, R_i = r_i)\right.$$

$$\left. - \frac{H(\tilde{U}_i|Q_i = q_i) + H(\tilde{V}_i|R_i = r_i)}{2}\right)$$

$$\geq \sum_{(q_i, r_i)\in G} p_{Q_i, R_i}(q_i, r_i)\left(H(\tilde{U}_i + \tilde{V}_i|Q_i = q_i, R_i = r_i)\right.$$

$$\left. - \frac{H(\tilde{U}_i|Q_i = q_i) + H(\tilde{V}_i|R_i = r_i)}{2}\right)$$

$$> \delta(\xi) \cdot \Delta(\xi) . \tag{23}$$

Indeed, the first equality is by stationarity; the first inequality is because the term in brackets is always non-negative[3]; the last inequality is by Lemma 11 and (22). Thus, recalling that we have taken $\theta(\xi) = \delta(\xi)\Delta(\xi)/2$, we have proved (11), under the assumptions (12) and (17).

*Case 2*: We now aim to prove (11), under the assumptions (12) and

$$\beta \geq \xi . \tag{24}$$

This will be shorter, informally because we are now assuming that the probability of $Q_i$ equalling a value for which the entropy of $\tilde{U}_i$ is "moderate" is "sufficiently high". We start by noticing that under the event $H(\tilde{U}_i|Q_i) \in (\xi, 1 - \xi)$ used to define $\beta$ in (15), we have that $(Q_i, \overline{R}_i) \in G$. Thus, the LHS of (22) is lower bounded by $\beta$. Next, we will show that $\beta > \delta(\xi)$, and hence (22) holds. Indeed, recall from the proof of Lemma 10 that $\delta(\xi) = \xi\left(1 - \sqrt{1-\xi}\right)/2 < \xi$. By this and (24) we deduce that (22) holds, and the proof continues as before. Hence, we have proved (11), under the assumptions (12) and (24). ∎

The following corollary to Lemma 12 shifts us back to $U_i$ and $V_i$ from $\tilde{U}_i$ and $\tilde{V}_i$.

**Corollary 13.** *For all $\xi > 0$ there exists $\theta(\xi) > 0$ such that*

$$H(U_i|Q_i) \in (3\xi, 1 - 3\xi) \text{ implies}$$
$$H(U_i + V_i|Q_i, R_i) - H(U_i|Q_i) > \theta(\xi) \tag{25}$$

*for a fraction of indices $i \in \{1, \ldots, N\}$ approaching 1 as $N \to \infty$.*

*Proof:* Let $\xi > 0$ be given and take $\theta(\xi)$ as in Lemma 12. Also, take $N_0$ as in Lemma 12. Fix $N > N_0$, and let $\mathcal{A}$ be the set of indices for which (9) holds, for $\epsilon = \theta(\xi)$. Note that by Corollary 9, the fraction of indices in $\mathcal{A}$ approaches 1 as $N$ tends to infinity. By assumption, for all indices $i$, and specifically for all $i \in \mathcal{A}$, we have that (11) holds. Our aim is to show that (25) holds for all $i \in \mathcal{A}$ as well. Indeed, let $i \in \mathcal{A}$. If $H(U_i|Q_i) \notin (3\xi, 1 - 3\xi)$, then (25) holds trivially.

[3]Note that $H(\tilde{U}_i + \tilde{V}_i|Q_i = q_i, R_i = r_i) \geq H(\tilde{U}_i + \tilde{V}_i|\tilde{V}_i, Q_i = q_i, R_i = r_i) = H(\tilde{U}_i|Q_i = q_i)$, and we can similarly lower bound by $H(\tilde{V}_i|R_i = r_i)$.

Thus, assume that $H(U_i|Q_i) \in (3\xi, 1 - 3\xi)$. By (10), this is equivalent to $H(\tilde{U}_i|Q_i) \in (3\xi, 1 - 3\xi)$. Thus, by assumption, the consequent in (11) holds. We deduce that

$$H(U_i + V_i|Q_i, R_i) - H(U_i|Q_i)$$
$$= H(U_i + V_i|Q_i, R_i) - H(\tilde{U}_i|Q_i)$$
$$= H(U_i + V_i|Q_i, R_i) - H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i)$$
$$\quad + H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i) - H(\tilde{U}_i|Q_i)$$
$$> -\theta(\xi) + H(\tilde{U}_i + \tilde{V}_i|Q_i, R_i) - H(\tilde{U}_i|Q_i)$$
$$> -\theta(\xi) + 2\theta(\xi)$$
$$= \theta(\xi) ,$$

where the first equality follows from (10); the first inequality follows from (9), recalling that $\epsilon = \theta(\xi)$; and the last inequality follows from our assumption that the consequent in (11) holds. Thus, the consequent in (25) holds. ∎

With Corollary 13 at hand, the proof of Theorem 1 is forthcoming. Indeed, we now essentially repeat the arguments in [1].

