

Fast Polarization for Processes with Memory

Boaz Shuval, Ido Tal

Department of Electrical Engineering,
Technion, Haifa 32000, Israel.

Email: {bshuval@campus, idotal@ee}.technion.ac.il

Abstract—Fast polarization is crucial for the performance guarantees of polar codes. In the memoryless setting, the rate of polarization is known to be exponential in the square root of the block length. A complete characterization of the rate of polarization for models with memory has been missing. We consider polar codes for processes with memory that are characterized by an underlying aperiodic and irreducible finite state Markov chain. We show that the rate of polarization for these processes is the same as in the memoryless setting, both to the high and to the low-entropy sets. Thus, polar codes achieve the Markov capacity in many information-theoretic applications.

I. INTRODUCTION

Memory is prevalent in many communication scenarios. In this research we show that polar codes can be used directly for a large class of scenarios with memory. This enables leveraging the attractive properties of polar codes — such as low complexity encoding and decoding, vanishing error performance, and versatility — to scenarios with memory.

Polar codes [1] were first developed for binary-input, symmetric, memoryless, channels. They employ *successive cancellation* (SC) decoding, which consists of N successive decoding operations. The polarization phenomenon is that for large enough N , the decoding operations polarize to two sets: a ‘low-entropy’ set and a ‘high-entropy’ set. The vanishing error performance of polar codes is due to polarization happening sufficiently fast. Fast polarization to the low-entropy set for the memoryless setting was established in [1], [2].

Polar codes were extended to many other memoryless scenarios, e.g., non-binary channels [3], [4], source coding [5], [6], asymmetric channels and sources [7], and more. Many of these applications are contingent upon fast polarization to the high-entropy set; for memoryless settings, this was established in [5] (see also [8] for a different proof that is closer to ours).

The study of polar codes for scenarios with memory began in [4, Chapter 5], which showed that polarization occurs for a certain class of processes with memory. In [9], polarization was established for a more general class of processes with memory. That paper further showed that polarization to the low-entropy set is fast even for processes with memory. Fast polarization to the high-entropy set was not addressed.

A practical decoding algorithm for polar codes for finite-state channels was suggested in [10]. This algorithm is an extension of SC decoding, taking into account an underlying state structure. Its increase in complexity relative to the complexity of SC decoding is polynomial with the number of states. The authors also showed [10, Theorem 3] that their elegant scheme

from [7] can be applied to models with memory. To this end, they required the additional (then unproved) assumption of fast polarization both to the low and high-entropy sets.

This paper completes the picture. We show that for a large class of processes with memory, polarization is fast both to the low-entropy and high-entropy sets. Fast polarization to the low-entropy set will follow from a specialization of [9]. Fast polarization to the high-entropy set, Theorem 10, is the main result of this paper. Consequently, polar codes can be used in settings with memory with vanishing error probability.

Specifically, we consider stationary processes whose memory can be encompassed by an underlying (hidden) aperiodic and irreducible finite-state Markov chain. This family of processes includes, as special cases, finite state Markov channels with an ergodic state sequence, discrete ergodic sources with finite memory, and many input-constrained systems (e.g., (d, k) -runlength limited (RLL) constraint, with and without noise).

Due to space constraints, some proofs are omitted/shortened. The full version of this paper [11] contains detailed proofs.

II. PRELIMINARIES

A. Distribution Parameters

Let random variables (U, Q) have joint distribution $P_{U,Q}(u, q) = P_Q(q)P_{U|Q}(u|q)$. For simplicity, we assume that U is binary. The random variable Q is some observation dependent on U that takes values in a finite alphabet \mathcal{Q} .

In the following equations, we define the *Bhattacharyya parameter* of U given Q , $\mathcal{Z}(U|Q)$; the *total variation distance* of U given Q , $\mathcal{K}(U|Q)$; and the *conditional entropy* of U given Q , $\mathcal{H}(U|Q)$.

$$\mathcal{Z}(U|Q) = \sum_q 2\sqrt{P_{U,Q}(0, q)P_{U,Q}(1, q)}, \quad (1)$$

$$\mathcal{K}(U|Q) = \sum_q |P_{U,Q}(0, q) - P_{U,Q}(1, q)|, \quad (2)$$

$$\mathcal{H}(U|Q) = \sum_q P_Q(q) \sum_u P_{U|Q}(u|q) \log_2 P_{U|Q}(u|q). \quad (3)$$

The three parameters take values in $[0, 1]$.

Lemma 1. *We have*

$$\begin{aligned} \mathcal{K}(U|Q) &\geq \sqrt{1 - \mathcal{H}(U|Q)} \geq \sqrt{1 - \mathcal{Z}(U|Q)}, \\ \mathcal{K}(U|Q) &\leq \sqrt{1 - \mathcal{Z}(U|Q)^2} \leq \sqrt{1 - \mathcal{H}(U|Q)^2}. \end{aligned}$$

Informally, as a consequence of the above lemma,

$$\begin{aligned} \mathcal{Z}(U|Q) \approx 0 &\iff \mathcal{H}(U|Q) \approx 0 \iff \mathcal{K}(U|Q) \approx 1, \\ \mathcal{Z}(U|Q) \approx 1 &\iff \mathcal{H}(U|Q) \approx 1 \iff \mathcal{K}(U|Q) \approx 0. \end{aligned} \quad (4)$$

The definitions in Equations (1) to (3) naturally extend to the case where there are multiple random variables related to U . For example, consider a triplet of random variables (U, Q, S) with joint distribution $P_{U,Q,S}(u, q, s)$ such that U is binary and Q, S take values in finite alphabets \mathcal{Q}, \mathcal{S} . Then, $\mathcal{K}(U|Q, S) = \sum_{q,s} |P_{U,Q,S}(0, q, s) - P_{U,Q,S}(1, q, s)|$; the remaining parameters are similarly extended.