*Proof of Theorem 1:* Recall that in Lemma 7, we proved that $H_n$ converges a.s. and in $L^1$ to $H_\infty \in [0, 1]$. We next show that $H_\infty$ converges a.s. to either 0 or 1. That is, we show that for all $\epsilon > 0$, $\mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon)) = 0$. Indeed, assume to the contrary that there exists $\epsilon > 0$ for which

$$\mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon)) > \rho , \tag{26}$$

where $\rho > 0$. Next, note that

$$\mathbb{P}(H_n \in (\epsilon/2, 1 - \epsilon/2))$$
$$\geq \mathbb{P}(H_n \in (\epsilon/2, 1 - \epsilon/2) \quad \text{and} \quad |H_n - H_\infty| < \epsilon/2)$$
$$\geq \mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon) \quad \text{and} \quad |H_n - H_\infty| < \epsilon/2)$$
$$= \mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon))$$
$$\quad - \mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon) \quad \text{and} \quad |H_n - H_\infty| \geq \epsilon/2)$$
$$\geq \mathbb{P}(H_\infty \in (\epsilon, 1 - \epsilon)) - \mathbb{P}(|H_n - H_\infty| \geq \epsilon/2)$$
$$> \rho - \mathbb{P}(|H_n - H_\infty| \geq \epsilon/2) ,$$

where the last inequality follows from (26). Since a.s. convergence implies convergence in probability [9, Theorem 4.1.2], we deduce from the above that

$$\liminf_{n\to\infty} \mathbb{P}(H_n \in (\epsilon/2, 1 - \epsilon/2)) \geq \rho .$$

Recall the definition of $H_{n+1}$ in (4), and further recall that $B_{n+1}$ equals 0 with probability 1/2. Now, take $\xi$ such that $3\xi = \epsilon/2$. We deduce from Corollary 13 that for $n$ large enough,

$$\mathbb{P}(|H_{n+1} - H_n| > \theta(\xi)) > \frac{\rho}{4} .$$

However, this implies that $H_n$ cannot converge in probability to $H_\infty$, a contradiction to what was stated earlier. We have proven that $H_\infty \in \{0, 1\}$ a.s.

We now show that

$$\lim_{n\to\infty} \mathbb{E}[H_n] = \mathbb{E}[H_\infty] . \tag{27}$$

Indeed,

$$\mathbb{E}[-|H_n - H_\infty|] \leq \mathbb{E}[H_n - H_\infty] \leq \mathbb{E}[|H_n - H_\infty|] ,$$

and by the $L^1$ convergence of $H_n$ to $H_\infty$ and the sandwich property, the limit of the middle term is 0. By definition, $\lim_{n\to\infty} \mathbb{E}[H_n] = \mathcal{H}_{X|Y}$. Hence, since $H_\infty \in \{0,1\}$ almost surely, we must have that $\mathbb{P}(H_\infty = 1) = 1 - \mathbb{P}(H_\infty = 0) = \mathcal{H}_{X|Y}$. Recalling that $H_n$ converges in probability to $H_\infty$, the claim in Theorem 5 follows. We end by noting that Theorem 5 is equivalent to Theorem 1.

■

## VI. PROOF OF THEOREM 2

Like most proofs of the speed of polarization, our proof of Theorem 2 relies on the following result by Arıkan and Telatar [11], although we need the more general form of the result given[4] in [3, Lemma 2.3].

**Lemma 14** ([11],[3]). *If $Z_n$ converges almost surely to a $\{0,1\}$-valued random variable $Z_\infty$ and if there exists $K < \infty$ such that*

$$Z_n \le K Z_{n-1}, \quad \text{if } B_n = 0 \tag{28}$$

$$Z_n \le K Z_{n-1}^2, \quad \text{if } B_n = 1 \tag{29}$$

*then*

$$\lim_{n\to\infty} \mathbb{P}\left(Z_n < 2^{-2^{n\beta}}\right) = \mathbb{P}(Z_\infty = 0)$$

*for all $\beta < 1/2$.*

Recall from the proof of Theorem 1 that $H_n$ converges almost surely to a $\{0,1\}$-valued random variable. It then follows from the relations [13, Proposition 2]

$$Z(A|B)^2 \le H(A|B)$$

$$H(A|B) \le \log(1 + Z(A|B))$$

that $Z_n$ also converges almost surely to a $\{0,1\}$-valued random variable $Z_\infty$. Indeed, $H_n \to 0$ implies $Z_n \to 0$ whereas $H_n \to 1$ implies $Z_n \to 1$. It then suffices to show that $Z_n$ satisfies inequalities (28) and (29).