Lemma 2. *The triplet (U, Q, S) satisfies $\mathcal{K}(U|Q) \leq \mathcal{K}(U|Q, S)$; $\mathcal{Z}(U|Q) \geq \mathcal{Z}(U|Q, S)$; and $\mathcal{H}(U|Q) \geq \mathcal{H}(U|Q, S)$.*

B. Polar Construction

Let (X_j, Y_j) , $j = 1, 2, \dots$ be a stationary process, such that X_j are binary and $Y_j \in \mathcal{Y}$, where \mathcal{Y} is a finite alphabet. Random variables X_j are to be estimated from observations Y_j .

We denote Arıkan's polarization matrix by G_N , where $N = 2^n$. Following [9], we define for $i = 1, 2, \dots, N$:

$$U_1^N = X_1^N G_N, \quad (5a)$$

$$V_1^N = X_{N+1}^{2N} G_N, \quad (5b)$$

$$Q_i = (U_1^{i-1}, Y_1^N), \quad (5c)$$

$$R_i = (V_1^{i-1}, Y_{N+1}^{2N}). \quad (5d)$$

Note that these definitions apply to two consecutive blocks of length N ; this will be useful in the sequel. By (5), we can write

$$(U_i, Q_i) = f(X_1^N, Y_1^N), \quad (V_i, R_i) = f(X_{N+1}^{2N}, Y_{N+1}^{2N}),$$

where function f depends solely on i . Due to stationarity, $P_{U_i, Q_i} = P_{V_i, R_i}$. Denoting $T_i = U_i + V_i$, we obtain

$$P_{T_i, V_i, Q_i, R_i}(t, v, q, r) = P_{U_i, V_i, Q_i, R_i}(t + v, v, q, r). \quad (6)$$

Let B_1, B_2, \dots be a sequence of independent and identically distributed (i.i.d.) Bernoulli-1/2 random variables. Set $i = 1 + (B_1 B_2 \dots B_n)_2$, where $(B_1 B_2 \dots B_n)_2 = \sum_{j=1}^n B_j 2^{n-j}$. Thus, i is a random variable that assumes any value in $\{1, 2, \dots, N\}$ with equal probability. Define the random variables

$$K_n = \mathcal{K}(U_i | U_1^{i-1}, Y_1^N) = \mathcal{K}(U_i | Q_i),$$

$$Z_n = \mathcal{Z}(U_i | U_1^{i-1}, Y_1^N) = \mathcal{Z}(U_i | Q_i),$$

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

whenever $i = 1 + (B_1 B_2 \dots B_n)_2$. They denote the relevant distribution parameters for a uniformly chosen index after n polarization steps.

By the properties of G_N [1, Section VII],

$$K_{n+1} = \begin{cases} \mathcal{K}(U_i + V_i | Q_i, R_i), & \text{if } B_{n+1} = 0 \\ \mathcal{K}(V_i | U_i + V_i, Q_i, R_i), & \text{if } B_{n+1} = 1. \end{cases} \quad (7)$$

Similar relationships hold for H_{n+1} and Z_{n+1} .

C. Polarization

Polarization occurs when the fraction of indices with moderate conditional entropy $|\{i : \mathcal{H}(U_i | Q_i) \in (\epsilon, 1 - \epsilon)\}|/N$ vanishes for large enough n , for any $\epsilon > 0$.

Definition 1. Let A_n , $n = 1, 2, \dots$ be a sequence of random variables that take values in $[0, 1]$.

- 1) The sequence A_n *polarizes* if it converges almost surely to a $\{0, 1\}$ -random variable A_∞ as $n \rightarrow \infty$.
- 2) The sequence A_n *polarizes fast to 0* with $\beta > 0$ if it polarizes and $\lim_{n \rightarrow \infty} \mathbb{P}(A_n < 2^{-2^{n\beta}}) = \mathbb{P}(A_\infty = 0)$.
- 3) The sequence A_n *polarizes fast to 1* with $\beta > 0$ if it polarizes and $\lim_{n \rightarrow \infty} \mathbb{P}(A_n > 1 - 2^{-2^{n\beta}}) = \mathbb{P}(A_\infty = 1)$.

The following lemma [2], [4], is an important tool for establishing fast polarization (see also [12]).

Lemma 3. [2], [4] *Let B_n , $n = 1, 2, \dots$ be an i.i.d. Bernoulli-1/2 process and A_n , $n = 1, 2, \dots$ be a $[0, 1]$ -valued process that polarizes to A_∞ . Assume that there exist $k \geq 1$ and $d_0, d_1 > 0$ such that $A_{n+1} \leq kA_n^{d_1}$ if $B_{n+1} = 0$, for $i = 0, 1$. Then, for any $0 < \beta < E = (\log_2 d_0 + \log_2 d_1)/2$, we have*

$$\lim_{n \rightarrow \infty} \mathbb{P}(A_n < 2^{-2^{n\beta}}) = \mathbb{P}(A_\infty = 0). \quad (8)$$

Arıkan showed in [1] that in the memoryless case (i.e., $P_{X_1^N, Y_1^N}(x_1^N, y_1^N) = \prod_{j=1}^N P_{X, Y}(x_j, y_j)$) the process H_n polarizes. Fast polarization to the low-entropy set was established in [2] by showing that $Z_{n+1} \leq 2Z_n$ if $B_{n+1} = 0$ and $Z_{n+1} = Z_n^2$ if $B_{n+1} = 1$, and using Lemma 3 and Equation (4).

Fast polarization to the high-entropy set is important for many applications of polar codes. For example, it is integral to source coding applications [5] and to channel coding without symmetry assumptions [7]. Lemma 3 will be useful for establishing fast polarization results to the high entropy set.

III. FINITE-STATE APERIODIC IRREDUCIBLE MARKOV PROCESSES

We now introduce a class of processes with memory, described using a hidden state sequence. We call them *Finite-state Aperiodic Irreducible Markov processes* (FAIM processes).

Our model applies to many problems in information theory that can be described using states. Examples include compression of finite-memory sources and coding for input constrained channels. Additionally, our model may be applied to finite-state channels and channels with intersymbol interference; in this case, the FAIM state sequence describes both the channel state and input state. That is, FAIM processes enable us to model non-i.i.d. input sequences.

A. Definition

Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a stationary process, where X_j is binary, $Y_j \in \mathcal{Y}$, and $S_j \in \mathcal{S}$. Alphabets \mathcal{Y} and \mathcal{S} are finite, and $\mathcal{S} = \{1, 2, \dots, |\mathcal{S}|\}$. We call $S_j, j \in \mathbb{Z}$ the *state sequence*. It encompasses the memory of the process.

The process (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ is described by the conditional distribution $P_{X_j, Y_j, S_j | S_{j-1}}$, which, by stationarity, is independent of j . We call it a FAIM process if

$$P_{X_j, Y_j, S_j | S_{j-1}} = P_{X_j, Y_j, S_j | S_{-\infty}^{j-1}, X_{-\infty}^{j-1}, Y_{-\infty}^{j-1}},$$

and the state sequence S_j , $j \in \mathbb{Z}$ is a finite-state, homogeneous, aperiodic, and irreducible Markov chain. For any $N > M > 0$, with b denoting the value of the middle state S_M ,

$$\begin{aligned} P_{X_1^N, Y_1^N, S_N | S_0} &= \sum_b P_{X_1^M, Y_1^M, S_M, X_{M+1}^N, Y_{M+1}^N, S_N | S_0} \\ &= \sum_b P_{X_{M+1}^N, Y_{M+1}^N, S_N | S_M} \cdot P_{X_1^M, Y_1^M, S_M | S_0}. \end{aligned} \quad (9)$$

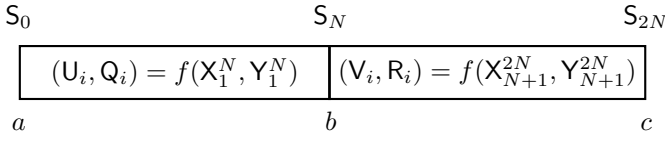


Fig. 1. Two adjacent length- N blocks of a FAIM process. Here, $S_0 = a$, $S_N = b$, and $S_{2N} = c$, where $a, b, c \in \mathcal{S}$.

We use the following shorthand for the stationary distribution of the state sequence: $\pi_N(a) = P_{S_N}(a)$, $\pi_{N|M}(b|a) = P_{S_N|S_M}(b|a)$, and $\pi_{N,M}(b, a) = P_{S_N, S_M}(b, a)$. By stationarity, $\pi_N(a) = \pi_0(a)$. By irreducibility and aperiodicity, $\pi_0(a) > 0$ for all $a \in \mathcal{S}$. Aperiodicity is assumed because periodic processes may not polarize [9, Theorem 3].

B. Blocks of a FAIM Process

Typically, the state sequence is not observed by the encoder or decoder. The joint distribution of (X_1^N, Y_1^N) is given by

$$P_{X_1^N, Y_1^N}(x_1^N, y_1^N) = \sum_{b,a} P_{X_1^N, Y_1^N, S_N|S_0}(x_1^N, y_1^N, b|a)\pi_0(a).$$

Definition 2 (Block). Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. We call (X_{L+1}^M, Y_{L+1}^M) a *block* of the FAIM process, with length $M - L$. State S_L is the *initial* state of the block and state S_M is the *final* state of the block.

We emphasize that for block (X_{L+1}^M, Y_{L+1}^M) , the initial state is S_L and *not* S_{L+1} .

The following lemma establishes that FAIM processes are a special case of the family of processes considered in [9].