We claim that this is indeed the case with $K = 2\psi_0$. To see this, let $\hat{X}_1^{2N}, \hat{Y}_1^{2N}$ be distributed as $P_{X_1^N Y_1^N} \cdot P_{X_{N+1}^{2N} Y_{N+1}^{2N}}$, and define the corresponding variables $\hat{U}_i, \hat{V}_i, \hat{Q}_i, \hat{R}_i$ as in (3). We know from [1, Proposition 5] that

$$Z(\hat{U}_i + \hat{V}_i | \hat{Q}_i, \hat{R}_i) \le 2 Z(\hat{U}_i | \hat{Q}_i), \tag{30}$$

$$Z(\hat{V}_i | \hat{Q}_i, \hat{R}_i, \hat{U}_i + \hat{V}_i) \le Z(\hat{U}_i | \hat{Q}_i)^2. \tag{31}$$

Now let $(A, B)$ and $(\hat{A}, \hat{B})$ be random variables that can be written as

$$(A, B) = f(X_1^{2N}, Y_1^{2N})$$

$$(\hat{A}, \hat{B}) = f(\hat{X}_1^{2N}, \hat{Y}_1^{2N})$$

for some function $f$. Observe that the assumption (1) implies $p_{A,B} \le \psi_0 \cdot p_{\hat{A},\hat{B}}$. Therefore, for binary $A$ we have

$$Z(A|B) = 2 \sum_b \sqrt{p_{A,B}(0,b) p_{A,B}(1,b)}$$

$$\le 2\psi_0 \sum_b \sqrt{p_{\hat{A},\hat{B}}(0,b) p_{\hat{A},\hat{B}}(1,b)}$$

$$= \psi_0 \cdot Z(\hat{A}|\hat{B}). \tag{32}$$

[4]See also [12] for a simpler proof.

Defining $A = U_i + V_i$ and $B = (Q_i, R_i)$ and combining (32) with (30) implies (28) with $K = 2\psi_0$. Similarly, defining $A = V_i$ and $B = (Q_i, R_i, U_i + V_i)$ and combining (32) with (31) implies (29) with $K = \psi_0$. This proves Theorem 2 since $\psi_0 < \infty$ by assumption.

## VII. PROOF OF THEOREM 4

Recall that the process we are considering is described in Figure 1. Let us start by defining the process exactly. The state of the process at time $t = 1, 2, \ldots$ is denoted $S_t$. Each such state has 4 possible values, $\{0, 1, 2, 3\}$. The initial state $S_1$ is picked uniformly at random. The value of $S_1$ determines the value of all $S_t$, specifically, $S_t = S_1 + t - 1 \pmod 4$. If $S_t \in \{0, 1\}$, then $X_t$, the output of the process at time $t$, is picked uniformly at random from $\{0, 1\}$. If $S_t \in \{2, 3\}$, then $X_t$ equals 0. Recall that for a given $N$, we have $U_1^N = X_1^N \mathsf{B}_N \mathsf{G}_N$.

The proof of Theorem 4 is divided into two parts. In the first part, we consider $H(U_i | U_1^{i-1}, S_1 = s_1)$. Namely, we consider a setting related to, yet distinct from, that of Theorem 4: we assume that the initial state $S_1$ is known to equal the fixed value $s_1$. As we will see, the case $N = 8$ is of particular importance. We refer the reader to Table II, which highlights key features of the distribution of $U_1^6$ when $N = 8$, for the 4 possible values of $s_1$. The entry "$U_6 \perp U_1^5$" denotes that $U_6$ is independent of $U_1^5$. The correctness of the Table II is easy to validate by using Table I.

**Lemma 15.** *Consider the stationary Markov process described in Figure 1. Then, for $N \ge 8$, the following holds.*

*For all $\dfrac{5N}{8} < i \le \dfrac{6N}{8}$ we have that*

$$H(U_i | U_1^{i-1}, S_1 = s_1) = \begin{cases} 0, & \text{if } s_1 \in \{1, 3\}, \\ 1, & \text{if } s_1 \in \{0, 2\}. \end{cases}$$

*Proof:* The correctness of the lemma is straightforward to validate for $N = 8$. Indeed, for $N = 8$ we must only consider $i = 6$, and the result follows from the last column of Table II. Namely, for $s_1 \in \{1, 3\}$ we have that $U_6$ is a function of $U_1^5$; for $s_1 \in \{0, 2\}$ we have that $U_6$ is independent of $U_1^5$ and is distributed Ber(1/2).