Lemma 4. *If (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ is a FAIM process, there exists a non-increasing sequence $\psi(N)$ with $\psi(0) < \infty$ and $\psi(N) \rightarrow 1$ as $N \rightarrow \infty$, such that for any $N > M \geq L \geq 1$,*

$$P_{X_1^L, Y_1^L, X_{M+1}^N, Y_{M+1}^N} \leq \psi(M-L) \cdot P_{X_1^L, Y_1^L} \cdot P_{X_{M+1}^N, Y_{M+1}^N}. \quad (10)$$

To see this, define

$$\psi(N) = \begin{cases} \max_{a,b} \pi_{N|0}(b|a)/\pi_0(b), & \text{if } N > 0 \\ \max_a 1/\pi_0(a), & \text{if } N = 0, \end{cases} \quad (11)$$

and use (9). A process satisfying (10) with $\psi(N) \rightarrow 1$ as $N \rightarrow \infty$ is called *ψ -mixing*.

Two adjacent blocks share a state: the final state of the first block is the initial state of the second block. By (9), for any $N > M \geq 1$,

$$\begin{aligned} P_{X_1^M, Y_1^M, X_{M+1}^N, Y_{M+1}^N | S_0, S_M, S_N} \\ = P_{X_1^M, Y_1^M | S_0, S_M} P_{X_{M+1}^N, Y_{M+1}^N | S_M, S_N}. \end{aligned} \quad (12)$$

We will use ascending letters to denote values of ordered states. In Figure 1 we illustrate a useful case. A block of length $2N$ comprises two adjacent blocks of length N . State S_0 , the initial state of the first block, takes value a . State S_N , at the end of the first block and the beginning of the second block, takes value b . State S_{2N} , at the end of the second block, takes value c .

C. Boundary State-Informed Parameters for FAIM Processes

Let (X_1^N, Y_1^N) be a block of a FAIM process and let $(U, Q) = f(X_1^N, Y_1^N)$, where function $f(\cdot, \cdot)$ is independent of the state sequence and U is binary. We denote

$$P_a^b(u, q) \triangleq P_{U, Q | S_N, S_0}(u, q | b, a) = \frac{P_{U, Q, S_N | S_0}(u, q, b | a)}{\pi_{N|0}(b|a)}. \quad (13)$$

and further define $P_a^b(q) = P_{Q | S_N, S_0}(q | b, a)$. We denote by $\mathcal{Z}_a^b(U|Q)$, $\mathcal{K}_a^b(U|Q)$, and $\mathcal{H}_a^b(U|Q)$ the results of replacing $P_{U, Q}(u, q)$ with $P_a^b(u, q)$ in equations (1) to (3), respectively. For example, $\mathcal{K}_a^b(U|Q) = \sum_q |P_a^b(0, q) - P_a^b(1, q)|$.

As $P_{U, Q, S_N, S_0}(u, q, b, a) = P_a^b(u, q) \cdot \pi_{N,0}(b, a)$, we respectively define the *boundary state-informed* (BSI) total variation distance, Bhattacharyya parameter, and conditional entropy, as

$$\begin{aligned} \mathcal{K}(U|Q, S_N, S_0) &= \sum_{a,b} \pi_{N,0}(b, a) \mathcal{K}_a^b(U|Q), \\ \mathcal{Z}(U|Q, S_N, S_0) &= \sum_{a,b} \pi_{N,0}(b, a) \mathcal{Z}_a^b(U|Q), \\ \mathcal{H}(U|Q, S_N, S_0) &= \sum_{a,b} \pi_{N,0}(b, a) \mathcal{H}_a^b(U|Q). \end{aligned}$$

BSI parameters are defined for blocks of the process; they depend on the initial and final states of the block.

IV. FAST POLARIZATION FOR FAIM PROCESSES

We use the notation of Section II-B for FAIM processes. That is, U_1^N, V_1^N, Q_i, R_i , $i = 1, \dots, N$ are defined using (5). The random variables B_1, \dots, B_n are Bernoulli-1/2 and i.i.d., and $i = 1 + (B_1 B_2 \dots B_n)_2$. Via i , we define the random variables $K_n = \mathcal{K}(U_i | Q_i)$, $H_n = \mathcal{H}(U_i | Q_i)$, and $Z_n = \mathcal{Z}(U_i | Q_i)$.

Let \hat{K}_n, \hat{H}_n , and \hat{Z}_n denote the BSI versions of K_n, Z_n , and H_n , respectively. That is, with i chosen randomly as above,

$$\begin{aligned} \hat{K}_n &= \mathcal{K}(U_i | Q_i, S_N, S_0), \\ \hat{Z}_n &= \mathcal{Z}(U_i | Q_i, S_N, S_0), \\ \hat{H}_n &= \mathcal{H}(U_i | Q_i, S_N, S_0). \end{aligned} \quad (14)$$

By Lemma 2, $K_n \leq \hat{K}_n$, $Z_n \geq \hat{Z}_n$, and $H_n \geq \hat{H}_n$ for any n . Similar to (7), recalling that $2N = 2^{n+1}$, we have

$$\hat{K}_{n+1} = \begin{cases} \mathcal{K}(U_i + V_i | Q_i, R_i, S_0, S_{2N}), & \text{if } B_{n+1} = 0 \\ \mathcal{K}(V_i | U_i + V_i, Q_i, R_i, S_0, S_{2N}), & \text{if } B_{n+1} = 1. \end{cases} \quad (15)$$

Relationships akin to (15) hold for \hat{Z}_{n+1} and \hat{H}_{n+1} , with \mathcal{K} replaced with \mathcal{Z} and \mathcal{H} , respectively.

A. Polarization of FAIM Processes

Let

$$\mathcal{H}_*(X|Y) \triangleq \lim_{N \rightarrow \infty} \frac{1}{N} \mathcal{H}(X_1^N | Y_1^N).$$

This limit exists due to stationarity [13, Section 4.2] and the identity $\mathcal{H}(X_1^N | Y_1^N) = \mathcal{H}(X_1^N, Y_1^N) - \mathcal{H}(Y_1^N)$. In [9], the following was shown:

Theorem 5. [9] *For a stationary ψ -mixing process (X_j, Y_j) , $j \in \mathbb{Z}$, with $\psi(0) < \infty$:*

- 1) H_n polarizes to H_∞ with $\mathbb{P}(H_\infty = 1) = \mathcal{H}_*(X|Y)$;
- 2) Z_n polarizes fast to 0 with any $\beta < 1/2$.

Since FAIM processes are ψ -mixing, we obtain:

Corollary 6. *Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. Then,*

- 1) *Its conditional entropy process H_n polarizes to H_∞ with $\mathbb{P}(H_\infty = 1) = \mathcal{H}_*(X|Y)$.*
- 2) *Its Bhattacharyya process Z_n polarizes fast to 0 with any $\beta < 1/2$.*

Proof: By Lemma 4, FAIM processes are ψ -mixing and satisfy the requirements of Theorem 5. ■

B. Polarization of the BSI Distribution Parameters

This section is concerned with proving that the BSI distribution parameters polarize.

Theorem 7. *Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. The BSI conditional entropy process \hat{H}_n polarizes to \hat{H}_∞ and $\hat{H}_\infty = H_\infty$ almost surely.*

Proof: Consider two adjacent blocks of length $N = 2^n$ and let $i = 1 + (B_1 B_2 \cdots B_n)_2$, as in Figure 1. By (12),

$$P_{U_i, V_i | Q_i, R_i, S_0, S_N, S_{2N}} = P_{U_i | Q_i, S_0, S_N} P_{V_i | R_i, S_N, S_{2N}}. \quad (16)$$

Thus,

$$\begin{aligned} \hat{H}_n &\stackrel{(a)}{=} \frac{1}{2} \left(\mathcal{H}(U_i | Q_i, S_0, S_N) + \mathcal{H}(V_i | R_i, S_N, S_{2N}) \right) \\ &\stackrel{(b)}{=} \frac{1}{2} \mathcal{H}(U_i, V_i | Q_i, R_i, S_0, S_N, S_{2N}) \\ &\stackrel{(c)}{=} \frac{1}{2} \mathcal{H}(U_i + V_i, V_i | Q_i, R_i, S_0, S_N, S_{2N}) \\ &\stackrel{(d)}{=} \frac{1}{2} \left(\mathcal{H}(U_i + V_i | Q_i, R_i, S_0, S_N, S_{2N}) \right. \\ &\quad \left. + \mathcal{H}(V_i | U_i + V_i, Q_i, R_i, S_0, S_N, S_{2N}) \right) \\ &\stackrel{(e)}{\leq} \frac{1}{2} \left(\mathcal{H}(U_i + V_i | Q_i, R_i, S_0, S_{2N}) \right. \\ &\quad \left. + \mathcal{H}(V_i | U_i + V_i, Q_i, R_i, S_0, S_{2N}) \right) \end{aligned}$$

where (a) is by stationarity, (b) is by (16), (c) is because the mapping $(U, V) \mapsto (U + V, V)$ is one-to-one and onto, (d) is by the chain rule for entropies, and (e) is by Lemma 2.

Thus, recalling (15) (applied to the BSI conditional entropy), \hat{H}_n is a submartingale sequence. It is also bounded, as $\hat{H}_n \in [0, 1]$ for any n . Hence, it converges almost surely to some random variable $\hat{H}_\infty \in [0, 1]$.

Consequently, $\Delta H_n = H_n - \hat{H}_n$ converges almost surely to the random variable $\Delta H_\infty = H_\infty - \hat{H}_\infty$. By Lemma 2, $\Delta H_n \geq 0$ for any n , implying that $\Delta H_\infty \geq 0$ almost surely.

The effect of knowing a block's initial and final states becomes negligible for sufficiently long blocks; thus, it can be shown that $\lim_{n \rightarrow \infty} \mathbb{E}[\Delta H_n] = 0$. For details, see [11, Lemma 10]. Since ΔH_n is a non-negative sequence that converges to ΔH_∞ almost surely, by Fatou's lemma,

$$\begin{aligned} 0 &\leq \mathbb{E}[\Delta H_\infty] = \mathbb{E} \left[\liminf_{n \rightarrow \infty} \Delta H_n \right] \\ &\leq \liminf_{n \rightarrow \infty} \mathbb{E}[\Delta H_n] = \lim_{n \rightarrow \infty} \mathbb{E}[\Delta H_n] = 0. \end{aligned}$$

In other words, $\mathbb{E}[\Delta H_\infty] = 0$. By Markov's inequality, $\mathbb{P}(\Delta H_\infty \geq \delta) \leq \mathbb{E}[\Delta H_\infty] / \delta = 0$ for any $\delta > 0$; consequently, $\mathbb{P}(\Delta H_\infty = 0) = \mathbb{P}(H_\infty = \hat{H}_\infty) = 1$. That is, $\hat{H}_\infty = H_\infty$ almost surely. ■

Corollary 8.