The general result is proved by induction on $N$. We have proved the basis $N = 8$ above. In order to prove the step, let us first tailor the notation (3) to our needs:

$$U_1^N = X_1^N \mathsf{B}_N \mathsf{G}_N, \tag{33a}$$

$$V_1^N = X_{N+1}^{2N} \mathsf{B}_N \mathsf{G}_N, \tag{33b}$$

$$Q_i = U_1^{i-1}, \tag{33c}$$

$$R_i = V_1^{i-1}. \tag{33d}$$

Proving the step is equivalent to proving that for all indices $\frac{5N}{8} < i \le \frac{6N}{8}$,

$$H(U_i + V_i | Q_i, R_i, S_1 = s_1)$$

$$= H(V_i | U_i + V_i, Q_i, R_i, S_1 = s_1)$$

$$= H(U_i | Q_i, S_1 = s_1). \tag{34}$$

Recall that $N$ is a power of 2 and $N \ge 8$. Thus, $N$ is a multiple of 4. Since the period of the process is 4, we have that

$$U_1 = X_1 + X_2 + X_3 + X_4 + X_5 + X_6 + X_7 + X_8$$
$$U_2 = \qquad\qquad\qquad\qquad\quad X_5 + X_6 + X_7 + X_8$$
$$U_3 = \qquad\quad X_3 + X_4 \qquad\qquad\quad X_7 + X_8$$
$$U_4 = \qquad\qquad\qquad\qquad\qquad\qquad\quad X_7 + X_8$$
$$U_5 = \quad X_2 \qquad + X_4 \qquad\quad + X_6 \qquad\quad + X_8$$
$$U_6 = \qquad\qquad\qquad\qquad\qquad\quad X_6 \qquad\quad + X_8$$

| | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ | $X_7$ | $X_8$ |
|---|---|---|---|---|---|---|---|---|
| $S_1 = 0$ | $B$ | $B$ | 0 | 0 | $B$ | $B$ | 0 | 0 |
| $S_1 = 1$ | $B$ | 0 | 0 | $B$ | $B$ | 0 | 0 | $B$ |
| $S_1 = 2$ | 0 | 0 | $B$ | $B$ | 0 | 0 | $B$ | $B$ |
| $S_1 = 3$ | 0 | $B$ | $B$ | 0 | 0 | $B$ | $B$ | 0 |

$S_1 = s_1$ iff $S_{N+1} = s_1$. Moreover, it is easily seen that given that $S_1 = s_1$, $(U_i, Q_i)$ and $(V_i, R_i)$ are identically distributed. Hence,

$$H(U_i|Q_i, S_1 = s_1)$$
$$= H(V_i|R_i, S_{N+1} = s_1) = H(V_i|R_i, S_1 = s_1) . \qquad (35)$$

Moreover, it is easily seen that given that $S_1 = s_1$, $(U_i, Q_i)$ and $(V_i, R_i)$ are independent.

We now prove (34) for the two cases of interest. Indeed, if $H(U_i|Q_i, S_1 = s_1) = 0$ then $U_i$ and $V_i$ are deterministic function of $Q_i$ and $R_i$, respectively, given that $S_1 = s_1$. Hence, the two equalities in (34) follow easily. If $H(U_i|Q_i, S_1 = s_1) = 1$, then by (35) and the independence of $(U_i, Q_i)$ and $(V_i, R_i)$ given $S_1 = s_1$ we deduce that

$$2 = H(U_i|Q_i, R_i, S_1 = s_1) + H(V_i|U_i, Q_i, R_i, S_1 = s_1)$$
$$= H(U_i, V_i|Q_i, R_i, S_1 = s_1)$$
$$= H(U_i + V_i, V_i|Q_i, R_i, S_1 = s_1)$$
$$= H(U_i + V_i|Q_i, R_i, S_1 = s_1)$$
$$\quad + H(V_i|U_i + V_i, Q_i, R_i, S_1 = s_1) .$$

Since the two terms on the RHS are at most 1, they must both equal 1, proving (34) for this case as well. ∎

An immediate corollary of Lemma 15 is that $H(U_i|U_1^{i-1}, S_1) = 1/2$, for $\frac{5N}{8} < i \leq \frac{6N}{8}$. To see this, note that all 4 states are equally likely as initial states. What remains is to prove that $S_1$ is essentially known from $U_1^{i-1}$.