- 1) *The sequences Z_n and \hat{Z}_n polarize to random variables Z_∞ and \hat{Z}_∞ , respectively. Moreover, $Z_\infty = \hat{Z}_\infty = H_\infty$ almost surely.*
- 2) *The sequences K_n and \hat{K}_n polarize to random variables K_∞ and \hat{K}_∞ , respectively. Moreover, $K_\infty = \hat{K}_\infty = 1 - H_\infty$ almost surely.*

This follows from Lemma 1, Corollary 6, and Theorem 7.

C. Fast Polarization to the High Entropy Set

In this section, we establish fast polarization to the high entropy set. We do this by proving that the total variation process K_n polarizes fast to 0, which implies that the Bhattacharyya process Z_n polarizes fast to 1.

Proposition 9. *Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. Then, with $\psi(0)$ as in (11),*

$$\hat{K}_{n+1} \leq \begin{cases} \psi(0) \hat{K}_n^2, & \text{if } B_{n+1} = 0 \\ 2\hat{K}_n, & \text{if } B_{n+1} = 1. \end{cases} \quad (17)$$

Proof of Proposition 9: Consider two adjacent blocks of length $N = 2^n$ and let $i = 1 + (B_1 B_2 \cdots B_n)_2$, as in Figure 1. By stationarity,

$$\hat{K}_n = \sum_{a,b \in \mathcal{S}} \pi_{N,0}(b, a) \mathcal{K}_a^b(U_i | Q_i) = \sum_{b,c \in \mathcal{S}} \pi_{2N,N}(c, b) \mathcal{K}_b^c(V_i | R_i). \quad (18)$$

As in (13), we denote $P_a^c(u, q) = P_{U_i, Q_i | S_N, S_0}(u, q | c, a) = P_{V_i, R_i | S_{2N}, S_N}(u, q | c, a)$, and $P_a^c(s) = P_a^c(0, s) + P_a^c(1, s)$; in particular, $\sum_s P_a^c(s) = 1$. Further denote

$$\mu(b) = \pi_{2N|N}(c|b) \pi_{N|0}(b|a) \pi_0(a) = \frac{\pi_{2N,N}(c, b) \pi_{N,0}(b, a)}{\pi_N(b)}.$$

For brevity, we omit the dependence of $\mu(b)$ on a, c . By (11),

$$\mu(b) \leq \psi(0) \cdot \pi_{2N,N}(c, b) \cdot \pi_{N,0}(b, a). \quad (19)$$

Also, since $\pi_N(b) = \sum_{a \in \mathcal{S}} \pi_{N,0}(b, a) = \sum_{c \in \mathcal{S}} \pi_{2N,N}(c, b)$,

$$\sum_a \mu(b) = \pi_{2N,N}(c, b), \quad \sum_c \mu(b) = \pi_{N,0}(b, a). \quad (20)$$

By (9) and (13),

$$\begin{aligned} \pi_{2N,0}(c, a) P_{U_i, V_i, Q_i, R_i | S_{2N}, S_0}(u, v, q, r | c, a) \\ &= \pi_0(a) \sum_b P_{U_i, Q_i, S_N | S_0}(u, q, b | a) P_{V_i, R_i, S_{2N} | S_N}(v, r, c | b) \\ &= \sum_b \mu(b) P_a^b(u, q) P_b^c(v, r). \end{aligned} \quad (21)$$

Set $T_i = U_i + V_i$. Using (7), a single-step polarization from \hat{K}_n to \hat{K}_{n+1} becomes

$$\hat{K}_{n+1} = \begin{cases} \sum_{a,c} \pi_{2N,0}(c, a) \mathcal{K}_a^c(T_i | Q_i, R_i), & \text{if } B_{n+1} = 0 \\ \sum_{a,c} \pi_{2N,0}(c, a) \mathcal{K}_a^c(V_i | T_i, Q_i, R_i), & \text{if } B_{n+1} = 1. \end{cases}$$

Here, $\mathcal{K}_a^c(\mathsf{T}_i|\mathsf{Q}_i, \mathsf{R}_i)$ and $\mathcal{K}_a^c(\mathsf{V}_i|\mathsf{T}_i, \mathsf{Q}_i, \mathsf{R}_i)$ are computed for a length- $2N$ block with initial state $\mathsf{S}_0 = a$ and final state $\mathsf{S}_{2N} = c$. The state at the middle of the block is $\mathsf{S}_N = b$. Denote by (6),

$$\begin{aligned}\bar{P}_a^c(t, v, q, r) &= P_{\mathsf{T}_i, \mathsf{V}_i, \mathsf{Q}_i, \mathsf{R}_i | \mathsf{S}_{2N}, \mathsf{S}_0}(t, v, q, r | c, a) \\ &= P_{\mathsf{U}_i, \mathsf{V}_i, \mathsf{Q}_i, \mathsf{R}_i | \mathsf{S}_{2N}, \mathsf{S}_0}(t + v, v, q, r | c, a)\end{aligned}$$

and $\bar{P}_a^c(t, q, r) = \sum_{v=0}^1 \bar{P}_a^c(t, v, q, r)$.