**Lemma 16.** *Consider the stationary Markov process depicted in Figure 1. Then, there exists an $\epsilon_N$ such that*

*for all* $\dfrac{5N}{8} < i \leq \dfrac{6N}{8}$ *we have that*

$$H(S_1|U_1^{i-1}) \leq \epsilon_N , \quad \text{and} \quad \lim_{N \to \infty} \epsilon_N = 0 . \quad (36)$$

*Proof:* We start by giving an informal explanation as to why the claim holds. Consider the first two columns of Table II, and suppose we had many i.i.d. realizations of $U_1^5$, all with the same initial state $s_1$. Hence, the first column would allow us to distinguish — with very high probability — between $s_1 = 0$, $s_1 = 2$, and $s_1 \in \{1, 3\}$:

- If $s_1 = 0$ then all the realizations of $U_4$ would equal 0.

| | $(U_2, U_4)$ | $(U_1, U_3, U_5)$ | $U_6$ vs. $U_1^5$ |
|---|---|---|---|
| $S_1 = 0$ | $U_4 = 0$ | | $U_6 \perp U_1^5$ |
| $S_1 = 1$ | i.i.d. | $U_5 = U_3$ | $U_6 = U_4$ |
| $S_1 = 2$ | $U_4 = U_2$ | | $U_6 \perp U_1^5$ |
| $S_1 = 3$ | i.i.d. | $U_5 = U_3 + U_1$ | $U_6 = U_4 + U_2$ |

- If $s_1 = 2$, all realizations would satisfy $U_2 = U_4$. In roughly half the realizations we would have $U_4 = 1$, since $U_4 \sim \mathrm{Ber}(1/2)$. Each such realization would rule out the previous case.
- If $s_1 \in \{1, 3\}$ then in roughly a quarter of the realizations we would have $U_4 = 1$ and $U_2 = 0$, since $U_2$ and $U_4$ are i.i.d. and $\mathrm{Ber}(1/2)$. Such an outcome would distinguishing this case from the two previous ones.

To distinguish between $s_1 = 1$ and $s_1 = 3$, we utilize the second column of Table II. Specifically, in both cases, $U_1 \sim \mathrm{Ber}(1/2)$. Thus, in roughly half of the realizations, $U_1 = 1$, and for each such realization we can distinguish between $s_1 = 1$ in which $U_5 = U_3$ and $s_1 = 3$ in which $U_5 \neq U_3$.

Lastly, we claim that such independent realization of $U_1^5$ can indeed be attained. Specifically, for $N \geq 8$ and $\frac{5N}{8} < i \leq \frac{6N}{8}$, the vector $U_1^{i-1}$ can be used to deduce the first 5 entries of each vector in the set $\{X_{1+8(j-1)}^{1+8j} \mathsf{B}_8 \mathsf{G}_8 : 1 \leq j \leq N/8\}$. Note that since the period of the process is 4, the state at time $1 + 8(j - 1)$ is equal to $s_1$, for all values of $j$. Also, given $S_1$, all the vectors in the above set are independent.

Let us move on to the formal proof. The statistical properties of $U_1^5$ detailed above are easy to validate using Table I. Suppose we have $N/8$ realizations of $U_1^5$, which are i.i.d. given $S_1$. The above description suggests an algorithm for guessing the value of $S_1$:

- If all the realizations of $U_4$ equal 0, set $\hat{S}_1 = 0$.
- Otherwise, if all realizations satisfy $U_2 = U_4$, set $\hat{S}_1 = 2$.
- Otherwise, if all realizations satisfy $U_5 = U_3$, set $\hat{S}_1 = 1$.
- Otherwise, set $\hat{S}_1 = 3$.

A straightforward calculation shows that the probability of misdecoding $S_1$ goes down to 0 exponentially in $N$. By Fano's inequality [7, Theorem 2.10.1], we have that

$$H(S_1|U_1^{i-1}) \leq h_2(p_e) + p_e \log_2 4 ,$$

where $p_e$ is the probability of misdecoding. Since $p_e$ tends to 0, the RHS of the above tends to 0 as well.

Recall the set $\{X_{1+8(j-1)}^{1+8j} \mathsf{B}_8 \mathsf{G}_8 : 1 \leq j \leq N/8\}$, and denote by $A$ the vectors obtained by taking the prefix of length 5 of each vector in the set. Obviously, the vectors in $A$ are i.i.d. given $S_1$, and have the same distribution as the $U_1^5$ discussed above. All that remains to prove is that we can deduce $A$ from $U_1^{i-1}$, when $\frac{5N}{8} < i \leq \frac{6N}{8}$. We prove this by induction on $N$. The case $N = 8$ is immediate. For the step, let the set $B$ be defined similarly to $A$, but with $j$ ranging as $N/8 + 1 \leq j \leq N/8 + N/8$. The induction step assumes that $A$ can be deduced from $U_1^{i-1}$. Hence, $B$ can be deduced from $V_1^{i-1}$, where we recall the shorthand (33). Recalling the definition of the polar

transform, we must prove that both $A$ and $B$ can be deduced from either $(U_1^{i-1} + V_1^{i-1}, V_1^{i-1})$ or $(U_1^{i-1} + V_1^{i-1}, V_1^{i-1}, U_i + V_i)$. Obviously, this is true. ∎

The proof of Theorem 4 is now a simple consequence of the above.

*Proof of Theorem 4:* By the chain rule applied in two ways to $H(U_i, S_1|U_1^{i-1})$ we deduce that

$$H(U_i|U_1^{i-1}) + H(S_1|U_i, U_1^{i-1}) = H(S_1|U_1^{i-1}) + H(U_i|U_1^{i-1}, S_1) .$$

As discussed, an immediate consequence of Lemma 15 is that $H(U_i|U_1^{i-1}, S_1) = 1/2$. Thus,

$$\left| H(U_i|U_1^{i-1}) - 1/2 \right| = \left| H(S_1|U_1^{i-1}) - H(S_1|U_i, U_1^{i-1}) \right| .$$

By Lemma 16, there exists an $\epsilon_N \to 0$ such that

$$0 \le H(S_1|U_i, U_1^{i-1}) \le H(S_1|U_1^{i-1}) \le \epsilon_N .$$

Hence,

$$\left| H(U_i|U_1^{i-1}) - 1/2 \right| \le \epsilon_N . \quad ∎$$

## VIII. Appendix

*Proof of Corollary 9:* By marginalizing (8) over $\tilde{v}_i$ we deduce that

$$p_{\tilde{U}_i|Q_i,R_i}(\tilde{u}_i|q_i, r_i) = p_{\tilde{U}_i|Q_i}(\tilde{u}_i|q_i) . \quad (37)$$

Similarly,

$$p_{\tilde{V}_i|Q_i,R_i}(\tilde{v}_i|q_i, r_i) = p_{\tilde{V}_i|R_i}(\tilde{v}_i|r_i) . \quad (38)$$

Thus, by (8) and the above we deduce that $\tilde{U}_i$ and $\tilde{V}_i$ are independent given $Q_i$ and $R_i$,

$$\begin{aligned} p_{\tilde{U}_i,\tilde{V}_i|Q_i,R_i}&(\tilde{u}_i, \tilde{v}_i|q_i, r_i) \\ &= p_{\tilde{U}_i|Q_i,R_i}(\tilde{u}_i|q_i, r_i) \cdot p_{\tilde{V}_i|Q_i,R_i}(\tilde{v}_i|q_i, r_i) . \end{aligned} \quad (39)$$

Define

$$h_2(\alpha) = -\alpha \log_2 \alpha - (1 - \alpha) \log_2(1 - \alpha) . \quad (40)$$

We start with the following simple claim: for $\alpha, \beta$ between 0 and 1,

$$|h_2(\beta) - h_2(\alpha)| \le h_2(|\beta - \alpha|) . \quad (41)$$

Indeed, assume w.l.o.g. that $\beta \ge \alpha$. Then,

$$\begin{aligned} h_2(\beta) - h_2(\alpha) &= \int_\alpha^\beta h_2'(t)\, dt \\ &\le \int_0^{\beta-\alpha} h_2'(t)\, dt = h_2(\beta - \alpha) , \end{aligned} \quad (42)$$

where the inequality follows from the concavity of $h_2$ (the derivative $h_2'$ is decreasing). Similarly,

$$\begin{aligned} h_2(\beta) - h_2(\alpha) &= \int_\alpha^\beta h_2'(t)\, dt \\ &\ge \int_{1-(\beta-\alpha)}^1 h_2'(t)\, dt \\ &= -h_2(1 - (\beta - \alpha)) = -h_2(\beta - \alpha) . \end{aligned} \quad (43)$$

We deduce (41) from (42) and (43).

For $q_i$ and $r_i$ fixed, let us adopt the shorthand $\alpha = p_{U_i+V_i|Q_i,R_i}(0|q_i, r_i)$ and $\beta = p_{\tilde{U}_i+\tilde{V}_i|Q_i,R_i}(0|q_i, r_i)$. We claim that

$$\begin{aligned} |H(\tilde{U}_i &+ \tilde{V}_i|Q_i, R_i) - H(U_i + V_i|Q_i, R_i)| \\ &= \left| \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i)\big(h_2(\beta) - h_2(\alpha)\big) \right| \\ &\le \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i)|h_2(\beta) - h_2(\alpha)| \\ &\le \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i) h_2(|\beta - \alpha|) \\ &\le h_2\left( \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i)|\beta - \alpha| \right) . \end{aligned} \quad (44)$$

The second inequality follows form (41) while the third inequality follows by applying Jensen's inequality [7, Theorem 2.6.2] with respect to the concave function $h_2$.