Consider first the case $\mathsf{B}_{n+1} = 0$:

$$\begin{aligned}\pi_{2N,0}(c, a) \mathcal{K}_a^c(\mathsf{T}_i | \mathsf{Q}_i, \mathsf{R}_i) &= \pi_{2N,0}(c, a) \sum_{q,r} |\bar{P}_a^c(0, q, r) - \bar{P}_a^c(1, q, r)| \\ &\stackrel{(a)}{=} \sum_{q,r} \left| \sum_b \mu(b) \sum_{v=0}^1 P_b^c(v, r) (P_a^b(v, q) - P_a^b(v+1, q)) \right| \\ &\stackrel{(b)}{\leq} \sum_{q,r,b} \mu(b) \left| \sum_{v=0}^1 P_b^c(v, r) (P_a^b(v, q) - P_a^b(v+1, q)) \right| \\ &= \sum_{q,r,b} \mu(b) |P_a^b(0, q) - P_a^b(1, q)| \cdot |P_b^c(0, r) - P_b^c(1, r)| \\ &= \sum_b \mu(b) \mathcal{K}_a^b(\mathsf{U}_i | \mathsf{Q}_i) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) \\ &\stackrel{(c)}{\leq} \psi(0) \sum_b \left(\pi_{2N,N}(c, b) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) \right) \left(\pi_{N,0}(b, a) \mathcal{K}_a^b(\mathsf{U}_i | \mathsf{Q}_i) \right) \\ &\stackrel{(d)}{\leq} \psi(0) \sum_b \pi_{2N,N}(c, b) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) \sum_{b'} \pi_{N,0}(b', a) \mathcal{K}_a^{b'}(\mathsf{U}_i | \mathsf{Q}_i),\end{aligned}$$

where (a) is by (21), (b) is by the triangle inequality, (c) is by (19), and (d) is by the inequality $\sum_j a_j b_j \leq \sum_j a_j \sum_{j'} b_{j'}$, which holds for $a_j, b_j \geq 0$. By (18), the sum over $a, c \in \mathcal{S}$ yields

$$\sum_{a,c} \pi_{2N,0}(c, a) \mathcal{K}_a^c(\mathsf{T}_i | \mathsf{Q}_i, \mathsf{R}_i) \leq \psi(0) \hat{\mathcal{K}}_n^2.$$

Next, let $\mathsf{B}_{n+1} = 1$. We have

$$\begin{aligned}\pi_{2N,0}(c, a) \mathcal{K}_a^c(\mathsf{V}_i | \mathsf{T}_i, \mathsf{Q}_i, \mathsf{R}_i) &= \pi_{2N,0}(c, a) \sum_{t,q,r} |\bar{P}_a^c(t, 0, q, r) - \bar{P}_a^c(t, 1, q, r)| \\ &\stackrel{(a)}{=} \sum_{t,q,r} \left| \sum_b \mu(b) (P_a^b(t, q) P_b^c(0, r) - P_a^b(t+1, q) P_b^c(1, r)) \right| \\ &\stackrel{(b)}{=} \frac{1}{2} \sum_{t,q,r} \left| \sum_b \mu(b) P_a^b(q) (P_b^c(0, r) - P_b^c(1, r)) \right. \\ &\quad \left. + \sum_b \mu(b) P_b^c(r) (P_a^b(t, q) - P_a^b(t+1, q)) \right| \\ &\stackrel{(c)}{\leq} \sum_{q,b} \mu(b) P_a^b(q) \left(\sum_r |P_b^c(0, r) - P_b^c(1, r)| \right) \\ &\quad + \sum_{r,b} \mu(b) P_b^c(r) \left(\sum_q |P_a^b(0, q) - P_a^b(1, q)| \right) \\ &= \sum_b \mu(b) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) + \sum_b \mu(b) \mathcal{K}_a^b(\mathsf{U}_i | \mathsf{Q}_i),\end{aligned}$$

where (a) is by (21), (b) is by the equality

$$(\alpha\beta - \gamma\delta) = ((\alpha + \gamma)(\beta - \delta) + (\beta + \delta)(\alpha - \gamma))/2,$$

which holds for any four numbers $\alpha, \beta, \gamma, \delta$, and (c) is by the triangle inequality. By (20),

$$\begin{aligned}\sum_{a,b,c} \mu(b) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) &= \sum_{b,c} \pi_{2N,N}(c, b) \mathcal{K}_b^c(\mathsf{V}_i | \mathsf{R}_i) = \hat{\mathcal{K}}_n, \\ \sum_{a,b,c} \mu(b) \mathcal{K}_a^b(\mathsf{U}_i | \mathsf{Q}_i) &= \sum_{a,b} \pi_{N,0}(b, a) \mathcal{K}_a^b(\mathsf{U}_i | \mathsf{Q}_i) = \hat{\mathcal{K}}_n.\end{aligned}$$

Thus,

$$\sum_{a,c} \pi_{2N,0}(c, a) \mathcal{K}_a^c(\mathsf{V}_i | \mathsf{T}_i, \mathsf{Q}_i, \mathsf{R}_i) \leq 2\hat{\mathcal{K}}_n. \quad \blacksquare$$

Theorem 10. Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. Then \mathcal{K}_n polarizes fast to 0 and \mathcal{Z}_n polarizes fast to 1 for any $\beta < 1/2$.

Proof: Fix $\beta < 1/2$. By Corollary 8 and (17), we can invoke Lemma 3 for $\hat{\mathcal{K}}_n$ with $E = 1/2$. Consequently, $\hat{\mathcal{K}}_n$ polarizes fast to 0, i.e.,

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\mathcal{K}}_n < 2^{-N^\beta}) &= \mathbb{P}(\hat{\mathcal{K}}_\infty = 0) \\ &= \mathbb{P}(\mathsf{H}_\infty = 1) = \mathcal{H}_*(X|Y).\end{aligned}$$

For any n , by Lemma 1 and Lemma 2,

$$1 - \mathcal{Z}_n \leq \sqrt{1 - \mathcal{Z}_n} \leq \mathcal{K}_n \leq \hat{\mathcal{K}}_n.$$

Thus, \mathcal{K}_n polarizes fast to 0. Moreover, $\mathbb{P}(\mathcal{Z}_n > 1 - 2^{-N^\beta}) \geq \mathbb{P}(\hat{\mathcal{K}}_n < 2^{-N^\beta})$. Taking limits, we obtain that $\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{Z}_n > 1 - 2^{-N^\beta}) \geq \mathcal{H}_*(X|Y)$.

On the other hand, by Corollary 6,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{Z}_n < 2^{-N^\beta}) = 1 - \mathcal{H}_*(X|Y).$$

Since $\mathbb{P}(\mathcal{Z}_n < 2^{-N^\beta}) + \mathbb{P}(\mathcal{Z}_n > 1 - 2^{-N^\beta}) \leq 1$ for any n , we must have

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathcal{Z}_n > 1 - 2^{-N^\beta}) = \mathcal{H}_*(X|Y). \quad \blacksquare$$

Corollary 11. Let (X_j, Y_j, S_j) , $j \in \mathbb{Z}$ be a FAIM process. Then $\mathsf{H}_n, \hat{\mathsf{H}}_n, \mathcal{Z}_n, \hat{\mathcal{Z}}_n, \mathcal{K}_n$, and $\hat{\mathcal{K}}_n$ polarize fast both to 0 and to 1 with any $\beta < 1/2$.

REFERENCES

- [1] E. Arkan, "Channel polarization: a method for constructing capacity-achieving codes for symmetric binary-input memoryless channels," *IEEE Trans. on Information Theory*, vol. 55, no. 7, pp. 3051–3073, July 2009.
- [2] E. Arkan and E. Telatar, "On the rate of channel polarization," in *Proc. IEEE Int. Sym. on Information Theory*, June 2009, pp. 1493–1495.
- [3] E. Şaşıoğlu, E. Telatar, and E. Arkan, "Polarization for arbitrary discrete memoryless channels," in *2009 IEEE Information Theory Workshop*, October 2009, pp. 144–148.
- [4] E. Şaşıoğlu, "Polar Coding Theorems for Discrete Systems," Ph.D. dissertation, IC, Lausanne, 2011.
- [5] S. B. Korada and R. L. Urbanke, "Polar codes are optimal for lossy source coding," *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1751–1768, April 2010.
- [6] E. Arkan, "Source polarization," in *2010 IEEE Int. Sym. on Information Theory*, June 2010, pp. 899–903.
- [7] J. Honda and H. Yamamoto, "Polar coding without alphabet extension for asymmetric models," *IEEE Transactions on Information Theory*, vol. 59, no. 12, pp. 7829–7838, December 2013.
- [8] M. Alsan, "Re-proving Channel Polarization Theorems," Ph.D. dissertation, IC, Lausanne, 2015. [Online]. Available: https://infoscience.epfl.ch/record/203886/files/EPFL_TH6403.pdf
- [9] E. Şaşıoğlu and I. Tal, "Polar coding for processes with memory," in *2016 IEEE Int. Sym. on Information Theory (ISIT)*. IEEE, 2016, pp. 225–229.
- [10] R. Wang, J. Honda, H. Yamamoto, R. Liu, and Y. Hou, "Construction of polar codes for channels with memory," in *2015 IEEE Information Theory Workshop*, October 2015, pp. 187–191.
- [11] B. Shuvail and I. Tal, "Fast polarization for processes with memory," 2017. [Online]. Available: arXiv:1710.02849
- [12] I. Tal, "A simple proof of fast polarization," *IEEE Transactions on Information Theory*, vol. 63, no. 12, pp. 7617–7619, Dec 2017.
- [13] T. Cover and J. Thomas, *Elements of Information Theory*, 2nd ed. Wiley, 2006.

V. ADDENDUM

There is a subtle error in the statement of Lemma 1. However, the consequence (4) and thus the rest of the paper are correct. For the correct statement of Lemma 1, see our journal paper.