Our aim now is to bound the argument of $h_2$ in the RHS of the above displayed equation. To this end, we use the shorthand $p = p_{U_i,V_i|Q_i,R_i}$ and $\tilde{p} = p_{\tilde{U}_i,\tilde{V}_i|Q_i,R_i}$. By (39),

$$I(U_i; V_i|Q_i, R_i) = \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i) D(p||\tilde{p}) ,$$

where $D(p||\tilde{p})$ is the relative entropy between $p$ and $\tilde{p}$, for $q_i$ and $r_i$ fixed,

$$D(p||\tilde{p}) = \sum_{u_i, v_i} p(u_i, v_i|q_i, r_i) \log_2 \frac{p(u_i, v_i|q_i, r_i)}{\tilde{p}(u_i, v_i|q_i, r_i)} .$$

Next, let us denote $p_+ = p_{U_i+V_i|Q_i,R_i}$ and $\tilde{p}_+ = p_{\tilde{U}_i+\tilde{V}_i|Q_i,R_i}$. Obviously, $p_+$ is gotten by quantizing $p$:

$$\begin{aligned} p_+(0|q_i, r_i) &= p(0, 0|q_i, r_i) + p(1, 1|q_i, r_i) , \\ p_+(1|q_i, r_i) &= p(1, 0|q_i, r_i) + p(0, 1|q_i, r_i) . \end{aligned}$$

The same quantization is used to derive $\tilde{p}_+$ from $\tilde{p}$. A simple consequence of the log-sum inequality [7, Theorem 2.7.1] is that such a quantization reduces the relative entropy. Namely, for $q_i, r_i$ fixed,

$$D(p||\tilde{p}) \ge D(p_+||\tilde{p}_+) .$$

Recalling that $\alpha = p_+(0|q_i, r_i)$ and $\beta = \tilde{p}_+(0|q_i, r_i)$, we get from Pinsker's inequality [7, Equation 11.147] that

$$D(p_+||\tilde{p}_+) \ge \frac{1}{2 \ln 2} \cdot 2(\beta - \alpha)^2 .$$

Aggregating the above inequalities yields

$$I(U_i; V_i|Q_i, R_i) \ge \frac{1}{\ln 2} \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i) \cdot (\beta - \alpha)^2 .$$

Now is the time to invoke Lemma 8. Namely, for an $\epsilon'$ which we will determine shortly, the fraction of indices $i$ for which $I(U_i; V_i|Q_i, R_i) \le \epsilon'$ approaches 1 as $N \to \infty$. Thus, for such an index $i$ we have that

$$\frac{1}{\ln 2} \sum_{q_i, r_i} p_{Q_i,R_i}(q_i, r_i) \cdot |\beta - \alpha|^2 \le I(U_i; V_i|Q_i, R_i) \le \epsilon' .$$

Since squaring is a convex function, we apply Jensen's inequality and deduce that

$$\sum_{q_i, r_i} p_{Q_i, R_i}(q_i, r_i) \cdot |\beta - \alpha| \leq \sqrt{\epsilon' \cdot \ln 2} \; .$$

Assuming the RHS of the above is less than $1/2$, we deduce from the above, the monotonicity of $h_2$ in $[0, 1/2]$, and (44) that

$$|H(\tilde{U}_i + \tilde{V}_i | Q_i, R_i) - H(U_i + V_i | Q_i, R_i)| \leq h_2(\sqrt{\epsilon' \cdot \ln 2}) \; .$$

Thus, taking $\epsilon'$ small enough so that $\sqrt{\epsilon' \cdot \ln 2} \leq 1/2$ and $h_2(\sqrt{\epsilon' \cdot \ln 2}) \leq \epsilon$ finishes the proof. ∎

*Proof of Lemma 11:* Denote the distributions of $A$ and $B$ as

$$A \sim \mathrm{Ber}(\alpha) \, , \quad B \sim \mathrm{Ber}(\beta) \; .$$

We will assume w.l.o.g. that $0 \leq \alpha \leq \beta \leq 1/2$ holds. Thus, according to our assumptions,

$$h_2(\alpha) \leq h_2(\beta) \, , \quad h_2(\alpha) \leq 1 - \xi \, , \quad h_2(\beta) \geq \xi \, ,$$

where $h_2$ is defined in (40). Since $h_2$ is strictly increasing when restricted to the domain $[0, 1/2]$, it is invertible and we conclude that

$$0 \leq \alpha \leq h_2^{-1}(1 - \xi) \, , \quad h_2^{-1}(\xi) \leq \beta \leq \frac{1}{2} \; .$$

We simplify the above to

$$0 \leq \alpha \leq \frac{1}{2} - \sigma \, , \quad \sigma \leq \beta \leq \frac{1}{2} \; . \tag{45}$$

where

$$\sigma = \sigma(\xi) = \min\left\{ h_2^{-1}(\xi), \frac{1}{2} - h_2^{-1}(1 - \xi) \right\} \; .$$

Define the random variable $D = (C, T)$ as follows,

$$D = (C, T) \, , \quad T \sim \mathrm{Ber}(1/2) \, , \quad C = \begin{cases} A & \text{if } T = 0 \, , \\ B & \text{if } T = 1 \, . \end{cases}$$

One easily gets that

$$H(A + B | D) = \frac{H(A) + H(B)}{2} \; .$$

Thus, we are interested in bounding the difference

$$H(A + B) - H(A + B | D) = I(A + B; D) \; .$$

We write $I(A + B; D)$ in terms of relative entropy [7, Equation (2.29)], and lower bound that with Pinsker's inequality [7, Equation 11.147]. Doing so results in a straightforward calculation which yields

$$H(A + B) - \frac{H(A) + H(B)}{2}$$
$$\geq \frac{2}{\ln 2} \big( \beta(1 - \beta)|1 - 2\alpha| + \alpha(1 - \alpha)|1 - 2\beta| \big)^2$$
$$\geq \frac{2}{\ln 2} \big( \beta(1 - \beta)|1 - 2\alpha| \big)^2$$
$$\geq \frac{2}{\ln 2} \big( \sigma(1 - \sigma)|2\sigma| \big)^2$$
$$= \frac{8}{\ln 2} \sigma^4 (1 - \sigma)^2 \, ,$$

where the last inequality follows from (45). Now, simply take $\Delta(\xi)$ as the RHS of the above. ∎

## REFERENCES

[1] E. Arıkan, "Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels," *IEEE Trans. Inform. Theory*, vol. 55, no. 7, pp. 3051–3073, July 2009.

[2] R. Wang, J. Honda, H. Yamamoto, R. Liu, and Y. Hou, "Construction of polar codes for channels with memory," in *Proc. IEEE Inform. Theory Workshop (ITW'2015)*, Jeju Island, Korea, 2015, pp. 187–191.

[3] E. Şaşoğlu, "Polar coding theorems for discrete systems," Ph.D. dissertation, Ecole Polytechnique Fédérale de Lausanne, 2011.

[4] B. Shuval and I. Tal, "Fast polarization for processes with memory," *IEEE Trans. Inform. Theory*, vol. 65, no. 4, pp. 2004–2020, April 2019.

[5] P. C. Shields, *The Ergodic Theory of Discrete Sample Paths*, ser. Graduate Studies in Mathematics. Providence (R.I.): American Mathematical Society, 1996, vol. 13.

[6] R. C. Bradley, *Introduction to Strong Mixing Conditions*. Heber City, Utah: Kendrick Press, 2007, vol. I.

[7] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Wiley, 2006.

[8] R. C. Bradley, "Equivalent mixing conditions for Markov chains," *Statis. Probab. Letters*, vol. 41, pp. 97–99, 1999.

[9] K. L. Chung, *A Course in Probability Theory*, 3rd ed. San Diego: Academic Press, 2001.

[10] E. Şaşoğlu, "Polarization and polar codes," in *Found. and Trends in Commun. and Inform. Theory*, vol. 8, no. 4, 2012, pp. 259–381.

[11] E. Arıkan and E. Telatar, "On the rate of channel polarization," in *Proc. IEEE Int'l Symp. Inform. Theory (ISIT'2009)*, Seoul, South Korea, 2009, pp. 1493–1495.

[12] I. Tal, "A simple proof of fast polarization," *IEEE Trans. Inform. Theory*, vol. 63, no. 12, pp. 7617–7619, December 2017.

[13] E. Arıkan, "Source polarization," in *Proc. IEEE Int'l Symp. Inform. Theory (ISIT'2010)*, Austin, Texas, 2010, pp. 899–903.

**Eren Şaşoğlu** received the B.Sc. degree in electrical engineering from Boğaziçi University in 2005, and the Ph.D. degree in communications systems from EPFL in 2011. He was a postdoctoral scholar at the University of California at San Diego and later at the University of California at Berkeley, an academic visitor at Technion, and a research scientist at Intel. He is now at Apple. He received the Best Doctoral Thesis Award at EPFL, and the STOC Best Paper Award in 2016.

**Ido Tal** (S'05–M'08–SM'18) was born in Haifa, Israel, in 1975. He received the B.Sc., M.Sc., and Ph.D. degrees in computer science from Technion — Israel Institute of Technology, Haifa, Israel, in 1998, 2003 and 2009, respectively. During 2010–2012 he was a postdoctoral scholar at the University of California at San Diego. In 2012 he joined the Electrical Engineering Department at Technion. His research interests include constrained coding and error-control coding. He received the IEEE Joint Communications Society/Information Theory Society Paper Award (jointly with Alexander Vardy) for the year 2017.